

# The Physics of Waves

*Physics 15c*

Matthew D. Schwartz

Department of Physics  
Harvard University

Spring 2016

# Lecture 1: Simple Harmonic Oscillators

## 1 Introduction

The simplest thing that can happen in the physical universe is nothing. The next simplest thing, which doesn't get too far away from nothing, is an oscillation about nothing. This course studies those oscillations. When many oscillators are put together, you get waves.

Almost all physical processes can be explained by breaking them down into simple building blocks and putting those blocks together. As we will see in this course, oscillators are the building blocks of a tremendous diversity of physical phenomena and technologies, including musical instruments, antennas, patriot missiles, x-ray crystallography, holography, quantum mechanics, 3D movies, cell phones, atomic clocks, ocean waves, gravitational waves, sonar, rainbows, color perception, prisms, soap films, sunglasses, information theory, solar sails, cell phone communication, molecular spectroscopy, acoustics and lots more. Many of these topics will be covered first in lab where you will explore and uncover principles of physics on your own.

The key mathematical technique to be mastered through this course is the **Fourier transform**. Fourier transforms, and **Fourier series**, play an absolutely crucial role in almost all areas of modern physics. I cannot emphasize enough how important Fourier transforms are in physics.

The first couple of weeks of the course build on what you've covered in 15a (or 16 or 11a or AP50) – balls and springs and simple oscillators. These are described by the differential equation for the **damped, driven oscillator**:

$$\frac{d^2x(t)}{dt^2} + \gamma \frac{dx(t)}{dt} + \omega_0^2 x(t) = \frac{F(t)}{m} \quad (1)$$

Here  $x(t)$  is the displacement of the oscillator from equilibrium,  $\omega_0$  is the natural angular frequency of the oscillator,  $\gamma$  is a damping coefficient, and  $F(t)$  is a driving force. We'll start with  $\gamma=0$  and  $F=0$ , in which case it's a simple harmonic oscillator (Section 2). Then we'll add  $\gamma$ , to get a damped harmonic oscillator (Section 4). Then add  $F(t)$  (Lecture 2).

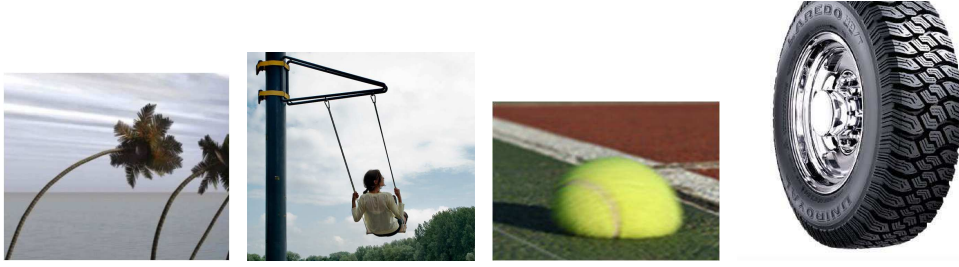
The damped, driven oscillator is governed by a **linear** differential equation (Section 5). Linear equations have the nice property that you can add two solutions to get a new solution. We will see how to solve them using complex exponentials,  $e^{i\alpha}$  and  $e^{-i\alpha}$ , which are linear combinations of sines and cosines (Section 6). A review of complex numbers is given in Section 7. Studying multiple coupled oscillators will lead to the concept of **normal modes**, which lead naturally to the **wave equation**, the Fourier series, and the Fourier transform (future lectures).

## 2 Why waves? Why oscillators?

Recall **Hooke's law**: if you displace a spring a distance  $x$  from its equilibrium position, the restoring force will be  $F = -kx$  for some constant  $k$ . You probably had this law told to you in high school or 15a or wherever. Maybe it's an empirical fact, deduced from measuring springs, maybe it was just stated as true. Why is it true? Why does Hooke's law hold?

To derive Hooke's law, you might imagine you need a microscopic description of a spring (what is it made out of, how does it bend, how are the atoms arranged, etc.). Indeed, if you hope to compute  $k$ , yes, absolutely, you need all of this. In fact you need so much detail that generally it's impossible to compute  $k$  in any real spring. But also generally, we don't care to compute  $k$ , we just measure it. That's not the point. We don't want to compute  $k$ . What we want to know is why is the force is proportional to displacement. Why is Hooke's law true?

First of all, it is true. Hooke's law applies not just to springs, but to just about everything:



**Figure 1.** The restoring force for pretty much anything (bending trees, swings, balls, tires, etc) is **linear** close to equilibrium.

You can move any of these systems, or pretty much anything else around you, a little bit away from its equilibrium and it will want to come back. The more you move it, the stronger the restoring force will be. And often to an excellent approximation, the distance and force are directly proportional.

To derive Hooke's law, we just need a little bit of calculus. Let's say we displace some system, a spring or a tire or whatever a distance  $x$  from its equilibrium and measure the function  $F(x)$ . We define  $x = 0$  as the equilibrium point, so by definition,  $F(0) = 0$ . Then, we can use Taylor's theorem

$$F(x) = F(0) + xF'(0) + \frac{1}{2}x^2F''(0) + \dots \quad (2)$$

Now  $F(0) = 0$  and  $F'(0)$  and  $F''(0)$  etc are just fixed numbers. So no matter what these numbers are, we can always find an  $x$  small enough so that  $F'(0) \gg \frac{1}{2}x^2F''(0)$ . Then we can neglect the  $\frac{1}{2}x^2F''(0)$  term compared to the  $xF'(0)$  term. Similarly, we can always take  $x$  small enough that all of the higher derivative terms are as small as we want. And therefore,

$$F(x) = -kx \quad (3)$$

with  $k = -F'(0)$ . We have just derived Hooke's law! Close enough to equilibrium, the restoring force for *anything* will be proportional to the displacement. Since  $y = -kx$  is the equation for a line, we say systems obeying Hooke's law are linear. Thus, everything is linear close to equilibrium. More about linearity in Section 5.

You might also ask, why does  $F$  depend only on  $x$ ? Well, what else could it depend on? It could, for example, depend on velocity. Wind resistance is an example of a velocity-dependent force. However, since we are assuming that the object is close to equilibrium, its speed must be small (or else our assumption would quickly be violated). So  $\dot{x}$  is small. Thus we can Taylor expand in  $\dot{x}$  as well

$$F(x, \dot{x}) = x \left. \frac{\partial F(x, \dot{x})}{\partial x} \right|_{x=\dot{x}=0} + \dot{x} \left. \frac{\partial F(x, \dot{x})}{\partial \dot{x}} \right|_{x=\dot{x}=0} + \dots \quad (4)$$

where the terms  $\dots$  are higher order in  $x$  or  $\dot{x}$ , so they are subleading close to equilibrium. Writing  $\left. \frac{\partial F(x, \dot{x})}{\partial \dot{x}} \right|_{x=\dot{x}=0} = -m\gamma$  we then have

$$F(x) = -kx - m\gamma\dot{x} \quad (5)$$

Then  $F = ma$  with  $a = \ddot{x}$  gives

$$\frac{d^2x(t)}{dt^2} + \gamma \frac{dx(t)}{dt} + \omega_0^2 x(t) = 0 \quad (6)$$

as in Eq. (1) with  $\omega_0 = \sqrt{\frac{k}{m}}$ .  $\gamma$  is called a damping coefficient, since the velocity dependence tends to slow the system down (as we will see).

The other piece of Eq. (1), labeled  $F(t)$ , is the driving force. It represents the action of something external to the system, like a woman pushing the swing with the girl on it, or the car tire being compressed by the car.

### 3 Simple harmonic motion

We have seen that Eq. (1) describes universally any system close to equilibrium. Now let's solve it. First, take  $\gamma = 0$ . Then Eq. (1) becomes

$$\frac{d^2}{dt^2}x(t) + \omega_0^2 x(t) = 0 \quad (7)$$

For a spring,  $\omega_0 = \sqrt{\frac{k}{m}}$ , for a pendulum  $\omega_0 = \sqrt{\frac{g}{L}}$ . Other systems have different expressions for  $\omega_0$  in terms of the relevant physical parameters.

We can solve this equation by hand, by plugging into Mathematica, or just by guessing. Guessing is often the easiest. So, we want to guess a function whose second derivative is proportional to itself. You know at least two functions with this property: sine and cosine. So let us write as an *ansatz* (ansatz is a sciency word for “educated guess”):

$$x(t) = A \sin(\omega t) + B \cos(\omega t) \quad (8)$$

This solution has 3 free parameters  $A$ ,  $B$  and  $\omega$ . Plugging in to Eq. (7) gives

$$-\omega^2[A \sin(\omega t) + B \cos(\omega t)] + \omega_0^2[A \sin(\omega t) + B \cos(\omega t)] = 0 \quad (9)$$

Thus,

$$\omega = \omega_0 \quad (10)$$

That is, the angular frequency  $\omega$  of the solution must be the parameter  $\omega_0 = \sqrt{\frac{k}{m}}$  in the differential equation. We get no constraint on  $A$  and  $B$ .

$\omega$  is called the **angular frequency**. It has units of radians per second. The **frequency** is

$$\nu = \frac{\omega}{2\pi} \quad (11)$$

units of 1/sec. The solution  $x(t)$  we found goes back to itself after  $t \rightarrow t + T$  where

$$T = \frac{1}{\nu} = \frac{2\pi}{\omega} \quad (12)$$

is the **period**.  $T$  has units of seconds. The function  $x(t) = A \sin(\omega t) + B \cos(\omega t)$  satisfies  $x(t) = x(t + nT)$  for any integer  $n$ . In other words, the solutions *oscillate*!

$A$  and  $B$  are the **amplitudes** of the oscillation. They can be fixed by boundary conditions. For example, you specify the position and velocity at any given time, you can determine  $A$  and  $B$ . To be concrete, suppose we start with  $x(0) = 1m$  and  $x'(0) = 2\frac{m}{s}$ . Then,

$$1m = x(0) = A \sin(\omega 0) + B \cos(\omega 0) = B \quad (13)$$

$$2\frac{m}{s} = x'(0) = \omega A \cos(\omega 0) - \omega B \sin(\omega 0) = \omega A \quad (14)$$

So we find  $A = \frac{2}{\omega} \frac{m}{s}$  and  $B = 1m$ .

Keep in mind that the angular frequency  $\omega$  is *not* fixed by boundary conditions. It is determined by the physical problem:  $\omega = \sqrt{\frac{k}{m}}$  where  $k = -F'(0)$  and  $m$  is the mass of the thing oscillating. That is why if you start a pendulum from any height and give it any sort of initial kick, it will oscillate with the same frequency.

Another representation of the general solution  $x(t) = A \sin(\omega t) + B \cos(\omega t)$  is often convenient. Instead of using  $A$  and  $B$  we can write

$$x(t) = C \sin(\omega t + \phi) \quad (15)$$

using trig identities, we find

$$C \sin(\omega t + \phi) = C \cos(\phi) \sin(\omega t) + C \sin(\phi) \cos(\omega t) \quad (16)$$

and so

$$A = C \cos(\phi) \quad B = C \sin(\phi) \quad (17)$$



Thus we can swap the amplitudes  $A$  and  $B$  for the sine and cosine components for a single amplitude  $C$  and a phase  $\phi$ .

## 4 Damped oscillators

A damped oscillator dissipates its energy, returning eventually to the equilibrium  $x(t) = \text{const}$  solution. When the object is at rest, the damping force must vanish. For small velocities, the damping force should be proportional to velocity:  $F = -\gamma \frac{dx}{dt}$  with  $\gamma$  some constant. Contributions to the force proportional to higher powers of velocity, like  $F = -\kappa \left(\frac{dx}{dt}\right)^2$  will be suppressed when the object is moving slowly. Thus the generic form for damped motion close to equilibrium is

$$\frac{d^2x}{dt^2} + \gamma \frac{dx}{dt} + \omega_0^2 x = 0 \quad (18)$$

Indeed, this equation describes a great many physical systems: vibrating strings, sound waves, etc. Basically everything we study in this course will have damping.

Neither  $\sin(\omega t)$  nor  $\cos(\omega t)$  solve the damped oscillator equation. Sines and cosines are proportional to their second derivatives, but here we have also a first derivative. Since  $\frac{d}{dt}\sin(\omega t) \propto \cos(\omega t)$  and vice versa, neither sines nor cosines alone will solve this equation. However, the exponential function is proportional to its *first* derivative. Thus exponentials are a natural guess, and indeed they will work.

So, let's try plugging

$$x(t) = Ce^{\alpha t} \quad (19)$$

into Eq. (18). We find

$$\alpha^2 Ce^{\alpha t} + \gamma \alpha Ce^{\alpha t} + \omega_0^2 Ce^{\alpha t} = 0 \quad (20)$$

Dividing out by  $Ce^{\alpha t}$  we have reduced this to an algebraic equation:

$$\alpha^2 + \gamma \alpha + \omega_0^2 = 0 \quad (21)$$

The solutions are

$$\alpha = -\frac{\gamma}{2} \pm \sqrt{\left(\frac{\gamma}{2}\right)^2 - \omega_0^2} \quad (22)$$

And therefore the general solution to the damped oscillator equation is

$$x(t) = e^{-\frac{\gamma}{2}t} \left( C_1 e^{t\sqrt{\left(\frac{\gamma}{2}\right)^2 - \omega_0^2}} + C_2 e^{-t\sqrt{\left(\frac{\gamma}{2}\right)^2 - \omega_0^2}} \right) \quad (23)$$

The cases when  $\gamma > 2\omega_0$ ,  $\gamma = 2\omega_0$  and  $\gamma < 2\omega_0$  give very different physical behavior.

### 4.1 Underdamping: $\gamma < 2\omega_0$

The case  $\gamma < 2\omega_0$  includes the case when  $\gamma = 0$ . For  $\gamma = 0$  the damping vanishes and we should regain the oscillator solution. Increasing  $\gamma$  from zero should slowly damp the oscillator. Let's see how this works mathematically.

Since  $\gamma < 2\omega_0$  then

$$\omega_u = \sqrt{\omega_0^2 - \left(\frac{\gamma}{2}\right)^2} \quad (24)$$

is a real number. In terms of  $\omega_u$ , the general solution is then

$$x(t) = e^{-\frac{\gamma}{2}t} (C_1 e^{i\omega_u t} + C_2 e^{-i\omega_u t}) \quad (25)$$

Since  $x(t)$  must be real, we must also have  $C_1 = C_2^*$ . Thus we can write

$$C_1 = \frac{1}{2} A e^{i\phi}, \quad C_2 = \frac{1}{2} A e^{-i\phi} \quad (26)$$

for two real constants  $A$  and  $\phi$ . This leads to

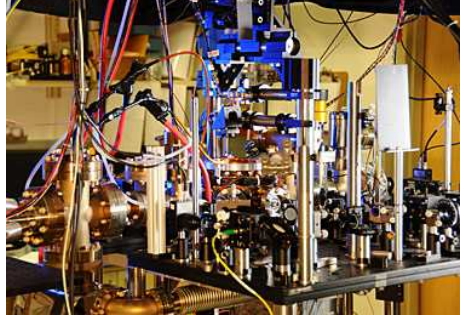


$$x(t) = A e^{-\frac{\gamma}{2}t} \cos(\omega_u t + \phi) \quad (27)$$

Thus we see that in the underdamped case, the object still oscillates, but at an angular frequency  $\omega_u = \sqrt{\omega_0^2 - (\frac{\gamma}{2})^2}$  and the amplitude slowly goes down over time.

Note that both  $\omega_0$  and  $\gamma$  have dimensions of  $\frac{1}{\text{seconds}}$ . Their relative size determines how much the amplitude gets damped in a single oscillation. To quantify this, we define the **Q-factor** (or **Q-value**) as

$$Q \equiv \frac{\omega_0}{\gamma} \quad (28)$$

The smaller the  $Q$  the more the damping.  $Q$  stands for quality. The higher  $Q$  is, the higher quality, and the less resistance/friction/damping is involved. For example, a tuning fork vibrates for a long time. It is a very high quality resonator with  $Q \sim 1000$ . Here are some examples:

		
Atomic clock: $Q \approx 10^{11}$	Tuning fork: $Q \approx 1000$	Silly putty: $Q \sim 0.01$

**Figure 2.** Some  $Q$ -factors

$Q$  is roughly the number of complete oscillations a system has gone through before it's amplitude goes down by a factor of around 20. To see this, note that due to the  $\cos(\omega_u t)$  factor, it takes a time  $t_Q = \frac{2\pi}{\omega_u} Q$  to go through  $Q$  cycles. Then due to the  $e^{-\frac{\gamma}{2}t}$  factor, the amplitude has decayed by a factor of

$$\exp\left(-\frac{\gamma}{2}t\right) = \exp\left(-\frac{\gamma}{2}Q \frac{2\pi}{\omega_u}\right) = \exp\left(-\frac{\omega_0}{\omega_u}\pi\right) \approx \exp(-\pi) = 0.043 \quad (29)$$

In the next-to-last step, we have used that  $\omega_u \approx \omega_0$  when  $Q \gg 1$ . (If  $Q$  is not large, then the system is highly damped and counting oscillations is not so useful). Since  $0.043 \approx \frac{1}{23}$  we get the  $\frac{1}{20}$  rule.

## 4.2 Overdamping: $\gamma > 2\omega_0$

In the over damped case,  $\gamma > 2\omega_0$ , then  $(\frac{\gamma}{2})^2 - \omega_0^2$  is positive so the roots in Eq. (22) are real. Thus the general solution is simply

$$x(t) = C_1 e^{-u_1 t} + C_2 e^{-u_2 t} \quad (30)$$

with

$$u_1 = \frac{\gamma}{2} + \sqrt{\left(\frac{\gamma}{2}\right)^2 - \omega_0^2}, \quad (31)$$

and

$$u_2 = \frac{\gamma}{2} - \sqrt{\left(\frac{\gamma}{2}\right)^2 - \omega_0^2} \quad (32)$$

Both solutions have exponential decay. Since  $u_1 > u_2$ , the  $u_1$  solution will die away first, leaving the  $u_2$  solution. Overdamped systems have  $Q < \frac{1}{2}$ .

### 4.3 Critical damping: $\gamma = 2\omega_0$

In the critically damped case, the two solutions in Eq. (23) reduce to one:

$$x(t) = Ce^{-\omega_0 t} \quad (33)$$

What happened to the other solution? That is, a second-order differential equation is supposed to have two independent solutions, but we have only found one. To find the other solution, let's look at the damped oscillator equation again, but set  $\gamma = 2\omega_0$  to begin with. Then the equation is

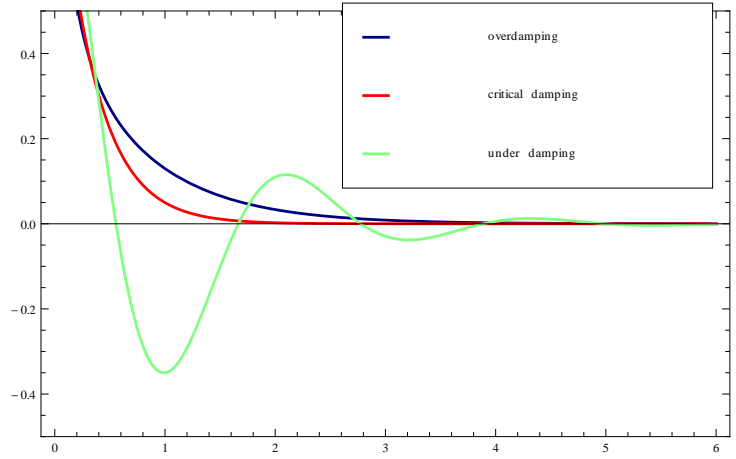
$$\frac{d^2x}{dt^2} + 2\omega_0 \frac{dx}{dt} + \omega_0^2 x = 0 \quad (34)$$

Solving this with Mathematica, we find the general solution is

$$x(t) = (C + Bt)e^{-\omega_0 t} \quad (35)$$

You should check yourself that this ansatz satisfies Eq. (34).

A comparison of over-damping, underdamping and critical damping is shown in Figure 3. One thing to note is that the critically damped curve goes to zero faster than the overdamped curve! Can you think of an application for which you'd want a critically damped oscillator?



**Figure 3.** Comparison of underdamping, overdamping, and critical damping. We have taken  $\omega_0 = 3$  and  $\gamma = 8, 2$  and  $6$ .

## 5 Linearity

The oscillator equation we have been solving has a very important property: **linearity**. Differential equations with at most single powers of  $x$  are **linear differential equations**. For example,

$$\frac{d^2x}{dt^2} + \omega^2 x = 0 \quad (36)$$

is linear. If there is no constant ( $x^0$  term), the differential equations are **homogeneous**.

Linearity is important because it implies that if  $x_1(t)$  and  $x_2(t)$  are solutions to the equations of motion for a homogenous linear system then

$$x(t) = x_1(t) + x_2(t) \quad (37)$$

is also a solution. Let's check this for Eq. (36). By assumption  $x_1(t)$  and  $x_2(t)$  satisfy:

$$\frac{d^2x_1}{dt^2} + \omega^2 x_1 = 0 \quad (38)$$

$$\frac{d^2x_2}{dt^2} + \omega^2 x_2 = 0 \quad (39)$$

Adding these equations, we find

$$\frac{d^2x}{dt^2} + \omega^2x = 0 \quad (40)$$

So that  $x$  satisfies Eq. (36) as well. So, *solutions to homogeneous linear differential equations add.*

## 5.1 Examples of linear systems

The damped oscillator in Eq. (18) is a linear system:  $\frac{d^2x}{dt^2} + \gamma\frac{dx}{dt} + \omega_0^2x = 0$ .

If you have a string of tension  $T$  and mass density  $\mu$  going along the  $x$  direction, then its displacement in a transverse direction  $y(x, t)$  satisfies

$$\mu \frac{\partial^2}{\partial t^2} y(x, t) - T \frac{\partial^2}{\partial x^2} y(x, t) = 0 \quad (41)$$

This is the **wave equation**. It is linear. (We'll derive this in a couple of weeks).

Electromagnetic waves are described by Maxwell's equations. With a little work, you can combine two of Maxwell's equations into an equation for the electric field  $\vec{E}(x, t)$  of the form (we'll derive this soon too):

$$c^2 \frac{\partial^2}{\partial t^2} \vec{E}(x, t) - \frac{\partial^2}{\partial x^2} \vec{E}(x, t) = 0 \quad (42)$$

with  $c$  the speed of light. Thus each component of  $\vec{E}$  satisfies the **wave equation**.

Sound waves, water waves, etc., all satisfy linear differential equations.

## 5.2 Forced oscillation

What happens if the equation is not homogeneous? For example, what if we had

$$\frac{d^2x}{dt^2} = F_1(t) \quad (43)$$

Here,  $F(t)$  represents some force from a motor or the wind or someone pushing you on a swing or the electromagnetic waves from the cell-phone tower transmitting that critical text message to you during class.

It is in general hard to solve this differential equation. But let's imagine we can do it and find a function  $x_1(t)$  which satisfies

$$\frac{d^2x_1}{dt^2} = F_1(t) \quad (44)$$

Say this is your motion on a swing when you are being pushed. Now say some friend comes and pushes you too. Then

$$\frac{d^2x}{dt^2} = F_1(t) + F_2(t) \quad (45)$$

The amazing thing about linearity is that if we can find a solution  $x_2(t)$  which satisfies

$$\frac{d^2x_2}{dt^2} = F_2(t) \quad (46)$$

Then  $x = x_1 + x_2$  satisfies

$$\frac{d^2x}{dt^2} = F_1(t) + F_2(t) \quad (47)$$

This is *extremely important*. It is the key to this whole course. Really complicated systems are solvable by simpler systems, as long as the equations are linear.

In contrast, we cannot add solutions to *nonlinear* equations. For example, suppose we want to solve

$$a \frac{d^2}{dt^2} x + b \frac{d}{dt} x^2 = F_1(t) + F_2(t) \quad (48)$$

Although we can solve for  $x_1$  and  $x_2$  produced by the two forces separately, it does not then follow that  $x = x_1 + x_2$  is a solution with the combined force is present. The  $x^2$  term couples the  $x_1$  and  $x_2$  solutions together, so there is interference.

All the systems we study in this course will be linear systems. An important example is electromagnetism (Maxwell's equations are linear). Suppose you are making some radio waves in your radio station with  $F_1 = \sin(7t)$  and some MIT student is making waves in her station with  $F_2 = \sin(4t)$ . Then  $x_1(t) = \frac{1}{49}\sin(7t)$  satisfies Eq. (44) and  $x_2(t) = \frac{1}{16}\sin(4t)$  satisfies Eq. (46). We then immediately conclude that

$$x(t) = \frac{1}{49}\sin(7t) + \frac{1}{16}\sin(4t) \quad (49)$$

must satisfy Eq. (45). We just add the oscillations! Thus if there are radio waves at frequency  $\nu = 89.9$  MHz flying around and frequency  $\nu = 90.3$  MHz flying around, they don't interfere with each other.

That explains why we can tune our radio – because electromagnetism is linear, we can add radio waves. There is no interference! The different frequencies don't mix with each other. All we have to do is get our radio to extract the coefficient of the  $\sin(7t)$  oscillation from  $x(t)$ . Then we will get only the output from our radio station. As we will see, you can always find out which frequencies are present with which amplitudes using **Fourier decomposition**. We'll come back to this soon.

### 5.3 Summary

Linearity is a really important concept in physics. The definition of linearity is that all terms in a differential equation for  $x(t)$  have at most one power of  $x(t)$ . So  $\frac{d^3}{dt^3}x(t) = 0$  is linear, but  $\frac{d}{dt}x(t)^2$  is nonlinear.

For linear systems, one can add different solutions and still get a solution. This lets us break the problem down to easier subproblems.

Linearity does not only let us solve problems simply, but it is also a universal feature of physical systems. Whenever you are close to a static solution  $x(t) = x_0 = \text{constant}$ , the equations for deviations around this solution will be linear. To see that, we again use Taylor theorem. We shift by  $x(t) \rightarrow x(t) - x_0$  so the equilibrium point is now  $x(t) = 0$ . Then no matter how complicated and nonlinear the exact equations of motion for the system are, when  $x - x_0 \ll 1$ , the linear term, proportional to  $x - x_0$  will dominate. For example, if we had

$$\frac{d^2}{dt^2} \frac{x}{x^2 - 2} e^{-x^4} + \frac{d}{dt} x^7 + (x^2 - 4)\sin^3(x) = 0 \quad (50)$$

then  $x(t) = x_0 = 0$  is a solution. For  $x \ll 1$  this simplifies to

$$-\frac{1}{2} \frac{d^2 x}{dt^2} - 4x = 0 \quad (51)$$

which is again linear (it's the oscillator equation again).

## 6 Solving general linear systems

At this point, we've defined linearity, argued that it should be universal for small deviations from equilibrium, and showed how it can help us combine solutions to a differential equation. Now we will see how to solve general linear differential equations.

### 6.1 Exponentials, sines and cosines

A general linear equation has a bunch of derivatives with respect to time acting on a single function  $x$ :

$$\cdots + a_3 \frac{d^3}{dt^3} x + a_2 \frac{d^2}{dt^2} x + a_1 \frac{d}{dt} x + a_0 x = F(t) \quad (52)$$

Let's first consider the case when  $F = 0$ . A really easy way to solve these equations for  $F = 0$  is to consider solutions for which all the derivatives are proportional to each other. What is a function with this property? Sines and cosines have derivatives proportional to themselves:  $\frac{d^2}{dt^2}\sin(\omega t) = -\omega^2\sin(\omega t)$ , but only *second* (or even numbers of) derivatives. A function with *all* of its derivatives proportional to itself is the exponential:  $x(t) = Ce^{\alpha t}$

$$\frac{d}{dt}Ce^{\alpha t} = C\alpha e^{\alpha t} \quad (53)$$

As an example, let's try this Ansatz into our oscillator equation, Eq. (8):

$$\frac{d^2x}{dt^2} = -\omega^2 x \quad (54)$$

Plugging in  $x(t) = Ce^{\alpha t}$  gives

$$\alpha^2 e^{\omega t} = -\omega^2 e^{\omega t} \quad (55)$$

which implies  $\alpha = \sqrt{-\omega^2}$  or  $\alpha = \pm i\omega$ ,

Thus the solutions are

$$x(t) = C_1 e^{i\omega t} + C_2 e^{-i\omega t} \quad (56)$$

Recalling that

$$\sin(\omega t) = \frac{e^{i\omega t} - e^{-i\omega t}}{2i}, \quad \cos(\omega t) = \frac{e^{i\omega t} + e^{-i\omega t}}{2} \quad (57)$$

We can also write

$$x(t) = i(C_1 - C_2)\sin(\omega t) + (C_1 + C_2)\cos(\omega t) \quad (58)$$

In summary

- Sines and cosines are useful if you have only 2nd derivatives
- Exponentials work for any number of derivatives

Now, if  $F \neq 0$ , then a simple exponential will not obviously be a solution. The key, however, is that we can *always* write any function  $F(t)$  on the interval  $0 < t \leq T$  as a sum of exponentials

$$F(t) = \sum_{n=-\infty}^{\infty} a_n e^{2\pi i n \frac{t}{T}} \quad (59)$$

for some coefficients  $a_n$ . This is called a **Fourier decomposition**. Since we can solve the equation as if  $F(t) = e^{2\pi i n \frac{t}{T}}$  for a fixed  $n$ , we can then add solutions using *linearity* to find a solution with the original  $F(t)$ . Don't worry about understanding this now – it's just a taste of what's to come.

## 6.2 Relating $e^{ix}$ to $\sin(x)$ and $\cos(x)$

Suppose you didn't know that  $\sin(\omega t) = \frac{e^{i\omega t} - e^{-i\omega t}}{2i}$ . How could you derive this?

One way is using the fact that we solved the oscillator equation

$$\frac{d^2x}{dt^2} + \omega^2 x = 0 \quad (60)$$

two ways. On the one hand we found

$$x(t) = A \sin(\omega t) + B \cos(\omega t) \quad (61)$$

and on the other hand we found

$$x(t) = C_1 e^{i\omega t} + C_2 e^{-i\omega t} \quad (62)$$

Since the differential equation is second order (two derivatives), the solution is given uniquely once two boundary conditions are set. Conversely, if we know the solution we can work out the boundary conditions. For example, if the solution were  $x(t) = \sin(\omega t)$  then  $x(0) = 0$  and  $x'(0) = \omega$ . Plugging in  $x(0) = 0$  to the exponential solution implies

$$C_1 + C_2 = 0 \quad (63)$$

plugging in  $x'(0) = \omega$  implies

$$i\omega C_1 - i\omega C_2 = \omega \quad (64)$$

The solution to these two equations is  $C_1 = -\frac{1}{2i}$  and  $C_2 = \frac{1}{2i}$ , as you can easily check. We thus conclude that the two solutions are exactly equal and so

$$\sin(\omega t) = \frac{e^{i\omega t} - e^{-i\omega t}}{2i} \quad (65)$$

Similarly

$$\cos(\omega t) = \frac{e^{i\omega t} + e^{-i\omega t}}{2} \quad (66)$$

You should have these relationships memorized – we will use them a lot.

## 7 Complex numbers (mathematics)

Complex numbers are a wonderful invention. They make complicated equations look really simple. Being able to take the square root of anything is unbelievably helpful.

To see how important complex numbers are for solving equations, consider how sophisticated mathematics needs to be to solve some equations. The equation

$$3x - 4 = 0 \quad (67)$$

has a solution  $x = \frac{4}{3}$  which is a simple **rational number** (rational numbers can be written as ratios of whole numbers  $0, 1, -1, 2, -2, \dots$ ).

To solve

$$x^2 - 2 = 0 \quad (68)$$

we need **irrational numbers**:  $x = \sqrt{2}$ . Such numbers cannot be written as ratios of whole numbers.

To solve

$$x^2 + 4 = 0 \quad (69)$$

we need **complex numbers**. The solutions are  $x = \pm 2i$ , with  $i = \sqrt{-1}$ .

Now the punch line: to solve

$$ax^3 + bx^2 + cx + d = 0 \quad (70)$$

we still need *only* complex numbers. Complex numbers are the end of the road. Any polynomial equation can be solved with complex numbers.

$$ax^3 + bx^2 + cx + d = (x - r_1)(x - r_2)(x - r_3) = 0 \quad (71)$$

for some  $r_i \in \mathbb{C}$ .

Exponentials are for linear differential equations what complex numbers are for algebraic equations. Any linear differential equation can be solved by exponentials. Say we had

$$a \frac{d^3}{dt^3} x(t) + b \frac{d^2}{dt^2} x(t) + c \frac{d}{dt} x(t) + d x(t) = 0 \quad (72)$$

We can factor this into

$$\left(\frac{d}{dt} - r_1\right)\left(\frac{d}{dt} - r_2\right)\left(\frac{d}{dt} - r_3\right)x(t) = 0 \quad (73)$$

Thus if

$$\left(\frac{d}{dt} - r_3\right)x(t) = 0 \quad (74)$$

Then we have a solution. The solution is therefore a product of factors like

$$x(t) = e^{ir_3 t} \quad (75)$$

So we're always going to have exponential solutions to linear equations.

## 7.1 Complex number arithmetic

I hope you're already familiar with complex numbers from your math classes. If not, here's a quick review.

We can write any complex number as

$$z = a + bi \quad (76)$$

with  $a$  and  $b$  real. Then

$$z_1 + z_2 = a_1 + a_2 + (b_1 + b_2)i \quad (77)$$

and

$$z_1 \cdot z_2 = (a_1 + b_1i)(a_2 + b_2i) = a_1a_2 + b_1a_2i + b_2a_1i + b_1b_2i^2 \quad (78)$$

$$= (a_1a_2 - b_1b_2) + (b_1a_2 + b_2a_1)i \quad (79)$$

It's helpful to define **complex conjugation**  $i \rightarrow -i$ . In fact, we could have used  $-i$  instead of  $i$  from the beginning. We define

$$\bar{z} = a - bi \quad (80)$$

as the **complex conjugate** of a complex number  $z = a + bi$ .

Then

$$z\bar{z} = (a + bi)(a - bi) = a^2 + b^2 \in \mathbb{R} \quad (81)$$

The trick to dividing complex numbers is to use that  $z\bar{z} \in \mathbb{R}$ :

$$\frac{1}{z} = \frac{\bar{z}}{z\bar{z}} = \frac{a - bi}{a^2 + b^2} \quad (82)$$

That is,

$$\frac{a_2 + b_2i}{a_1 + b_1i} = (a_2 + b_2i) \frac{1}{a_1 + b_1i} = (a_2 + b_2i) \frac{a_1 - b_1i}{a_1^2 + b_1^2} = \frac{a_1a_2 + b_1b_2}{a_1^2 + b_1^2} + \frac{a_1b_2 - a_2b_1}{a_1^2 + b_1^2}i \quad (83)$$

For functions, we usually write  $f^*$  instead of  $\bar{f}$  for complex conjugation. For any function  $f(x) \in \mathbb{C}$ , we have

$$[f(x)][f(x)]^* \in \mathbb{R} \quad (84)$$

This is easy to see using that the conjugate of a product of complex numbers is the product of the conjugates:

$$(ff^*)^* = f^*f^{**} = f^*f = ff^* \quad (85)$$

Since  $ff^*$  is invariant under complex conjugation it must be real.

Any complex number can be also written as

$$z = re^{i\theta} = a + bi \quad (86)$$

to relate  $r$  and  $\theta$  to  $a$  and  $b$  we use

$$\bar{z} = re^{-i\theta} \quad (87)$$

$$a = \frac{z + \bar{z}}{2} = r \left( \frac{e^{i\theta} + e^{-i\theta}}{2} \right) = r \cos \theta \quad (88)$$

$$b = \frac{z - \bar{z}}{2i} = r \left( \frac{e^{i\theta} - e^{-i\theta}}{2i} \right) = r \sin \theta \quad (89)$$

$$z\bar{z} = r^2 = a^2 + b^2 \quad (90)$$

$r$  is sometimes called the **modulus** of a complex number and  $\theta$  the **phase** of a complex number.



# Lecture 2: Driven oscillators

## 1 Introduction

We started last time to analyze the equation describing the motion of a damped-driven oscillator:

$$\frac{d^2x}{dt^2} + \gamma \frac{dx}{dt} + \omega_0^2 x = F(t) \quad (1)$$

For small damping  $\gamma \ll \omega_0$ , we found solutions for  $F(t) = 0$  of the form

$$x(t) = A e^{-\frac{\gamma}{2}t} \cos(\omega_0 t + \phi) \quad (2)$$

where the amplitude  $A$  and the phase  $\phi$  are determined by initial conditions. Now we will see how to deal with  $F(t)$ .

We found the damped solution by guessing that an exponential  $x(t) = A e^{\alpha t}$  should work, since its derivatives are all proportional to itself. Plugging this ansatz in with  $F(t)$  we find

$$A e^{\alpha t} (\alpha^2 + \gamma \alpha + \omega_0^2) = F(t) \quad (3)$$

This will clearly not be solved for constant  $\alpha$  unless  $F(t)$  happens to be of the form  $e^{\alpha t}$ . The trick to solving this equation is to use linearity.

Let us suppose that we can write

$$F(t) = \sum_j c_j \cos(\omega_j t) \quad (4)$$

where  $c_j$  are real numbers. It may seem that only a handful of functions can be written this way, but actually *any periodic function* can be written as in Eq. (4) or as

$$F(t) = \sum_j [a_j \sin(\omega_j t) + b_j \cos(\omega_j t)] \quad (5)$$

with real numbers  $a_j$  and  $b_j$ . This truly remarkable fact is known as Fourier's theorem, and we will study it soon. For now, let us just take Eq. (4) as given.

Once  $F(t)$  is written as a sum of cosines, we can solve the differential equation for each cosine separately then add them. By linearity we then get a solution to the original equation. That is, if we can find functions  $x_j(t)$  satisfying

$$\frac{d^2x_j}{dt^2} + \gamma \frac{dx_j}{dt} + \omega_0^2 x_j = \cos(\omega_j t) \quad (6)$$

Multiplying this equation by  $c_j$  and summing over  $j$  gives

$$\sum_j c_j \left[ \frac{d^2x_j}{dt^2} + \gamma \frac{dx_j}{dt} + \omega_0^2 x_j \right] = \sum_j c_j [\cos(\omega_j t)] \quad (7)$$

Therefore, if we define

$$x(t) = \sum_j c_j x_j(t) \quad (8)$$

we immediately get

$$\frac{d^2x}{dt^2} + \gamma \frac{dx}{dt} + \omega_0^2 x = F(t) \quad (9)$$

as desired.

In summary, assuming Eq. (4) holds (we will come back to this soon), we have reduced solving Eq. (1) to solving Eq. (6). That is, if we can solve the equation with a  $\cos(\omega_d t)$  driving force, we can solve the equation for *any* driving force.

## 2 Driven oscillator

Our first task is to solve

$$\frac{d^2x}{dt^2} + \gamma \frac{dx}{dt} + \omega_0^2 x = \frac{F_0}{m} \cos(\omega_d t) \quad (10)$$

Here we have made the normalization more physical by adding  $F_0$ , for the strength of force with units of force, and dividing by the oscillator mass  $m$  to get an acceleration (the left hand side has units of acceleration, as in  $\frac{d^2x}{dt^2}$ ). What's a good guess for a solution? Trying  $x(t) = \cos(\omega_d t)$  or  $x(t) = \sin(\omega_d t)$  will not work since there are first *and* second derivatives in the equation. We need exponentials.

The key to turning the problem from cosines into exponentials is to recall that

$$e^{-i\omega t} = \cos(\omega t) - i \sin(\omega t) \quad (11)$$

so that

$$\cos(\omega_d t) = \operatorname{Re}(e^{-i\omega_d t}) \quad (12)$$

Now suppose we find a solution to

$$\frac{d^2}{dt^2} z + \gamma \frac{d}{dt} z + \omega_0^2 z = \frac{F_0}{m} e^{-i\omega_d t} \quad (13)$$

with a complex function  $z(t)$ . Then we define

$$x(t) \equiv \operatorname{Re}[z(t)] \quad (14)$$

Taking the real part of Eq. (13) then gives

$$\operatorname{Re}\left[\frac{d^2}{dt^2} z + \gamma \frac{d}{dt} z + \omega_0^2 z\right] = \operatorname{Re}\left[\frac{F_0}{m} e^{-i\omega_d t}\right] \quad (15)$$

which is exactly Eq. (10). So we have reduced the problem to using an exponential driving force instead of a cosine driving force.

Plugging in a guess  $z(t) = C e^{-i\omega_d t}$  into Eq. (15) gives

$$C e^{-i\omega_d t} [-\omega_d^2 - i\gamma\omega_d + \omega_0^2] = \frac{F_0}{m} e^{-i\omega_d t} \quad (16)$$

Now the  $e^{i\omega_d t}$  factors drop out and we have a simple algebraic relation

$$C = \frac{F_0}{m} \frac{1}{\omega_0^2 - i\gamma\omega_d - \omega_d^2} \quad (17)$$

Thus

$$z(t) = \frac{F_0}{m} \frac{1}{\omega_0^2 - i\gamma\omega_d - \omega_d^2} e^{-i\omega_d t} \quad (18)$$

To get the solution to the original equation with a real function  $x(t)$  we use Eq. (14):

$$x(t) = \operatorname{Re}\left[\frac{F_0}{m} \frac{1}{\omega_0^2 - i\gamma\omega_d - \omega_d^2} e^{-i\omega_d t}\right] \quad (19)$$

Now we just have to simplify this using algebra.

First, we get the  $i$ 's to the numerator by writing

$$\frac{1}{\omega_0^2 - i\gamma\omega_d - \omega_d^2} = \frac{\omega_0^2 - \omega_d^2 + i\gamma\omega_d}{\omega_0^2 - \omega_d^2 + i\gamma\omega_d} \frac{1}{\omega_0^2 - \omega_d^2 - i\gamma\omega_d} = \frac{\omega_0^2 - \omega_d^2 + i\gamma\omega_d}{(\omega_0^2 - \omega_d^2)^2 + (\gamma\omega_d)^2} \quad (20)$$

$$\equiv A + Bi \quad (21)$$

where

$$A = \frac{\omega_0^2 - \omega_d^2}{(\omega_0^2 - \omega_d^2)^2 + (\gamma\omega_d)^2} \quad B = \frac{\gamma\omega_d}{(\omega_0^2 - \omega_d^2)^2 + (\gamma\omega_d)^2} \quad (22)$$

Then

$$x(t) = \operatorname{Re}\left[\frac{F_0}{m} (A + Bi) e^{-i\omega_d t}\right] = \frac{F_0}{m} \operatorname{Re}[(A + Bi)(\cos(\omega_d t) - i \sin(\omega_d t))] \quad (23)$$

$$= \frac{F_0}{m} (A \cos \omega_d t + B \sin \omega_d t) \quad (24)$$

In summary, we found an exact solution to Eq. (10):

$$x(t) = \frac{F_0}{m} \left\{ \frac{\omega_0^2 - \omega_d^2}{(\omega_0^2 - \omega_d^2)^2 + (\gamma\omega_d)^2} \cos\omega_d t + \frac{\gamma\omega_d}{(\omega_0^2 - \omega_d^2)^2 + (\gamma\omega_d)^2} \sin\omega_d t \right\} \quad (25)$$

## 2.1 Transients

We found a single exact solution. What happened to the boundary conditions? The dependence on boundary conditions is entirely determined by solutions to the **homogeneous** equation, with  $F = 0$ :

$$\frac{d^2 x_0}{dt^2} + \gamma \frac{dx_0}{dt} + \omega_0^2 x_0 = 0 \quad (26)$$

Solutions to this equation are called **homogeneous solutions**. The solution  $x(t)$  in Eq. (25) is called the **inhomogeneous solution**. Note that  $x_0(t) + x(t)$  will also satisfy the inhomogeneous Eq. (10), due to linearity. Thus we can always add a homogeneous solution to an inhomogeneous solution. We saw before that the homogeneous solutions all have  $e^{-\frac{\gamma}{2}t}$  factors, plus possibly some oscillatory component. Thus they die off at late time. For this reason, they are called **transient**. Transients are determined by boundary conditions. If you have a driving force for long enough time, then the transient is irrelevant.

## 2.2 Phase lag

A good way to see the physics hidden in the solution  $x(t)$  is to take limits. First, consider the limit with no damping,  $\gamma = 0$ . Then,

$$x(t) = \frac{F_0}{m} \frac{1}{\omega_0^2 - \omega_d^2} \cos\omega_d t \quad (27)$$

We can compare this to our driving force  $F(t) = F_0 \cos\omega_d t$ . For  $\omega_d < \omega_0$  the sign of the position and the force are the same so they are exactly in phase. Now say we crank up the driving frequency  $\omega_d$  until it reaches then surpasses  $\omega_0$ . For  $\omega_d > \omega_0$ , the sign of the solution flips and the oscillator is out of phase with the driver. Physically, the oscillator can't keep up with the driving force: it experiences **phase lag**.

## 2.3 Power and energy

We see from Eq. (25) there is a part of  $x(t)$  which is exactly proportional to the driving force  $F(t) = F_0 \cos\omega_d t$  and a part which is out of phase. We call the in-phase part the **elastic amplitude**. It is proportional to

$$A = \frac{\omega_0^2 - \omega_d^2}{(\omega_0^2 - \omega_d^2)^2 + (\gamma\omega_d)^2} \quad (28)$$

The out-of-phase part is the **absorptive amplitude**. Its magnitude is

$$B = \frac{\gamma\omega_d}{(\omega_0^2 - \omega_d^2)^2 + (\gamma\omega_d)^2} \quad (29)$$

Thus for  $\gamma = 0$ , no damping, there is no absorptive part. Since the absorptive part is proportional to  $\gamma$  it should have to do with energy being lost from the oscillator into the system. To see how this works, we need to compute the energy and the power.

Recall that work is force times displacement  $W = F\Delta x$  and power is work per unit time:

$$P = \frac{W}{\Delta t} = F \frac{\Delta x}{\Delta t} \quad (30)$$

For small displacements and small times, this becomes

$$P = F \frac{dx}{dt} \quad (31)$$

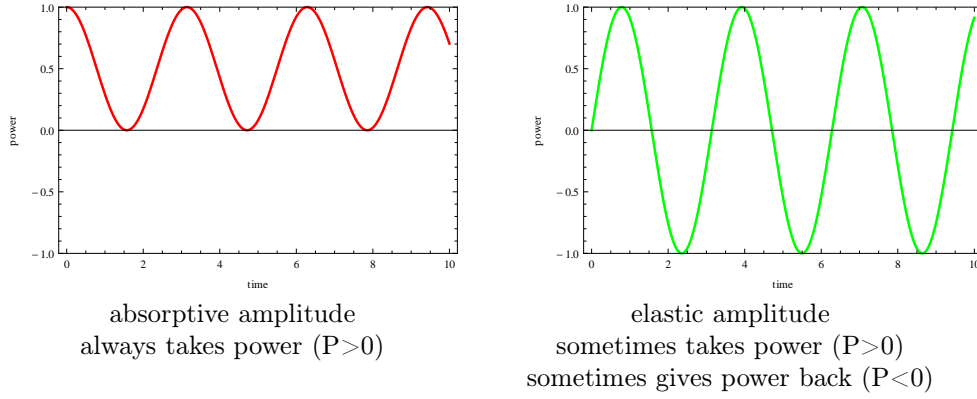
Plugging in our solution  $x(t) = \frac{F_0}{m}[A \cos(\omega_d t) + B \sin(\omega_d t)]$

$$P = F_0 \cos(\omega_d t) \left[ -\omega_d \frac{F_0}{m} A \sin(\omega_d t) + \omega_d \frac{F_0}{m} B \cos(\omega_d t) \right] \quad (32)$$

$$= -\frac{F_0^2}{2m} \omega_d A \sin(2\omega_d t) + \frac{F_0^2}{m} B \omega_d \cos^2(\omega_d t) \quad (33)$$

where  $2\sin\theta\cos\theta = \sin(2\theta)$  has been used. This is the **power put into the system by the driving force**.

We see that the absorptive part is proportional to  $\cos^2 \omega_d t$  which is positive for all times. Thus it always takes (absorbs) power. On the other hand, the elastic amplitude is proportional to  $\sin(2\omega_d t)$  which is sometimes positive and sometimes negative. These are shown here



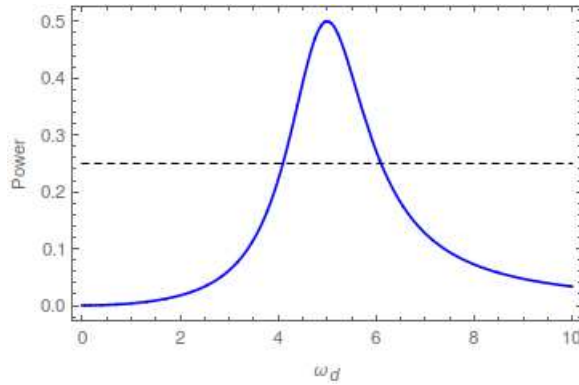
**Figure 1.** Absorptive and elastic amplitudes

When the power is negative, as in the elastic amplitude, the oscillator is returning power to the driver. The elastic amplitude averages to zero. Since  $\gamma = 0$  implies that the absorptive amplitude vanishes so the entire solution is elastic, we draw the logical conclusion that with no damping ( $\gamma = 0$ ) no net power is needed to drive the system (a little power is needed to get it started, but once it's moving, the driver no longer does work).

The average power put into the system is over a period  $T = \frac{2\pi}{\omega_d}$  is

$$\langle P \rangle = \frac{1}{T} \int_0^T dt P(t) = \frac{F_0^2}{2m} B \omega_d = \left( \frac{F_0^2}{2\gamma m} \right) \frac{(\gamma \omega_d)^2}{(\omega_0^2 - \omega_d^2)^2 + (\gamma \omega_d)^2} \quad (34)$$

Here is a plot of this average power as a function of  $\omega_d$  for fixed  $\gamma$  and  $\omega_0$ .



**Figure 2.** Power absorbed for  $\gamma = 2$  and  $\omega_0 = 5$  as a function of the driving frequency  $\omega_d$ . The maximum is when  $\omega_d = \omega_0$  known as **resonance**. The dashed line is half of the maximum power. The length of the dashed line between the points where it hits the curve (width at half-maximum) is  $\gamma$ .

This power absorption curve has a maximum at  $\omega_d = \omega_0$  (you can check this) where  $\langle P \rangle = \frac{F_0^2}{2\gamma m}$ . This is known as a **resonance**. One way to find the resonance frequency  $\omega_0$  of a system is by varying the driving force until maximum power is absorbed. The power is half the resonant power,  $\langle P \rangle = \frac{F_0^2}{4\gamma m}$ , when

$$\omega_d = \frac{1}{2} \sqrt{4\omega_0^2 + \gamma^2} \pm \frac{1}{2} \gamma \quad (35)$$

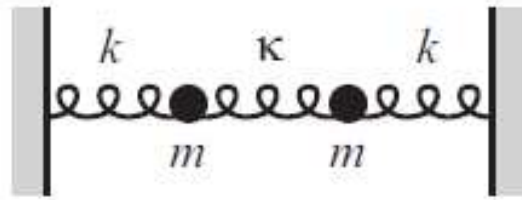
The difference between these two driving frequencies is  $\gamma$ . Thus, one can also read  $\gamma$  off of the plot in Fig. 2: it is the value of the width at half-maximum. This kind of curve is called a **Lorentzian**. Its maximum is at  $\omega_0$  and its width is  $\gamma$ .

## Lecture 3: Coupled oscillators

### 1 Two masses

To get to waves from oscillators, we have to start coupling them together. In the limit of a large number of coupled oscillators, we will find solutions while look like waves. Certain features of waves, such as resonance and normal modes, can be understood with a finite number of oscillators. Thus we start with two oscillators.

Consider two masses attached with springs



(1)

Let's say the masses are identical, but the spring constants are different.

Let  $x_1$  be the displacement of the first mass from its equilibrium and  $x_2$  be the displacement of the second mass from its equilibrium. To work out Newton's laws, we first want to know the force on  $x_1$  when it is moved from its equilibrium while holding  $x_2$  fixed. This is

$$F_{\text{on 1 from moving 1}} = F = -kx_1 - \kappa x_1 \quad (2)$$

The signs are both chosen so that they oppose the motion of the mass. There is also a force on  $x_1$  if we move  $x_2$  holding  $x_1$  fixed. This force is

$$F_{\text{on 1 from moving 2}} = \kappa x_2 \quad (3)$$

To check the sign, note that if  $x_2$  is increased, it pulls  $x_1$  to the right. There is no contribution to this force from the spring between the second mass and the wall, since we are moving the mass by hand and just asking how it affects the first mass. Thus

$$m \ddot{x}_1 = -(k + \kappa)x_1 + \kappa x_2 \quad (4)$$

similarly,

$$m \ddot{x}_2 = -(k + \kappa)x_2 + \kappa x_1 \quad (5)$$

One way to solve these equations is to note that if we add them, we get

$$m(\ddot{x}_1 + \ddot{x}_2) = -k(x_1 + x_2) \quad (6)$$

This is just  $m \ddot{y} = -ky$  for  $y = x_1 + x_2$ , so the solutions are sines and cosines, or cosine and a phase:

$$x_1 + x_2 = A_s \cos(\omega_s t + \phi_s), \quad \omega_s = \sqrt{\frac{k}{m}} \quad (7)$$

Another way solve them is taking the difference

$$m(\ddot{x}_1 - \ddot{x}_2) = (-k - 2\kappa)(x_1 - x_2) \Rightarrow x_1 - x_2 = A_f \cos(\omega_f t + \phi_f), \quad \omega_f = \sqrt{\frac{k + 2\kappa}{m}} \quad (8)$$

We write  $\omega_s$  for  $\omega_{\text{slow}}$  and  $\omega_f$  for  $\omega_{\text{fast}}$ , since  $\omega_f > \omega_s$ . Thus we have found two solutions each of which oscillate with fixed frequency. These are the **normal modes** for this system. A general solution is a linear combination of these two solutions. Explicitly, we have:

$$x_1 = \frac{1}{2}[(x_1 + x_2) + (x_1 - x_2)] = \frac{1}{2}[A_s \cos(\omega_s t + \phi_s) + A_f \cos(\omega_f t + \phi_f)] \quad (9)$$

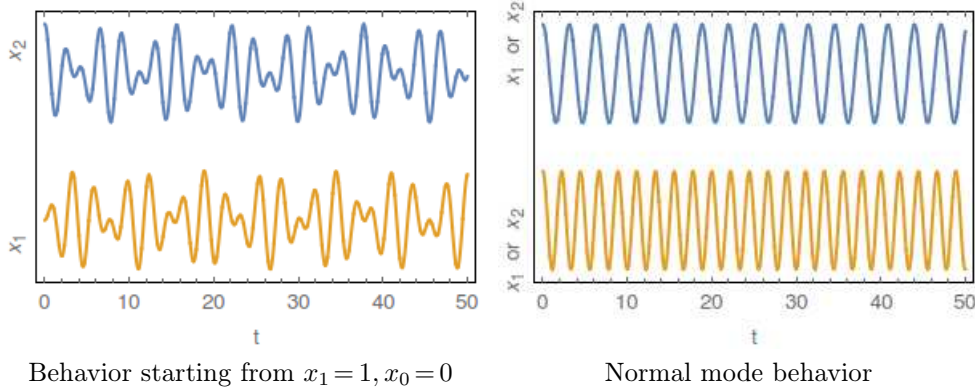
$$x_2 = \frac{1}{2}[(x_1 + x_2) - (x_1 - x_2)] = \frac{1}{2}[A_s \cos(\omega_s t + \phi_s) - A_f \cos(\omega_f t + \phi_f)] \quad (10)$$

If we can excite the masses so that  $A_f = 0$  then the masses will both oscillate at the frequency  $\omega_s$ . In practice, we can do this by pulling the masses to the right by the same amount, so that  $x_1(0) = x_2(0)$  which implies  $A_f = 0$ . The solution is then  $x_1 = x_2$  and both oscillate at the frequency  $A_s$  for all time. This is the **symmetric oscillation mode**. Since  $x_1 = x_2$  at all times, both masses move right together, then move left together.

If we excite the masses in such a way that  $A_s = 0$  then  $x_1 = -x_2$  and both oscillate at frequency  $\omega_f$ . We can set this up by pulling the masses in opposite directions. In this mode, when one mass is right of equilibrium, the other is left, and vice versa. So this is an **antisymmetric mode**.

## 2 Beats

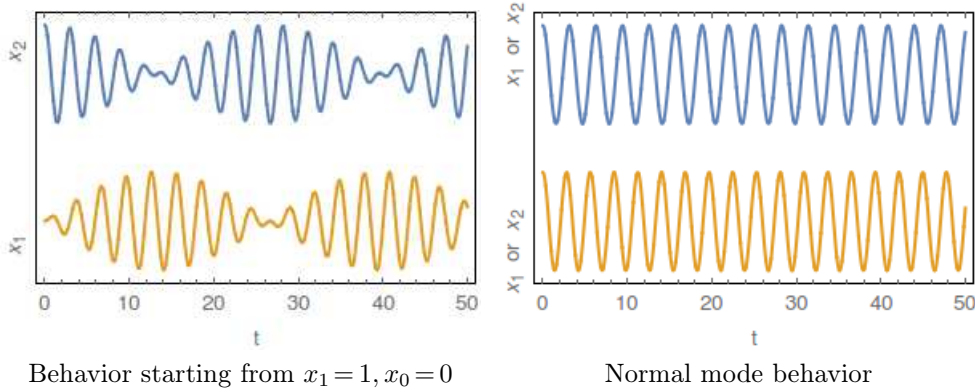
You should try playing with the coupled oscillator solutions in the Mathematica notebook `oscillators.nb`. Try varying  $\kappa$  and  $k$  to see how the solution changes. For example, say  $m = 1$ ,  $\kappa = 2$  and  $k = 4$ . Then  $\omega_s = 2$  and  $\omega_f = 2\sqrt{2}$ . Here are the solutions:



**Figure 1.** Left shows the motion of masses  $m = 1, \kappa = 2$  and  $k = 4$  starting with  $x_1 = 1$  and  $x_2 = 0$ . Right shows the normal modes, with  $x_1 = x_2 = 1$  (top) and  $x_1 = 1, x_2 = -1$  (bottom).

If you look closely at the left plot, you can make out two distinct frequencies: the normal mode frequencies, as shown on the right.

Now take  $\kappa = 0.5$  and  $k = 4$ . Then  $\omega_s = 2$  and  $\omega_f = 2.2$ . In this case



**Figure 2.** Motion of masses and normal modes for  $k = 0.5$  and  $\kappa = 4$

Now we can definitely see two distinct frequencies in the positions of the two masses. Are these the two frequencies  $\omega_s$  and  $\omega_f$ ? Comparing to the normal mode plots, it is clear they are not. One is much slower. However, we do note that  $\omega_s \approx \omega_f$ . What we are seeing here is the emergence of **beats**. Beats occur when two normal mode frequencies get close.

Beats can be understood from the simple trigonometric relation

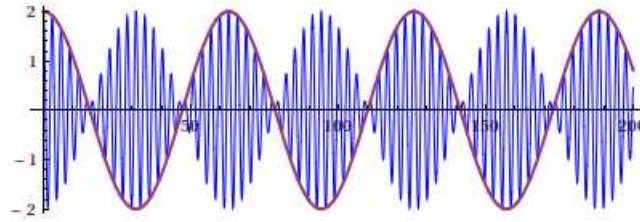
$$\cos(\omega_1 t) + \cos(\omega_2 t) = 2\cos\left(\frac{\omega_1 + \omega_2}{2}t\right)\cos\left(\frac{\omega_1 - \omega_2}{2}t\right) \quad (11)$$

When you excite two frequencies  $\omega_1$  and  $\omega_2$  at the same time, the solution to the equations of motion is the sum of the separate oscillating solutions (by linearity!). Eq. (11) shows that this sum can also be written as the *product* of two cosines. In particular, if  $\omega_1 \approx \omega_2$  then

$$\omega = \frac{\omega_1 + \omega_2}{2} \approx \omega_1 \approx \omega_2 \quad \varepsilon = \frac{\omega_1 - \omega_2}{2} \ll \omega_1, \omega_2 \quad (12)$$

So the sum looks like an oscillation whose frequency  $\omega$  is the *average* of the two normal mode frequencies modulated by an oscillation with frequency  $\varepsilon$  given by half the difference in the frequencies.

Beats are important because they can generate frequencies well below the normal mode frequencies. For example, suppose you have two strings which are not quite in tune. Say they are supposed to both be the note A4 at 440 Hz, but one is actually  $\nu_1 = 442\text{Hz}$  and the other is  $\nu_2 = 339\text{Hz}$ . If you pluck both strings together you will hear the average frequency  $\Omega = 440.5\text{Hz}$ , but also there will be an oscillation at  $\varepsilon = \frac{1}{2}(442 - 339)\text{Hz} = 1.5\text{Hz}$ . This oscillation is the enveloping curve over the high frequency (440.5 Hz) oscillations



**Figure 3.** The red curve is  $\cos\left(2\pi\frac{\nu_1 - \nu_2}{2}t\right)$ . When hearing beats, the observed frequency is the frequency of the extrema  $\nu_{\text{beat}} = \nu_1 - \nu_2$  which is twice the frequency of this curve.

As you can see from the figure, due to the high frequency oscillations, there are peaks in the amplitude twice as often as peaks in  $\cos\left(2\pi\frac{\nu_1 - \nu_2}{2}t\right)$ . Thus what we hear are beats at the **beat frequency**

$$\nu_{\text{beat}} = |\nu_1 - \nu_2| \quad (13)$$

We use an absolute value since we want a frequency to be positive (it's the same frequency whether  $\nu_1 > \nu_2$  or  $\nu_2 > \nu_1$ ). Note that there is no factor of 2 in the conventional definition of  $\nu_{\text{beat}}$ , since we only ever hear the modulus of the oscillation not the phase.

Thus with  $\nu_f = 442\text{Hz}$  and  $\nu_s = 339\text{Hz}$  the beat frequency is  $\nu_{\text{beat}} = 3\text{Hz}$ . Thus you hear something happening 3 times a second. This is a regular beating in off-tune notes which is audible by ear. In fact, it is a useful trick for tuning – change one string until the beating disappears. Then the strings are in tune. We will see numerous examples of beats as the course progresses.

### 3 Two masses with matrices

We solved the two coupled mass problem by looking at the equations and noting that their sum and difference would be independent solutions. For more complicated systems (more masses, different couplings) we should not expect to be able to guess the answer in this way. Can you guess the solution if the two oscillators have different masses?



To develop a more systematic procedure, suppose we have lots of masses with lots of different springs connected in a complicated way. Then the equations of motion are

$$m_1 \ddot{x}_1 = k_{11} x_1 + k_{12} x_2 + \cdots + k_{1n} x_n \quad (14)$$

$$\dots \quad (15)$$

$$m_n \ddot{x}_n = k_{n1} x_1 + k_{n2} x_2 + \cdots + k_{nn} x_n \quad (16)$$

where  $k_{ij}$  are constants, representing the strength of the spring between masses  $i$  and  $j$ . Note that all of these equations are linear. What are the solutions in this general case? This is an algebra problem involving linear equations. Hence we should be able to solve it with **linear algebra**.

To connect to linear algebra, let's return to our two mass system. Since the equations of motion are linear, we expect them to be solved by exponentials  $x_1 = c_1 e^{i\omega t}$  and  $x_2 = c_2 e^{i\omega t}$  for some  $\omega$ ,  $c_1$  and  $c_2$ . As with the driven oscillator from the last lecture, we are using complex solutions to make the math simpler, then we can always take the real part at the end. Plugging in these guesses, Eqs. (4) and (5) become

$$-m_1 \omega^2 c_1 = -(k + \kappa) c_1 + \kappa c_2 \quad (17)$$

$$-m_2 \omega^2 c_2 = -(k + \kappa) c_2 + \kappa c_1 \quad (18)$$

We have let the masses be different for generality.

Next, we will write these equations in matrix form. To do so, we define a vector  $\vec{c}$  as

$$\vec{c} = \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} \quad (19)$$

Then the equations of motion become

$$M \cdot \vec{c} = \begin{pmatrix} \frac{-k - \kappa}{m_1} & \frac{\kappa}{m_1} \\ \frac{\kappa}{m_2} & \frac{-k - \kappa}{m_2} \end{pmatrix} \cdot \vec{c} = -\omega^2 \vec{c} \quad (20)$$

where  $M$  is defined by this equation.

You might recognize this as an eigenvalue equation. An  $n \times n$  matrix  $A$  has  $n$  **eigenvalues**  $\lambda_i$  and  $n$  associated **eigenvectors**  $\vec{v}_i$  which satisfy

$$A \cdot \vec{v}_i = \lambda_i \vec{v}_i \quad (21)$$

The eigenvalues don't all have to be different. Note that the left hand side is a matrix multiplying a vector while the right-hand side is just a number multiplying a vector. So studying eigenvalues and eigenvectors lets us turn matrices into numbers! Eigenvalues and eigenvectors are *the* fundamental mathematical concept of quantum mechanics. I cannot emphasize enough how important it is to master them.

Let's recall how to solve an eigenvalue equation. The trick is to write it first as

$$(A - \lambda \mathbb{1}) \vec{v} = 0 \quad (22)$$

where  $\mathbb{1}$  is the  $n \times n$  identity matrix. For  $n = 2$ ,  $\mathbb{1} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$ . For most values of  $\lambda$ , the matrix  $(A - \lambda \mathbb{1})$  has an inverse. Multiplying both sides of Eq. (22) by that inverse, we find  $\vec{v} = 0$ . This is the trivial solution (it obviously satisfies Eq. (21) for any  $A$ ). The nontrivial solutions consequently must correspond to values of  $\lambda$  for which  $(A - \lambda \mathbb{1})$  does **not** have an inverse. When does a matrix not have an inverse? A result from linear algebra is that a matrix is not invertible if and only if its determinant is zero. Thus the equation  $\det(A - \lambda \mathbb{1}) = 0$  is an algebraic equation for  $\lambda$  whose solutions are the eigenvalues  $\lambda_i$ .

It is useful to know that determinant of a  $2 \times 2$  matrix is

$$\det \begin{pmatrix} a & b \\ c & d \end{pmatrix} = ad - bc \quad (23)$$

You should have this memorized. For a  $3 \times 3$  matrix, the determinant is:

$$\det \begin{pmatrix} a & b & c \\ d & e & f \\ g & h & i \end{pmatrix} = a(ei - fh) - b(di - fg) + c(dh - eg) \quad (24)$$

You should know how to compute this, but don't need to memorize the formula. Beyond  $3 \times 3$ , you probably want to take determinants with Mathematica rather than by hand.

So, returning to Eq. (20), the eigenvalues  $-\omega^2$  must satisfy

$$0 = \det(M + \omega^2 \mathbb{1}) = \det \begin{pmatrix} \frac{-k - \kappa}{m_1} + \omega^2 & \frac{\kappa}{m_1} \\ \frac{\kappa}{m_2} & \frac{-k - \kappa}{m_2} + \omega^2 \end{pmatrix} \quad (25)$$

$$= \left( \frac{-k - \kappa}{m_1} + \omega^2 \right) \left( \frac{-k - \kappa}{m_2} + \omega^2 \right) - \frac{\kappa^2}{m_1 m_2} \quad (26)$$

This is a quadratic equation for  $\omega^2$ , with two roots: the two eigenvalues.

Let's set  $m_1 = m_2 = m$  now to check that we reproduce our old result. Multiplying Eq. (26) by  $m^2$ , it reduces to

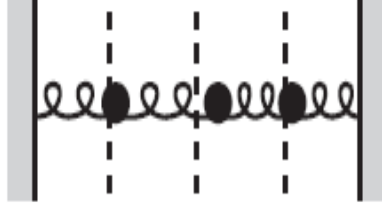
$$(k + \kappa - m\omega^2)^2 = \kappa^2 \quad (27)$$

Thus  $k + \kappa - m\omega^2 = \pm \kappa$ . Or in other words

$$\omega = \omega_s = \sqrt{\frac{k}{m}}, \quad \omega = \omega_f = \sqrt{\frac{k + 2\kappa}{m}} \quad (28)$$

These are the two normal mode frequencies we found above. Note that we didn't have to take the real part of the solution to find the normal mode frequencies. We only need to take the real part to find the solutions  $x(t)$ .

Now let's try three masses. We can couple them all together and to the walls in any which way



(29)

The equations of motion for this system will be of the form

$$m_1 \ddot{x}_1 = k_{11} x_1 + k_{12} x_2 + k_{13} x_3 \quad (30)$$

$$m_2 \ddot{x}_2 = k_{21} x_1 + k_{22} x_2 + k_{23} x_3 \quad (31)$$

$$m_3 \ddot{x}_3 = k_{31} x_1 + k_{32} x_2 + k_{33} x_3 \quad (32)$$

Some of these  $k_{ij}$  are probably zero, but we don't care. Writing  $x_1 = c_1 e^{i\omega t}$ ,  $x_2 = c_2 e^{i\omega t}$  and  $x_3 = c_3 e^{i\omega t}$ , these equations become algebraic:

$$-\omega^2 c_1 = \frac{k_{11}}{m_1} c_1 + \frac{k_{12}}{m_1} c_2 + \frac{k_{13}}{m_1} c_3 \quad (33)$$

$$-\omega^2 c_2 = \frac{k_{21}}{m_2} c_1 + \frac{k_{22}}{m_2} c_2 + \frac{k_{23}}{m_2} c_3 \quad (34)$$

$$-\omega^2 c_3 = \frac{k_{31}}{m_3} c_1 + \frac{k_{32}}{m_3} c_2 + \frac{k_{33}}{m_3} c_3 \quad (35)$$

In other words,

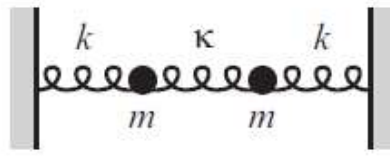
$$(M + \omega^2 \mathbb{1}) \vec{x} = 0 \quad (36)$$

with  $M$  the matrix whose entries are  $M_{ij} = \frac{k_{ij}}{m_i}$ . So to find the normal mode frequencies  $\omega$ , we need to solve  $\det(M + \omega^2 \mathbb{1}) = 0$ . For a  $3 \times 3$  matrix, there will be 3 eigenvalues and hence three normal-mode frequencies.

# Lecture 4: Oscillators to Waves

## 1 Review two masses

Last time we studied how two coupled masses on springs move



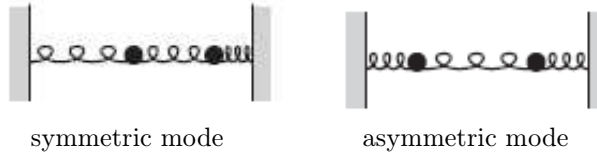
If we take  $\kappa = k$  for simplicity, the two normal modes correspond to

$$\begin{pmatrix} x_1(t) \\ x_2(t) \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \text{Re}[e^{i\omega_s t}], \quad \omega_s = \sqrt{\frac{k}{m}} \quad (1)$$

and

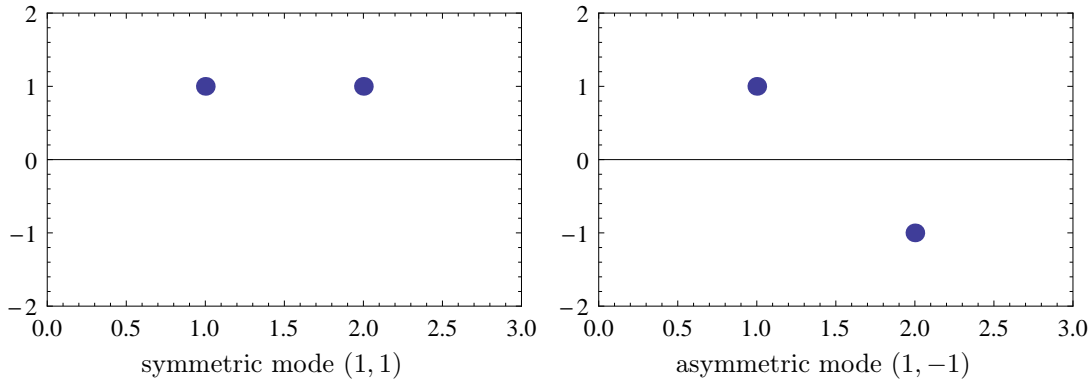
$$\begin{pmatrix} x_1(t) \\ x_2(t) \end{pmatrix} = \begin{pmatrix} 1 \\ -1 \end{pmatrix} \text{Re}[e^{i\omega_f t}], \quad \omega_f = \sqrt{\frac{3k}{m}} \quad (2)$$

One way to make sure we only excite these modes is to displace the masses so that their initial conditions are either  $x_1 = x_2$ , for the symmetric (slow) mode or  $x_1 = -x_2$  if we want to excite only the antisymmetric (fast mode). In other words, the normal mode solutions are in 1-to-1 correspondence with initial conditions. We can draw the initial conditions as



**Figure 1.** Initial conditions to excite normal modes

These pictures are going to be a little hard to look at if there are multiple masses. Thus it is helpful to draw the initial conditions as points in the  $y$  direction. That is, we write



**Figure 2.** Initial conditions to excite normal modes. The x-axis in these plots is the index of the mass (only two masses,  $m_1$  and  $m_2$  in this case). The y-axis is the displacement from equilibrium that you can release the masses from to get them oscillating with a normal mode frequency. In other words, these are plots of the eigenvectors.

In these figures, the displacement is still **longitudinal** (in the direction of the spring), we are just drawing it in the  $y$  axis because it is easier to see.

## 2 Three masses

Now consider 3 identical masses with all identical spring constants



The equations of motion for the first mass are

$$m \frac{d^2 x_1}{dt^2} = -2kx_1 + kx_2 \quad (4)$$

As before, you should think of first term on the right as the force generated on mass 1 when it is moved by a distance  $x_1$ . The sign on  $-2kx_1$  is negative since it always wants to go back to equilibrium. The second term  $+kx_2$  is the force that is exerted on mass 1 when mass 2 is moved *holding everything else fixed*. It has a  $+$  sign, since if I move mass 2, then mass 1 wants to leave its equilibrium position. If we move mass 3 holding everything else fixed, no force is exerted on mass 1.

Similarly,

$$m \frac{d^2 x_2}{dt^2} = -2kx_2 + kx_1 + kx_3 \quad (5)$$

and

$$m \frac{d^2 x_3}{dt^2} = -2kx_3 + kx_2 \quad (6)$$

Writing

$$\vec{x} = \begin{pmatrix} x_1^0 \\ x_2^0 \\ x_3^0 \end{pmatrix} e^{i\omega t} \quad (7)$$

The equations of motion become

$$-\omega^2 \vec{x} = \omega_0^2 \begin{pmatrix} -2 & 1 & 0 \\ 1 & -2 & 1 \\ 0 & 1 & -2 \end{pmatrix} \vec{x} \quad (8)$$

where

$$\omega_0 = \sqrt{\frac{k}{m}} \quad (9)$$

is the frequency associated with a single mass.

The normal mode frequencies are the eigenvalues of this matrix. Plugging in to Mathematica we find eigenvalues and associated eigenvectors

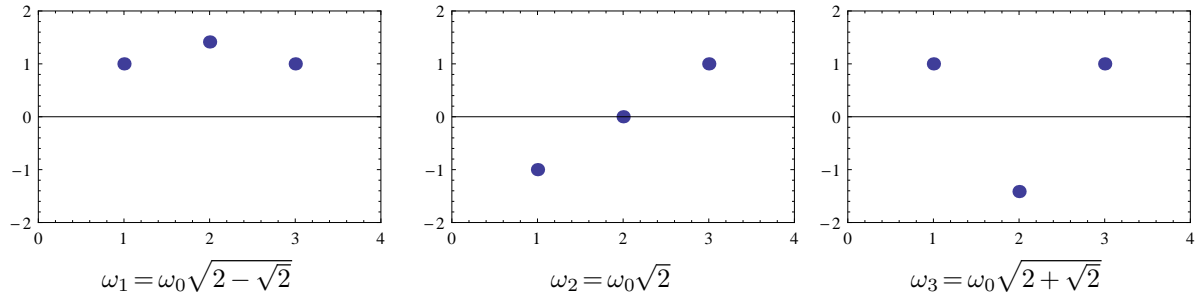
$$\omega_1 = \omega_0 \sqrt{2 - \sqrt{2}}, \quad \vec{x}(0) = \begin{pmatrix} 1 \\ \sqrt{2} \\ 1 \end{pmatrix} \quad (10)$$

$$\omega_2 = \omega_0 \sqrt{2}, \quad \vec{x}(0) = \begin{pmatrix} -1 \\ 0 \\ 1 \end{pmatrix} \quad (11)$$

$$\omega_3 = \omega_0 \sqrt{2 + \sqrt{2}}, \quad \vec{x}(0) = \begin{pmatrix} 1 \\ -\sqrt{2} \\ 1 \end{pmatrix} \quad (12)$$

We have ordered these from slowest to fastest.

To think about these solutions, it is helpful to plot the eigenvectors (the initial displacements)



**Figure 3.** Initial conditions to excite normal modes with 3 masses

From these pictures, it is not hard to see why the frequencies are higher for the third solution - the masses are more stretched apart, so there is more force between them, causing faster oscillation.

### 3 Completeness of eigenvectors

Before going on to  $N$  modes, it's worth making one very important point. Recall from linear algebra that the set of eigenvectors of any matrix is **complete**, meaning that any vector can be written as a linear combination of eigenvectors. For oscillators, this means that any solution to the equations of motion can be written as

$$\vec{x}(t) = \sum_j c_j \vec{x}_j(t) \quad (13)$$

where the sum is over normal modes  $n$ . The normal modes are solutions  $\vec{x}_j(t)$  to the equation of motion with frequencies  $\omega_j$ . That is

$$\vec{x}_j(t) = \vec{x}_j^0 \cos(\omega_j t) \quad (14)$$

where  $\vec{x}_j^0$  is a constant vector.

In summary, normal modes oscillate with a single frequency. A general solution can always be written as a sum of normal modes.

## 4 $N$ modes

Now we'll solve the  $N$  mass system. You should think of lots of springs put together as being simply one long spring where the masses are pieces of the spring itself. We'll see the wave equation result from this system. Solutions to the wave equation describe not just normal modes, but also waves, such as pulses sent down the spring (like pulses sent down a slinky). These pulses are called traveling waves, which are actually linear combinations of normal modes. Understanding waves will occupy the rest of the course, but first we have to solve the  $N$  spring system.

We'll construct the equations of motion for the  $N$  springs. Then we'll solve the coupled equations numerically for finite  $N$  to get a sense for what the answer should look like. Then we'll solve the system exactly for any  $N$ .

You may find this section quite abstract. It is not critical that you follow all the details here and be able to reproduce it all on your own. This is one of the most complicated derivations we will do in the course. Do your best. You should understand the result though, as summarized in Section 4.4.

### 4.1 Equations of motion

Ok, so, we want to string  $N$  masses together. Adding more masses to the right of mass 3 does not affect the equations of motions for masses 1 and 2. So their equations are

$$m \frac{d^2 x_1}{dt^2} = -2kx_1 + kx_2 \quad (15)$$

and

$$m \frac{d^2 x_2}{dt^2} = kx_1 - 2kx_2 + kx_3 \quad (16)$$

as before. Mass 3 is now like mass 2 – it has masses to the right and left of it. Thus,

$$m \frac{d^2 x_3}{dt^2} = kx_2 - 2kx_3 + kx_4 \quad (17)$$

In fact, it is easy to see that the generalization for any of the middle masses is

$$\boxed{m \frac{d^2 x_n}{dt^2} = kx_{n-1} - 2kx_n + kx_{n+1}} \quad (18)$$

The last mass has no mass on its right, so it gets an equation like mass 1:

$$m \frac{d^2 x_N}{dt^2} = kx_{N-1} - 2kx_N \quad (19)$$

Eq. (18), (15) and (19) are what we want to solve.

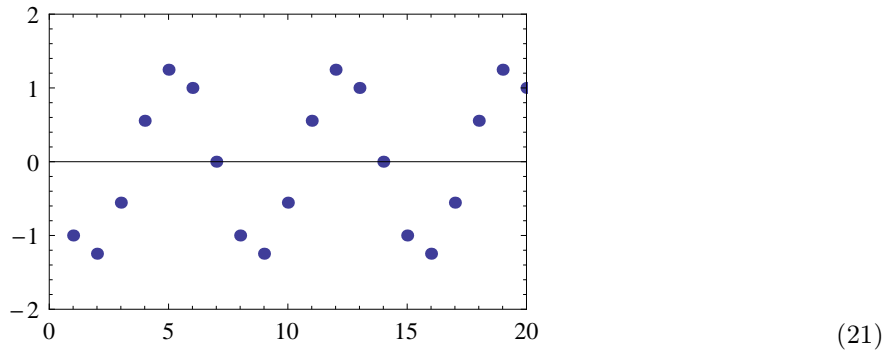
Putting all of these equations together with time dependence  $e^{i\omega t}$  for all masses leads to the matrix equation

$$-\omega^2 \vec{x} = \omega_0^2 \begin{pmatrix} -2 & 1 & & & \\ 1 & -2 & 1 & & \\ & 1 & -2 & \cdots & \\ & & \cdots & \cdots & \cdots \\ & & & \cdots & -2 & 1 \\ & & & & 1 & -2 \end{pmatrix} \vec{x} \quad (20)$$

with  $\omega_0 = \sqrt{\frac{k}{m}}$  as usual. All the entries not shown are zero.

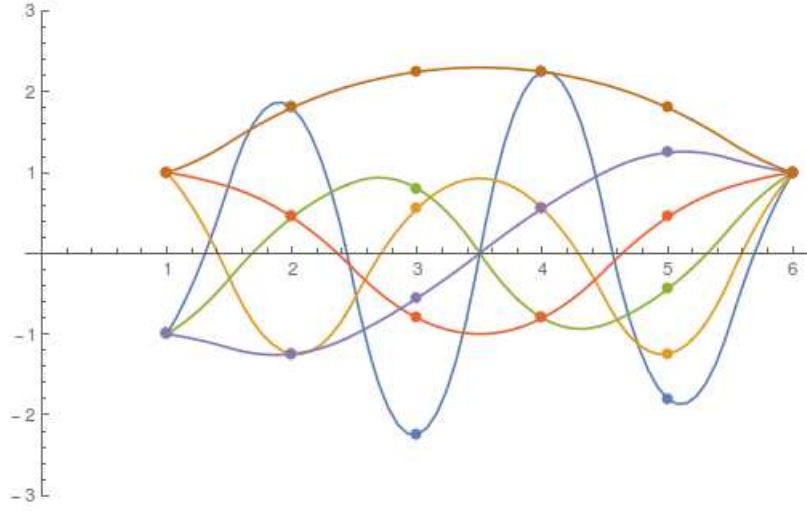
### 4.2 Numerical solutions

First, let's solve for the eigenvalues and eigenvectors of this system numerically using Mathematica. With 20 masses, the displacements associated with the 15th eigenvalue is



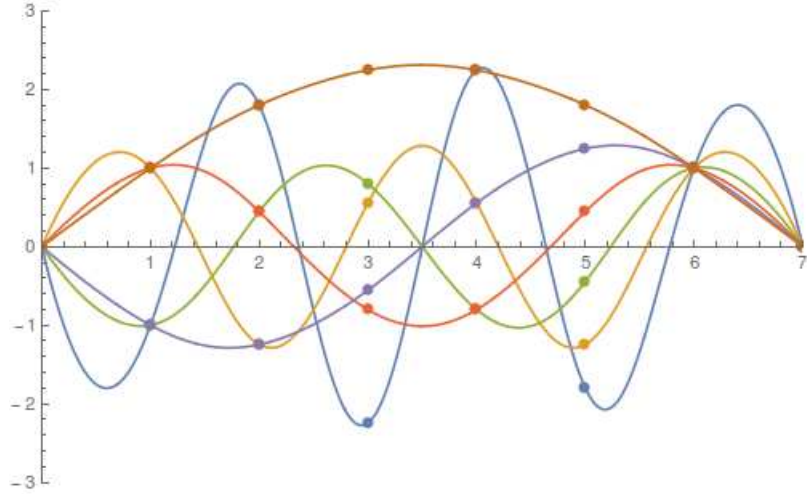
You may notice that it looks a lot like a cosine curve.

For the 6 mass system, we can plot the displacements for all the normal modes at the same time. Here they are, with the dots for clarity:



**Figure 4.** Normal modes displacements for the 6 mass system. These curves look like sine curves.

Already with six masses, in Figure 4, we see that the normal modes look like sine and cosine curves. They are not complete periods though – they stop abruptly. This is due to the equations of motion for masses 1 and  $N$  being determined by Eqs. (15) and (19) rather than the equation (18) that the rest of the masses satisfy. What is special about 1 and  $N$  is that they are attached to rigid walls, while all the other masses are attached to springs only. These rigid walls correspond to fixed boundary conditions at  $n=0$  and  $n=7$ :



**Figure 5.** Same as Figure 4, but with interpolations extended to  $n=0$  and  $n=7$  to show the boundary conditions.

To be clear, there is really no mass at  $n=0$ , but we are just pretending one is there. Now we can see that the solutions look like

$$x_n = B \sin(pn) e^{i\omega t} \quad (22)$$

for some  $p$ . The boundary conditions imply that  $p = \frac{\pi}{N+1}j$  for some  $j = 1, 2, 3, \dots$ . These  $p$  are called **wavenumbers**. In the continuum limit, we will see that wavenumber  $p = \frac{2\pi}{\lambda}$  with  $\lambda$  the wavelength. In the discrete case,  $p$  is dimensionless so it's harder to think of it as related to a wavelength. The fact that the wavenumbers are quantized by the boundary conditions is extremely important, both classically and in quantum mechanics. We will be revisiting this quantization at length throughout the course.



### 4.3 Exact solution

With the numerical solution giving a hint of where to look, let us just solve the system. We want to find vectors  $x_n$  that satisfy Eq. (18):

$$\frac{d^2 x_n}{dt^2} = \omega_0^2 [x_{n-1} - 2x_n + x_{n+1}] \quad (23)$$

Let's take a guess

$$x_n = B e^{ipn} e^{i\omega t} \quad (24)$$

That the time dependence is exponential follows from linearity – we always guess this. Here we are also guessing that since our numerical solution looks like sine functions the  $n$  dependence should be oscillatory.

Plugging our guess in we get

$$-\omega^2 e^{ipn} = [e^{ip(n-1)} - 2e^{ipn} + e^{ip(n+1)}] \omega_0^2 \quad (25)$$

Dividing both sides by  $-e^{ipn}$  gives

$$\omega^2 = [2 - e^{-ip} - e^{ip}] \omega_0^2 \quad (26)$$

$$= [2 - 2\cos(p)] \omega_0^2 \quad (27)$$

Thus, we have found solutions for any  $B$  as long as  $\omega$  and  $p$  are related by

$$\boxed{\omega(p) \equiv \pm \sqrt{2(1 - \cos(p))} \omega_0} \quad (28)$$

this is a type of **dispersion relation**. We will come back to dispersion relations later (once we have talked about dispersion).

To fix  $p$ , we need to use the boundary conditions, that is, the equations of motion for the end masses. The dispersion relationship doesn't tell us the sign of  $p$ . In fact, both  $p$  and  $-p$  lead to the same  $\omega$  so we can add solutions with  $p$  and  $-p$  and still have a solution. Considering that the numerical solutions looked like sine curves, let's guess

$$x_n = B \sin(pn) e^{i\omega t} \quad (29)$$

Now, mass 1 satisfies Eq. (15),  $m \frac{d^2 x_1}{dt^2} = -2kx_1 + kx_2$ . Plugging in our guess gives

$$-B\omega^2 \sin p = B\omega_0^2 [-2\sin p + \sin 2p] = -B\omega_0^2 \sin p [2 - 2\cos p] \quad (30)$$

where  $\sin(2x) = 2\sin(x)\cos(x)$  was used. Using Eq. (27) we then find that our guess works. You should make sure that you agree that Eq. (29) satisfies both Eq. (15) and Eq. (18) at this point.

Finally, we need to use the equation for mass  $N$ . It satisfies  $m \frac{d^2 x_N}{dt^2} = kx_{N-1} - 2kx_N$  as in Eq. (19). So

$$-B\omega^2 \sin Np = B\omega_0^2 [\sin(N-1)p - 2\sin Np] \quad (31)$$

Substituting Eq. (27) on the left-hand side, canceling the  $-2\sin Np$  terms, and then expanding  $\sin(Np-p)$  using  $\sin(\alpha - \beta) = \sin\alpha\cos\beta - \sin\beta\cos\alpha$  gives

$$2\sin Np \cos p = \sin(N-1)p = \sin Np \cos p - \sin p \cos Np \quad (32)$$

Thus

$$0 = \sin Np \cos p + \sin p \cos Np = \sin((N+1)p) \quad (33)$$

This equation is only satisfied if

$$\boxed{p = \frac{\pi}{N+1} j, \quad j = 1, 2, 3, \dots} \quad (34)$$

This are the same eigenvalues we found by guessing after Eq. (22). We have now derived rigorously that the equations of motion for masses 1 and  $N$  correspond to boundary conditions where we hold masses  $n=0$  and  $n=N+1$  fixed.

Thus the normal mode frequencies are

$$\boxed{\omega^2 = 2 \left( 1 - \cos \frac{\pi}{N+1} j \right) \omega_0^2, \quad j = 1, 2, 3, \dots, N} \quad (35)$$

and the solutions are  $x_n = B \sin(pn)e^{\pm i\omega t}$ . Let's check that this is the right answer for  $N = 3$ . We find

$$j = 1, \quad \omega^2 = 2\left(1 - \cos\frac{\pi}{4}\right)\omega_0^2 = (2 - \sqrt{2})\omega_0^2 \quad (36)$$

$$j = 2, \quad \omega^2 = 2\left(1 - \cos\frac{\pi}{2}\right)\omega_0^2 = 2\omega_0^2 \quad (37)$$

$$j = 3, \quad \omega^2 = 2\left(1 - \cos\frac{3\pi}{4}\right)\omega_0^2 = (2 + \sqrt{2})\omega_0^2 \quad (38)$$

These are exactly the frequencies we found in Eqs. (10) to (12).

For large  $N$ , the lowest frequencies have  $j \ll N$  thus using

$$\cos(x) = 1 - \frac{1}{2}x^2 + \dots, \quad x \ll 1 \quad (39)$$

we find

$$\omega^2 = 2\left(1 - \cos\frac{\pi}{N+1}j\right)\omega_0^2 = \left(\frac{\pi}{N+1}j\right)^2\omega_0^2 = p^2\omega_0^2 \quad (40)$$

That is

- For a large number of modes,  $\omega = \omega_0 p$ : the frequency is proportional to the wavenumber

In other words, the dispersion relation becomes linear:

$$\omega(p) = p\omega_0 \quad (41)$$

This linearity will be important when we discuss dispersion.

#### 4.4 Summary

In summary, we found the following solution for a large  $N$  number of masses connected by springs. For each integer  $j = 1, 2, 3, \dots$  there is a single normal mode solution. The position of mass  $n$  during the oscillation of normal mode  $j$  is given by

$$x_n^{(j)}(t) = \sin\left(\frac{\pi j}{N+1}n\right)\cos(\omega_j t + \phi_j) \quad (42)$$

We have chosen to write the solutions in manifestly real form (meaning we use sines and cosines rather than exponentials). The phase  $\phi_j$  is arbitrary. The frequencies are given by

$$\omega_j = \omega_0 \sqrt{2\left(1 - \cos\frac{\pi}{N+1}j\right)}, \quad (43)$$

The normal mode solutions are periodic with frequencies  $\omega_j$ .

For small  $j$

$$\omega_j \approx \frac{\pi j}{N+1}\omega_0 \quad (44)$$

An arbitrary solution can be written as a sum over normal modes as

$$(\vec{x})_n(t) = \sum_j a_j \sin\left(\frac{\pi j}{N+1}n\right)\cos(\omega_j t + \phi_j) \quad (45)$$

for some real constants  $a_j$  and  $\phi_j$

These solutions all satisfy the boundary conditions  $(\vec{x}_j)_0(t) = (\vec{x}_j)_{N+1}(t) = 0$ . It is straightforward to work out the general solution for other boundary conditions.

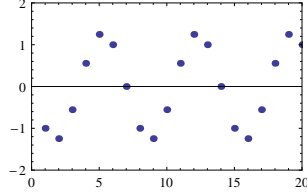
## 5 Continuum limit

We will now take the limit  $N \rightarrow \infty$ . This will turn our discrete problem into a continuous problem, and our differences into derivatives.

With  $N$  masses, we called the displacement of each mass from its equilibrium point  $x_n$ . Since all the springs have the same constant, at equilibrium, all the masses are a distance  $\Delta x$  apart. Let us define a function  $A(\Delta x, t)$  as the **amplitude** of the displacement from equilibrium at a point  $x$ . So,

$$A(n\Delta x, t) = x_n(t) \quad (46)$$

Thus, plots like this



(47)

are plots of  $A(x, t)$ . To be clear, these displacements are still longitudinal (in the direction of the springs), but we are drawing  $A(x, t)$  in the transverse direction.  $A(x, t)$  so far only has values at discrete points given by  $x = n\Delta x$ . Its value at those points is  $x_n(t)$ .

In terms of  $A(x, t)$ , the equations of motion for the coupled system, Eq. (18)

$$\frac{d^2 x_n}{dt^2} = \frac{k}{m} [x_{n-1} - 2x_n + x_{n+1}] \quad (48)$$

become

$$\frac{\partial^2}{\partial t^2} A(n\Delta x, t) = \frac{k}{m} [A((n+1)\Delta x, t) - 2A(n\Delta x, t) + A((n-1)\Delta x, t)] \quad (49)$$

$$= \frac{k}{m} \Delta x \left[ \frac{A(n\Delta x + \Delta x, t) - A(n\Delta x, t)}{\Delta x} - \frac{A(n\Delta x, t) - A(n\Delta x - \Delta x, t)}{\Delta x} \right] \quad (50)$$

Writing  $x = n\Delta x$  this becomes

$$\frac{\partial^2}{\partial t^2} A(x, t) = \frac{k}{m} \Delta x \left[ \frac{A(x + \Delta x, t) - A(x, t)}{\Delta x} - \frac{A(x, t) - A(x - \Delta x, t)}{\Delta x} \right] \quad (51)$$

Starting to look like calculus...

As  $\Delta x \rightarrow 0$ , this becomes

$$\frac{d^2}{dt^2} A(x, t) = \frac{k}{\mu} \left[ \frac{\partial A(x, t)}{\partial x} - \frac{\partial A(x - \Delta x, t)}{\partial x} \right] \quad (52)$$

where  $\mu = \frac{m}{\Delta x}$  is the mass per unit length or **mass density**. We also define  $E = k\Delta x$  as the **elastic modulus** to get

$$\frac{d^2}{dt^2} A(x, t) = \frac{E}{\mu} \frac{1}{\Delta x} \left[ \frac{\partial A(x, t)}{\partial x} - \frac{\partial A(x - \Delta x, t)}{\partial x} \right] \quad (53)$$

Now take  $\Delta x \rightarrow 0$  turns the first derivatives into second derivatives. Writing

$$v \equiv \sqrt{\frac{E}{\mu}} \quad (54)$$

and we find

$$\boxed{\frac{\partial^2}{\partial t^2} A(x, t) = v^2 \frac{\partial^2}{\partial x^2} A(x, t)} \quad (55)$$

This is the **wave equation**.

## 6 Solving the wave equation

The wave equation is linear, so we can solve it with exponentials. Writing

$$A(x, t) = e^{i\omega t} e^{ikx} \quad (56)$$

we get

$$\omega^2 = v^2 k^2 \quad (57)$$

So

$$\omega(k) = v|k| \quad (58)$$

This is a linear dispersion relation. Since we have taken  $N \rightarrow \infty$ , all modes have  $j \ll N$ , thus, the linearity of the dispersion relation is consistent with what we found for finite  $N$ ,

Since the wave equation (without damping) has only second derivatives, its easy to see that the general solution for a fixed frequency  $\omega$  is

$$A_k(x, t) = a_k \cos(kx) \cos(\omega t) + b_k \sin(kx) \cos(\omega t) + c_k \cos(kx) \sin(\omega t) + d_k \sin(kx) \sin(\omega t) \quad (59)$$

exactly as in the discrete case. The only difference is that now

$$\omega(k) = vk \quad (60)$$

instead of the more complicated  $\omega(p) = \sqrt{2(1 - \cos(p))} \omega_0$  we found before.

Note that in the continuum case  $k$  has dimensions of  $\frac{1}{\text{length}}$ . We call  $k$  the wavenumber, which is equal to  $\frac{2\pi}{\text{wevelength}}$ . In the discrete case  $p$  was dimensionless. Note that this  $k$  has nothing to do with the spring constant – we just use the same letter “ $k$ ” for both.

Which of  $a_k, b_k, c_k$  or  $d_k$  vanish depends on boundary conditions. Let’s consider some interesting cases. First, we note that one solution is

$$A(x, t) = \cos(kx) \cos(\omega t) - \sin(kx) \sin(\omega t) = \cos(kx - \omega t) \quad (61)$$

$$= \cos\left(\frac{\omega}{v}(x - vt)\right) \quad (62)$$

This solution has the property that  $A(x, t + \Delta t) = A(x - v\Delta t, t)$  meaning that the amplitude at  $x$  in the future is given by the amplitude at position to the left at current time. In other words, the curve is moving to the right. This is a **right-moving traveling wave**.

More generally, we note that for any function  $f(z)$  the amplitude

$$A(x, t) = f(x - vt) \quad (63)$$

will satisfy the wave equation. Indeed, it is easy to check that

$$\frac{\partial^2}{\partial t^2} f(x - vt) = v^2 f''(x - vt) \quad (64)$$

$$\frac{\partial^2}{\partial x^2} f(x - vt) = f''(x - vt) \quad (65)$$

So that

$$\left[ \frac{\partial^2}{\partial t^2} - v^2 \frac{\partial^2}{\partial x^2} \right] f(x - vt) = 0 \quad (66)$$

Thus  $A(x, t) = f(x - vt)$  is a general right-moving traveling wave. In general it will not be associated with a fixed frequency. However, since any solution can be written as a sum over normal modes, any traveling wave can be written as a sum over solutions of fixed frequency. How this is done is known as the Fourier decomposition, which we study next time.

Waves of the form

$$A(x, t) = f(x + vt) \quad (67)$$

are also solutions for any  $f$ . These are **left-moving** traveling waves.

Another solution of fixed frequency is

$$A(x, t) = \cos(kx - \omega t) + \cos(kx + \omega t) = 2\cos(kx)\cos(\omega t) \quad (68)$$

This solution has the property that the amplitudes at any two points  $x_1$  and  $x_2$  always have the same ratio at any time

$$\frac{A(x_1, t)}{A(x_2, t)} = \frac{2\cos(kx_1)}{2\cos(kx_2)} \quad (69)$$

These are **standing waves**.

We see that

- Whether a traveling wave or a standing wave is produced depends on initial conditions.
- Standing waves are the sum of a left-moving and right-moving wave.

# Lecture 5: Fourier series

## 1 Fourier series

When  $N$  oscillators are strung together in a series, the amplitude of that string can be described by a function  $A(x, t)$  which satisfies the **wave equation**:

$$\left[ \frac{\partial^2}{\partial t^2} - v^2 \frac{\partial^2}{\partial x^2} \right] A(x, t) = 0 \quad (1)$$

We saw that electromagnetic fields satisfy this same equation with  $v = c$  the speed of light.

We found normal mode solutions of the form

$$A(x, t) = A_0 \cos\left(\frac{\omega}{v}(x \pm vt) + \phi\right) \quad (2)$$

for any  $\omega$  which are **traveling waves**. Solutions of the form

$$A(x, t) = A_0 \cos(kx) \cos(\omega t) \quad (3)$$

with  $\omega^2 = v^2 k^2$  are called **standing waves**. Whether traveling waves or standing waves are relevant depends on the boundary condition.

More generally, we found traveling wave solutions could come from any function  $f(x + vt)$ :

$$\left[ \frac{\partial^2}{\partial t^2} - v^2 \frac{\partial^2}{\partial x^2} \right] f(x + vt) = 0 \quad (4)$$

Similarly  $f(x - vt)$  is a solution. Functions  $f(x - vt)$  are right-moving traveling waves and functions  $f(x + vt)$  are left-moving traveling waves.

Now, since any vector can be written as a sum of eigenvectors, any solution can be written as a sum of normal modes. This is true both in the discrete case and in the continuum case. Thus we must be able to write

$$f(x + vt) = \sum_k a_k \cos(kx) \cos(\omega t) + b_k \sin(kx) \cos(\omega t) + c_k \cos(kx) \sin(\omega t) + d_k \sin(kx) \sin(\omega t) \quad (5)$$

where the sum is over wavenumbers  $k$ . In particular, at  $t = 0$  any function can be written as

$$f(x) = \sum_k a_k \cos(kx) + b_k \sin(kx) \quad (6)$$

We have just proved Fourier's theorem!

(Ok, we haven't really proven it, we just assumed the result from linear algebra about a finite system applies also in the continuum limit. The actual proof requires certain properties about the smoothness of  $f(x)$  to hold. But we are physicists not mathematicians, so let's just say we proved it.)

## 2 Fourier's theorem

Fourier's theorem states that any square-integrable function<sup>1</sup>  $f(x)$  which is periodic on the interval  $0 < x \leq L$  (meaning  $f(x + L) = f(x)$ ) can be written as

$$f(x) = a_0 + \sum_{n=1}^{\infty} a_n \cos\left(\frac{2\pi n}{L}x\right) + \sum_{n=1}^{\infty} b_n \sin\left(\frac{2\pi n}{L}x\right) \quad (7)$$

---

1. A function is square-integrable if  $\int_0^L dx f(x)^2$  exists.

with

$$a_0 = \frac{1}{L} \int_0^L dx f(x) \quad (8)$$

$$a_n = \frac{2}{L} \int_0^L dx f(x) \cos\left(\frac{2\pi n}{L}x\right) \quad (9)$$

$$b_n = \frac{2}{L} \int_0^L dx f(x) \sin\left(\frac{2\pi n}{L}x\right) \quad (10)$$

This decomposition is known as a **Fourier series**. Fourier series are useful for periodic functions or functions on a fixed interval  $L$  (like a string). One can do a similar analysis for non-periodic functions or functions on an infinite interval ( $L \rightarrow \infty$ ) in which case the decomposition is known as a Fourier transform. We will study Fourier series first.

It is easy to verify these formulas for  $a_n$  and  $b_n$ . For  $a_0$ , we just integrate  $f(x)$ . Since  $\cos\left(\frac{2\pi}{L}nx\right)$  goes through  $n$  cycles of the complete cosine curve as  $x$  goes from 0 to  $L$ , we have

$$\int_0^L dx \cos\left(\frac{2\pi n}{L}x\right) = 0, \quad n > 0 \quad (11)$$

Similarly,

$$\int_0^L dx \sin\left(\frac{2\pi n}{L}x\right) = 0, \quad n > 0 \quad (12)$$

Thus,

$$\int_0^L dx f(x) = a_0 \int_0^L dx + \sum_{n=1}^{\infty} a_n \int_0^L dx \cos\left(\frac{2\pi n}{L}x\right) + \sum_{n=1}^{\infty} b_n \int_0^L dx \sin\left(\frac{2\pi n}{L}x\right) \quad (13)$$

$$= a_0 L \quad (14)$$

in agreement with Eq. (8).

For  $a_n$  we can use the cosine sum formula to write

$$\int_0^L dx \cos\left(\frac{2\pi m}{L}x\right) \cos\left(\frac{2\pi n}{L}x\right) = \int_0^L dx \left[ \frac{1}{2} \cos\left(\frac{n+m}{L}2\pi x\right) + \frac{1}{2} \cos\left(\frac{n-m}{L}2\pi x\right) \right] \quad (15)$$

Now again we have that these integrals all vanish over an integer number of periods of the cosine curve. The only way this wouldn't vanish is if  $n - m = 0$ . So we have for  $n > 0$

$$\int_0^L dx \cos\left(\frac{2\pi m}{L}x\right) \cos\left(\frac{2\pi n}{L}x\right) = \frac{1}{2} \delta_{mn} \int_0^L dx = \frac{L}{2} \delta_{mn} \quad (16)$$

where  $\delta_{mn}$  is the **Kronecker  $\delta$ -function**

$$\delta_{mn} = \begin{cases} 0, & m \neq n \\ 1 & m = n \end{cases} \quad (17)$$

Similarly,

$$\int_0^L dx \cos\left(\frac{2\pi m}{L}x\right) \sin\left(\frac{2\pi n}{L}x\right) = 0 \quad (18)$$

$$\int_0^L dx \sin\left(\frac{2\pi m}{L}x\right) \sin\left(\frac{2\pi n}{L}x\right) = \frac{L}{2} \delta_{mn} \quad (19)$$

Thus, for  $n > 0$

$$\int_0^L dx f(x) \cos\left(\frac{2\pi n}{L}x\right) \quad (20)$$

$$= \int_0^L dx \left[ a_0 + \sum_{m=1}^{\infty} a_m \cos\left(\frac{2\pi m}{L}x\right) + \sum_{m=1}^{\infty} b_m \sin\left(\frac{2\pi m}{L}x\right) \right] \cos\left(\frac{2\pi n}{L}x\right) \quad (21)$$

$$= \frac{L}{2} \sum_{m=1}^{\infty} a_m \delta_{mn} \quad (22)$$

$$= \frac{L}{2} a_n \quad (23)$$

as in Eq. (9). In the same way you can check the formula for  $b_n$ .

We use

- Fourier cosine series for functions which are **even** on the interval ( $f(x) = f(L - x)$ )
- Fourier sine series for functions which are **odd** on the interval ( $f(x) = -f(L - x)$ )
- For functions that are neither even nor odd on the interval, we need both sines and cosines

### 3 Example

Find the Fourier series for the sawtooth function:

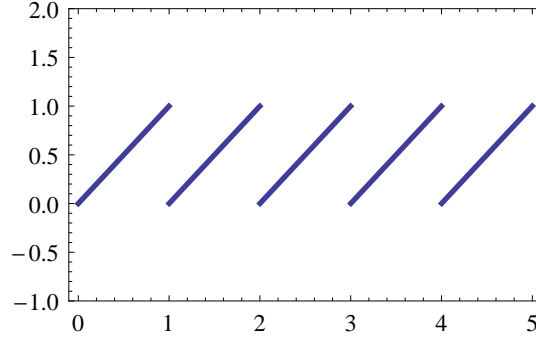


Figure 1. Sawtooth function

This function is clearly periodic. It is equal to  $f(x) = x$  on the interval  $0 < x \leq 1$ . Thus we can compute the Fourier series with  $L = 1$ . We get

$$a_0 = \frac{1}{L} \int_0^L f(x) dx = \int_0^1 dx x = \frac{1}{2} \quad (24)$$

Next,

$$a_n = \frac{2}{L} \int_0^L dx f(x) \cos\left(\frac{2\pi n}{L}x\right) = 2 \int_0^1 dx x \cos(2\pi n x) \quad (25)$$

This can be done by integration by parts

$$a_n = 2 \left. \frac{x}{2\pi n} \sin(2\pi n x) \right|_0^L - 2 \int_0^1 dx \sin(2\pi n x) = 0 \quad (26)$$

Finally,

$$b_n = 2 \int_0^1 dx x \sin(2\pi n x) \quad (27)$$

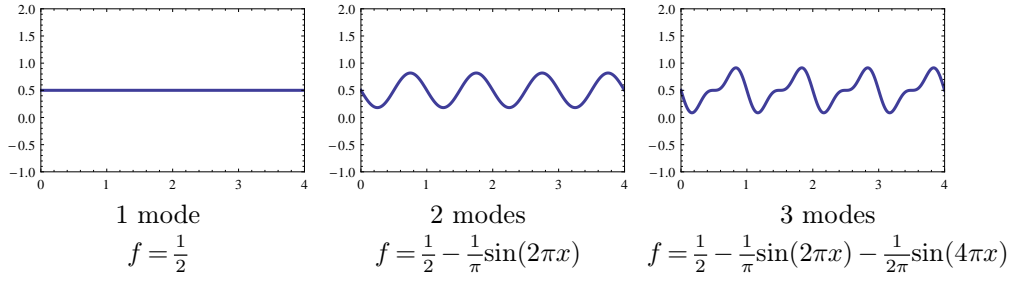
$$= -2 \left. \frac{x}{2\pi n} \cos(2\pi n x) \right|_0^L + 2 \int_0^1 dx \cos(2\pi n x) \quad (28)$$

$$= -\frac{1}{\pi n} \quad (29)$$

Thus we have

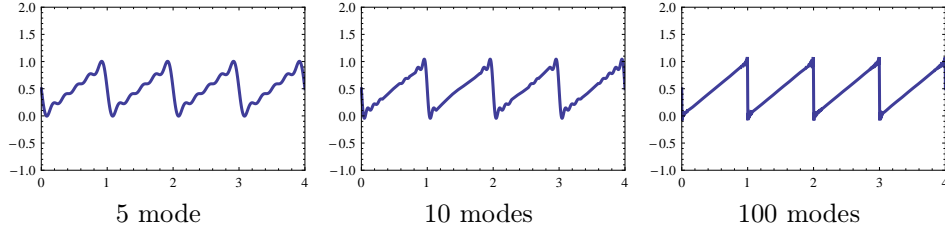
$$\boxed{f(x) = \frac{1}{2} + \sum_{n=1}^{\infty} -\frac{1}{\pi n} \sin \frac{2\pi n x}{L}} \quad (30)$$

Let's look at how well the series approximates the function when including various terms. Taking 0, 1 and 2 terms in the sum gives



**Figure 2.** Approximations to the sawtooth function

Already at 3 modes, it's looking reasonable. For 5, 10 and 100 modes we find

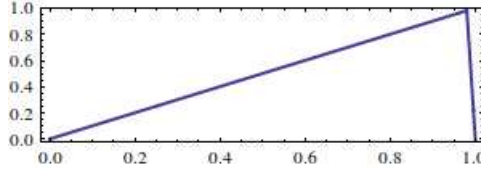


**Figure 3.** More approximations to the sawtooth function

For 10 modes we find excellent agreement.

## 4 Plucking a string

Let's apply the Fourier decomposition we worked out to plucking a string. Suppose we pluck a string by pulling up one end:



What happens to the string? To find out, let us do a Fourier decomposition of the  $x$ -dependence of the pluck. We start by writing

$$A(x, t) = \sum_{n=0}^{\infty} \left[ a_n \cos\left(\frac{2n\pi}{L}x\right) \cos(\omega_n t) + b_n \sin\left(\frac{2n\pi}{L}x\right) \cos(\omega_n t) \right], \quad \omega_n = \frac{2n\pi}{L}v \quad (31)$$

Here  $v$  is the speed of sound in the string. For a given wavenumber,  $k_n = \frac{2n\pi}{L}$ , we know that  $\omega_n = k_n v$  to satisfy the wave equation. We could also have included components with  $\sin(\omega_n t)$ ; however since the string starts off at rest (so that  $\partial_t A(x, t) = 0$ ), then the coefficients of the  $\sin(\omega_n t)$  functions must all vanish.

At time  $t = 0$ , the amplitude is

$$A(x, 0) = \sum_{n=0}^{\infty} \left[ a_n \cos\left(\frac{2n\pi}{L}x\right) + b_n \sin\left(\frac{2n\pi}{L}x\right) \right] \quad (32)$$

This is just the Fourier decomposition of the function described by our pluck shape. If we approximate the pluck as the sawtooth function from the previous section, then we already know that

$$a_n = 0, \quad b_n = -\frac{1}{\pi n} \quad (33)$$

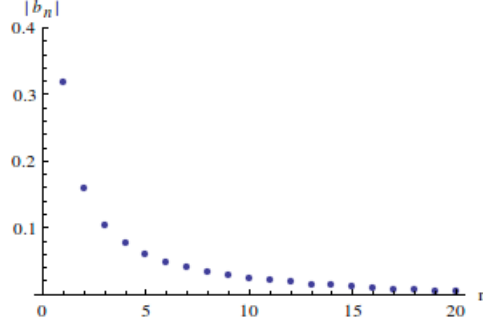


So that, setting  $L = 1$

$$A(x, t) = \sum_{n=1}^{\infty} -\frac{1}{\pi n} \sin(2\pi x) \cos(2\pi n t) \quad (34)$$

This gives the motion of the string for all time.

The relative amplitudes of each mode are



**Figure 4.** Amplitudes of the relative harmonics of a string plucked with a sawtooth plucking.

The  $n = 1$  mode is the largest. This the **fundamental** frequency of the string. Thus the sound that comes out of the string will be mostly this frequency:  $\omega_1 = \frac{2\pi}{L}v$ . The modes with  $n > 1$  are the **harmonics**. Harmonics have frequencies which are integer multiples of the fundamental.

## 5 Exponentials

Fourier series decompositions are even easier with complex numbers. There we can replace the sines and cosines by exponentials. The series is

$$f(x) = \sum_{n=-\infty}^{\infty} c_n e^{in x \frac{2\pi}{L}} \quad (35)$$

where

$$c_n = \frac{1}{L} \int_0^L dx f(x) e^{-in x \frac{2\pi}{L}} \quad (36)$$

To check this, we substitute in

$$\int_0^L dx f(x) e^{-in x \frac{2\pi}{L}} = \sum_{m=-\infty}^{\infty} c_m \int_0^L dx e^{im x \frac{2\pi}{L}} e^{-in x \frac{2\pi}{L}} = \sum_{m=-\infty}^{\infty} c_m \int_0^L dx e^{i(m-n) x \frac{2\pi}{L}} \quad (37)$$

If  $n \neq m$  then,

$$\int_0^L dx e^{i(m-n) x \frac{2\pi}{L}} = \frac{L}{2\pi(m-n)} e^{i(m-n) x \frac{2\pi}{L}} \Big|_0^L \quad (38)$$

$$= \frac{L}{2\pi(m-n)} [e^{i2\pi(m-n)} - 1] = 0 \quad (39)$$

If  $m = n$ , then the integral is just

$$\int_0^L dx = L \quad (40)$$

Thus,

$$\int_0^L dx e^{i(m-n) x \frac{2\pi}{L}} = L \delta_{mn} \quad (41)$$

and so

$$\int_0^L dx f(x) e^{-in x \frac{2\pi}{L}} = L c_n \quad (42)$$

If  $f(x)$  is real, then

$$f(x) = \sum_{n=-\infty}^{\infty} \operatorname{Re}(c_n + c_{-n}) \cos\left(\frac{2\pi n x}{L}\right) + \operatorname{Im}(c_{-n} + c_n) \sin\left(\frac{2\pi n x}{L}\right) \quad (43)$$

So  $a_n = \operatorname{Re}(c_n + c_{-n})$  and  $b_n = \operatorname{Im}(c_{-n} + c_n)$ . Thus, the exponential series contains all the information in both the sine and cosine series in an efficient form.

## 6 Orthogonal functions (optional)

In verifying Fourier's theorem, we found the relevant integral equations

$$\frac{1}{L} \int_0^L dx e^{i(m-n)x \frac{2\pi}{L}} = \delta_{mn} \quad (44)$$

$$\frac{1}{L} \int_0^L dx \cos\left(\frac{2\pi m}{L}x\right) \sin\left(\frac{2\pi n}{L}x\right) = 0 \quad (45)$$

$$\frac{1}{L} \int_0^L dx \cos\left(\frac{2\pi m}{L}x\right) \cos\left(\frac{2\pi n}{L}x\right) = \frac{1}{2} \delta_{nm} \quad (46)$$

$$\frac{1}{L} \int_0^L dx \sin\left(\frac{2\pi m}{L}x\right) \sin\left(\frac{2\pi n}{L}x\right) = \frac{1}{2} \delta_{nm} \quad (47)$$

These are examples of orthogonal functions. The integral is a type of **inner product**. The dot-product among vectors is another example of an inner product. We can write the inner product in various ways

$$\langle v | w \rangle \equiv \vec{v} \cdot \vec{w} = \sum_i v_i w_i \quad (48)$$

The integral inner product is a generalization of this from vectors of numbers to functions.

We can define the inner product of two functions as

$$\langle f | g \rangle = \frac{1}{2\pi} \int_0^{2\pi} dx f^*(x) g(x) \quad (49)$$

where  $f^*(x)$  is the complex conjugate of  $f(x)$ . For example,

$$\langle e^{imx} | e^{inx} \rangle = \delta_{mn} \quad (50)$$

This is the analog of

$$\langle x_i | x_j \rangle = \delta_{ij} \quad (51)$$

where  $|x_i\rangle = (0, \dots, 0, 1, 0, 0)$  with the 1 in the  $i^{\text{th}}$  component. That is the  $|x_i\rangle$  are the unit vectors. When a set of functions satisfy

$$\langle f_i | f_j \rangle = \delta_{ij} \quad (52)$$

we say that they are **orthonormal**. The **ortho** part means they are **orthogonal**:  $\langle f_i | f_j \rangle = 0$  for  $i \neq j$ . The **normal** part means they are normalized,  $\langle f_i | f_j \rangle = 1$  for  $i = j$ .

If any function can be written as a linear combination of functions  $f_i$  we say that the set  $\{f_i\}$  is **complete**. Then

$$f(x) = \sum_i a_i f_i(x) \quad (53)$$

We can extract  $a_i$  via

$$\langle f(x) | f_i \rangle = \sum_j a_j \langle f_j | f_i \rangle = \sum_j a_j \delta_{ij} = a_i \quad (54)$$

This is exactly what we did with the Fourier decomposition above. It is also what we do with vectors

$$\vec{v} = \sum c_i \vec{x}_i \quad (55)$$

Then  $c_i = \langle v | x_i \rangle$  which is just the  $i^{\text{th}}$  component of  $\vec{v}$ .

We will see various sets of orthonormal function bases with different inner products come up in physics. Other examples are:

**Bessel functions:**  $\mathcal{J}_n(x)$ . These functions are the solutions to the differential equation

$$x^2 f''(x) + x f'(x) + (x^2 - n^2) f(x) = 0 \quad (56)$$

They satisfy the orthonormality condition

$$\langle \mathcal{J}_n | \mathcal{J}_m \rangle = \int_0^1 dx x \mathcal{J}_n(x) \mathcal{J}_m(x) = \delta_{nm} \quad (57)$$

Bessel functions come up in 2 dimensional problems. You will start to see them all over the place in physics.

**Legendre polynomials**  $P_n$ . These satisfy  $P_0(x) = 1$ ,  $P_1(x) = x$  and

$$(n+1)P_{n+1}(x) = (2n+1)xP_n(x) - nP_{n-1}(x) \quad (58)$$

Their inner product is

$$\langle P_n | P_m \rangle = \int_0^1 dx P_n(x) P_m(x) = \frac{2}{2n+1} \delta_{nm} \quad (59)$$

Legendre polynomials come up in problems with spherical symmetry. You will study them to death in quantum mechanics.

**Hermite polynomials**

$$H_n(x) = (-1)^n e^{-\frac{x^2}{2}} \frac{d}{dx^2} e^{-\frac{x^2}{2}} \quad (60)$$

So  $H_0(x) = 1$ ,  $H_1(x) = x$ ,  $H_2(x) = x^2 - 1$  and so on. These satisfy

$$\langle H_n | H_m \rangle = \int_0^1 dx e^{-\frac{x^2}{2}} H_n(x) H_m(x) = \delta_{nm} \quad (61)$$

Hermite polynomials play a critical role in the quantum harmonic oscillator.

# Lecture 6: Waves in strings and air

## 1 Introduction

In Lecture 4, we derived the wave equation for two systems. First, by stringing together masses and springs and taking the continuum limit, we found

$$\left[ \begin{array}{c} k \quad k \quad k \quad k \\ \text{---} \text{---} \text{---} \text{---} \\ m \quad m \quad m \\ x_1 \quad x_2 \quad x_3 \end{array} \right] \Rightarrow \left[ \frac{\partial^2}{\partial t^2} - v^2 \frac{\partial^2}{\partial x^2} \right] A(x, t) = 0 \quad (1)$$

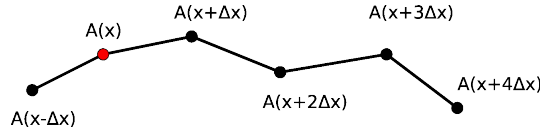
where  $A(x, t)$  is the displacement from equilibrium of the mass at position  $x$ . These are longitudinal waves. In this equation, for waves in a solid, the wave speed is

$$v = \sqrt{\frac{E}{\mu}} \quad (2)$$

where  $E$  is the elastic modulus and  $\mu$  is the density per unit length. Now we consider two more cases: transverse oscillations on a string and longitudinal motion of air molecules (sound waves).

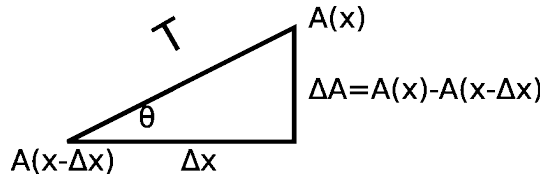
## 2 Transverse oscillations

Consider a string of tension  $T$ . We define the amplitude of the string at a point  $x$  at time  $t$  as  $A(x, t)$ . In this section, we'll sometimes write  $A(x, t)$  just as  $A(x)$  to avoid clutter. Let us treat the string as a bunch of massless test probes connected by a elastic strings. Then we can draw the picture as



What is the force acting on the test mass at position  $x$  (in red)?

First, consider the *downward* component of the force pulling on the test mass at  $x$  from the mass to the left (at  $x - \Delta x$ ). We can draw a triangle:



The force is given by

$$F_{\text{downwards, from left mass}} = T \sin \theta = T \frac{\Delta A}{\sqrt{\Delta A^2 + \Delta x^2}} \quad (3)$$

If the system is close to equilibrium, then the slope will be small. That is,  $\Delta A \ll \Delta x$ . In this case, we can approximate  $\sqrt{\Delta A^2 + \Delta x^2} \approx \Delta x$  and so

$$F_{\text{downwards, from left mass}} = T \frac{\Delta A}{\Delta x} = T \frac{A(x) - A(x - \Delta x)}{\Delta x} = T \frac{\partial A}{\partial x} \quad (4)$$

where we have taken  $\Delta x \rightarrow 0$  in the last step turning the difference into a derivative.

Similarly, the *downward* force from the mass on the right is

$$F_{\text{downwards, from right mass}} = T \frac{A(x) - A(x + \Delta x)}{\Delta x} = -T \frac{\partial A(x + \Delta x)}{\partial x} \quad (5)$$

Thus,

$$F_{\text{total downwards}} = -T \left[ \frac{\partial A(x + \Delta x)}{\partial x} - \frac{\partial A}{\partial x} \right] \quad (6)$$

Now we use  $F = ma$ , where  $a = -\frac{\partial^2 A}{\partial t^2}$  is the downward acceleration. So we should have

$$F_{\text{total downwards}} = -m \frac{\partial^2 A}{\partial t^2} = -\mu \Delta x \frac{\partial^2 A}{\partial t^2} \quad (7)$$

Plugging this into Eq. (6) we find

$$\frac{\partial^2 A}{\partial t^2} = \frac{T}{\mu} \left[ \frac{\frac{\partial A(x + \Delta x)}{\partial x} - \frac{\partial A}{\partial x}}{\Delta x} \right] = \frac{T}{\mu} \frac{\partial^2 A}{\partial x^2} \quad (8)$$

Thus,

$$\left[ \frac{\partial^2}{\partial t^2} - v^2 \frac{\partial^2}{\partial x^2} \right] A(x, t) = 0 \quad \text{with} \quad v = \sqrt{\frac{T}{\mu}} \quad (9)$$

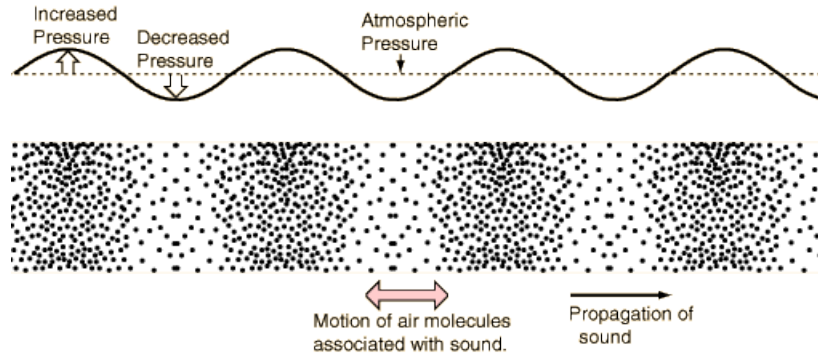
So the wave equation is again satisfied with a wave speed  $v = \sqrt{\frac{T}{\mu}}$ .

Note that the acceleration is due to a **difference of forces**. The force pulling up from the right has to be different from the force pulling down from the left to get an acceleration. Each force is proportional to a first derivative, thus the acceleration is proportional to a second derivative.

### 3 Sound waves

Waves in air are just like waves in a solid: the air molecules are like little masses and the forces between them act like springs. Thus we have already derived the wave equation. What's left is to think about what's actually going on when a wave propagates through the air.

Sound waves are longitudinal density waves, which look like



**Figure 1.** Visualization of sound waves

What is the amplitude for a sound wave? As always, the amplitude  $A(x, t)$  measures the displacement from equilibrium. In fact, in a sound wave, each individual molecule is just oscillating back and forth around an equilibrium position, and the wave appears as a collective phenomenon among these moving molecules. This is easiest to see in an animation. Try the animation on this web page <http://www.acs.psu.edu/drussell/demos/waves/wavemotion.html> under longitudinal waves. In this animation,  $A(x, t)$  is the displacement from equilibrium at time  $t$  of the red dot whose equilibrium position is at  $x$ . It has the form  $A(x, t) = A_0 \cos(kx - \omega t + \phi_0)$  for some overall amplitude  $A_0$  and some phase  $\phi_0$ . From the point of view of the molecules strung together, this system is identical to the masses and springs strung together that we discussed in lecture 4. So the derivation of the wave equation for a gas is identical.

In a snapshot of the wave in Fig. 1, it's hard to see which molecules are at their equilibrium position and which are not. That is, it's hard to see  $A(x, t)$ . Instead, we see the density of the gas  $\rho(x, t)$ . In fact the two are related. Looking the animations for a long time, you can see that in fact there is a close relation between the amplitude (displacement from equilibrium)  $A(x, t)$  and the density  $\rho(x, t)$ . As with any oscillator, the molecules move fastest as they pass through their equilibrium going left or right, and stop when they are farthest from equilibrium. Now note that when the molecules (red dots on the web animation) are moving fastest to the right, they are in the most dense region, and when they moving fastest to the left, they are in the least dense region of the gas. That is, the maximal velocity in the  $+x$  direction corresponds to the maximal density and the minimal velocity in the  $+x$  direction to the minimal density. Therefore density agrees with velocity:  $\rho \propto \rho_0 + \frac{dA}{dt}$ . In other words, if the displacement is  $A(x, t) = A_0 \cos(kx - \omega t)$  then the density is  $\rho(x, t) = \rho_0 + (\Delta\rho) \sin(kx - \omega t)$ . Another way to say this is that the density lags behind the amplitude by  $90^\circ$ .

### 3.1 Speed of sound in a gas

Let us consider the case where the sound wave is excited by a large membrane like a drum or a speaker. If we are interested in wavelengths much less than the size of the membrane, and much larger than the distance between air molecules, then waves in air become exactly like waves in a solid or waves on a string. We simply have to divide by the unit area:

$$\frac{\mu}{A} \frac{\partial^2 A}{\partial t^2} = \frac{T}{A} \frac{\partial^2 A}{\partial x^2} \quad (10)$$

Now,  $\frac{\mu}{A}$  is the mass per unit length per unit area, also known as mass per unit volume or density  $\rho$ . Also,  $\frac{T}{A}$  is the force per unit area or the pressure, so we get

$$\rho \frac{\partial^2 A}{\partial t^2} = p \frac{\partial^2 A}{\partial x^2} \quad (11)$$

Thus  $v = \sqrt{\frac{p}{\rho}}$  for a gas. It turns out that this is only correct at constant temperature.

At constant temperature, the gas doesn't heat up. You may remember the ideal gas law from chemistry:

$$pV = nRT \quad (12)$$

where  $V$  is the volume,  $n$  is the number of molecules,  $R$  is the ideal gas constant and  $T$  is the temperature. Dividing by  $V$  and using  $\rho = \frac{n}{V}m$  with  $m$  the molecular weight of the gas molecules, we get

$$p = \rho \frac{RT}{m} \quad (13)$$

This means that

$$\left( \frac{dp}{d\rho} \right)_T = \frac{RT}{m} = \frac{p}{\rho} \quad (14)$$

where the subscript  $T$  means "constant temperature". Thus we would have  $v = \sqrt{\left( \frac{\partial p}{\partial \rho} \right)_T}$  for a gas which could not heat up. There is unfortunately no such gas.

The correct velocity for a wave in air is

$$v = \sqrt{\left( \frac{\partial p}{\partial \rho} \right)_S} \quad (15)$$

where the subscript  $S$  means "constant entropy". It should be constant entropy since when a wave passes through some air, it leaves the air in the same state it started in, without increasing the entropy. If you take physical chemistry (physics 181 or chem 161), you can study these constrained partial derivatives to death. I will just summarize the important result

$$\left( \frac{\partial p}{\partial \rho} \right)_S = \gamma \frac{p}{\rho} \quad (16)$$

where

$$\gamma = \frac{C_P}{C_V} \quad (17)$$

where  $C_P$  is the specific heat at constant pressure and  $C_V$  is the specific heat at constant volume.

A more useful form of  $\gamma$  is

$$\gamma = \frac{f+2}{f} \quad (18)$$

where  $f$  is the number of degrees of freedom in the gas. For a monatomic gas, like argon, the only degrees of freedom are from translations. For the  $x$ ,  $y$  and  $z$  directions, we get  $f=3$ . So

$$\gamma = \frac{3+2}{3} = \frac{5}{3} = 1.67 \quad (\text{monatomic gas}) \quad (19)$$

For a diatomic gas, like  $N_2$  or  $O_2$  (which is mostly what air is), both atoms can move, so we would get  $f=6$ , however, if they rotate around the bond axis, the molecule is unchanged, so in fact  $f=5$ . You can think of 5 as 3 translations, one rotation and one vibration along the bond axis. Thus

$$\gamma = \frac{5+2}{5} = \frac{7}{5} = 1.4 \quad (\text{diatomic gas like air}) \quad (20)$$

To match to the notation for waves in solids we sometimes define a **bulk modulus**

$$B \equiv \gamma p \quad (21)$$

Then the speed of sound in air is

$$c_s = \sqrt{\frac{\gamma p}{\rho}} = \sqrt{\frac{B}{\rho}} \quad (22)$$

Note that  $B$  and  $c_s$  are properties of the gas, not the wave. All waves have the same velocity in the same type of air.

Another useful formula is that, using the ideal gas law,

$$c_s = \sqrt{\frac{\gamma p}{\rho}} = \sqrt{\gamma \frac{RT}{m}} \quad (23)$$

This tells us that the speed of sound only depends on the temperature of a gas, not on its density or pressure separately. It also tells us the speed of sound is different in two gases with the same temperature but different molecular masses  $m$ .

You may also recall from chemistry that the root-mean-square velocity of gas is determined by its temperature:  $v_{\text{rms}} = \sqrt{\frac{3RT}{m}}$ . Again, this is something you will show in a physical chemistry class. So

$$c_s = \sqrt{\frac{\gamma}{3}} v_{\text{rms}} \quad (24)$$

Thus the speed of sound is proportional to, but not greater than, the speed of the molecules in the gas. This makes sense – how could it be sound travel faster than the molecules transmitting it?

### 3.2 Summary

For sound waves, the amplitude of the wave  $A(x, t)$  is the displacement from  $x$  of the molecule whose equilibrium position is at  $x$ . The density of the wave  $\rho(x, t)$  also oscillates and lags in phase behind the amplitude by a quarter wavelength,  $\frac{\pi}{2}$ . Sound waves satisfy the wave equation with a sound speed

$$v = \sqrt{\gamma \frac{p}{\rho}} = \sqrt{\gamma \frac{RT}{m}} \quad (25)$$

where  $p$  is the average pressure,  $\rho$  the average density and  $T$  the average temperature. Also,

$$\gamma = \frac{C_P}{C_V} = \frac{f+2}{f} \quad (26)$$

where  $C_P$  is the specific heat at constant pressure and  $C_V$  is the specific heat at constant volume and  $f$  is the number of degrees of freedom of the gas molecules. For a monotonic gas like Ar,  $f=3$  and  $v = \sqrt{1.67 \frac{P}{\rho}}$ . For a diatomic gas like  $N_2$  or  $O_2$ ,  $f=5$  and  $v = \sqrt{1.4 \frac{P}{\rho}}$ .

## 4 Standing waves

Now lets talk about standing wave solutions in more detail. Again, we consider the wave equation

$$\left[ \frac{\partial^2}{\partial t^2} - v^2 \frac{\partial^2}{\partial x^2} \right] A(x, t) = 0 \quad (27)$$

and we would like solutions of fixed frequency  $\omega$ . These are solutions which are periodic in time. We can write the general such solution as a sum of terms of the form

$$A(x, t) = A_0 \sin(kx + \phi_1) \sin(\omega t + \phi_2) \quad (28)$$

In this solution,  $A_0$  is the **amplitude** and  $k$  the **wavenumber**. The frequency determined from the wavenumber through the dispersion relation

$$\omega = vk \quad (29)$$

There are two phases  $\phi_1$  and  $\phi_2$ . Instead of using phases, we could write the general solution as

$$A(x, t) = A_0 \sin(kx) \sin(\omega t) + A_1 \sin(kx) \cos(\omega t) + A_2 \cos(kx) \cos(\omega t) + A_3 \cos(kx) \sin(\omega t) \quad (30)$$

The two forms are equivalent and we will go back and forth between them as convenient.

Consider first the case where one of the boundary conditions is that the string is fixed at  $x = 0$ . That is

$$\boxed{A(0, t) = 0} \quad (31)$$

This is known as a fixed, closed, or **Dirichlet** boundary condition. If there were a  $A_3 \cos(kx) \sin(\omega t)$  component, then the  $x = 0$  point would oscillate as  $x(0, t) = A_3 \sin(\omega t)$  meaning it is not fixed. Thus  $A_3 = 0$ . Similarly,  $A_2 = 0$ . Thus the general solution with  $A(0, t) = 0$  is

$$A(x, t) = A_0 \sin(kx) \sin(\omega t + \phi) \quad (32)$$

If we fix the other end of the string at  $x = L$  then we must have  $\sin(kL) = 0$  which implies

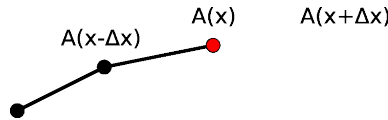
$$k = \frac{\pi}{L} n, \quad n = 1, 2, 3, \dots \quad (33)$$

This tells us which frequencies can be produced

$$\boxed{\omega_n = vk_n = v \frac{\pi}{L} n, \quad n = 1, 2, 3, \dots} \quad \text{both ends fixed} \quad (34)$$

This is the spectrum for 2 Dirichlet boundary conditions.

Next, consider having the end at  $x = 0$  fixed but the end at  $x = L$  free. To figure out what happens if the end is free, we have to go back to our picture



Now there is no line connected  $A(x)$  to  $A(x + \Delta x)$ . Then,  $F_{\text{right}} = 0$  and Eq. (6) becomes

$$F_{\text{total}} = F_{\text{left}} + F_{\text{right}} = T \frac{\partial A}{\partial x} \quad (35)$$

So, Eq. (8) becomes

$$\frac{T}{\mu} \frac{\partial A}{\partial x} = \Delta x \frac{\partial^2 A}{\partial t^2} \quad (36)$$



In this case if we take  $\Delta x$  to 0 we see that  $\frac{\partial A}{\partial x} \rightarrow 0$ . Thus a free end must satisfy

$$\boxed{\frac{\partial A(L, t)}{\partial x} = 0} \quad (37)$$

This is known as a free, open, or **Neumann** boundary condition.

Now using the  $x=0$  fixed solution, Eq. (32), the Neumann condition at  $x=L$  implies

$$0 = \frac{\partial A(L, t)}{\partial x} = k A_0 \cos(kL) \sin(\omega t + \phi) \quad (38)$$

For this to hold at all times,  $\cos(kL)$  must be at a zero of the cosine curve. Now,  $\cos(x) = 0$  when  $x = (n + \frac{1}{2})\pi$ . Thus,

$$\boxed{\omega_n = v \frac{n + \frac{1}{2}}{L} \pi, \quad n = 0, 1, 2, 3}, \quad \text{one fixed end, one free end} \quad (39)$$

This solution says that the lowest frequency is

$$\nu_0 = \frac{\omega_1}{2\pi} = \frac{1}{2} v \frac{\frac{1}{2}}{L} = \frac{1}{4} \frac{v}{L} \quad (40)$$

the next frequency up is

$$\nu_1 = \frac{1}{2} v \frac{1 + \frac{1}{2}}{L} = \frac{3}{4} \frac{v}{L} = 3\nu_0 \quad (41)$$

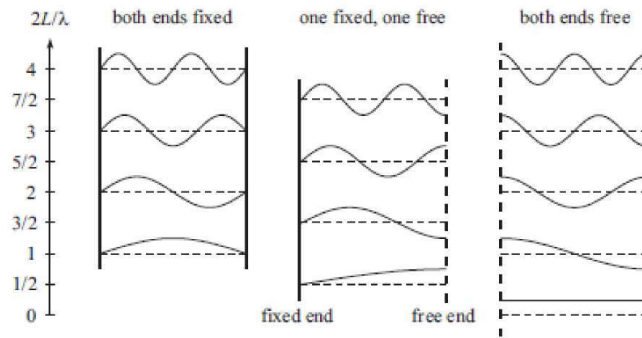
and so on. Thus the even harmonics are missing!! This has dramatic consequences for instruments like the trumpet and the clarinet.

Finally, if  $x=0$  is free, we must have  $A(x, t) = A_0 \cos(kx) \sin(\omega t + \phi)$ . Then, if  $x=L$  is also free, we find  $\sin(kL) = 0$  which implies

$$\boxed{\omega_n = v \frac{n}{L} \pi, \quad n = 0, 1, 2, 3}, \quad \text{both ends free} \quad (42)$$

The only difference between both free ends solution and both fixed end solution is that for free ends  $n=0$  is allowed. However, the  $n=0$  solution is  $A(x, t) = \text{const}$  which has  $k = \omega = 0$  thus it is not physically interesting.

Here are the lowest harmonics with different boundary conditions



**Figure 2.** Frequencies allowed for different boundary conditions

If the fundamental note (lowest frequency) is  $\nu$ , then we find

	lowest	next	second	third
both fixed	$\nu$	$2\nu$	$3\nu$	$4\nu$
one fixed, one free	$\nu$	$3\nu$	$5\nu$	$7\nu$
both free	$\nu$	$2\nu$	$3\nu$	$4\nu$

(43)

What are the implications for this? Well for one thing different instruments correspond to different boundary conditions. String instruments have both ends fixed. Woodwinds and brass have one end open. A flute has both ends free. The absence of the even harmonics is one of the reasons that clarinets tend to sound eerie. Many of the complications in the designs of brass instruments help restore the even harmonics. This is explained well in Rick Heller's book. See for example this figure from page 317: (don't worry too much about understanding this picture now – it will make more sense after the next lecture, on music, but it naturally falls in to this lecture):

### Strategy for resonance placement in modern brass instruments

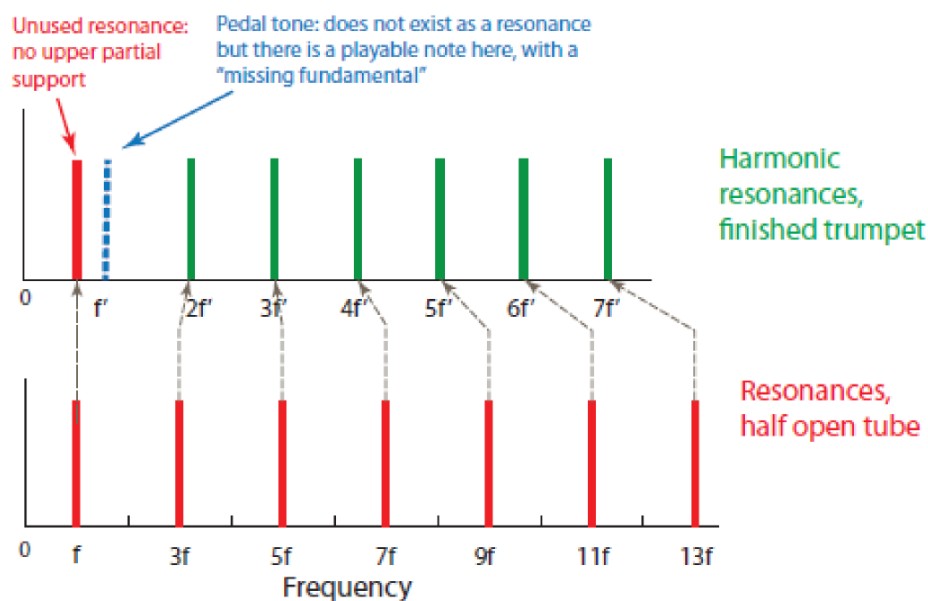
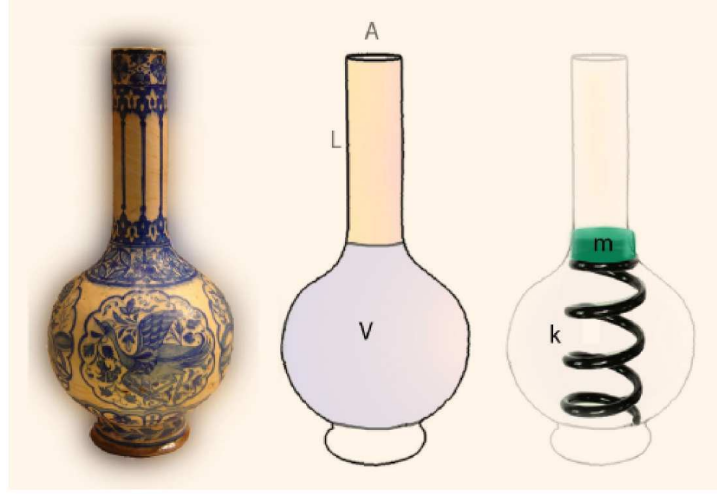


Figure 16.6: The shifting of the even-numbered resonances of a half open tube (bottom, red), based on a fundamental  $f$ , to a set of equally spaced odd and even harmonic resonances based on a fundamental “pedal tone”  $f'$ , a frequency which does not exist as a resonance. This is accomplished with a bell and mouthpiece. The pedal tone is playable, due to the harmonic support it gets, but is normally not used. The lowest, quarter-wave resonance (red) is not part of the harmonic series (green) and is also unused.

## 5 Helmholtz resonators (optional)

An important object in the physics of sound is the resonance chamber or **Helmholtz resonator**. Helmholtz resonators are hollow cavities with small openings, like a bottle or a violin body. They work because the volume of the air in the body cannot change, thus pushing down on the air in the neck forces the air in the body to push back with essentially a linear restoring force, like a spring. The setup is like this



**Figure 3.** In a Helmholtz resonator, the air in the neck acts like a mass and the air in the base acts like a spring. Figure taken from Fig 13.1 of Heller.

To work out the resonant frequency for a Helmholtz resonator we can use  $\omega = \sqrt{\frac{k}{m}}$ . We can extract the spring constant  $k$  from  $F = -k\Delta x$ . For pressure,  $F = A \cdot dp$ , where  $A$  is the area, in this case the cross sectional area of the neck. Now,  $\rho = \frac{m}{V}$  so

$$d\rho = \frac{d}{dV} \left( \frac{m}{V} \right) dV = -\frac{m}{V^2} dV = -\rho \frac{dV}{V} \quad (44)$$

Also using Eq. (16),  $\frac{dp}{d\rho} = \gamma \frac{p}{\rho}$  for sound waves, we have

$$dp = \frac{dp}{d\rho} d\rho = \gamma \frac{p}{\rho} d\rho = -\gamma \frac{p}{V} dV \quad (45)$$

Now,  $dV = A \Delta x$  and so

$$F = A \cdot dp = -\gamma A^2 \frac{p}{V} \Delta x \quad (46)$$

Thus

$$k = \gamma A^2 \frac{p}{V} = A^2 c_s^2 \frac{\rho}{V} \quad (47)$$

The mass on which the spring acts is the air in the neck. It has mass  $m = \rho A L$ , thus

$$\omega = \sqrt{\frac{k}{m}} = \sqrt{\frac{A^2 c_s^2 \rho / V}{\rho A L}} = c_s \sqrt{\frac{A}{V L}} \quad (48)$$

Thus Helmholtz resonators resonate at a single frequency

$$\boxed{\nu = \frac{c_s}{2\pi} \sqrt{\frac{A}{V L}}} \quad (49)$$

where  $A$  is the area of the opening,  $L$  is the length of the neck, and  $V$  is the volume of the cavity.

For example, consider a 10 cm wide jar with a 10 cm long neck. Using  $v = 343 \frac{m}{s}$ ,  $A = 1\text{cm}^2$ ,  $L = 1\text{cm}$ ,  $V = 1\text{L} = 1000\text{cm}^3$ , we find  $\nu = 172\text{Hz}$ . The associated wavelength in air is  $\lambda = \frac{c_s}{\nu} = 2m$ . Note that the wavelength of sound in the resonator is much larger than the size of the resonator.

Since Helmholtz resonators have only one frequency, they have no harmonics (no overtones). However, they can have low  $Q$  values. Indeed, if you blow on a bottle, you see that the sound does not resonate for long at all. This is good, if you are building an instrument, since you want all the audible frequencies to resonate. On a violin, the vibrations are produced on the strings, transmitted to the wooden body of the violin through the bridge (the part of the violin which connects the strings to the body). The body then vibrates, exciting the air in the body which emits sound through the holes. I can't describe the function of the body of a violin better than Heller. Here's his description [Heller, p. 267]

Helmholtz resonators can be used as transducers, turning mechanical energy into sound energy. A prime example is the violin. The violin body is basically a box containing air, with the f-holes opening to the outside. It functions deliberately as a Helmholtz resonator, enhancing the low frequency response of the violin, giving it much of its richness of tone...

The violin body's broad Helmholtz resonance peaks around 300 Hz. No doubt the shape thin but long holes serve to increase air friction and thus lower the  $Q$  of the Helmholtz mode, spreading the resonance over a broader frequency range. This props up the transduction of string vibrations into sound down to the frequency of the open  $D$  string [ $\nu \sim 300\text{Hz}$ ].

# Lecture 7: Music

## 1 Why do notes sound good?

In the previous lecture, we saw that if you pluck a string, it will excite various frequencies. The amplitude of each frequency which is excited will be proportional to the coefficient in the Fourier decomposition. In this lecture we will start to understand how different frequencies combine to produce music. This lecture is best studied alongside the Mathematica notebook `music.nb` on the isite.

In the first section “playing notes” of the notebook, you can listen to a pure frequency (C4 = middle C = 261 Hz). It sounds pleasant, but not particularly interesting. Now play the “square wave” version of middle C – you should notice that it sounds somewhat tinny and unpleasant. These are the same notes, but different sounds. Why do they sound different? One way to understand the difference is to compare the Fourier decomposition of the sine wave and the square wave:

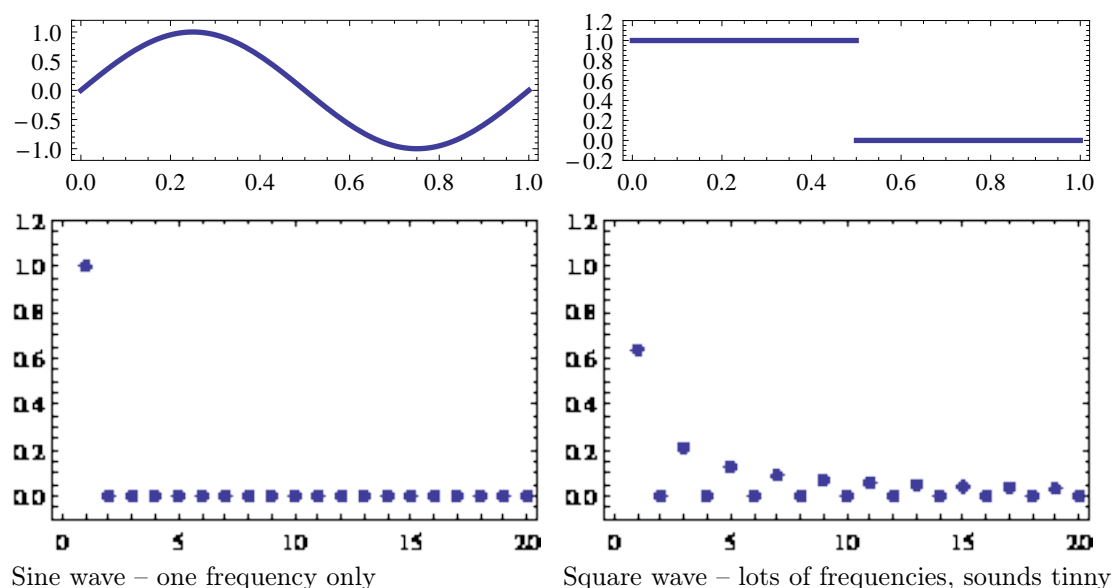


Figure 1. Comparison of Fourier modes for sin and square wave

For both the sine wave and the square wave, the dominant mode is the fundamental. However, the square wave has lots of other modes which make the note sound less pure. One way to understand why the square wave sounds worse is that it has many high frequency notes with significant amplitude. It is hard for our brains to process all these high frequency notes, so we find it jarring. In fact, if all frequencies are present at once we get so-called **white noise**. White noise is perhaps as unmusical as you can get.

Now consider playing two notes at once. In the “playing pairs of notes” section, play the 250 Hz and 270 Hz notes at the same time. It doesn’t sound great. Why?

The problem is that you hear a rattling around 20 times a second. This 20 Hz rattling is the beat frequency between the 250 Hz and 270 Hz. Indeed,

$$\cos(270 \text{ Hz } 2\pi t) + \cos(250 \text{ Hz } (2\pi t)) = 2 \cos(10 \text{ Hz } 2\pi t) \cos(260 \text{ Hz } 2\pi t) \quad (1)$$

This the sum of those two notes oscillates at 10 Hz and at 260 Hz. The 10 Hz oscillation (which we hear as a 20 Hz beat frequency) is jarring – your mind tries to process it consciously. Frequencies as high as 260 Hz do not have this effect.

Thus there seem to be two reasons sounds appear unmusical:

- Too many frequencies are present at once
- Beating occurs at frequencies we can consciously process.

This is a physics class, not a biology class, so we will not try to explain why these facts hold. We merely observe that whenever each criteria is satisfied, sounds appear unmusical. However, now that we have defined the problem, we can start to study music scientifically.

## 2 Dissonant and consonant note pairs

Now we're ready to study music. If we had a pure sine wave at 300 Hz, then it sounds nasty when played at the same time as a sine wave of 320 Hz. However, if we play it with a sine wave of 580 Hz it does not sound so bad. That is because

$$\cos(300 \text{ Hz } 2\pi t) + \cos(580 \text{ Hz } 2\pi t) = 2 \cos(140 \text{ Hz } 2\pi t) \cos(440 \text{ Hz } 2\pi t) \quad (2)$$

The beat frequency  $2 \times 140 \text{ Hz} = 280 \text{ Hz}$  is not low enough to be harsh – it is just a note (try the Mathematica notebook). On the other hand, if we played 300 Hz and 580 Hz on an actual instrument it would sound horrible.

We can see why from studying our plucked string example. Recall that for a string plucked near the end, the relative Fourier coefficients scale like  $\frac{1}{n}$ . So the dominant frequency  $n = 1$  (the fundamental) has only twice the amplitude of the first harmonic ( $n = 2$ ). Thus playing 580 Hz along side a plucked string would give

$$f(t) = \cos(580 \text{ Hz } 2\pi t) + \sum_{n=1}^{\infty} \frac{1}{n} \cos(300n \text{ Hz } 2\pi t) \quad (3)$$

Writing  $T = \text{Hz } 2\pi t$  to clean up the equation and expanding the sum

$$f(t) = \cos(580 T) + \cos(300 T) + \frac{1}{2} \cos(600 T) + \dots \quad (4)$$

Let us combine the 580 Hz oscillation with the 600 Hz oscillation using trig sum rules. Using

$$\cos(580 T) + \frac{1}{2} \cos(600 T) = \frac{3}{2} \cos(10 T) \cos(590 T) + \frac{1}{2} \sin(10 T) \sin(590 T) \quad (5)$$

we find

$$f(t) = \cos(300 T) + \frac{3}{2} \cos(10 T) \cos(590 T) + \frac{1}{2} \sin(10 T) \sin(590 T) + \sum_{n=3}^{\infty} \frac{1}{n} \cos(300n T) \quad (6)$$

Now we see beating at  $2 \times 10 \text{ Hz} = 20 \text{ Hz}$ , which is audible and jarring. There is beating between the 580 Hz note and the first harmonic of the plucked string. The point is that with pure sine waves, no harmonics are excited, but with real instruments they are.

In general, instruments will have significant amplitudes for many harmonics. Here is the Fourier spectrum of a flute

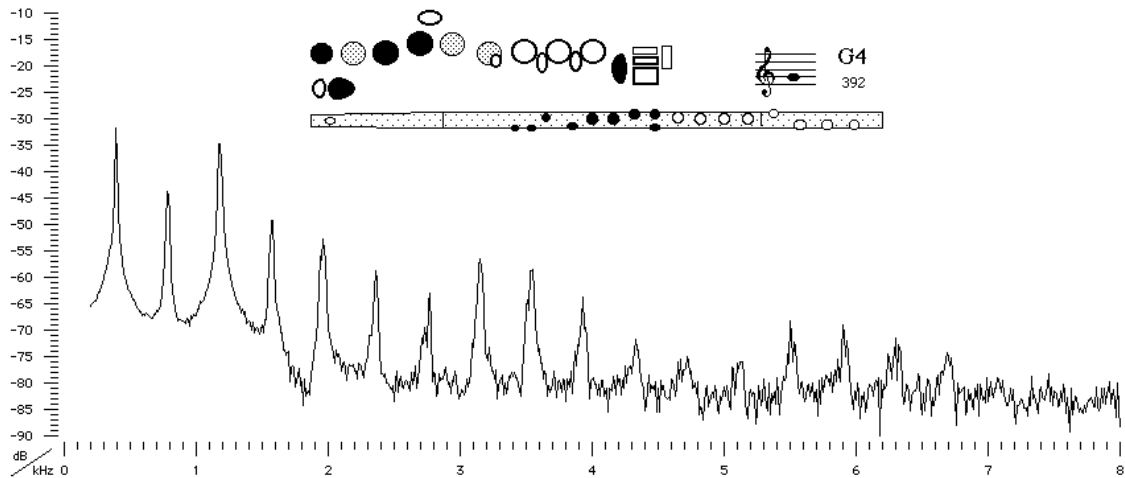


Figure 2. Spectrum of the note G4 on a flute.

This is a flute playing the note G4 which is 392 Hz. You see that not only is the fundamental frequency (the G4) largely excited, but many other modes have significant amplitudes. These modes are all the higher harmonics. The harmonics determine what an instrument sounds like – its **timbre**. Timbre is what a note sounds like when played, as distinguished from its **pitch**, which is the fundamental frequency played, and **intensity**, which is the power going into sound from the instrument. We will study the timbres of different instruments in lecture, in section and on psets. The pitch of a note is the frequency with highest intensity, which is usually the lowest frequency where there is a peak. If the first peak is not the highest, what the pitch is can be somewhat subjective. The relative heights of the different peaks and their  $Q$  values is the timbre, and the absolute scale is the intensity. All of this can be read off the Fourier spectrum, as in Figure 2 for a flute. We'll discuss timbre more later, but for now, the main point is that most instruments will have integer multiples of the fundamental frequency excited with significant amplitudes; it is these harmonics which are the key to the scale in Western music.

Key points are:

- On a real instrument, there will be unmusical beating whenever an integer multiple of one harmonic is close but not equal to an integer multiple of another harmonic.
- Conversely, the most consonant notes will have some harmonics which exactly agree.

For example, let's start with middle  $C$ . This note is called  $C4$  (the  $C$  is the note and 4 is the octave) and has a frequency of  $\nu_0 = 261$  Hz. Which notes sound good along side  $C4$ ? Well, the 261 Hz note has harmonics of  $\nu_0, 2\nu_0, 3\nu_0, 4\nu_0$ , etc.:

$$261 \text{ Hz}, \quad 522 \text{ Hz}, \quad 783 \text{ Hz}, \quad 1044 \text{ Hz}, \quad 1305 \text{ Hz}, \dots \quad (7)$$

Thus if we play any of those notes along with  $C4$  it will sound harmonic. If the fundamental frequency is  $\nu_0$  then

$$2\nu_0 = 1 \text{ octave} = 1\text{st harmonic} = C5 \quad (8)$$

Are there more notes which are harmonious? Yes. Consider the note with  $\nu_5 = 391\text{Hz}$ . This note has  $2\nu_5 = 783\text{Hz}$ . Thus the second harmonic of  $\nu_5$  matches the 3rd harmonic of  $\nu_0$ . We call the note  $\nu_5 = \frac{3}{2}\nu_0$  the **perfect fifth**

$$\nu_5 = \frac{3}{2}\nu_0 = \text{perfect fifth} = G4 \quad (9)$$

This is the *G* above middle *C*. In the same way, consider  $\nu_4 = 348\text{Hz}$ . The 3rd harmonic of  $\nu_4$  agrees with the 4th harmonic of  $\nu_0$ . We call this the perfect fourth

$$\nu_4 = \frac{4}{3}\nu_0 = \text{perfect fourth} = F4 \quad (10)$$

And so on.

It is easy to see that any rational number ratio of frequencies will be consonant. Many of these ratios have names

$\frac{\nu}{\nu_0}$	1	2	$\frac{3}{2}$	$\frac{4}{3}$	$\frac{5}{4}$	$\frac{6}{5}$	$\frac{5}{3}$	$\frac{8}{5}$
name	fundamental	octave	perfect fifth	perfect fourth	Major third	Minor third	Major sixth	Minor sixth
example	C4	C5	G4	F4	E4	E♭	A5	A♭5

**Table 1.** Notes names and ratios in the **just intonation scale**

There are an infinite number of rational numbers. So where do we stop? The answer is that the lower the numbers in the ratio (that is, the 3 and 2 in  $\frac{3}{2}$  are lower than the 8 and 5 in  $\frac{8}{5}$ ), the more consonant they will be. That's because for something like  $\frac{11}{17}$ , one would need the 17th harmonic of one note to match the 11th harmonic of another note. By such high harmonics, the amplitudes are no longer large, and the spectrum is messy (as you can see in Figure 2). Also, it is more likely for frequencies with a ratio of  $\frac{11}{17}$  to give harmonics which are close but not equal, generating dissonant beating, before generating the harmonic consonance. Thus, for numbers large than about 5 in the ratio, notes are no longer appreciated as harmonic.

### 3 Scales

If we have a given note, say *C4*, we can define all the other notes so that they will be harmonic with *C4*.

#### 3.1 Just intonation scale

The most harmonic notes will have the smallest integers in the ratio, as in Table 1. This is a particular choice of tuning known as the **just intonation scale**. The just intonation scale is in a sense the most harmonic choice for the frequencies of notes in a scale (it is default tuning for some non-Western instruments, such as the Turkish Baglama). But it is just a choice.

Note that if we pick a set of notes that sound harmonic with *C4*, the same set of notes will generically not sound harmonic with another note, like *D4*. Thus if we're playing a song, the set of notes we want to use is determined relative to some starting note. This starting note is called the **key**. For example, if you are in the key of C, the notes C, G and F sound good. But if you are in the key of D, the notes C, G and F will generally not sound as good.



On some instruments, such as a violin which has no frets, there are no predefined notes. Thus on a violin, if you work in the key of  $C$ , you can play all the harmonics in exactly the right place (if you have a skilled enough ear and hand). Thus you can play the just intonation scale in any key. There are an *infinite* number of notes you can play on a violin. The same is true on most instruments actually. For example, even though an oboe has fixed holes corresponding to notes, oboe players can easily move the notes up or down by manipulating the reed. Controlling the precise frequency of a note with your mouth is critical to playing any woodwind instrument well.

On other instruments, like a piano or a stringed instrument with frets like a guitar, the notes are essentially built in to the instrument. You can sometimes tweak the notes if you are skilled, or tune the instrument to a different key, but there are a finite number of notes which can be played in a given tuning. Unfortunately, it's impossible to have an instrument with a finite number of notes be capable of playing the most harmonic notes in every key.

To see the problem, suppose you want your piano to be in the just intonation scale in the key of  $C$ . That means that you want all the other notes to be defined so that they are related to  $C$  by rational numbers with low integers. For example, the whole notes can be defined as

note	C	D	E	F	G	A	B	C
$\frac{\nu}{\nu_0}$	1	$\frac{9}{8}$	$\frac{5}{4}$	$\frac{4}{3}$	$\frac{3}{2}$	$\frac{5}{3}$	$\frac{15}{8}$	2
decimal	1	1.125	1.25	1.333	1.5	1.666	1.875	2

**Table 2.** Notes names and ratios in **just intonation**

This defines the note  $D$  as having the frequency  $\frac{9}{8}$  times the frequency of the  $C$ . Now, where is the 5th of  $D$ ? This should be at  $\frac{3}{2} \times \frac{9}{8} = \frac{27}{16}$  times  $\nu_0$ . This note is somewhere between the  $A$  and the  $B$ , but it is not exactly a note in the key of  $C$ . It is not hard to see that to get an instrument which could play any note in any key, you would need an enormous number of available notes. Please make sure you understand this point, as it is key to understanding scales.

So what can we do? There are two options: the first is you can tune your instrument to the key you want to play in. Stringed instruments can do this. But it is not so appealing of an option if we want to play music in different keys without retuning every time. The other option is to compromise. Most of our ears are not sensitive enough to distinguish close but slightly different notes. Thus we can choose scales which are not exactly correct in any key, but close to correct in all keys.

### 3.2 Pythagorean scale

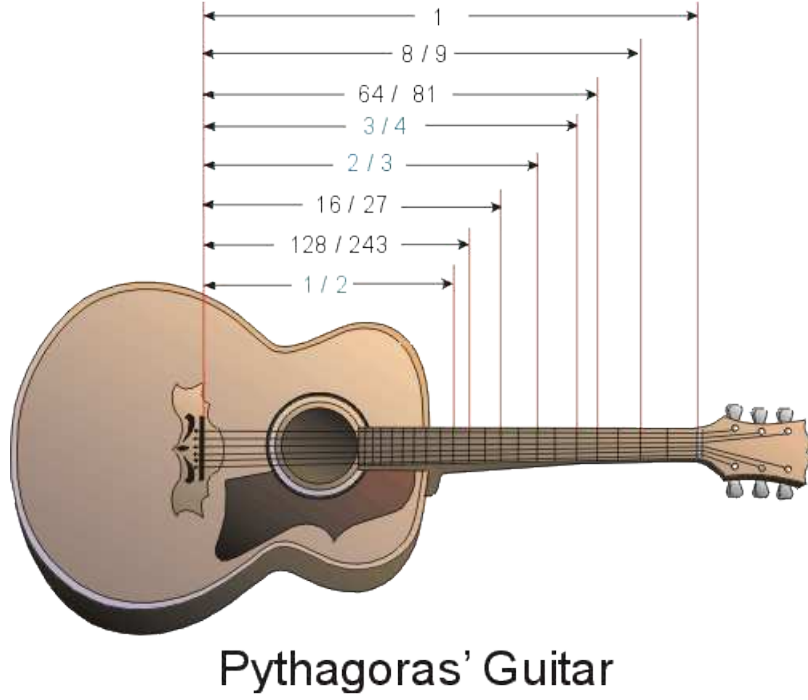
One way to approximate the scale is by choosing the notes to be related by powers of  $\frac{3}{2}$  and octaves to be related by factors of 2. For example,  $\frac{3}{2}\nu_0$  is the perfect fifth ( $G4$  when  $C4=\nu_0$ ). Then  $(\frac{3}{2})^2\nu_0$  which is the fifth of the fifth, or the fifth of  $G4$  which is  $D5$ . The next note has  $(\frac{3}{2})^3\nu_0$  or the fifth of  $D5$  which is  $A6$  and so on. We can bring any power of  $\frac{3}{2}$  back to the interval between 0 and 1 by dividing by 2 to some power. For example, since  $D5 = \frac{9}{4}\nu_0$  then  $D4 = \frac{9}{8}\nu_0$ . We then get

note	C	D	E	F	G	A	B	C
$\frac{\nu}{\nu_0}$	1	$\frac{9}{8}$	$\frac{81}{64}$	$\frac{4}{3}$	$\frac{3}{2}$	$\frac{27}{16}$	$\frac{243}{128}$	2
decimal	1	1.125	1.266	1.333	1.5	1.688	1.898	2

**Table 3.** Notes names and ratios in Pythagorean tuning

This is called the **Pythagorean tuning**. Note that the octave, perfect fifth ( $G$ ) and perfect fourth ( $F$ ) agree with their values in Table 2. Some notes do not agree: for example,  $E$  is  $\frac{81}{64}\nu_0 = 1.266\nu_0$  in this tuning. This ratio is close to  $\frac{5}{4}$  but not exactly. Thus if we play  $C$  and  $E$  it will be close to a consonant sounding note, but not exactly.

The advantage of the Pythagorean tuning is that the perfect fifth and perfect fourth of every note is included in the scale.



**Figure 3.** Pythagorean guitar: frets positions are related by powers of  $\left(\frac{3}{2}\right)^n 2^m$

### 3.3 Equal-tempered scale

An interesting feature of the Pythagorean scale is that the 12th fifth is very close to 8 octaves:

$$\left(\frac{3}{2}\right)^{12} = 129.748 \approx 128 = 2^7 \quad (11)$$

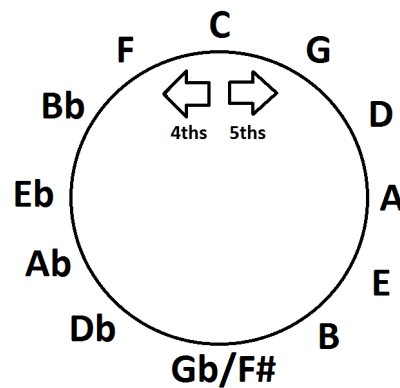
Thus another compromise is to say that when we go 12 steps around the circle of fifths, we get back to the note we started at. Then we can devise a scale which does not choose any key. We simply relate the notes by powers of  $2^{1/12}$ . Each half-step gives another factor of  $2^{1/12}$ . The result are the frequencies in Table 4.

note	C	C $\sharp$	D	D $\sharp$	E	F	F $\sharp$	G	G $\sharp$	A	A $\sharp$	B	C
$\frac{\nu}{\nu_0}$	1	$2^{\frac{1}{12}}$	$2^{\frac{2}{12}}$	$2^{\frac{3}{12}}$	$2^{\frac{4}{12}}$	$2^{\frac{5}{12}}$	$2^{\frac{6}{12}}$	$2^{\frac{7}{12}}$	$2^{\frac{8}{12}}$	$2^{\frac{9}{12}}$	$2^{\frac{10}{12}}$	$2^{\frac{11}{12}}$	2
decimal	1	1.059	1.122	1.189	1.260	1.335	1.414	1.498	1.587	1.682	1.782	1.888	2

**Table 4.** Notes names and ratios in the **equal-tempered scale**.

This tuning is called the **equal-tempered scale**. The equal-tempered scale is the standard tuning for all of Western music.

Because of Eq. (11) if you keep going up by perfect fifths, and normalizing by octaves, you will get back to the note you started at. We can see this in the **circle of fifths**



**Figure 4.** Circle of fifths. Each note going clockwise is a perfect fifth above the previous note. Going counterclockwise, each note is a perfect 4th above the previous note. The circle only closes in the equal-tempered scale.

In this circle, each note is 1 fifth above the note clockwise. So  $G$  is a fifth above  $C$ ,  $D$  a fifth above  $G$  and so on. Going counterclockwise, the intervals are fourths:  $C$  is a fourth above  $G$  and  $F$  is a fourth above  $C$ . Going up a fifth is the same as going down a fourth and adding an octave,

If the notes are defined with the Pythagorean scale, the circle doesn't close: going up by 12 steps, and normalizing back down to the original octave leaves you  $(\frac{3}{2})^{12}2^{-8} = 1.014$  times where you started. Thus the circle doesn't close by 1.4%. In the equal-tempered scale, it *does* exactly close, however, none of the notes have frequency ratios of exactly  $\frac{2}{3}$ .

### 3.4 Summary

We discussed 3 scales. The just intonation scale chooses notes to be related by rational number ratios with integers as small as possible in the numerator and denominator. The Pythagorean scale has all notes related by  $3^n2^m$  for some  $m$  and  $n$ . Both just intonation and the Pythagorean scale require a key to start in. The third scale is the equal-tempered scale. Notes in the equal-tempered scale are related by  $2^{\frac{n}{12}}$  for some  $n$ .

Here is a comparison between the relative frequencies of the 3 scales in the key of  $C$ :

note	C	D	E	F	G	A	B	C
just-intonation	1	$\frac{9}{8}$	$\frac{5}{4}$	$\frac{4}{3}$	$\frac{3}{2}$	$\frac{5}{3}$	$\frac{15}{8}$	2
Pythagorean	1	$\frac{9}{8}$	$\frac{81}{64}$	$\frac{4}{3}$	$\frac{3}{2}$	$\frac{27}{16}$	$\frac{243}{128}$	2
equal-tempered	1	$2^{\frac{2}{12}}$	$2^{\frac{4}{12}}$	$2^{\frac{5}{12}}$	$2^{\frac{7}{12}}$	$2^{\frac{9}{12}}$	$2^{\frac{11}{12}}$	2

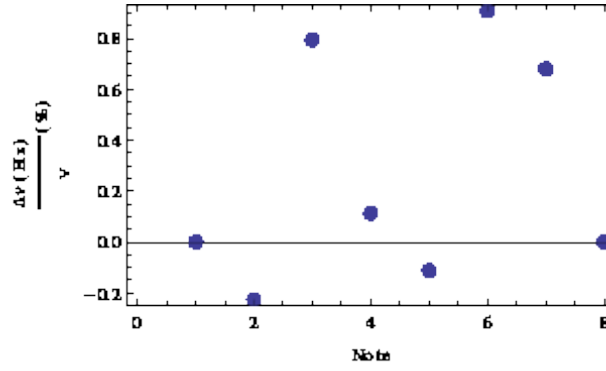
**Table 5.** Comparison of scales: exact ratios.

In decimals

note	C	D	E	F	G	A	B	C
just-intonation	1	1.125	1.25	1.333	1.5	1.666	1.875	2
Pythagorean	1	1.125	1.266	1.333	1.5	1.688	1.898	2
equal-tempered	1	1.122	1.260	1.335	1.498	1.682	1.888	2

**Table 6.** Comparison of scales: decimal approximations.

Here is a graphical comparison of how far off the frequency is in the equal-tempered scale from the frequency in the just intonation scale



**Figure 5.** Difference between the equal-tempered frequencies  $\nu_{\text{WT}}$  and the just intonation frequencies for whole notes  $C, D, E, \dots, C$  labeled as 1 to 8.

One thing we can see is that the perfect 4th and perfect 5th are very close to their optimal values, while the 6th and 7ths are not so close.

# Lecture 8: Fourier transforms

## 1 Strings

To understand sound, we need to know more than just which notes are played – we need the shape of the notes. If a string were a pure infinitely thin oscillator, with no damping, it would produce pure notes. In the real world, strings have finite width and radius, we pluck or bow them in funny ways, the vibrations are transmitted to sound waves in the air through the body of the instrument etc. All this combines to a much more interesting picture than pure frequencies. For example, the spectrum of a violin looks like this:

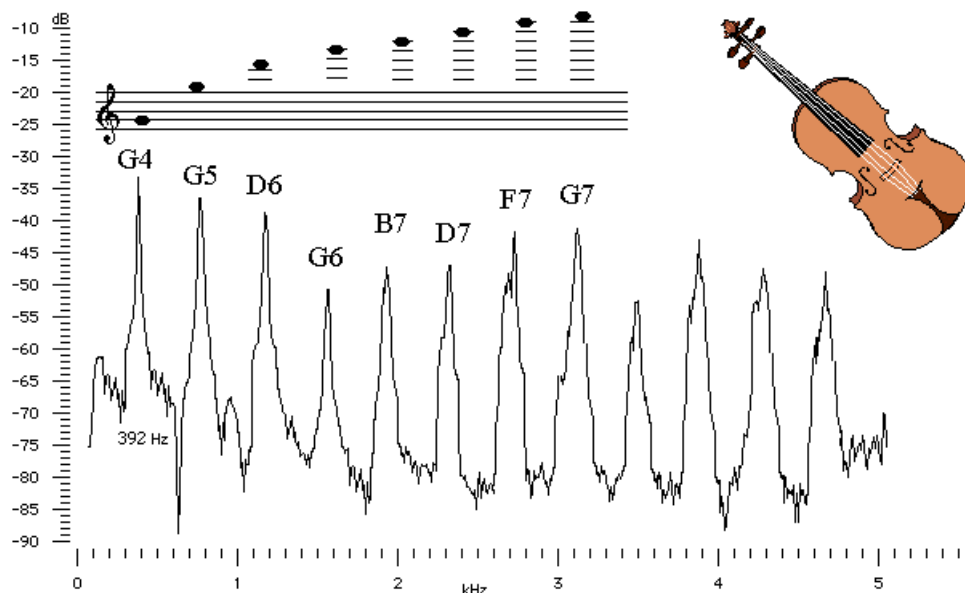


Figure 1. Spectrum of a violin

This figure shows the intensity of each frequency produced by the violin (the vertical axis is in decibels, which is a logarithmic measure of sound intensity; we'll discuss this scale in Lecture 10). We know the basics of this spectrum: the fundamental and the harmonics are related to the Fourier series of the note played. Now we want to understand where the shape of the peaks comes from. The tool for studying these things is the Fourier transform.

## 2 Fourier transforms

In the violin spectrum above, you can see that the violin produces sound waves with frequencies which are arbitrarily close. The way to describe these frequencies is with **Fourier transforms**.

Recall the Fourier exponential series

$$f(x) = \sum_{n=-\infty}^{\infty} c_n e^{i \frac{2\pi n x}{L}} \quad (1)$$

where

$$c_n = \frac{1}{L} \int_{-\frac{L}{2}}^{\frac{L}{2}} dx f(x) e^{-i \frac{2\pi n x}{L}} \quad (2)$$

To check this, we plug Eq. (1) into Eq. (2) giving

$$c_n = \frac{1}{L} \int_{-\frac{L}{2}}^{\frac{L}{2}} dx \left[ \sum_{m=-\infty}^{\infty} c_m e^{i \frac{2\pi m x}{L}} \right] e^{-i \frac{2\pi n x}{L}} = \frac{1}{L} \sum_{m=-\infty}^{\infty} c_m \int_{-\frac{L}{2}}^{\frac{L}{2}} dx e^{i \frac{2\pi m x}{L}} e^{-i \frac{2\pi n x}{L}} \quad (3)$$

Then using the mathematical identity

$$\int_{-\frac{L}{2}}^{\frac{L}{2}} dx e^{i(m-n)x \frac{2\pi}{L}} = L \delta_{mn} \quad (4)$$

we get

$$c_n = \frac{1}{L} \sum_{m=-\infty}^{\infty} c_m L \delta_{nm} = c_n \quad (5)$$

as desired. That is, we have checked Eq. (2).

To derive the Fourier transform, we write

$$k_n = \frac{2\pi n}{L} \quad (6)$$

where  $n$  is still an integer going from  $-\infty$  to  $+\infty$ . For arbitrary  $L$ ,  $k_n$  can get arbitrarily big in the positive or negative direction. However, at fixed  $L$ , the lowest non-zero  $k_n$  cannot be arbitrarily small:  $|k_n| > \frac{2\pi}{L}$ . Then, we define

$$\tilde{f}(k_n) = \frac{L c_n}{2\pi} = \frac{1}{2\pi} \int_{-\frac{L}{2}}^{\frac{L}{2}} dx f(x) e^{-i k_n x} \quad (7)$$

The factor of  $2\pi$  in this equation is just a convention. Now we can take  $L \rightarrow \infty$  so that  $k_n$  can get arbitrarily close to zero. This gives

$$\tilde{f}(k) = \frac{1}{2\pi} \int_{-\infty}^{\infty} dx f(x) e^{-i k x} \quad (8)$$

where now  $k$  can be any real number. This is the Fourier transform. It is a continuum generalization of the  $c_n$ 's of the Fourier series.

The inverse of this comes from writing Eq. (1) as an integral. From Eq. (6), we find  $dk_n = \frac{2\pi}{L} \Delta n$ . This leads to

$$f(x) = \sum_{n=-\infty}^{\infty} c_n e^{i k_n x} \Delta n = \sum_{n=-\infty}^{\infty} c_n e^{i k_n x} \frac{L}{2\pi} dk_n = \int_{-\infty}^{\infty} dk \tilde{f}(k) e^{i k x} \quad (9)$$

where we have used Eq. (7) and taken  $L \rightarrow \infty$  in the last step.

So we have

$$\boxed{\tilde{f}(k) = \frac{1}{2\pi} \int_{-\infty}^{\infty} dx f(x) e^{-ikx} \iff f(x) = \int_{-\infty}^{\infty} dk \tilde{f}(k) e^{ikx}} \quad (10)$$

We say that  $\tilde{f}(k)$  is the **Fourier transform** of  $f(x)$ . The factor of  $2\pi$  is just a convention. We could also have defined  $f(x)$  with the  $2\pi$  in it. The sign on the phase is also a convention (that is, we could have defined  $\tilde{f}(k) = \frac{1}{2\pi} \int_{-\infty}^{\infty} dx f(x) e^{ikx}$  instead). Keep in mind that different conventions are used in different places and by different people. There is no universal convention for the  $2\pi$  factors. All conventions lead to the same physics.

The Fourier transform of a function of  $x$  gives a function of  $k$ , where  $k$  is the **wavenumber**. The Fourier transform of a function of  $t$  gives a function of  $\omega$  where  $\omega$  is the angular frequency:

$$\tilde{f}(\omega) = \frac{1}{2\pi} \int_{-\infty}^{\infty} dt f(t) e^{-i\omega t} \quad (11)$$

### 3 Example

As an example, let us compute the Fourier transform of the position of an underdamped oscillator:

$$f(t) = e^{-\gamma t} \cos(\omega_0 t) \theta(t) \quad (12)$$

where the **unit-step function** is defined by

$$\theta(t) = \begin{cases} 1, & t > 0 \\ 0, & t \leq 0 \end{cases} \quad (13)$$

This function insures that our oscillator starts at time  $t = 0$ . If didn't include, the amplitude would blow up as  $t \rightarrow -\infty$ .

We first write

$$f(t) = e^{-\gamma t} \cos(\omega_0 t) \theta(t) = \frac{1}{2} e^{-\gamma t} e^{i\omega_0 t} \theta(t) + \frac{1}{2} e^{-\gamma t} e^{-i\omega_0 t} \theta(t) \quad (14)$$

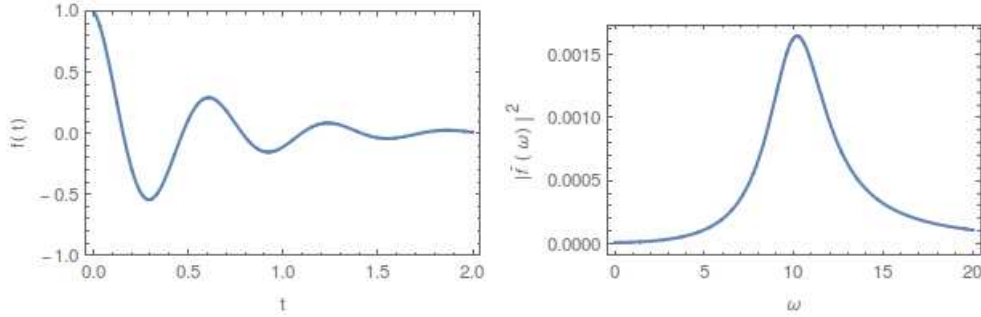
So we can Fourier transform the simpler exponential function. Starting with the first term, we find

$$\begin{aligned} \tilde{f}_{+\omega_0}(\omega) &= \frac{1}{4\pi} \int_{-\infty}^{\infty} dt e^{-\gamma t} e^{-i(\omega - \omega_0)t} \theta(t) \\ &= \frac{1}{4\pi} \int_0^{\infty} dt e^{(-\gamma - i\omega + i\omega_0)t} \\ &= \frac{1}{4\pi} \frac{1}{-\gamma - i(\omega - \omega_0)} e^{(-\gamma - i\omega + i\omega_0)t} \Big|_0^{\infty} \\ &= \frac{1}{4\pi} \frac{1}{\gamma + i(\omega - \omega_0)} \end{aligned}$$

In the last step we have used that the  $t = \infty$  endpoint vanishes due to the  $e^{-\gamma t}$  factor and that at the  $t = 0$  endpoint the exponential is 1. The second term in Eq. (14) is the first term with  $\omega_0 \rightarrow -\omega_0$ . Thus the full Fourier transform is

$$\tilde{f}(\omega) = \frac{1}{4\pi} \left[ \frac{1}{\gamma + i(\omega - \omega_0)} + \frac{1}{\gamma + i(\omega + \omega_0)} \right] = \frac{1}{2\pi i} \frac{\omega - i\gamma}{(\omega - i\gamma)^2 - \omega_0^2} \quad (15)$$

As mentioned before, the spectrum plotted for an audio signal is usually  $|\tilde{f}(\omega)|^2$ . Let's see what this looks like. We'll take  $\omega_0 = 10$  and  $\gamma = 2$ . The function and the modulus squared  $|\tilde{f}(\omega)|^2$  of its Fourier transform are then:



**Figure 2.** An underdamped oscillator and its power spectrum (modulus of its Fourier transform squared) for  $\gamma = 2$  and  $\omega_0 = 10$ .

We now can also understand what the shapes of the peaks are in the violin spectrum in Fig. 1. The widths of the peaks give how much each harmonic damps with time. The width at half maximum gives the damping factor  $\gamma$ .

## 4 Fourier transform is complex

For a real function  $f(t)$ , the Fourier transform will usually not be real. Indeed, the imaginary part of the Fourier transform of a real function is

$$\text{Im}[\tilde{f}(k)] = \frac{\tilde{f}(k) - \tilde{f}(k)^*}{2i} = \frac{1}{2i} \frac{1}{2\pi} \left[ \int_{-\infty}^{\infty} dx f(x) e^{-ikx} - \int_{-\infty}^{\infty} dx f(x) e^{ikx} \right] \quad (16)$$

$$= \frac{1}{2\pi} \int_{-\infty}^{\infty} dx f(x) \sin(kx) \equiv \tilde{f}_s(k) \quad (17)$$

This is a Fourier sine transform. Thus the imaginary part vanishes only if the function has no sine components which happens if and only if the function is even. For an odd function, the Fourier transform is purely imaginary. For a general real function, the Fourier transform will have both real and imaginary parts. We can write

$$\tilde{f}(k) = \tilde{f}_c(k) + i\tilde{f}_s(k) \quad (18)$$

where  $\tilde{f}_s(k)$  is the Fourier sine transform and  $\tilde{f}_c(k)$  the Fourier cosine transform. One hardly ever uses Fourier sine and cosine transforms. We practically always talk about the complex Fourier transform.

Rather than separating  $\tilde{f}(k)$  into real and imaginary parts, which amounts to Cartesian coordinates, it is often helpful to write it as a magnitude and phase, as in polar coordinates. So we write

$$\tilde{f}(k) = A(k) e^{i\phi(k)} \quad (19)$$

with  $A(k) = |\tilde{f}(k)|$  the **magnitude** and  $\phi(k)$  the **phase**.

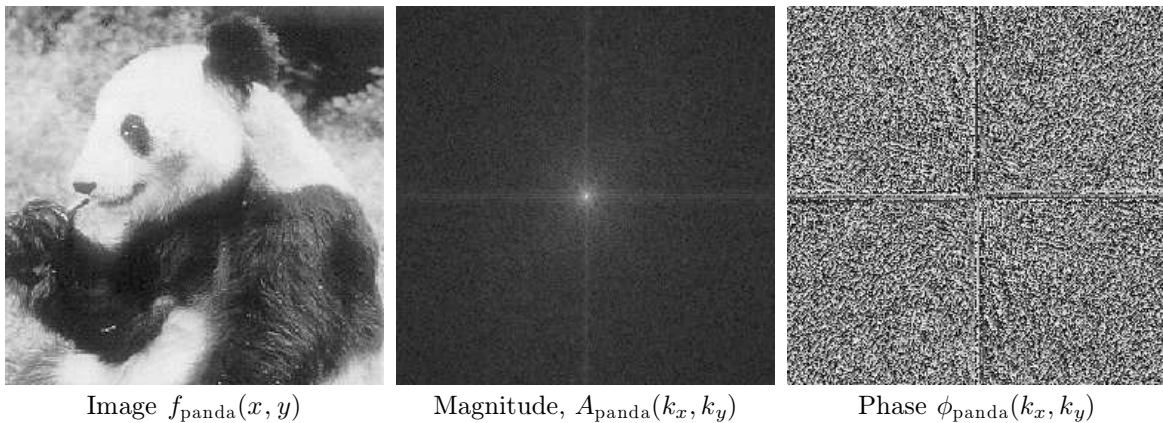


The energy in a frequency mode only depends on the amplitude:  $I = A(\omega)^2$ . When one plots the spectrum as in audacity, what is being shown is  $A(\omega)^2$ . This corresponds to the intensity or power in a particular mode, as we will see in Lecture 10. Power is useful in doing a frequency analysis of sound since it tells us how loud that frequency is. But looking at the amplitude is not the only thing one can do with a Fourier transform. Often one is also interested in the phase.

For a visual example, we can take the Fourier transform of an image. Suppose we have a grayscale image that is  $640 \times 480$  pixels. Each pixel is a number from 0 to 255, going from black (0) to white (255). Thus the image is a function  $f(x, y)$  with  $0 \leq x < 640$ ,  $0 \leq y < 480$  which takes values from 0 to 255. We can then Fourier transform this function to a function  $\tilde{f}(k_x, k_y)$ :

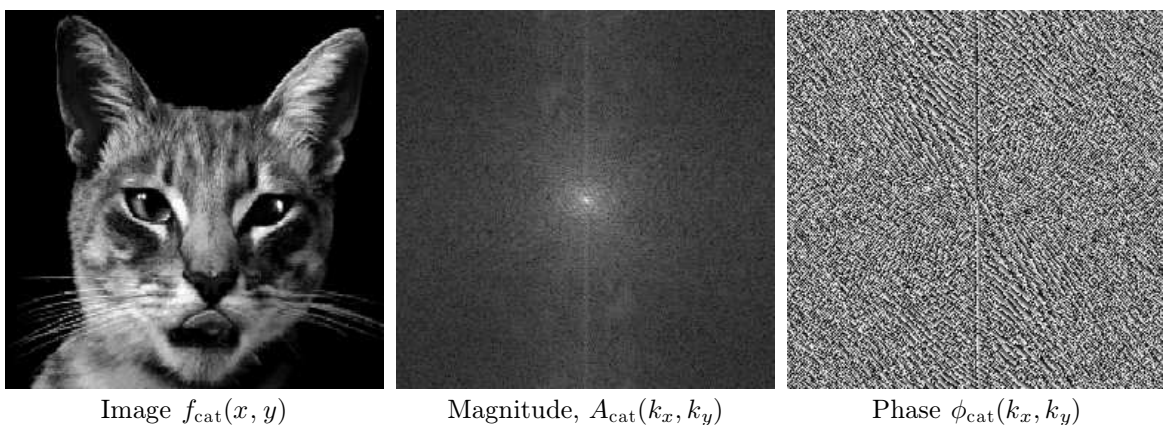
$$\tilde{f}(k_x, k_y) = \frac{1}{2\pi} \int_{-\infty}^{\infty} dx \int_{-\infty}^{\infty} dy f(x, y) e^{-ik_x x} e^{-ik_y y} \quad (20)$$

The 2D Fourier transform is really no more complicated than the 1D transform – we just do two integrals instead of one. So what we do we get? Here's an example



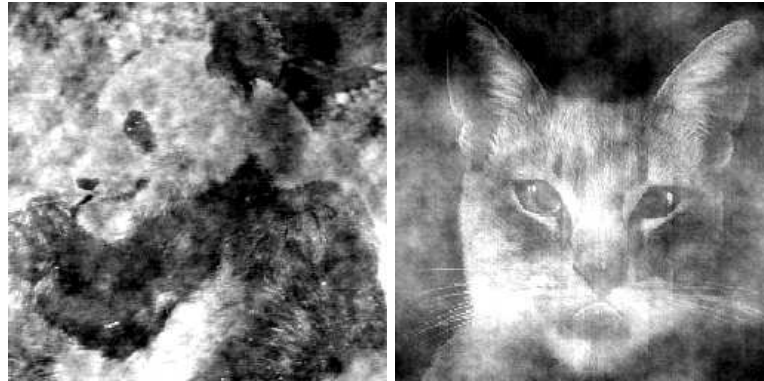
**Figure 3.** Fourier transform of a panda. The magnitude is concentrated near  $k_x \sim k_y \sim 0$ , corresponding to large-wavelength variations, while the phase looks random.

We can do the same thing for a picture of a cat:



**Figure 4.** Fourier transform of a cat. The magnitude is concentrated near  $k_x \sim k_y \sim 0$ , but maybe not as much as the panda, since that cat has smaller wavelength features. Phase still looks random.

Now let's Fourier transform back. Of course for the cat and panda we get back the original image. But what happens if we combine the magnitude for the panda with the phase for the cat, and vice versa?



$A_{\text{cat}}(k_x, k_y)$  and  $\phi_{\text{panda}}(k_x, k_y)$      $A_{\text{panda}}(k_x, k_y)$  and  $\phi_{\text{cat}}(k_x, k_y)$

**Figure 5.** We take the inverse Fourier transform of function  $A_{\text{cat}}(k_x, k_y)e^{i\phi_{\text{panda}}(k_x, k_y)}$  on the left, and  $A_{\text{panda}}(k_x, k_y)e^{i\phi_{\text{cat}}(k_x, k_y)}$  on the right.

It looks like the phase is more important than the magnitude for reconstructing the original image. The importance of phase is critical for many engineering applications, such as signal analysis. It is also relevant for image compression technologies.

## 5 Filtering

One thing we can do with the Fourier transform of an image is remove some components. If we remove low frequencies, less than some  $\omega_f$  say, we call it a **high-pass filter**. A lot of background noise is at low frequencies, so a high-pass filter can clean up a signal. If we throw out the high frequencies, it is called a **low-pass filter**. A low pass filter can be used to smooth data (such as a digital photo) since it throws out high frequency noise. A filter that cuts out both high and low frequencies is called a **band-pass filter**.

Here are some examples

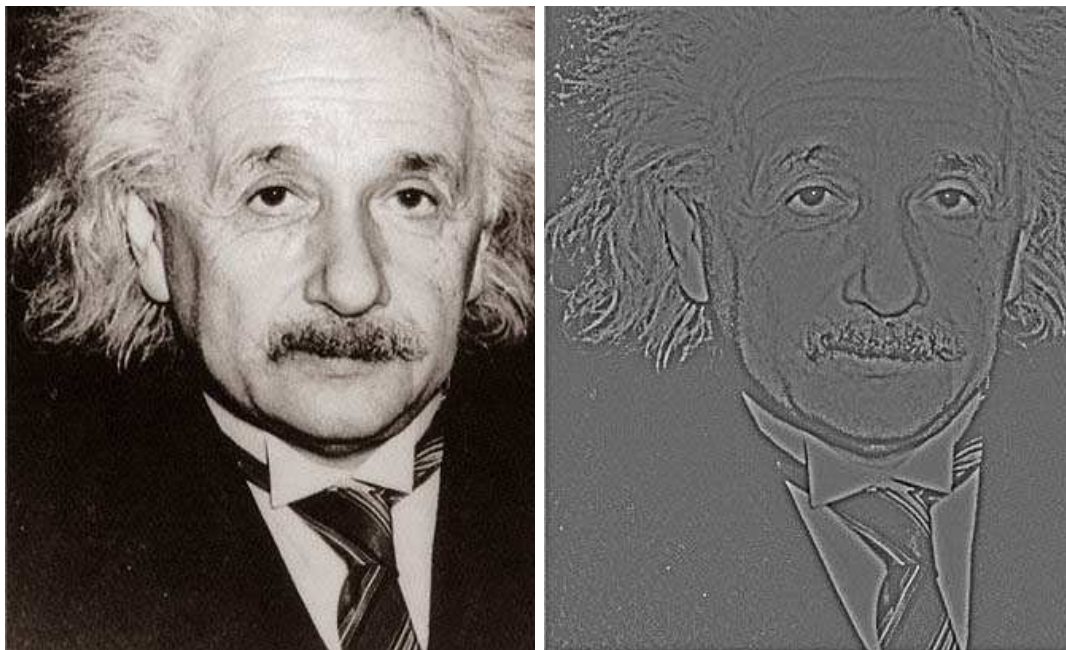


photo of Einstein

Photo after high-pass filter

**Figure 6.** What a high-pass filter does to Albert Einstein.



photo of Einstein

Photo after low-pass filter

**Figure 7.** What a low-pass filter does to Marylyn Monroe.

Now let's combine the two

high-pass Einstein  
+low pass Marylynlow-pass Einstein  
+high-pass Marylyn**Figure 8.** Combining filtered images

Take a look at these last two images from up close and from far away. What do you see? Why?

## 6 Dirac $\delta$ function

Another extremely important example is the Fourier transform of a constant:

$$\delta(\omega) \equiv \frac{1}{2\pi} \int_{-\infty}^{\infty} dt e^{-i\omega t} \quad (21)$$

Its Fourier inverse is then

$$1 = \int_{-\infty}^{\infty} d\omega \delta(\omega) e^{i\omega t} \quad (22)$$

This object  $\delta(\omega)$  is called the **Dirac  $\delta$  function**. It is enormously useful in a great variety of physics problems, especially in quantum mechanics, but also in waves.

To figure out what  $\delta(\omega)$  looks like, we use the fact that the Fourier transform of the inverse Fourier transform gives a function back. That is, for any smooth function  $f(x)$

$$f(x) = \int_{-\infty}^{\infty} dk e^{ikx} \tilde{f}(k) = \frac{1}{2\pi} \int_{-\infty}^{\infty} dk e^{ikx} \int_{-\infty}^{\infty} dy f(y) e^{-iky} \quad (23)$$

$$= \int_{-\infty}^{\infty} dy \frac{1}{2\pi} \int_{-\infty}^{\infty} dk e^{-ik(y-x)} f(y) \quad (24)$$

$$= \int_{-\infty}^{\infty} dy \delta(y-x) f(y) \quad (25)$$

where we used Eq. (21) in the last step. Setting  $x=0$ , we see that the  $\delta$ -function satisfies

$$\int_{-\infty}^{\infty} dx \delta(x) f(x) = f(0) \quad (26)$$

for any smooth function  $f(x)$ .  $\delta(x)$  also has the property that  $\delta(x) = 0$  for  $x \neq 0$  (see Section 6.1 below), so that

$$\int_{-x_0}^{x_0} dx \delta(x) f(x) = f(0) \quad (27)$$

for any  $x_0$ .

Eq. (26) and (27) uniquely define the  $\delta$ -function. Indeed, the  $\delta$ -function is no ordinary function. It is instead a member of a class of mathematical objects called **distributions**. While functions take numbers and give numbers (like  $f(x) = x^2$ ), distributions only give numbers after being integrated.

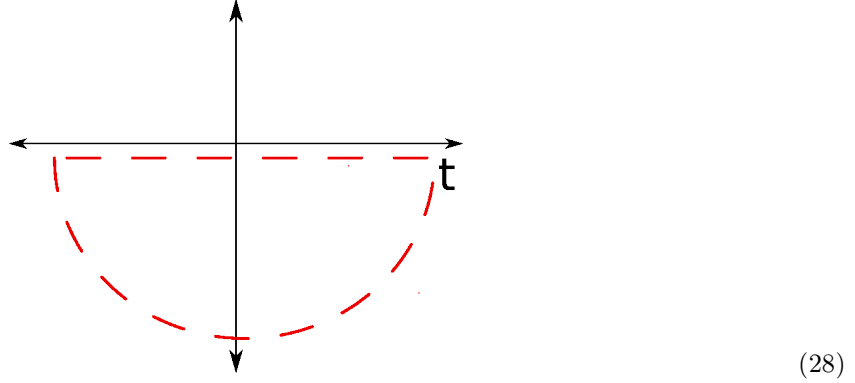
You should think of  $\delta(x)$  as zero everywhere except at  $x = 0$  where it is infinite. However, the infinity is integrable:  $\int_{-x_0}^{x_0} \delta(x) = 1$  for any  $x_0 > 0$ .

From the physics point of view, we showed that if we have an amplitude which is constant in time  $f(t) = 1$  then the only frequency mode supported has 0 frequency. This makes sense – a constant has an infinite wavelength and never repeats. Conversely, if  $\tilde{f}(\omega) = 1$  it says that all frequencies are excited. This corresponds to **white noise**. The Fourier transform of  $\tilde{f}(\omega) = 1$  gives a function  $f(t) = \delta(t)$  which corresponds to an infinitely sharp pulse. For a pulse has no characteristic time associated with it, no frequency can be picked out. That's why white noise has all frequencies equally.

### 6.1 Some mathematics of $\delta(\omega)$ (optional)

For  $\omega \neq 0$  the quickest way to evaluate  $\delta(\omega)$  integral is by contour integration. If you've never seen any complex analysis, just ignore this section. If you have, consider the integral in the com-

plex  $\omega$  plane along the red contour:



The integral along the contour is equal to  $2\pi i$  times the residues of poles within the contour.

$$\int_{-\infty}^{\infty} dt e^{-i\omega t} f(t) + \int_{\text{curve}} dt e^{-i\omega t} f(t) = 2\pi i \sum_{\text{poles } \omega_j} \text{Res}[f, \omega_j] \quad (29)$$

For the curved part of the contour,  $t$  has a negative imaginary part. Thus  $e^{-i\omega t} \rightarrow 0$  as  $|t| \rightarrow \infty$  and the integral along the curved part vanishes. There are no poles in  $e^{-i\omega t}$ , thus the right hand side of Eq. (29) vanishes. Therefore

$$\delta(\omega) = 0, \quad \omega \neq 0 \quad (30)$$

On the other hand, for  $\omega = 0$ ,

$$\delta(0) = \frac{1}{2\pi} \int_{-\infty}^{\infty} dt = \infty \quad (31)$$

So

$$\delta(\omega) = \begin{cases} 0, & \omega \neq 0 \\ \infty, & \omega = 0 \end{cases} \quad (32)$$

Clearly  $\delta(\omega)$  is no ordinary function. It is a distribution.

A practical way to define  $\delta(x)$  is as a limit. There are lots of ways to do this. Here are three:

$$\delta(x) = \lim_{\varepsilon \rightarrow 0} \frac{1}{\pi} \frac{\varepsilon}{x^2 + \varepsilon^2}, \quad \delta(x) = \lim_{\varepsilon \rightarrow 0} \varepsilon \left( \frac{1}{x} \right)^{1-\varepsilon}, \quad \delta(x) = \lim_{\varepsilon \rightarrow 0} \frac{1}{2\sqrt{\pi\varepsilon}} e^{-\frac{x^2}{4\varepsilon}}, \quad \dots \quad (33)$$

To check these definitions, try integrating any of them against any test function  $g(x)$  to see that Eq. (27) is reproduced.

# Lecture 9:

## Reflection, Transmission and Impedance

### 1 Boundary conditions at a junction

Suppose we take two taut strings, one thick and one thin and knot them together. What will happen to a wave as it passes through the knot? Or, instead of changing the mass density at the junction, we could change the tension (for example, by tying the string to a ring on a fixed rod which can absorb the longitudinal force from the change in tension). What happens to a sound wave when it passes from air to water? What happens to a light wave when it passes from air to glass? In this lecture, we will answer these questions.

Let's start with the string with varying tension. Say there is a knot at  $x = 0$  and the tension changes abruptly between  $x < 0$  and  $x > 0$ . To be concrete, imagine we have a left-moving traveling wave coming in at very early times, hitting the junction around  $t = 0$  (obviously all parts of the wave can't hit the junction at the same time). We would like to know what the wave looks like at late times. Let us write the amplitude of the wave as  $\psi_L(x, t)$  to the left of the knot at  $\psi_R(x, t)$  to the right of the knot.

$$\psi(x, t) = \begin{cases} \psi_L(x, t), & x < 0 \\ \psi_R(x, t), & x \geq 0 \end{cases} \quad (1)$$

To the left of the knot, the wave must satisfy one wave equation

$$\left[ \frac{\partial^2}{\partial t^2} - v_1^2 \frac{\partial^2}{\partial x^2} \right] \psi_L(x, t) = 0, \quad v_1 = \sqrt{\frac{T_1}{\mu_1}} \quad (2)$$

and to the right of the knot, another wave equation must be satisfied

$$\left[ \frac{\partial^2}{\partial t^2} - v_2^2 \frac{\partial^2}{\partial x^2} \right] \psi_R(x, t) = 0, \quad v_2 = \sqrt{\frac{T_2}{\mu_2}} \quad (3)$$

Recalling that the Heaviside step function (or theta-function) is defined by  $\theta(x) = 0$  if  $x < 0$  and  $\theta(x) = 1$  for  $x \geq 0$ , we can also write Eq. (1) as

$$\psi(x, t) = \psi_L(x, t)\theta(-x) + \psi_R(x, t)\theta(x) \quad (4)$$

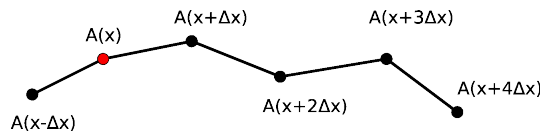
This way of writing  $\psi(x, y)$  makes it clear that it is just some function of position and time. We need to determine what the boundary conditions are at the junction, and then find the full solution  $\psi(x, t)$  for all times.

Obviously  $\psi(x, t)$  should be continuous. So

$$\boxed{\psi_L(0, t) = \psi_R(0, t)} \quad (5)$$

This is one boundary condition at the junction.

Recall from Lecture 6 that a point on the string of mass  $m$  gets a force from the parts of the string to the left and to the right:





The force from the part to the left is  $T \frac{\Delta\psi}{\Delta x} \approx T \frac{\partial\psi(x,t)}{\partial x}$ . This form makes sense, since if the string has no slope, it is flat and there is no force. From the right, the force is  $-T \frac{\partial\psi(x,t)}{\partial x}$ . The sign has to be opposite so that if there is no difference in slope there is no force (with equal tensions). So if there are different tensions to the right and left, as at  $x=0$ , we have

$$m \frac{\partial^2\psi(0,t)}{\partial t^2} = T_1 \frac{\partial\psi_L(0,t)}{\partial x} - T_2 \frac{\partial\psi_R(0,t)}{\partial x} \quad (6)$$

Now  $m$  is the mass of an infinitesimal point of string at  $x=0$ . But  $T_1$  and  $T_2$  as well as the slopes  $\frac{\partial\psi_L(0,t)}{\partial x}$  and  $\frac{\partial\psi_R(0,t)}{\partial x}$  are macroscopic quantities. Thus, if the right hand side doesn't vanish, we would find  $\frac{\partial^2\psi(0,t)}{\partial t^2} \rightarrow \infty$  as  $m \rightarrow 0$ . Equivalently, we can write  $m = \mu\Delta x$  then this becomes

$$\mu\Delta x \frac{\partial\psi(0,t)}{\partial t^2} = T_1 \frac{\partial\psi_L(0,t)}{\partial x} - T_2 \frac{\partial\psi_R(0,t)}{\partial x} \quad (7)$$

Taking  $\Delta x \rightarrow 0$  we find

$$\boxed{T_1 \frac{\partial\psi_L(0,t)}{\partial x} = T_2 \frac{\partial\psi_R(0,t)}{\partial x}} \quad (8)$$

So the slope must be discontinuous at the boundary to account for the different tensions.

Now we have the boundary conditions. What is the solution?

## 2 Reflection and transmission

Suppose we have some incoming traveling wave. Before it hits the junction it has the form of a right-moving traveling wave

$$\psi_L(x,t) = \psi_i(x - v_1 t), \quad t < 0 \quad (9)$$

To be clear,  $\psi_L(x,t)$  is the part of  $\psi(x,t)$  with  $x < 0$ .  $\psi_i(x)$  is some function describing the wave's shape in this region. It is easy to check that  $\psi_L(x,t)$  satisfies the wave equation in the  $x < 0$  region:  $\left[ \frac{\partial^2}{\partial t^2} - v_1^2 \frac{\partial^2}{\partial x^2} \right] \psi_L(x,t) = 0$ . The  $i$  subscript on  $\psi_i(t)$  refers to the **incident wave**. Let  $t=0$  be the time when the first part of the wave hits the knot at  $x=0$ .

To be concrete, think of  $\psi_i(t)$  as a square wave. For example  $\psi_i(z) = 2\text{mm}$  for  $-1\text{ cm} < z \leq 0\text{cm}$  and  $\psi_i(z) = 0$  otherwise. At  $t=0$ ,  $\psi_L(x,0)$  is zero outside of  $-1\text{ cm} < x < 0$ , so it just starts to hit  $x=0$ . At earlier times, say  $t_1 = -\frac{5\text{cm}}{v_1}$ , then  $\psi_L(x,t_1)$  is zero outside of  $-6\text{ cm} < x < -5\text{cm}$ . So as time goes on, it approaches the junction, and hits it just at  $t=0$ . So  $\psi(x,t) = \psi_L(x,t)\theta(-x)$  is a perfectly good solution of the wave equation for  $t < 0$ . The real wave doesn't have to be a square wave, it can have any shape.

Actually, it will be extremely helpful to make a cosmetic change and write  $\psi_i\left(t - \frac{x}{v_1}\right)$  instead of  $\psi_i(x - v_1 t)$ . Clearly these functions carry the same information, because we just rescaled the argument. The new form is nicer since at the boundary  $x=0$ ,  $\psi_i$  doesn't depend on  $v$  (so Eq. (12) below has a simple form). So let's pretend we wrote  $\psi_i\left(t - \frac{x}{v_1}\right)$  from the start of this section (I didn't want to actually write it that way from the start to connect more clearly to what we did before).

Now, after  $t=0$   $\psi_L$  can have left and a right moving components, so we can more generally write

$$\psi_L(x,t) = \psi_i\left(t - \frac{x}{v_1}\right) + \psi_r\left(t + \frac{x}{v_1}\right) \quad (10)$$

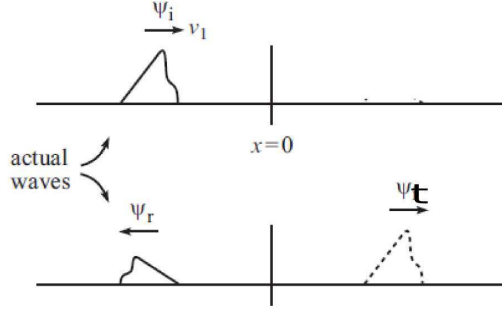
where  $\psi_r$  is the **reflected wave**. Recall that any wave can be written as a sum of left and right moving waves. So writing  $\psi_L$  this way does not involve any assumptions, it is just convenient to solve the wave equation including boundary conditions at the junction.

For  $t > 0$  there will also be some  $\psi_R$  (the part at  $x > 0$ ). This part will always be right-moving. We call this the **transmitted wave** and write

$$\psi_R(x, t) = \psi_t\left(t - \frac{x}{v_2}\right) \quad (11)$$

That we can write the wave for  $x > 0$  in this form follows from the assumption that for  $t < 0$  then  $\psi = 0$  for  $x > 0$ . If there were a left-moving component on the right side, then as  $t \rightarrow -\infty$  it would always be there. Note that the transmitted wave has wave speed  $v_2$ , since it is in the string on the right. Note that we are not assuming that the incident, transmitted and reflected waves all have the same shape.

The picture is as follows



**Figure 1.** Incident, reflected and transmitted waves.

Now we impose our boundary conditions. Continuity at  $x = 0$ , Eq. (5) implies

$$\psi_i(t) + \psi_r(t) = \psi_t(t) \quad (12)$$

For the other boundary condition, Eq. (8), we have

$$T_1 \frac{\partial \psi_L(0, t)}{\partial x} = T_1 \frac{\partial}{\partial x} \left[ \psi_i\left(t - \frac{x}{v_1}\right) + \psi_r\left(t + \frac{x}{v_1}\right) \right]_{x=0} = \frac{T_1}{v_1} [-\psi'_i(t) + \psi'_r(t)] \quad (13)$$

and

$$T_2 \frac{\partial \psi_R(0, t)}{\partial x} = T_2 \frac{\partial}{\partial x} \left[ \psi_t\left(t - \frac{x}{v_2}\right) \right]_{x=0} = -\frac{T_2}{v_2} \psi'_t(t) \quad (14)$$

Thus,

$$\frac{T_1}{v_1} [-\psi'_i(t) + \psi'_r(t)] = -\frac{T_2}{v_2} \psi'_t(t) \quad (15)$$

In other words

$$\frac{d}{dt} \left[ -\frac{T_1}{v_1} \psi_i(t) + \frac{T_1}{v_1} \psi_r(t) + \frac{T_2}{v_2} \psi_t(t) \right] = 0 \quad (16)$$

Since a function whose derivative vanishes must be constant, we then have

$$\frac{T_1}{v_1} [-\psi_i(t) + \psi_r(t)] = -\frac{T_2}{v_2} \psi_t(t) + \text{const} \quad (17)$$

If the constant were nonzero, it would mean that the wave on the righthand side,  $\psi_t$  has a net displacement at all times. There is nothing particularly interesting in such a displacement, so we set the integration constant to zero.

Substituting Eq. (12) into Eq. (17) we get

$$\frac{T_1}{v_1} [-\psi_i(t) + \psi_r(t)] = -\frac{T_2}{v_2} [\psi_i(t) + \psi_r(t)] \quad (18)$$

or

$$\left( \frac{T_1}{v_1} + \frac{T_2}{v_2} \right) \psi_r = \left( \frac{T_1}{v_1} - \frac{T_2}{v_2} \right) \psi_i(t) \quad (19)$$



which implies

$$\psi_r = \frac{\frac{T_1}{v_1} - \frac{T_2}{v_2}}{\frac{T_1}{v_1} + \frac{T_2}{v_2}} \psi_i \quad (20)$$

We have found that the reflective wave has exactly the same shape as the incident wave, but with a different overall magnitude. By Eq. (12) the transmitted wave also has the same shape. The relevant amplitudes are the main useful formulas coming out of this analysis.

Defining

$$Z_1 = \frac{T_1}{v_1}, \quad Z_2 = \frac{T_2}{v_2} \quad (21)$$

we have

$$\boxed{\psi_r = \frac{Z_1 - Z_2}{Z_1 + Z_2} \psi_i} \quad (22)$$

Substituting back in to Eq. (12) we get

$$\boxed{\psi_t = \frac{2Z_1}{Z_1 + Z_2} \psi_i} \quad (23)$$

Sometimes this solution is written as

$$\psi_r = R\psi_i, \quad \psi_t = T\psi_i \quad (24)$$

where

$$R = \frac{Z_1 - Z_2}{Z_1 + Z_2} \quad (25)$$

is the **reflection coefficient** and

$$T = \frac{2Z_1}{Z_1 + Z_2} \quad (26)$$

is the **transmission coefficient**.

$Z$  is known as an **impedance**. In this case it's tension over velocity, but more generally

### Impedance is force divided by velocity

That is, impedance tells you how much force is required to impart a certain velocity. Impedance is a property of a medium. In this case, the two strings have different tensions and different velocities. Using  $v = \sqrt{\frac{T}{\mu}}$  we can write

$$Z = \frac{T}{v} = \sqrt{T\mu} \quad (27)$$

Note that as  $Z_1 = Z_2$  there is no reflection and complete transmission. If we want no reflection we need to **match impedances**. For example, if we want to impedance-match across two strings with different mass densities  $\mu_1$  and  $\mu_2$  we can choose the tensions to be  $T_2 = \frac{\mu_1}{\mu_2} T_1$  so that

$$Z_2 = \sqrt{T_2 \mu_2} = \sqrt{\frac{\mu_1}{\mu_2} T_1 \mu_2} = \sqrt{T_1 \mu_1} = Z_1 \quad (28)$$

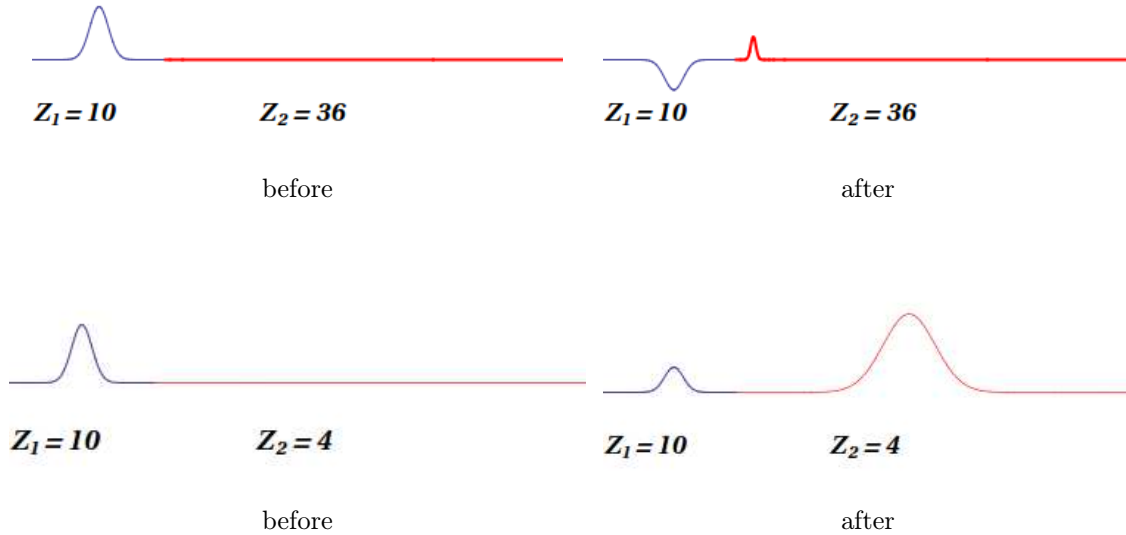
Thus the impedances can agree in strings of different mass density.

Note that the transmission coefficient is greater than 1 if  $Z_1 < Z_2$ . That means the amplitude increases when a wave travels from a medium of lower impedance to a medium of higher impedance. This is an important fact. We'll discuss a consequence in Section 7.1 below.

## 3 Phase flipping

What happens when a wave hits a medium of higher impedance, such as when the tension or mass density of the second string is very large? Then  $Z_2 > Z_1$  and so,  $R = \frac{Z_1 - Z_2}{Z_1 + Z_2} < 0$ . Thus, if  $\psi_i > 0$  then  $\psi_r < 0$ . That is, the wave flips its sign. This happens in particular if the wave hits a wall, which is like  $\mu = \infty$ .

On the other hand if a wave passes to a less dense string then  $Z_2 < Z_1$  and there is no sign flip. This can happen if  $Z_2 = 0$ , for example, if the second string is massless or tensionless – as in an open boundary condition.

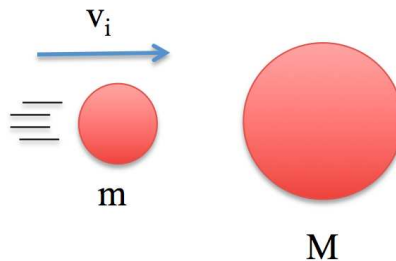


**Figure 2.** Phase shift of reflected pulses on a string. Top has pulse going from lower to higher impedance. Bottom has pulse going from higher to lower impedance.

This phase flipping has important consequences due to interference between the reflected pulse and other incoming pulses. There will be constructive interference if the phases are the same, but destructive interference if they are opposite. We will return to interference after discussing light.

## 4 Impedance for masses

To get intuition for impedance, it is helpful to go back to a more familiar system: masses. Suppose we collide a block of mass  $m$  with a larger block of mass  $M$

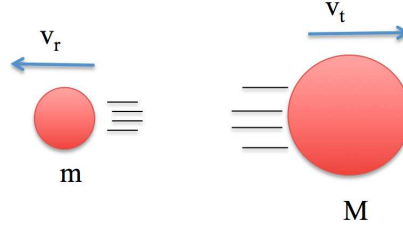


**Figure 3.** Mass  $m$  starts with velocity  $v_i$ , with  $M$  at rest.

Say  $m$  has velocity  $v_i$ . To find out the velocity of  $M$  we solve Newton's laws, or more easily, use conservation of momentum and energy. The initial momentum and energy are

$$p_i = m v_i, \quad E_i = \frac{1}{2} m v_i^2 \quad (29)$$

After the collisions,  $m$  bounces off  $M$  and goes back the way it came with “reflected velocity”  $v_r$  and  $M$  moves off to the right with “transmitted velocity”  $v_t$ :



**Figure 4.** After the collision,  $m$  has the reflected velocity  $v_r$  and  $M$  the transmitted velocity  $v_t$ .

The final momentum and energy are

$$p_f = -v_r m + v_t M, \quad E_f = \frac{1}{2} m v_r^2 + \frac{1}{2} M v_t^2 \quad (30)$$

Conservation of momentum implies

$$v_t = \frac{m}{M} (v_i + v_r) \quad (31)$$

then conservation of energy implies

$$\frac{1}{2} m v_i^2 = \frac{1}{2} m v_r^2 + \frac{1}{2} M \left[ \frac{m}{M} (v_i + v_r) \right]^2 \quad (32)$$

After a little more algebra we find

$$v_r = \frac{M - m}{M + m} v_i, \quad v_t = \frac{2m}{m + M} v_i \quad (33)$$

These equations have exactly the same form as Eqs. (22) and (23) with  $Z_1 = m$  and  $Z_2 = M$ . Thus for masses, **impedance is mass**. This makes sense – the bigger the mass, the less force you can impart with a given velocity.

Let's take a concrete example. Suppose  $m = 1$ ,  $M = 3$  and the incoming velocity is  $v$ . Then the final velocity of  $M$  is

$$v_t = \frac{2m}{m + M} v = \frac{2(1)}{1 + 3} v = \frac{1}{2} v \quad (34)$$

Thus the mass  $M$  gets half the velocity of  $m$ . Now say we put a mass  $m_2 = 2$  in between them. When  $m$  bangs into  $m_2$  it gives it a velocity

$$v_2 = \frac{2m}{m + m_2} v = \frac{2(1)}{1 + 2} v = \frac{2}{3} v \quad (35)$$

Then  $m_2$  bangs into  $M$  and gives it a velocity

$$v_t = \frac{2m_2}{m_2 + M} v_2 = \frac{2(2)}{2 + 3} \left( \frac{2}{3} v \right) = \frac{4}{5} \frac{2}{3} v = \frac{8}{15} v = 0.533 v \quad (36)$$

Thus  $M$  goes faster. Thus inserting a mass between the two masses helps impedance match. Similarly inserting lots of masses can make the impedance matching very efficient.



## 5 Complex impedance

It is sometimes useful to generalize impedance to complex numbers. For example, suppose we have a driven oscillator satisfying

$$m\ddot{x} + kx = F_0 e^{i\omega t} \quad (37)$$

First consider the case where  $k \approx 0$ . Then  $m\ddot{x} = F_0 e^{i\omega t}$ . Integrating this gives

$$\dot{x} = \frac{F_0}{i\omega m} e^{i\omega t} \quad (38)$$

Then

$$Z_m = \frac{\text{force}}{\text{velocity}} = \frac{F_0 e^{i\omega t}}{\frac{F_0}{i\omega m} e^{i\omega t}} = i\omega m \quad (39)$$

Thus at fixed driving frequency  $\omega$ ,  $Z_m \propto m$  as with the masses.

In the other case, when  $m \approx 0$ ,  $kx = F_0 e^{i\omega t}$  and so

$$\dot{x} = i\omega \frac{F_0}{k} e^{i\omega t} \quad (40)$$

Then

$$Z_k = \frac{\text{force}}{\text{velocity}} = \frac{F_0 e^{i\omega t}}{i\omega \frac{F_0}{k} e^{i\omega t}} = -i \frac{k}{\omega} \quad (41)$$

The impedance of the whole system is the sum of the impedances

$$Z_{\text{total}} = Z_m + Z_k = i \left( \omega m - \frac{k}{\omega} \right) \quad (42)$$

So at high frequencies, the mass term dominates. This is called **mass-dominated impedance**. Physically, when the driver is going very fast, the mass has no time to react: a lot of force at high frequency has little effect on velocity. At low frequencies, the  $k$  term dominates. For slow motion, how much velocity the mass gets for a given force depends very much on how stiff the spring is. This is called **stiffness dominated impedance**.

Note that  $Z_{\text{total}} = 0$  when  $\omega = \sqrt{\frac{k}{m}}$ , that is, no resonances. At the resonant frequency, nothing impedes the motion of the oscillator: a small force gives a huge velocity.

With complex impedances you can add a damping term.

$$\gamma \dot{x} = F_0 e^{i\omega t} \Rightarrow \dot{x} = \frac{F_0}{\gamma} e^{i\omega t} \quad (43)$$

Thus,

$$Z_\gamma = \frac{F}{v} = \gamma \quad (44)$$

This makes perfect sense: damping impedes the transfer of energy from the driver to the oscillator.

With all 3 terms,

$$Z_{\text{total}} = \gamma + i \left( \omega m - \frac{k}{\omega} \right) \quad (45)$$

Now the impedance is always nonzero, for any frequency.

## 6 Circuits (optional)

An important use of complex impedances is in circuits. Recall that the equation of motion for an LRC circuit is just like a damped harmonic oscillator. For a resistive circuit:

$$V = IR = \dot{Q}R \quad (46)$$

where  $Q$  is the charge,  $I$  is the current,  $R$  is the resistance and  $V$  is the voltage. For a capacitor

$$V = \frac{Q}{C} \quad (47)$$

For an inductor

$$V = L\dot{I} = L\ddot{Q} \quad (48)$$

Putting everything together, the total voltage is

$$V_{\text{tot}} = L\ddot{Q} + \frac{Q}{C} + \dot{Q}R \quad (49)$$

This is the direct analog of

$$F = m\ddot{x} + kx + \gamma\dot{x} \quad (50)$$

Instead of driving the mass with an external force  $F = F_0 e^{i\omega t}$ , we drive the circuit with an external voltage  $V = V_0 e^{i\omega t}$ . That is we find the simple correspondence

mass/spring	$F$	$x$	$\dot{x}$	$\ddot{x}$	$\gamma$	$k$	$m$	$Z = \frac{F}{\ddot{x}}$
circuit	$V$	$Q$	$I = \dot{Q}$	$\dot{I} = \ddot{Q}$	$R$	$\frac{1}{C}$	$L$	$Z = \frac{V}{\ddot{Q}}$

(51)

Thus instead of being  $Z = \frac{F}{\ddot{x}}$ , impedance for a circuit is

$$Z = \frac{V}{\ddot{Q}} = \frac{V}{I} \quad (52)$$

A resistor has

$$Z_R = \frac{V}{I} = R \quad (53)$$

A capacitor has

$$Z_C = \frac{V}{I} = \frac{1}{i\omega C} \quad (54)$$

and an inductor has

$$Z_L = i\omega L \quad (55)$$

Impedance of an AC circuit plays the role that resistance does for a DC circuit. We can add impedances in series or in parallel just like we do for resistance. Impedance has the units of resistance, that is Ohms. In practice, impedances are more easily measured than calculated.

Matching the impedances of two different wave carrying media is of critical importance in electrical engineering. Say one wishes to drive an antenna, such as the wifi antenna on your router. The maximum power we can couple into the antenna occurs when the impedances of the power supply and antenna are equal in magnitude. This is pretty important in high power applications, where can waves which are reflected from your antenna can come back and destroy your amplifying equipment. It's also critical if you are a receiver. All modern radios have impedance matching circuits in them. This is because antennas are resonant devices, and as we just saw, tuning away from resonances causes some impedance. Thus you would need to match your radio input impedance to your antenna as you pick up different wavelengths.

## 7 Impedance for other stuff

For air, we recall  $v = \sqrt{\frac{B}{\rho}}$  with  $B = \gamma p = \rho v^2$  the bulk modulus and  $v$  the speed of sound in the gas. Then

$$Z_0 = \frac{B}{v} = \rho v \quad (56)$$

$Z_0 = \rho v$  is called the **specific impedance**.  $Z_0$  is a property of the medium. For example, in air

$$\rho = 1.2 \frac{\text{kg}}{\text{m}^3}, \quad v = 343 \frac{\text{m}}{\text{s}} \Rightarrow Z_{\text{air}} = 420 \frac{\text{kg}}{\text{m}^2 \text{s}} = 420 \frac{\text{Pa} \cdot \text{s}}{\text{m}} \quad (57)$$

for water

$$\rho = 1000 \frac{\text{kg}}{\text{m}^3}, \quad v = 1480 \frac{\text{m}}{\text{s}} \Rightarrow Z_{\text{water}} = 1.48 \times 10^6 \frac{\text{kg}}{\text{m}^2 \text{s}} = 1.48 \times 10^6 \frac{\text{Pa} \cdot \text{s}}{\text{m}} \quad (58)$$

Thus if you try to yell at someone under water, you find that the amount reflected is

$$R = \frac{Z_{\text{air}} - Z_{\text{water}}}{Z_{\text{air}} + Z_{\text{water}}} = -0.9994 \quad (59)$$

So almost all of the sound is reflected (and there is a phase flip).

If the wavelength of the sound waves is smaller than the size of the cavity holding the waves (for example in a pipe) then one must account for this finite size in the impedance. For air in a finite size cavity, the relevant quantity is not the specific impedance (which is a property of the gas itself), but the impedance per area

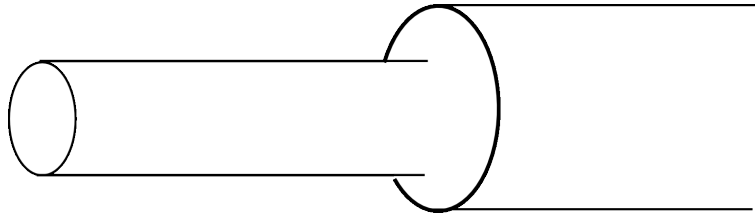
$$Z = \frac{Z_0}{A} = \frac{B}{v \cdot A} = \frac{B}{\text{volume flow rate}} = \frac{\rho v}{A}, \quad \lambda > \sqrt{A} \quad (60)$$

This is relevant when  $\lambda > \sqrt{A}$  where  $\lambda$  is the wavelength of the sound wave and  $A$  is the cross sectional area of the pipe.

For air of the same density, the impedance is effectively  $\frac{1}{\text{area}}$ . Thus the reflection coefficient going between pipes of different radii is

$$R = \frac{\frac{1}{A_1} - \frac{1}{A_2}}{\frac{1}{A_1} + \frac{1}{A_2}} = \frac{A_2 - A_1}{A_1 + A_2} \quad (61)$$

So a situation like this will have bad impedance matching:



On the other hand, a megaphone is designed to impedance match much better:



Now you know why megaphones are shaped this way!

## 7.1 Solids

For liquids or solids, impedance is also  $Z = \rho v$ . The nice thing about a formula like this is that both  $\rho$  (density of the solid) and  $v$  (speed of sound in the solid) are easy to measure, in contrast to the bulk modulus (what is that?) and the pressure (what is pressure for a solid?). For example,

material	density ( $\text{kg}/\text{m}^3$ )	speed of sound ( $\text{m}/\text{s}$ )	specific impedance ( $\text{MPa}\cdot\text{s}/\text{m}$ )
brick	2,200	4,200	9.4
concrete	1,100	3,500	3.8
steel	7,900	6,100	48
water	1,000	1,400	1.4
wood	630	3,600	2.3
rubber	1,100	100	0.11
rock	2,600	6,000	16
diamond	3,500	12,000	42
dirt	1,500	100	0.15

**Table 1.** Properties of various liquids and solids.  $1 \text{ MPa} = 10^6 \text{ Pascals} = 10^6 \frac{\text{kg}}{\text{m}\cdot\text{s}^2}$ .

It's good to have a little intuition for speeds of sounds and densities, which you can get from this table. For example, sound goes very fast in diamond. That's because diamond is very hard and rigid, so the the atoms move back to their equilibrium very quickly as the wave passed through (spring constant is high). Steel is also hard and has a fast sound speed. Rubber and dirt are soft, so waves propagate slowly through them. Dirt is denser than concrete, but sound goes much slower since it is not rigid.

Regarding the impedance, because impedance is  $\rho \cdot v$ , soft stuff generally has small  $\rho$  and small  $v$ , so it has much lower impedance. The highest impedances are for the hardest substances: steel and diamond, the lowest for the softest stuff, water and dirt.

As an application, recall from Eq. (26) that when the impedance goes down  $T > 1$  and the amplitude increases. Now, consider an earthquake as it travels from rock ( $Z_1=16 \text{ MPa}\cdot\text{s}/\text{m}$ ) into dirt or landfill ( $Z_2=0.15 \text{ MPa}\cdot\text{s}/\text{m}$ ). Then  $T = \frac{2Z_1}{Z_1 + Z_2} = 1.98$ . So the amplitude of the shaking will double in amplitude! That's why you shouldn't build houses on landfill in an earthquake zone.

# Lecture 10: Energy and Power in Waves

## 1 Energy in a string

The kinetic energy of a mass  $m$  with velocity  $v$  is  $\frac{1}{2}mv^2$ . Thus if we have a oscillating wave in a string, the kinetic energy of each individual bit of the string is

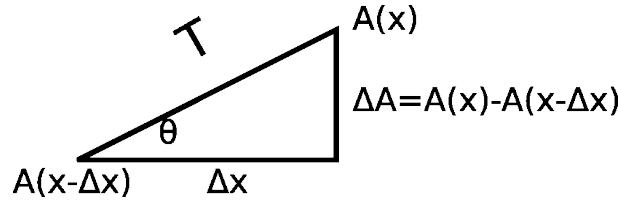
$$\text{KE} = \frac{1}{2}mv^2 = \frac{1}{2}(\mu\Delta x)\left(\frac{\partial A(x,t)}{\partial t}\right)^2 \quad (1)$$

Thus the kinetic energy per unit length is

$$\frac{\text{KE}}{\text{length}} = \frac{1}{2}\mu\left(\frac{\partial A(x,t)}{\partial t}\right)^2 \quad (2)$$

The potential energy depends on how stretched the string is. Of course, having a string with some tension  $T$  automatically has some potential energy due to the stretching, even if there are no waves passing through the string. We are instead interested in the potential energy stored in the string as it is stretched further, due to the propagation of transverse waves. For simplicity, we use potential energy to refer to only this additional potential energy due to the extra strength.

If the string is in equilibrium, so  $\frac{\partial A}{\partial x} = 0$ , then by definition the potential energy is zero. The amount the string is stretched at point  $x$  is given by the difference between the length of the hypotenuse of the triangle and the base. Recall our triangle:



The amount the string is stretched is

$$\Delta L = \sqrt{(\Delta x)^2 + [A(x) - A(x - \Delta x)]^2} - \Delta x \quad (3)$$

$$= \Delta x \sqrt{1 + \left(\frac{\Delta A}{\Delta x}\right)^2} - \Delta x \quad (4)$$

Since the string is close to equilibrium  $\frac{\Delta A}{\Delta x} = \frac{\partial A}{\partial x} \ll 1$ , so we can Taylor expanding the square-root

$$\Delta L = \Delta x \left( 1 + \frac{1}{2} \left( \frac{\partial A}{\partial x} \right)^2 - \frac{1}{8} \left( \frac{\partial A}{\partial x} \right)^4 + \dots \right) - \Delta x \quad (5)$$

and drop subleading terms

$$\Delta L = \frac{1}{2} \Delta x \left( \frac{\partial A}{\partial x} \right)^2 \quad (6)$$

Thus the potential energy is

$$\text{PE} = \text{force} \times \text{distance} = \frac{1}{2}T\Delta x \left( \frac{\partial A(x,t)}{\partial x} \right)^2 \quad (7)$$



And so

$$\frac{\text{PE}}{\text{length}} = \frac{1}{2}T \left( \frac{\partial A(x, t)}{\partial x} \right)^2 \quad (8)$$

Note that this is proportional to the first derivative squared, not the second derivative. Indeed, even if there were no net force on the test mass at position  $x$  (so that  $\frac{\partial^2 A}{\partial x^2} = 0$ ), there would still be potential energy stored in the stretched string. Even though the wave is transverse, the energy comes from stretching the string both longitudinally *and* transversely.

So the total energy per unit length is

$$\frac{E_{\text{tot}}}{\text{length}} = \frac{1}{2}\mu \left( \frac{\partial A}{\partial t} \right)^2 + \frac{1}{2}T \left( \frac{\partial A}{\partial x} \right)^2 \quad (9)$$

Now consider the special case of a traveling wave  $A(x, t) = f(x \pm vt)$ . Then,

$$\frac{\partial A}{\partial t} = \pm v \frac{\partial A}{\partial x} \quad (10)$$

Thus

$$\left( \frac{\partial A}{\partial x} \right)^2 = \frac{1}{v^2} \left( \frac{\partial A}{\partial t} \right)^2 = \frac{\mu}{T} \left( \frac{\partial A}{\partial t} \right)^2 \quad (11)$$

So the total energy per unit length for a traveling wave

$$\frac{E_{\text{tot}}}{\text{length}} = \frac{1}{2}\mu \left( \frac{\partial A}{\partial t} \right)^2 + \frac{1}{2}T \left( \frac{\partial A}{\partial x} \right)^2 = \mu \left( \frac{\partial A}{\partial t} \right)^2 \quad (12)$$

Recalling that impedance for a string is

$$Z = \frac{\text{force}}{\text{velocity}} = \frac{T}{v} = \sqrt{T\mu} = v\mu \quad (13)$$

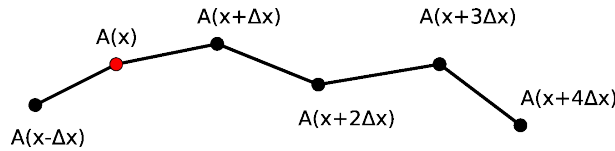
we can write the energy as

$$\frac{E_{\text{tot}}}{\text{length}} = \frac{Z}{v} \left( \frac{\partial A}{\partial t} \right)^2 \quad (14)$$

## 2 Power

An extremely important quantity related to waves is power. We want to use waves to do things, such as transmit sound or light, or energy in a wire. Thus we want to know the rate at which work can be done using a wave. For example, if you have an incoming sound wave, how much power can be transmitted by the wave to a microphone?

For an incoming traveling wave, let us return to this figure



We want to know how much power can be transmitted from the test mass at  $A(x - \Delta x)$  to the test mass at  $A(x)$ . Now, power = force  $\times$  velocity. But we don't want the net force on  $A(x)$ , only the force from the left to compute power transmitted.

We calculated that the force from the left is  $F = T \frac{\partial A}{\partial x}$ . For  $\frac{\partial A}{\partial x} > 0$  this force pulls downward. To see the power transmitted, we need the force which moves the string away from equilibrium, which is the upward force  $F = -T \frac{\partial A}{\partial x}$ . Then

$$P = F \cdot v = -T \left( \frac{\partial A(x, t)}{\partial x} \right) \left( \frac{\partial A(x, t)}{\partial t} \right) \quad (15)$$

For a traveling wave  $A(x, t) = f(x \pm vt)$  we can use Eq. (10) to write this as

$$P = \mp v T \left( \frac{\partial A(x, t)}{\partial x} \right)^2 = \mp v \mu \left( \frac{\partial A(x, t)}{\partial t} \right)^2 = \mp Z \left( \frac{\partial A(x, t)}{\partial t} \right)^2 \quad (16)$$

The sign is + for a right-moving wave (power goes to the right) and – for a left-moving wave.

Now recall that if we have a wave going from a medium with impedance  $Z_1$  into a medium with impedance  $Z_2$ , the amplitude of the transmitted and reflective waves are related to the amplitude of the incoming wave by

$$A_T = \frac{2Z_1}{Z_1 + Z_2} A_I, \quad A_R = \frac{Z_1 - Z_2}{Z_1 + Z_2} A_I, \quad (17)$$

Thus if the power in the incoming wave is

$$P_I = Z_1 \left( \frac{\partial A_I}{\partial t} \right)^2 \quad (18)$$

Then the power in the reflected wave is

$$P_R = Z_1 \left( \frac{Z_1 - Z_2}{Z_1 + Z_2} \frac{\partial A_I}{\partial t} \right)^2 = \left( \frac{Z_1 - Z_2}{Z_1 + Z_2} \right)^2 P_I \quad (19)$$

and the power in the transmitted wave is

$$P_T = Z_2 \left( \frac{2Z_1}{Z_1 + Z_2} \frac{\partial A_I}{\partial t} \right)^2 = \frac{Z_2}{Z_1} \left( \frac{2Z_1}{Z_1 + Z_2} \right)^2 P_I \quad (20)$$

Thus, the fraction of power reflected is

$$\frac{P_R}{P_I} = \left( \frac{Z_1 - Z_2}{Z_1 + Z_2} \right)^2 = \frac{Z_1^2 - 2Z_1Z_2 + Z_2^2}{(Z_1 + Z_2)^2} \quad (21)$$

and the fraction of power transmitted is

$$\frac{P_T}{P_I} = \frac{Z_2}{Z_1} \left( \frac{2Z_1}{Z_1 + Z_2} \right)^2 = \frac{4Z_1Z_2}{(Z_1 + Z_2)^2} \quad (22)$$

Note that

$$\frac{P_T + P_R}{P_I} = 1 \quad (23)$$

so that, overall, power is conserved.

### 3 Sound intensity (decibels)

**Intensity** is defined as power per unit area:  $I = \frac{P}{A}$ . Sound intensity is measured in **decibels**. A logarithmic scale is used for sound intensity because human hearing is logarithmic. For example, if something has an intensity 1000 times larger, you will perceive it as being 3 times as loud.

Decibel is a logarithmic scale normalized so that 0 dB is  $10^{-12} \frac{W}{m^2}$ . That is, by definition

$$0 \text{ dB} \equiv 10^{-12} \frac{\text{Watts}}{\text{meter}^2} = I_0 \quad (24)$$

A decibel is 1 tenth of a bel. The number of bels can be computed from a given intensity by

$$\text{loudness in bels} = 100 \log_{10} \frac{I}{I_0} \quad (25)$$

The number of decibels is therefore

$$\text{loudness in decibels} = 10 \log_{10} \frac{I}{I_0} \quad (26)$$

Some reference intensities are

sound	decibels
Threshold for human hearing	0
Breathing at 3 meters	10
Rustling of leaves	20
...	
Music at 1 meter	70
vacuum cleaner	80
...	
rock concert	120
threshold for pain	130
jet engine at 30 meters	150

**Table 1.** Reference decibel intensities

It is easy to compute the decibel intensity. For example, suppose you are 3 meters away from a 50 Watt speaker. Watts are a unit of power, so at 3 meters, the power is distributed across a sphere of surface area  $A = 4\pi r^2$ . Thus, if all the power went into sound, the intensity would be

$$I = \frac{50 W}{4\pi(3m)^2} = 0.44 \frac{W}{m^2} \quad (27)$$

Thus,

$$L = 10 \log_{10} \frac{0.44 \frac{W}{m^2}}{10^{-12} \frac{W}{m^2}} = 116 \text{dB} \quad (28)$$

This is not actually how loud a speaker is (it's more like the loudness of a rock concert). In reality, the energy of the speaker is only transmitted into sound very inefficiently. We can define the efficiency as the power coming out in sound divided by the power which the speaker draws from the battery. The efficiency is around  $10^{-5}$  for a typical speaker. So,

$$L = 10 \log_{10} \frac{10^{-5} 0.44 \frac{W}{m^2}}{10^{-12} \frac{W}{m^2}} = 66 \text{dB} \quad (29)$$

The efficiency is so low because the speaker and the air have very different impedances.

It takes about 150 mW of power to bow a violin and about 6mW of power comes out in sound. So a violin has an efficiency of  $\varepsilon = 0.04 = 4\%$ . This is much greater than a speaker, but still most of the energy used in bowing a violin is mechanical and not transmitted into sound.

How does the loudness change with distance? Since  $I = \frac{P}{4\pi r^2}$  we have

$$L = 10 \log_{10} \frac{P}{4\pi r^2 I_0} = 10 \log_{10} \frac{P}{4\pi I_0} - 20 \log_{10} r \quad (30)$$

Thus loudness only drops logarithmically with distance. Say you measure a loudness  $L_0$  at a distance  $r_0$ . Then the loudness at a distance  $2r_0$  would be

$$L = 10 \log_{10} \frac{P}{4\pi I_0} - 20 \log_{10}(2r_0) = L_0 - 20 \log_{10} 2 = L_0 - 6.021 \quad (31)$$

That is, it's 6 decibels lower.

Or we can ask at what distance compared to  $r_0$  we should have to go for the loudness to drop by 10 decibels? That would be

$$20 \log_{10} \frac{r}{r_0} = 10 \quad \Rightarrow \quad r = 3.162 r_0 \quad (32)$$

So if you go 3 times farther away, you are down by 10 decibels.

## 4 Plane waves

The waves will propagate in 3 dimensions, so we need the 3-dimensional version of the wave equation:

$$\left[ \frac{\partial^2}{\partial t^2} - v^2 \left( \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} + \frac{\partial^2}{\partial z^2} \right) \right] A(x, y, z, t) = 0 \quad (33)$$

This is the obvious generalization of the 1D wave equation. It is invariant under rotations of  $x$ ,  $y$  and  $z$  (in fact, it is invariant under a larger group of symmetries, Lorentz transformations, which mix space and time).

There are many solutions to this 3D wave equation. Important solutions are **plane waves**

$$A(x, y, z, t) = A_0 \cos(\vec{k} \cdot \vec{x} - \omega t + \phi) \quad (34)$$

for some amplitude  $A_0$ , frequency  $\omega$  and fixed vector  $\vec{k}$  called the **wavevector**. For a plane wave to satisfy the wave equation, its frequency and wavevector must be related by

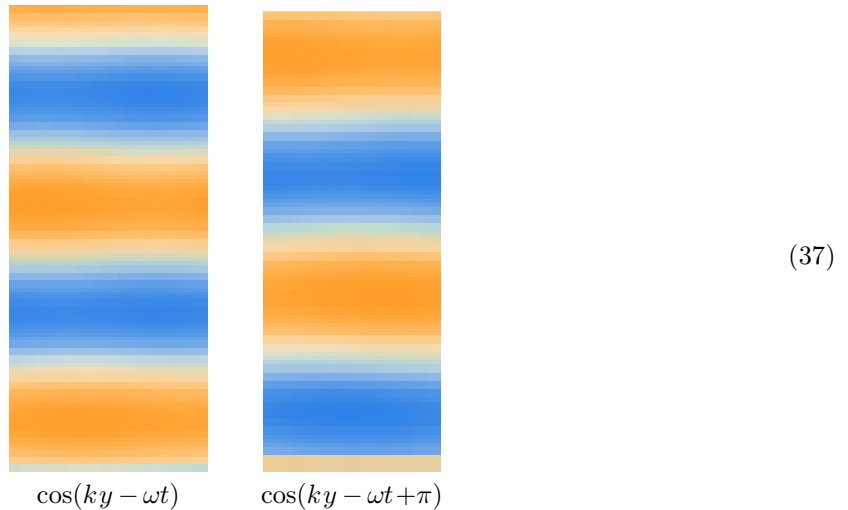
$$\omega = v |\vec{k}| \quad (35)$$

Thus norm of  $\vec{k}$  is fixed by  $\omega$  and  $v$ , but  $\vec{k}$  can point in any direction. The direction  $\vec{k}$  points is the direction the plane wave is traveling. For example, if  $\vec{k} = (0, k, 0)$  the wave is

$$A(x, y, z, t) = A_0 \cos(k(y - vt) + \phi) \quad (36)$$

which is a wave traveling in the  $y$  direction with angular frequency  $\omega = kv$ . This plane wave is constant in the  $x$  and  $z$  directions.

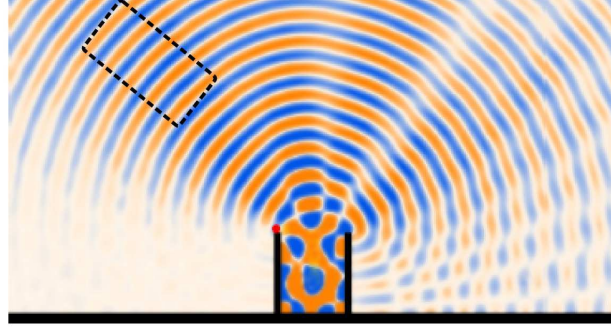
Planes waves have **direction** and **phase**. Falstad's ripple simulation (see link on isite) program gives a nice way to visualize waves. It shows the amplitude of the wave by different by colors. Plane waves going in the  $y$  direction look like



These are both plane waves, but they have different phases. These happen to be exactly out of phase.

Plane waves form a basis of all possible solutions to the wave equation. They are the normal modes of the 3D wave equation. For each frequency  $\omega$  there are plane waves in any direction  $\vec{k}$  with  $|\vec{k}| = \frac{\omega}{v}$  with any possible phase.

Another important feature of plane waves is that if you are far enough away from sources, *everything reduces to a plane wave*. For example, if we have some messy source in a cavity, the solution to the wave equation might look like



Inside the dashed box, the solution is very similar to the solution of the plane wave.

How much power is in a plane wave? At a time  $t$  the part of the wave in Eq. (36) at position  $y$  has power

$$P(t, y) = Z \left( \frac{\partial A}{\partial t} \right)^2 = Z A_0^2 \omega^2 \sin^2(ky - \omega t + \phi) \quad (38)$$

This power is always positive but it oscillates from 0 to its maximum value as the wave oscillates. We don't care so much about these fluctuations. A more useful quantity is the average power. Averaging the power over a wavelength  $\lambda = \frac{2\pi}{k}$  gives

$$\langle P \rangle = P_{\text{avg}} = \frac{k}{2\pi} \int_0^{\frac{2\pi}{k}} dy P(t, y) = \frac{1}{2} Z \omega^2 A_0^2 \quad (39)$$

where  $Z = \rho c_s$  is the impedance for air. The average power is in principle a function of time. In the case of a plane wave, the average power is time-independent.

## 5 Interference

Now we are ready to discuss one of the most important concepts in waves (and perhaps all of physics) constructive and destructive interference..

Suppose we have a speaker emitting sound a frequency  $\omega$ . If the speaker is at  $y=0$  it and we are at large enough distances, the sound will appear as a plane wave

$$A_1 = A_0 \cos(\omega t - ky + \phi_1) \quad (40)$$

Now say we have another speaker directly behind the first speaker producing the same sound at the same volume. We also find a plane wave solution are at large enough distances, the sound will appear as a plane wave

$$A_2 = A_0 \cos(\omega t - ky + \phi_2) \quad (41)$$

We know the frequencies are the same, and  $k$  is the same, but the phases can be different. The total wave is then

$$A_{\text{tot}} = A_1 + A_2 = A_0 \cos(\omega t - ky + \phi_1) + A_0 \cos(\omega t - ky + \phi_2) \quad (42)$$

$$= 2A_0 \cos\left(\omega t - ky + \frac{\phi_1 + \phi_2}{2}\right) \cos\left(\frac{\Delta\phi}{2}\right) \quad (43)$$

where  $\Delta\phi = \phi_1 - \phi_2$  is the phase difference. This last step is just some trigonometry, which you can check by applying the Mathematica command `TrigReduce[]` to Eq. (43).

Thus the average power is now

$$\langle P_2 \rangle = \frac{1}{2} Z \omega^2 (2A_0)^2 \cos^2\left(\frac{\Delta\phi}{2}\right) = 4\langle P_1 \rangle \cos^2\left(\frac{\Delta\phi}{2}\right) \quad (44)$$

where  $\langle P_1 \rangle = \frac{1}{2} Z \omega^2 A_0^2$  is the average power from a single source.

In a generic situation, say where multiple frequencies are produced by different uncorrelated speakers, the phases will have nothing to do with each other. Then, over time the phase difference can change and we should average the phase difference too. Replacing  $\cos^2\left(\frac{\Delta\phi}{2}\right)$  by  $\frac{1}{2}$  then gives

$$\langle P_2 \rangle = 4\langle P_1 \rangle \frac{1}{2} = 2\langle P_1 \rangle \quad (45)$$

So two speakers produce twice the power of one speaker. This makes perfect sense.

Now suppose instead that the two speakers are exactly out of phase, so that  $\Delta\phi = \pi$ , we find

$$\langle P_2 \rangle = 0 \quad (46)$$

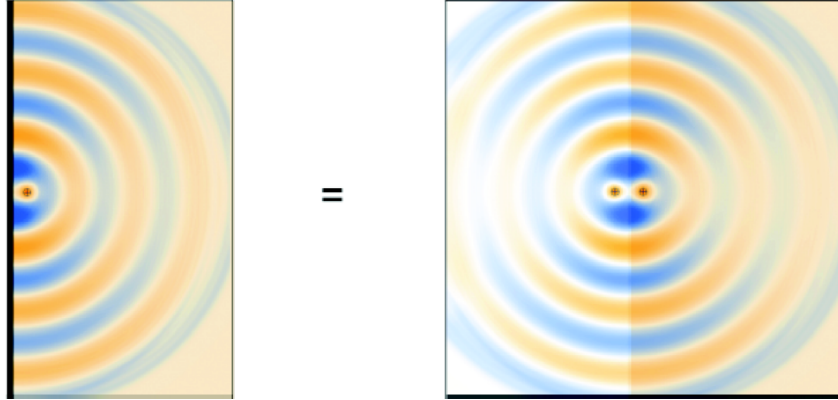
Thus no power is emitted. This is **destructive interference**. Conversely, if  $\Delta\phi = 0$ , then

$$\langle P_2 \rangle = 4\langle P_1 \rangle \quad (47)$$

This is **constructive interference**. Thus two speakers working in phase can produce *four times the power* of a single speaker.

How is this possible? Where is the power coming from? If one looks at how much power is being drawn from the currents driving the speakers, will you find twice as much power is being drawn when the speakers are coherent as when they are incoherent? What is actually happening is that one speaker is pushing down on the other speaker, forcing it to work harder. This is called **source loading**. Thus (in principle) more power is being used by the speakers. However, you won't see this by looking at the power being drawn, since speakers are very inefficient – only 0.01% of the power goes to sound. Instead, the source loading is actually making the speaker more efficient, so that more sound comes out with the same power.

One way to get a produce two coherent speakers is having one speaker a distance  $d$  from a wall:



**Figure 1.** A source near a wall looks like two sources

A wall has infinite impedance ( $Z_2 = \infty$ ), so the sound will reflect off completely. Recall that the amplitude for a sound wave  $A(x, t)$  describes the longitudinal displacement at time  $t$  of molecules whose equilibrium position is at  $x$ . If  $A > 0$  is a displacement to the right, then when the wave reflects it will have  $A < 0$ , since  $R = -1$ . Thus the reflected wave will be displaced to the left. Thus its displacement looks like the mirror image of what the displacement would be if there were no wall. Thus, at a point  $x$  there will be sound coming directly from the source and also sound from the reflection. Say the angular frequency is  $\omega$  and the wavelength is  $\lambda = \frac{2\pi}{k} = 2\pi \frac{v}{\omega}$ . What is the intensity picked up by a microphone a distance  $L$  from the wall?

The easiest way to compute the effect of the reflection is using the method of images: the reflection will act exactly like a source a distance  $d$  on the other side of the wall. Far enough from the source and the wall both the source and its image will produce plane waves. The original source is a distance  $d$  from the wall and so a distance  $L - d$  from the microphone it produces

$$A_1 = A_0 \cos(\omega t - k(L - d)) \quad (48)$$

where we have set the  $\phi = 0$  for the source for simplicity. Then the image source will produce

$$A_2 = A_0 \cos(\omega t - k(L + d)) \quad (49)$$

Thus the phase difference is

$$\Delta\phi = 2kd = 4\pi \frac{d}{\lambda} \quad (50)$$

For  $d \ll \lambda$  (the speaker is close to the wall), then  $\Delta\phi \approx 0$  and we have complete constructive interference. Thus by putting a speaker near a wall we get **four times the power**. This is called a **proximity resonance** or **self-amplification**.

You might have expected there to be twice the power. Since all the power is going into half the space, by conservation of energy this is perfectly logical. Indeed, if the source and the image were incoherent, there would be twice the power. However, we get another factor of 2 from source loading so in fact the power goes up by 4.

It is natural to try to add more proximity resonances. For example, if we have 4 walls or a 30% wedge we would get

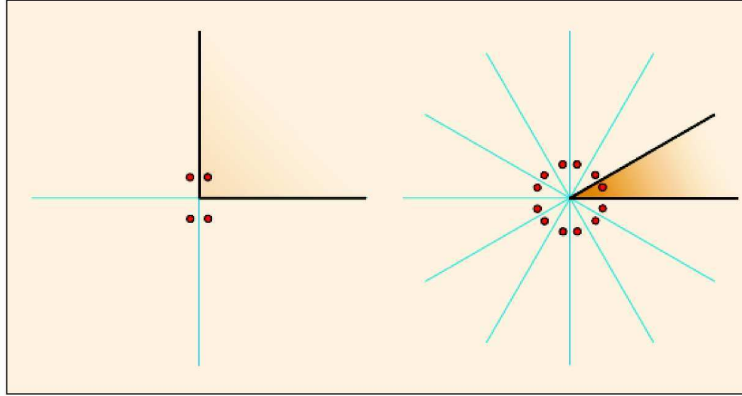


Figure 7.4: Method of images for a single source near a 90° corner and in a 30° wedge. The physical walls are shown as the black lines, and the physical region is shaded. The extra reflection planes for determining the location of the images are shown in light blue. The single source with the walls present or the multiple sources with the walls removed give the same result in the physical region.

This figure and caption are taken from Heller.

With four walls, the area goes to one fourth the area, so naively you would expect a factor of 4 increase in intensity. However, due to source loading the enhancement is a factor of 16. Try standing near a corner and you can hear this yourself. For the 30 degree wedge, source loading gives a factor of 144 enhancement. For a 30 degree wedge in 3 dimensions, the enhancement is around a factor of 200.

The amount of enhancement depends on the frequencies involved. When  $d \sim \lambda$  there is as much destructive interference as constructive interference. Thus, there is no source loading or proximity resonances for high frequencies. One still gets the enhancement in intensity from confining the sound to a smaller volume, but this is not as dramatic as when source loading is relevant. When the size of the wedge is of order the distance to the source, boundary effects become important and one cannot use the plane wave approximations. You can play with numerical solutions with arbitrary configurations using Falstad's ripple.

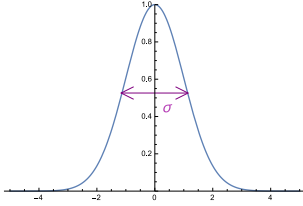




# Lecture 11: Wavepackets and dispersion

## 1 Wave packets

The function

$$g(x) = e^{-\frac{1}{2}\left(\frac{x}{\sigma_x}\right)^2} = \quad (1)$$


is called a **Gaussian**. For a Gaussian, note that  $g(\pm\sigma_x) = \frac{1}{\sqrt{e}}g(0) \approx 0.6g(0)$ , so when  $x = \pm\sigma_x$ , the Gaussian has decreased to about 0.6 of its value at the top. Alternatively, the Gaussian is at half its maximal value at  $x = \pm 1.1\sigma_x$ . Either way,  $\sigma_x$  indicates the width of the Gaussian. The plot above has  $\sigma_x = 1$ . (You may recall that the power of a driven oscillator is given by a Lorentzian function  $l(x) = \frac{\gamma}{x^2 + \gamma^2}$ , which has roughly similar shape to a Gaussian and decays to half of its value at the top at  $x = \pm\gamma$ . Try not to get the functions confused.)

The Fourier transform of the Gaussian is

$$\tilde{g}(k) = \frac{1}{2\pi} \int_{-\infty}^{\infty} dx e^{-ikx} g(x) = \frac{\sigma_x}{\sqrt{2\pi}} e^{-\frac{1}{2}\sigma_x^2 k^2} = \frac{1}{\sqrt{2\pi}\sigma_k} e^{-\frac{1}{2}\left(\frac{k}{\sigma_k}\right)^2}, \quad \text{where } \sigma_k = \frac{1}{\sigma_x} \quad (2)$$

This is also a Gaussian, but with width  $\sigma_k = \frac{1}{\sigma_x}$ . Thus, the narrower the Gaussian is in position space ( $\sigma_x \rightarrow 0$ ), the broader its Fourier transform is ( $\sigma_k \rightarrow \infty$ ), and vice versa

When  $\sigma = \infty$ , the Gaussian is infinitely wide: it takes the same value at all  $x$ . Then  $\tilde{g}(k)$  becomes a  $\delta$ -function at  $k=0$ . That is, to construct a constant, one only needs the infinite wavelength mode (recall  $\lambda = \frac{2\pi}{k}$ ). To construct something narrower than a constant, one needs more and more wavenumbers. To construct a very sharp Gaussian in  $x$  ( $\sigma_x \rightarrow 0$ ) the Fourier transform flattens out: one needs an infinite number of wavenumbers to get infinitely sharp features.

As you know, if we shift the Gaussian  $g(x + x_0)$ , then the Fourier transform rotates by a phase. Conversely, if we shift the Fourier transform, the function rotates by a phase. Even with these extra phases, the Fourier transform of a Gaussian is still a Gaussian:

$$\boxed{f(x) = e^{-\frac{1}{2}\left(\frac{x-x_0}{\sigma_x}\right)^2} e^{ik_c x} \iff \tilde{f}(k) = \frac{\sigma_x}{\sqrt{2\pi}} e^{-\frac{\sigma_x^2}{2}(k-k_c)^2} e^{-ix_0(k-k_c)}} \quad (3)$$

The Gaussian is called a **wavepacket** because of its Fourier transform: it is a packet of waves with frequencies/wavenumbers clustered around a single value  $k_c$  (the subscript “c” is for “carrier”, as we explain below).

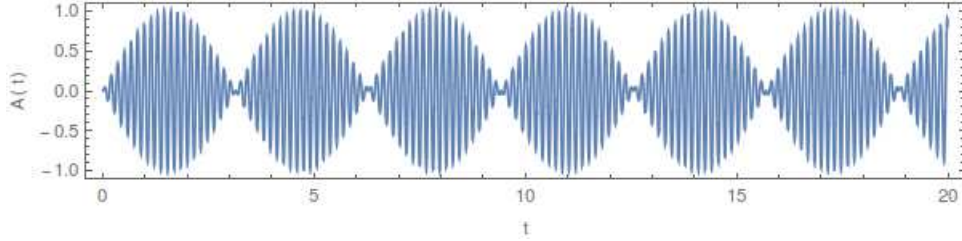
## 2 Amplitude modulation

One of the most important applications of wavepackets is in communication. How do we encode information in waves?

The simplest way is just to play a single note. For example, we produce a simple sin wave,  $\sin(2\pi\nu_c t)$  for some  $\nu_c$ , say 40 Hz for a low  $E$ . If this is all that ever happens, then no information is actually being transferred from one point to another. To transmit a signal, we can start and stop the note periodically. For example, suppose we modulate our note by turning it on and off once a second (think whole notes). Then we would have something like

$$A(t) = f(t) \sin(2\pi\nu_c t) \quad (4)$$

where  $f(t)$  has a frequency of  $\nu_m \sim 1$  Hz. So let's say  $f(t) = \sin(2\pi\nu_m t)$ . Since  $\nu_m \ll \nu_c$ , the curves will look like what we had for beats:

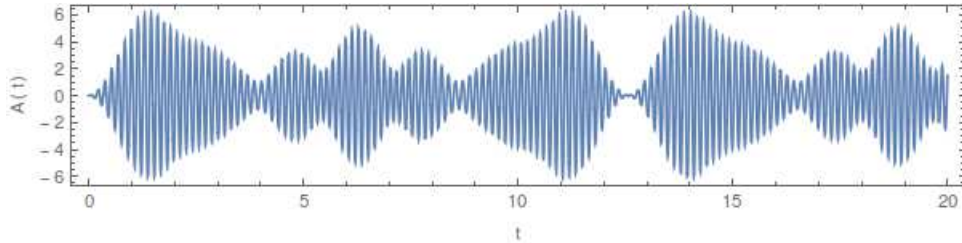


**Figure 1.** A note encoded with high frequency oscillations

We know that since these curves look like the beat curves, they are really the sum of Fourier modes with  $\nu = \nu_c \pm \nu_m$ . In other words  $A(t) = \frac{1}{2}[\cos(39 \text{ Hz} \times 2\pi t) - \cos(41 \text{ Hz} \times 2\pi t)]$ . There is still not any information carried in the signal. But by adding a few more frequencies, we can get something more interesting. For example, consider

$$A(t) = \cos(39t) - \cos(41t) + 0.5 \cos(38t) + 2 \cos(43t) - 2.5 \cos(41.5t) \quad (5)$$

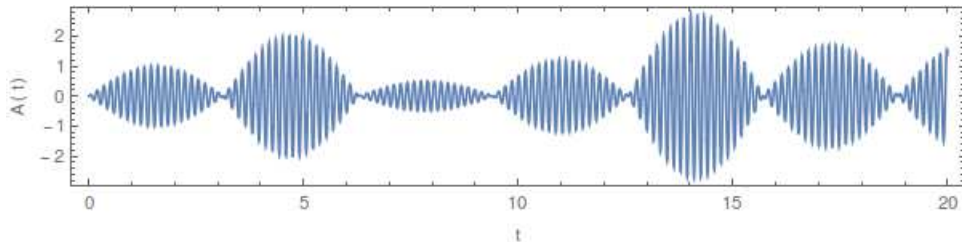
(I didn't write the  $2\pi\text{Hz}$  everywhere to avoid clutter). This looks like this



**Figure 2.** Combining frequencies close to the carrier frequency of 40Hz we can encode information in the signal.

Note that this signal is constructed using only frequencies within 3Hz of the carrier frequency of 40Hz.

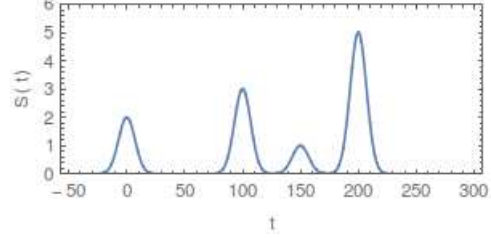
Rather than combining particular frequencies, it's somewhat easier to think about writing the amplitude as in Eq. (4) with  $f(t) = F(t)\sin(2\pi\nu_m t)$  and the function  $F(t)$  having a constant value which changes after each note of the modulated signal. For example, something like this



**Figure 3.** Varying the amplitude at a frequency of  $\nu_m = 1$  Hz using a  $\nu_c = 40$  Hz carrier frequency. This is an amplitude-modulated signal.

Finally, we observe that we can separate the pulses quite cleanly if we construct them with wavepackets, as long as the width of the packets is smaller than the distance between them. For example, we can add Gaussians with different widths and amplitudes:

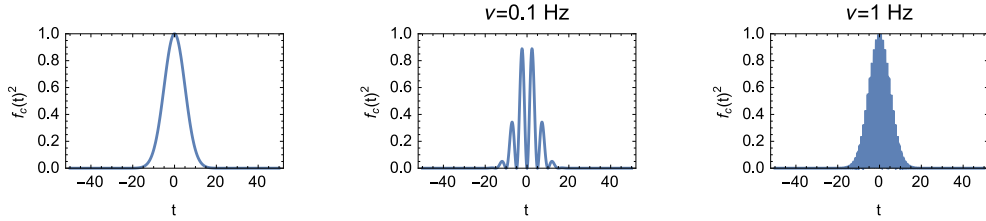
$$S(t) = 2f(t) + 3f(t - 100s) + f(t - 150s) + 5f(t - 200s) =$$



Now, we would like to construct these pulses with a carrier of frequency  $\nu_c$ . Think of this as trying to draw little hills using a pen which wiggles up and down at a rate  $\nu_c$ . The high-frequency pen changes each packet from  $f(t)$  to

$$f_c(t) = e^{-\frac{1}{2}\left(\frac{t}{\sigma_t}\right)^2} \cos(2\pi\nu_c t) = \text{Re}\left[e^{-\frac{1}{2}\left(\frac{t}{\sigma_t}\right)^2} e^{2\pi i \nu_c t}\right] \quad (6)$$

As long as the carrier frequency is larger than the width of the wavepacket,  $\nu_c \gtrsim \frac{1}{\sigma_t}$ , the wiggles in the carrier will be imperceptible and the packet will be faithfully reconstructed. For example, in  $S(t)$  above, the pulses are separated by  $100s$ , so taking  $\sigma_t = 10s$  should do. The corresponding width of the packets is then  $\gamma = \frac{1}{\sigma} = 0.1\text{Hz}$ . The following plots show the amplitude squared, centered around  $\nu_c = 0.1\text{Hz}$  and  $\nu_c = 1\text{Hz}$ .



**Figure 4.** The Gaussian wave-packet (left) with  $\gamma = \frac{1}{\sigma} = 0.1\text{Hz}$  is well approximated by varying the amplitude of a  $\nu_c = 1\text{Hz}$  signal (right). Using  $\nu_c = 0.1\text{Hz}$  (middle) it's not that well constructed.

This example shows that information can be conveyed in  $S(t)$  at the rate of  $\nu_m = \frac{1}{100s} = 0.01\text{Hz}$  using a carrier frequency of  $\nu_c = 1\text{Hz}$ .

More generally, this is how **AM (Amplitude Modulated)** radio works. In radios, the information is conveyed at the **information rate** of  $\nu_m \sim \text{Hz}$  on the **carrier frequency**  $\nu_c$  typically in the  $100\text{MHz}$  range. For cell phones and wireless,  $\text{GHz}$  frequencies are used as carrier frequencies.

In terms of time and frequency, Eq. (3) becomes

$$f(t) = e^{-\frac{1}{2}\left(\frac{t-t_0}{\sigma_t}\right)^2} e^{i\omega_c t} \iff \tilde{f}(\omega) = \frac{\sigma_t}{\sqrt{2\pi}} e^{-\frac{\sigma_t^2}{2}(\omega-\omega_c)^2} e^{-it_0(\omega-\omega_c)} \quad (7)$$

From this, we see that to construct a signal  $f(t)$  with width  $\sigma_t$ , we can use frequencies within a range  $\sigma_\omega = \frac{1}{\sigma_t}$  centered around *any*  $\omega_c$ . The central frequency (the carrier frequency  $\omega_c$ ) can be anything. The key is that enough frequencies around  $\omega_c$  be included. More precisely, we need a band of width  $\sigma_\omega = \frac{1}{\sigma_t}$  to construct pulses of width  $\sigma_t$ . The pulses should be separated by, at minimum,  $\sigma_t$ . Thus the feature which limits how much information can be transmitted is the **bandwidth**. To send more information (smaller distance  $\sim \sigma_t$  between pulses) a larger bandwidth is needed.

### 3 Dispersion relations

An extremely important concept in the study of waves and wave propagation is *dispersion*. Recall the dispersion relation is defined as the relationship between the frequency and the wavenumber:  $\omega(k)$ . For non-dispersive systems, like most of what we've covered so far,  $\omega(k) = vk$  is a linear relation between  $\omega$  and  $k$ . An example of a dispersive system is a set of pendula coupled by springs (see Problem Set 3), where the wave equation is modified to

$$\frac{\partial^2 A(x, t)}{\partial t^2} - \frac{E}{\mu} \frac{\partial^2 A(x, t)}{\partial x^2} + \frac{g}{L} A(x, t) = 0 \quad (8)$$

The dispersion relation can be derived by plugging in  $A(x, t) = A_0 e^{i(kx + \omega t)}$ , leading to the relation  $\omega = \sqrt{\frac{E}{\mu} k^2 + \frac{g}{L}}$ , with  $k = |\vec{k}|$ .

Here is a quick summary of some physical systems and their dispersion relations

- Deep water waves,  $\omega = \sqrt{gk}$ , with  $g = 9.8 \frac{m}{s^2}$  the acceleration due to gravity. Here, the phase and group velocity (see below) are  $v_p = \sqrt{\frac{g\lambda}{2\pi}}$ ,  $v_g = \frac{1}{2}v_p$  and the longer wavelength modes move faster. This regime applies if  $\lambda \gg d$  with  $d$  the depth of the water
- Shallow water waves  $\omega = \sqrt{gd}k$ , where  $d$  is the depth of the water. This is a dispersionless system with  $v_p = v_g = \sqrt{gd}$ .
- Surface waves (capillary waves), like ripples in a pond:  $\omega^2 = k^3 \sigma \rho$ , with  $\sigma$  the surface tension and  $\rho$  density. Thus,  $v_p = \sqrt{\frac{2\pi\sigma}{\rho\lambda}}$ ,  $v_g = \frac{3}{2}v_p$  and shorter wavelength modes move faster. These involve surface tension so can be seen when the disturbance is small enough not to break the water's surface.
- Light propagation in a plasma:  $\omega = \sqrt{\omega_p^2 + ck^2}$ , with  $\omega_p$  the *plasma frequency* and  $c$  the speed of light. This is the same functional form as for the pendula/spring system above.
- Light in a glass  $\omega = \frac{c}{n}k$ .  $n$  is the index of refraction, which can be weakly dependent on wavenumber. In most glass, it is well described by  $n^2 = 1 + \frac{a}{k_0^2 - k^2}$ .

We'll talk about the water waves in Lecture 12 and light waves in later lectures.

### 4 Time evolution of modes: phase velocity

Now we will understand the importance of dispersion relations (and their name) by studying the time-evolution of propagating wavepackets.

To begin, let's think about how to solve the wave equation in a dispersive system with initial condition

$$A(x, t=0) = f(x) \quad (9)$$

Think about setting up a pulse of this form in a medium like a string and then sending down the string. For a non-dispersive wave, with  $\omega(k) = vk$ , the solution is easy

$$A(x, t) = f(x \pm vt) \quad (10)$$

with the sign determined by initial conditions.

Now say we want to solve the pendula/spring wave equation, Eq. (8) with  $A(x, 0) = f(x)$ . So far, we have only solved Eq. (8) for solutions with fixed  $k$ .

$$A_k(x, t) = A_0 e^{i(kx - \sqrt{\frac{E}{\mu} k^2 + \frac{g}{L}} t)} \quad (11)$$

This is indeed of the form  $f(x - vt)$  for  $v = \frac{\sqrt{\frac{E}{\mu} k^2 + \frac{g}{L}}}{k}$ . However, since  $v_p$  depends on  $k$ , this only works for if only one  $k$  is present in the Fourier transform. But if  $A(x, 0)$  is not of the form of a monochromatic (fixed frequency/wavenumber) plane wave, then this solution doesn't apply and we have to think a little harder.

Before thinking harder, we note that a fixed  $k$  solution is possible with any dispersion relation, not just this one. For a dispersion relation  $\omega(k)$  the amplitude  $A_0 \exp\left[ik\left(x - \frac{\omega(k)}{k}t\right)\right]$  is a solution to the corresponding wave equation. We call the speed of this particular solution the phase velocity

$$\text{phase velocity: } v_p(k) = \frac{\omega(k)}{k} \quad (12)$$

Thus  $A(x, t) = A_0 \exp[ik(x - v_p(k)t)]$  will always be a solution.

So what happens to  $A(x, t)$  when  $A(x, 0) \neq e^{ikx}$  for some  $k$ ? The easiest way to solve the wave equation is through Fourier analysis. We know we can write

$$A(x, t=0) = f(x) = \int dk e^{ikx} \tilde{f}(k) \quad (13)$$

where

$$\tilde{f}(k) = \frac{1}{2\pi} \int dx e^{-ikx} f(x) = \frac{1}{2\pi} \int dx e^{-ikx} A(x, 0) \quad (14)$$

Eq. (13) writes the initial condition as a sum of plane wave (fixed  $k$ ) modes. Then, since we know that each mode evolves by replacing  $x \rightarrow x - v_p(k)t$ , we have

$$\boxed{A(x, t) = \int dk e^{ik(x - v_p(k)t)} \tilde{f}(k) = \int dk e^{i(kx - \omega(k)t)} \tilde{f}(k)} \quad (15)$$

It's that simple. This is the exact solution to Eq. (8) with initial condition  $A(x, 0) = f(x)$ .

Let us check that Eq. (15) satisfies Eq. (8) with  $\omega(k) = \sqrt{\frac{E}{\mu}k^2 + \frac{g}{L}}$ . Plugging in we get

$$\left[ \frac{\partial^2}{\partial t^2} - \frac{E}{\mu} \frac{\partial^2}{\partial x^2} + \frac{g}{L} \right] A(x, t) = \left[ \frac{\partial^2}{\partial t^2} - \frac{E}{\mu} \frac{\partial^2}{\partial x^2} + \frac{g}{L} \right] \int dk e^{i(kx - \sqrt{\frac{E}{\mu}k^2 + \frac{g}{L}}t)} \tilde{f}(k) \quad (16)$$

$$= \left[ -\left( \frac{T}{\mu}k^2 + \frac{g}{L} \right) + \frac{T}{\mu}k^2 + \frac{g}{L} \right] \int dk e^{i(kx - \sqrt{\frac{E}{\mu}k^2 + \frac{g}{L}}t)} \tilde{f}(k) \quad (17)$$

$$= 0 \quad (18)$$

The boundary condition  $A(x, 0) = f(x)$  is also satisfied.

Another check is that in the special case of a dispersionless medium, where  $\omega(k) = vk$  and so  $v_p(k) = v$  constant, the solution is exactly what we expect:

$$A(x, t) = \int e^{ik(x - vt)} \tilde{f}(k) dk = f(x - vt) \quad (19)$$

which we already knew.

## 5 Time evolution of signals: group velocity

In this section, we take  $\omega(k)$  to be arbitrary and take the initial signal shape to be our beautiful Gaussian wavepacket constricted with a carrier wave of wavenumber  $k_c$ . So

$$f(x) = e^{-\frac{1}{2}\left(\frac{x-x_0}{\sigma_x}\right)^2} e^{ik_c x} \quad (20)$$

where  $k_0 = k_c$  is the carrier wavenumber. Here, we let the signal be complex to efficiently encode phase information. One can always take the real part at the end, as we have done before. We again want to solve the general wave equation with dispersion relation  $\omega(k)$  for  $A(x, t)$  with initial condition  $A(x, 0) = f(x)$ .

The Fourier transform of this packet is

$$\tilde{f}(k) = \frac{\sigma_x}{2\sqrt{\pi}} e^{-\frac{\sigma_x^2}{2}(k-k_c)^2} e^{ix_0(k-k_c)} \quad (21)$$

as in Eq. (3). In Fourier space, the time evolution is easy to compute:

$$A(x, t) = \int dk e^{i(kx - \omega(k)t)} \tilde{f}(k) \quad (22)$$

As noted above, it is impossible to solve this in general. But since in our case  $\tilde{f}(k)$  is exponentially suppressed away from  $k = k_c$ , we can Taylor expand the dispersion relation

$$\omega(k) = \omega(k_c) + (k - k_c)\omega'(k_c) + \cdots \quad (23)$$

$$= k_c v_p + (k - k_c)v_g + \cdots \quad (24)$$

$v_p = v_p(k_c)$  is the phase velocity at  $k_c$  and  $v_g = v_g(k_c)$  is called the group velocity

$$\text{group velocity} \quad v_g(k) = \frac{d\omega(k)}{dk} \quad (25)$$

In general, both the phase and group velocities depend on  $k$ . Here, because of the Taylor expansion, we are only interested in the special value  $v_g = v_g(k_c)$ .

If we truncate the Taylor expansion to order  $(k - k_c)$ , then the solution for  $A(x, t)$  is:

$$A(x, t) = \int dk e^{i(kx - k_c v_p t - (k - k_c)v_g t)} \tilde{f}(k) \quad (26)$$

$$= e^{-ik_c t(v_g - v_p)} \int dk e^{ik(x - v_g t)} \tilde{f}(k) \quad (27)$$

$$= e^{-ik_c t(v_g - v_p)} f(x - v_g t) \quad (28)$$

Thus we have found that the wave-packet moves at the velocity  $v_g$ .

Note that for a non-dispersive wave  $v_p = v_g$  and we get back our original solution. Also note that in deriving this, we didn't need to use the exact form of the wavepacket, just that it was exponentially localized around  $k_c$ .

Stating our results in terms of time dependence and frequency, we have found that

- A pulse can be constructed with a group of wavenumbers in a band  $k_c - \sigma_k < k < k_c + \sigma_k$  or equivalently with a group of frequencies in a band  $\nu_c - \sigma_\nu < \nu < \nu_c + \sigma_\nu$ .
- To send a pulse which lasts  $\sigma_t$  seconds using a carrier frequency  $\nu_c$ , one needs frequencies within  $\sigma_\nu = \frac{1}{\sigma_t}$  of  $\nu_c$ .
- The pulse travels with the group velocity  $v_g = \left. \frac{d\omega}{dk} \right|_{k=k_c}$  evaluated at the carrier wavenumber/frequency.

Note that because  $\sigma_k \ll k_c$  ( $\sigma_\nu \ll \nu_c$ ), the group velocity is roughly constant for all of the relevant wavenumbers,  $k_c - \Delta k < k < k_c + \Delta k$ . But it may be very different from the phase velocity. For example, if  $\omega(k) = 5k^4$ , then at  $k_c = 100$ ,  $v_p = 5 \times 10^8$  while  $v_g = 20k^3 = 2 \times 10^7$ . Again, for non-dispersive media,  $v_g = v_p$ . We will contrast group and phase velocity more in the next lecture when we have some concrete examples of dispersive systems.

## 6 Dispersion

Now we come to where dispersion relations got their name.

We just saw that to the first approximation, a wave-packet moves with velocity  $v_g$ . Of course, in the first order approximation in the Taylor expansion, the dispersion relation might as well be linear (non-dispersive). So let's add the second term to see the dispersion. Then

$$\omega(k) = \omega(k_c) + (k - k_c)\omega'(k_c) + \frac{1}{2}(k - k_c)^2\omega''(k_c) + \cdots \quad (29)$$

$$= k_c v_p + (k - k_c)v_g + \frac{1}{2}(k - k_c)^2\Gamma + \cdots \quad (30)$$

where  $\Gamma = \omega''(k_c)$  is a new parameter. Note that if the wave is non-dispersive, so  $\omega(k) = vk$ , then  $\omega_p = \omega_g$  and  $\Gamma = 0$ .

With this expansion, let's go back to our Gaussian. We start with

$$A(x, t=0) = f(x) = e^{-\frac{1}{2\sigma^2}(x-x_0)^2} e^{ik_c x} \quad (31)$$

Then,

$$\tilde{f}(k) = \frac{\sigma}{2\sqrt{\pi}} e^{-\frac{\sigma^2}{2}(k-k_c)^2} e^{ix_0(k-k_c)} \quad (32)$$

So

$$A(x, t) = \frac{\sigma}{2\sqrt{\pi}} \int e^{i(kx - [k_c v_p + (k-k_c)v_g + \frac{1}{2}(k-k_c)^2 \Gamma]t)} e^{-\frac{\sigma^2}{2}(k-k_c)^2} e^{ix_0(k-k_c)} dk \quad (33)$$

If you stare at the exponent, you will see that it is still quadratic in  $k$  – still a Gaussian – so in this special case we can actually perform the inverse Fourier transform. And of course, we will get a new Gaussian. The result is

$$A(x, t) = \exp\left[-\frac{1}{2}\left(\frac{x - (x_0 + v_g t)}{\sqrt{\sigma_x^2 - i\Gamma t}}\right)^2\right] e^{ik_c x} e^{-ik_c t(v_g - v_p)} \quad (34)$$

This is the exact solution for the time dependence if  $\omega(k) = k_c v_p + (k - k_c)v_g + \frac{1}{2}(k - k_c)^2 \Gamma$  exactly. It is helpful to also pull the  $i$  out of the denominator, writing the solution

$$A(x, t) = \exp\left[-\frac{1}{2}\left(\frac{x - (x_0 + v_g t)}{\sigma(t)}\right)^2\right] e^{i\phi(x, t)} \quad (35)$$

where

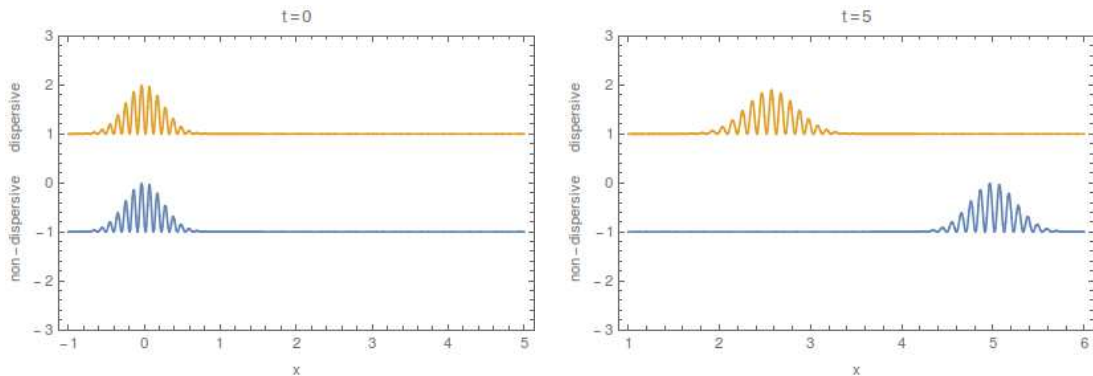
$$\sigma(t) = \sigma_x \sqrt{1 + \frac{\Gamma^2}{\sigma_x^4} t^2} \quad (36)$$

and

$$\phi(x, t) = k_c x - k_c t(v_g - v_p) - \frac{t\Gamma}{t^2 \Gamma^2 + \sigma_x^4} \quad (37)$$

How do we interpret this solution? It has a magnitude and a phase. The phase just causes the real part to oscillate between  $-1$  and  $1$ , which is not that interesting. So let's concentrate on the magnitude. The magnitude of a Gaussian only has three parameters, its overall normalization, its center and its width. We have written Eq. (35) in a way so that it is easy to read of that at time  $t$  the packet is centered at  $(x_0 + v_g t)$ . This is consistent with what we found above: the center of the Gaussian moves with the group velocity. We can also read off from Eq. (35) that the width at time  $t$  is given by the function  $\sigma(t)$  in Eq. (36). Notice that the **width is increasing with time**. That is, the wave-packet is broadening. This is why we call it a dispersion relation. Recall that a non-dispersive wave has  $\Gamma = 0$ , so with non-dispersive dispersion relations, wavepackets don't disperse.

Here's a comparison of a nondispersive pulse, with  $v = 1$  to one with dispersion relation  $\omega(k) = \sqrt{k^2 + 50^2}$ . We construct a wavepacket of width  $\sigma_x = 0.5$  with a carrier wavenumber of  $k_c = 30$ .



**Figure 5.** Pulses at  $t=0$  and  $t=10$ . Dispersive packet is on top. Note that the dispersive one is moving at  $v=0.5$  and the non-dispersive one at  $v=1$ .

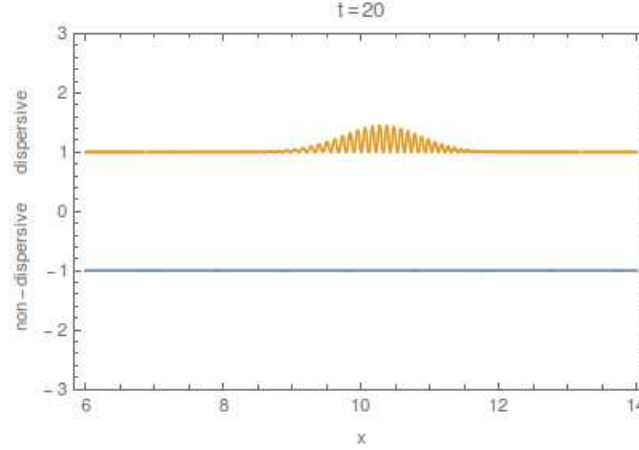


For the non-dispersive pulse, the phase and group velocity are  $v_p = v_g = 1$ . Thus, after 5 seconds it is at  $x = 5$ , consistent with the figure. For the dispersive pulse, the phase and group velocities are

$$v_p = \frac{\omega(k_c)}{k_c} = 1.94, \quad v_g = \omega'(k_c) = 0.51 \quad (38)$$

One can see from the figure that at  $t = 5$ , the dispersive packet has gone half as far as the non-dispersive one, which is consistent with it traveling at the group velocity of  $v_g = 0.51$ .

At longer times, you can really see the pulse flatten out.



**Figure 6.** At  $t=20$ , the dispersive packet is significant broadened. The non-dispersive pulses is off the plot, near  $x = 20$ . It has the same shape it did originally.

Dispersion in optical media are critical to modern optics and to telecommunications. For example, high speed internet and long distance telephone communication are now done through fiber optic cables. Fiber optic cable contains a glass core surrounded by a lower index-of-refraction cladding. This allows light to be transported along the cable via total internal reflection. A key figure in telecommunication is the rate at which data can be communicated. In fiber optic telecommunications, information is transmitted via optical wavepackets in a glass fiber. Due to dispersion in the glass, pulses too close together begin to overlap, destroying the information. This sets a fundamental limit to the speed internet communications. Luckily, silica-based glass has very low dispersion (and absorption) in the near IR region (1.3-1.5 micron). This frequency band, now known as the telecomm band, has seen extensive technological development in the last 20 years due to its use for fiber optic communication. Fiber optics will be revisited when we discuss light.

By the way, in quantum mechanics, an electron is often effectively treated as a wavepacket. We will see that in Lecture 20 that non-relativistic dispersion relation for an electron is  $\omega(k) = \frac{\hbar}{2m} k^2$ . So  $v_g = 2v_p = \hbar \frac{k}{m} = \frac{p}{m}$  with  $p = \hbar k$  the momentum and  $\Gamma = \frac{\hbar}{m}$  the width. Thus the width becomes  $\sigma(t) = \sigma_x \sqrt{1 + \left( \frac{\lambda_c c t}{2\pi \sigma_x^2} \right)^2}$  with  $\lambda_c = \frac{h}{mc} = 2.42 \times 10^{-12} m$  called the Compton wavelength of the electron. The phase becomes  $\phi(x, t) = k_c \left( x - \frac{p}{m} \right) + \frac{\lambda_c c t}{\lambda_c^2 c^2 t^2 + \sigma_x^4}$ . At late times  $\sigma(t) = \frac{\lambda_c}{2\pi \sigma_x} c t$  and  $\phi(x, t) = k_c \left( x - \frac{p}{m} \right)$ . Thus the center moves with the velocity  $\frac{p}{m}$ , as expected, and the width grows very rapidly: at the speed of light for an electron localized to within its Compton wavelength. The more you try to pin down an electron (smaller  $\sigma_x$ ), the faster the wavepacket grows!



Matthew Schwartz

## Lecture 12: Wave phenomena

This lecture is from Muller's new book *Physics and Technology for Future Presidents*. The attached chapter on Waves is freely available on the web

<http://muller.lbl.gov/teaching/Physics10/PffP.html>

Feel free to browse other chapters, and buy the book if you like it.

Much of what is in here will be review for you, but it may be nice to hear about things from a different perspective. Fewer equations, more stories.

Don't worry about anything past page 7-26. We'll discuss interference, diffraction and the Doppler effect later on. We already discussed music. You can skim through these final pages though, as you might appreciate Muller's perspective.

# 7. Waves

including UFOs, earthquakes, and music

## Two strange but true stories

The following two anecdotes, *Flying Saucers* and *Rescuing Pilots*, are actually closely related, as you will see later in this chapter. They both will lead us into the physics of waves.

### Flying saucers crash near Roswell, New Mexico

In 1947, devices that the U.S. government called "flying disks" crashed in the desert of New Mexico. The debris was collected by a team from the nearby Roswell Army Air Base, which was one of the most highly classified locations in the United States. The government put out a press release announcing that flying disks had crashed, and the story made headlines in the respected local newspaper, The Roswell Daily Record. Take a moment to look at the headlines for July 8, 1947:

Leased Wire  
Associated Press

**Roswell Daily Record**

RECORD PHONES  
Business Office 2288  
News Department  
2287

CL. 47, NUMBER 39, ESTABLISHED 1888

ROSWELL, NEW MEXICO, TUESDAY, JULY 8, 1947

14 PAGES

**Grand**  
Grand Tower, Ill., but while the manager of the movie theater sweeps out the water that has entered the lobby these youngsters are standing in line for tickets for the night performance. (AP Wirephoto).

**Claims Army Is Stacking Courts Martial**  
Indiana Senator Lays Protest Before Patterson  
Washington, July 8 (AP)—Indiana Senator Richard J. Roth, today, charged that the army is stacking courts martial in order to avoid a full investigation of the crash of a flying saucer near Roswell, N. M., last night.

**House Passes Tax Slash by Large Margin**  
Defeat Amendment by Demos to Remove Many from Rolls  
Washington, July 8 (AP)—The House today passed a bill to cut income taxes by 10 percent, 340-100.

**Security Council Paves Way to Talks On Arms Reductions**  
Lake Success, July 8 (AP)—The United Nations Security Council today approved a resolution calling for a general conference to discuss arms reduction.

**No Details of Flying Disk Are Revealed**  
Roswell Hardware Man and Wife Report Disk Seen  
The Roswell Hardware Man and Wife today reported that they had seen a flying saucer.

**Ex-King Carol Weds Mme. Lupescu**  
Former King Carol of Romania and Mme. Elena Lupescu today wedded in Bucharest.

**Some of Soviet Satellites May Attend Paris Meeting**  
Paris, July 8 (AP)—Indications pointed today that at least some of the satellite states of the Soviet Union would attend the Paris conference on the Marshall plan.

**Roswellians Have Differing Opinions On Flying Saucers**  
Roswell, July 8 (AP)—Indications pointed today that at least some of the Roswellians who saw the flying saucer were not convinced that it was a flying saucer.

**American League Wins All-Star Game**  
Chicago, July 8 (AP)—The American League today won the All-Star game, 4-3, over the National League.

**RAAF Captures Flying Saucer On Ranch in Roswell Region**

House Passes Tax Slash by Large Margin

Security Council Paves Way to Talks On Arms Reductions

No Details of Flying Disk Are Revealed

Ex-King Carol Weds Mme. Lupescu

Some of Soviet Satellites May Attend Paris Meeting

Roswellians Have Differing Opinions On Flying Saucers

American League Wins All-Star Game

Serious newspaper headlines from the respected Roswell Daily Record, RAAF stands for "Roswell Army Air Force".

The next day, the U.S. government retracted the press release, and said their original announcement was mistaken. There were no flying disks, they claimed. It was only a weather balloon that had crashed. Anybody who had

seen the debris knew it wasn't a weather balloon. It was far too large, and it appeared to be made from some exotic materials. In fact, the object that crashed was *not* a weather balloon. The government was lying, in order to protect a highly classified program. And most people could tell that the government was lying.

The story I have just related sounds like a fantasy story from a supermarket tabloid--or maybe like the ravings of an anti-government nut. But I assure you, everything I said is true. The story of the events of Roswell, New Mexico is fascinating, and not widely known, since many of the facts were classified until recently. In this chapter I'll fill in the details so that the Roswell story makes sense.

Incidentally, if you are unfamiliar with the name Roswell, that means you have not watched the TV program "The X Files" or read any of the other voluminous literature about flying saucers and UFOs. Try doing an Internet search on Roswell in 1947 and see what you find. Be prepared to be astonished.

Now for the second anecdote.

## **Rescuing Pilots in World War II**

The true story of the flying disks began with an ingenious invention made by the physicist Maurice Ewing near the end of World War II. His invention involved small objects called "**sofar**" spheres that could be placed in the emergency kits of pilots flying over the Pacific Ocean. If a pilot was shot down, but he managed to inflate and get on to a life raft, then he was instructed to take one of these spheres and drop it into the water. If he wasn't rescued within 24 hours, then he should drop another.

What was in these miraculous spheres? If the enemy had captured one and opened it up, they would have found that the spheres were hollow with nothing inside. How could hollow spheres lead to rescue? How did they work?

Here's the answer to the **sofar** question: Ewing had been studying the ocean, and he was particularly interested in the way that sound travels in water. He knew that the temperature of the water got colder as it got deeper--and that should make sound travel slower. But as you go deeper, the pressure gets stronger, and that should make the sound travel faster. The two effects don't cancel. When he studied it in detail, he concluded that the sound velocity would vary with depth. His most interesting conclusion was that at a depth of about 1 km, the sound travels slower than at any other depth. As we will discuss later, this implies the existence of a "sound channel" at this depth, a layer that tends to concentrate and focus sound and keep it from escaping to other depths. Ewing did some experiments off the coast of New Jersey and verified that this sound channel existed, just as he had predicted.

The **sofar** spheres were hollow and heavier than an equal volume of water. They sank but were strong enough to hold off water pressure until they reached the depth of the sound channel. At that depth the sphere suddenly collapsed with a bang. That sent out a pulse of sound that could be heard thousands of kilometers away. From these sounds, the Navy could figure out the approximate location of the downed pilot, and send out a rescue team.

It turns out (this wasn't known back then) that Ewing's little spheres used the same phenomena that whales use to communicate with other

whales: the focusing of sound in the sound channel. We'll discuss this shortly.

At the end of World War II, the same Maurice Ewing proposed a second project based on the same idea. This project was eventually given the name **Project Mogul**. It used "flying disks" for a highly classified purpose: to detect nuclear explosions. It made use of a sound channel in the atmosphere. But the flying disks crashed in Roswell, New Mexico in 1947, made headlines, and became part of a modern legend.

To explain these stories, we have to get into the physics of sound. And to understand sound, we have to talk about waves.

## Waves

All waves are named after water waves. Think for a moment about how strange water waves are. Wind pushes up a pile of water, and the pile creates a wave. The wave moves and keeps on moving, carrying energy far from the place where the wave was created. Waves at the coast are frequently an indicator of a distant storm. But the water from that distant storm didn't move very far, just the wave. The wind pushed the water and the water pushed other water and the energy traveled for thousands of miles, even though the water only moved a few feet.

You can make waves on a rope or with a toy called a *slinky*. (If you've never played with a slinky, you should go to a toy store as soon as possible and buy one.) Take a long rope or a slinky, stretch it across a room, shake one end, and watch the wave move all the way to the other end and then bounce back. (Water waves, when they hit a cliff, also bounce.) The rope jiggles, but no part of it moves very far. Yet the wave does travel, and with remarkable speed.

Sound is also a wave. When your vocal cords vibrate they shake the air. The air doesn't move very far, but the shaking does. The shaking moves as far as the ear can hear and further. The initial shaking air around your vocal cords makes the air nearby shake also, and so on. If the shaking reaches someone else, then it causes his eardrums to shake, which sends signals to his brain and causes him to hear you.

For a nice animation of a sound wave, showing how the molecules bounce back and forth but create a wave that moves forward only, see <http://www.kettering.edu/~drussell/Demos/waves/wavemotion.html>

If the sound wave hits a wall, it bounces. That's what gives rise to echoes. Sound waves bounce just like water waves and rope waves.

A remarkable thing about all these kinds of waves is that the shaking leaves the location where it started. Shake some air and you create a sound, but the sound doesn't stay around. A wave is a way of transporting energy long distances without actually transporting matter. It is also a good way to send a signal.

It turns out that light, radio, and TV signals also consist of waves. We'll get to that in the next chapter. What is waving for these? The traditional answer is "nothing" but that is really misleading. A much better answer is that there is a "field" that is shaking – the electric and magnetic fields.

Another correct answer is that “the vacuum” is what is shaking. We’ll discuss this further in the chapter on quantum mechanics.<sup>1</sup>

## Wave packets

Waves can be long with many vibrations, as when you hum, or they can be short, as in a shout. We call such short waves “wave packets.” You may have noticed water waves often travel in packets. Splash a rock into a pool and you’ll see a bunch of waves moving out, forming a ring that contains several up and down oscillations. That’s a packet. A shout contains many oscillations of the air, but these oscillations are confined to a relatively small region. So that too is a wave packet.

Now think about this: short waves act in a way very similar to particles. They move and they bounce. They carry energy. If the packet were extremely short, maybe you wouldn’t notice that it was really a wave. Maybe you would think it was a small particle.

In fact, the theory of quantum mechanics is really a fancy name for the theory that all particles are really little packets of waves. The packets for an electron and proton are so small that we don’t normally see them. What is waving in an electron? We think it is the same thing that is waving for light: the vacuum.

So when you are studying sound, water, and earthquakes, you are really learning the properties of waves. That will be most of what you need to understand quantum mechanics.

## Sound

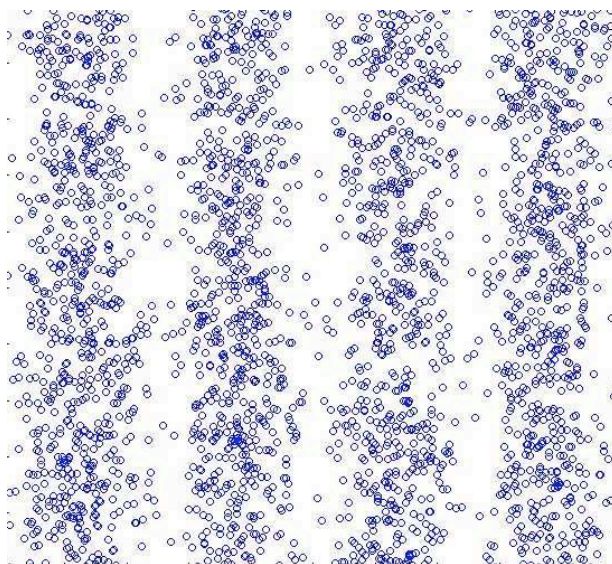
Sound in air results when air is suddenly compressed, for example by a moving surface (such as a vibrating vocal cord or bell). The compression pushes against adjacent air, and that pushes against the air in front of it, and so on. The amazing thing about sound is that the disturbance travels, and the shaking of the original air stops. The energy is carried away very effectively.

Sound is generated in air when something compresses it in a local region. This could be the vibrating of vocal cords, a violin string, or a bell. The compressed air expands, and compresses the air next to it. The air never moves very far, but the compression is passed on from one region to the next. This is depicted in the following diagram. Each little circle represents a

---

<sup>1</sup> Here is a brief summary of the answer: when it was discovered that light is a wave, physicists didn’t know what was waving, but they gave it a name: the “Aether.” (I spell it this way to distinguish it from the chemical “ether” which is totally different.) Most modern physicists believe that the Aether was shown not to exist, but that isn’t true. The distinguished theoretical physicist Eyvind Wichmann pointed out (in a class I took from him at Berkeley) that the Aether was only shown to be invariant under the laws of special relativity, and therefore was unnecessary. But then quantum mechanics started giving it properties: it can be polarized, and it carries dark energy. Wichmann says that the Aether never went away from physics; it was made more complex, and simply was reborn with a new name: the vacuum.

molecule. The wave consists of compression and expansion of regions of the gas.



Each molecule shakes back and forth, and doesn't travel very far. But the waves travel forward. Look at the diagram, and imagine that you are looking at a series of water waves from an airplane. But the waves in sound don't come from up and down motion, but from compression and dilation. When these compressions reach your eardrum, they make it vibrate. Those vibrations are then passed on through the rest of your ear to nerves and then to the brain, where the vibrations are interpreted as sound.

To understand this, it is easiest to watch a movie. A very nice one is posted at [www.kettering.edu/~drussell/Demos/waves/wavemotion.html](http://www.kettering.edu/~drussell/Demos/waves/wavemotion.html). A wave is moving from the left to the right. But if you watch one molecule, you'll see that it is shaking back and forth, and never travels very far. It bangs into a nearby molecule, and transfers its energy.

That is the key aspect of waves. No individual molecule travels very far, but the energy is transferred. They pass on the energy, from one to the next. It is the energy that travels long distances, not the particles. Waves are means for sending energy without sending matter.

Sound waves can travel in rock, water, or metal. All those materials compress slightly, and this compression travels and carries away the energy. If you hit a hammer on a railroad rail, then the metal rail is momentarily distorted and the distortion travels down the rail. If someone puts her ear to the rail a mile away, she will hear the sound. The best way to hear the sound is to put your head against the rail. The vibration in the rail will make your skull vibrate, and this will make the nerves in your ear respond—even if none of the sound is actually in air.

Because steel is so stiff, it turns out that sound travels 18 times faster in steel than in air. In air, sound takes 5 s to go 1 mi; in steel, sound will go that same distance in less than 1/3 s. In the olden days, when people lived near railroad tracks, they could listen to the track to hear if a train was coming, and they could even estimate the distance to the train by the loudness of the sound.

For sound to travel, the molecules of air have to hit other molecules of air. That's why the speed of sound is approximately equal to the speed of



molecules. We discussed this fact in Chapter 2. But in steel, the molecules are already touching each other. That's why sound in steel can move much faster than the thermal velocity of the atoms in the steel.

Sound travels in any material that is springy, i.e. which returns to its original shape when suddenly compressed and then released. The faster it springs back, the faster the wave moves. The speed of sound in water is about 1 mi/s, but it varies slightly depending on the temperature and depth of the water.

Note that a sound wave in water is a different kind of wave than the water wave that moves on the surface. In water, sound travels *under* the surface, in the bulk of the water. It consists of a compression of the water. Water waves on the surface are not from compression, but from movement of the water up and down, changing the shape of the surface. So although they are both in water, they are really very different kinds of waves. You can see surface waves easily. You usually cannot see sound waves. Surface waves are slow and big. Sound waves are microscopic and fast.

The speed of sound in air doesn't depend on how hard you push, that is, on how intense the sound is! No matter how loud you shout, the sound doesn't get there any faster. That's surprising, isn't it?

Why is that true? Remember, at least for air, the speed of sound is approximately the speed of molecules. The signal has to go from one molecule to the next, and it can't do that until the air molecule moves from one location to another. (The added motion from the sound vibration is actually very small compared to the thermal motion of the molecules.) When you push on the air, you don't speed up the molecules very much; you just push them closer to each other.

But the speed of sound does depend on the temperature of the air. That's because the speed depends on the velocity of the air molecules, and when air is warmer, the velocity is greater.

The table below gives the speed of sound in several materials:

material and temperature	speed of sound
air at 0 C = 32 F	331 m/s = 1 mi for every 5 s
air at 20 C = 68 F	343 m/s
water at 0 C	1402 m/s = 1.4 km/s
water at 20C	1482 m/s = almost 1 mi/s
Steel	5790 m/s = 3.6 mi/s
Granite	5800 m/s

There is no need to memorize this table. But you should remember that sound moves faster in solids and liquids than in air. And you should know that the speed of sound in air is about one mile every five seconds.

Sound traveling in rock gives us very interesting information about distant earthquakes. We'll come back to that later in this chapter. Observations of the surface of the sun show sound waves arriving from the other side, traveling right through the middle of the Sun. Much of our knowledge of the interior of the Sun comes from the study of these waves. (We detect them by sensitive measurements of the surface of the Sun.) Sound has been detected traveling through the Moon, created by meteorites

hitting the opposite side. On the Moon we use instruments that were left behind by the Apollo astronauts.

There is no sound in space because there is nothing to shake. A famous tag line from the science fiction movie *Alien* (1979) is, "In space, nobody can hear you scream." Astronauts on the moon had to talk to each other using radios. Science fiction movies that show rockets roaring by are not giving the sound that you would hear if you were watching from a distance--since there would be no sound.<sup>2</sup>

## Transverse and longitudinal waves

When you shake the end of a rope, the wave travels down its length, from one end to the other. However, the shaking is sideways, i.e. the rope vibrates sideways even though the direction that the wave is moving is along the rope. This kind of wave is called a *transverse* wave. In a transverse wave, the motion of the particles is along a line that is perpendicular to the direction the wave is moving.

For an illustration of a transverse wave, go back to:

<http://www.kettering.edu/~drussell/Demos/waves/wavemotion.html>

and look at the second illustration on that page.

A sound wave is different. The vibration of the air molecules is back and forth, in the same direction that the wave is moving. This kind of compressional wave is called a *longitudinal* wave. In such a wave, the motion and direction of the wave are both along the same line.

This may seem peculiar, but water waves are even stranger.

## Water surface waves

Water waves (the term we will use when we mean the ordinary surface water waves--as opposed to water sound waves) gave all waves their name. If you are swimming or floating and a water wave passes by, you move slightly back and forth as well as up and down. It is worthwhile to go swimming in the ocean just to sense this. In fact, for most water waves, the sideways motion is just as big as the up and down, and you wind up moving in a circle! But when the wave is past, you and the water around you are left in the same place. The wave, and the energy it carries, passes by you.

For a nice illustration of the motion of particles in a water wave, take a look again at

<http://www.kettering.edu/~drussell/Demos/waves/wavemotion.html>

but this time scroll down to the third animation. Look at one particle, maybe one of the blue ones, and watch how that particle moves. Does it move in a circle?

When there is a series of waves following each other, we call that a wave packet. The distance between the crests (the high points of the waves) is called the **wavelength**. Waves with different wavelengths travel at very different speeds. Those with a short wavelength go slower, and those with a

---

<sup>2</sup> To enjoy the movie, I always assume that the microphone is located on the spacecraft, so although we are watching the rocket pass, we are hearing sound as if we were on the rocket.



long wavelength go faster. In deep water (when the depth is greater than the wavelength), the equation is as follows:<sup>3</sup>

$$v \approx \sqrt{L}$$

In this equation,  $v$  is the velocity in meters per second (m/s), and  $L$  is the wavelength in meters (m), and the squiggly equals sign  $\approx$  means “approximately equal to.” So, for example, if the wavelength (distance between crests) is  $L = 1$  m, then the velocity is about  $v = 1$  m/s. If the wavelength is 9 m, the velocity is 3 m/s. Does that agree with your image of ocean waves? Next time you swim in the ocean, check to verify that long waves move faster.

That equation is remarkably simple, but it is correct only for deep water, that is, for water that is much deeper than a wavelength.

### Shallow water waves

When the water is “shallow” (the depth  $D$  is much less than the wavelength  $L$ ) then the equation changes to

$$v = 3.13\sqrt{D} \\ \approx \pi\sqrt{D}$$

where  $D$  is the depth in meters.<sup>4</sup> Note that all shallow water waves travel at the same velocity, determined only by the depth of the water, regardless of the wave’s wavelength. The speed of shallow water waves depends only on the depth of the water. This might match your experience when you surf on relatively long waves in shallow water.

If the wavelength is very long, then we have to regard even the deep ocean as shallow. This is often the case for tsunamis.

### Tsunamis (tidal waves)

A tsunami is a giant wave that hits the coast and washes far up on the shore, often destroying buildings that are within a few hundred meters of the beach. Tsunamis were traditionally called tidal waves, but a few decades ago scientists (and newspapers) decided to adopt the Japanese word, and now it is more commonly used.

Underwater earthquakes and landslides often generate tsunamis. These waves usually have a very high velocity and a very long wavelength. In the deep ocean, they may have a very low amplitude, so they can travel right under a ship without anyone on board even noticing. But as they approach land, they are slowed down, and the energy is spread out over a smaller depth of water. As a result, the height of the wave rises. The rise can be enormous, and that is what causes the damage near the coast.

---

<sup>3</sup> *For the physics major:* The standard physics equation for deep water waves is  $v = \sqrt{gL/(2\pi)}$ , where  $g = 9.8 \text{ m/s}^2$  is the acceleration of gravity (from Chapter 3). Putting in  $g = 9.8$ , gives  $v \approx 1.2 \sqrt{L} \approx \sqrt{L}$ .

<sup>4</sup> The second equation is only approximate. I wrote it using the symbol  $\pi$  to make it easier to remember, even though you don’t have to remember it.

In Pacific islands (such as Hawaii) you'll see sirens mounted on poles near the beaches. If an earthquake is generated within a few thousand miles, these sirens will be sounded to warn the residents to evacuate. A tsunami could arrive within a few hours.

If a very large earthquake fault moves underneath deep water, the wave it creates can be very long. For a large tsunami, a typical wavelength is 10 km, although some have been seen with wavelengths of 100 km and more. That means even in water with a depth of 1 km = 1000 meters, a tsunami is a *shallow* water wave! (Recall that a "shallow water wave" is one in which the wavelength is greater than the depth.)

The velocity of the tsunami can be calculated from the shallow water equation. In water, 3 km deep,  $D = 3000$  meters, so the velocity is  $v = 3.13 \sqrt{3000} \approx 171$  meters per second. That's 386 miles per hour, about half as fast as the speed of sound in air. A tsunami that is generated by an earthquake 1000 mi away will take 2.6 hours to arrive. That's enough time to give warning to coastal areas that a tsunami is on its way.

### You can outrun a tsunami

Imagine a tsunami with that velocity, with a wavelength of 30 km. Imagine that one crest of the wave passes you. The next one is approaching you from 30 km away. Even with its speed of 313 m/s, it will take  $t = d/v = 30000/313 = 100$  s to reach you. The water will fall for the first 50 of these seconds, and then rise for the next 50. Thus, although these waves travel fast, they are slow to rise and fall. That's why tsunamis were called tidal waves. If you are in a harbor, and there is a small tsunami, it might take 100 s for the water to rise and fall, and it gives the appearance of a rapid tide. The image of a huge breaking wave hitting the shore is largely fictional; most tsunamis are just very high tides that come and wash away everything close to the shore.<sup>5</sup> That's how they do their damage. If the ocean rises 10 m, it destroys everything, even if it takes 50 s to reach its peak. If you are young and healthy, you can usually outrun the rising water as it comes in. If you are not fast enough, then you get swept up in a very large volume of water, and dragged out to sea when the wave recedes. Small tidal waves are frequently observed as slow (100 s) rises and falls in harbors. Boats tied to docks are often damaged by these slow waves as they rise above the dock and get thrown into other boats. Many captains take their boats out into the harbor or out to sea when they are alerted that a tsunami is coming. In Japanese, the word "tsunami" means "harbor wave."

## The equation for waves

Recall from the last section that if the wavelength is  $L$ , and the velocity is  $v$ , then the time it takes between crests hitting you is  $T = L/v$ . The time  $T$  is called the **period** of the wave. This calculation is true for all waves, sound,

---

<sup>5</sup> The tsunami in the movie *Deep Impact* (1998) is particularly inaccurate. It shows a giant wave breaking over Manhattan Island. But the harbor of New York City is relatively shallow; there is no place for that much water to come from, unless a giant wave broke far out to sea.

tsunamis, deep water waves, even light waves. It is the fundamental relationship between velocity, period, and wavelength.

$$T = L/v$$

If the period is less than one second, it is usually more convenient to refer to the number of crests that pass by every second. That is called the *frequency*  $f$  of the wave, and it is given by  $f = 1/T$ . Putting these into the above equation, we get  $1/f = L/v$ , or

$$v = fL$$

You don't have to memorize this equation, but we will use it a lot, especially when we discuss light. Light in vacuum has a speed  $v = 3 \times 10^8$  m/s, a number we usually call  $c$ . Since we know  $c$ , the equation allows us to calculate the frequency whenever we know the wavelength, or the other way around.

## Sound doesn't always travel straight

Sound waves, whether in air or in ocean, often do not travel in straight lines. They will bend upwards or downwards, to the left or to the right, depending on the relative sound speed in the nearby material. Here is the key rule:

**Waves tend to change their direction by bending their motion towards the side that has a slower wave velocity.**

To understand why this is so, imagine that you are walking arm-in-arm with a friend. If your friend is on your left side and slows down, that pulls your left side backwards and turns you towards the left. If your friend speeds up, that pulls your left arm forward and turns you to the right (and also turns your friend to the right). The same phenomenon happens with waves. A more complete description of this is given in the optional section at the end of the chapter about Huygens's principle.

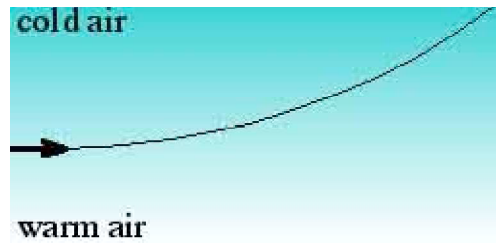
This principle can be demonstrated in a large classroom by having students raise their hands as soon as their neighbor raises a hand. The location of students' hand-raising moves throughout the classroom like a wave. This procedure is popular among fans in sporting events, where it is also called a "wave." If the students in one part of the room are told to be a little slower, then the wave will bend towards them as it spreads across the room.

The direction changing rule is true for all kinds of waves, including sound, water surface waves, and even earthquakes and light.

### Example: "normal atmosphere"

Here is an example from the atmosphere. At high altitude, the air is usually colder. That means that the velocity of sound at high altitude is slower than it is at low altitude.

Now imagine a sound wave that is initially traveling horizontally, near the surface of the earth. Above it, the velocity is lower, so it will tend to bend upward. This is shown in the following diagram.

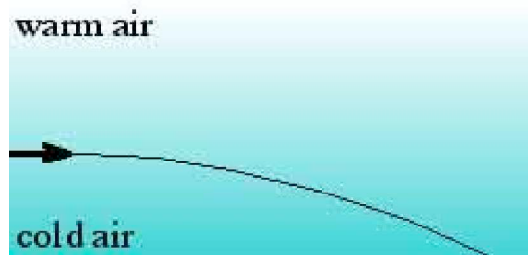


Cool atmosphere near the warm ground. Sound bends upward.

Notice that the sound bends away from the ground towards higher altitude. It bends upward. That's because the air above it has a slower sound velocity.

### Sound in the evening

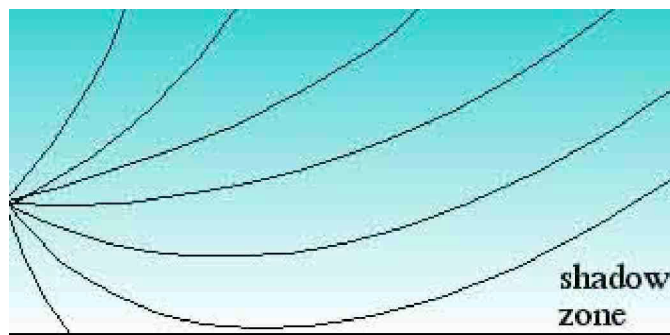
When the sun sets, the ground cools off rapidly. (It does this by emitting infrared radiation; we'll discuss this further in the Chapter 9 "Invisible Light.") The air does not cool so quickly, so in the evening, the air near the ground is often cooler than it is up higher. This phenomenon is called a "temperature inversion" because it is opposite to the normal pattern of the daytime. When there is a temperature inversion, sound tends to bend down towards the ground, as shown in the following figure.



Sound near the ground when there is an inversion

### Sound during the day, again

Now let's look at the morning situation again, with warm air near the ground and cold air up high. But let's draw many sound paths, all coming from the same point. This is done in the diagram below.

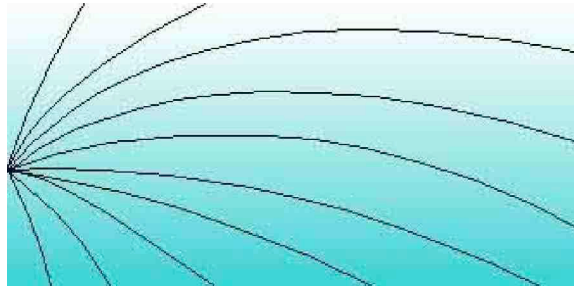


The solid line at the bottom represents the ground. Note that it blocks certain paths--the ones that drop too steeply. In the lower right corner is a small region that none of the paths can reach, since to reach this region the sound waves would have to go through the ground. (We'll assume for now that the ground absorbs or reflects sound, and does not transmit it, at least, not very well.) If the sound were coming from the point on the left, and you were standing in the shadow zone, then you wouldn't hear any sound at all. You are in the sound shadow of the ground.

This diagram shows why mornings tend to be quiet. Sounds bend up toward the sky, and if you are near the ground, there is no way that most of them can reach you. You won't hear distant automobiles, birds, waves, lion roars...

### Sound in the evening, again

In the diagram below, I've redrawn the evening situation, with the inverted temperature profile (cold at the bottom, warm at the top).



Note that there is no shadow zone. No matter where you stand, there are paths by which the sound can reach you.

Have you ever noticed that you can hear more distant sounds in the evening than in the morning? I've noticed that in the evening I can often hear the sound of distant traffic, or of a train; I rarely hear such sounds in the morning. (This phenomenon first mystified me when I was a teenager living 1/4 mile from the beach. I noticed that I could hear the waves breaking in the evening, but almost never in the morning.)

The explanation is in the diagrams above: in the evening, sound that is emitted upward bends back down, and you can hear sound from distant places. There is no shadow zone.

If you happen to be a wild beast, then the evening would be a good time to search for prey, since you could hear it even when it was far away.

### Forecasting a hot day

There are times when I wake up in the morning and hear distant traffic. Then I know that it will probably be a hot (and maybe smoggy) day. I learned this from experience long before I figured out the reason why.

The reason is that hearing distant sounds means there is an inversion, i.e. the high air is warmer than the low air. The sound diagram for an inversion in the morning is identical to the sound diagram shown above for the evening.

Inversions are unusual in the morning, but they do happen. The presence of an inversion in the morning leads to a special weather condition. On normal (no inversion) days, hot air is near the ground, and cold air is above it. Hot air is less dense than cold air, so it tends to float upward. (In the same way, wood floats on water if it is less dense than water.) So the hot air tends to leave the ground, replaced by cooler air from above.

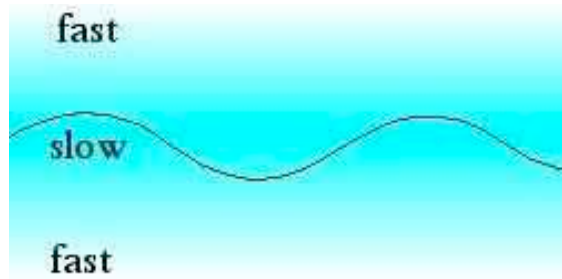
But if there is an inversion, i.e. there is hot air above and cold below, then the air above is less dense than the air at the ground. So convection, the floating of the ground air upward, doesn't take place. With no place to rise to, the hot air accumulates near the ground, making for a hot day. Smog and other pollutants also accumulate. The weather forecast on the radio or TV will often announce that there is an "inversion." Now you know what that means: the normal temperature profile is inverted, i.e. it is upside down with the cool air near the ground and warm air above.

Inversions frequently happen at the end of a hot day. The ground cools more rapidly than does the air (that's because it emits IR radiation; we'll discuss that in the next chapter). The air near the ground is cooled by contact with the cool ground, and the air above remains warm (unless there is turbulence from wind). For people who are sensitive to smog, the announcement of an inversion is bad news. For people who love hot weather, it is good news.

## The sound channel explained: focused sound

Now let's get back to the mysterious sound channel in the ocean that Maurice Ewing exploited for his sonar system in World War II.

In the ocean, the temperature of the water gets cooler as we go further down. This would make the speed of sound less. But, as mentioned earlier, the water is also getting more and more compressed (i.e. denser) because of the increasing pressure. This tends to make sound go faster. When these two effects are combined, we get a gradual decrease in sound velocity as we go from the surface to about 1 km of depth, and then the sound velocity increases again. This is illustrated in the diagram below. Darker means slower sound (just as it did in the atmosphere diagrams).



I've also drawn the path of a ray of sound. Notice that it always bends towards the slow region. The path I drew starts with an upward tilt, bends downward, passes through the slow region, and then bends upward. The path oscillates up and down, but never gets very far from the slow region, the 1 km-deep sound channel.

**Exercise:** Draw some other paths, starting at different angles. What happens if the ray starts out horizontally? Vertically?

## How Sofar saved downed pilots

Let's return now to the magic of Ewing's sofar spheres. As I stated earlier, they were hollow, and yet they were made of heavy material. Since they weighed more than an equal volume of water, they didn't float, but sank. Ewing designed the spheres to be strong enough to withstand the pressure of water down to a depth of 1000 m. At this depth, the spheres were suddenly crushed. (Like an egg, the round surface provides lots of strength, but when it breaks, it breaks suddenly.) The water and metal collapsed, and banged against the material coming in from the other side. It's like a hammer hitting a hammer, it generates a loud sound. The energy released from a sphere with radius of 1 in at a depth of 1 km is approximately the same as in 60 mg of TNT. That doesn't sound like a lot--but it is about the same amount you might find in a very large firecracker.

In the air, the sound of a firecracker doesn't go far, perhaps a few kilometers. But at a depth of 1000 m, the ocean sound channel focuses the sound. Moreover, the sound channel is quiet. Sound doesn't get trapped unless it originates within the sound channel itself. (Can you see why?) Any sounds created in the sound channel by whales or submarines stay in there, so the sound doesn't spread out as much as it would otherwise. Microphones placed within the sound channel can hear sounds that come from thousands of kilometers away.

During World War II, the Navy had arranged for several such microphones placed at important locations, where they could pick up the ping of the imploding Ewing spheres. They could locate where the implosion had taken place by the time of arrival of the sound. If the sound arrived simultaneously at two microphones (for example), then they knew the sound had been generated somewhere on a line that is equally distant from the two microphones. With another set of microphones they could draw another line, and the intersection of the two lines gave the location of the downed pilot.

**Historical note:** "Sofar" supposedly stands for "SOund Fixing And Ranging." Fixing and ranging was Navy terminology for determining the direction to a source (that's the fixing part) and its distance (ranging). Despite all this, I suspect that the acronym was forced, and the real name came about because the channel enabled you to hear things that were *so far* away. Some people still refer to the sound channel as the sofar channel. I learned about the sofar spheres from Luis Alvarez, who knew about them from his scientific work during World War II. I have spoken to several other people who remember them, including Walter Munk and Robert C. Spindel. Spindel believes that the spheres contained a small explosive charge to enhance the sound. We have not yet found any historical documentation that verifies this.

## Whale songs

What does the sound channel look like? The word "channel" can be misleading, since it brings up a vision of a narrow corridor. It is not like a tube. It is a flat layer, existing about a kilometer deep, spreading over most of

the ocean. Sound that is emitted in the sound channel tends to stay in the sound channel. It still spreads out, but not nearly as much as it would if it also spread vertically. That's why the sound can be heard so far away from its source. It tends to get focused and trapped in that sheet.

In fact, the sound channel is like one floor in a very large building, with ceiling and floor but without walls. Sound travels horizontally, but not vertically. If sound is emitted at the surface of the ocean, then it does not get trapped. So the sound of waves and ships does not pollute the sound channel. The sound channel is a quiet place for listening to sofar spheres and other sounds that are generated in the sound channel.

Whales discovered this, probably millions of years ago. We now know that whales like to sing when they are at the sound channel depth. These songs are hauntingly beautiful. If your computer has the right software, you can listen here to the recorded song of the humpback whale at

[www.muller.lbl.gov/teaching/physics10/whale\\_songs/humpback.wav](http://www.muller.lbl.gov/teaching/physics10/whale_songs/humpback.wav)

and of the gray whale at

[www.muller.lbl.gov/teaching/physics10/whale\\_songs/gray.wav](http://www.muller.lbl.gov/teaching/physics10/whale_songs/gray.wav)

You can find other recordings on the Internet, and you can buy recordings on CDs. Nobody knows what the whales are singing about. Some unromantic people think they are saying nothing more than "I am here."

## Global Positioning System (GPS)

A favorite gadget for hikers, boaters, and travelers is a Global Positioning System (GPS) receiver. This is a small device that will tell you where you are on the earth, to an accuracy of a few feet. You can buy one at a sporting goods store for about \$100. If you rent a car, for a small charge you can get one with GPS built in--to help keep you from getting lost. Many new cars now come equipped with GPS and a built in map system.

Why am I talking about GPS? Because GPS uses the same idea that Maurice Ewing used for locating pilots. For GPS, however, the signals are sent using radio waves rather than sound. And instead of using microphones set on the edges of the ocean, it uses radio receivers orbiting the Earth.

The GPS system works because there are now about 24 GPS satellites in space that are emitting signals. Each signal contains the time when it was emitted and the position of the satellite when it was emitted. The GPS receiver has a small computer and an accurate clock.

When the GPS receiver detects a signal, it looks at the time, reads the message saying when the signal was emitted, compares it with its own clock, and determines how long the signal was traveling. It multiplies that by the speed of light, and that gives it the distance to the first satellite. Of course, it also knows exactly where that satellite was when it emitted the signal. Once the GPS knows its distance from three different satellites, it can use geometry to calculate where it is. Can you see why that works?

## GPS geometry

How does the GPS system get its location by knowing the distance to three satellites? It's easy to see by analogy. Suppose you didn't know where you were in the US, but knew that you were 1000 mi from Denver and 1500 mi from San Francisco. To pin point your positions, you could first get a map and draw a circle around Denver with a 1000-mile radius. Then draw a 1500-



mile circle around San Francisco. The circles intersect at two points. If you knew your distance to one other city, you would know which of those two points represents your location.

If the GPS receiver gets a signal from four satellites, then it can see if the distance to that satellite is exactly what it expected. It should come out right--since the GPS already knows its own position. Suppose it turns out to be wrong? The only explanation can be that the clock in your inexpensive GPS receiver has drifted and is no longer accurate. So the receiver can use the fourth satellite to adjust its clock! The result is that if it can pick up four satellites, the receiver does not need an accurate clock.

## The Cold War and SOSUS

During World War II, the part of the military that used submarines was called “The Silent Service.” This reflected the fact that any sound emitted by a submarine could put it in danger of detection, so submariners trained themselves to be very quiet. Someone in a sub who drops a wrench makes a sound that is unlike any other in the ocean. (Fish don’t drop wrenches.) The wrench clatters against the hull, and the hull carries the sound to the water, and the vibrations of the hull send the sound into the ocean. Ships on the surface, and other submarines, had sensitive microphones to listen to possible sounds emitted from submarines.

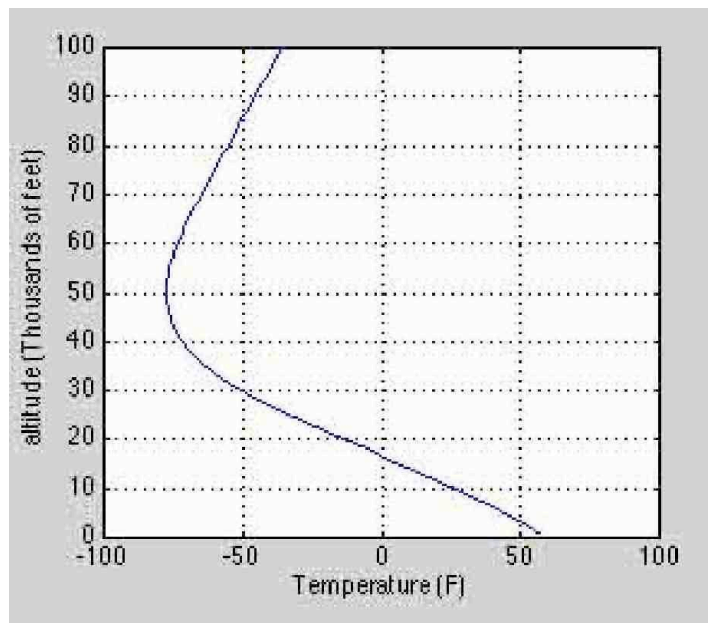
The presence of the sound channel did not remain secret for long, but its properties did. In the period from the 1950s to 1990s, the United States spent billions of dollars to put hundreds of microphones into the channel at locations all around the world. These microphones carried the signals back to an analysis center, and then the world’s best computers analyzed them. The system was called SOSUS, an acronym for “SOund SURveillance System.” The magnitude of the SOSUS effort was one of the best kept secrets of the Cold War. Effective use of SOSUS required the Navy to make extensive measurements of the ocean and its properties, and to update the temperature profile of the ocean all around the world. (The ocean has weather fronts analogous to those in the atmosphere.)

**Optional--Book to read:** If you really want to know more about this subject, one of the best introductions is the novel *The Hunt for Red October* by Tom Clancy. (Not the movie. The movie skips all the interesting technology.) When this novel came out in 1984, much of the material in it was still classified. Clancy had a talent for reading documents, talking to people, and figuring out from what they said and what was really true. The book was so detailed and so accurate (although it does have some fiction in it and some errors) that new people joining the submarine service were told to read the book in order to get a good picture of how operations worked! Many of the details of the SOSUS system were finally declassified in 1991, seven years after Clancy’s book was published. The SOSUS system was one of the largest and most expensive secret systems any nation ever built.

## Back to UFOs: a sound channel in the atmosphere!

Soon after he did his work in the ocean, Maurice Ewing realized that there should be a sound channel in the atmosphere! His reasoning was simple: as you go higher, everyone knows the air gets colder. Mountain air is colder than sea level air. The temperature of the air drops about 4 F for every 1000 ft of altitude gain.

That means that the velocity of sound decreases with altitude. But he also knew that when you get to very high altitudes, the temperature begins to rise again. Starting at about 40-50,000 ft, the air starts getting warmer. This is shown in the following figure.



Remember that the speed of sound depends on the temperature of the air. When the temperature is low, so is the speed of sound. That means that the speed of sound is fast at both high and low altitudes, and slower at about 50,000 ft.

Look at the diagram above the previous one, the diagram that showed sound moving in a wiggly line through the ocean. Exactly the same diagram can be used for sound in the atmosphere. That means that there is a sound channel in the atmosphere, centered at about 50,000 ft. (The exact altitude depends on latitude, as well as on the season of the year.) This is what Ewing figured out. He had an important US National Security application in mind to take advantage of this realization.

But first, we need a little more physics. Why does the atmosphere get warmer above 50,000 ft?

### Ozone: The cause of the high altitude heating

Why is high altitude air hot? The reason is the famous ozone layer. At about 40-50,000 ft, there is an excess of ozone, and this ozone absorbs much of the ultraviolet radiation from the sun. Ultraviolet is that part of sunlight that is

more violet than violet. This light is there, but invisible to the human eye. The ozone layer protects us, since ultraviolet light can induce cancer if absorbed on the skin. We'll talk more about the ultraviolet radiation in Chapter 9 "Invisible Light."

At the end of the 20th century, scientists began to fear that the ozone layer could be destroyed by human activity, and that would let the cancer-causing, ultraviolet radiation reach the ground with greater intensity. In particular, the scientists worried about the release of certain chemicals into the atmosphere called CFCs (chlorofluorocarbons, used in refrigerators and air conditioners). CFCs release chlorine and fluorine, and these catalyze the conversion of ozone  $O_3$  into ordinary  $O_2$ . (To balance the equation, 2 molecules of  $O_3$  turn into 3 molecules of  $O_2$ .)

The use of CFCs was outlawed internationally, and that was expected to solve the problem. For this reason, the human destruction of the ozone layer is no longer considered an urgent problem. For more, see chapter 9.

### **Looking at the ozone layer (and the sound channel): thunderhead tops**

On a day where there are large thunderstorms, you can see where the ozone layer is--right at the top of the tallest thunderheads. A thunderstorm grows from hot air at the ground, rising up through the colder (and denser) air above it. When the warm air hits the warm air of the ozone layer, it no longer rises. The cloud spreads out, making the "anvil head" shape that people associate with the biggest storms. So when you see the flat top of a large thunderhead cloud, you are looking at the ozone layer, and at the middle of the sound channel.

Think of it this way: there is a permanent "inversion" of atmospheric temperatures, if you go all the way up to 50,000 ft. That inversion prevents rising air from going any further, just as the low altitude inversion can prevent smog and hot air from rising away from the ground.



Image of an "anvil head" thunderstorm. The ozone layer and the sound channel are at the top of the cloud. (NOAA photo)

## Ewing's Project Mogul and his flying disks

Maurice Ewing had an urgent application for his predicted atmospheric sound channel: the detection of nuclear tests in Russia. In the late 1940s, the Cold War had begun and there was growing fear in many countries of the totalitarian communism represented by Russia. The Russians had great scientists and there was widespread belief that they would be building an atomic bomb soon. At that time, Russia was a very secret and closed society. In fact, Stalin was starving to death 30 million “kulak” farmers, and he could get away with it because he controlled information going in and out of the nation. In 1948, George Orwell wrote *1984*, expressing his fears of such a government.

Ewing realized that as the fireball from a nuclear explosion rose through the atmospheric sound channel, it would generate a great deal of noise that would travel around the world in the channel. (Not all of the sound *bang* is generated when the bomb detonates. The roiling fireball continues to generate sound as it reaches the atmospheric sound channel.) Ewing argued that we should send microphones up into the sound channel to detect and measure any such sound. That way we could detect Soviet nuclear tests, even with microphones in the United States!

The microphones that he used were called “disk microphones.” You can see them in photographs of old radio shows. For an old photo of a disk microphone, see

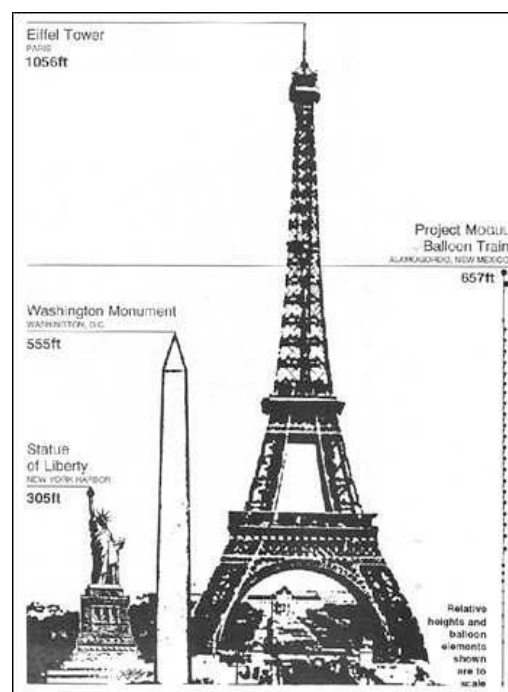
[www.muller.lbl.gov/photos/DiskMicrophone.jpg](http://www.muller.lbl.gov/photos/DiskMicrophone.jpg)

There is also a photograph of a somewhat smaller disk microphone used by Orson Wells in his famous 1938 broadcast of War of the Worlds, when he actually convinced many listeners that the world was under attack by Martians! For a photo, see [www.muller.lbl.gov/photos/DiskMicWelles3.jpg](http://www.muller.lbl.gov/photos/DiskMicWelles3.jpg)

You'll also see lots of disk microphones in the movie *The Aviator* (2004).

Ewing's idea was to string the microphones under a high altitude balloon, have them pick up the sounds in the sound channel, and then radio the sounds back to the ground. The disk microphones were called “flying disks.” (The word flying was not confined to airplanes; it was equally used by ballooners when they went up.) The balloons were huge, and the string of microphones was 657 ft long, longer than the Washington Monument is high.

The project was a success. The system detected American nuclear explosions, and on August 29, 1949, it detected the first Russian test.



The Project Mogul chain of microphones and balloons, compared in size to monuments. (From the US Air Force “The Roswell Report”. 1995.)

## The Roswell crash of 1947

One of the Project Mogul balloon flights crashed near the Roswell Army Air Force base on July 7, 1947. It was recovered by the U.S. Army, who issued a press release stating that “flying disks had been recovered.” The Roswell

Daily Record had headlines the next day. We referred to these at the beginning of this chapter: “RAAF Captures Flying Saucer“.

The fallen object was not a flying saucer, it was a complex balloon project that carried flying disk microphones to pick up Russian nuclear explosions. The program was highly classified, and the press release said more than the security people considered acceptable, so the next day the press release was “retracted.” A new press release stated that what had crashed was a “weather balloon.” It wasn’t a weather balloon. The US Government was lying.

### The government finally tells the truth

In 1994, at the request of a congressman, the U.S. government declassified the information they had on the Roswell incident, and prepared a report. Popular articles appeared in The New York Times, Popular Science (June 1997). But the Official U.S. government report (synopsis only) on Project Mogul is also available at

<http://muller.lbl.gov/teaching/physics10/Roswell/USMogulReport.html>

See also the official U.S. government report on the Roswell Incident at:

<http://muller.lbl.gov/teaching/physics10/Roswell/RoswellIncident.html>

Of these articles, you should read at least the New York Times article. The Popular Science articles give interesting background. The official U.S. government reports give details that you might find interesting.

Should the US Government ever lie? This is just the sort of issue that you should confront *before* you become president! It would make a good discussion question.

### How do we know the government isn’t lying now?



Richard Muller examines something interesting in the UFO Museum, Roswell, New Mexico

Many people believe the official government report on Project Mogul is just an elaborate cover-up. They believe that a flying saucer really did crash, and the government doesn’t want the public to know. Maybe I am part of this conspiracy, and part of my job is to mislead you into believing that flying saucers don’t exist! (According to the movie *Men in Black* (1997), the job of the Men in Black is to make sure the public never finds out.)

I suggest the following answer: the people who continue to believe that Project Mogul never happened probably don’t understand the remarkable science of the ocean and atmosphere sound channels. I could not have invented such a wonderful story. It has too many amazing details. In contrast, it is relatively easy to make up stories about flying saucers. Those don’t require much imagination. So here is my hypothesis: it is possible to distinguish the truth by the fact that it is more imaginative and more fascinating!

Of course, I might be lying. At left is a photograph taken of me at the UFO museum in Roswell, New Mexico in 2007.

## Earthquakes

When a fault in the Earth suddenly releases energy, it creates a wave in the ground. The location where an earthquake starts is called the epicenter. Most people who experience an earthquake are far from the epicenter, and are shaken by the wave that starts at the epicenter and shakes them as it passes by.

The epicenter of an earthquake can be located by noting when the earthquake wave arrived at several different locations--just as the sofar disks were used to locate downed pilots in World War II. Moreover, the epicenter is often deep underground, so even someone who is standing at the latitude and longitude of the epicenter can be standing over 15 miles away from it (i.e. above it).<sup>6</sup>

Huge amounts of energy are released in earthquakes, often greater than in our largest atomic weapons. That shouldn't surprise you. If you are making mountains shake over distances of tens or hundreds of miles, it takes a lot of energy. In 1935, Charles Richter found a way to estimate the energy from the measured shaking. His scale, originally called the Magnitude Local, became known as the Richter scale. An earthquake with magnitude 6 is believed to release the energy equivalent of about 1 million tons of TNT. That is the energy of a large nuclear weapon. Go up to magnitude 7 (roughly that of the Loma Prieta earthquake that shook San Francisco and the World Series in 1989) and the earthquake releases energy 10 to 30 times greater.

Why do I say a factor of 10 to 30? Which is it? The answer is that we don't really know. Magnitude is not exactly equivalent to energy. For some earthquakes, a magnitude difference of 1 unit will be a factor of 10, and for others it will be a factor of 30. It is easier to determine magnitude than it is to determine energy, and that's why magnitude is so widely used.

In the table on the next page I give the approximate magnitudes of some historical earthquakes in the US. I rounded them off to the nearest integer.

Earthquake	Approximate magnitude	Megatons of TNT
	6	1
San Francisco area 1989	7	10 to 30
San Francisco 1906	8	100 to 1000
Alaska 1999	8	100 to 1000
Alaska 1964	9	1000 to 30,000
New Madrid Missouri 1811	9	1000 to 30,000

Waves transport energy from one location to another. The velocity of an earthquake wave depends on many things, including the nature of the rock or soil in which it is traveling (granite? limestone?), and its temperature (particularly for earthquakes traveling in deep rock).

An especially deadly effect occurs when a wave moves from high velocity material into low velocity material, such as from rock to soil. When a wave slows down, its wavelength (the spacing between adjacent crests) decreases. But the energy is still there, but now squeezed into a shorter distance. That increases the amplitude of the shaking. Even though the

---

<sup>6</sup> "Shallow" earthquakes are defined to be those less than 70 km deep.

energy carried by the wave is unchanged, the effect on buildings becomes much stronger. This is what happened in downtown Oakland in the 1989 Loma Prieta quake. The earthquake wave passed right through much of Oakland without causing great damage, until it reached the area near the freeway. This region had once been part of the bay, and had been filled in. Such soft ground called “landfill” has a slow wave velocity, so the amplitude of the earthquake increased when it reached this ground. The most dangerous areas in an earthquake are regions of landfill. The Marina District in San Francisco is also landfill, and that is why it was so extensively damaged.

**Personal story from the author:** My daughters were at the Berkeley WMCA when the 1989 Loma Prieta earthquake hit. One of them told me that she was thrown up against the wall by the earthquake. I said to her, “No, Betsy, that was an illusion. You weren’t thrown against the wall. The wall came over and hit you.”

## Locating the epicenter of an earthquake

You already know that you can measure the distance to a lightning flash by counting the seconds and dividing by 5. The result is the distance to the lightning in miles. But here is another trick: as soon as you feel the ground shaking, and as you are ducking for cover, start counting seconds. When the bigger shaking finally arrives, take the number of seconds and *multiply* by 5. That will give you the distance to the epicenter (the place where the earthquake started) in miles.

Why does that work? To understand it, you should know that in rock, there are three important kinds of seismic waves. These are the P wave, the S wave, and the L wave.

### **The P wave (primary, pressure, push)**

P stands for “primary” because this wave arrives first. This is a longitudinal (compressional) wave, as is ordinary sound. That means that the shaking is back and forth in the same direction as the direction of propagation. So, for example, if you see that the lamppost is shaking in the east-west direction, that means that the P wave is coming from either the East or the West. Some people like to use the memory trick that the P wave is a Pressure wave, i.e. it is like sound because it is a compression and rarefaction, rather than a transverse motion. The P wave travels at about  $6 \text{ km/s} = 3.7 \text{ mi/s}$ . That is a lot faster than the speed of sound in air (which is  $300 \text{ m/s} = 0.3 \text{ km/s}$ ).

### **The S wave (secondary, shear)**

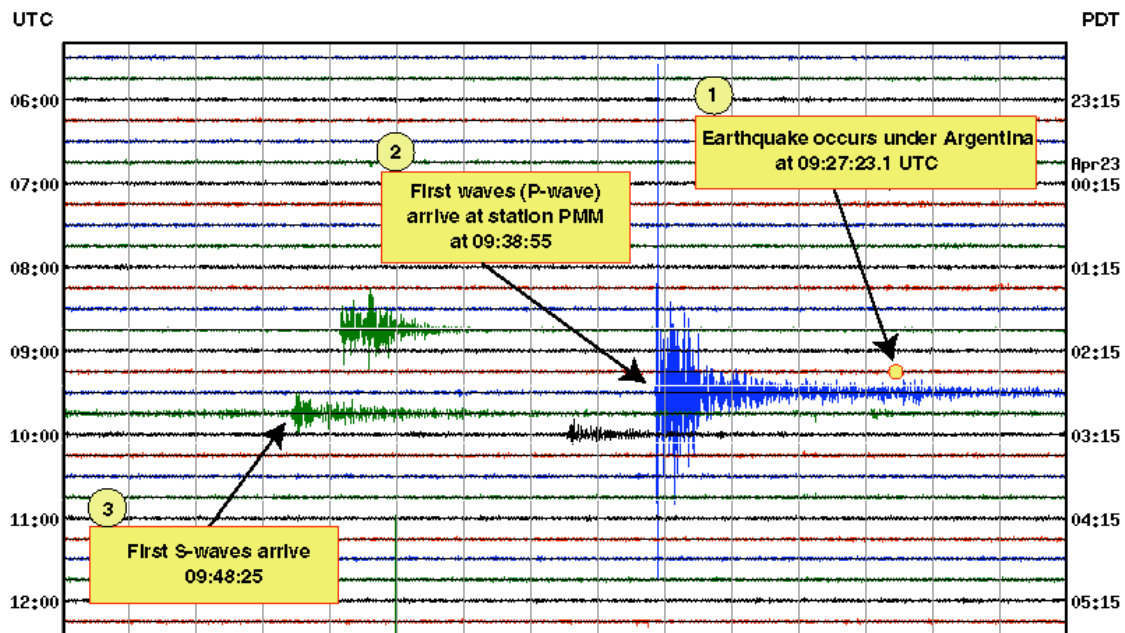
S stands for “secondary” because this wave arrives second. This is a transverse wave. That means that the shaking is perpendicular to the direction of propagation. If the wave is traveling from the east, then this implies that the shaking is either north-south, up-down, or some angle in between. Some people like to use the memory trick that the S wave is a shear wave, i.e. it can only propagate in a stiff material which does not allow

easy shear motion (sideways slipping). Liquids do not carry shear waves. We know there is a liquid core near the center of the Earth because shear waves do not go through it. The S wave travels at about  $3.5 \text{ km/s} = 2.2 \text{ mi/s}$ .

### The L wave (long, last)

L stands for “long.” These are waves that travel only on the surface of the Earth. Like water waves, they are a combination of compression and shear. They are created near the epicenter when the P and S waves reach the surface. They are called long because they tend to have the longest wavelength of the three kinds of seismic waves. It is the L wave that usually does the most damage, because the wave traveling on the surface often retains the biggest amplitude since it is not spreading out into three dimensions. The L wave travels at about  $3.1 \text{ km/s} = 2 \text{ mi/s}$ . Some people like to use the memory trick that the L wave is the last to arrive. (Careful with memory tricks. The L wave is NOT a pure “longitudinal” wave!)

The image below shows the shaking of the ground caused by a distant earthquake. Look at the wiggly line that crosses the image near the top. That is the first line, and it shows the shaking measured by the seismograph. Later lines are below this one. The little circle shows when the earthquake actually took place, at 9:27:23 UT. (UT stands for “universal time”, and it is the time at Greenwich UK.) The first shaking, due to the P wave, actually reaches the seismograph about 11 minutes later (shown as point 2). The S wave arrives about 10 minutes after that. There is no evidence for an L wave.



For a nice animation of the L wave (also known as the Rayleigh wave) go to: <http://www.kettering.edu/~drussell/Demos/waves/wavemotion.html> and look at the fourth animation on the page. If you look at the blue dots, you'll see they move in circular-like patterns (they are actually ellipses). At the top of the wave, the circle moves backwards – that is, opposite to the



direction the wave is moving! That's the opposite of what you saw in a water wave. And, if you go deep enough, the motion is forwards. Very strange.

### distance to the epicenter, again

Let's return now to the method of estimating the distance to the quake. As you are ducking under a table, start counting seconds from when you felt the first tremor, i.e. the P wave. (You can get very good at doing this if you live in California long enough.) When the S wave arrives then you know:

**For every second, the epicenter is about 8.4 km = 5 mi away.**

This is the rule I mentioned earlier. Thus if there is a 5-second gap between the waves, the epicenter was  $5 \times 5 = 25$  mi away. You may even be able to estimate the direction from the P wave shaking--the back and forth motion is in the same direction as the source. If we are lucky enough to have an earthquake during class, then you can watch me do this. (This equation is not true for travel through the deep earth, where the velocities are faster.)

For those of you who like math, can you see how I got the value of 8.4? It is based on the P and S velocities. *Hint:* the distance a wave travels is equal to the velocity multiplied by the time. This calculation is optional (not required) and relegated to a footnote.<sup>7</sup>

There are small earthquake waves passing by all the time, just as there are small waves everywhere you look on the ocean surface. To see the waves recorded for the last few hours, look at the University of California Berkeley Seismograph record at

<http://quake.geo.berkeley.edu/ftp/outgoing/userdata/quicklook/BKS.LHZ.current.gif>

This is an extremely interesting link to keep on your computer; it is something you can check any time you think you might have felt a quake. (Wait a little while before checking, the online plot is only updated every few minutes.) You'll see it there, even when it is not reported on the news. The quakes that occur every day in this region are also available on a map: <http://quake.wr.usgs.gov/recenteqs/>

## The liquid core of the Earth

Halfway to the center of the Earth, about 2900 km deep (1800 mi) is a very thick layer of liquid. (The distance to the center of the Earth is 6378 km.) You could say that the entire earth is "floating" on this liquid layer. The layer is mostly liquid iron, and the flow of this liquid creates the Earth's magnetic field (as discussed in Chapter 6). The liquid is so hot, that if we didn't have

---

<sup>7</sup> Suppose an earthquake is at a distance  $d$  from where you are standing. The P wave moves with a velocity  $v_p$ . The time it takes the wave to reach you is  $T_p = d/v_p$ . The S wave moves with a velocity  $v_s$ . The time it takes to reach you is  $T_s = d/v_s$ . First you feel the P wave, and you start counting seconds. Then the S wave arrives. The time difference that you measured is  $T = T_s - T_p$ . According to our equations, this is  $T = T_s - T_p = d/v_s - d/v_p = d(1/v_s - 1/v_p) = d(1/2.2 - 1/3.7) = d(0.184)$ . Solving for  $d$  gives:  $d = T/0.184 = 5.4 T$ . We approximate this as  $d = 5 T$ .

the rock blanket between it and us, the heat radiation from the core would quickly burn us to a crisp.

So much for curious facts--the real question for now is: how could we possibly know all this? The deepest we can drill is only a few miles. Nobody has ever gone to the core. Volcanoes don't come from regions that deep. How could we possibly know?

The interesting answer is that we know from watching signals from earthquakes. Thousands of these happen every year, and they are studied by earthquake detectors all around the Earth. The largest earthquakes send strong waves that travel down through the bulk of the Earth, and are detected on the opposite side.

An interesting aspect of the earthquakes is that **only the P waves pass through the core**. The S waves are all reflected! That is the wonderful clue. P waves are longitudinal "pressure" waves, and they travel through rock, air, or liquids. But S waves are transverse "shear" waves. Shear waves travel through solids, but they don't go through liquids or gases. That's because liquids and gases moving in the transverse direction can just slip past the rest of the liquid or gas; it doesn't exert much shear. So the fact that the P waves pass but the S waves don't gave one of the clues that there is a liquid core. Scientists also measured the speed at which the waves travel, and from this they can rule out gas and many kinds of liquids. They measure the density of the core from its contribution to the mass of the Earth, and they also see the magnetic field that the core creates. From all this, they were able to rule out every possible liquid except iron, although there could be liquid nickel mixed in with it.

We believe that the iron melted on the Earth when the Earth first formed. Most of the iron sank to the core, since it was denser than the other rock. The liquid iron is still in the core and it hasn't yet completely cooled off. The very center of the core, called the inner core, is under great pressure. Even though it too is hot, the inner core has been compressed into a solid. If the pressure of all the weight of the Earth were removed, it would turn into a liquid, or possibly into a gas.

**Discussion question:** how do we know the liquid core has a solid center? (Or rather, how did scientists figure that out?) For the answer, see the footnote.<sup>8</sup>

## Bullwhips

In a bullwhip,<sup>9</sup> the thickness of the whip is tapered towards the end. When the whip is snapped, a wave begins to travel down the whip to the end. Because the end is thin, the velocity of the wave increases near the end. The

---

<sup>8</sup> When a compressional wave hits the depth of the inner solid core it breaks up into two waves. From the behavior of these waves, we know that one of them is a shear wave. So although shear wave didn't travel through the outer core, shear waves are generated in the inner core. That means that the inner core must be made of a solid.

<sup>9</sup> If you don't know what a bullwhip is, then you might watch the opening scene in the movie *Indiana Jones and the Raiders of the Lost Ark* (1981) in which Indiana Jones uses a bullwhip to "whip" a gun out of the hand of a bad guy.

loud “crack” that you hear from the bullwhip is a sonic boom that occurs when the velocity of the wave exceeds the speed of sound.

Note this difference: in earthquakes and tsunamis, the added danger comes because the wave enters a region in which it slows. In the bullwhip, the crack comes because the wave speeds up.

## Waves can cancel (or reinforce)

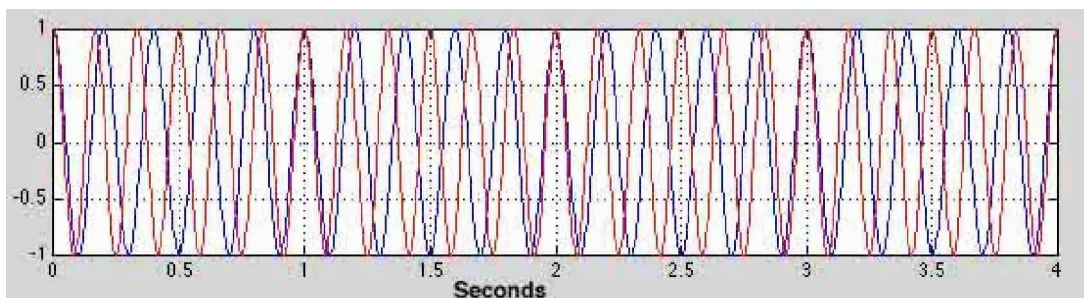
Suppose you are very unlucky, and are standing right in the middle of two earthquakes. One is to the north, and it takes you up, down, up, down, up, down, etc. The other earthquake arrives from the south, and it shakes you down, up, down, up, down, up, etc., exactly the opposite of the shaking of the first wave. What will happen? Will the up from one be canceled by the down of the other?

The answer is yes! If you are unlucky enough to be between two such waves, then try to be lucky enough to be at just the right place for them to cancel. You are depending on the fact that the two waves arrive with exactly opposite shakings.

Of course, if you were standing at a different location, the waves would arrive at different times, and they might not cancel. Suppose the first wave gave you up, down, up, down... and so did the second wave. Then the ups would arrive together, as well as the downs, and you would be shaken twice as much.

This circumstance is not as unlikely as you might think. Even if there is only one earthquake, parts of the wave can bend, and so you can be hit by the same earthquake but from two different directions. If you are lucky, the two waves will cancel, but a short distance away, they can add. This phenomenon was seen in the 1989 Loma Prieta earthquake that shook Berkeley, Oakland, and San Francisco. There were buildings where one side was shaken badly (causing that side to fall down) and the other side of the building was undamaged. This was probably due to the arrival of the wave from two directions at once, and the cancellation of the wave at the lucky end of the building.

If two waves are traveling together in the same direction but have different wavelengths (or frequencies), then the same kind of cancellation can happen. Take a look at the two different waves shown in the plot on the next page, one in lighter in color and one darker. The curves show the amount that the ground moves up and down (in centimeters) at different times due to the light earthquake and the dark earthquake. Zero represents the original level. The dark earthquake shakes the ground upward (to 1 cm), and downward (to -1 cm). So does the light earthquake. So far, we have not considered the effects when added together.



First look at the darker wave. At zero seconds, it starts at the maximum value of 1. It oscillates down and up, and by the time it reaches 1 s it has gone through 5 cycles. (Verify this. Try not to be distracted by the lighter wave.) We say that the frequency of the dark wave is 5 cycles per second = 5 Hz.

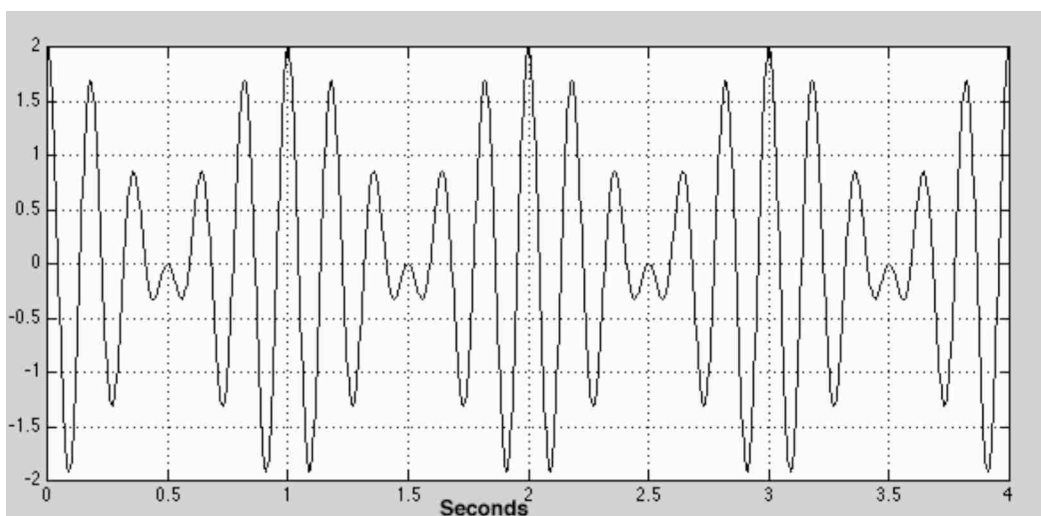
Now look at the light curve. In 1 s, it oscillates up and down 6 times. The frequency of the light wave is 6 Hz.

Suppose you are shaken by both waves at the same time. At time zero, you are shaken in the upward direction by both the dark and light waves; their effects add, and you will move up by  $1 + 1 = 2$  cm. Look at what happens at 0.5 s. The red wave is pushing you up by 1 cm, and the dark wave is pushing you down by 1 cm, so the two effects cancel, and at that instant you will be at level ground.

Note that there are also times when both waves are pushing you down. There's no place when they are both exactly at their minima, but they come pretty close at about 0.1 s. At this time both light and dark waves are down near -1 cm, so the sum effect will be to lower the ground by a total of 2 cm.

## Beats

If we add the light and the dark waves, point by point, we get the oscillation shown below.



The curve is taller because it ranges between 2 and -2. The shaking is not as regular, because of the alternating reinforcement and cancellation. Try counting cycles, and see what you get.

You probably got a frequency of 6 Hz (that's what I got). But some of the cycles are much bigger than others. Mathematically, we would not try to characterize this oscillation by a single frequency; it is a superposition (sum) of two frequencies.

If you felt this combination of waves under you, would say that the shaking was modulated, with the biggest shaking taking place every 1 s (at 0, 1, 2, 3, 4, ...). These are called the **beats**. The beat frequency is given by this elegant equation:

$$f_{\text{BEATS}} = f_1 - f_2$$

where  $f_1$  and  $f_2$  are the two frequencies that make up the signal (i.e. they are the frequencies of the light and dark waves). If the number comes out negative, ignore the sign; that's because beats look the same if they are upside down.

To demonstrate beats, you can listen to two tuning forks with slightly different frequencies. The demonstration that we use at Berkeley is described at [www.mip.berkeley.edu/physics/B+35+20.html](http://www.mip.berkeley.edu/physics/B+35+20.html)

For a very nice computer demonstration of how water waves can interfere, look at the UCLA site

<http://ephysics.physics.ucla.edu/physlets/einterference.htm>.

This site needs a fairly up-to-date browser with Java installed. You may already have that without knowing it, so it is worth trying. Move the red dots around, and then click on "calculate." Waves will come out of the two spots. These waves will add ("reinforce") at some locations, and cancel at others.

## Music: notes and intervals

A musical note usually consists of sound waves that have one dominant frequency. The middle white key on a piano, known as middle C, has a frequency of 256 Hz (at least when the piano is tuned to the "Just scale"). The white keys are designated A, B, C, D, E, F, G, A, B...with the 8 different letters which repeat in cycles. They repeat because, to most people, 2 consecutive Cs sound similar. They are said to be an "octave" apart. In fact, when you go 1 octave (8 notes), the frequency is exactly doubled. So the C above middle C has a frequency of 512 Hz. The next C has a frequency of 1024 Hz. Normal human hearing is quite good up to 10,000 Hz, and some people can hear tones as high as 15 to 20,000 Hz.

If two notes are played at the same time, and their frequency differs by just a little bit, then you will hear beats. Suppose you have a tuning fork that you know has a frequency of 256 Hz. You play the C string on a guitar, and listen to it and the tuning fork together. If you hear 1 beat per second, then you know the guitar is mistuned by 1 Hz; it is either 255 Hz or 257 Hz. You adjust the tension on the string until the frequency of the beats gets lower and lower. When there no longer are beats, the string is "in tune."

The interval between the A note and the higher E note is called a "fifth" because there are five notes: A, B, C, D, E. Likewise, middle C and the higher G make a fifth: C, D, E, F, G.

A violin is tuned so that the fifth has 2 frequencies with a ratio of exactly 1.5. So, with the middle C tuned to 256 Hz, the G above middle C has a

frequency of 379 Hz. This combination is also considered particularly pleasant, so many chords (combinations of notes played simultaneously, or in rapid sequence) contain this interval, as well as octaves.

Another pleasant interval is called the “third.” C and E make a third. The ratio of notes for a perfect third is  $1.25 = 5/4$ . The pleasant reaction of the sound is believed to be related to the fact that these frequencies have ratios equal to those of small, whole numbers.

A particularly unpleasant interval is the *tritone*, in which the frequencies have the ratio of  $7/5$ . The tritone is used in music to make the listener temporarily uncomfortable, and so it is considered dissonant. It is also used in ambulance sirens to make a sound that you can’t easily ignore.

## Vibrations and the sense of sound

As I said, the middle C on a piano vibrates 256 times per second. The C below that is 128 Hz. The next lower C is 64 Hz, and the one below that is 32 Hz. That’s pretty slow. If you can find a piano, play that note. Try singing it. Can you sense that your vocal cords are only vibrating 32 times per second? You almost feel that you can count the vibrations, but you perceive the tone as a tone, not as a collection of vibrations. If a light flickers at 32 times per second, you sense it as flickering, but your eye is more sensitive to the rapid changes than your ear. TV sets in Europe flicker at 50 Hz, and many people notice that. In the United States, TVs flicker at 60 Hz; most people do not perceive this! It is strange that the eye responds so differently to 50 Hz than to 60 Hz.

Ordinary house electricity oscillates 60 times per second, from positive to negative and back. Sometimes this causes a buzz in electronics, or in a faulty light bulb. The buzz is actually 120 times per second, since both the positive and the negative excursions of the current make sound. Do you remember hearing such a buzz? Can you hum the buzz, approximately? That is 120 Hz.

Remember the sound of the “light saber” in the Star Wars movies? That sound is 120 Hz. It sounds like a faulty fluorescent light bulb. In fact, it was made by picking up the buzz from electrical wires.

### Noise-canceling earphones

Because sound is a wave, it can be cancelled just like the shaking of an earthquake. So some smart people have made earphones that have a built-in microphone on the outside. This microphone picks up noise, reverses it, and then puts it into the earphone speakers. If done correctly, the reversed sound exactly cancels the noise, and the wearer hears “the sound of silence.” On top of this quiet, the electronics can put music into the earphones. Since the music does not reach the outside microphone, it is not cancelled.

I have a set of Bose noise-canceling earphones and I use them mostly on airplane flights. The result is that I can listen to high quality classical music, or to a typical airplane movie, and hear it as clearly as I would in a movie theater, without distracting noise.

There are even more expensive versions of noise-canceling earphones that are used by professional pilots and others who work in very noisy environments. It would be very nice to be able to cancel noise over a much larger region, e.g. in an entire room. However that is probably not possible, at least not from a single small speaker. The reason is that the wavelength of

sound (see next section) is typically 1 m. If the noise is not coming from the same location as the speaker, then although the sound could be cancelled in one location, it would probably be reinforced in a different location. That is not a problem for earphones, since the entire earphone is so small. Noise cancellation for an entire room might be possible if the walls of the room were made out of loudspeakers, or if they otherwise could be caused to vibrate to cancel any noise that might otherwise pass through.

## Wavelength of sound

Let's apply our wave equation to sound. Recall that the equation is:

$$v = fL$$

Let's use this equation to figure out the wavelength of sound for middle C on a piano. That has  $f = 256$  Hz. You learned in Chapter 1 that the velocity of sound in air is about 330 m/s. So the wavelength is  $L = v/f = 330/256 = 1.3$  m.

Does that seem long to you? It is large compared to the typical size of a head, so the wave is moving your two ear drums together.

Suppose we go up by 3 octaves. That means the frequency is doubled 3 times, i.e. increased by a factor of 8. Since the sound velocity  $v$  is the same in the wave equation, that means that the wavelength will be reduced by a factor of 8, from 1.3 meters, to 0.16 meters = 16 cm. That's smaller than the distance between your ears. So for this frequency, the eardrums on the opposite sides of your head may be vibrating opposite to each other.

## Doppler shift

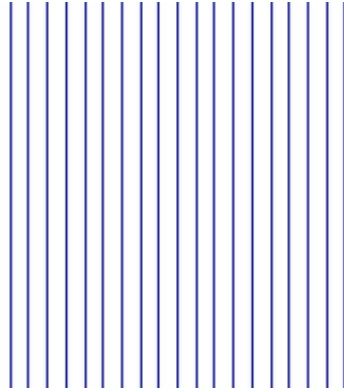
When an object is approaching you, you'll hear a higher frequency than the one they emit. That's because each time a crest is emitted, or a trough, the object is closer to you than it was for the previous cycle. So you hear them closer together. Likewise, if the object is moving away, you'll hear a lower frequency. This effect is called the Doppler shift, and it is extremely important in radar and in cosmology since it allows us to detect the velocity of very distant objects.

When a car or truck goes by, listen to the sound. I'm not sure how to describe it in words--something like "shhhhh--ooooo." (Sorry. That's the best I could come up with. Suggestions for better ways to write this would be appreciated.) But the important thing to note is that the pitch of the sound drops just as the car goes by (that's the change from the shhhhh to the ooooo). That's the Doppler shift.

The Doppler shift is seen in all waves, not just sound. The Doppler shift in light means that an object moving away from you has a lower frequency. In astronomy, this is referred to as the red shift. It was from the red shift that Edwin Powell Hubble (1889-1953) discovered that the Universe is expanding away from us.

## Huygens's Principle -- why waves bend toward the slow side

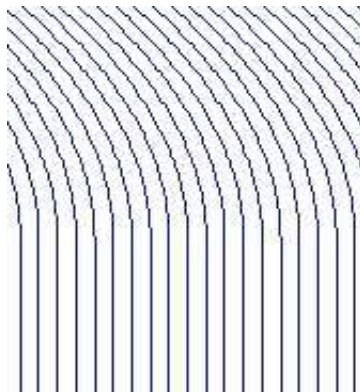
Imagine you are in an airplane, and you are watching waves on the ocean. Draw lines on the crests of the waves, i.e. on the highest points. Suppose the waves are moving to the right. The image will look like this:



Look carefully at this image. The lines are the crests, i.e. they are the high points of the waves. The waves are all moving towards the right. That means that if we had a movie, each crest (each line) would move towards the right. In between the lines are the low points of the water waves, called the troughs. They move too.

Recall that the distance between the crests is called the wavelength. In the figure, the wavelength is the separation of the lines.

Now imagine the waves moving to the right, but the ones near the top of the picture moving slower than the ones at the bottom. The lines would have to distort for this to be true. This is what it would look like:

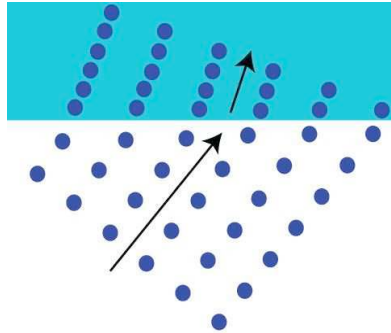


The waves near the top are moving to the right, but slower than the ones near the bottom. They will arrive at the right edge later. Notice how the slowing tends to bend their direction. But also notice that the waves near the top are becoming diagonal. The crests are no longer straight up and down. The direction of the wave is perpendicular to the crest. So the wave is no longer moving from left to right, but is also moving slightly upward. The direction has changed towards the side that has slower velocity.

The same thing would happen with a marching band (assuming that adjacent band members held hands), seen from above, if the field near the



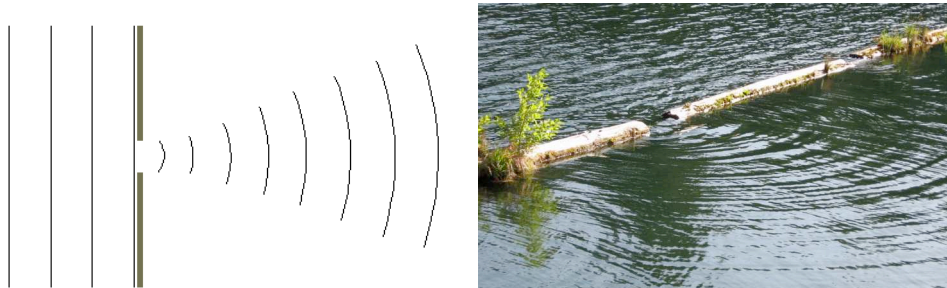
top was muddy and the members marched slower than the ones near the bottom. This is illustrated in the figure below.



Note how the direction of the marching band (the arrows) changes once the band enters the muddy field – presuming that they try to stay lined up with each other. For waves, they do stay lined up, because each wave is generating the next one. That is a big abstraction, but many people find the above diagram helpful anyway. This method of explaining wave direction change is called “Huygens principle”.

## Spreading of waves

Any wave, when passing through an opening, spreads. If not for such spreading, we usually would not hear people shout when they talk from behind a corner. The figure below left shows a wave (coming from the left) going through a hole and spreading out. On the right is a photo of waves passing through a gap in a floating log (taken by Michael Leitch).



There is a simple formula for this spreading of waves. The only thing you need to know is the wavelength of the wave  $L$  and the diameter of the opening  $D$ . Then the size  $S$  of the wave when it goes a distance  $R$  will be given by the approximate formula

$$S = \frac{L}{D} R$$

You don't need to learn this formula, but it will be useful for us to calculate wave spreading. It becomes very important for light, because it limits the ability of telescopes to “resolve” objects. As we will show in the next

chapter, this spreading is what prevents spy satellites from being able to read license plates.

This spreading equation is true for all kinds of waves, including sound waves and earthquake waves. The same equation works for them all. Let's take sound as an example. We showed above that for the tone of middle C, the wavelength is  $L = 1.3$  meters. Suppose this sound passes through a doorway that is 1.3 meters in diameter. If not for the spreading, the wave would be only 1.3 meters in diameter even after it went 10 meters past the door. But from the equation, there will be spreading. The amount of spreading will be large:

$$\begin{aligned} S &= \frac{L}{D} R \\ &= \frac{1.3}{1.3} 10 \\ &= 10 \text{ meters} \end{aligned}$$

The spreading is so great that you can hear a person on the other side of the door even if you can't see him. The spreading of light will be much less, because  $L$  will be much smaller. We'll talk about the spreading of light in the next chapter.

## Chapter Review

Waves travel in many materials, such as water, air, rock, and steel. Even though the material only shakes and none of the molecules move very far, the wave moves and carries energy over long distances. Waves are longitudinal when the direction of vibration is along the direction of the wave. Longitudinal waves include sound waves and the P earthquake wave. Waves can also be transverse. This means that the shaking is perpendicular to the direction of motion. An example is a wave on a rope. Water waves are both transverse and longitudinal. Light waves travel in a vacuum or in a material such as glass. They consist of a shaking electric and magnetic field. Light waves are transverse. Electrons and other particles are actually waves too, but they are so short ("wave packets") that this was not discovered until the 20<sup>th</sup> century. The fact that particles are waves is called the theory of quantum mechanics.

If a wave repeats, then the number of repeats per second is called the frequency. For sound, frequency is the tone, i.e. high pitch or low pitch means high frequency or low frequency. For light, frequency is color. Blue is high frequency and red is low frequency. The wavelength is the distance between crests of the wave.

The velocity of the wave depends on the material it is passing through. Sound travels about 1 mi/5 s through air, but 1 mi/s in water, and even faster in rock and steel. Light travels at 1 ft every billionth of a second, i.e. 1 ft per typical computer cycle. That is 186,000 mi/s.

The speed of sound depends on the temperature of the air. In hot air, sound travels faster. If sound is traveling horizontally, but the air above or below has a different temperature, then the direction of the sound will bend towards the side that is slower. This phenomenon causes sound to get trapped beneath the ocean, and is exploited by whales to send sounds thousands of miles. It also was used by the military for sofar (locating downed pilots) and

for SOSUS (to locate submarines). If four different microphones can pick up the same sound, then the source of the sound can be located. The same principle is used using radio waves for GPS.

A sound channel in the atmosphere is created because of the high altitude heating caused by the ozone layer. Project Mogul took advantage of the sound channel in the atmosphere. It was designed to detect Soviet nuclear tests. When the flying microphone disks crashed near Roswell, New Mexico in 1947, stories began to spread about flying saucers.

When the ground is cool, sound bends downward, and that lets us hear distant sounds. When the ground is warm, sound bends upward, and we do not hear distant sounds.

The velocity of sound waves does not depend on their frequency or wavelength. If it did, it would be hard to understand speech from someone standing far away. But the velocity of water waves does depend on the frequency and wavelength. Long wavelength water waves travel faster than short wavelength ones. Very long wavelength water waves, usually triggered by earthquakes, are called tsunamis or tidal waves.

Earthquakes begin when a fault ruptures at the epicenter, but they then travel as waves to distant places. The Richter scale gives a rough idea of the energy released. One point in the Richter scale is about a factor of 10 to 30 in energy released. The P wave is a compressional wave that travels fastest. Next comes the S wave (transverse), and finally the L wave. The time between the P and the S wave can be used to tell the distance to the epicenter. The fact that S waves do not travel through the center of the Earth, enables us to deduce that there is liquid there, probably (from the velocity we measure) liquid iron.

Waves can cancel, and that gives rise to beats (in music) and to strange effects, such as buildings that feel no shaking because of the fact that two canceling earthquake waves approached the building from different directions.

# Lecture 13:

## Electromagnetic waves

### 1 Light as waves

In previous classes (15a, high school physics), you learned to think about light as rays. Remember tracing light as lines through a curved lens and seeing whether the image was inverted or right side up and so on? This ray-tracing business is known as **geometrical optics**. It applies in the limit that the wavelength of light is much smaller than the objects which are bending and moving the light around. We're not going to cover geometrical optics much in 15c (except for what you do in lab). Approaching light as waves will allow us to understand topics of polarization, interference, and diffraction which you may not have seen in a physics class before.

Electromagnetic waves can have any wavelength. Electromagnetic waves with wavelengths in the visible spectrum ( $\lambda \sim 10^{-6} \text{ m}$ ) are called **light**. Red light has a longer wavelength than blue light. Longer wavelengths than red go into infrared, then microwaves, then radio waves. Smaller wavelengths than blue are ultraviolet, then  $x$  rays than  $\gamma$  rays. Here is a diagram of the spectrum:

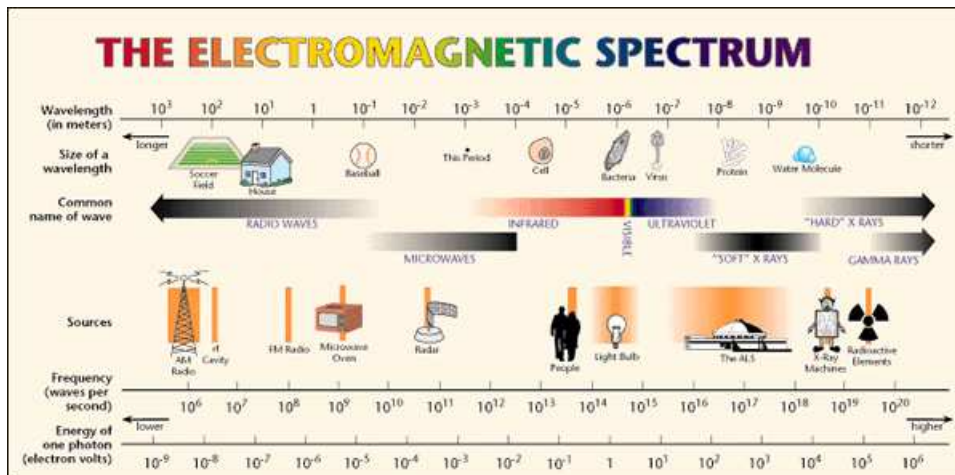


Figure 1. Electromagnetic spectrum

It is helpful to go back and forth between wavelengths  $\lambda$ , frequency  $\nu = \frac{c}{\lambda}$  and energy  $E = h\nu$ , where  $h = 6.6 \times 10^{-34} \text{ J} \cdot \text{s}$  is Planck's constant. The fact that energy and frequency are linearly related is a shocking consequence of quantum mechanics. In fact, although light is wavelike, it is made up of particles called **photons**. A photon of frequency  $\nu$  has energy  $E = h\nu$ . The intensity of light (recall intensity is power/area,  $I = \frac{P}{A}$ ) is also the number density  $n$  of photons per unit volume, times their rate (the speed of light  $c$ ), times their energy:  $I = nch\nu$ . Thus, given a frequency of light, the intensity of light at that frequency tells you how many photons are present. It's usually a lot. For example, light from the sun has intensity  $I = 1.2 \frac{\text{kW}}{\text{m}^2}$ . Assuming this peaks in the visible ( $\sim 500 \text{ nm}$ ), we get  $n = \frac{I\lambda}{c^2 h} = 10^{22} \text{ m}^{-3}$ . This is roughly Avogadro's's number per cubic meter.

We use the words soft and hard to mean lower energy and higher energy respectively. That's because higher energy photons just plow through things, like a bullet, while lower energy photons are easy to deflect. Let's start by briefly going through the spectrum from hard to soft (small  $\lambda$  to big  $\lambda$ ), discussing some of the relevant physics.

Very energetic photons are called **gamma rays**. Gamma rays are a type of radiation coming from radioactive decay. Their high energy makes these rays very dangerous. Gamma rays can also be seen coming in from outer space with energies of  $10^{19}$  eV or wavelengths as small as  $10^{-25}$  m!! These are called **cosmic rays**.

Typical atomic spacings in a solid are of order 1 Angstrom =  $10^{-10}$  m. This is not surprisingly the wavelength of x-rays; their wavelength, also of order  $10^{-9}$  m, can probe these solids. Hard (meaning shorter wavelength/larger energy) x-rays, as used in x-ray machines scatter off of electrons in solids. They usually scatter inelastically, knocking the electron out of its orbit and ionizing the atom. Thus they are called **ionizing radiation**. Ionizing radiation can be dangerous. The more dense a solid, the more electrons are there to scatter off of, so the x-rays go through less before scattering. That's why x-rays machines can distinguish dense materials like metal or bone from less dense ones like skin. Soft x-rays have wavelengths of multiple atomic spacings. These can be used to probe how a solid is organized through the wave nature of light. These softer x-rays are used in x-ray crystallography.

There is an important relationship between temperature and the wavelength where the spectrum of a blackbody peaks:  $\lambda_{\max} T = 2.898 \text{ mm K}$ . This is known as **Wein's displacement law**. We will discuss it more in Lecture 15. For  $T = 5778 \text{ K}$  the surface temperature of the sun,  $\lambda_{\max} = 501 \text{ nm}$ , which is yellow light. Thus sunlight peaks as yellow.

Visible light has wavelengths  $\lambda \sim 10^{-7} \text{ m}$  which are a few orders of magnitude longer than typical atomic spacings. However the energy of visible light,  $\sim$  few electron volts, is the same as the spacing between energy levels in atoms. More precisely, the key unit is a Rydberg,

$$1 \text{ Ry} = \frac{1}{2} m_e c^2 \alpha^2 = 13 \text{ eV} \quad (1)$$

with  $m_e$  the electron mass, and  $\alpha \approx \frac{1}{137}$  the fine-structure constant. Visible light has energy of a few Rydbergs. Materials have different colors because of resonances which excite electronic transitions in materials absorbing the light. This is a quantum phenomenon. We will discuss color in Lecture 15.

Infrared light ( $\lambda = 10^{-4} \text{ m} \sim 10^{-6} \text{ m}$ ) has lower energy than visible light. It cannot excite electron transitions in molecules, but it can excite vibrational modes. These excitations are broader (lower  $Q$ ) and more sensitive to the thermal motion. They do not usually appear as sharp resonances. By Wein's displacement law, a blackbody at temperature  $310 \text{ K}$  (like mammals) has a peak wavelength of  $\lambda_{\max} = 10^{-5} \text{ m}$  which is in the infrared. That's why infrared cameras which see heat maps in the dark. Some snakes are also sensitive to infrared radiation for the same reason.

Microwaves ( $\lambda = 10^{-3} \text{ m} \sim 10^{-1} \text{ m}$ ) are longer wavelength/lower energy than infrared. They can excite rotational modes of water, which is how microwave ovens work. Microwaves are also used in full body airport scanners: they can see skin under clothing because they see the water in the skin! (Alternative airport scanners use ionizing x-rays at low intensity and look at the photons which backscatter).

Longer wavelength radiation are generically called radio waves. They are very low energy, non-ionizing, and not dangerous. The higher frequency the radio waves, the more information can be carried (see Lecture 12). For example, you need about 10 kHz bandwidth for audio and 5 MHz bandwidth for television. AM radio has frequencies around 1 MHz ( $\lambda \sim 100 \text{ m}$ ), FM radio around 100 MHz ( $\lambda \sim 1 \text{ m}$ ). Cell phones and wireless data transfer are usually done in GHz bands ( $\lambda \sim 1 \text{ cm}$ ). 3G technology has around 1.7 GHz carrier frequencies, and 4G is 2.5 GHz. Although higher frequencies than this could transmit more data, we don't want cell phones to transmit in the microwave region (for obvious reasons!)

## 2 Wave equation from Maxwell's equations

With no charges around ( $\rho = \vec{J} = 0$ ), as in empty space, Maxwell's equations are

$$\boxed{\vec{\nabla} \cdot \vec{E} = 0 \quad \vec{\nabla} \times \vec{E} = -\frac{\partial}{\partial t} \vec{B} \quad \vec{\nabla} \cdot \vec{B} = 0 \quad \vec{\nabla} \times \vec{B} = \mu_0 \epsilon_0 \frac{\partial}{\partial t} \vec{E}} \quad (2)$$

Taking a time derivative of the 4th equation gives

$$\mu_0\epsilon_0\frac{\partial^2}{\partial t^2}\vec{E} = \frac{\partial}{\partial t}[\vec{\nabla} \times \vec{B}] = \vec{\nabla} \times \frac{\partial}{\partial t}\vec{B} = -\vec{\nabla} \times (\vec{\nabla} \times \vec{E}) \quad (3)$$

where the 2nd Maxwell equation was used in the last step. Now we can use the identity

$$\vec{A} \times (\vec{B} \times \vec{C}) = \vec{B}(\vec{A} \cdot \vec{C}) - (\vec{A} \cdot \vec{B})\vec{C} \quad (4)$$

For  $\vec{A} = \vec{B} = \vec{\nabla}$  and  $\vec{C} = \vec{E}$  gives

$$\vec{\nabla} \times (\vec{\nabla} \times \vec{E}) = \vec{\nabla}(\vec{\nabla} \cdot \vec{E}) - \vec{\nabla}^2 \vec{E} = -\vec{\nabla}^2 \vec{E} \quad (5)$$

where the 1st Maxwell equation was used in the last step. We therefore have

$$\mu_0\epsilon_0\frac{\partial^2}{\partial t^2}\vec{E} = \vec{\nabla}^2 \vec{E} \quad (6)$$

or more simply

$$\left[ \frac{\partial^2}{\partial t^2} - c^2 \vec{\nabla}^2 \right] \vec{E}(\vec{x}, t) = 0 \quad (7)$$

where  $c \equiv \frac{1}{\sqrt{\mu_0\epsilon_0}}$  is the **speed of light**. Thus the electric field satisfies the wave equation. With similar manipulations, you can also show that

$$\left[ \frac{\partial^2}{\partial t^2} - c^2 \vec{\nabla}^2 \right] \vec{B}(\vec{x}, t) = 0 \quad (8)$$

Thus, the magnetic field must also satisfy the wave equation. Moreover, since the second Maxwell equation is  $\vec{\nabla} \times \vec{E} = -\frac{\partial}{\partial t}\vec{B}$  the solutions for the electric and magnetic fields are not independent. Thus we call the **electromagnetic waves**. We will next see how to solve these coupled equations.

### 3 Plane waves

Recall from Lecture 10 that for a sound wave, satisfying

$$\left[ \frac{\partial^2}{\partial t^2} - v^2 \vec{\nabla}^2 \right] p(\vec{x}, t) = 0 \quad (9)$$

with  $p(\vec{x}, t)$  the pressure at position  $\vec{x}$  and time  $t$ , plane waves looked like

$$p(x, y, z, t) = p_0 \cos(\vec{k} \cdot \vec{x} - \omega t + \phi) \quad (10)$$

Here,  $\vec{k}$  is the wavevector. The wave equation implies  $\omega = v|\vec{k}|$ . These waves are constant in the plane perpendicular to  $\vec{k}$ . For example, if  $\vec{k} = (0, 0, k_z)$  then

$$p(x, y, z, t) = p_0 \cos(k_z z - \omega t + \phi) \quad (11)$$

with  $\omega = vk_z$ . For a given frequency  $\omega$ , there can be plane waves with any wavevector  $\vec{k}$  satisfying  $|\vec{k}| = \frac{\omega}{v}$  and any phase  $\phi$ . Any wave can be written as a sum over plane waves

$$p(\vec{x}, t) = \int dk_x dk_y dk_z \tilde{p}(\vec{k}) e^{i(\vec{k} \cdot \vec{x} - \omega t)} \quad (12)$$

with  $\tilde{p}$  in general complex. Or for real  $p$  we can write

$$p(\vec{x}, t) = \int d^3k \tilde{p}_s(\vec{k}) \sin(\vec{k} \cdot \vec{x} - \omega t) + \int d^3k \tilde{p}_c(\vec{k}) \cos(\vec{k} \cdot \vec{x} - \omega t) \quad (13)$$

where  $d^3k$  is shorthand for  $dk_x dk_y dk_z$ .

For electromagnetic waves, the plane wave solutions to Eq. (7) are similar

$$\vec{E}(\vec{x}, t) = \text{Re} \left[ \vec{E}_0 e^{i(\vec{k} \cdot \vec{x} - \omega t)} \right] \quad (14)$$

with

$$\omega = c |\vec{k}| \quad (15)$$

We usually leave the taking of the real part implicit and work with complex electric and magnetic fields for simplicity. We can always take the real part at the end of the day.

A critical difference with sound or other waves is that the coefficient of the wave  $\vec{E}_0$  is itself a vector. Thus, in addition to a direction  $\vec{k}$  the wave is moving, there is a direction  $\vec{E}_0$  describing which way the electric field is pointing. We call  $\vec{E}_0$  the **polarization vector for the electric field**. As with any plane wave, a plane wave electric field is constant in the direction perpendicular to  $\vec{k}$ . For example, if  $\vec{k} = (0, 0, k_z)$  then  $\vec{E}$  is independent of  $x$  and  $y$  at a fixed  $z$  and  $t$ .

Now, the Maxwell's equation  $\vec{\nabla} \cdot \vec{E}(x, t) = 0$  implies that for a plane wave solution with wavevector  $\vec{k}$  that

$$0 = \vec{\nabla} \cdot \vec{E}(\vec{x}, t) = \vec{\nabla} \cdot \vec{E}_0 e^{i(\vec{k} \cdot \vec{x} - \omega t)} = i(\vec{k} \cdot \vec{E}_0) e^{i(\vec{k} \cdot \vec{x} - \omega t)} \quad (16)$$

The only way this can be satisfied at all times is if

$$\boxed{\vec{k} \cdot \vec{E}_0 = 0} \quad (17)$$

So the **polarization is orthogonal to the direction of propagation of the plane wave**.

Since the magnetic field satisfies the same equations  $\left[ \frac{\partial^2}{\partial t^2} - c^2 \vec{\nabla}^2 \right] \vec{B}(\vec{x}, t) = 0$  and  $\vec{\nabla} \cdot \vec{B} = 0$ , its plane wave solutions are similar

$$\vec{B}(\vec{x}, t) = \vec{B}_0 e^{i(\vec{k} \cdot \vec{x} - \omega t)}, \quad \vec{k} \cdot \vec{B}_0 = 0 \quad (18)$$

Finally, we can use that  $\vec{\nabla} \times \vec{E} = -\frac{\partial}{\partial t} \vec{B}$ . For an electric plane wave

$$\vec{\nabla} \times \vec{E} = \vec{\nabla} \times \vec{E}_0 e^{i(\vec{k} \cdot \vec{x} - \omega t)} = i(\vec{k} \times \vec{E}_0) e^{i(\vec{k} \cdot \vec{x} - \omega t)} \quad (19)$$

For a magnetic plane wave

$$-\frac{\partial}{\partial t} \vec{B}(\vec{x}, t) = -\frac{\partial}{\partial t} \vec{B}_0 e^{i(\vec{k} \cdot \vec{x} - \omega t)} = i\omega \vec{B}_0 e^{i(\vec{k} \cdot \vec{x} - \omega t)} \quad (20)$$

Setting  $\vec{\nabla} \times \vec{E} = -\frac{\partial}{\partial t} \vec{B}$  we see that **the wavevectors  $\vec{k}$  must be the same for  $\vec{E}$  and  $\vec{B}$** . We also find

$$\omega \vec{B}_0 = \vec{k} \times \vec{E}_0 \quad (21)$$

Since the vector  $\vec{A} \times \vec{B}$  is perpendicular to both  $\vec{A}$  and  $\vec{B}$ , Eq. (21) implies that  $\vec{B}_0$  is perpendicular to both  $\vec{k}$  (which we already knew) and  $\vec{E}_0$ . Also, since  $\omega = c|\vec{k}|$  by Eq. (15) we have

$$|\omega \vec{B}_0| = |\vec{k} \times \vec{E}_0| = |\vec{k}| |\vec{E}_0| = \frac{\omega}{c} |\vec{E}_0| \quad (22)$$

where we have used that  $\vec{k}$  is perpendicular to  $\vec{E}_0$  to write  $|\vec{k} \times \vec{E}_0| = |\vec{k}| |\vec{E}_0|$ . Thus we find

$$\boxed{|\vec{E}_0| = c |\vec{B}_0|} \quad (23)$$

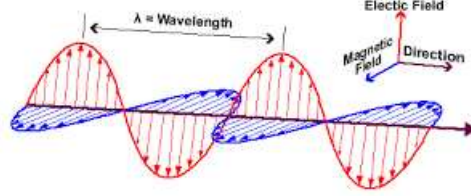
So the amplitude for the electric and magnetic fields in a plane wave are related by the speed of light.

Keep in mind that  $|\vec{E}| = c|\vec{B}|$  is a special condition for plane waves. It certainly does not hold for most solutions to Maxwell's equations. For example, we can certainly have pure electric fields (or pure magnetic fields). These satisfy Maxwell's equations in empty space but do not have  $|\vec{E}| = c|\vec{B}|$ . However, because any solution can be written as a sum of plane waves, by Fourier's theorem, discussing plane waves where  $|\vec{E}| = c|\vec{B}|$  will be very productive.

In summary, for plane waves, we have found for a given wavevector  $\vec{k}$  we have

$$\vec{k} \perp \vec{E}, \quad \vec{k} \perp \vec{B}, \quad \vec{E} \perp \vec{B}, \quad |\vec{E}_0| = c|\vec{B}_0| \quad (24)$$

We can draw a plane wave like this



**Figure 2.** Propagation of a linearly polarized plane wave where  $\vec{k}$  points to the right. In this wave, the  $\vec{E}$  field always points in vertically, the  $\vec{B}$  field horizontally, and their magnitudes vary with position in sync.

This picture illustrates some things but hides others. The  $\vec{E}$  field is constant in magnitude and direction in the entire plane perpendicular to  $\vec{k}$ . It points, at every position, in the same direction. Taking the real part (to find the actual electric field):

$$\text{Re}(\vec{E}) = \vec{E}_0 \cos(\vec{k} \cdot \vec{x} - \omega t + \phi) \quad (25)$$

we see that  $\vec{E}$  varies in time and in the  $\vec{k}$  direction. At any point, the  $\vec{B}$  field is exactly perpendicular to the  $\vec{E}$  field. The magnitude of the  $\vec{B}$  field is always  $|\vec{B}_0| = \frac{1}{c}|\vec{E}_0|$  at all points at all times.  $\vec{B}$  and  $\vec{E}$  are always in phase.

Finally, note that for a given  $\vec{k}$  and phase, there are two orthogonal directions that  $\vec{E}$  can point. Once  $\vec{E}$  is fixed,  $\vec{B}$  is fixed too. These are the two linearly independent linear polarizations of light. Any other solution to the wave equation in free space can be written as a sum of plane waves of these two polarizations.

## 4 Energy and power

In physics 15b you derived that the energy density in the electromagnetic field is

$$\mathcal{E} = \frac{1}{2}\epsilon_0|\vec{E}|^2 + \frac{1}{2}\frac{1}{\mu_0}|\vec{B}|^2 \quad (26)$$

This holds for *any type of electromagnetic field*, not just plane waves. I'm not going to re-derive this equation, but we can quickly check that it makes sense – we saw for other waves that the energy is proportional to the square of the amplitude, as it is in this case. The factors of  $\epsilon_0$  and  $\mu_0$  can then be put in by dimensional analysis.

What we are more interested in here is the energy transported by a wave. We can write

$$\begin{aligned} \frac{\partial \mathcal{E}}{\partial t} &= \epsilon_0 \vec{E} \cdot \frac{\partial \vec{E}}{\partial t} + \frac{1}{\mu_0} \vec{B} \cdot \frac{\partial \vec{B}}{\partial t} \\ &= \frac{1}{\mu_0} (\vec{E} \cdot (\vec{\nabla} \times \vec{B}) - \vec{B} \cdot (\vec{\nabla} \times \vec{E})) \\ &= \frac{1}{\mu_0} \vec{\nabla} \cdot (\vec{B} \times \vec{E}) \end{aligned}$$

The total energy in a volume  $V$  is  $E_V = \int_V \mathcal{E}$  so that

$$\frac{dE_V}{dt} = \frac{1}{\mu_0} \int_V \vec{\nabla} \cdot (\vec{B} \times \vec{E}) = \frac{1}{\mu_0} \oint (\vec{B} \times \vec{E}) \cdot d\vec{A} \quad (27)$$



where  $\vec{A}$  is the area element on the surface of  $V$ . This is just the divergence theorem. This says that

$$\vec{S} \equiv -\frac{dE_V}{dt \text{ Area}} = \frac{1}{\mu_0} \vec{E} \times \vec{B} = \frac{\text{watts}}{\text{meter}^2} \quad (28)$$

This is called the **Poynting vector**. It gives the power flowing *into* a region (the minus sign in Eq. (28) makes it *into* the region rather than *out of* the region). Note that since  $\vec{E} \times \vec{B} \propto \vec{k}$  for a plane wave, the Poynting vector points in the direction of the wave motion. So power flows in the  $\vec{k}$  direction. The Poynting vector measures the **intensity of a wave**.

Using  $\omega \vec{B}_0 = \vec{k} \times \vec{E}_0$ , for a plane wave the magnitude of the intensity is

$$I = |\vec{S}| = \frac{1}{\mu_0 c} \vec{E}^2 = c \epsilon_0 \vec{E}^2 \quad (29)$$

So the intensity is given by the square of the electric (or magnetic) fields.

In a sense, light is a transverse wave since  $\vec{E}$  and  $\vec{B}$  are transverse to the direction  $\vec{k}$  of the wave propagation. On other hand, light acts like a longitudinal wave in that it carries momentum. To see that, consider an electric charge which the wave hits. The electric field will push the charge in the  $\vec{E}$  direction. Then once the charge is moving, the magnetic field will push the charge in the  $\vec{v} \times \vec{B}$  direction. Since  $\vec{v} \propto \vec{E}$ , this is the  $\vec{E} \times \vec{B} \propto \vec{k}$  direction. That is, the combination of the electric and magnetic fields has the net effect of pushing the charge along the direction  $\vec{k}$  that the wave is going. Thus the charge ends up starting to spin in a little helix going in the  $\vec{k}$  direction. The effect is that light exerts pressure on matter. The pressure is called **radiation pressure**. The radiation pressure is given by the Poynting vector divided by  $c$ :

$$\vec{p} = \frac{1}{c} \vec{S} \quad (30)$$

By conservation of momentum, the light wave must also carry momentum.

So light carries momentum in the direction of  $\vec{S} \propto \vec{k}$ . Since light is made up of photons, each photon must also carry momentum proportional to  $\vec{k}$ . We give the proportionality constant a name, Planck's constant  $h$ :

$$\vec{p} = h \vec{k} \quad (31)$$

In other words, knowing that light is made up of photons is enough to derive that the momentum of the photons is proportional to  $\vec{k}$ . Since  $\vec{k}$  is the Fourier conjugate of  $\vec{x}$ , if a particle is well localized in  $\vec{x}$ , it must be poorly localized in  $\vec{p}$  (as we learned from our study of wavepackets). This is the origin of the Heisenberg uncertainty principle. Of course, you probably can't make any sense out of what it means for an electron to be "well-localized" in  $x$ . To understand this concept requires quantum mechanics.

Of course, mostly when electromagnetic radiation passes through an object, not much of the momentum is transferred. If the radiation is completely absorbed, as visible light is on black paper, the momentum is completely transferred.

By the way, Kepler postulated radiation pressure way back in 1619 to explain why the tails of comets are always pointing away from the sun.

## 5 Historical note (optional)

Since the electric and magnetic fields satisfy the wave equation, it is natural to ask what are the little oscillators which are strung together to allow the waves to propagate? One case is that they might be air molecules. This guess is quickly disproved since the propagation is independent of the density of air. In particular, electromagnetic waves propagate perfectly well in a vacuum. Thus the medium is something new which we call the **aether**.

Maxwell proposed a model of the aether based on gears and wheels. These gears have to be able to move very rapidly to support high frequencies. A number of other models were proposed in the 19th century. The aether has never been discovered. It is also not needed since Maxwell's equations correctly predict all of the phenomenology associated with electromagnetic waves without the need for an aether.

Aether theories were disproved by the famous Michelson-Morley experiment. This experiment showed that the speed of light is the same in different reference frames through **interferometry**. We will cover interferometry later, and in lab. The idea is that there were some microscopic system whose oscillations generated electromagnetic waves, then the speed of the waves should depend on the speed of the aether. For example, if you put the shive wave machine on a truck, then the speed that the waves propagate in the machine as observed by someone on the ground depend on the speed of the truck. If you put light on a truck, it still goes at the speed of light.

We now understand that the speed of light is the same in all frames follows from a symmetry principle. That is, we start with an axiom of frame-independence, called **Lorentz invariance**, and use it to derive the constancy of the speed of light. Lorentz invariance is the generalization of spatial rotations to include time. For example, with one spatial direction  $x$  and time  $t$ , a Lorentz transformation is a boost (change of frames) which has the form

$$x' = x \cosh \beta + ct \sinh \beta \quad (32)$$

$$t' = t \cosh \beta + \frac{x}{c} \sinh \beta \quad (33)$$

where

$$\sinh \beta \equiv \frac{e^x - e^{-x}}{x} \quad \cosh \beta \equiv \frac{e^x + e^{-x}}{x} \quad (34)$$

These satisfy

$$\cosh^2 \beta - \sinh^2 \beta = 1 \quad (35)$$

Then

$$\frac{\partial}{\partial t} = \frac{\partial t'}{\partial t} \frac{\partial}{\partial t'} + \frac{\partial x'}{\partial t} \frac{\partial}{\partial x'} = \cosh \beta \frac{\partial}{\partial t'} + c \sinh \beta \frac{\partial}{\partial x'} \quad (36)$$

$$\frac{\partial}{\partial x} = \frac{\partial t'}{\partial x} \frac{\partial}{\partial t'} + \frac{\partial x'}{\partial x} \frac{\partial}{\partial x'} = \frac{1}{c} \sinh \beta \frac{\partial}{\partial t'} + \cosh \beta \frac{\partial}{\partial x'} \quad (37)$$

so

$$\frac{\partial^2}{\partial t^2} - c^2 \frac{\partial^2}{\partial x^2} = [\cosh^2 \beta - \sinh^2 \beta] \frac{\partial^2}{\partial t'^2} + c^2 [\cosh^2 \beta - \sinh^2 \beta] \frac{\partial^2}{\partial x'^2} \quad (38)$$

$$= \frac{\partial^2}{\partial t'^2} - c^2 \frac{\partial^2}{\partial x'^2} \quad (39)$$

So the wave equation is Lorentz invariant. That is, it is the same equation, with the same speed  $c$  in any frame.

# Lecture 14: Polarization

## 1 Polarization vectors

In the last lecture, we showed that Maxwell's equations admit plane wave solutions

$$\vec{E} = \vec{E}_0 e^{i(\vec{k} \cdot \vec{x} - \omega t)}, \quad \vec{B} = \vec{B}_0 e^{i(\vec{k} \cdot \vec{x} - \omega t)} \quad (1)$$

Here,  $\vec{E}_0$  and  $\vec{B}_0$  are called the **polarization vectors** for the electric and magnetic fields. These are complex 3 dimensional vectors. The wavevector  $\vec{k}$  and angular frequency  $\omega$  are real and in the vacuum are related by  $\omega = c|\vec{k}|$ . This relation implies that electromagnetic waves are dispersionless with velocity  $c$ : the speed of light. In materials, like a prism, light can have dispersion. We will come to this later.

In addition, we found that for plane waves

$$\vec{B}_0 = \frac{1}{\omega} (\vec{k} \times \vec{E}_0) \quad (2)$$

This equation implies that the magnetic field in a plane wave is completely determined by the electric field. In particular, it implies that their magnitudes are related by

$$|\vec{E}_0| = c|\vec{B}_0| \quad (3)$$

and that

$$\vec{k} \cdot \vec{E}_0 = 0, \quad \vec{k} \cdot \vec{B}_0 = 0, \quad \vec{E}_0 \cdot \vec{B}_0 = 0 \quad (4)$$

In other words, the polarization vector of the electric field, the polarization vector of the magnetic field, and the direction  $\vec{k}$  that the plane wave is propagating are all orthogonal.

To see how much freedom there is left in the plane wave, it's helpful to choose coordinates. We can always define the  $\hat{z}$  direction as where  $\vec{k}$  points. When we put a hat on a vector, it means the unit vector pointing in that direction, that is  $\hat{z} = (0, 0, 1)$ . Thus the electric field has the form

$$\vec{E} = \vec{E}_0 e^{i\omega(\frac{z}{c} - t)} \quad (5)$$

which moves in the  $z$  direction at the speed of light. Since  $\vec{E}_0$  is orthogonal to  $\vec{k}$  we can write also write

$$\vec{E}_0 = (E_x, E_y, 0) \quad (6)$$

with  $E_x$  and  $E_y$  complex numbers.

The two complex amplitudes  $E_x$  and  $E_y$  each have magnitudes and phases. We can write these explicitly as  $E_x = |E_x|e^{i\phi_x}$  and  $E_y = |E_y|e^{i\phi_y}$ . So one needs four real numbers to specify the polarization vector. We usually separate the phases into an overall phase and the difference in phase between  $E_x$  and  $E_y$ ,  $\phi = \phi_x - \phi_y$ . The reason for doing this is that as  $t$  and  $z$  change the phases of  $E_x$  and  $E_y$  change in the same way while  $\phi$  is unaffected. We don't usually care about this overall phase, or about the overall magnitude  $E = |E_x|^2 + |E_y|^2$ . Thus, to specify a polarization vector, we talk about the relative size of  $E_x$  and  $E_y$  and the phase difference  $\phi = \phi_x - \phi_y$ .

Don't let the complex polarization vectors confuse you. We can just as well

$$\vec{E} = E_x \hat{x} \cos(kz - \omega t + \phi_x) + E_y \hat{y} \cos(kz - \omega t + \phi_y) \quad (7)$$

where you see the four real numbers specifying the polarization explicitly. Exponentials just make a lot of the algebra easier.

## 2 Linear polarization

We say a plane wave is **linearly polarized** if there is no phase difference between  $E_x$  and  $E_y$ . We can write linear polarizations as

$$\vec{E}_0 = (E_x, E_y, 0) \quad (8)$$

and choose the overall phase so that  $E_x$  and  $E_y$  are real numbers. If  $E_y = 0$  but  $E_x \neq 0$ , we have

$$\vec{E} = E_0 \hat{x} e^{i(kz - \omega t)} \quad (9)$$

with  $E_0 = |\vec{E}_0|$  just a number now. Then, from Eq. (2), since  $\hat{z} \times \hat{x} = \hat{y}$ ,

$$\vec{B} = \frac{1}{c} E_0 \hat{y} e^{i(kz - \omega t)} \quad (10)$$

This configuration is said to be **linear polarized in the x direction**.

Similarly we can have

$$\vec{E} = E_0 \hat{y} e^{i(kz - \omega t)} \quad (11)$$

and using  $\hat{z} \times \hat{y} = -\hat{x}$ :

$$\vec{B} = -\frac{1}{c} E_0 \hat{x} e^{i(kz - \omega t)} \quad (12)$$

This configuration is said to be **linear polarized in the y direction**. Note that in both cases, the magnetic field is given by rotating the electric field  $90^\circ$  counterclockwise in the  $x$ - $y$  plane and dividing by  $c$ .

More generally, for any unit vector  $\hat{v}$  we can have

$$\vec{E} = E_0 \hat{v} e^{i(\vec{k} \cdot \vec{x} - \omega t)} \quad (13)$$

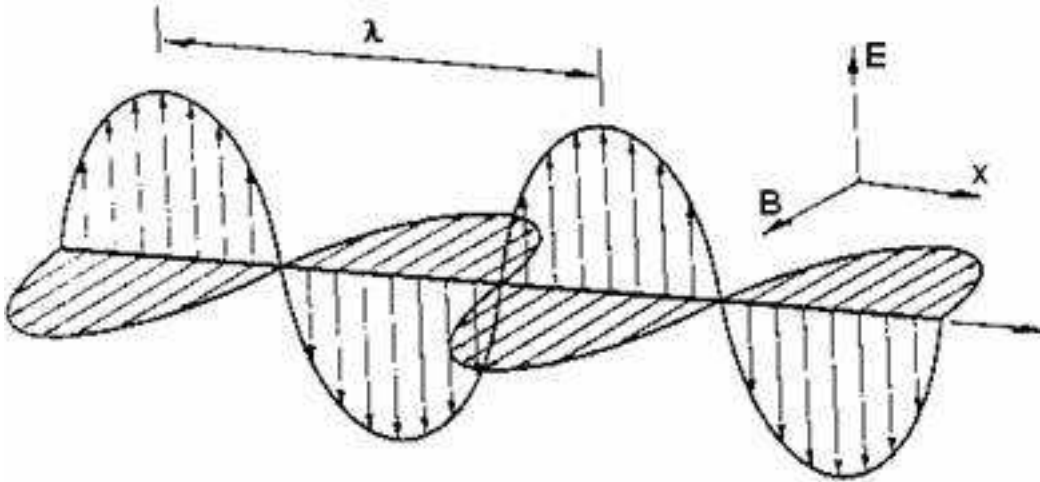
which is linearly polarized in the  $\hat{v}$  direction. The magnetic field will be polarized in direction  $90^\circ$  behind the electric field.

Remember, we always implicitly want to have the real part of these fields. So, for linearly polarized light in the  $x$  direction, the fields are actually

$$\text{Re}[\vec{E}] = (E_0 \cos(kz - \omega t), 0, 0), \quad \text{Re}[\vec{B}] = \left(0, \frac{E_0}{c} \cos(kz - \omega t), 0\right) \quad (14)$$

Note that there is no  $x$  or  $y$  dependence in these solutions: the fields are completely uniform in  $x$  and  $y$  – they are plane waves. At each point on the plane, the electric field points in the  $\hat{x}$  direction with the same magnitude. This magnitude varies as we move in  $z$  and in  $t$  but is always uniform in the plane.

Here's a picture of how the fields vary as they move along for a wave moving in the  $x$  direction:



### 3 Circular polarization

What if the components of the electric field are not in phase? First suppose they have the same magnitude but are a quarter wavelength out of phase, so  $\phi_x - \phi_y = \frac{\pi}{2}$ . Then,

$$\vec{E}_0 = (E_0, E_0 e^{i\frac{\pi}{2}}, 0) = (E_0, iE_0, 0) \quad (15)$$

Thus,

$$\vec{E} = (E_0 e^{i(kz - \omega t)}, iE_0 e^{i(kz - \omega t)}, 0) \quad (16)$$

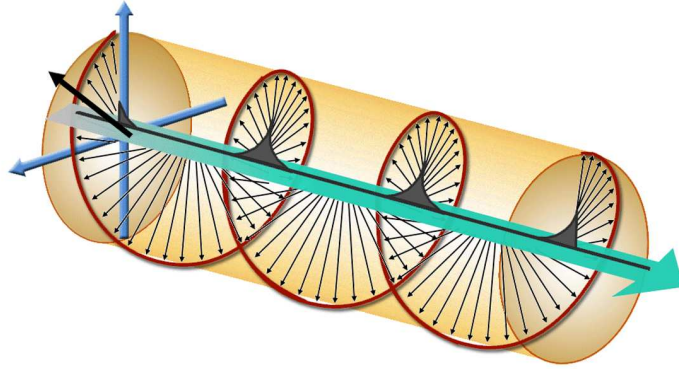
Taking the real part gives the actual electric field

$$\text{Re}[\vec{E}] = \begin{pmatrix} E_0 \cos(kz - \omega t), & -E_0 \sin(kz - \omega t), & 0 \end{pmatrix} \quad (17)$$

This is called **left-handed circularly polarized light**.

What does it look like? Well, at  $t = z = 0$ , the field is  $\text{Re}(\vec{E}) = (E_0, 0, 0)$  pointing in the  $x$  direction. A little farther along, when  $kz = \frac{\pi}{2}$ , still at  $t = 0$ , then  $\text{Re}(\vec{E}) = (0, -E_0, 0)$  which points in the negative  $y$  direction. Farther along still, when  $kz = \pi$ , it points along  $-\hat{x}$ , and then  $\hat{y}$  and so on. Finally, after a full wavelength, it goes back to  $\hat{x}$ . Thus, as we move along  $z$ , the polarization rotates clockwise in the  $x-y$  plane. Equivalently, at a given  $z$ , as time progresses it also rotates in the  $x-y$  plane.

Here is a picture



**Figure 1.** Circularly polarized light changes the direction of its polarization as it moves.

Similarly, taking  $\phi_x - \phi_y = -\frac{\pi}{2}$  gives  $\vec{E}_0 = (E_0, -iE_0, 0)$

$$\text{Re}[\vec{E}] = (E_0 \cos(kz - \omega t), E_0 \sin(kz - \omega t), 0) \quad (18)$$

which is **right-handed circularly polarized light**. In this case, at  $t = 0$  and  $kz = 0$ , the polarization points in the  $\hat{x}$  direction. A quarter wavelength farther, it points in the  $\hat{y}$  direction, and so on. So this field rotates counterclockwise in the  $x-y$  plane.

Note that when we add left and right handed polarizations we get

$$\vec{E}_0 = (E_0, iE_0, 0) + (E_0, -iE_0, 0) = (2E_0, 0, 0) \quad (19)$$

which is linearly polarized in the  $\hat{x}$  direction. Similarly, subtracting them gives linear polarization in the  $\hat{y}$  direction. Thus circular and linear polarizations are not linearly independent. Indeed, any possible polarization can be written as a linear combination of left-handed and right-handed circularly polarized light.

Circularly polarized light has angular momentum  $\pm e_0 |\vec{E}|^2 \hat{k}$ , with the positive sign for left-handed and the negative sign for right-handed. Linearly polarized light carries no angular momentum. This is easy to understand because it is a sum of left and right handed light. As an analogy, suppose you have two tops, one spinning clockwise and one spinning counterclockwise; the two-top system has no net angular momentum.

As mentioned before, light cannot have arbitrarily small intensity. The smallest intensity light can have is a single photon. Thus the photon itself must be polarized. A single circularly polarized photon has angular momentum  $\vec{J} = \pm \frac{h}{2\pi} \hat{k}$ . It may seem surprising that photons which are pointlike particles with no substructure can have angular momentum on their own. It is a very odd fact, but true nonetheless. Photons have **spin**. We say photons are particles of spin 1. The two signs for the angular momentum correspond to the two helicities of the photon.

To complete the discussion of polarization, we can consider also varying the magnitudes of the different components and the phase. The most general parameterization is

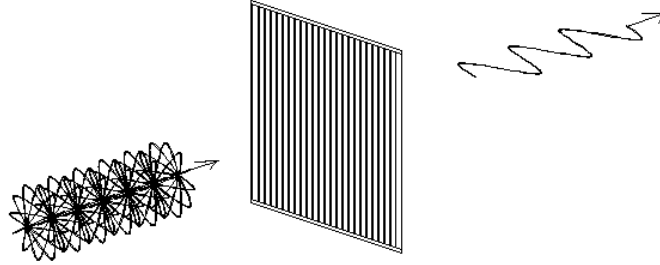
$$\vec{E}_0 = E_x \hat{x} \cos(kz - \omega t) + E_y \hat{y} \cos(kz - \omega t + \phi) \quad (20)$$

When  $E_x$  and  $E_y$  are different and the relative phase  $\phi$  is nonzero, the polarization changes in magnitude as it rotates in the  $x$ - $y$  plane. It thereby describes an ellipse, and this is called **elliptical polarization**. Linear polarization corresponds to  $\phi = 0$ . Circular polarization corresponds to  $\phi = \pm \frac{\pi}{2}$  and  $E_x = E_y$ .

## 4 Polaroid film

Now that we understand the mathematics of polarizations, what is the physics? How do we actually produce polarized light?

One way to polarize light is using a **polaroid film**. The first such film was invented by Edwin Land in 1928, while an undergraduate at Harvard. He came up with a way to align polymer molecules in a thin sheet into long needlelike strands. When an electric field acts on the electrons in those strands, it can only move the electrons up and down in the strand direction, but not perpendicular to the strands. Thus the only polarization which can pass through is linearly polarized **in the direction perpendicular to the strips**. Here is a figure



**Figure 2.** Polarizing film absorbs the electric field in the direction of the strips.

Land called his invention **Polaroid film**, patented it, and then founded the Polaroid corporation here in Cambridge. One of the primary applications at the time was sunglasses (see below). His next most famous invention was the Land instant camera, now known as the polaroid camera. The polaroid camera has nothing to do with polarization!

Back to polarizing film. Say the film has the strips in the  $y$  direction and a plane wave comes in along the  $z$  direction with an arbitrary polarization

$$\vec{E}_{\text{init}} = E_x \hat{x} e^{i(kz - \omega t)} + E_y \hat{y} e^{i(kz - \omega t + \phi)} \quad (21)$$

Since the polarizer has strips in the  $y$  direction, it absorbs the  $\hat{y}$  components of the electric field. Thus, what exits the polarizer must be

$$\vec{E}_{\text{final}} = E_x \hat{x} e^{i(kz - \omega t)} \quad (22)$$

Thus the polarizer takes any initial polarization and turns it into linear polarization in the  $x$  direction.

For example, say the field is linearly polarized at an angle  $\theta$  to the direction where the field would just go through. That is,

$$\vec{E}_{\text{init}} = E_0 (\cos\theta \hat{x} + \sin\theta \hat{y}) e^{i(kz - \omega t)} \quad (23)$$

Then if we put it through the polarizer with strips in the  $y$  direction we would find

$$\vec{E}_{\text{init}} = E_0 \cos\theta \hat{x} e^{i(kz - \omega t)} \quad (24)$$

So that the result is still linearly polarized but now in the  $\hat{x}$  direction with magnitude  $E_0 \cos\theta$ . Since the intensity of radiation is proportional to the square of the field, we would find that a polarizer oriented at an angle  $\theta$  to the direction of propagation reduces the intensity

$$I_{\text{final}} = I_{\text{initial}} \cos^2\theta \quad (25)$$

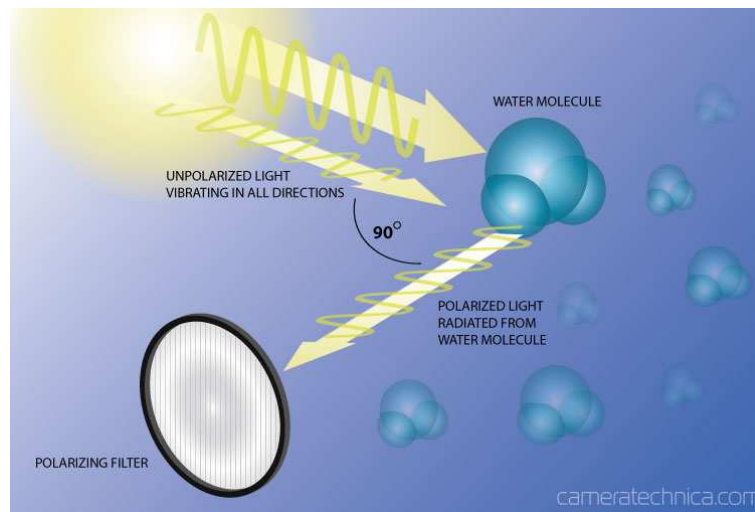
This is known as **Malus' law**.

## 5 Polarization in the real world

The key facts which let us understand why natural light is polarized are 1) that the electric field makes charged particles move in the direction of  $\vec{E}$  and 2) moving charged particles produce electric fields along their direction of motion. Fact 1) and the consequent absorption of energy of the electric field is the mechanism behind the polaroid film we just discussed.

During the day, sunlight is constantly being absorbed and emitted by molecules (mostly water) in the atmosphere. Let's say  $z$  is the vertical direction (towards the sky), so the sky above us is the  $x$ - $y$  plane. Let's say the sun is in the  $x$  direction. Thus plane waves from the sun have  $\vec{k}$  in the  $\hat{x}$  direction and are polarized in the  $y$ - $z$  plane. Thus, the sky molecules can only get pushed in the  $y$  and  $z$  directions from sunlight. Molecules moving in the  $z$  direction emit light in the  $x$ - $y$  plane, that is, off into the sky, but not to us. Molecules excited in the  $y$  direction however can emit light going directly off in  $z$  down to us. Since all this light comes from motion in the  $y$ - $z$  plane, the light will all be linearly polarized in this plane. Thus the sky directly above us, at right angles to the sun, is polarized.

Similar arguments show that light from some other angle will have some polarization, but not be completely polarized. Looking through the sky directly at the sun, the light should not be polarized at all.



**Figure 3.** Polarization of the sky

Another important example is light from a reflection. Here the story is nearly identical, as shown in Fig 4 with the maximum polarization coming if the source is at a right angle to the viewer. In this case, the index of refraction of the two materials involved (air and lake for example) play a role and the angle is not exactly 90 degrees. Instead it is called Brewster's angle. We will derive a formula for Brewster's angle in Lecture 16 on reflection.



**Figure 4.** Polarization from reflection.

So reflected light is polarized. Most light around us is not polarized. So for example, if you are on a boat and light is coming in from the sky (not just the sky directly above you), and also from the water, only the light reflected off the water will be polarized. Thus if you put on polarizing sunglasses, you can filter out the reflected light, and see more clearly the things around you. Similarly, reflected light coming off the ground will be polarized horizontally. Polarizing sunglasses generally try to filter out horizontally polarized light.

## 6 Index of refraction

The speed of light  $c$  is the speed of light in the vacuum. In air, light goes slightly slower than  $c$ . More generally, light travels at a speed

$$v = \frac{1}{n}c \quad (26)$$

with  $n$  called the **index of refraction**. The physical origin of the index of refraction is interesting. Light comes in and excites molecules. These molecules then radiate more light which interferes with the incoming light. The phases of the incoming and radiated light align in such a way as to make the light appear to be slower. We will work this out in detail in Lecture 17. For now, let's accept the index of refraction as a phenomenological fact.

When light hits a boundary where the index of refraction changes, we have to solve the wave equation with the boundary conditions. We will work out the transmission and reflection coefficients (which are polarization dependent) in Lecture 16. All we need to know for now is that the incident and transmitted wave have the same frequency (this is always true – the boundary knows about the frequency but not the wavelength of the incoming wave).

Say light with frequency  $\omega$  goes from a vacuum with  $n = 1$  to a region with  $n = 2$  and then back to the vacuum with  $n = 1$ . Since the frequency is fixed and

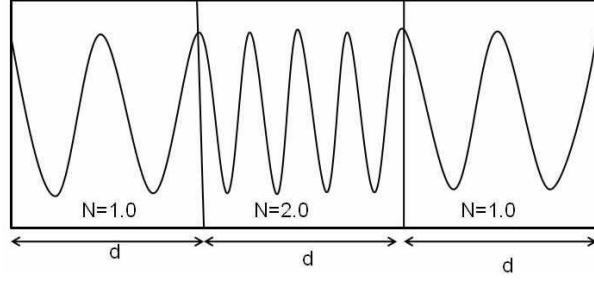
$$\omega = vk = \frac{c}{n} \frac{2\pi}{\lambda} \quad (27)$$

we must have

$$n_1 \lambda_1 = n_2 \lambda_2 \quad (28)$$

So that as  $n$  goes up (at fixed frequency),  $\lambda$  must go down. The picture at a fixed time looks like this:



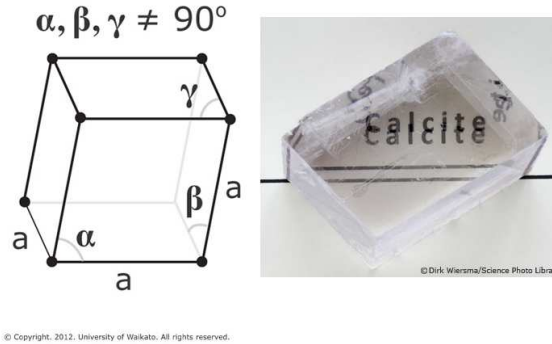


**Figure 5.** Plane wave entering and emerging from a medium with different index of refraction. This figure illustrates how light going slower means it goes through more oscillations over the same distance.

Many materials are **birefringent**, meaning they have different indices of refraction in different directions. A linear polarizer is an extreme example: it has  $n = \infty$  in one direction, so that light simply does not propagate for one polarization. The mineral calcite is a naturally occurring birefringent material. Calcite's two indices of refraction are

$$n_1 = 1.66, \quad n_2 = 1.49 \quad (29)$$

These indices of refraction are determined by the atomic structure of the mineral and associated with the primary axes of the crystal structure. You can see the birefringence by looking at some writing through a chunk of calcite. You will see two images. The two images are due to the different polarizations of light refracting differently in the crystal (refraction is the topic of Lecture 16). As you rotate the crystal, you can see how one polarization is refracting more than the other.



**Figure 6.** Calcite is naturally birefringent. When looking through a calcite crystal you see two images because the two polarizations of light refract differently.

## 7 Quarter and half-wave plates

Suppose we have a plane wave moving in the  $\hat{z}$  direction whose electric field is polarized in the  $\hat{x} + \hat{y}$  direction. Say this wave passes through a calcite crystal of thickness  $L$  whose optical axes are aligned with the  $\hat{x}$  and  $\hat{y}$  directions. The initial electric field is

$$\vec{E}_{\text{init}}(z, t) = E_0 (\hat{x} + \hat{y}) e^{i\omega(t - \frac{n}{c}z)} \quad (30)$$

where  $n$  is the index of refraction in air. Inside the crystal, the two components propagate with different speeds, so (assuming zero reflection)

$$\vec{E}_{\text{inside}}(z, t) = E_0 \hat{x} e^{i\omega(t - \frac{n_x}{c}z)} + E_0 \hat{y} e^{i\omega(t - \frac{n_y}{c}z)} \quad (31)$$

where  $n_x$  and  $n_y$  are the two indices of refraction of the calcite.

In particular, at  $z = L$  we have

$$\vec{E}_{\text{inside}}(L, t) = E_0 \hat{x} e^{i\omega(t - \frac{n_x L}{c})} + E_0 \hat{y} e^{i\omega(t - \frac{n_y L}{c})} \quad (32)$$

Thus there is a net phase difference of

$$\Delta\phi = \frac{\omega}{c} L (n_x - n_y) \quad (33)$$

between the  $x$  and  $y$  components of a field. Then when the light exits and both indices of refraction are  $n$  again, the outgoing wave has electric field

$$\vec{E}_{\text{final}} = E_0 (\hat{x} + \hat{y} e^{i\Delta\phi}) e^{i\omega(t - \frac{n}{c} z)} \quad (34)$$

From this calculation, we can now deduce how to use calcite to turn linearly polarized light into circularly polarized light. To do so, we must choose the length  $L$  so that  $\Delta\phi = \frac{\pi}{2}$ . This is easy to do if we know  $n_x$  and  $n_y$  using Eq. (33). A material which rotates the phase by  $\frac{\pi}{2}$ , which is a quarter of a circle, is known as a **quarter-wave plate**.

A few things to note about quarter wave plates or more generally using birefringent materials to polarize light:

- The phase  $\Delta\phi$  depends linearly on the frequency  $\omega$ . So you need different lengths for different frequencies. In most materials,  $n_x$  and  $n_y$  also depend on  $\omega$ .
- Light which is linearly polarized along one of the axes of the crystal will remain linearly polarized. That is, the orientation of the quarter wave plate is important.

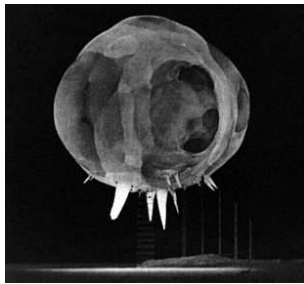
In practice, if one has a coherent, monochromatic, polarized source, such as a laser, one knows the frequency and the direction of polarization. Then one can rotate between linear and circular polarization using wave plates.

If one places two quarter-wave plates together, or equivalently takes a crystal or length  $2L$ , one gets a **half-wave plate**. The first  $L$  of it adds a phase  $e^{i\frac{\pi}{2}} = i$ , rotating linear polarization to circular polarization; the second  $L$  adds another phase  $e^{i\frac{\pi}{2}} = i$ , rotating circular back to linear. However, the net phase changes is  $e^{i\pi} = -1$  so  $(E_x, E_y) \rightarrow (E_x, -E_y)$  and the final polarization is now rotated  $90^\circ$  from the original polarization. Thus half-wave plates can rotate the polarization of linearly polarized light. A crystal of length  $4L$  will have  $\Delta\phi = 2\pi$  and return the polarization to its initial value. This is called a **full-wave plate**.

One application of the use of circularly polarized light is in 3D movies. Ever wonder what those glasses they give you at IMAX do? The IMAX theater projects two movies at once, one with left-hand circular polarization and the other with right-handed circular polarization. Your 3D glasses have one lens which can see transmit the left-handed light and the other lens transmits only right-handed. Thus you see different images with each eye which gives the illusion of depth. Why do you think they don't use linearly polarized light?

For another application, consider that the index of refraction along axes in some materials can be controlled by applying electric and magnetic fields. In the case of electric fields, this effect is called **electro-optic effect**. Then the index of refraction is  $n(\vec{E})$ . In the case of magnetic fields, this effect is called the Faraday effect, and the index of refraction is  $n(\vec{B})$ . These effects are exploited in electro-optic devices which allow fast and precise control of the polarization of light.

An example of an electro-optic device is the rapatronic camera. In an ordinary camera, where the shutter is mechanical, the shutter speed is limited to about 1 millisecond. In a rapatronic camera, the shutter is composed of two crossed polarizers with a birefringent material exhibiting the electro-optic effect in between the polarizers. With the voltage on, the material acts as a half-wave plate so that all the light is transmitted through the shutter. If the voltage is switched off, which can be done on the order of microseconds or faster, the material has no effect on the polarization. This causes the shutter to be closed. Rapatronic cameras were heavily developed during the 1940's to take pictures of atomic bomb detonations very shortly ( $10^{-6}$  s) after the detonation. This allowed scientists to learn about how well the bomb, in particular its detonation, was working. Here's a picture from a rapatronic camera:

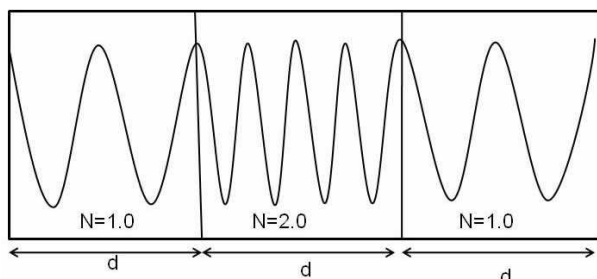


**Figure 7.** Picture of the detonation of an atomic bomb exposed with a rapatronic camera for 3 millionths of a second.

# Lecture 15: Refraction and Reflection

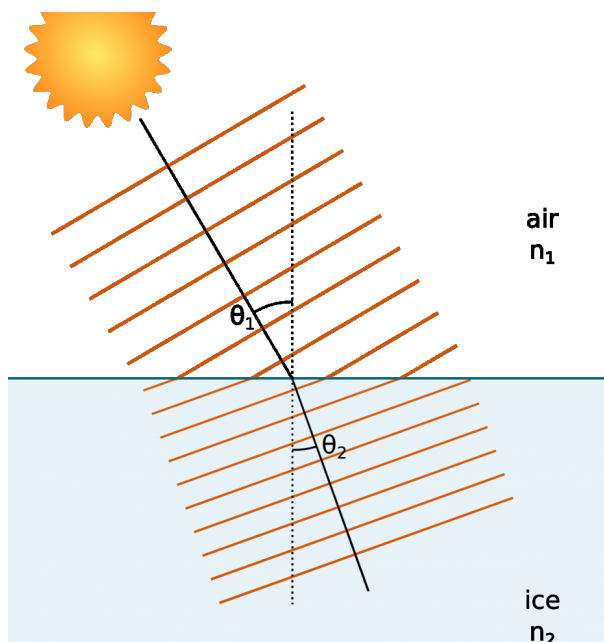
## 1 Refraction

When we discussed polarization, we saw that when light enters a medium with a different index of refraction, the frequency stays the same but the wavelength changes. Using that the speed of light is  $v = \frac{c}{n}$  we deduced that  $\lambda_1 n_1 = \lambda_2 n_2$ , so that as the index of refraction goes up, the wavelength goes down. The picture looks like this



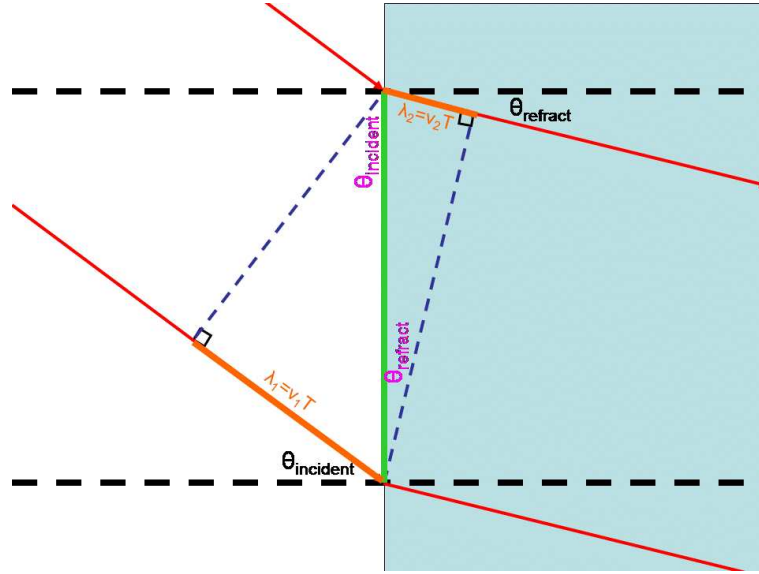
**Figure 1.** Plane wave entering and emerging from a medium with different index of refraction.

Now let us ask what happens when light enters a medium with a different index of refraction at an angle. Since we know the wavelength of light in the two media, we can deduce the effect with pictures. The key is to draw the plane waves as the location of the maximum field values. These crests will be straight lines, but spaced more closely together in the medium with higher index of refraction. For example, if sunlight hits ice (or water), the picture looks like this



**Figure 2.** Matching wavefronts demonstrates refraction. Orange lines represent the crests of waves (or the maximum amplitude)

The bending of light when the index of refraction changes is called **refraction**. To relate the angles  $\theta_1$  (the **angle of incidence**) to  $\theta_2$  (the **angle of refraction**) we draw a triangle



**Figure 3.** Light comes in from the air on the left with the left dashed blue line indicating one wavecrest with the previous wavecrest having just finished passing into the water. Thus the wavelength  $\lambda_1$  in medium 1 is the solid thick orange line on the bottom left, and the wavelength in the second medium  $\lambda_2$  is the thick solid orange line on the top right.

Call the distance between the places where the wave crest hits the water along the water  $R$  (the thick green vertical line in the picture). The distance between crests is  $\lambda_1 = R \sin \theta_1$  in the air and  $\lambda_2 = R \sin \theta_2$  in the water. Since  $R$  is the same, trigonometry and  $n_1 \lambda_1 = n_2 \lambda_2$  imply

$$\boxed{\frac{n_1}{n_2} = \frac{\lambda_2}{\lambda_1} = \frac{\sin \theta_2}{\sin \theta_1}} \quad (1)$$

This is known as **Snell's law**.

The same logic holds for reflected waves:  $R$  is the same and  $\lambda$  is the same (since  $n_1 = n_2$  for a reflection) therefore  $\theta_1 = \theta_2$ . This is usual the **law of reflection**: the angle of reflection is equal to the angle of incidence.

For a fast-to-slow interface (like air to water), the angle gets smaller (the refracted angle is less than the incident angle). For a slow-to-fast interface (like water to air), the angle gets larger. Since the angle cannot be larger than  $90^\circ$  while remaining in the second medium, there is a largest incident angle for refraction when  $n_2 < n_1$ . In equations, since the refracted angle satisfies  $\sin \theta_2 = \frac{n_1}{n_2} \sin \theta_1$ , we see that if  $\frac{n_1}{n_2} \sin \theta_1 > 1$  there is no solution. The **critical angle** beyond which no refraction occurs is therefore

$$\boxed{\theta_c = \sin^{-1} \frac{n_2}{n_1}} \quad (2)$$

For air  $n \approx 1$  and for water  $n = 1.33$  so  $\theta_c = 49^\circ$ . For incident angles larger than the critical angle, there is no refraction: all the light is reflected. We call this situation **total internal reflection**. Total internal reflection can only happen if  $n_2 < n_1$ . Thus, light can be confined to a material with higher index of refraction but not a lower one.

Total internal reflection is the principle behind fiber optics. A fiber optical cable has a solid silica core surrounding by a cladding with an index of refraction about 1% smaller. For example, the core might have  $n_1 \approx 1.4475$  and the cladding  $n_2 = 1.444$ , so the critical angle is  $\theta_c = 86^\circ$  from normal incidence, or  $4^\circ$  from the direction of propagation. As long as the cable is not bent too much (typically the cable is thick enough so that this is very difficult), the light will just bounce around in the cable, never exiting, with little loss. Why might you want to send signals with visible light rather than radio waves?

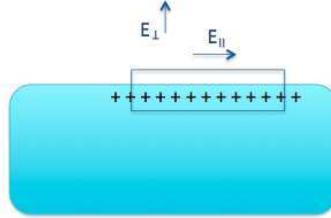
In the fiber optic cable, the high index of refraction is surrounded by a lower index of refraction, or equivalently the lower wave speed medium is surrounded by higher wave speed medium. You can always remember this through Muller's analogy with the people holding hands and walking at different speeds – if slow people are surrounded by fast people, the fast will bend in to the slow. The same principle explains the SOFAR sound channel in the ocean and the sound channel in the atmosphere. Recall that for a gas, the speed of sound is  $c_s = \sqrt{\gamma \frac{RT}{m}}$ , thus as the temperature goes down, the speed of sound goes down. However when you go high enough to hit the ozone layer, the temperature starts increasing. This is because ozone absorbs UV light, converting it to heat. Thus speed of sound has a minimum, and there is a sound channel.

## 2 Boundary conditions

Snell's law holds for any polarization. It determines the direction of the transmitted fields. It does not determine the magnitude. In fact, the magnitude depends on the polarization (as we already know, since reflections are generally polarized). To determine the magnitude(s), we need the boundary conditions at the interface.

You might think  $\vec{E}_1 = \vec{E}_2$  and  $\vec{B}_1 = \vec{B}_2$  are the right boundary conditions. However, if charge accumulates on the boundary, as in a conductor, the fields outside and inside the material will have to be different. If current accumulates, for example through a growing number of little swirling eddies of charge, the magnetic field will be different.

How much charge or current accumulates is determined by the electric permittivity  $\epsilon$  and magnetic permeability  $\mu$  of a material. Recall that in a vacuum, these reduce to  $\epsilon_0$  and  $\mu_0$  and that the speed of light in the vacuum is  $c = \sqrt{\epsilon_0 \mu_0}$ . In a medium, we have to replace all the  $\epsilon_0$  and  $\mu_0$  factors in Maxwell's equations with  $\epsilon$  and  $\mu$ . So,  $v = \sqrt{\mu\epsilon}$  in a material and  $|\vec{B}| = \frac{1}{v}|\vec{E}|$ . Moreover, since one of Maxwell's equations is  $\vec{\nabla} \times \frac{\vec{B}}{\mu} = \frac{\partial}{\partial t}(\epsilon \vec{E})$ , it is natural to work with the rescaled fields  $\vec{D} = \epsilon \vec{E}$  and  $\vec{H} = \frac{1}{\mu} \vec{B}$ .



**Figure 4.** Charge accumulating on a boundary can affect  $E_{\perp}$  not  $E_{\parallel}$ .

To figure out the effect of the charge, we can use Gauss's law. Drawing a little pillbox around the charges as in Fig 4, we see the only the field perpendicular to the interface can be affected. Let's call this  $E_{\perp}$ . You should remember from 15b that the way electric permittivity works is that it lets you use Gauss's law as if there is no charge, provided you integrate  $\vec{D} = \epsilon \vec{E}$  over the surface rather than  $\vec{E}$ . Therefore the boundary condition is  $D_{\perp}^{(1)} = D_{\perp}^{(2)}$  which implies

$$\epsilon_1 E_{\perp}^{(1)} = \epsilon_2 E_{\perp}^{(2)} \quad (3)$$

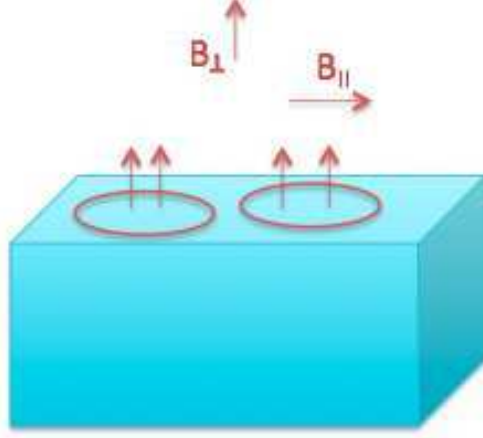
where  $\epsilon_1$  is the dielectric constant in the medium where the transverse electric field is  $E_{\perp}^{(1)}$  and  $\epsilon_2$  is the dielectric constant in the medium where the transverse electric field is  $E_{\perp}^{(2)}$ .

Since the parallel electric field is unaffected by the charges, we also have

$$E_{\parallel}^{(1)} = E_{\parallel}^{(2)} \quad (4)$$

There are no factors of  $\epsilon$  here because the accumulated charge has no effect on  $E_{\parallel}$ .

For the magnetic field, the boundary conditions can be affected due to the magnetic moment of the particles in the medium, as encoded in  $\mu$ . This picture looks like



**Figure 5.** Current on the boundary can affect  $B_{\parallel}$  not  $B_{\perp}$ .

Since a current induces a field perpendicular to the current, only  $B_{\parallel}$  can be affected, not  $B_{\perp}$ . Thus,

$$B_{\perp}^{(1)} = B_{\perp}^{(2)} \quad (5)$$

Since the part of the magnetic field which is sensitive to an accumulated current in a material is treated using  $\vec{H} = \frac{1}{\mu} \vec{B}$  the condition is then that

$$\frac{1}{\mu_1} B_{\parallel}^{(1)} = \frac{1}{\mu_2} B_{\parallel}^{(2)} \quad (6)$$

In summary, the boundary conditions are

$\epsilon_1 E_{\perp}^{(1)} = \epsilon_2 E_{\perp}^{(2)}$	$B_{\perp}^{(1)} = B_{\perp}^{(2)}$
$E_{\parallel}^{(1)} = E_{\parallel}^{(2)}$	$\frac{1}{\mu_1} B_{\parallel}^{(1)} = \frac{1}{\mu_2} B_{\parallel}^{(2)}$

(7)

In the vacuum,  $\mu = \mu_0 = 1.25 \times 10^{-6} \frac{H}{m}$  (Henries per meter). If  $\mu$  is much larger than this, the material can acquire a large current with a small magnetic field. That is, it conducts. For example, iron has  $\mu = 6.3 \times 10^{-3}$ . Conductors are opaque and therefore not of much interest for the study of refraction. In most transparent materials,  $\mu \approx \mu_0$ . For example, water has  $\mu = 1.256 \times 10^{-6} \frac{H}{m}$  which is the same as for air up to 6 decimal places. In transparent materials, the electric permittivity can vary significantly. So we will assume that  $\mu \approx \mu_0$  and  $\epsilon$  varies. In this case, the boxed equations reduce to  $\epsilon_1 E_{\perp}^{(1)} = \epsilon_2 E_{\perp}^{(2)}$ ,  $E_{\parallel}^{(1)} = E_{\parallel}^{(2)}$  and  $\vec{B}^{(1)} = \vec{B}^{(2)}$ .

Next, we'll solve the wave equation with these boundary conditions, much like we solved the transmission and reflection problem for waves in a string in Lecture 9 which led to the concept of impedance.

### 3 Normal incidence

Let's start with the case of normal incidence (perpendicular to the interface). To find out how much light is reflected, we need to work out the impedance, which determines the reflection and transmission coefficients.

For normal incidence the incident angle to the normal  $\theta_1 = 0$  (see Fig. 8 below). Thus  $E_\perp = B_\perp = 0$ . That is, the electric and magnetic fields are both polarized in the plane of the interface so there simply is no  $\perp$  component. Without loss of generality, let's take a plane wave of frequency  $\omega$  moving in the  $\hat{z}$  direction with  $\vec{E}$  in the  $\hat{y}$  direction. Then since  $\vec{B} = \frac{1}{\omega} \vec{k} \times \vec{E}$ , the magnetic field points in the  $\hat{x}$  direction. So the incident fields are

$$\vec{E}_I = E_I \hat{y} e^{i\omega(t - \frac{1}{v_1}z)}, \quad \vec{B}_I = B_I \hat{x} e^{i\omega(t - \frac{1}{v_1}z)} \quad (8)$$

The transmitted and reflected waves are also moving normal to the surface (by Snell's law), so we can write

$$\vec{E}_T = E_T \hat{y} e^{i\omega(t - \frac{1}{v_1}z)} \quad \vec{B}_T = B_T \hat{x} e^{i\omega(t - \frac{1}{v_1}z)} \quad (9)$$

$$\vec{E}_R = E_R \hat{y} e^{i\omega(t + \frac{1}{v_1}z)} \quad \vec{B}_R = -B_R \hat{x} e^{i\omega(t + \frac{1}{v_1}z)} \quad (10)$$

with  $E_T$ ,  $E_R$ ,  $B_T$  and  $B_R$  the transmission and reflection coefficients to be determined. Note that we have flipped the sign on the  $z$  term in the phase since the reflected wave is moving in the  $-\hat{z}$  direction (that is,  $\vec{k}$  flips). Since  $\omega \vec{B} = \vec{k} \times \vec{E}$ , this means that  $\vec{B}$  flips sign too, which explains the minus sign in Eq. (10).

Since there is no  $\perp$  component the boundary condition  $E_{||}^{(1)} = E_{||}^{(2)}$  implies that

$$E_I + E_R = E_T \quad (11)$$

Similarly,  $\frac{1}{\mu_1} B_{||}^{(1)} = \frac{1}{\mu_2} B_{||}^{(2)}$  implies

$$\frac{1}{\mu_1} B_I - \frac{1}{\mu_1} B_R = \frac{1}{\mu_2} B_T \quad (12)$$

Since  $|\vec{B}| = \frac{1}{v} |\vec{E}|$  for any plane wave,  $B_I = \frac{1}{v_1} E_I$  and  $B_T = \frac{1}{v_2} E_T$  and Eq. (12) becomes

$$\frac{1}{\mu_1 v_1} (E_I - E_R) = \frac{1}{\mu_2 v_2} E_T \quad (13)$$

Eqs. (11) and (13) look just very much like the equations for transmission and reflection for a wave that we studied in Lecture 9. Solving them gives

$$E_R = E_I \frac{Z_2 - Z_1}{Z_2 + Z_1}, \quad E_T = E_I \frac{2Z_2}{Z_2 + Z_1} \quad (14)$$

where

$$Z = \mu v = \mu \frac{c}{n} = \sqrt{\frac{\mu}{\epsilon}} \quad (15)$$

For most materials,  $\mu$  is pretty constant, so the form  $Z = \mu c \frac{1}{n}$ , which says  $Z \propto \frac{1}{n}$ , is most useful

What is the power reflected and transmitted? The incident power in a medium is given by the Poynting vector  $\vec{P} = \frac{1}{\mu} \vec{E} \times \vec{B}$  times area  $A$ . So

$$|\vec{P}_I| = \frac{1}{\mu v} |\vec{E}_I|^2 A = \frac{1}{Z} |\vec{E}_I|^2 A \quad (16)$$

Thus the reflected power is

$$P_R = \frac{E_R^2}{Z_1} A = \left( \frac{Z_2 - Z_1}{Z_2 + Z_1} \right)^2 \frac{E_I^2}{Z_1} A = \left( \frac{Z_1 - Z_2}{Z_1 + Z_2} \right)^2 P_I \quad (17)$$

and the transmitted power is

$$P_T = \frac{E_T^2}{Z_2} A = \left( \frac{2Z_2}{Z_2 + Z_1} \right)^2 \frac{E_I^2}{Z_2} A = \frac{4Z_1 Z_2}{(Z_1 + Z_2)^2} \frac{E_I^2}{Z_1} A = \frac{4Z_1 Z_2}{(Z_1 + Z_2)^2} P_I \quad (18)$$

These satisfy  $P_T + P_R = P_I$  as expected. Note that these equations hold for **normal incidence only**.



By the way, Eq. (15) implies that empty space has impedance of  $Z_0 = \sqrt{\frac{\mu_0}{\epsilon_0}} = 376.7\Omega$ . This number is very useful in broadcasting, since you want to impedance-match your antenna to air to get efficient signal transmission. Coaxial cables usually have  $Z = 50\Omega$  so some circuitry is required to get the antenna impedance higher. Antennas and interference patterns are the subject of Lecture 18.

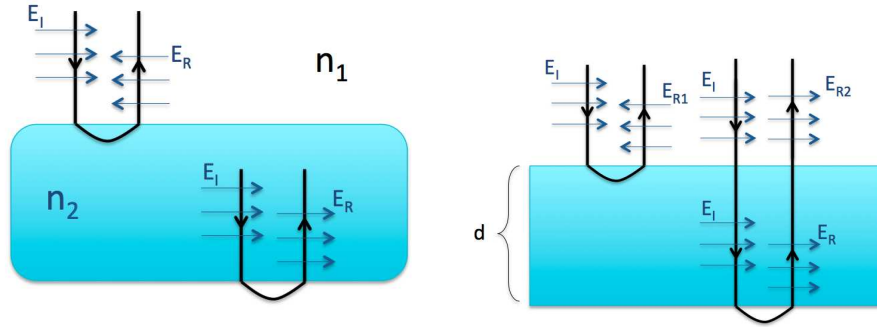
## 4 Thin film interference

Eq. (14) is  $E_R = E_I \frac{Z_2 - Z_1}{Z_2 + Z_1}$ . Using  $Z \propto \frac{1}{n}$ , this becomes

$$\frac{E_R}{E_I} = \frac{n_1 - n_2}{n_1 + n_2} \quad (19)$$

Thus for  $n_2 > n_1$  the reflected electric field has a phase flip compared to the incoming field, and for  $n_1 < n_2$  there is no phase flip.

So what happens when light hits a thin film? The film will generally have  $n_2 = n_{\text{film}} > n_{\text{air}} = n_1$ . Thus we get a phase flip at the top surface, but no phase flip at the bottom surface. The picture looks like this



**Figure 6.** There is a  $\pi$  phase flip from reflections off the top surface, where  $n_2 > n_1$ , but not off the bottom surface, where  $n_2 < n_1$ .

So what happens if light of wavelength  $\lambda$  hits a film of thickness  $d$ ? There will be two reflected waves, one off the top surface ( $A$ ) and one off the bottom ( $B$ ) which will interfere. Say the incoming electric field is  $E_I \cos(kz - \omega t)$ , pointing to the right. Then the wave  $A$  which reflects off the top surface ( $z=0$ ) will be

$$E_A = R E_I \cos(-\omega t - \pi) = -R E_I \cos(-\omega t) \quad (20)$$

with  $R$  the reflection coefficient. This now points to the left. The wave which gets through will be  $E T \cos(\frac{2\pi}{\lambda} z - \omega t)$ , with  $T$  the transmission coefficient. This wave then goes to the bottom surface, reflects with no phase flip, and comes back and exits with no more phase flips. So when it exits, it is back to  $z=0$  after having traversed a distance  $\Delta z = 2d$ . Thus it is

$$E_B = E_I T^2 R \cos\left(-\omega t + 2\pi \frac{2d}{\lambda}\right) \quad (21)$$

The total wave at  $t=0$  therefore

$$E_{\text{tot}} = E_I R \left[ -1 + T^2 \cos\left(\frac{4\pi d}{\lambda}\right) \right] \quad (22)$$

where  $\lambda$  is the wavelength in the second medium.

In the limit that  $d \ll \lambda$  the two reflections will be exactly out of phase. For  $T \approx 1$ , there will be complete destructive interference and no reflection. But there will also be complete destructive interference whenever  $\cos\left(\frac{4\pi d}{\lambda}\right) = 1$ , which is when

$$d = \frac{\lambda}{2}, \frac{2\lambda}{2}, \frac{3\lambda}{2}, \frac{4\lambda}{2}, \dots \quad (\text{complete destructive interference}) \quad (23)$$

On the other hand if the two waves are completely in phase there will be constructive interference. This happens when  $\cos\left(\frac{4\pi d}{\lambda}\right) = -1$  which is when

$$d = \frac{\lambda}{4}, \frac{3\lambda}{4}, \frac{5\lambda}{4}, \dots \quad (\text{complete constructive interference}) \quad (24)$$

If the material is much thicker than the wavelength of light, and not of completely uniform thickness, then there will be some constructive and some destructive interference and we won't see much interesting. However, if the material has a well-defined thickness which is of the same order of magnitude as the wavelength of visible light, we will see different wavelengths with different intensities. This happens in a soap film.

If we put a soap film vertically, then gravity will make it denser at the bottom. The result is the following:



Note at the very top, the film is black because there is complete destructive interference for all wavelengths. Similar patterns can be seen in soap bubbles.



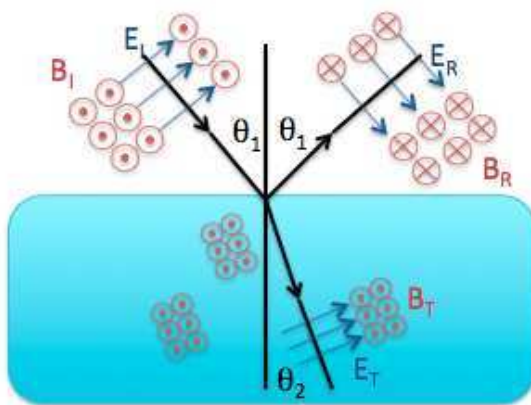
Color due to thin-film interference is known as **iridescence**. The color of many butterflies and the gorgets of hummingbirds is due to iridescence.



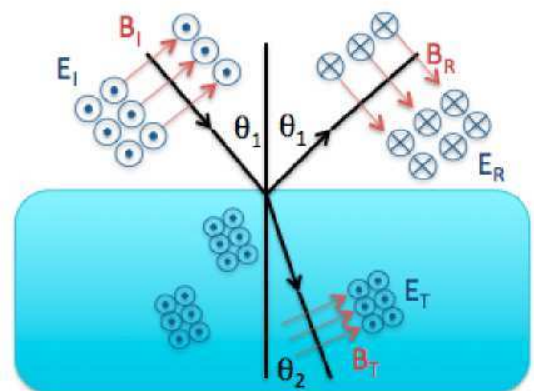
**Figure 7.** Hummingbirds and butterflies get some of their color from thin-film interference.

## 5 Fresnel coefficients

Now let's consider the more general case. Suppose we have a plane wave moving in the  $\vec{k}$  direction towards a surface with normal vector  $\vec{n}$ . Thus the angle  $\theta_1$  that the wave is coming in at satisfies  $\vec{n} \cdot \vec{k} = |\vec{k}| |\vec{n}| \cos \theta_1$ . The two vectors  $\vec{k}$  and  $\vec{n}$  form a plane. There are two linearly independent possibilities for the polarization: the electric field can be polarized in the  $\vec{k}$ - $\vec{n}$  plane of the boundary, or transverse to that plane. We call these **vertical** and **horizontal polarizations** respectively. Snell's law then tells us the directions of the transmitted and reflected electric fields. The magnetic field is always determined by the electric field through  $\omega \vec{B} = \vec{k} \times \vec{E}$ . What we need to solve for is the amplitude of the transmitted and reflected fields. The two cases look like:



Vertical polarization  
(E in plane)



Horizontal polarization  
(B in plane)

**Figure 8.** Two linearly independent linear polarizations are vertical and horizontal. Crosses indicate vectors point into the page, and dots that the vectors come out.

Let's start with vertical polarization. From Fig. 8 we see that  $\vec{B}$  is always parallel to the surface, so  $B_{\perp} = 0$ . The electric field has  $E_{\parallel}^I = E^I \cos \theta_1$ ,  $E_{\parallel}^R = E^R \cos \theta_1$  and  $E_{\parallel}^T = E_T \cos \theta_2$  as well as  $E_{\perp}^I = E^I \sin \theta_1$ ,  $E_{\perp}^R = -E^R \sin \theta_1$  and  $E_{\perp}^T = E_T \sin \theta_2$ . You can check that when  $\theta_1 = \theta_2 = 0$ , this reduces to the normal incidence case, so we have the sines and cosines right. Note the relative sign between  $E_{\perp}^I$  and  $E_{\perp}^R$ . You can check this sign because when  $\theta = \frac{\pi}{2}$   $E_I$  points up and  $E_R$  points down, i.e. they have opposite signs. There is a similar sign flip between  $B_I$  and  $B_R$ , as can be seen in the figure, since  $\vec{B}$  always lags behind  $\vec{E}$  by  $\frac{\pi}{2}$ .

Now we simply plug into our boundary conditions as solve. Eqs. (3) and (4) imply

$$\epsilon_1 E_I \sin \theta_1 - \epsilon_1 E_R \sin \theta_1 = \epsilon_2 E_T \sin \theta_2 \quad (25)$$

$$E_I \cos \theta_1 + E_R \cos \theta_1 = E_T \cos \theta_2 \quad (26)$$

Eq. (5) is trivially satisfied since  $B_{\perp} = 0$ . Finally, Eq. (6) implies

$$\frac{1}{\mu_1} (B_I - B_R)^{(2)} = \frac{1}{\mu_2} B_T \quad (27)$$

Using  $B = \frac{1}{v} E$  the last equation becomes

$$\frac{1}{v_1 \mu_1} (E_I - E_R) = \frac{1}{v_2 \mu_2} E_T \quad (28)$$

Dividing Eq. (25) by Eq. (28) we get

$$\epsilon_1 v_1 \mu_1 \sin \theta_1 = \epsilon_2 v_2 \mu_2 \sin \theta_2 \quad (29)$$

Using  $v = \frac{1}{\sqrt{\mu \epsilon}}$  and  $n = \frac{c}{v} = \sqrt{\frac{\mu \epsilon}{\mu_0 \epsilon_0}}$  we then get

$$n_1 \sin \theta_1 = n_2 \sin \theta_2 \quad (30)$$

which is Snell's law. It is reassuring that our derivation here reproduces Snell's law.

To simplify the solution it is helpful to define one parameter which depends on the angles but not on the materials

$$\alpha \equiv \frac{\cos \theta_2}{\cos \theta_1} \quad (31)$$

and another which depends only on the materials.

$$\beta = \frac{Z_1}{Z_2} = \frac{n_2 \mu_1}{n_1 \mu_2} = \frac{n_1 \epsilon_2}{n_2 \epsilon_1} = \frac{\epsilon_2 v_2}{\epsilon_1 v_1} = \sqrt{\frac{\epsilon_2 \mu_1}{\epsilon_1 \mu_2}} \quad (32)$$

Then the equations reduce to

$$E_I - E_R = \beta E_T \quad (33)$$

$$E_I + E_R = \alpha E_T \quad (34)$$

with solutions

$$\boxed{E_R^{\text{vert}} = \frac{\alpha - \beta}{\alpha + \beta} E_I, \quad E_T^{\text{vert}} = \frac{2}{\alpha + \beta} E_I, \quad n_1 \sin \theta_1 = n_2 \sin \theta_2}, \quad \text{vertical polarization} \quad (35)$$

The solution for horizontal polarization is similar,

$$\boxed{E_R^{\text{horiz}} = \frac{\alpha \beta - 1}{\alpha \beta + 1} E_I, \quad E_T^{\text{horiz}} = \frac{2}{\alpha \beta + 1} E_I, \quad n_1 \sin \theta_1 = n_2 \sin \theta_2}, \quad \text{horizontal polarization} \quad (36)$$

These are known as **Fresnel coefficients**. The Fresnel coefficients tell us how much of each polarization is reflected and transmitted. Since any polarization vector can be written as a linear combination of vertical and horizontal polarizations, we can use the Fresnel coefficients to understand how any polarizations reflect.

## 6 Transmitted power

To interpret the Fresnel coefficients, we would like to know not just the amplitude of the wave transmitted, but the intensity of the light that passes through, or equivalently the power. Recall that the power in a plane wave in the vacuum is  $P = c\epsilon_0 |\vec{E}|^2$ . For a plane wave in a medium,  $\epsilon$  changes and only the component of the velocity moving into the medium is relevant, so this becomes  $P = v \cos\theta \epsilon |\vec{E}|^2 = \cos\theta \sqrt{\frac{\epsilon}{\mu}} |\vec{E}|^2 = \frac{\cos\theta}{Z} |\vec{E}|^2$ . Thus the fraction of power reflected for vertical and horizontal polarizations is:

$$\frac{P_R^{\text{vert}}}{P_I^{\text{vert}}} = \left( \frac{\alpha - \beta}{\alpha + \beta} \right)^2, \quad \frac{P_R^{\text{horiz}}}{P_I^{\text{horiz}}} = \left( \frac{\alpha\beta - 1}{\alpha\beta + 1} \right)^2, \quad (37)$$

For the fraction of power transmitted,

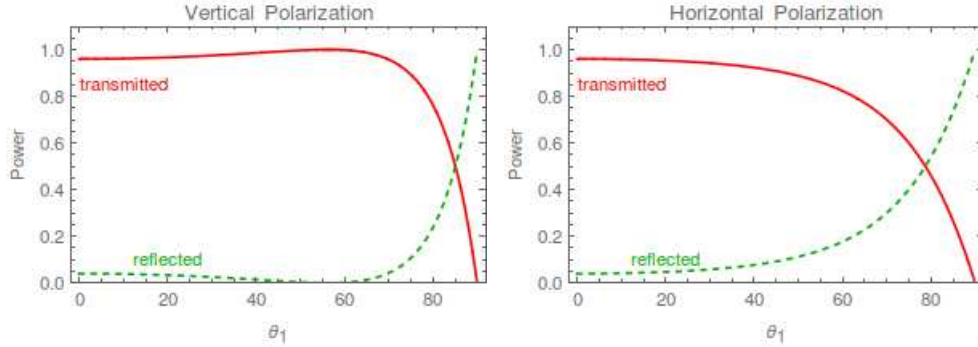
$$\frac{P_T^{\text{vert}}}{P_I^{\text{vert}}} = \frac{\cos\theta_2}{\cos\theta_1} \frac{Z_1}{Z_2} \left( \frac{2}{\alpha + \beta} \right)^2 = \alpha\beta \left( \frac{2}{\alpha + \beta} \right)^2 \quad (38)$$

and

$$\frac{P_T^{\text{horiz}}}{P_I^{\text{horiz}}} = \frac{\cos\theta_2}{\cos\theta_1} \frac{Z_1}{Z_2} \left( \frac{2}{\alpha\beta + 1} \right)^2 = \alpha\beta \left( \frac{2}{\alpha\beta + 1} \right)^2 \quad (39)$$

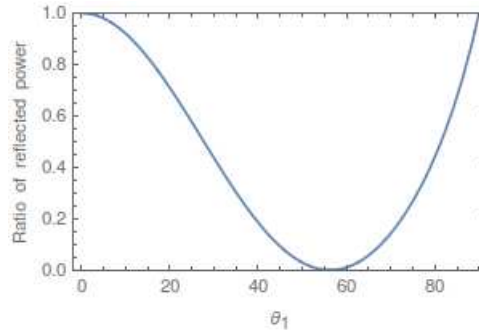
One can check that for either polarization,  $P_T + P_R = P_I$  and that these equations agree with Eqs. (17) and (18) for normal incidence ( $\alpha = 1, \beta = \frac{Z_1}{Z_2}$ ).

For example, the air-glass interface has  $\beta = 1.5$ . Then the transmitted and reflected power in vertical and horizontal polarizations are



**Figure 9.** Transmitted and reflected power as a function of incident angle for the two polarizations.

We can also plot the ratio of power reflected in vertical to horizontal polarizations:



**Figure 10.** Ratio  $P_R^{\text{vert}}/P_R^{\text{horiz}}$  of power reflected vertically polarized to horizontally polarized light, .

From these plots we see there is an angle where the vertical polarization exactly vanishes. This angle is called **Brewster's angle**,  $\theta_B$ . We can see from the plots that for the glass-air interface  $\theta_B \approx 56^\circ$ . At this angle, the reflected light is completely polarized.

What is the general formula for  $\theta_B$ ? From Eq. (37) we see that  $P_R^{\text{vert}} = 0$  when  $\alpha = \beta$ . That is

$$\frac{\cos\theta_2}{\cos\theta_B} = \sqrt{\frac{\epsilon_2\mu_1}{\epsilon_1\mu_2}} \quad (40)$$

For most materials  $\mu_1 \sim \mu_2 \approx \mu_0$ . Then we can use  $n = \sqrt{\mu\epsilon}$  to get  $\frac{\cos\theta_2}{\cos\theta_B} = \frac{n_2}{n_1}$ . Using also  $n_1\sin\theta_B = n_2\sin\theta_2$  we can solve for  $\theta_B$  giving

$$\boxed{\tan\theta_B = \frac{n_2}{n_1}} \quad (41)$$

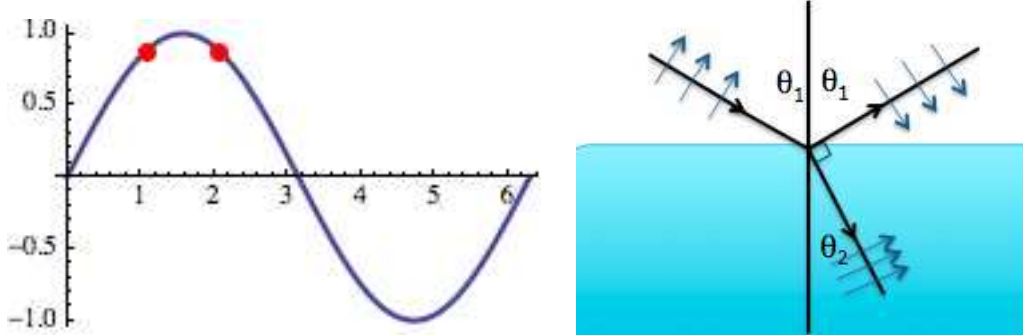
For the air-water interface,  $\theta_B = \tan^{-1}1.5 = 56.3^\circ$ . One can see from Figure 10 that one does not have to be exactly at this angle to have little reflected vertical polarization. Angles close to  $\theta_B$  work almost as well.

What is going on physically at Brewster's angle? We know that Snell's law  $n_1\sin\theta_1 = n_2\sin\theta_2$  is always satisfied. At Brewster's angle, when  $\theta_1 = \theta_B$ , we found also  $n_1\cos\theta_2 = n_2\cos\theta_1$ . Dividing these two equations gives  $\cos\theta_1\sin\theta_1 = \cos\theta_2\sin\theta_2$ , or equivalently

$$\sin(2\theta_1) = \sin(2\theta_2) \quad \text{when } \theta_1 = \theta_B \quad (42)$$

Now, by definition  $\theta_1$  and  $\theta_2$  are both between 0 and  $\frac{\pi}{2}$ , so  $2\theta_1$  and  $2\theta_2$  are between 0 and  $\pi$ . Thus,  $\sin(2\theta_1) = \sin(2\theta_2)$  has two solutions. Either  $\theta_1 = \theta_2$ , which corresponds to  $n_1 = n_2$  so the light just passes through, or  $\frac{\pi}{2} - 2\theta_1 = 2\theta_2 - \frac{\pi}{2}$  (see Fig 11, left) which simplifies to

$$\theta_1 + \theta_2 = \frac{\pi}{2} \quad (43)$$



**Figure 11.** The Brewster's angle happens when  $\sin(2\theta_1) = \sin(2\theta_2)$  which corresponds to  $\frac{\pi}{2} - 2\theta_1 = 2\theta_2 - \frac{\pi}{2}$  as can be seen in the left figure. This means that the transmitted and reflected waves are perpendicular. Thus the reflected vertically polarized light (shown on the right) cannot be produced by the motion of particles in the surface. Hence the reflected vertically polarized light vanishes at Brewster's angle.

Working out the geometry, as in Fig. 11, we see that the transmitted and reflected waves are perpendicular at Brewster's angle. Since the reflected wave has to be produced by the motion of particles in the surface we can understand why there is no reflection at Brewster's angle: particles moving in the surface can only produce light polarized transverse to their direction of motion. In Lecture 17, we will understand in more detail how accelerating charges produce electromagnetic fields and waves.

The famous Harvard dropout Edwin Land started a corporation called Polaroid that made its first fortune with polarizing sunglasses. These glasses remove glare from reflection by removing the horizontally polarized light. His second fortune was made with the Polaroid camera. In 1973 Land donated the money for the construction of the Science Center (which some people say looks like a camera...).

# Lecture 16:

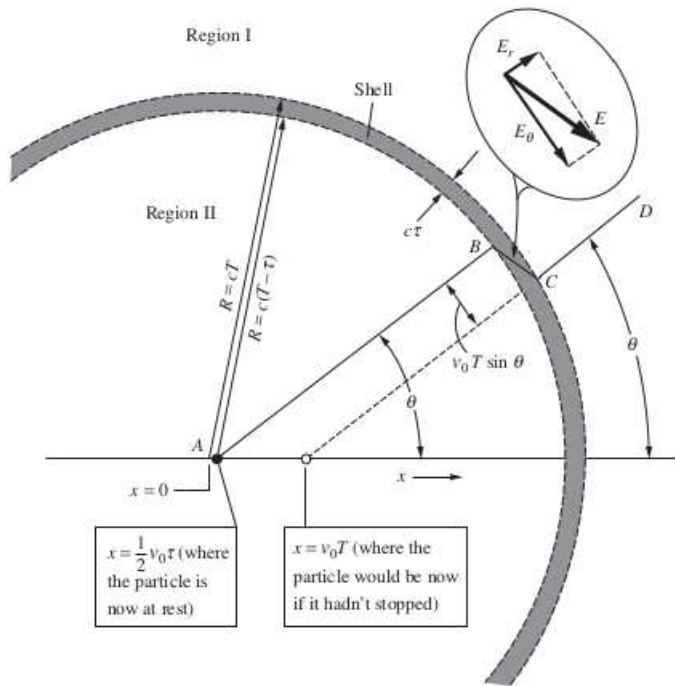
## Accelerating charges and prisms

### 1 Larmor formula

Now that we know how electromagnetic waves propagate, we need to know how they are produced. To produce oscillating electric fields, we need oscillating charges. In particular, we need charges which accelerate. It is easy to see why acceleration is necessary: a charge at rest produces only a static electric field. A charge moving at constant velocity, as in a current, produces a static magnetic field. To have the fields change with time, as in a electromagnetic wave, the charges must not be at rest or moving at constant velocity. That is, they must be accelerating.

A useful example of acceleration is a charge moving at velocity  $v$  that suddenly changes to velocity  $v_1$ . This example is useful because any acceleration can be built out of small little bits of acceleration: going from  $v$  to  $v_1$  to  $v_2$  etc.. Let's take  $v_1 = 0$  for simplicity and say that stopping occurs between times  $t = 0$  and  $t = \tau$ . Then the entire acceleration happens in this time interval. What is the field due to the accelerating charge? The relevant calculation is beautifully explained by Purcell. It's essentially just geometry I've put the relevant section, Appendix H of Morin/Purcell on openrev. Here I'll summarize the calculation.

The key insight is that after a time  $T$  one can use Gauss's law to determine the electric field for  $r < c(T - \tau)$  and  $r > cT$ . So the field due to the acceleration has to be confined to a shell of thickness  $\Delta r = c\tau$ . The situation looks like this



**Figure H.2.**

Space diagram for the instant  $t = T \gg \tau$ , a long time after the particle has stopped. For observers in region I, the field must be that of a charge located at the position  $x = v_0 T$ ; for observers in region II, it is that of a particle at rest close to the origin. The transition region is a shell of thickness  $c\tau$ .

**Figure 1.** Figure from Purcell/Morin Appendix H. A charge at  $x=0$  has acceleration  $a$  for a time  $\tau$ .



What Purcell shows is that the electric field line has to flow along  $ABCD$ . So from  $B$  to  $C$ , which is within the shell where the acceleration affects the field, it has a radial component  $E_r$  and a tangential component  $E_\theta$ . Looking at the geometry, it is not hard to see that

$$\frac{E_\theta}{E_r} = \frac{aR}{c^2} \sin\theta \quad (1)$$

where  $R$  is the distance to the shell. Now,  $E_r$  is determined by Gauss's law to be

$$E_r = \frac{q}{4\pi\epsilon_0 R^2} \quad (2)$$

where  $q$  is the charge of the thing moving (for an electron,  $\frac{q^2}{4\pi\epsilon_0} \approx \frac{1}{137}$ ). Note that  $E_r$  only depends on the net charge, not the acceleration.  $E_\theta$  is given by combining these two equations

$$E_\theta = E_r \frac{aR}{c^2} \sin\theta = \frac{qa}{4\pi\epsilon_0 c^2 R} \sin\theta \quad (3)$$

Note that  $E_\theta$  is proportional to the acceleration. Also note that at a fixed angle  $E_\theta$  dies with  $R$  only as  $\frac{1}{R}$  while  $E_r$  dies as  $\frac{1}{R^2}$ . When we have an AC current in an antenna, the net charge is 0, but the electrons are accelerating. Thus  $E_r = 0$  for antennas, but  $E_\theta$  is not.  $E_\theta$  has the information about the outgoing electromagnetic radiation

In a situation where  $E_r = 0$ , as in a current carrying wire or antenna, the energy density in the electric field is

$$\mathcal{E} = \frac{1}{2} \epsilon_0 \vec{E}^2 = \frac{q^2}{32\pi^2 \epsilon_0 c^4 R^2} \sin^2\theta \quad (4)$$

The total energy in the shell is twice this (because the magnetic field has the same energy density) integrated over the shell:

$$E_{\text{tot}} = 2 \int \mathcal{E} dV = 2 \int_0^\pi \sin\theta d\theta \int_0^{2\pi} d\phi \int_R^{R+\Delta r} r^2 dr \frac{q^2}{32\pi^2 \epsilon_0 c^4 R^2} \sin^2\theta \quad (5)$$

$$= \frac{q^2}{6\pi \epsilon_0 c^4} \Delta r \quad (6)$$

Now,  $\Delta r = c\tau$  as we observed above, where  $\tau$  is the time of the acceleration. So the power emitted, which is energy per time, is

$$P = \frac{q^2}{6\pi \epsilon_0 c^3} \frac{a^2}{c^3} \quad (7)$$

This is the **Larmor formula** for the power radiated by an accelerating charge.

The key results from this calculation are the boxed equations. Note that:

- An accelerating charge produces an electric field at a point  $\vec{R}$  in the direction perpendicular to the line between 0 and  $\vec{R}$  that scales as  $\frac{1}{R}$ .
- This field is also proportional to  $\sin\theta$ , where  $\theta$  is the angle between the line between the charge and  $\vec{R}$  and the direction of the acceleration  $\vec{a}$ . So it is maximal along the plane normal to the acceleration, and minimal in the direction of the acceleration.
- Power radiated is proportional to acceleration squared.

By Gauss's law, which follows from conservation of charge, one expects the electric field integrated over shell to be independent of the radius of the shell, so that the field dies as  $\frac{1}{R^2}$  at large  $R$ . But Gauss's law holds for static charges, not accelerating ones. The right conservation law in this case is conservation of energy – the energy in the field just moves outward. Indeed we see that the energy dies as  $\frac{1}{R^2}$  as in Eq. (4), which is what we expect for a conserved quantity. For energy to be conserved the field must scale like  $\frac{1}{R}$  which is only possible if the field is not uniform over the sphere. In particular, because of the  $\sin\theta$  factor, the  $\frac{1}{R}$  component of the field is zero at  $\theta = 0$  and  $\theta = \pi$ . Thus the radiation is essentially in a circle rather than a sphere.



To make this more concrete, think about an antenna. An antenna is just a set of charges moving back and forth in one direction. This antenna produces fields which decay only as  $\frac{1}{R}$  in the direction perpendicular to the antenna. In particular, antennas have directionality and can be more efficient at generating signals (or picking up signals), than a point charge would be. That is, a point charge has a field which dies like  $\frac{1}{R^2}$  so you have to be very close to sense its electric field. An antenna produces a field which dies like  $\frac{1}{R}$  so you can sense it from much farther away. We'll study antennas in detail in the next lecture.

## 2 Rayleigh and Mie scattering

When sunlight hits the sky it causes air molecules to vibrate. These vibrating molecules then radiate electromagnetic field down to us. Thus the light scatters off of the molecules.

There are two important limits. First, the wavelength of the radiation  $\lambda$  can be much larger than the size  $d$  of the molecule. For example, a water molecule has  $d \sim 1\text{nm}$  and visible light has  $400\text{ nm} < \lambda < 700\text{nm}$ . When light scatters off of small molecules in the atmosphere, such as  $H_2O$  or  $O_2$  or  $N_2$  then  $\lambda \gg d$ . In this limit light acts coherently and sets the molecule vibrating and emitting. This limit is known as **Rayleigh scattering**.

In Rayleigh scattering, the amplitude of the moving molecule becomes proportional to the electric field. For a plane wave, we then have the position of the molecule is

$$A(\vec{x}, t) = A_0 e^{i(\vec{k} \cdot \vec{x} - \omega t)} \quad (8)$$

The acceleration is

$$a = \frac{d^2 A}{dt^2} = -\omega^2 A(x, t) \quad (9)$$

By the Larmor formula, the power radiated is

$$P = \frac{q^2}{6\pi\epsilon_0} \frac{a^2}{c^3} \propto \omega^4 \propto \frac{1}{\lambda^4} \quad (10)$$

So the power radiated is **inversely proportional to the fourth power of the wavelength**. This is known as **Rayleigh's law**.

Now, red light has  $\lambda_{\text{red}} \sim 700\text{nm}$  and blue light has  $\lambda_{\text{blue}} \sim 400\text{nm}$ . Since blue light is shorter wavelength, more power is radiated in blue than in red. The ratio of power emitted is

$$\frac{P_{\text{blue}}}{P_{\text{red}}} = \frac{\lambda_{\text{red}}^4}{\lambda_{\text{blue}}^4} = 9.4 \quad (11)$$

So that's why the sky overhead is blue! When we talk about color, we'll compute what shade of blue it is.

When the sun is setting, we look at much shallower angles towards the sun. Then the blue scattered light goes off sideways and we see mostly what is left over, which is the red light. So that's why sunsets are red!

In the opposite limit the wavelength of light is much less than the size of the scatterer:  $\lambda \ll d$ . For example, dust particles have  $d \sim 1\mu\text{m}$  or larger. In this limit, light is more particle-like – it bounces off the particles like a mirror. This limit is called **Mie scattering**. For Mie scattering, all wavelengths of light are equally well reflected. For example, a cloud is made when lots of water molecules coalesce into droplets of sizes  $d \sim 1\text{mm}$  or larger (think about a raindrop). Then light just bounces off of them. That's why clouds are white! Similarly, fat globules in milk are  $10\mu\text{m}$  or so, much larger than the wavelengths of light. That's why milk is white, and why fatty milk is whiter than skim milk.

### 3 Transmission lines

Another consequence of the  $\omega^4$  frequency dependence of the power radiated is that an alternating current (AC) passing through a wire can radiate a lot. Power in the United States is AC at a frequency of 60 Hz (most of the rest of the world uses 50 Hz). Either way, this is a fairly low frequency compared to say the frequency of radio or television which is in the MHz ( $10^6$  Hz) to GHz ( $10^9$  Hz) range, or higher. An ordinary power cord looks like this

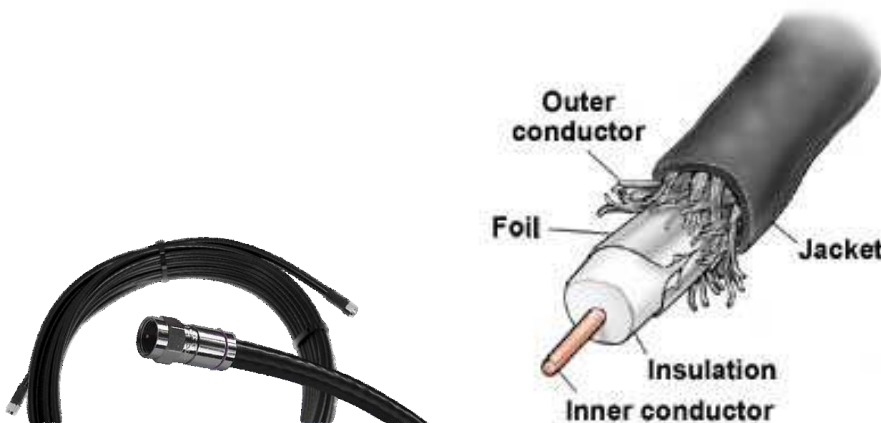


(12)

It has a pair of wires in parallel with current running through both to form a closed circuit. If the currents are exactly out of phase, as they would be if they go in opposite directions, then the field outside the wire largely cancels and not much power is radiated. Indeed, the wavelength associated with the 60 Hz oscillation is around 5000 kilometers, so the spacing  $d$  between the wires is much less than the wavelength  $d \ll \lambda$  and there is nearly complete destructive interference everywhere.

On the other hand, if you have an antenna receiving a radio signal at 1 GHz, with  $\lambda = 1\text{ m}$ . Then, there won't be perfect destructive interference everywhere. Moreover, since the power emitted goes like  $\omega^4$ , it can be enormous for 1 GHz frequencies. In fact, if you tried to connect your TV signal through a regular parallel electrical cord, it would probably catch fire or melt.

To transmit high frequency signals, we use **coaxial cables**.



**Figure 2.** Coaxial cables are used for high frequency transmissions

Coaxial cables have the current going one way in the middle and the return current going through an outer conductor which forms a cylinder around the inner conductor. In this way, the symmetry guarantees that the power will exactly cancel (to a very good approximation) even at very high frequency.

## 4 Microscopic origin of the index of refraction

Light moves at the speed of light  $c$ . So how can it move at a speed  $v = \frac{c}{n} < c$  in a medium with an index of refraction? What we will show is that an incoming plane wave excites charged particles in a material which then radiate. The interference between the original plane wave and the radiation from the accelerated charges conspire to make light propagate slower than  $c$ . In other words,  $v$  can be less than  $c$  in materials due to interference. I personally find this to be a very deep and satisfying result. Hopefully you will too.

When an electric field enters a medium like water with an index of refraction, it acts on the charged particles in the medium. The electric field pushes the charged particles up and down in the direction of the polarization of the field. A charged particle of mass  $m$  in a material would satisfy the wave equation in the absence of the external field

$$m \frac{d^2 x}{dt^2} + kx = 0 \quad (13)$$

where  $x(t)$  is the displacement of the particle from equilibrium and  $\omega_0 = \sqrt{\frac{k}{m}}$  is its characteristic oscillation frequency. Since the force from an electric field is  $\vec{F} = q\vec{E}$ , a plane wave with frequency  $\omega$  and amplitude  $E_0$  modifies this equation to

$$m \left[ \frac{d^2 x}{dt^2} + kx \right] = qE_0 e^{i\omega t} \quad (14)$$

This is are old friend the driven oscillator, whose solution is

$$x(t) = \frac{qE_0}{m(\omega_0^2 - \omega^2)} e^{i\omega t} \quad (15)$$

This solution is for a single charge. When a plane wave passes through a material, it acts on all the charges in an entire plane all at once. Each charged particle in the plane will be in phase (since the incoming wave is in phase) and will have displacement  $x(t)$ .

So now we have a plane of charges of thickness  $dz$  all moving coherently. Next, we need to work out the field produced by this plane. The “large”  $E_\theta$  component of this field from one charge given by Eq. (3):

$$E_\theta^{\text{one charge}}(t) = \frac{qa(t)}{4\pi\epsilon_0 c^2 R} \sin\theta \quad (16)$$

where  $R$  is the distance to the charge. How is the field of a plane of charges related to the field from a single charge? Consider what the field is at a distance  $z$  from the plane. This field gets a contribution from all the charges in the plane. It’s not a terribly easy calculation, since one must account for the different phases from points along the plane. The result is that the field is given by the *velocity* that the charges had at the time  $t_{\text{emit}} = t - \frac{z}{c}$  when the charges emitted the radiation:

$$E_{\text{plane}}(z, t) = -\frac{q\sigma}{2\epsilon_0 c} v\left(t - \frac{z}{c}\right) \quad (17)$$

where  $v(t) = \frac{dx}{dt}$  and  $\sigma$  is the number of particles per unit area. You can find a detailed derivation of this formula in Chapter 30-12 of the Feynman lectures (see Eq. 30.19). Plugging Eq. (15) into Eq. (17) we get

$$E_{\text{plane}}(z, t) = -i\omega \frac{\sigma q^2 E_0}{2\epsilon_0 c m(\omega_0^2 - \omega^2)} e^{i\omega(t - \frac{z}{c})} \quad (18)$$

This is the electric field produced by an infinitesimally thin plane of charges which have been accelerated due to an incoming plane wave.

Now, the total electric field at  $z$  is given by the sum of the incoming electric field and the field produced from the plane of charges

$$E_{\text{tot}}(z, t) = E_0 e^{i\omega(t - \frac{z}{c})} - i\omega \frac{\sigma}{2\epsilon_0 c} \frac{qE_0}{m(\omega_0^2 - \omega^2)} e^{i\omega(t - \frac{z}{c})} \quad (19)$$

$$= E_0 e^{i\omega(t - \frac{z}{c})} \left( 1 - i \frac{\omega}{c} \delta z \right) \quad (20)$$

where

$$\delta = \frac{1}{2\epsilon_0} \frac{q^2 \rho}{m(\omega_0^2 - \omega^2)} \quad (21)$$

and  $\rho = \frac{\sigma}{dz}$  is the density of charges per unit volume. Now,  $dz \ll 1$ , since the plane of charges is infinitesimally thin. So  $1 + i\delta dz = e^{i\delta dz}$  and we can therefore write in the limit  $dz \rightarrow 0$  that

$$E_{\text{tot}}(z, t) = E_0 e^{i\omega\left(t - \frac{z}{c} - \delta \frac{dz}{c}\right)} \quad (22)$$

So each little infinitesimally thin plane that the wave passes through forces the wave's phase to shift by  $\delta \frac{\omega}{c} dz$ .

Finally, once the wave has passed through a distance  $z$  of the material, this phase shift turns into  $\int_0^z dz \frac{\omega \delta}{c} = \frac{\omega}{c} z$  so that

$$E_{\text{tot}}(z, t) = E_0 e^{i\omega\left(t - z\frac{(1+\delta)}{c}\right)} \quad (23)$$

and we can identify the index of refraction as

$$n = 1 + \delta \quad (24)$$

The final result is the same as if light simply had the velocity  $v = \frac{c}{n} = \frac{c}{1+\delta}$  to begin with. Thus, the slowing down of light is just interference!!

## 5 Prisms

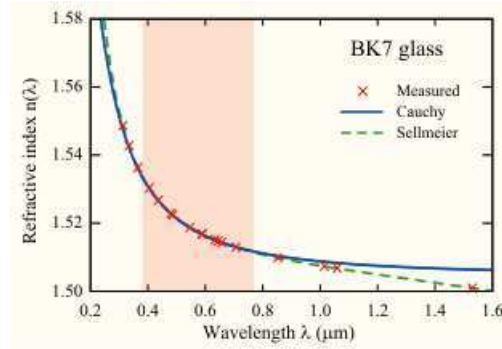
One can actually use this calculation for something. It not only tells us that  $n$  is related to interference, but also tells us that  $n$  depends on  $\omega$ . We found that

$$n = 1 + \frac{1}{2\epsilon_0} \frac{q^2 \rho}{m(\omega_0^2 - \omega^2)} = 1 + \frac{1}{8\pi^2 \epsilon_0} \frac{q^2 \rho \lambda^2 \lambda_0^2}{m(\lambda^2 - \lambda_0^2)} \quad (25)$$

Here,  $\omega_0$  is some characteristic wavelength of oscillation of the glass or whatever it is. Since light does not usually have enough energy to disrupt the glass, we expect it to be lower frequency:  $\omega \ll \omega_0$ . Expanding in this limit, or equivalently  $\lambda \gg \lambda_0$  gives

$$n = A + \frac{B}{\lambda^2} + \dots \quad (26)$$

with  $A = \left(1 + \frac{1}{2\epsilon_0} \frac{q^2 \rho}{m\omega_0^2}\right)$  and  $B = \frac{2\pi^2 q^2 \rho}{m\omega_0^4}$ . This dependence of the index of refraction on wavelength is known as **Cauchy's formula**.



**Figure 3.** Index of refraction of as a function of wavelength for a certain glass known as BK7 glass and a comparison to Cauchy's formula.

The fact that the index of refraction in glass depends on wavelength is the reason that prisms can spread the colors of the rainbow. Since the angle of refraction from air into glass is  $\sin\theta_1 = n(\lambda)\sin\theta_2$ , we see that incoming light at different angles refracts a different amount. The result looks like this



Figure 4. White light dispersed by a prism

## 6 Faraday cages

You may have noticed that the door of your microwave oven has grid lines on it. This grid is made of a conducting material and helps screen the microwaves from getting out. It is an example of a **Faraday cage**. But why don't the microwaves just pass through the grid?

A similar grid can be seen in many radio wave antennas, like this one in Arecibo, Puerto Rico:

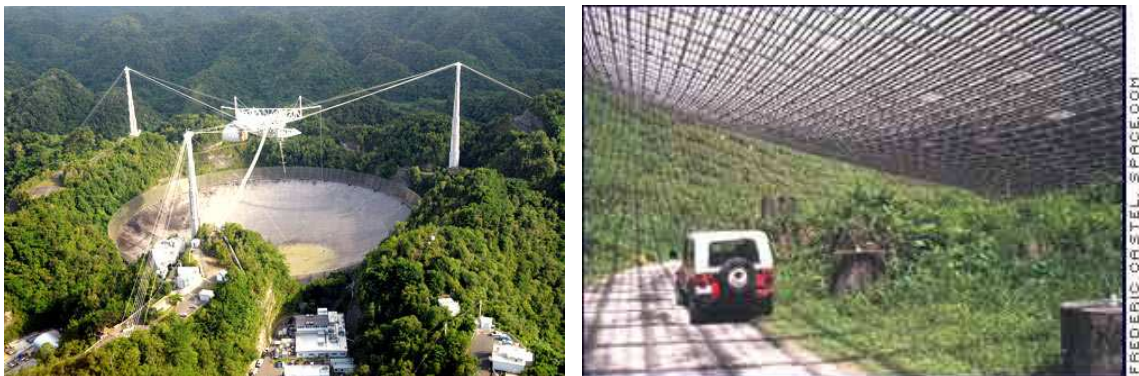


Figure 5. Arecibo telescope is a 1 km diameter dish antenna. The dish reflects radio waves to the receiver dangling above. The right shows the dish from close up. You can see it's just a grid of metal.

In this case, how does the grid reflect all the radio waves? Why don't most of them pass through the holes between the metal?

Your intuition is probably based on thinking of light like a stream of particles. If this were true, then most of the intensity would indeed go through the grid. Indeed in the limit that the grid spacing  $d$  is much bigger than the wavelength,  $d \gg \lambda$ , the grid does not block much light and the particle picture gives the right answer. On the other hand if  $d \ll \lambda$  or  $d \sim \lambda$  then we really need to think of light as waves. In this limit, the light comes in and excites electrons in the grid. These electrons then produce an electromagnetic wave which is exactly out of phase with the incoming wave. The two then destructively interfere on the opposite side of the grid. Thus the transmitted wave is zero and the wave is entirely reflected.

The key fact that lets this work is that a grid of conductors produces plane waves. Of course, they don't produce exactly plane waves. If the grid spacing is too large, each conductor will produce waves which go in circles. But if  $d \lesssim \lambda$  then the curvature is small (on the scale of the wavelength), and when the waves from all the conductors are summed coherently the net effect will be a plane wave which exactly cancels the incoming wave.

A typical microwave oven heats water using frequencies around 2.5 GHz (a wavelength of 12 cm). The spacing of the Faraday cage on the door is typically around 0.5 cm. So it is in the  $d \lesssim \lambda$  limit. For the Arecibo telescope, typical frequencies are 400 MHz, with wavelengths around 1 m. So as long as the spacing is a bit less than a meter (it looks like  $d \sim 10$  cm in the picture), all the radiation will be reflected.

By the way, the fact that the produced field is exactly out of phase with the incoming field in a conductor is a consequence of conductors accumulating charge on their surfaces. They can do this because the electrons in a conductor are only weakly bound and flow essentially freely. In particular, the model with spring constants as in Section 4 does not work for conductors; that derivation only works for systems where the electrons are bound near atoms and the effective spring constant picture can be applied. Such materials are insulators, not conductors.

# Lecture 17: Color

## 1 History of Color

You already know that wavelengths of light have different colors. For example, red light has  $\lambda \approx 650\text{nm}$ , blue light has  $\lambda \approx 480\text{nm}$  and purple light has  $\lambda \approx 420\text{nm}$ . But there is much more to color than pure monochromatic light. For example, why does mixing red paint and blue paint make purple paint? Surely monochromatic plane waves can't suddenly change wavelength, so what is going on when we mix colors? Also, there are colors (like cyan or brown) which do not appear in the rainbow. What wavelengths do they have? As we will see, the wavelength of monochromatic light is only the starting point for thinking about color. What we think of as a color depends on the way our brains and our visual system processes the light coming into our eyes.

The earliest studies of color were done by Newton. He understood that white light was a combination of lots of wavelengths. He performed some ingenious experiments to show this, described in his book *Optiks* (1704). For example, even though prisms turned white light into a rainbow, people thought that maybe the prism was producing the rainbow colors. Newton showed that this wasn't true. He could use lenses and mirrors to shine only red light into a prism, then only red light came out. He also separated white light into different colors, then used lenses and mirrors to put them into another prism which made white light again.

In his classic *Optiks*, Newton compared the colors going in a circle to the cycle of 5ths in music. His color wheel was

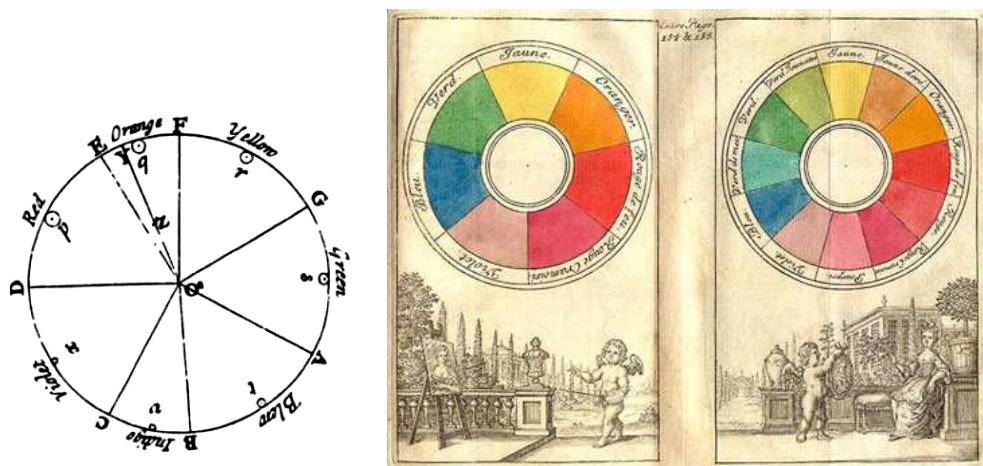


Figure 1. Newton's color circle 1704 and a contemporary rendition by Boutet (1708)

If the spectrum of visible light spans wavelengths from 350 nm to 750 nm, why should colors make a circle? Read on...

The next big advance came in 1853 by the mathematician Hermann Grassmann (Grassmann is well known in mathematical physics for his work on anti-commuting numbers which are now used to describe spinors like the electron in quantum mechanics). Grassmann was intrigued by the idea that two colors could mix to produce a different color, such as red+blue=purple. He showed that when you have an equation like this, you can add any color to either side, and it will still be a match, for example, red+blue+yellow=purple+yellow. This is known as Grassmann's law, which we will explore more quantitatively in the next section.

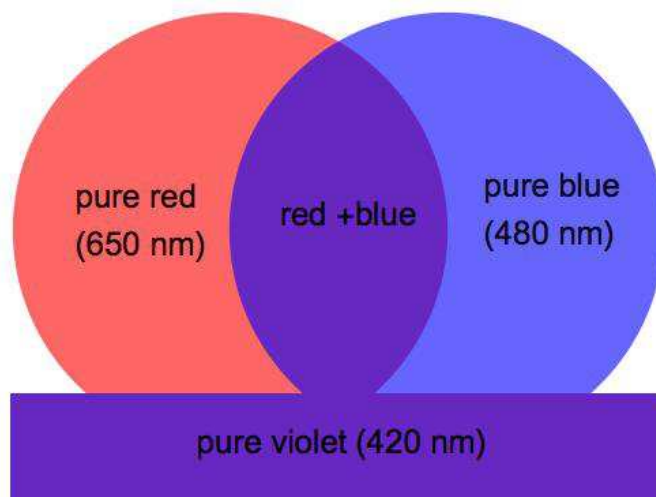


Another great contributor to our understanding of color was James Clerk Maxwell. He did a cool demonstration where he took black and white photos of the same colorful object using red, green and blue filters. Then when he projected those photos back again, through the same filters, the original multicolored object could be seen (this is still the way that color photography works, both film-based and digital). The amazing thing was that it wasn't just the red, blue and green colors that showed up, but the oranges and yellows and purples as well – all the colors. How is that possible? If the filters just let through pure red, blue, and green wavelengths, then none of the orange wavelengths would pass through. So how could they be reproduced? To understand this, we need to understand the perception of color.

## 2 Combining color

The modern theory of color was not laid down until 1931 by a set of classic experiments by W. D. Wright and John Guild in France, which built on the insights of many other people, including Newton, Grassmann and Maxwell. Their work led to the International Commission on Illumination (CIE) and their now standard color space.

What Wright and Guild did was to systematically quantify when two colors made up of different frequencies look the same. For example, if I show you some combination of red light ( $\lambda = 650\text{nm}$ ) and blue light ( $\lambda = 480\text{nm}$ ), it looks the same as pure violet light ( $\lambda = 420\text{nm}$ ):



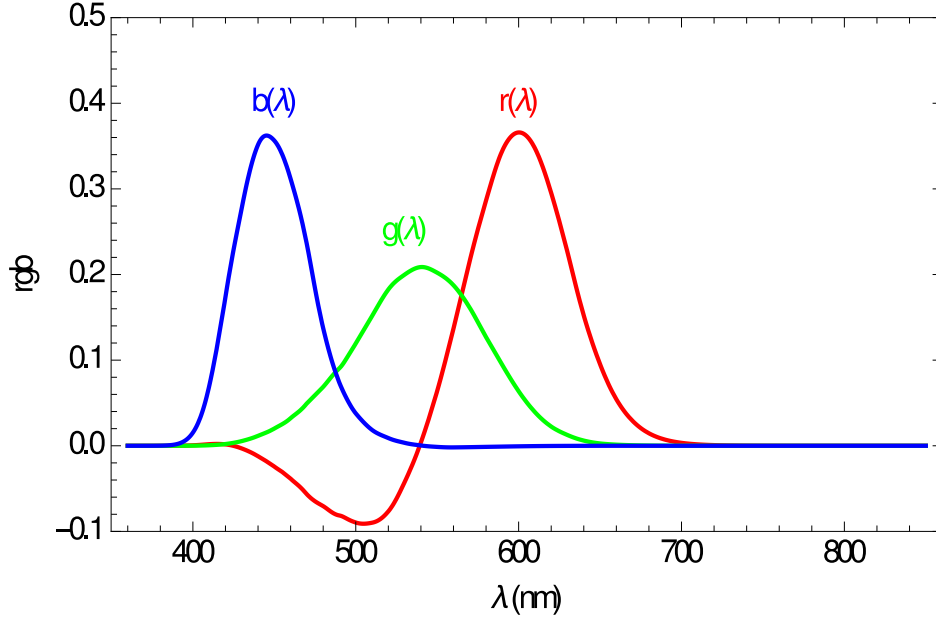
**Figure 2.** Mixed colors can look indistinguishable from a pure wavelength.

What Wright and Guild did was show people combinations of different wavelengths of light (the basis colors) and ask which pure wavelength of light it looked the same as. When two spectra of monochromatic light look the same, the colors are called **metamers**.

How many pure basis wavelengths do we need to combine to produce every color? Well, if we use just one color, like green, all we can do is vary the intensity – we always have green. So we need at least 2 colors. With two colors, the answer is not obvious. In principle, two colors could have worked. Two colors do work for colorblind people, and for most animals. But it's kind of obvious that 2 won't work for people, since by mixing red and blue we can only get shades of purple, never something like green. So Wright and Guild tried 3.

Wright and Guild showed people combinations of red (650 nm), blue (480nm) and green (530nm) and asked how much of the three colors looked the same pure light of wavelength  $\lambda$ . For example, they found that light of wavelength  $\lambda = 420\text{ nm}$ , which is a violet, can be matched with 0.97 parts blue, 0.20 red and 0.02 green. Their results are summarized in three functions  $r(\lambda)$ ,  $g(\lambda)$  and  $b(\lambda)$ , representing the intensities of red, green and blue light they needed to match light of a pure frequency  $\lambda$ . The functions look like this:





**Figure 3.** Original Wright and Guild RGB functions. These curves encode how we perceive color.

You should notice immediately something weird about these curves: the red one is less than zero in places. What does that mean? Well, in reality, Wright and Guild found that some colors could not be matched by mixing red, green and blue, for example, yellow with  $\lambda = 580$ . They got around this using Grassmann’s law:

- **Grassmann’s Law:** If two colors are indistinguishable (metamers), you can add any other color to them and they will still be metamers.

This essentially says that our color perception is linear. For example, Wright and Guild found that

$$I_{\text{violet}} \cong 0.97 I_{\text{blue}} + 0.20 I_{\text{red}} + 0.02 I_{\text{green}} \quad (1)$$

The “ $\cong$ ” sign in this equation means the colors on the two sides are metamers. Violet on the left hand side here means a monochromatic violet,  $\lambda = 420\text{nm}$ , while the right hand side is a combination of 3 pure wavelengths. Since metamerism is a linear equivalence, Grassmann’s law says that if we add more red, they will still look identical. Then

$$I_{\text{violet}} + 0.3 I_{\text{red}} \cong 0.97 I_{\text{blue}} + 0.50 I_{\text{red}} + 0.02 I_{\text{green}} \quad (2)$$

This is obvious when we write it as an equation like this, but it is very non-trivial.

One thing the linearity lets us do is take a color like yellow, which has  $\lambda = 580\text{nm}$ , and no matching to positive values of R, G, and B, and add to it something else, like violet to make the matching positive. For example, if

$$I_{\text{yellow}} + I_{\text{violet}} \cong 1.3 I_{\text{blue}} + 0.1 I_{\text{red}} + 0.3 I_{\text{green}} \quad (3)$$

we conclude that

$$I_{\text{yellow}} \cong 0.33 I_{\text{blue}} - 0.1 I_{\text{red}} + 0.28 I_{\text{green}} \quad (4)$$

This has a negative coefficient for red, just like in Wright and Guild’s curves.

To summarize, Grassmann, Wright and Guild’s observations were that our perception of light is a *linear system*. We can study perception with linear algebra. To repeat the “ $=$ ” sign in the above equation indicates metameric equivalence: “the combinations of pure wavelengths on the two sides are perceived identical by human beings”.

Negative intensities are fine mathematically, but what happens physically? Since we can see all the colors of the rainbow, are we somehow seeing negative intensities? Remarkably, the answer seems to be yes. Somehow we appreciate negative intensities as a result of higher-level cognitive processing. If this sounds nuts to you, try looking at

- Lilac Chaser: [http://www.michaelbach.de/ot/col\\_lilacChaser](http://www.michaelbach.de/ot/col_lilacChaser)

We will have more discussion of color *perception* in Sections 7 and 8 below.

### 3 CIE color space

Negative intensities are great mathematically, but not so practical for real world applications – you can't have an RGB monitor with negative intensities. But it would still be useful to have a representation of the total amount of color we can see, and to know which of those colors can be reproduced by a given RGB basis. A French committee called the International Commission on Illumination (CIE) got together in 1931 and came up with a very useful representation of the Wright and Guild data called the CIE color space.

The CIE commission realized that no three monochromatic colors could be mixed in positive amounts to match the full spectrum. However, they found that they could take a basis XYZ which are *sums* of monochromatic colors which would work. To make this precise, it is helpful first to write the Wright and Guild RGB values as a vector

$$\vec{R}(\lambda) = (r(\lambda), g(\lambda), b(\lambda)) \quad (5)$$

We can think of  $\vec{R}(\lambda)$  as a parametrized 3-dimensional curve. The new basis will be related to  $\vec{R}$  by a linear transformation

$$\vec{X}(\lambda) = M \cdot \vec{R}(\lambda) \quad (6)$$

with  $M$  a matrix. The transformation chosen by the CIE is

$$\vec{X}(\lambda) = \begin{pmatrix} X(\lambda) \\ Y(\lambda) \\ Z(\lambda) \end{pmatrix} = \frac{1}{0.17} \begin{pmatrix} 0.49 & 0.31 & 0.20 \\ 0.17 & 0.81 & 0.01 \\ 0.00 & 0.01 & 0.99 \end{pmatrix} \vec{R}(\lambda) \quad (7)$$

This particular transformation is more-or-less arbitrary. It is chosen so that  $X(\lambda)$ ,  $Y(\lambda)$  and  $Z(\lambda)$  are always positive. Lots of other choices would have worked just as well.

The functions  $X(\lambda)$ ,  $Y(\lambda)$  and  $Z(\lambda)$  are called the **tristimulus curves** or **CIE curves**. They look like this

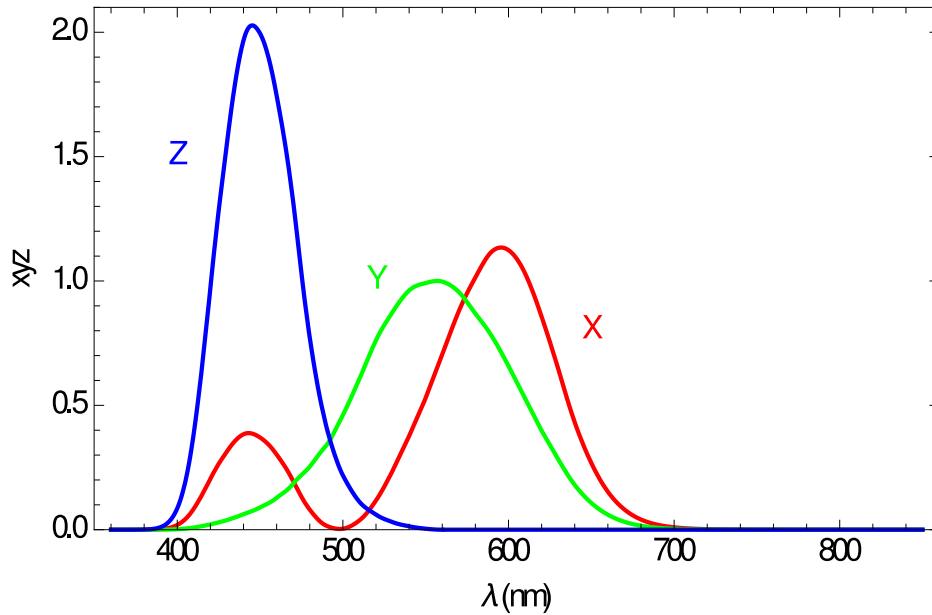


Figure 4. Standard tristimulus curves (CIE 1931).

As you can see from the figure any pure wavelength light will look identical to us as linear combination of  $X, Y, Z$  with positive coefficients. As we will see, tristimulus curves are essential to the design of film, printers, monitors, cameras, etc.

Note the small bump in the red tristimulus curve, and around 450 nm. This tells that violet looks like red + blue. It also explains the “line of purples”, which goes between red and purple and why Newton’s circle is a circle, as opposed to a line. Perhaps “explain” is too strong of a word. Of course, the CIE curves don’t actually explain why we see things that way – we could have just seen the rainbow, with red not looking anything like purple like a robot does. But somehow our mind wants to make a circle out of it. Maybe there’s an evolutionary explanation for this, I don’t know.

### 3.1 Normalizing the intensity

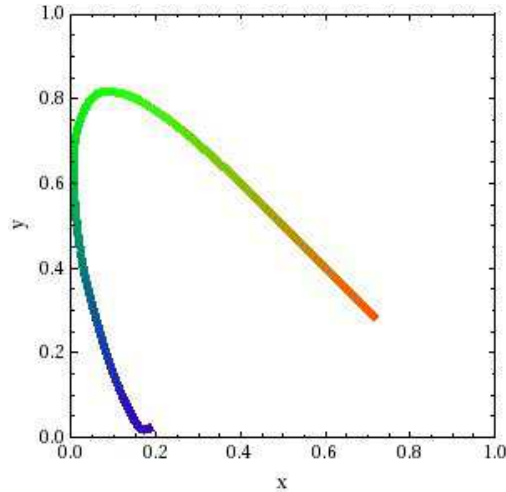
Because of the linearity, there is a well-defined notion of intensity of light: just sum the intensities in all the frequencies. In terms of the representation of a single frequency in terms of the CIE  $X, Y$  and  $Z$  values, this says that

$$I = \text{Intensity} = X + Y + Z \quad (8)$$

Then we only need two numbers to describe all the colors we can see. We call these two numbers lowercase  $x$  and  $y$ :

$$x = \frac{X}{X + Y + Z}, \quad y = \frac{Y}{X + Y + Z} \quad (9)$$

We can also define  $z = \frac{Z}{X + Y + Z}$ , but  $z = 1 - x - y$  so it is not independent. Each frequency of pure light gives a value of  $x(\lambda)$  and  $y(\lambda)$ . Plotting these curves gives a contour:



**Figure 5.** CIE curve:  $x$  and  $y$  values for **pure monochromatic light**, with  $z = 1 - x - y$ . That is, the colors of the rainbow fall along this curve.

Here’s where it gets interesting. Remember, the Wright and Guild experiments determined how much of their red, green and blue were needed to match a pure frequency  $\lambda$ . What happens when you mix just red and green? Or two other pure colors? Since the intensities are additive, we know that if we add two colors, the result would still match. For example, suppose we added some yellow light, with 0.3 intensity. Then we would find

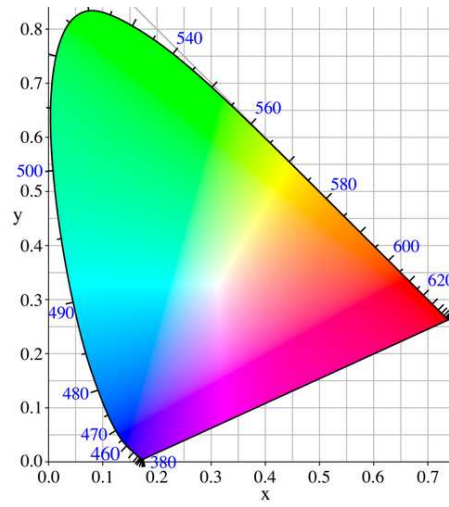
$$I_{\text{violet}} + 0.3I_{\text{yellow}} \cong 0.97I_{\text{blue}} + 0.20I_{\text{red}} + 0.02I_{\text{green}} + 0.3I_{\text{yellow}} \quad (10)$$

The two sides of this equation show that two weird linear combinations of different colors are the same, even though neither side is monochromatic. What is this color? It's not monochromatic, so it's not in the rainbow.

Whatever this color is, since we can perceive it, there will be values of  $x$ ,  $y$  and  $z$  which will match the perception. Which values? We can read them off of Fig. 5. Say yellow has a wavelength of 600nm and violet 420nm. Then the new  $x$ ,  $y$  values will be just the linear combination of the violet and yellow above

$$x_{\text{mixed}} \cong \alpha x_{\text{violet}} + (1 - \alpha)x_{\text{yellow}}, \quad y_{\text{mixed}} = \alpha y_{\text{violet}} + (1 - \alpha)y_{\text{yellow}} \quad (11)$$

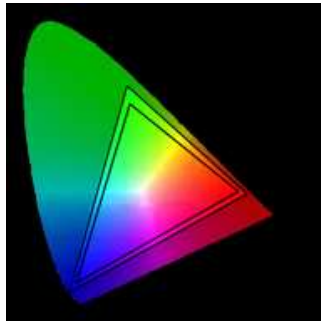
These are the equations for a line. Thus this color is somewhere on the line connecting yellow and violet on the curve in Fig. 5. Working out all the linear combinations, we get the full CIE color space, which is the gamut of colors we can perceive. It looks like this



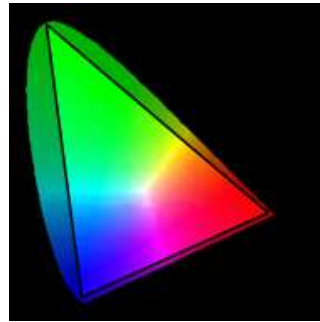
**Figure 6.** Approximate color space for the “CIE standard observer”. The wavelengths of monochromatic colors are marked on the boundary. The vertical line is a limitation of this RGB image.

For example, using the CIE functions, color in Eq. (10) corresponds to  $x = 0.46$  and  $y = 0.22$ , which is a pink.

A very important point here is that if you start with any three monochromatic colors, corresponding to three points on the boundary of this curve, the **gamut** of colors you will be able to produce with that basis will be just a triangle. For example, here are two triangles using the RGB color basis of TVs and computer monitors (left) and film (right):



**Figure 7.** CIE for computer monitors.



**Figure 8.** CIE triangle for film.

As we observed before, any monochromatic color can be made from any basis of 3 pure colors if negative intensities are allowed. If only positive intensities are allowed, then you are stuck in a triangle. What are the best colors we can choose? Those would be the colors that give the biggest triangle. This is pretty close to what they use for film. For computer monitors, the engineering requirements force them to use colors for which there are cheap LEDs, which gives a smaller triangle. In any case, we can see **the best basis will be red, green and blue**. A basis of red, yellow and blue, like Newton wanted, would get many fewer colors. In particular, RYB would have trouble making green.

Note that the “colors”  $X$ ,  $Y$ , and  $Z$  correspond to the points  $(x, y) = (1, 0)$ ,  $(x, y) = (0, 1)$  and  $(x, y) = (0, 0)$  respectively. Thus the triangle the CIE colors make is the bottom-left half of Fig. 6. This contains all the visible colors. It also contains lots of other linear combinations with negative intensities that don’t correspond to anything we can perceive at all. We can’t make a monitor with  $X$ ,  $Y$  and  $Z$  LEDs, since  $X$ ,  $Y$  and  $Z$  are not colors. The best we can do with actual colors is draw a triangle touching the color-space boundary. The point of inventing  $X$ ,  $Y$  and  $Z$  is it lets us embed all perceivable colors in a 2D triangle.

## 4 Using the CIE values

As an example application, we can compute what color the sky is. In Lecture 16, we showed that the intensity of scattered light depends on wavelength like  $I(\lambda) = c \frac{1}{\lambda^4}$  for some constant  $c$ . If you put all these scattered colors together with these relative intensities, what color do you get? By Grassmann’s linearity principle, the net color off such a signal would be equivalent to amounts of pure red, green and blue given by

$$r = \int r(\lambda) I(\lambda) d\lambda, \quad g = \int g(\lambda) I(\lambda) d\lambda, \quad b = \int b(\lambda) I(\lambda) d\lambda \quad (12)$$

Equivalently, we can use the standard  $x$ ,  $y$ ,  $z$  values to compute the integrals, then convert to RGB using the transfer matrix in Eq. (7) above. Performing the integrals, then normalizing the intensities so that  $r + g + b = 1$ , we find

$$(r, g, b) = (0.117, 0.295, 0.587) \quad (13)$$

This looks like, with various levels of saturation with white for cloud color (saturation is discussed in the next section):



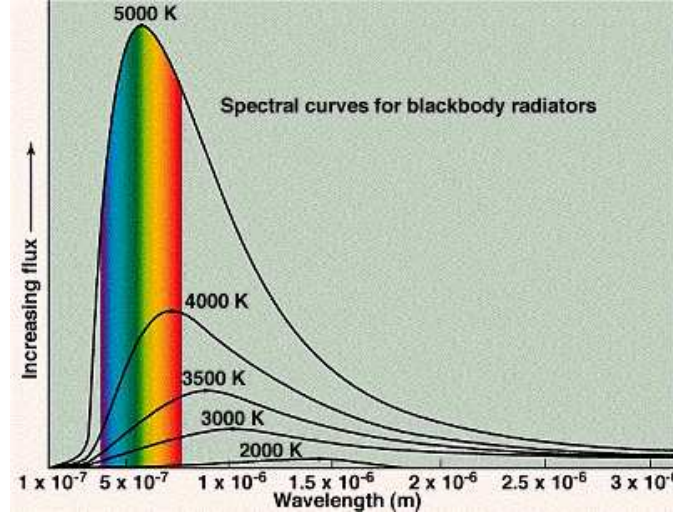
**Figure 9.** Sky color calculated using CIE observer functions

So, the simple answer to why sky is blue is that the formula  $I(\lambda) = \frac{1}{\lambda^4}$  is largest at smallest wavelengths, which is on the blue end. The more sophisticated answer is that we can figure out *what* shade of blue it is by integrating against the CIE curves. To be more careful, we should also include the fact that the incoming light is not completely flat in intensity because the sun has a spectrum.

To work out the color of the sun, we use that it is a blackbody at 5000 K. Blackbody radiation is thermal: it comes from highly excited molecules in a heat bath with their surroundings. This bath contains photons which are also thermal. The amount of each photon depends on the ratio of the photon energy  $h\nu$  to the thermal energy  $kT$ . A beautiful result from quantum statistical mechanics is that the intensity of radiation at a wavelength  $\lambda$  from an object at temperature  $T$  is

$$I(\lambda) = \frac{2hc}{\lambda^5} \left( e^{\frac{hc}{\lambda kT}} - 1 \right)^{-1} \quad (14)$$

where  $h$  is Planck's constant,  $k$  is Boltzmann's constant and  $c$  is the speed of light. These blackbody curves look like



**Figure 10.** Blackbody radiation curves for different temperature. The sun is around 5000 K, incandescent light bulbs are around 3000K.

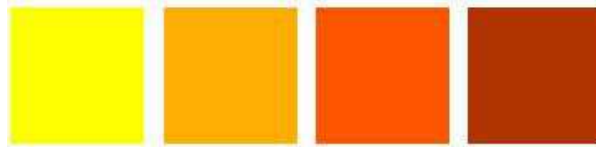
Note that the sun does not have much intensity in violet, which partially explains why the sky is blue and not violet (the main reason, however, is that the  $\lambda^{-4}$  Rayleigh spectrum pushes the peak intensity away from violet).

A useful feature of the blackbody spectrum is that the wavelength where the spectrum is maximal scales inversely with the temperature, as given by **Wein's displacement law**:

$$\lambda_{\max} = \frac{2.898 \text{ mm K}}{T} \quad (15)$$

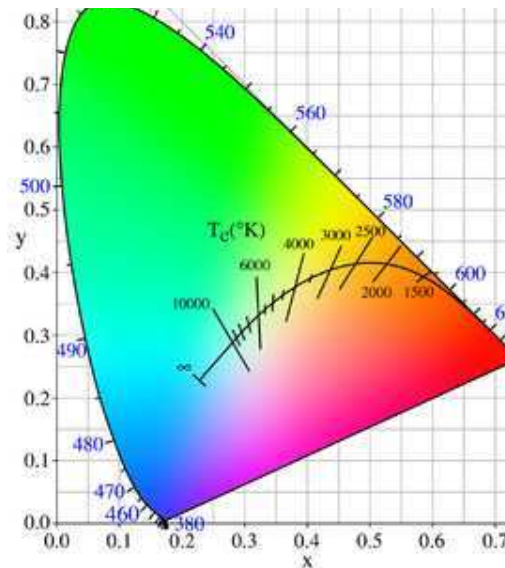
Plugging in room temperature  $T = 300 \text{ K}$  gives  $\lambda_{\max} = 9.6 \mu\text{m}$ , which is in the infrared. That's why people show up in infrared cameras. Plugging in  $T = 5000 \text{ K}$  gives  $\lambda_{\max} = 501.4 \text{ nm}$  which is green.

To find the color of the sun, we integrate the Planck blackbody distributions against the CIE color functions, as in Eq. (12). The result is basically white. The reason the sun looks yellow is because blue light is scattered away from the sun by the sky. In fact, we can figure out the colors of the sun and sunset by subtracting the sky spectrum from the sun spectrum. Doing this, gives:



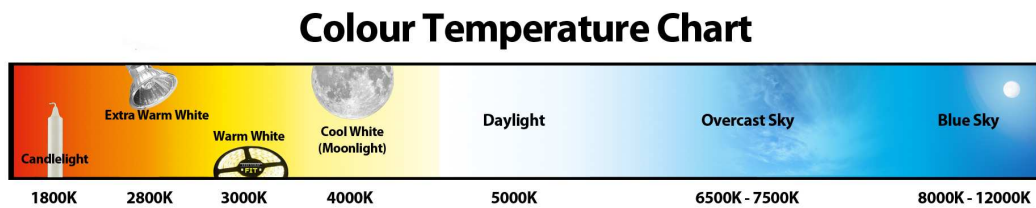
**Figure 11.** Sunset colors – sun color with scattering from the sky subtracted out.

More generally, each temperature translates to a particular color when the blackbody spectrum is integrated against the CIE curves. So as we vary the temperature, we trace out a curve of colors, or equivalently a path in the CIE color space. This curve has a name: the **Planckian locus**. It looks like this



**Figure 12.** Color of blackbody radiation: The line is the **Planckian Locus**.

If you look at the color along this locus, you get something like



**Figure 13.** Colors along the Planckian locus and associated temperatures. The sun at 5700K is white.

This says that the hotter a color the more blue it is, the cooler, the more yellow or red, and in the middle it's white. Red stars, like red giants such as Betelgeuse (Orion's shoulder) are cooler, Yellow stars, like our sun are next, then white stars, such as white dwarfs, then blue stars. Sirius, the dog star, is the brightest star in the sky and is whitish-blue in color.

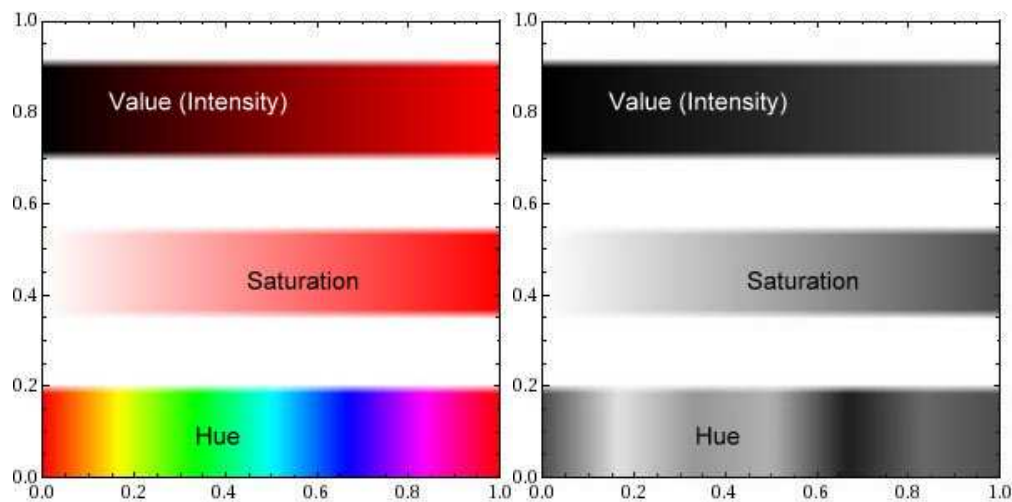
## 5 RGB and HSV

If all the colors of the rainbow are on the boundary of the CIE shape, what are all the colors in the middle? Our gamut includes a lot more than monochromatic colors.

We saw that we have a 3-dimensional basis of color perception – every color we can see can be represented by  $x$ ,  $y$ ,  $z$  values, representing (roughly) the intensities of red, green and blue light. One linear combination gives something equivalent to a pure monochromatic rainbow color (called **hue**). This corresponds to moving around the boundary CIE curve. One linear combination is the intensity (called **value**). This is normalized to 1 in the CIE color space. What is the remaining degree of freedom? We call it **saturation**. It is the amount of white a color has. It corresponds to moving inwards from the CIE boundary towards the white spot in the center.

Here are the three degrees of freedom in this basis

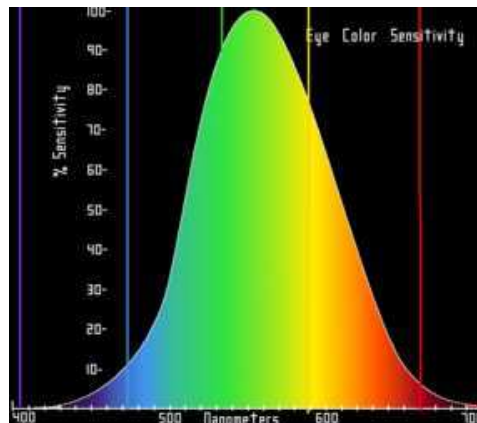




**Figure 14.** The HSV coordinates – using the Hue[h,s,v] command in mathematica.

One advantage of HSV is it lets us easily understand what colors are not in the rainbow. They are all the colors that you can get from a rainbow color (monochromatic, fixed hue), by varying the saturation (moving away from the boundary). If you vary the intensity, you just get a different value, which is normalized away in the CIE color space. Note from black-and-white part of Fig.14 that saturation, not value, characterizes how black-or-white a color is. That is, black-or-whiteness is different from brightness.

Another point you might notice is that the yellow hue seems brighter than the blue hue, despite the fact that they have the same intensity. This is because our eyes are more sensitive to yellow and green light than to blue or red light. Our overall color sensitivity curve looks like this:



**Figure 15.** Color sensitivity of human eyes.



## 6 Additive versus Subtractive colors

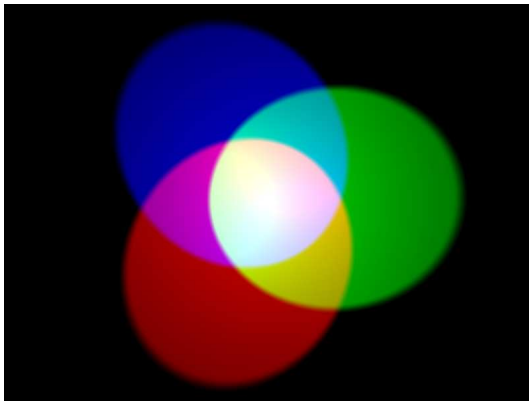
The color mixing we have been talking about so far has been mixing of light. This is called **additive** mixing. For example,

$$\text{blue} + \text{red} = \text{purple} \quad (16)$$

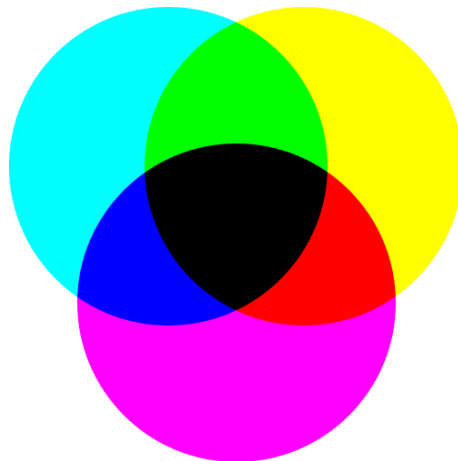
This is in contrast to what would happen if we mixed paint: blue paint + red paint = black. This is because blue paint *reflects* blue. It absorbs red and green. red paint absorbs blue and green. Put them together and you absorb everything, giving black. This is called **subtractive mixing**. Here's an easy way to remember the difference

- Additive mixing makes things whiter.
- Subtractive mixing makes things black.

For example, here are some standard color mixing diagrams



**Figure 16.** Additive color mixing



**Figure 17.** Subtractive color mixing

How can we figure out the subtractive color mixing rules using the CIE functions? Subtractive color mixing is multiplicative, since pigments *remove* a fraction of the incoming light. When you mix two colors, you simply multiply their intensities. For example, consider magenta and cyan. Their RGB values are

$$\text{magenta} = (1, 0, 1), \quad \text{cyan} = (0, 1, 1) \quad (17)$$

If we mix them additively we get

$$\text{magenta} + \text{cyan} = (1, 1, 2) \sim \text{lavender} \quad (18)$$

If we mix them subtractively, we need to write

$$\text{magenta} \times \text{cyan} = (1, 0, 1) \times (0, 1, 1) = (0, 0, 1) = \text{blue} \quad (19)$$

We can see the relationships between the additive and subtractive systems in the color wheels above.

The **complimentary color**  $\bar{C}$  of a color  $C$  is one which you must add to get white:  $C + \bar{C} = \text{white}$ . Colors on opposite sides of white (or black) are complimentary. Note that the subtractive primaries are white minus the additive primaries *e.g.*  $\text{cyan}(0,1,1) = \text{white}(1,1,1) - \text{red}(1,0,0)$ , which is why they are called subtractive. When you combine many colors additively, they give white and when you combine them subtractively they give black. Try also this useful additive and subtractive mixing applet

- Mixing applet [http://www.michaelbach.de/ot/col\\_mix/index.html](http://www.michaelbach.de/ot/col_mix/index.html).

## 6.1 CMYK basis – printing

Why are subtractive colors important? Printing, for one thing. If you mix red, green and blue you get white light. Printers work by printing the same thing in three colors. But if you put down red blue and green ink, you will just get black, not white. This is a little subtle. You might think that if you painted a bunch of thin green lines, red lines, and blue lines next to each other, all the colors would get reflected, making the combination white. However, only 1/3 of the light is reflected, so it is a very dark white. Also known as black.

The translation from RGB colors, which are good on your screen, to colors which print well is not so easy. In fact, it is very complicated. First of all, what is the best basis for ink colors we can choose? For light colors, we saw it was the basis with the biggest gamut triangle, which was some kind of red, green and blue. For printing, we need the complimentary CIE space.

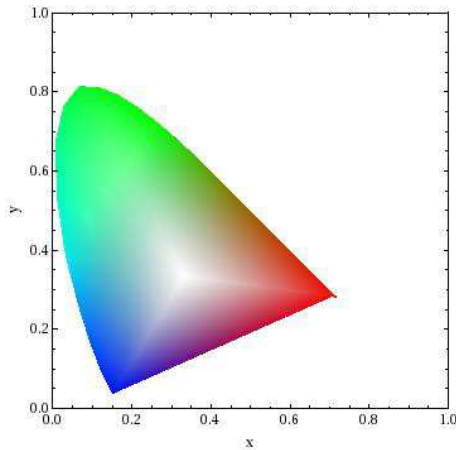


Figure 18. Additive color space

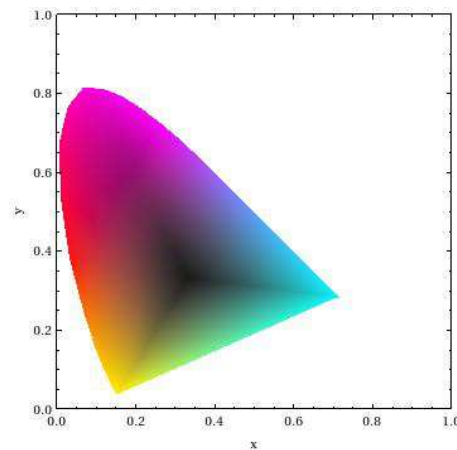


Figure 19. Subtractive color space

From this we can see that cyan, yellow and magenta make good subtractive colors. This is why printers use these three in their ink. The CMYK stands for these three colors and the K

stands for “Key”, which means the key color black. If the 3 dimensional basis is enough to produce all the colors, why also use black? Because black ink is cheaper than color ink!

This is only part of the story. To get various saturations of color, printers have to mix in white. This is done with **half-toning**, a more sophisticated version of the Ben-Day dot technique which inspired pop artist Roy Lichtenstein. In half-toning, colors are mixed, with some black, and various angles so that at large distances they look like intermediate colors

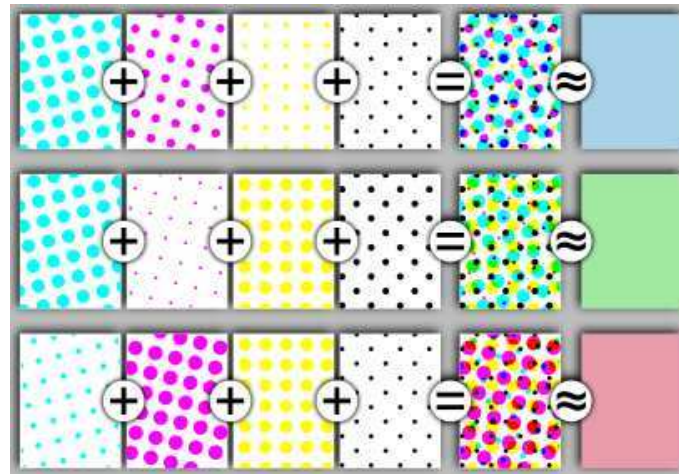


Figure 20. Half-toning

## 7 The Eye

Next, we’ll talk about perception. Color cannot be understood without understanding more of how we humans process it with our cognitive system. To begin, let’s look the way the input is processed – through our eye.

The eye looks like this

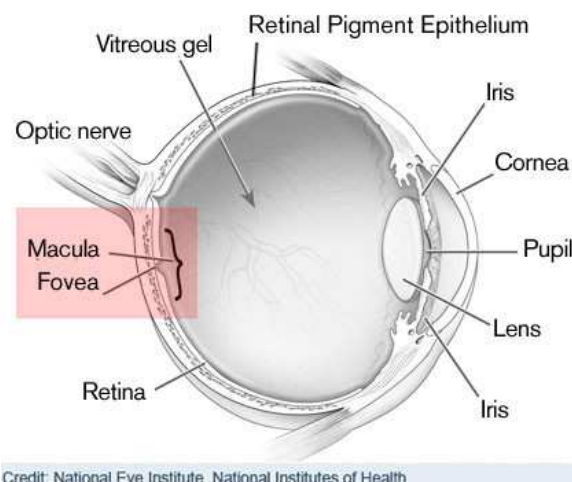


Figure 21. The Eye. The boxed area contains most of the cones, which are how we sense color.

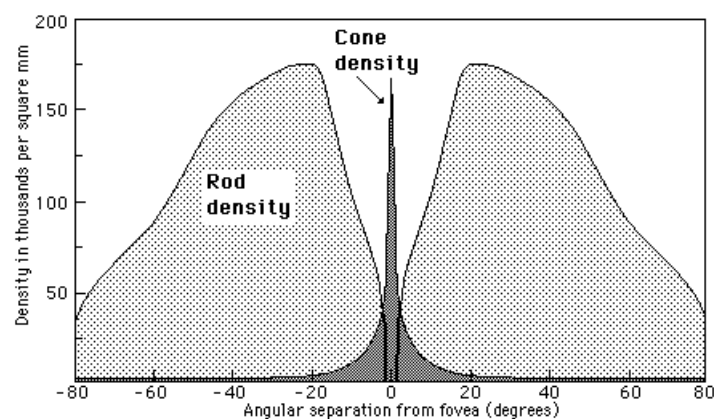
There are many parts of the eye which affect how we see:

- The pupil: hole where the light comes in.
- The iris: opens and closes around the pupil to let more or less light in.
- The retina: surface of back of the eye. Has the rods and cones
- The lens: focuses the light on the the retina.
- The cornea: Also helps focus light, onto the lens.
- Vitreous humor: fluid in the eye. Index of refraction  $\sim 1.333$
- The macula: part of retina with the best resolution.
- The fovea: Small part of the macula. Consists only of cones
- The optic nerve: pipes information out of the eye. Leads to *blind spot*.

There's a lot of physics in all these parts. For now, we're going to concentrate on the parts relating to color vision. The important facts about the eye relevant to our discussion are mostly about rods and cones.

## 7.1 Rods and cones

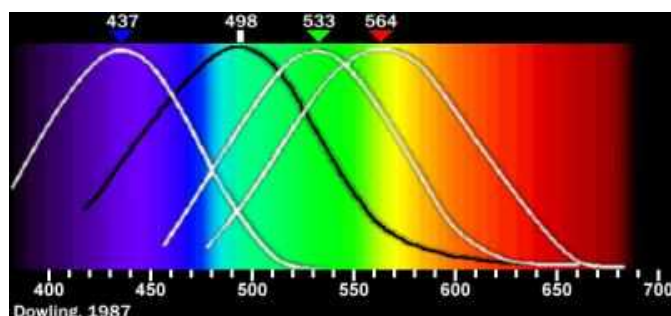
Rods and cones are the photoreceptors in the eye. There are many more rods ( $\sim 120$  million) than cones ( $\sim 6$  million) in the eye. The rods are scattered all about, except in the fovea, which is the central part of the macula, which is the central part of the retina. The fovea is all cones. There are a few cones outside of the fovea, but none towards the edges of the retina. The angular distribution looks like this:



**Figure 22.** Rod and cone density in the eye

There are 3 kinds of **cones**, called short (blue) medium (green) and long (red), which correspond to different sensitivities in wavelength. The reason we have a 3-dimensional space of color perception (we are *trichromats*) is because we have three cones. The measured sensitivities look

like this



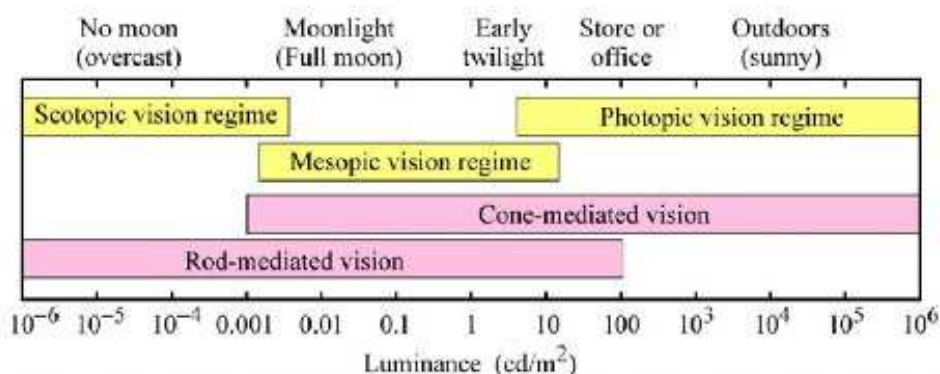
**Figure 23.** Sensitivity of cones (white curves) and rods (black curve)

Note that the “red” cone, is closer to having it’s maximal sensitivity in yellow than in red. You might have expected these to look like the  $X, Y, Z$  CIE curves, Fig. (4), but they don’t – there is no bump in the red curve at small wavelength. These are basically just the power spectra for forced simple linear systems: resonance with damping, so you wouldn’t expect a second bump. So where does the second bump in the red CIE curve come from? It must be a result of our higher level cognitive processing, which is still not well understood. Since the cones are in the central part of the eye, we have the best color perception in the center of our field of view.

Most animals, such as dogs, only have 2 cones and are *dichromats*. It’s not that dogs are color blind – they can distinguish all the colors of the rainbow, all the hues. They just can’t measure hue, saturation and value, only hue and value. So many colors that we can distinguish, dogs can’t.

Now lets talk about **rods**. There are 20 times as many rods as cones in the eye. Rods have greater sensitivity to low light conditions and for seeing motion. Since the rods are mostly scattered around the periphery of the retina, we have better motion perception at the edge of our field of view. This makes sense evolutionarily – we want to see predators coming in from the side. When we are looking right at something, we want to know more details, like its color.

Since there are so many more rods than cones, they have different absolute light sensitivities. Overall, our eyes are sensitive to roughly 12 orders of magnitude in intensity, with the rods dominating at low intensity and the cones at high intensity.



**Figure 24.** Sensitivity of our eyes. Cd is a candela, which is the SI unit of luminance. Scotopic, mesopic and photopic vision means vision under low-light, medium-light and well-lit conditions.

Since the rods are also in different places in the eye than cones, and since there are so many more of them, our rods do not just act like another cone with blueish sensitivity. Some input from rods is used in the interpretation of color, but we do not have a 4 dimensional space of color perception, only 3 dimensional. Basically, rods are used to interpret intensity, not color.

Rods and cones are types of neurons. The way neurons work is building up potential gradients and then discharging electrical or chemical signals. They produce an all-or-none response. When you combine lots of neurons you can get fairly smooth response curves, but the basic point is that there is some threshold sensitivity for them to work. When there is a lot of light coming in to your eye, you want these thresholds to be low, so that you can tell intensity differences between say  $1000 \frac{\text{cd}}{\text{m}^2}$  and  $10,000 \frac{\text{cd}}{\text{m}^2}$  (cd is a **candela**. 1 cd is about the amount of light coming off a candle). At low light, you want them to go between  $10^{-4} \frac{\text{cd}}{\text{m}^2}$  and  $10^{-5} \frac{\text{cd}}{\text{m}^2}$ . They can do this, but it takes a little while to recalibrate. It takes longer to calibrate at lower light, simply because there's less light coming in to do the calibration with. Rods typically take 7-10 minutes to calibrate at low light. This is why it takes around 10 minutes to really see stars at night. It takes much less time to calibrate to bright light, ~30 seconds.

Let's go back to the rod color sensitivity curve, Fig.(23). Notice that **rods do not see red**. This is why darkrooms are red. It's not that film isn't sensitive to red, it's all about the adjustment time of your eyes. You need to have the lights all the way off when you develop the film. So if your eyes have adapted to the dark, if you put red light on, it will not spoil your rods' adaptation – the rods will still be on low sensitivity settings – low ISO. Your cones will see the red light, and they can adapt quickly. If you shined bright lights, your rods would have to readapt and it would take 10 minutes. This is also why instrument panels on boats are often illuminated with red lights, so the captains can see at night.

Another consequence of the cones not seeing red is that in dim light, red looks black. This is the origin of the Purkinje effect. Purkinje was a 19th century biologist, who observed that the red geraniums in his garden looked striking during the day, but in the early morning, before it was light, they were much darker even their leaves. The contrast is something like this:



**Figure 25. Purkinje shift.** Red flower during daylight (left) and at dusk (right). The right image is simulated, not a photo.

If we just took a photo of the flower at night, it would look like a dimmer version of the flower, not like the black thing on the right. That's because cameras have color sensors, which are like cones, not rods. The Purkinje effect is due to rods taking over at night. That's part of what makes color photography so tricky! More about this below.

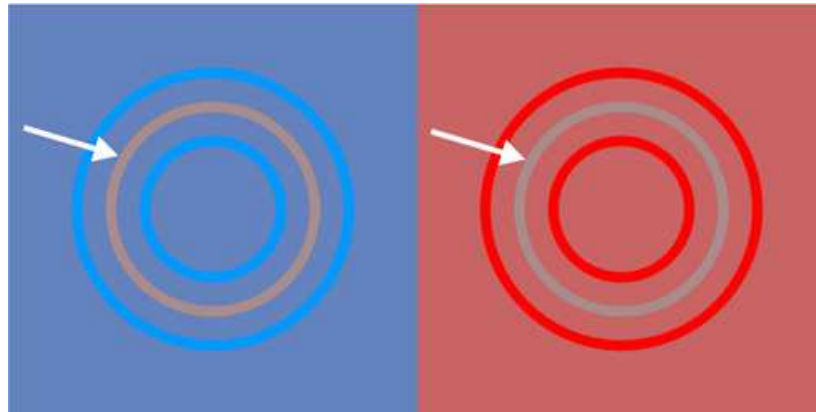
## 8 Vision and human perception

Now that we know a little about the eye, we can talk about processing color. A surprising amount of color is context based.



## 8.1 Color adaptation

Look at this figure:

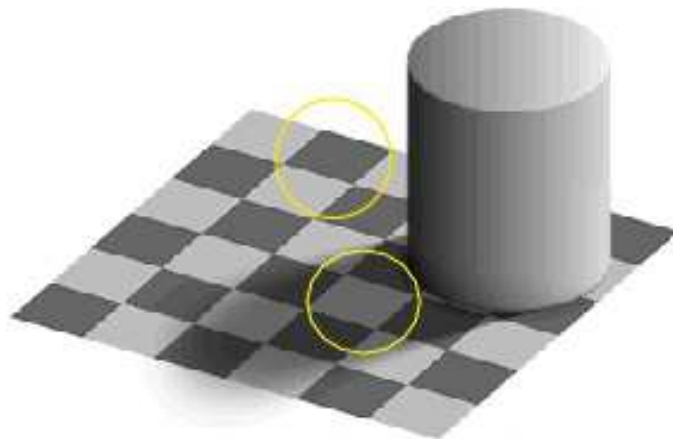


**Figure 26.** The two rings indicated by the arrow are the same color

The rings pointed at by the arrow look to be different colors, right? In fact, they are the same color. The issue here is that our interpretation of a color is determined in part by its surroundings. Our cones are being saturated by the backgrounds, so that the foregrounds look like they are different hues. It's also possible for our eyes to make color out of pure black and white:

- Color from nothing: [http://www.michaelbach.de/ot/col\\_benham](http://www.michaelbach.de/ot/col_benham)

The sensitivity to context can be seen also for black and white images. Perhaps the best illustration of this is the Adelson's Checkerboard:



**Figure 27.** Do the two circled squares look the same? In fact they are exactly the same shade of gray!

Here's another example. This dress looks like a different color, depending on what's in the background:



**Figure 28.** What color is the dress? It depends on what's around it.

## 8.2 Blind Spot

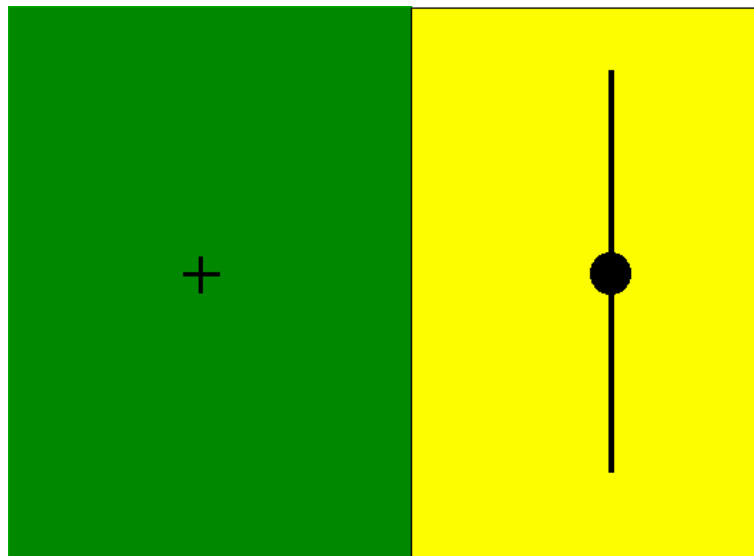
Your mind can do all kinds of other things to you as well. The signals of what you are seeing have to exit the eye somehow. They do this through the optic nerve. Unfortunately, the region where the optic nerve exits your eye can't have rods or cones. So it is a blind spot. It's pretty amazing what our minds do to pretend this blind spot doesn't exist.

First, let's prove there is a blind spot. Close your left eye and stare at the cross with your right eye. Move your eye closer to the page. At some point you will see the dot disappear. This is because the dot is right in line with the blind spot



**Figure 29.** Blind spot

That wasn't the interesting thing (although it is interesting). Now try it again with this figure.

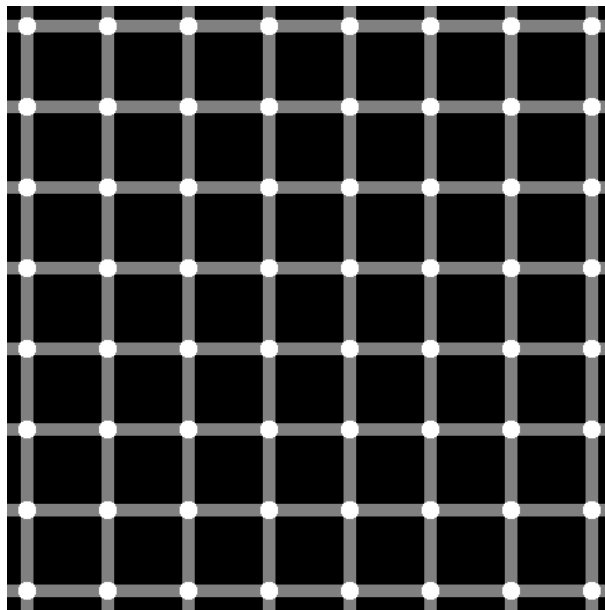


**Figure 30.** Blind spot #2

You should have seen a number of things: the dot went away. But it was replaced by a line! And a yellow image! Why a line? You think that's the most natural thing to be there.

Our mind tries to fill in information from context in other situations too, not just over our blind spot. For example, look at this image





**Figure 31.** Count the black dots!

There are no black dots, but your mind thinks they should be there, so it puts them in. Perception is complicated!

## 9 Photography and imaging (optional, but interesting)

We have seen that there is a mathematical way to describe color. All the colors we can see can be written as a linear combination of intensities of a standard set of  $X$ ,  $Y$  and  $Z$  “colors”. However, the color we see depends a lot on how we process that signal. This processing involves converting from a linear to a logarithmic scale, by adapting our rods and cones to the dark, and taking clues from context. This makes it very difficult to present color on text or on a computer that look like what we see in real life.

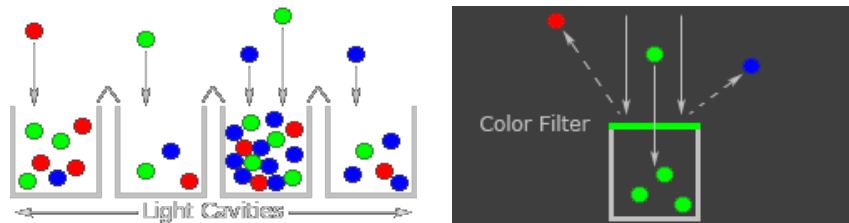
First of all, since any 3 RGB values can produce only a subset of all the colors we can see, it’s impossible to get a true color image of anything. That’s why photographs of rainbows never look like the real thing – they can’t.

RGB values are designed to simulate the color we can see. This does great in bright light, where cones dominate, but not for dim light, where the rods are critical. Recall the Purkinje effect, Figure (25). If you took a photograph of a flower in dim light, it would look like a dimmer version of the flower, not what we see. Of course, if you look at the photo in dim light, it would start to look like the dim flower, but who wants to look at photos in the dark? Instead, we look in bright light, and see the reds perfectly well, even though they are at low intensity. So how can we get the photo to look like what we see, where the reds appear black? That is one of the main challenges in color imaging: from photography to TV to movies to computers. The main point: if you want to make the image simulate our perception, you have to do more than simply reproduce the RGB values.

### 9.1 Digital Cameras

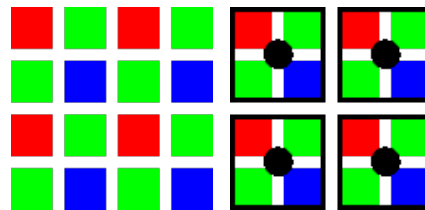
It is helpful at this point to discuss some elements of how cameras work. Since no one uses film anymore, we will talk about digital cameras. We will also only concentrate on the parts having to do with color. There are a lot of interesting things we could say about the physics of a camera’s lens and aperture, but this lecture is long enough already.

Digital cameras work just like you’d think they would. They have a bunch of little cavities where light enters. These are the pixels. The pixels get covered with a filter which only lets in certain light. Something like this:



**Figure 32.** Components of the sensor of a digital camera

So each cavity can only take one color. Since we want a basis of 3 colors, how do we arrange the pixels? What is normally done is to use a **Bayer array**, which contains twice as many green pixels as red or blue. It look like this:

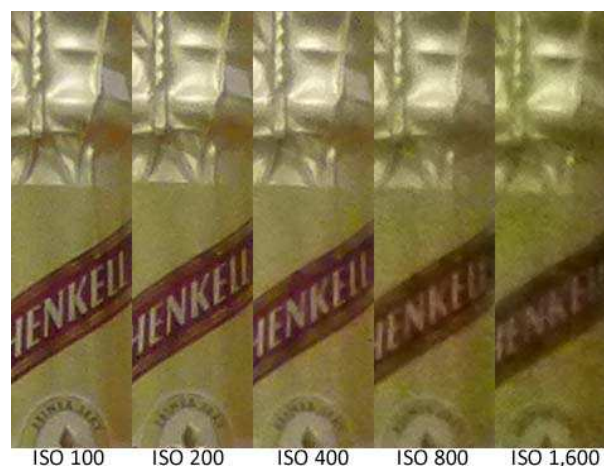


**Figure 33.** Bayer array. The 4 squares of pixels are turned into RGB values at each black point.

The colors are then combined into RGB values. Digital cameras generally have twice as much luminosity resolution in green than in red or blue, so the red and blue colors appear twice as noisy. This was chosen because we are maximally sensitive to green.

The output of the picture is R, G, and B intensity readings for each of the millions of pixels in your camera. In an image, as in an uncompressed .tiff file, each pixel has 8-bits (a number from 0 to 255) of information in each color channel. This is the standard **bit depth** of the image. Recall that we are sensitive to 12 orders of magnitude in intensity. So the first thing the camera has to do is to pick the range of intensities which it wants to map.

The **ISO** (International Organization for Standardization) setting of a camera tells which intensity range to use. If you are in dim light, you should use a high ISO, like 1600 ISO, which multiplies the intensities measured by roughly a factor of 1600. So going from 100 ISO to 1600 ISO, you can increase the sensitivity dramatically. Even in dim light, at 1600 ISO, there may be so few photons coming in that not all the pixels will get lit up. In any case, there will probably not be a nice smooth distribution of intensities, so the picture will often look grainy. On the other end of the range, if it is very bright out, you want to use something like ISO 100, so you will have good resolution of the bright light coming in. Here is a comparison of different ISOs



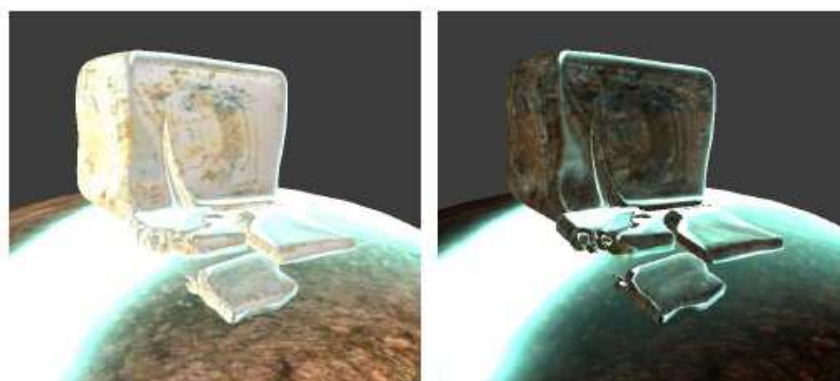
**Figure 34.** ISO comparison. Note how the 1600 ISO is very grainy.

Even with an ISO setting, you often still want to record more than a factor of 256 in intensity. This can be done, if we map

$$I(\lambda) \rightarrow I(\lambda)^\gamma \quad (20)$$

for some  $\gamma$ . The number  $\gamma$  is called the **gamma** factor. It is recorded into your image file when you save the photo. Because our eyes see light on essentially a logarithmic scale, in the same way we hear sound on a logarithmic scale, this gamma factor is critical to getting an image to look right.

If the  $\gamma$  factor is recorded in the image, and then interpreted correctly by your computer screen, or your printer, it just acts like a compression parameter for the spectrum. However cameras are very poor at figuring out the right  $\gamma$ . Moreover, because our minds interpret logarithmic separations in intensity on an essentially linear scale, we have a built in  $\gamma$  factor. So finding the right  $\gamma$  is not just a matter of recording the image intensity, we also need to know about our perception. A practical result is that changing the  $\gamma$  value can make an image look much more realistic. For example, look at this figure



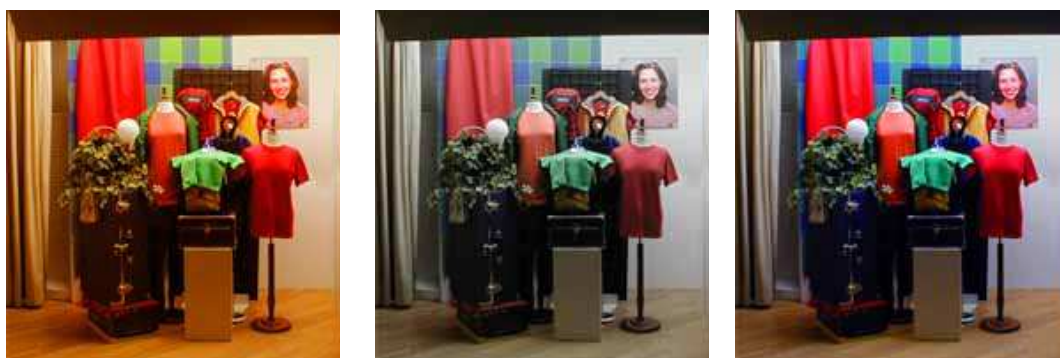
**Figure 35.** Original image (left) image with gamma adjusted (right)

This image has a wide range of intensities, so it is very sensitive to  $\gamma$ . Adjusting  $\gamma$  appropriately can make the image appear much more realistic.

Monitors and TVs have  $\gamma$  correction factors. The intensity of light is produced by the monitor can range over more than 256 bits, due to this  $\gamma$  factor. You can usually calibrate the  $\gamma$  correction by hand so that the images look proper. It is very difficult to get a monitor to look right both during the day, with lots of ambient light, and during the night, without ambient light.

## 9.2 Light Bulbs

Now let's talk about illumination. We can light a scene with all kinds of different light bulbs. Here is the same image lit up with three different bulbs:



**Figure 36.** Same picture, with an incandescent bulb (left), GE cool white bulb (middle), and GE Chroma 50 (right).

The first picture clearly looks yellowish. But which of the other two bulbs looks more like “the real thing”. Unfortunately, there’s no exact answer for that. The colors seem brighter on the right. So we might say the colors in the middle correspond to a dimmer scene, where the reds are muted. But it’s not dimmer – it’s just different lighting. In fact, the scene *is* yellowish in incandescent light. We just don’t usually notice in a room with regular light bulbs because we compensate in our minds using facts we already know about what the scene *should* look like.

To see why light bulbs color things, we can compare the spectrum of daylight and the spectrum of an incandescent bulb

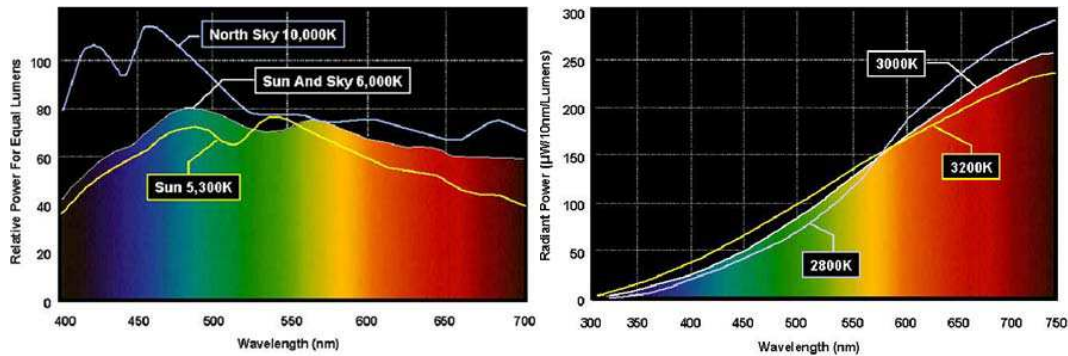


Figure 37. Spectra for daylight (left) and an Incandescent bulb (right).

**Incandescent bulbs** are the traditional lightbulbs. They usually have a tungsten filament which simply heats up and radiates thermally. That is, incandescent bulbs are black bodies, and therefore radiate power smoothly across the entire visible spectrum. They are generally around 3000K, which is less than the sun ( $\sim 5000K$ ), so they usually peak in the infrared. That’s why they are so energy inefficient – most of the power goes into invisible heat radiation. Thus, the visible spectrum is just the tail of the blackbody distribution (see Fig. (10)), so they tend to tint things reddish-yellow.

The other two bulbs lighting the scene in Fig. 36 are CFLs (Compact Fluorescent Lamps). Their power spectra look like this:

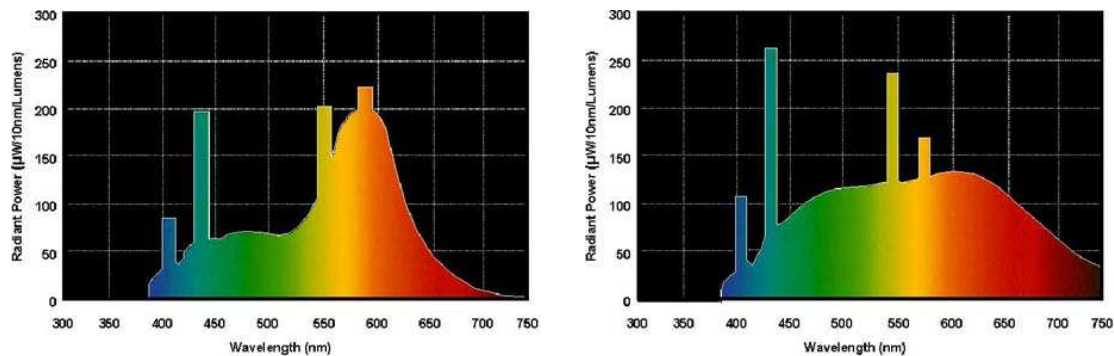


Figure 38. Spectra for a cool white (left) CFL light bulb and the GE Chroma 50 (right) CFL light bulb.

Flourescent bulbs, including the compact kind, are usually filled with mercury. Electrical discharges excite the mercury, which then de-excites, emitting UV radiation. The UV radiation hits phosphors (complicated chemicals – not phosphorous) on the walls of the bulb which convert the UV to visible light. The result of this conversion is the smooth part of the spectrum. The peaks are the visible light emissions from mercury itself. The difference in spectra is due to different kinds of phosphor on the walls of the bulb.



One other thing about lightbulbs is that they are sometimes described as “full spectrum”. That means close-to-a-blackbody. There’s something called the **Color Rendering Index** (CRI), which is defined as the closest distance a light source is to the Planckian Locus on the CIE diagram (see Figure 12). So if you know the temperature of a lightbulb and the CRI, you can place the spectrum on the CIE the diagram. A full-spectrum bulb has CRI above 90%. An incandescent bulb is full-spectrum since it is a blackbody (CRI=100%).

### 9.3 White balance

Now that we understand a bit about lighting, we can ask what a picture is supposed to look like. If we compare the differently-lit photos in Fig. 36, they are all good renditions of what is seen. The incandescent photo looks yellow because it is yellow! But somehow we want it to look as if it weren’t yellow. That’s one of the hardest challenges in photography.

If we use incandescent light, things look yellow. So we want to compensate for this so things look “normal”. What is normal? The generally accepted definition is sunlight:

- **We expect things to look like they were illuminated with sunlight**

So a lot of the art of photography is trying to correct for the way color is portrayed. The best tool we have for this is **white balance**.

The white balance setting on your camera tries to **correct back to sunlight**. This is a very difficult problem, and most cameras suck at it. Here’s an example scene with incandescent lighting and various white balance settings on the camera.



**Figure 39.** awb is auto white balance, then tungsten and manual white balance. The bottom right is an 80a filter put directly on the lens (an old-school solution).

So you see white balance is more than just white, it corrects the whole spectrum. But white is the hardest color to correct for, since it is a careful balance of all the others. One way to correct for it is to use a pure white piece of paper and photograph it with your image. This is still commonly done by professional photographers. An alternative is to load your photo into your favorite editing program and click on a region which you know is white. Then it can try to compensate.

As a final important point, a camera actually takes the picture with more than 8 bits, usually 12. Then it applies your white balance setting and  $\gamma$ -correction to encode the larger range of colors and intensities it sees into a fewer number of bits. This mapping cannot be undone since information is lost in the process. It can be compensated for in post-processing, but the image will be degraded. Some point-and-shoot cameras, and all professional cameras, have a setting called RAW. Raw records all the information that the camera actually sees (the *raw* Bayer array intensities) into a file, with 12 bit depth. Then you can load up your image into photoshop and choose the white balance and  $\gamma$  *afterwards*. If you really want your pictures to look right, this is what you have to do. Keep in mind, if you edit a picture on a computer, the monitor has its own gamma correction, which must be calibrated, then you also have to make sure your printer is calibrated too.

# Lecture 18:

## Antennas and interference

### 1 Visualizing radiation patterns

In the Lecture 16, we showed that an accelerating charge produces a field which decays with distance differently in different directions. In the plane transverse to the acceleration, the field dies like  $E_\theta \sim \frac{1}{r}$  while in the direction parallel to the acceleration it decays much faster like  $E_r \sim \frac{1}{r^2}$ . In this lecture, we will be considering more and more elaborate arrangements of sources placed near the origin ( $r=0$ ) and looking at how large the field is very far away from those sources. We call this the **far-field limit**. In the far-field limit, we only care about the transverse component of the electric field  $E_\theta$  since it is parametrically larger than the parallel component.

We will be calculating expressions where this field varies in space and time

$$E_\theta = E_0 e^{i(kr - \omega t + \delta)} \quad (1)$$

We will try to write this complex representation of the electric field in a form where  $E_0$  is real, so that the actual electric field is  $\text{Re}(E_\theta) = E_0 \cos(kr - \omega t)$ . The corresponding intensity averaged over time is then  $I = \epsilon_0 \text{Re}(E_\theta)^2 = \frac{1}{2} \epsilon_0 E_0^2$ .

A useful observation is that when you add two fields, the time-averaged intensity only depends on their amplitudes and phase difference, not the phases separately. There are many ways to see this, but the most effective might simply be direct calculation. Say we have a field

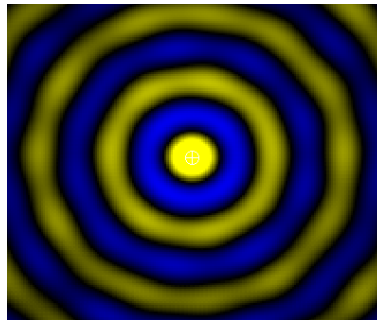
$$E = E_1 e^{i(kx - \omega t + \phi_1)} + E_2 e^{i(kx - \omega t + \phi_2)} \quad (2)$$

Then averaging the intensity over a period  $T = \frac{2\pi}{\omega}$  we get

$$\langle I \rangle = \frac{1}{T} \int_0^T dt \frac{1}{2} \epsilon_0 (\text{Re}[E])^2 = \frac{\epsilon_0}{2} \left\{ \frac{(E_1 + E_2)^2}{2} \cos^2 \frac{\Delta\phi}{2} + \frac{(E_1 - E_2)^2}{2} \sin^2 \frac{\Delta\phi}{2} \right\} \quad (3)$$

where  $\Delta\phi = \phi_1 - \phi_2$ . The reason the separate phases drop out is that we can always write the sum of two waves as having an overall, average phase, and a phase difference. The average phase combines with the time oscillation and gets averaged out, but the phase difference does not. Thus, in the following we will concentrate on phase differences among sets of waves. We will go back and forth between talking about fields and intensities, with the time-averaging of the intensities always being implicit.

To begin, say we just have a group of charges exactly at the origin moving up and down in the  $z$  direction. Such an arrangement is called a **monopole antenna**. Now, consider how the amplitude of the electric field produced looks in the  $x$ - $y$  plane. It will be  $E_\theta \sim \frac{1}{r} \cos(kr - \omega t)$ . We can draw this pattern as follows

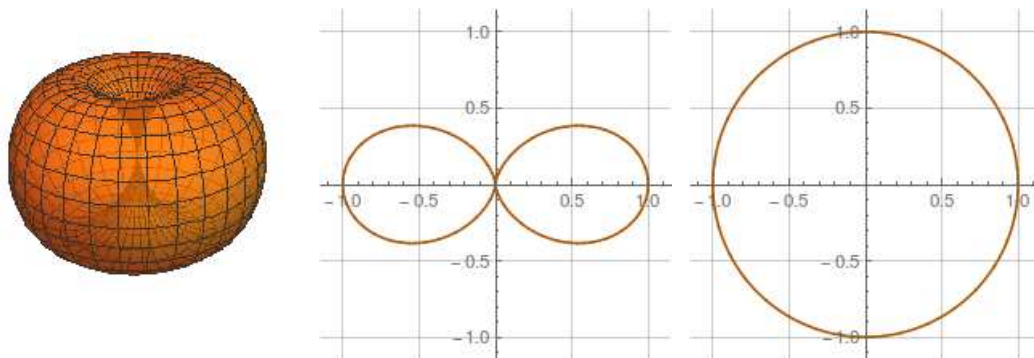


**Figure 1.** A monopole antenna produces the radiation pattern as generated in Falstad's ripple program.

In this visualization, the blue means  $E_\theta > 0$  the yellow means  $E_\theta < 0$  and the black is  $E_\theta = 0$ . The brightness of the color represents  $|E_\theta|$ .

Now, pictures like Fig. 1 show how the amplitude of the electric field varies in two dimensions. To understand how much power can be received from a transmitter, we instead want to know how large the intensity is at different points in three dimensions. So let's imagine a spherical shell of radius  $R$  surrounding the antenna with the antenna direction pointing from the north pole to the south pole and consider how large the generated field is around the shell. The amplitude and intensity are given by functions of the latitude and longitude angles  $\theta$  and  $\phi$  around the sphere. We can then plot a surface where the radial distance from the origin in the plot is the intensity  $I(\theta, \phi)$ .

For example, for the monopole antenna, the surface looks like



**Figure 2.** Radiation pattern for a monopole antenna. Left is the 3D pattern. Middle and right show vertical and equatorial slices, respectively.

This is known as a **3D radiation pattern**.

To make sure this is clear, let's review what is plotted. In Lecture 16 we showed that an accelerating charge produces a field whose components are  $E_r = \frac{q}{4\pi\epsilon_0 R^2}$  and  $E_\theta = \frac{qa}{4\pi\epsilon_0 c^2 R} \sin\theta$ . For an antenna, the net charge is neutral so  $E_r = 0$  and only  $E_\theta$  matters. At a fixed  $R$  very far away from the source,  $E_\theta \sim \sin\theta$ , where  $\theta$  is the angle to the direction of the accelerating charges. Thus the 3D pattern above is a spherical plot of the surface  $R = \sin^2\theta$ ; the middle pattern above is a polar plot of the contour  $R = \sin^2\theta$  which is a vertical slice through the 3D plot; the right pattern is a polar plot along the equator  $\theta = \frac{\pi}{2}$ , so  $R = 1$ . There is a notebook `Interference.nb` on the canvas site which generates these patterns.

Since the field dies off so fast in the vertical direction, we are often only interested in the 2D pattern in the plane of the equator, as in the right panel of Figure 2. In this case, since the intensity is constant along the equator at a distance  $r$  from the antenna, the equatorial cross section of the 3D radiation pattern is just a circle. When we put more antennas together, the resulting interference patterns will be more interesting, as we will now see. So we will work with two different 2D visualizations: the Falstad ripple type, as in Fig. 1 which attempts to show the phase of the amplitude, and the 2D projections as in the right panel of Fig. 2 which shows the relative intensity as a function of angle.

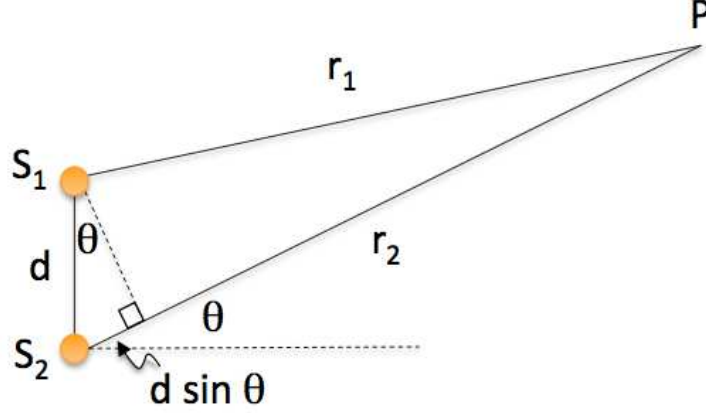
## 2 Two sources

What happens if we put two sources with the same amplitude and in phase a distance  $d$  apart?

Let's start with a very small separation. If the distance  $d \ll \lambda$ , where  $\lambda$  is the wavelength, then the waves will be entirely in phase everywhere. This radiation pattern will look just like the single source but more intense. Say the field from one source is  $E_0$ ; the intensity from one source is therefore  $I_0 = \frac{1}{2}\epsilon_0 E_0^2$ . Two sources at the same spot will produce a field  $2E_0$  and an intensity  $4I_0$ . Since power is conserved, one factor of 2 comes from having two sources instead of 1, and the other factor of two comes from source loading, as we discussed in the context of sound.



Now consider moving the sources farther apart. For a general  $d$ , the field at a distant point is given by the sum of the fields produced from the two charges. The picture is like this



**Figure 3.** Two antennas spaced  $d$  apart produce fields which go to  $P$ .

Let us denote by  $r_1$  and  $r_2$  the distances from  $S_1$  and  $S_2$  to the point  $P$  where we want to evaluate the field, as shown in the picture. Then the field at  $P$  produced from the two sources are

$$E_1 = E_0 e^{i(kr_1 - \omega t)}, \quad E_2 = E_0 e^{i(kr_2 - \omega t)}, \quad (4)$$

So, by linearity, the total field at  $P$  is

$$E_P = E_1 + E_2 = E_0 e^{-i\omega t} (e^{ikr_1} + e^{ikr_2}) \quad (5)$$

Even though the sources are in phase, since the distance to the far point is different from the two sources, the fields may or may not add coherently.

Let us call  $\theta$  the angle between the line connecting  $S_2$  to  $P$  and the  $x$  axis. In terms of  $\theta$ , the difference in the distance the light travels from  $S_2$  and  $S_1$  to get to  $P$  is (see Fig. 3):

$$\Delta r = r_2 - r_1 = d \sin \theta \quad (6)$$

The corresponding phase difference is therefore

$$\Delta \phi = 2\pi \frac{\Delta r}{\lambda} = 2\pi \frac{d}{\lambda} \sin \theta \quad (7)$$

To see the effect on the amplitude, define

$$r = \frac{r_1 + r_2}{2} \quad (8)$$

Then

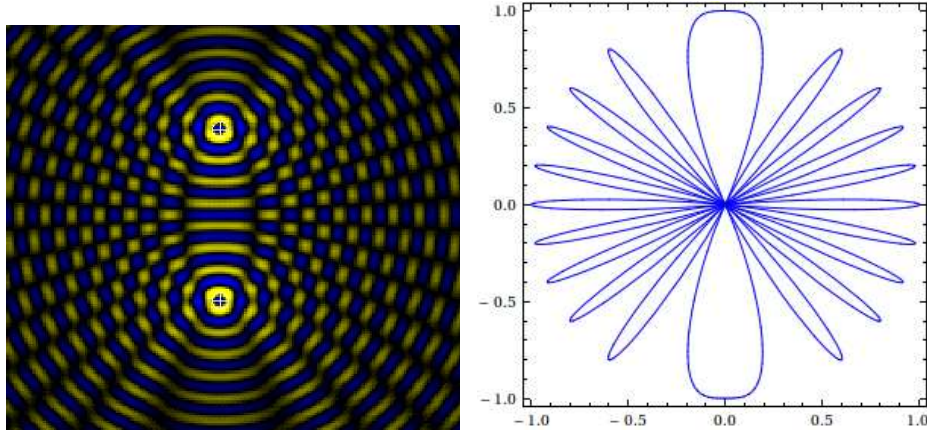
$$\begin{aligned} E_P &= E_0 e^{-i\omega t} \left[ e^{ik\left(r + \frac{\Delta r}{2}\right)} + e^{ik\left(r - \frac{\Delta r}{2}\right)} \right] \\ &= E_0 e^{-i\omega t} e^{ikr} \left[ e^{i\frac{k\Delta r}{2}} + e^{-i\frac{k\Delta r}{2}} \right] \\ &= 2E_0 e^{-i\omega t} e^{ikr} \cos\left(k\frac{\Delta r}{2}\right) \\ &= 2E_0 e^{-i\omega t} e^{ikr} \cos\left(\pi \frac{d}{\lambda} \sin \theta\right) \\ &= 2E_0 e^{-i\omega t} e^{ikr} \cos\left(\frac{\Delta \phi}{2}\right) \end{aligned}$$

So the intensity (averaged over time) is

$$I = \frac{1}{2}\epsilon_0 \text{Re}(E_P)^2 = 4I_0 \cos^2\left(\pi \frac{d}{\lambda} \sin\theta\right) = 4I_0 \cos^2\left(\frac{\Delta\phi}{2}\right) \quad (9)$$

where  $I_0 = \frac{1}{2}\epsilon_0 E_0^2$  is the intensity from a single source, as above. Taking the limit  $d \ll \lambda$  we find that  $I = 4I_0$ , which is in agreement with total coherence, as discussed above.

In limit  $d \gg \lambda$ , the radiation pattern has a bunch of maxima and minima



**Figure 4.** Two coherent sources separated by  $d \gg \lambda$ . Left shows the amplitude from Ripple, and the right shows the radiation pattern.

In the  $d \gg \lambda$  case, the argument of the cosine,  $\pi \frac{d}{\lambda} \sin\theta$  goes around the circle many times. The average intensity along a circle is

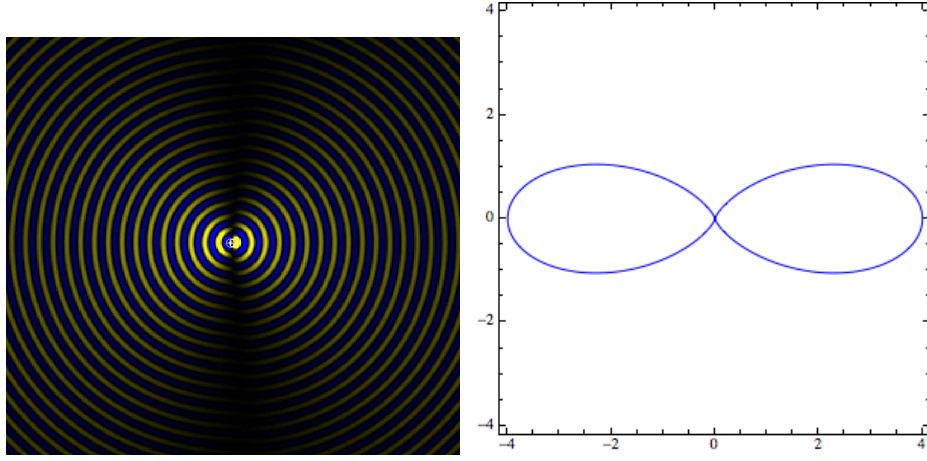
$$\langle I \rangle = \frac{1}{2\pi} \int_0^{2\pi} 4I_0 \cos^2\left(\pi \frac{d}{\lambda} \sin\theta\right) d\theta \approx 2I_0 \quad (d \gg \lambda) \quad (10)$$

The integral has been evaluated by replacing  $\cos^2\theta$  by its average  $\frac{1}{2}$  (if you don't trust this approximation, go ahead and can check it numerically yourself). Thus when the sources are farther than  $\lambda$  apart, all the constructive interference effects cancel on average – the total power emitted is then just the sum of the power begin emitted by the sources.

Ok, so neither  $d \ll \lambda$  or  $d \gg \lambda$  are particularly interesting. What happens if the distance between the sources is half a wavelength:  $d = \frac{\lambda}{2}$ ? In this case, when the wave from one source gets to the other, there will be complete destructive interference. Thus, along the direction of the line between the sources, even far away from the sources, the intensity will be zero. On the other hand, on the line which goes perpendicular to the sources, there must be constructive interference. So we get  $I = 4I_0$  along that line. We can confirm these assessments with a direct evaluation of Eq. (9) with  $\lambda = \frac{d}{2}$ :

$$I = 4I_0 \cos^2\left(\frac{\pi}{2} \sin\theta\right) = I_0 \times \begin{cases} 4, & \theta = 0 \\ 0, & \theta = \frac{\pi}{2} \\ 4, & \theta = \pi \\ 0, & \theta = \frac{3\pi}{2} \end{cases} \quad (11)$$

So it has maxima at  $\theta = \frac{\pi}{2}$  and  $\frac{3\pi}{2}$  and minima at 0 and  $\pi$ . The pattern looks like this



**Figure 5.** Amplitude and radiation pattern for two antennas separated in the  $y$  direction by  $d = \frac{\lambda}{2}$ .

So the field is focuses more in the  $x$  direction than the  $y$  direction.

This is useful: if we have a broadcasting antenna, we can arrange it to broadcast only in the East and West directions, with little power going North and South. That is, we can use interference to direct our transmission.

Unlike the  $d \ll \lambda$  case, the power does not go out uniformly. Recall that in the  $d \ll \lambda$  case the intensity along a circle was constant and the average intensity was just  $\langle I \rangle = 4I_0$ . In this case, the average intensity around the circle is

$$\langle I \rangle = \frac{1}{2\pi} \int_0^{2\pi} d\theta 4I_0 \cos^2\left(\frac{\pi}{2} \sin\theta\right) \approx 1.4I_0 \quad (12)$$

So we are able to achieve a local intensity of  $4I_0$  in the  $\pm x$  direction using only 1.4 times the power of a single source.

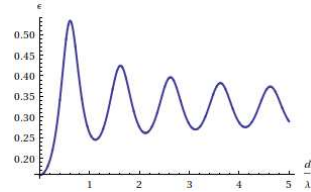
Let us define the transmission efficiency as amount of power going into a particular direction, say  $\theta = 0$ , divided by the power averaged over the whole circle. That is

$$\varepsilon = \frac{I(\theta=0)}{\langle I \rangle} \quad (13)$$

For example, if we had a set of transmitting antennas at the origin and another receiving set far away at  $\theta = 0$ , this efficiency would tell us how much of the emitted power would get to the receiver. With one source  $\varepsilon = 1$ . With two sources with  $d \ll \lambda$ ,  $\varepsilon = 1$  as well. In this case, because the gain in intensity is due to source loading, although more power is being converted into radiation, the transfer of power is not actually more efficient than for a single source. For  $d = \frac{\lambda}{2}$  we found  $\varepsilon = \frac{4}{1.4} = 2.9$ . So an antenna array with two sources set a half wavelength apart is nearly 3 times as efficient as a single source for broadcasting in a particular direction.

The next obvious question is how well can we do? That is, how can we maximize  $\frac{I_0(\theta=0)}{\langle I \rangle}$ ? With two sources spaced  $d$  apart, we want to maximize

$$\varepsilon = \frac{\cos^2\left(\pi \frac{d}{\lambda} \sin(\theta=0)\right)}{\frac{1}{2\pi} \int_0^{2\pi} d\theta' \cos^2\left(\pi \frac{d}{\lambda} \sin\theta'\right)} = \frac{2\pi}{\int_0^{2\pi} d\theta' \cos^2\left(\pi \frac{d}{\lambda} \sin\theta'\right)} = \quad (14)$$



Plotting this function it's not hard to see that the maximum is at  $d = \frac{\lambda}{2}$ . The efficiency at that point is  $\varepsilon = 2.9$ , as above.

Can we do better?

### 3 Phased arrays

We found that if we have two sources, the field at a point  $P$  is determined by the phase shift  $\Delta\phi = 2\pi\frac{d}{\lambda}\sin\theta$  from the different distance light has to go to get to  $P$  from the two sources. Then,

$$E_P = 2E_0e^{-i\omega t}e^{ikr}\cos\left(\frac{\Delta\phi}{2}\right) \quad (15)$$

How does the pattern change if we decide to produce the sources out of phase. Say the sources differ by a phase  $\delta$ .

This phase shift then add to the phase shift  $\Delta\phi$  from the path lengths at the point  $P$  and we get

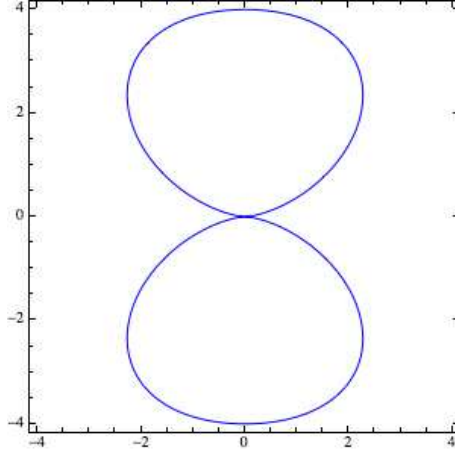
$$E_P = 2E_0e^{-i\omega t}e^{ikr}\cos\left(\frac{\Delta\phi + \delta}{2}\right) \quad (16)$$

So that

$$I = 4I_0\cos^2\left(\frac{\Delta\phi + \delta}{2}\right) = 4I_0\cos^2\left(\pi\frac{d}{\lambda}\sin\theta + \frac{\delta}{2}\right) \quad (17)$$

How does  $\delta$  affect the radiation pattern?

For example, consider again the case with two antennas spaced  $d = \frac{\lambda}{2}$  apart in the  $y$  direction, but now have them be broadcasting exactly out of phase  $\delta = \pi$ . The radiation pattern from Eq. (17) looks like



**Figure 6.** Radiation pattern with  $d = \frac{\lambda}{2}$  and the two sources exactly out of phase ( $\delta = \pi$ ).

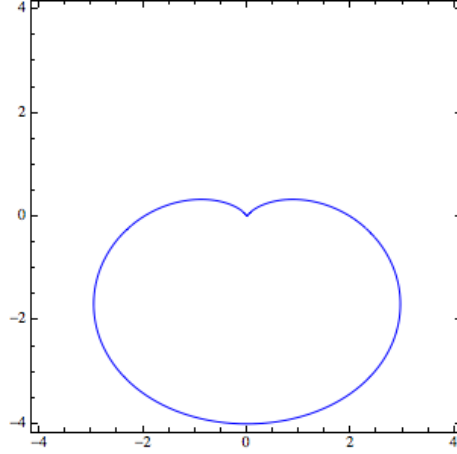
This looks just like Fig. 5, but rotated  $90^\circ$ . We can check that this makes sense. When a wave leaves the bottom antenna and goes towards the top one, it will have rotated by  $\pi$  since the two antennas are exactly half a wavelength apart. This phase shift of  $\pi$  adds to the phase shift  $\delta = \pi$  we put in by hand so that the wave going upward from the bottom antenna is now exactly in phase with the wave from the second antenna. Thus they will continue to be in phase as we move upward, which explains the upward lobe (and similarly the downward lobe).

This result is very practical: if we want to broadcast North-South instead of East-West we don't have to go up on the roof and reorient the antenna. Instead, we can simply insert a delay into the current driving one antenna to change its phase by  $\pi$ . If you play with the plots in the Interference.nb Mathematica notebook, you can see how the pattern changes as  $\delta$  goes from 0 to  $2\pi$ .

Now consider the case with  $d = \frac{\lambda}{4}$  and  $\delta = \frac{\pi}{2}$ . Then,

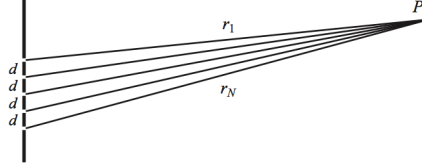
$$I(\theta) = 4I_0\cos^2\left(\frac{\pi}{4}\sin\theta + \frac{\pi}{4}\right) = I_0 \times \begin{cases} 2, & \theta = 0 \\ 0, & \theta = \frac{\pi}{2} \\ 2, & \theta = \pi \\ 4, & \theta = \frac{3\pi}{2} \end{cases} \quad (18)$$

Which looks like this



So now the radiation points mostly in a single direction. This is clearly a much more efficient way of transmitting a signal from point to point. Now our antenna array can radiate South only, with little going North, East or West.

To do better, we can put a bunch of sources in a row. Say we have  $N$  sources in a line, each separated from the previous by a distance  $d$  with a phase shift of  $\delta$ . This is known as a **phased array**:



Then far away, at an angle  $\theta$  to the array, each one has a phase shift of

$$\Delta = 2\pi \frac{d}{\lambda} \sin\theta + \delta \quad (19)$$

from the previous one. If  $E_0$  is the field from one source, then the net field at  $P$  is

$$E_P = E_0 e^{-i\omega t} e^{ikr} (1 + e^{i\Delta} + e^{2i\Delta} + \dots + e^{(N-1)i\Delta}) \quad (20)$$

To sum this, we use the mathematical formula

$$\sum_{n=0}^{N-1} x^n = \frac{x^N - 1}{x - 1} \quad (21)$$

So that

$$\begin{aligned} E_P &= E_0 e^{-i\omega t} e^{ikr} \sum_{n=0}^{N-1} (e^{i\Delta})^n \\ &= E_0 e^{-i\omega t} e^{ikr} \frac{e^{iN\Delta} - 1}{e^{i\Delta} - 1} \\ &= E_0 e^{-i\omega t} e^{ikr} e^{iN\frac{\Delta}{2}} e^{-i\frac{\Delta}{2}} \left[ \frac{\sin\left(N\frac{\Delta}{2}\right)}{\sin\left(\frac{\Delta}{2}\right)} \right] \end{aligned}$$

Averaging over time the phases drop out and the intensity is then

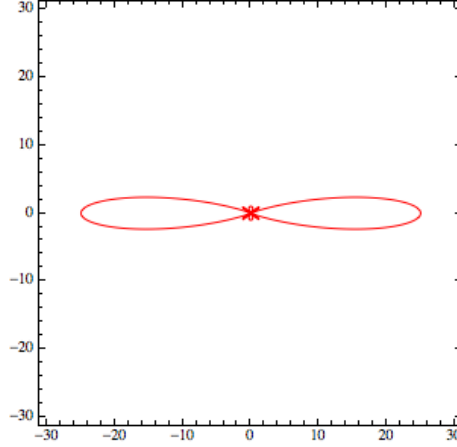
$$I = I_0 \frac{\sin^2\left(N\frac{\Delta}{2}\right)}{\sin^2\left(\frac{\Delta}{2}\right)} \quad (22)$$

We should check this against our previous results. For  $N = 1$  we find  $I = I_0$ . For  $N = 2$  we can use  $\frac{\sin(2x)}{\sin(x)} = 2\cos(x)$  to see that  $I = 4I_0\cos^2\frac{\Delta}{2}$  as in Eq. (17).

Since  $\Delta = 2\pi\frac{d}{\lambda}\sin\theta + \delta$ , our general formula for  $N$  antennas with a phase shift of  $\delta$  between each pair is:

$$I = I_0 \frac{\sin^2\left(N\left(\pi\frac{d}{\lambda}\sin\theta + \frac{\delta}{2}\right)\right)}{\sin^2\left(\pi\frac{d}{\lambda}\sin\theta + \frac{\delta}{2}\right)} \quad (23)$$

For  $N = 5$ ,  $d = \frac{\lambda}{2}$  and  $\delta = 0$  this looks like



So the power is more narrowly focused in the  $x$  direction. The efficiency at  $\theta = 0$  is now

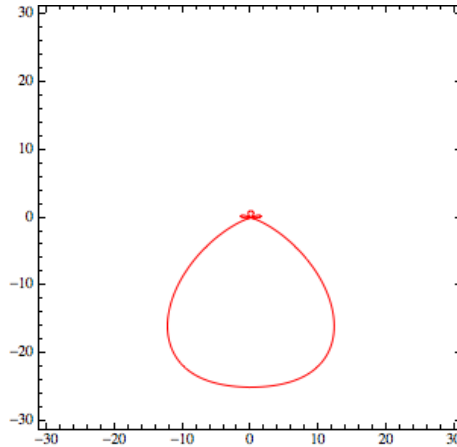
$$\varepsilon = \frac{I(0)}{\langle I \rangle} = \frac{25}{3.5} = 8.18 \quad (24)$$

More generally, the intensity at  $\theta = 0$  with  $\delta = 0$  is

$$I = I_0 \lim_{\theta \rightarrow 0} \frac{\sin^2\left(N\left(\pi\frac{d}{\lambda}\sin\theta\right)\right)}{\sin^2\left(\pi\left(\frac{d}{\lambda}\sin\theta\right)\right)} = I_0 \lim_{\theta \rightarrow 0} \frac{\left(N\pi\frac{d}{\lambda}\theta\right)^2}{\left(\pi\frac{d}{\lambda}\theta\right)^2} = N^2 I_0 \quad (25)$$

On the other hand  $\langle I \rangle$  grows linearly with  $N$ , since when one averages over all directions there is as much destructive as constructive broadcast. Thus, for  $\delta = 0$ , the efficiency  $\varepsilon \sim N$  which grows linearly with the number of antenna.

We can also adjust the phases so that the antenna points in one direction. With  $N = 5$   $d = \frac{\lambda}{4}$  and  $\delta = \frac{\pi}{2}$  we find



What does a phased array of  $N$  evenly spaced sources look like? It looks like



**Figure 7.** Yagi antenna array.

This is called a **Yagi antenna array**.

Phased arrays are very common for radar in the military, since radar can scan very fast in different directions by changing phases electronically, rather than rotating a radar dish. Ever heard of a patriot missile? Did you know that PATRIOT stands for "Phased Array Tracking Radar to Intercept On Target". These missiles are guided by ground-based phased-array radar. Ever wonder what's in the nose of a fighter jet?

Here's a MIG-31 with the nose cone removed:



**Figure 8.** That orange grid under the nose cone of a fighter jet is a phased array.



Fighter jets, and many modern aircraft are equipped with phased array radar. The orange lines you see under the cone are antennas whose phases can be adjusted electronically to point the radar in any direction. There are horizontal antennas in front and vertical ones in back, so the antenna can point in any direction in three-dimensions.

## 4 Other antennas

There are a number of other types of antennas you might want to know about.

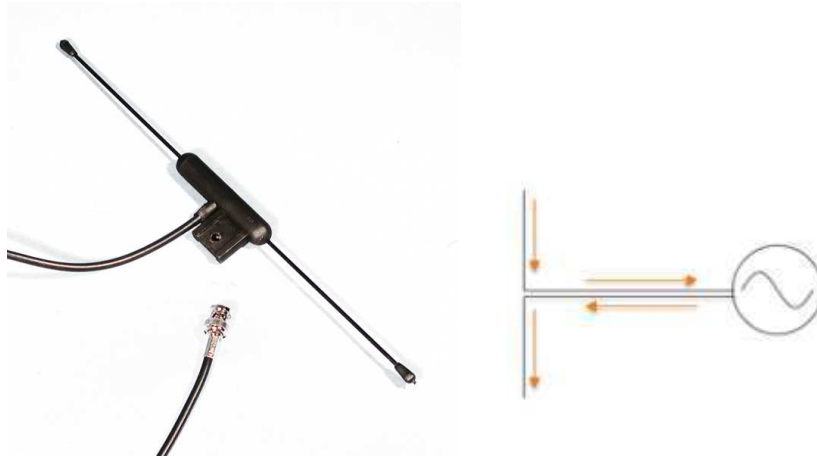
A car antenna is called a **whip antenna** or **monopole antenna**. It is not direction specific, since the car is driving around, there would be no reason to point it in a particular direction



Figure 9. Whip antenna on a car



Next, there are dipole antennas, which have two arms



**Figure 10.** In a dipole antenna charge moves up and down in the two arms coherently. This produces waves which are mostly confined to the plane transverse to the antenna.

Walkie-talkies have antennas that are much shorter. They curl around to pick up the  $\vec{B}$  field rather than the  $\vec{E}$  field. They are also non-directional



**Figure 11.** Rubber ducky antenna

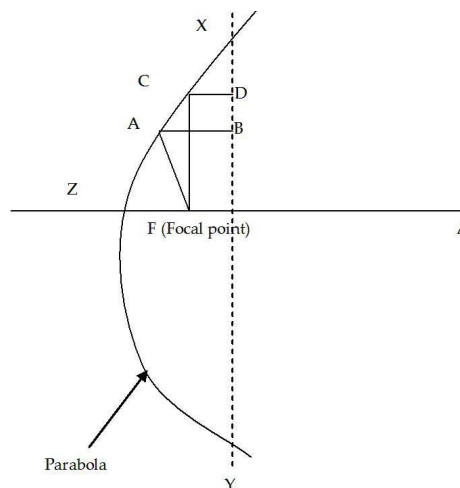
Finally, there are the good old rabbit ears



**Figure 12.** Rabbit ear antenna.

Fiddling with the straight antennas (the whips) can pick up greater signal. When signals come into a house they can bounce off walls and the polarizations can change. Thus changing the angle of the whips can help pick up particular polarizations. The middle round ring is for picking up magnetic fields. Usually it can rotate to pick up different polarizations as well. The whips pick up VHF (very-high frequency) stations, with frequencies 3 MHz - 300 MHz (100 m-1m). The ring picks up UHF (ultra-high frequency) stations which have frequencies 300 MHz – 1 MHz or wavelengths of 1m to 10cm.

As you probably already know, dish antennas are usually parabolic. The parabolic shape focuses incoming plane waves to a point, the focus of a parabola. The picture looks like this

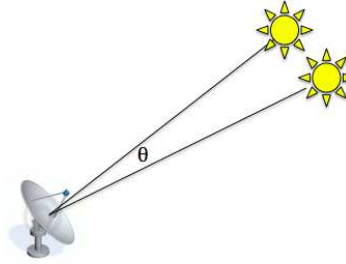


A plane wave has the same phase at  $D$  and  $B$ . Using Snell's law and some calculus you can then show that the shape of the dish must be a parabola for the reflected waves all to be in phase when they reach a single point (the focal point). There is a nice simulation of this in Ripple.

## 5 Interferometry

Sometimes we want to use an antenna to do more than receive a signal as intensely as possible. For example, we might be interested in angular resolution. Suppose two stars are separated by an angle  $\theta$  in the sky. When  $\theta$  is very small, what kind of antenna would be able to make out that there are two stars instead of just one?

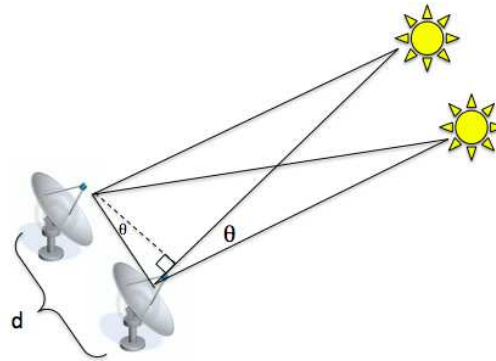
Suppose we only had one antenna to receive the signal.



This antenna would measure the sum of the amplitudes of the waves from the two stars. If the antenna focuses the signal to just a point, then it just receives some amplitude  $A(t) = A_0 \cos(\omega t)$ . One can get the frequency and intensity from such an amplitude, but there is nothing more. A single star or two nearby stars are indistinguishable. A bigger dish collects more light and hence is more **sensitive**. But its angular **resolution** is no better than a simple monopole antenna. Since the wave is still measured at just a single point and has no angular information.

Of course, you can point the dish in a different direction to get some angular resolution. But there's actually a much smarter and more powerful way to increase the angular resolution.

Consider an array of antennas now, each separated from the next by a distance  $d$



(26)

To an excellent approximation, the light coming from both stars will be plane waves in the direction between the stars and the earth. This is true no matter what produces the light in the stars, or how far they are from us. Let us write

$$A_{\text{dish } 1} = A_{\text{star } 1} + A_{\text{star } 2} \quad (27)$$

where  $A_{\text{star } 1}$  and  $A_{\text{star } 2}$  are the amplitudes for the two waves, including their phases, when they hit dish 1.

For simplicity, let us first consider the case where the antennas are in a line perpendicular to the direction to star 1. Then the light from star 1 will hit all the antennas with exactly the same phase. Since star 2 is offset from star 1 by an angle  $\theta$ , its light will not have the same

phase at both antennas. The phase difference will between the two waves at dish 2 is the usual  $\Delta\phi = 2\pi\frac{d}{\lambda}\sin\theta$ . So

$$A_{\text{dish 2}} = A_{\text{star 1}} + A_{\text{star 2}} \cos(\Delta\phi) \quad (28)$$

Similarly if we for a third dish equally spaced from the other two, its amplitude is

$$A_{\text{dish 3}} = A_{\text{star 1}} + A_{\text{star 2}} \cos(2\Delta\phi) \quad (29)$$

Now we have three equations. Assuming that these amplitudes are measured perfectly, we can solve them for the three unknowns  $A_{\text{star 1}}$ ,  $A_{\text{star 2}}$  and  $\Delta\phi$ .

Note that if we only had the intensity instead of the amplitude at each antenna, the phase difference would simply get time-averaged away and each dish would measure the same intensity. This way of resolving small angles is an example of **interferometry**. It requires amplitude-level phase information.

What is the angular resolution of an antenna array? Well, if the phase difference is too small, the signal will get washed out by noise. The biggest the phase difference can possibly be is  $\Delta\phi = \pi$ . Basically, if  $\Delta\phi \ll \pi$  we are not going to see much interference. For small angles  $\theta$ ,  $\Delta\phi = 2\pi\frac{d}{\lambda}\sin\theta \approx 2\pi\frac{d}{\lambda}\theta$ . Thus setting  $\Delta\phi = \pi$  we find

$$\theta_{\text{resolvable}} = \frac{\lambda}{2d} \quad (30)$$

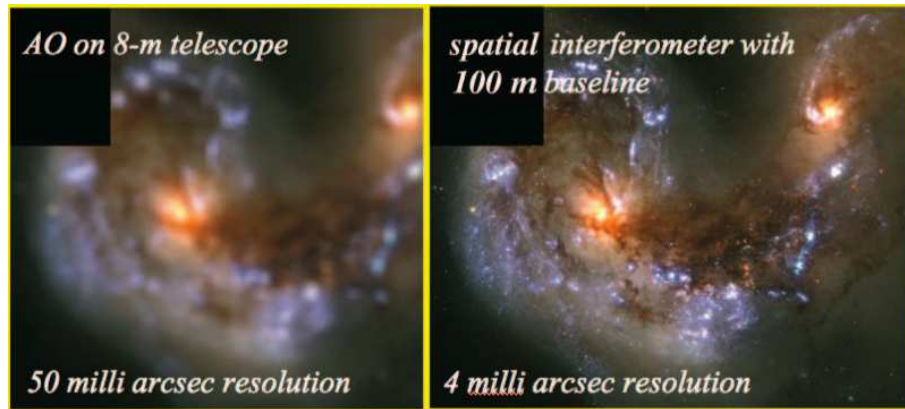
If you can resolve phase shifts smaller than  $\pi$ , you might do a little better than this. The important point is the scaling of the resolvable angles with wavelength and the distance  $d$  between the the antennas. The bigger  $d$  is the smaller the angles are that you can resolve. For this reason, you want very large interferometers. Also, the smaller the wavelengths, the smaller angles you can resolve.

There is an interferometer in New Mexico called the VLA (Very Large Array) which looks like this:



**Figure 13.** VLA array in New Mexico has 27 antennas each 25 meters in diameter. The antennas can be moved up to a maximum baseline of 22 miles

Here is an example of the improvement in angular resolution which you can get with an interferometer.

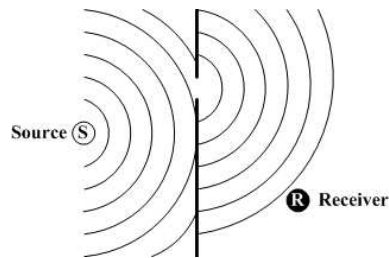


**Figure 14.** The left is how well one can see with a single 8m diameter telescope. The second shows what you can see with 2 such telescopes separated by 100 m.

# Lecture 19: Diffraction and resolution

## 1 Huygens' principle

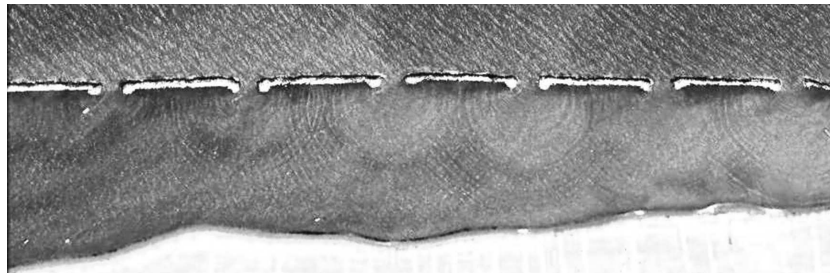
Diffraction refers to what happens to a wave when it hits an obstacle. The key to understanding diffraction is a very simple observation first due to Huygens in 1678. Say a wave arrives at an opaque screen with a little hole in it. On the other side of the screen, the wave equation must still be satisfied with boundary conditions given by the motion of the wave in the hole. That is, the solution is identical to a situation where there was a source in the hole. The wavefront diagram looks like this:



(1)

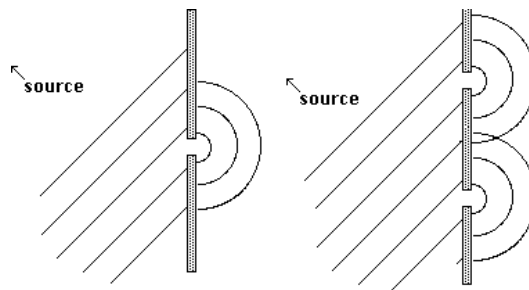
The key point is that the receiver does not care about whatever happened on the other side of the screen. All that matters at is what the wave was like as it passed through the hole.

Here is a picture of real wavefronts near a beach in Italy, taken from a satellite, showing this phenomenon:



**Figure 1.** Waves near a shore somewhere in Italy. This image was taken from Google Earth.

The great thing about this way of thinking about diffraction is that, since the wave equation is linear, you can use this trick for any number of holes. You simply add the amplitude for the waves produced from a “source” at each hole:



(2)

Computing the amplitude by adding point sources in this way is known as **Huygens' principle**.

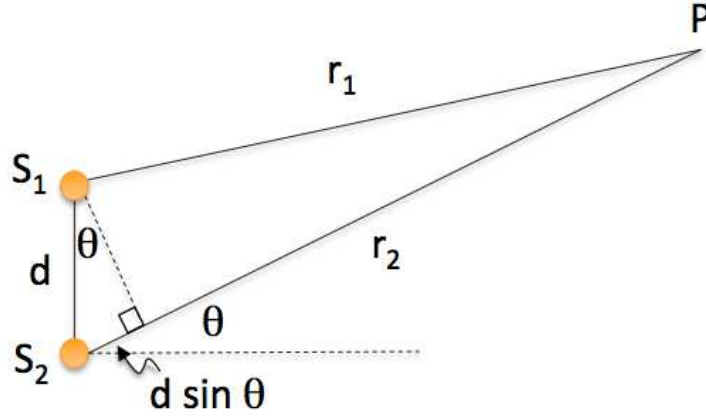
Huygen's principle works even if the holes are very close together. In fact, it works if they are connected, so instead of a hole, it's a slit. One just has to integrate over the holes rather than sum over them. We'll begin with a discrete set of holes as in a diffraction grating. Then

we'll take the continuum limit and talk about slits.

## 2 Multiple hole diffraction

Using Huygens' principle, we can easily compute the diffraction pattern from a plane wave passing through any number of holes. Say there are  $N$  holes in a row separated by a distance  $d$ . The solution will be as if there are  $N$  sources separated by a distance  $d$ . Thus the pattern will be the same as the antenna pattern for this source configuration we computed in Lecture 18!

Say one source has field  $E = E_0 e^{-i\omega t} e^{ikr}$ . Recall the diagram for two sources separated by a distance  $d$ :



**Figure 2.** Phase difference from two holes is the same as from two sources.

We found that at a distant point  $P$  at an angle  $\theta$  to the sources the field is

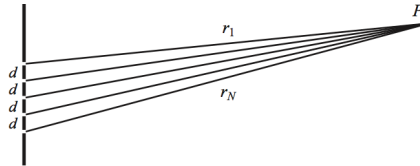
$$E_P = 2E_0 e^{-i\omega t} e^{ikr} \cos\left(\frac{\Delta}{2}\right) \quad (3)$$

where

$$\Delta = 2\pi \frac{d}{\lambda} \sin \theta + \delta \quad (4)$$

with  $\delta$  the difference in phase between the two sources. This phase difference will be  $\delta = 0$  if the sources are illuminated by a plane coming perpendicular to the separation between the sources. If the wave comes on an angle,  $\delta \neq 0$ .

For  $N$  slits, we get  $N$  sources, and the picture looks like this



This is called a **diffraction grating**. In this case, the amplitude at  $P$  is (again, see Lecture 18)

$$E_P = E_0 \left[ \frac{e^{iN\frac{\Delta}{2}}}{e^{i\frac{\Delta}{2}}} \right] \left[ \frac{\sin\left(N\frac{\Delta}{2}\right)}{\sin\left(\frac{\Delta}{2}\right)} \right] \quad (5)$$

So the intensity is

$$I = I_0 \frac{\sin^2\left(N\frac{\Delta}{2}\right)}{\sin^2\left(\frac{\Delta}{2}\right)} \quad (6)$$

where  $I_0$  is the intensity from a single source.

Since the slits are spaced  $d$  apart, the total size of the diffraction grating is  $a = Nd$ . Now say we place a screen at a distance  $L$  from the slits with  $L \gg a$ . Then the height  $y$  up this screen is  $y = L \tan \theta \approx L \theta \approx L \sin \theta$ , and so

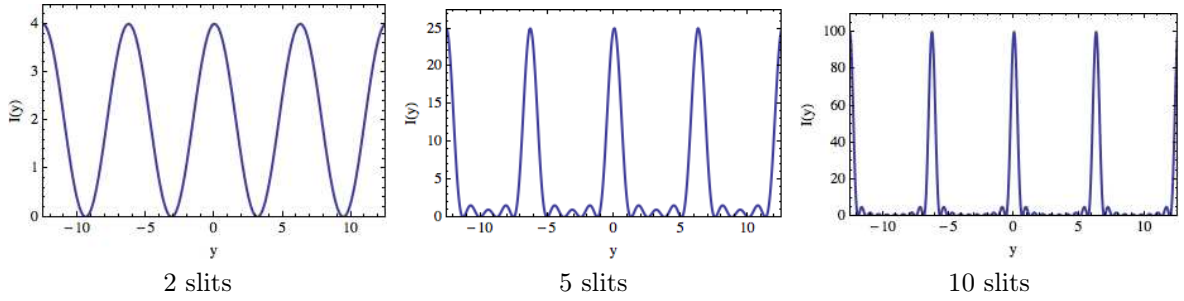
$$\Delta = 2\pi \frac{d}{\lambda} \frac{y}{L} + \delta \quad (7)$$

If the plane wave hits the grating straight on, so  $\delta = 0$  then

$$I(y) = I_0 \frac{\sin^2\left(\pi N \frac{dy}{\lambda L}\right)}{\sin^2\left(\pi \frac{dy}{\lambda L}\right)} \quad (8)$$

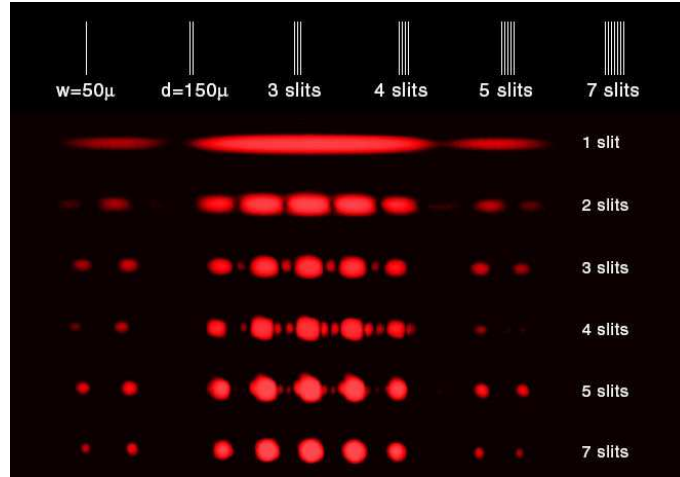
This is the intensity pattern from a diffraction grating. Note that for  $N = 2$  this simplifies to  $I(y) = 4I_0 \cos^2\left(\pi \frac{dy}{\lambda L}\right)$ , which we computed in the previous lecture.

Here are some plots of  $I(y)$  for  $N = 2$ ,  $N = 5$  and  $N = 10$  slits:



**Figure 3.** Patterns from diffraction gratings with  $N = 2, 5$  and  $10$  slits. These plots show the intensity at a distance  $y$  up a screen placed at a distance  $L$  from the wall.

Here's what this actually looks like with a laser passing through a diffraction grating



**Figure 4.** Pattern from a diffraction grating. As the number of slits is increased, the peaks become clearly separated. The extra dots on the right and left side after the gaps are due to the finite width of the slits, as discussed in Section 4.

You can see from the plots that each intensity pattern has a maximum at  $y = 0$ . Indeed, we find

$$I(0) = \lim_{y \rightarrow 0} I(y) = I_0 N^2 \quad (9)$$

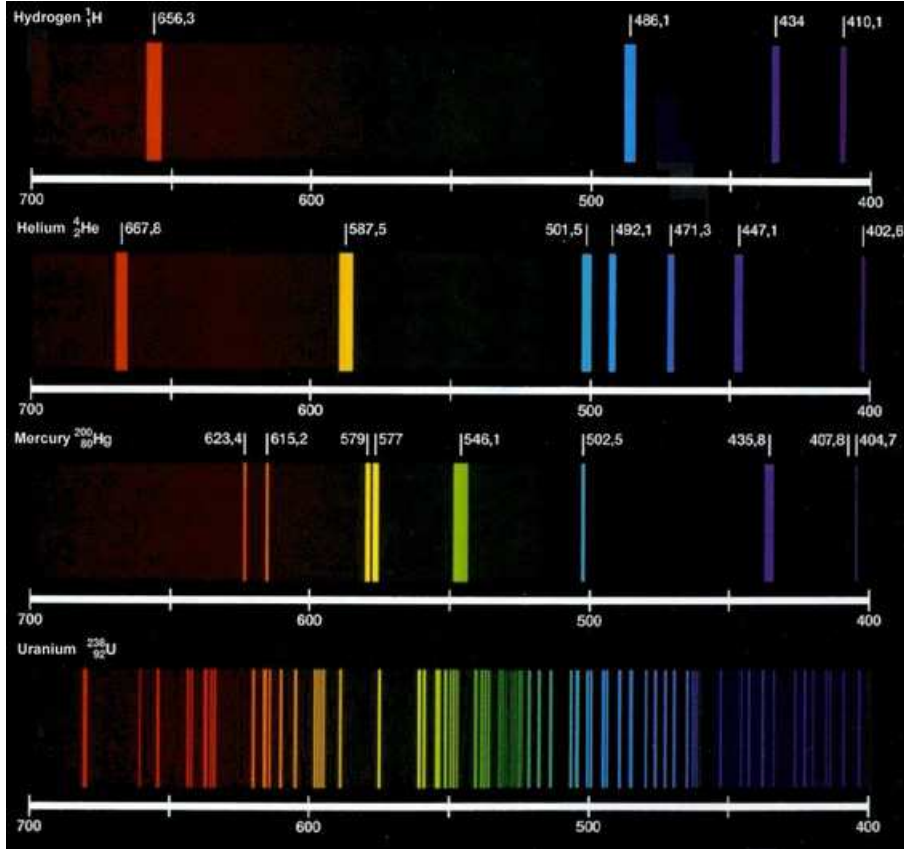


which agrees with the vertical axes on the plots. So the more slits in the grating, the larger the central peak will be. The small bumps are due to the oscillation of the  $\sin^2\left(\pi N \frac{yd}{\lambda L}\right)$  factor in the numerator. Indeed, the factor of  $N$  makes these oscillations have very high frequency. The overall periodicity is due to the denominator factor  $\sin^2\left(\pi \frac{dy}{\lambda L}\right)$ . This vanishes when

$$y_{\max} = j \frac{L\lambda}{d} = j \frac{L}{d} \frac{c}{\nu}, \quad j = 0, 1, 2, \dots \quad (10)$$

with  $\nu$  the frequency.

Since the peaks are spaced by  $\frac{L}{d}\lambda$ , different wavelengths will have peaks in different places. Thus, diffraction gratings can be used to characterize the spectra of various things. For example, gases have interesting spectra which can be resolved with diffraction gratings.



**Figure 5.** Spectra of hydrogen, helium, mercury and uranium as viewed through a diffraction grating.

The wavelengths in the hydrogen spectrum have the interesting pattern  $\frac{1}{\lambda} = \frac{1}{91\text{nm}}\left(\frac{1}{4} - \frac{1}{j^2}\right)$  for  $j = 3, 4, 5$  and  $6$ . This phenomenological observation was known in the 19th century. We will explain where it comes from in Lecture 21, on quantum mechanics.

### 3 Resolving power

An important use of diffraction gratings is to separate light into different wavelengths, as in Figure 5. Say a gas emits two colors, with wavelengths  $\lambda$  and  $\lambda'$ . How do we know if a given diffraction grating can separate the colors? Obviously, if the two wavelengths are infinitesimally close together, we will not be able to separate them. A reasonable criteria for whether frequencies can be separated is if the maximum of the intensity pattern from one wavelengths lies on top of the minimum of the intensity pattern of the other:

$$\max \lambda' = \min \lambda \quad (11)$$

This is known as the **Rayleigh Criterion**.

All frequencies have maxima at  $y=0$ , so we must look to the second peak to apply the criterion. The second maximum for  $\lambda'$  is at

$$y_{\max} = \frac{L}{d} \lambda' \quad (12)$$

The minima for  $\lambda$  happen only when the  $\sin^2\left(\pi N \frac{yd}{\lambda L}\right)$  numerator in Eq. (8) vanishes, which is when

$$y_{\min} = m \frac{L \lambda}{d N}, \quad m = 1, 2, \dots, N-1, N+1, \dots \quad (13)$$

Values  $m=0$  and  $m=N$  are excluded from this list since for these values the denominator in Eq. (8) vanishes as well and they are associated with maxima not minima.

For the second maximum of  $\lambda'$  to overlap with a minimum of  $\lambda$  we then want

$$\frac{L}{d} \lambda' = m \frac{L \lambda}{d N} \quad (14)$$

Or

$$\frac{\lambda'}{\lambda} = \frac{m}{N} \quad (15)$$

Note that this is independent of  $L$  and  $d$ . We want the value of  $m$  where this ratio is as close to 1 as possible (so that  $\lambda'$  is as close to  $\lambda$  as possible). So we take  $m = N-1$ , referring to the minimum right next to the first maximum. Then  $\frac{\lambda'}{\lambda} = \frac{N-1}{N} = 1 - \frac{1}{N}$ . Thus  $\frac{1}{N} = 1 - \frac{\lambda'}{\lambda} = \frac{\Delta \lambda}{\lambda}$  or

$$\boxed{\frac{\lambda}{\Delta \lambda} = N} \quad (16)$$

where  $\Delta \lambda = \lambda' - \lambda$ . In general  $\frac{\lambda}{\Delta \lambda}$  with  $\Delta \lambda$  is the smallest shift in  $\lambda$  which can be resolved is known as the **resolving power** of an optical system. Eq. (16) gives the resolving power of a diffraction grating.

For example, Hydrogen has wavelengths  $\lambda_{\text{red}} = 656.2 \text{ nm}$ ,  $\lambda_{\text{blue}} = 486.1 \text{ nm}$ ,  $\lambda_{\text{violet}-1} = 434.0 \text{ nm}$  and  $\lambda_{\text{violet}-2} = 410.1 \text{ nm}$ . As you can see from Figure 5 the closest are 434.0 nm and 410 nm with

$$\frac{410}{434 - 410} = 18 \quad (17)$$

So we would need 18 slits to see this separation. With 18 slits the lines would be barely separable. Typical diffraction gratings have 100s slits, so separations of  $\frac{\Delta \lambda}{\lambda} \sim 1\%$  can be resolved. High quality diffraction gratings can have  $\sim 10,000$  slits.

## 4 Wide slits

What happens if instead of having  $N$  holes separated  $d$  apart, we have one big hole? An easy way to derive the result is to take the limit of our previous result, Eq. (8)  $N \rightarrow \infty$  holding the overall size  $a = Nd$  fixed. Writing Eq. (8) in terms of  $N$  and  $I$  we have

$$I(y) = I_0 \frac{\sin^2\left(\pi \frac{Ndy}{\lambda L}\right)}{\sin^2\left(\pi \frac{dy}{\lambda L}\right)} \quad (18)$$

Note that  $I(0) = \lim_{y \rightarrow 0} I(y) = I_0 N^2$  which seems to blow up as  $N \rightarrow \infty$ . However, the relative intensity at a distance  $y$  compared to the intensity at  $y=0$  is finite

$$\frac{I(y)}{I(0)} = \lim_{N \rightarrow \infty} \frac{1}{N^2} \frac{\sin^2\left(\pi \frac{Ndy}{\lambda L}\right)}{\sin^2\left(\pi \frac{dy}{\lambda L}\right)} = \lim_{N \rightarrow \infty} \frac{1}{N^2} \frac{\sin^2\left(\pi \frac{ay}{\lambda L}\right)}{\sin^2\left(\pi \frac{ay}{N\lambda L}\right)} = \left(\frac{\sin\left(\pi \frac{ay}{\lambda L}\right)}{\pi \frac{ay}{\lambda L}}\right)^2 \quad (19)$$

This function is known as the **sinc** function

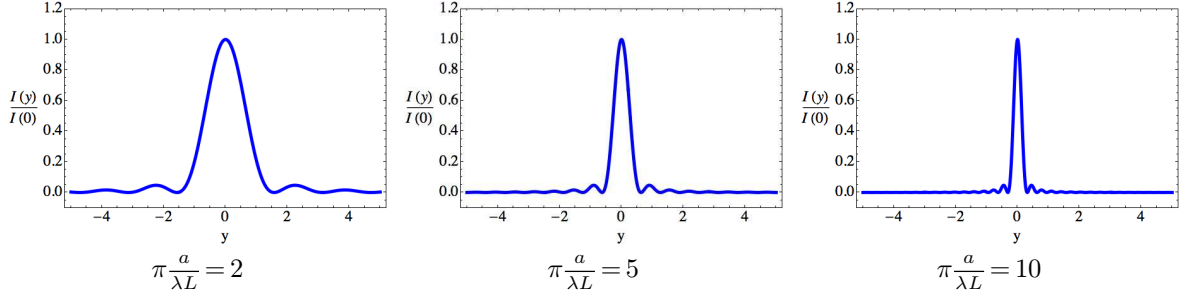
$$\text{sinc}(\beta) \equiv \frac{\sin \beta}{\beta} \quad (20)$$

So

$$\frac{I(y)}{I(0)} = \text{sinc}^2\left(\pi \frac{ay}{\lambda L}\right) \quad (21)$$

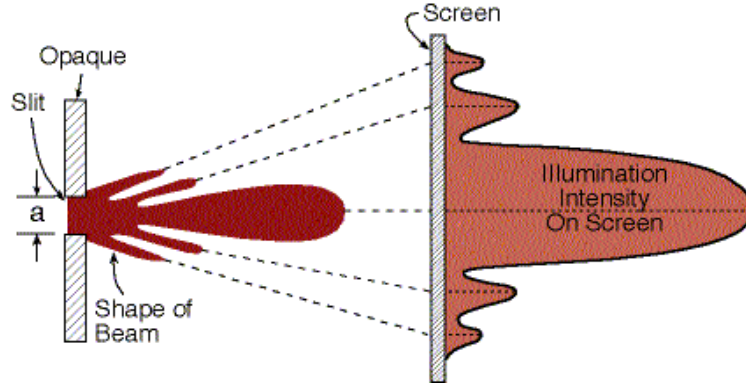
where  $a$  is the size of the slit,  $\lambda$  the wavelength,  $L$  the distance to the screen and  $y$  the distance up the screen.

This pattern looks like



**Figure 6.** The  $\text{sinc}^2(\frac{\pi a}{\lambda L}y)$  function gives the intensity from a wide slit. The wider the slit, the narrower the peak at fixed  $\lambda$ .

Or we can draw the picture like this

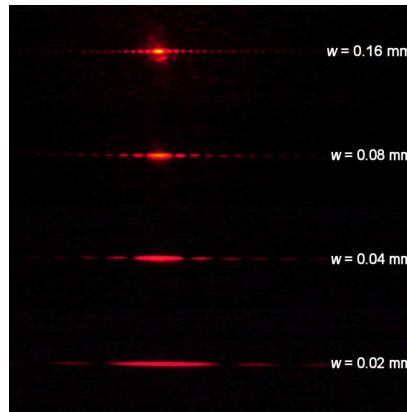


**Figure 7.** Diagram of wide-slit diffraction pattern.

The location of the first minimum in  $\text{sinc}^2(\beta)$  is at  $\beta = \pi$  or  $\pi \frac{ay}{\lambda L} = \pi$ , thus the width of the first peak is

$$\Delta y = \frac{2\lambda L}{a} \quad (22)$$

Note that as  $a$  gets smaller,  $y$  gets bigger. In the limit that  $a$  goes to a point, the outgoing wave is in phase everywhere and  $\Delta y \rightarrow \infty$ . Here's a photo of what this actually looks like



**Figure 8.** Wide slit interference patterns with variable width slits.

Now say we have two far away sources (like stars) and the light passes through a slit, making an image on a screen. What angular separation of the stars will we be able to distinguish on the screen? To find out, we apply the Rayleigh criterion: the maximum of one should hit the minimum of the other. For light coming in at an angle  $\Delta\theta$ , its plane-wave light will have slightly different phases at different parts of the slit. The net effect is to shift the pattern so that the maximum is not at  $y=0$  but at  $y=L\sin\theta$ . So

$$\frac{I(y)}{I(0)} = \text{sinc}^2\left[\pi \frac{a}{\lambda L}(y - L\sin\theta)\right] \quad (23)$$

Say one star's light comes in at  $\theta=0$  and the other at  $\Delta\theta$ . The one at  $\theta=0$  is directly normal to the screen, so the *minimum* of its intensity pattern is at  $y = \frac{\lambda L}{a}$ . If the other one comes in at an angle  $\Delta\theta$ , so its *maximum* is at  $y = L\sin\Delta\theta \approx L\Delta\theta$  for small angles. Setting these equal gives

$$\boxed{\Delta\theta_R = \frac{\lambda}{a}} \quad (24)$$

This is the resolving power of a slit. For a circular aperture, the formula is

$$\Delta\theta_R = 1.22 \frac{\lambda}{a} \quad (25)$$

where the 1.22 comes from the circular geometry. If you have 20/20 vision, you can only resolve objects at around 12 times farther apart than the Rayleigh criterion. That is, ordinary vision is not diffraction limited.

If we have two wide slits, the interference pattern is a combination of the  $\text{sinc}^2\left(\pi \frac{ay}{\lambda L}\right)$  behavior of a wide slit and the  $\cos^2\left(\pi \frac{dy}{\lambda L}\right)$  behavior of the double slit. This is shown in Figure 9.

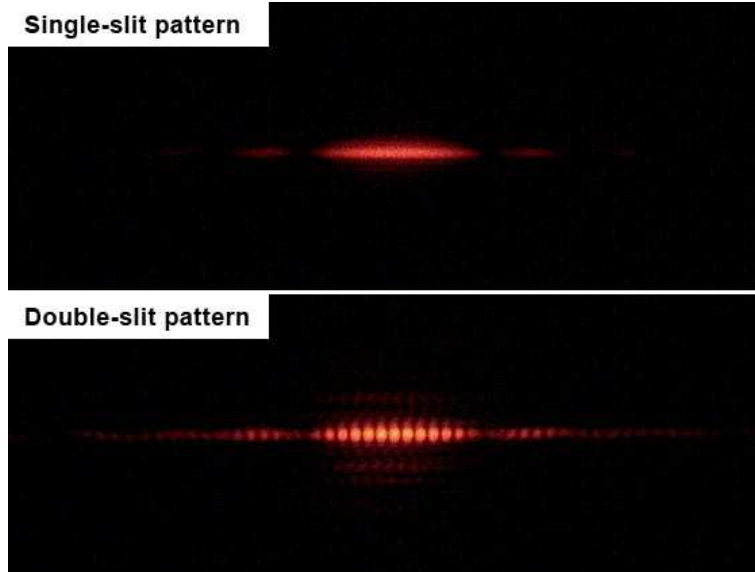


Figure 9. Comparison of the interference pattern for a single wide slit and two wide slits.

## 5 Fourier optics

What would happen if instead of a slit, we put a film in the slit which had a transparency  $T(x)$  which varies along the length of the slit. Let's define  $T=0$  as opaque and  $T=1$  as totally transparent. Thus, the wide slit is  $T(x)=1$ . A half transparent slit would have  $T(x)=\frac{1}{2}$ . The diffraction grating is a bunch of small regions with  $T=1$  surrounded by regions with  $T=0$ . We can model this with a set of  $\delta$ -functions:

$$T(x) = \sum_{j=1}^N \delta(x - jd) \quad (26)$$

What is the intensity pattern on the screen from a general  $T(x)$ ? We can compute the answer using Huygens' principle!

If a plane wave comes in with a field  $E_0$  at the slit, then each point will come with a field given by  $E_0 T(x)$ . At a distance  $y$  on a screen, we can then use our formula that the phase shift is

$$\Delta = 2\pi \frac{x}{\lambda} \frac{y}{L} \quad (27)$$

So the sum of the phases gives

$$E(y) \propto \int dx T(x) e^{2\pi i x \frac{y}{L\lambda}} = \tilde{T}\left(2\pi \frac{y}{L\lambda}\right) \quad (28)$$

where

$$\tilde{T}(k) = \int dx T(x) e^{ikx} \quad (29)$$

is the Fourier transform of  $k$ . That is, the image on the screen is the **Fourier transform of the transparency**. So finally, we understand Fourier optics. You have been using this result in lab all semester, and we have finally derived it.

To make sure we're not nuts, let's check on our wide slit. There we expect the field  $E$  to be the square root of the intensity in Eq. (21):

$$E(y) \propto \text{sinc}\left(\pi \frac{ay}{\lambda L}\right) \quad (30)$$

The relevant  $T(x)$  is given in Eq. (23). Its Fourier transform is

$$\tilde{T}(k) = \int_{-\infty}^{\infty} dx T(x) e^{ikx} = \int_{-\frac{a}{2}}^{\frac{a}{2}} dx e^{ikx} = \frac{1}{ik} \left[ e^{i \frac{ka}{2}} - e^{-i \frac{ka}{2}} \right] = \frac{2}{k} \sin\left(\frac{ka}{2}\right) = a \text{sinc}\left(\frac{ka}{2}\right) \quad (31)$$

So

$$E(y) \propto \tilde{T}\left(\pi \frac{ay}{L\lambda}\right) \quad (32)$$

exactly as expected. We can also check the diffraction grating using Eq. (26),

$$\tilde{T}(k) = \int_{-\infty}^{\infty} dx T(x) e^{ikx} = \sum_{j=1}^N \int_{-\infty}^{\infty} dx \delta(x - jd) e^{ikx} = \sum_{j=1}^N e^{ikjd} = e^{i(N+1)\Delta} \frac{\sin\left(N \frac{\Delta}{2}\right)}{\sin\left(\frac{\Delta}{2}\right)} \quad (33)$$

where  $\Delta = kd = 2\pi \frac{d}{\lambda}$ . This agrees with our result for the diffraction grating in Eq. (5).

In summary, if we put a screen over the slit with transparency  $T(x)$  then the image on the screen will be proportional to  $\tilde{T}\left(\pi \frac{ay}{L\lambda}\right)$ . Note that this holds even for monochromatic (fixed frequency) incoming light, as in a laser. Say the wavelength  $\lambda$  of light and the distance  $L$  to the screen are fixed. Then the intensity at a distance  $y$  up the screen tells us the strength of the Fourier component of the slit corresponding to a wavenumber of

$$k = 2\pi \frac{y}{L\lambda} \quad (34)$$

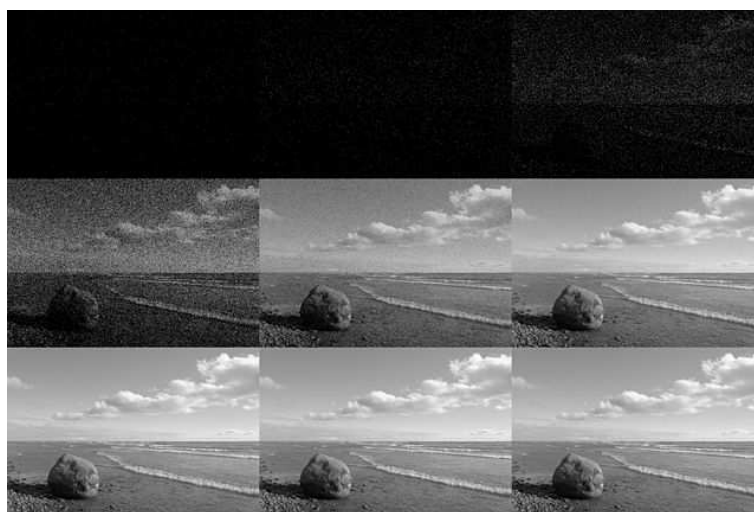
To be absolutely clear, this  $k$  is **not**  $\frac{2\pi}{\lambda}$ . The Fourier transform on the screen has nothing to do with the wavelength of the light entering the slit. The farther up the slit you look, the larger  $k$ , and hence smaller wavelength modes you are seeing. This is why the spot gets smaller as the slit gets bigger as in Figure 6.

## Lecture 20: Quantum mechanics

### 1 Motivation: particle-wave duality

This has been a course about waves. We have seen how light acts like a wave in many contexts – it refracts, diffracts, interferes, disperses, etc. Light also has a particle-like nature. Newton was the first to propose this, based mostly on observations about ray-tracing in geometrical optics. For example, that shadows had sharp edges were evidence to Newton that light was made of particles. This was in the 17th century. However, all of Newton’s observations are perfectly consistent with light as a wave. Geometric optics is simply the regime when the wavelength of light is much smaller than the relevant scales of the apparatus. Over the next 200 years, as the wave theory was developed by physicists such as Fresnel, Huygens and Maxwell, the particle-like nature of light was basically accepted as wrong.

However, light is made of particles. An easy way to see this is with a digital camera. Digital cameras are wonderfully more sensitive than film cameras. As we discussed in Lecture 17 on color, as you pump up the “ISO” on a digital camera, you are pumping up the gain on the CCD. That is, you change the mapping from intensity to the number 0-255 recorded at each pixel. At very high ISO such as 6400, you can be sensitive to single photon depositions on your sensor. Suppose we take a light source and add more and more filters in front of it to reduce the intensity and take pictures of it. The result looks like this



**Figure 1.** Going from bottom right to top left, nine photos of the same scene are taken with more and more light filtered out. The graininess of the images is mostly due to shot noise: there are a finite number of photons hitting the camera.

The type of graininess that results from the finite number of photons in light is called **shot noise** (There are other sources of graininess, such as sensor imperfections, but shot noise dominates at low light).

If you thought light was just a wave, when you lower the intensity, you would expect the photo to just get dimmer and dimmer at lower light. That there are dots in discrete places is evidence that the intensity is produced by particles hitting the sensor. If you take the same picture again at low light, the dots show up in different places. So there is an element of randomness to the locations of the photons.

We have seen that classical electromagnetic waves carry power. This is obvious from the fact that they can induce currents. By setting electrons into motion, the energy of the wave goes into kinetic energy in the electron motion. Electromagnetic waves also carry momentum. Again, this is easy to see classically: if a wave hits an electron at rest, the oscillating electric field alone would only get it moving up and down, like a cork on water when a wave passes. However, once the electron is moving, the magnetic field in the wave pushes it forward, through the  $F = eB \times v$  force law. The electron is pushed in the direction of the wavevector  $\vec{k}$ , thus the momentum in a plane wave is proportional to  $\vec{k}$ .

So if light is made of particles, those particles must be carrying momentum  $\vec{p}$ , and this  $\vec{p}$  must be proportional to the wavevector  $\vec{k}$ . We write  $\vec{p} = \hbar \vec{k}$ . By dimensional analysis, this constant  $\hbar$  must have units of  $J \cdot s$ . We write  $\hbar = \frac{h}{2\pi}$  and the value observed is

$$h = 6.26 \times 10^{19} J \cdot s \quad (1)$$

known as **Planck's constant**. Since  $\omega = c|\vec{k}|$  for electromagnetic radiation, and since light is massless,  $mc^2 = \sqrt{E^2 - c^2\vec{p}^2} = 0$  we find that  $E = c|\vec{p}| = \hbar c|\vec{k}| = \hbar\omega = h\nu$ , so energy is proportional to frequency for light.

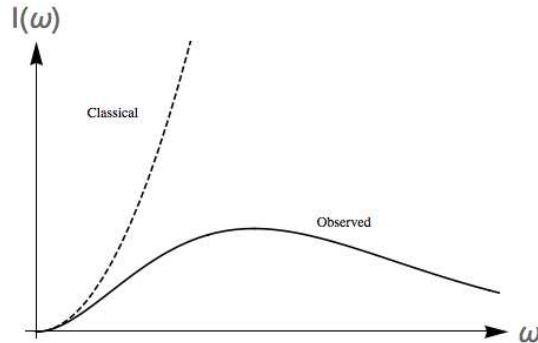
Of course, there were no digital cameras when quantum mechanics was invented. The key experimental evidence for photons came from careful studies of shining light on metals. By 1900, it was known that shining light could induce a current in a conductor. Interestingly though, the induced current depended not just on the intensity of the light, but also on its frequency. For example, in potassium, it was found that low intensity blue light would induce a current but high-intensity red light would not. The current is induced in a metal because light hits electrons, which are then freed to move around. In some situations, an even more dramatic effect occurs: the electrons get ejected from the metal. This is called the **photoelectric effect**. In 1902, Philipp Lenard found that the energy of the ejected electrons did not depend on the intensity of the light, but only on its frequency. A high frequency light source would free lots of electrons, but they would all have the same energy, independent of the intensity of the light. A low-frequency light source would never free electrons, even at very high intensity.

These observations are explained by quantum mechanics. The basic observation is that electrons have a binding energy  $E_{\text{bind}}$  in a metal. One needs a minimal amount of energy to free them. This binding energy is itself not a quantum effect. For example, we are bound to the earth; it takes a minimum amount of energy  $E_{\text{min}} = \frac{1}{2}mv_{\text{esc}}^2$  with  $v_{\text{esc}} = 11.2 \frac{\text{km}}{\text{s}}$  the escape velocity to set us free. Quantum mechanics plays a role through the identification of frequency with energy: because photons have energy  $E = h\nu$  it takes a minimum frequency  $\nu_{\text{min}} = \frac{E_{\text{bind}}}{h}$  to free the electron. After that, the electron will have kinetic energy  $E_{\text{kin}} = E_{\text{bind}} - h\nu$ . So one can make a precise prediction for how the velocity of the emitted electron depends on energy  $\frac{1}{2}mv^2 = E_{\text{bind}} - h\nu$ . This relation between  $v$  and  $\nu$  was postulated by Albert Einstein in 1905 and observed by Robert Millikan in 1914. Both were later awarded Nobel prizes (as was Lenard).

Two other important historical hints that classical theory was wrong were the classical instability of the atom and the blackbody paradox. The classical instability is the problem that if atoms are electrons orbiting nuclei, like planets orbiting the sun, then the electrons must be radiating power since they are accelerating. If you work out the numbers (as you did on the problem set), you will see that the power emitted is enormous. So this classical atomic model is untenable and must be wrong.

The blackbody paradox is that classically, one expects a hot gas to spread its energy evenly over all available normal modes. This evenness can actually be proved from statistical mechanics; the general result is known as the **equipartition theorem**. Now say you have light in a finite volume, such as the sun. It will have normal mode excitations determined by boundary conditions. Treating the volume like a box whose sides have length  $L$ , these will have wavevectors  $\vec{k} = \frac{2\pi}{L}\vec{n}$ , where  $\vec{n}$  is a vector of integers, e.g.  $\vec{n} = (4, 5, 6)$ . Since  $\omega = c|\vec{k}|$  in three dimensions there are many more modes with higher frequency than lower frequency. If these all

have equal energy, the total intensity at high frequency would grow without bound. More precisely, the classical result is that  $I(\omega) \propto \omega^4$ . This is known as the **ultraviolet catastrophe**. One can see this  $\omega^4$  behavior of the spectrum of a blackbody for low frequencies, but experimentally it was known by 1900 at least that the spectrum turned over, as shown in Figure 2.

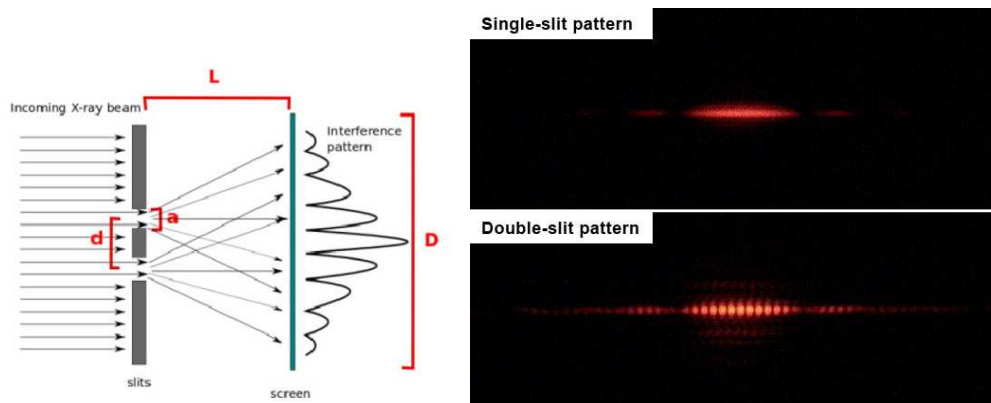


**Figure 2.** The classical prediction for the blackbody spectrum diverges at high frequency.

In 1900, Max Planck noted that the observed spectrum could be fit by the curve now known as the blackbody radiation formula (see Lecture 17). He later justified this fit by postulating that a normal mode of frequency  $\omega$  has energy  $E = \hbar\omega$  and redoing the statistical mechanical calculation of the intensity spectrum to get his curve.

## 2 Interference

So light is made of particles, but has wavelike properties. How can this be? Let's reconsider the double slit interference pattern, as in Figure 3.

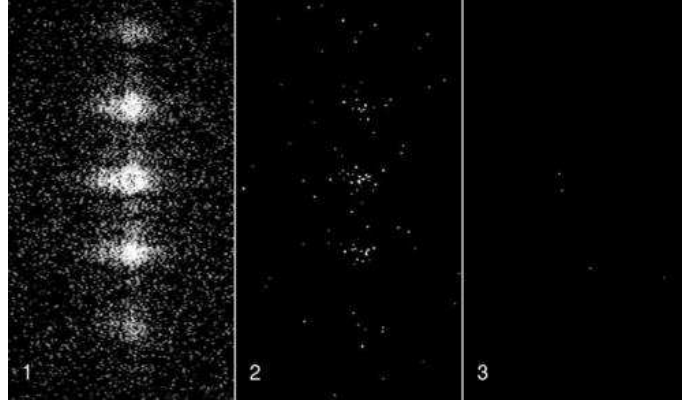


**Figure 3.** Double slit setup. Light passes through two slits and we see an interference pattern. The top image on the right is the single-slit pattern, which follows if one of the slits is covered up. The bottom image on the right is the result if both slits are open.

Recall from Lecture 19 that when a plane wave passes through two slits, there will be an interference pattern produced. More precisely, if we have two slits separated by a distance  $d$  and place a screen a distance  $L$  from the slits, the intensity pattern on the screen will be  $I(y) = 4I_0 \cos^2\left(\pi \frac{dy}{\lambda L}\right)$  where  $y$  is the distance up the screen. The oscillations with  $y$  are the small black stripes you see on the bottom right of Figure 3 (the longer scale oscillation is the  $\text{sinc}\left(\pi \frac{ay}{\lambda L}\right)$  pattern, due to the slit having finite width  $a$ ).



Since intensity is just photon counts, we expect this pattern will slowly build up if we shine the light at low intensity. That is just what happens, as can be seen in Figure 4:



**Figure 4.** Double slit interference pattern at low intensity. The pattern accumulates slowly, going from panel 3, to panel 2, to panel 1. Eventually, the granularity in the intensity pattern would wash out and it would match what we get with high intensity light over a shorter time.

This figure is equivalent to zooming in on the central region in the bottom right panel of Figure 3 and lowering the intensity. It is also true that if you cover up one of the slits, you find the single-slit interference pattern accumulating over time, as in the top right panel of Figure 3.

Now if you start to think about it a little, this result is very weird. The intensity pattern accumulates even if the intensity is very low, say one photon per minute. But if only one photon is going through at a time, what is it interfering with? Certainly not other photons. The answer is, it is interfering with itself! The photon must be going through both slits at once. Otherwise the interference pattern couldn't possibly be affected by closing off one slit.

To make this even clearer, you can do the same experiment with electrons. Exactly the same thing happens: the electrons produce an interference pattern if both slits are open, but no interference pattern if one slit is closed. For the photon case, the pattern was  $I(y) = 4I_0 \cos^2\left(\pi \frac{dy}{\lambda L}\right)$  with  $\lambda$  the photon's wavelength.

What is the “wavelength” of an electron? Well, for light momentum was  $\vec{p} = \hbar \vec{k}$  and  $|\vec{k}| = \frac{2\pi}{\lambda}$ . Thus, the wavelength of light can be extracted from its momentum as  $\lambda = \frac{2\pi}{|\vec{k}|} = \frac{h}{|\vec{p}|}$ . It is natural to expect that this equivalence will also work for electrons. This definition is called the **de Broglie wavelength**

$$\lambda \equiv \frac{h}{|\vec{p}|} \quad (2)$$

This de Broglie wavelength correctly determines the scale for the interference pattern in the electron double slit experiment.

Continuing this logic, the electron's wavevector is  $\vec{k} = \frac{1}{\hbar} \vec{p}$  and the electron's angular frequency is  $\omega = \frac{1}{\hbar} E$ , with  $E$  the electron's energy. From special relativity, we know that the electron's energy is  $E = \sqrt{m^2 c^4 + p^2 c^2}$ . Thus electrons satisfy the dispersion relation

$$\hbar \omega = \sqrt{m^2 c^4 + \hbar^2 c^2 \vec{k}^2} \quad (3)$$

This has a low-frequency cutoff (like in a plasma): the frequency of the electron satisfies  $\omega \geq \frac{m c^2}{\hbar}$ .

Another useful quantity with units of length is the **Compton wavelength** of the electron

$$\lambda_c = \frac{h}{m c} = 2.42 \times 10^{-12} m \quad (4)$$

Unlike the de Broglie wavelength, the Compton wavelength does not depend on momentum. The Compton wavelength is just a translation of the electron mass into length units and represents, roughly, the length scale where quantum mechanics becomes relevant. Calling it a wavelength is somewhat inappropriate. It should more appropriately be called the Compton length scale. The de Broglie wavelength really is the wavelength of the wave that a particle with momentum  $p$  has. So all the interference effects and so on depend on the de Broglie wavelength. The Compton wavelength does not directly characterize the scale for interference; it is merely a length scale that comes up in a lot of quantum mechanics calculations, as in Eq. (17) below, so it is handy to give it a name.

### 3 The Schrödinger equation

So light is made of particles that act like waves. Electrons are also particles which act like waves. We will now focus mostly on electrons. The reason is that to we want to focus on just one particle at a time, and with photons it's hard to do that. Due to  $E = mc^2$  and that  $m = 0$  for a photon, it is very easy to produce a lot of light. Indeed, there is no good regime where we can study single photons without having to consider multiple photons. **Quantum field theory** is quantum mechanics in regimes where energies are greater than the threshold for producing pairs of new particles with  $E = mc^2$ . For photons, this threshold energy is zero. For electrons, it is twice the rest mass,  $2m_e \approx 1.5 \times 10^{-13} \text{ J} \approx 1 \text{ MeV}$ . This is an enormous amount of energy, not generated in the laboratory until the 1930's.

What is the wave equation for the electron? To deduce it, we can work backwards from the dispersion relation. For electrons, plane waves  $\psi(x, t) = \exp(i\omega t - ikx)$  with the dispersion relation in Eq. (3) would satisfy the wave equation

$$\left[ \hbar^2 \frac{\partial^2}{\partial t^2} + \hbar^2 \vec{\nabla}^2 + m^2 c^4 \right] \psi(x, t) = 0 \quad (5)$$

with  $\psi(x, t)$  the amplitude of the electron, more commonly called its **wavefunction**. Eq. (5) is called the **Klein-Gordon equation**.

The solutions to the Klein-Gordon equation are plane waves, or linear combinations of plane waves. However, we are not usually interested in what happens to an electron pummeling through free space. Instead, we like to study systems with lots of positively and negatively charged particles, like atoms, molecules and solids. For example, around a proton, the electron is influenced by the background  $V(r) = \frac{q^2}{4\pi\epsilon_0 r}$  Coulomb potential.

So, let's adapt the Klein-Gordon equation a little so it's more practical. First of all, in most situations, the kinetic energy of the electron is much less than its rest mass (if this weren't true, we'd need quantum field theory). So it makes sense to Taylor expand the mass-energy relation:  $E = \sqrt{m^2 c^4 + p^2 c^2} = mc^2 + \frac{p^2}{2m} + \dots$ . To include the stuff around the electron, we can include the fact that the energy of the electron can also be modified in the presence of some potential  $V(x)$ , like around a proton or a nucleus. So  $E = mc^2 + \frac{p^2}{2m} + V(x) + \dots$ . In fact, since only energy differences are ever measurable, we might as well absorb the rest mass into this potential. So we now have  $E = \frac{\vec{p}^2}{2m} + V(x)$ . Note that since  $\vec{p} = m \vec{v}$  the  $\frac{\vec{p}^2}{2m} = \frac{1}{2} m \vec{v}^2$  is just the ordinary non-relativistic kinetic energy. In terms of the frequency and wavevector,  $\omega = \frac{E}{\hbar}$  and  $\vec{k} = \frac{\vec{p}}{\hbar}$ , this becomes  $\hbar\omega = \hbar^2 \frac{\vec{k}^2}{2m} + V(x)$  and the wave equation becomes

$$\left[ i\hbar \frac{\partial}{\partial t} + \frac{\hbar^2}{2m} \vec{\nabla}^2 - V(x) \right] \psi(x, t) = 0 \quad (6)$$

This is called the **Schrödinger equation**.

## 4 The hydrogen atom

As an example of the kind of things that can be calculated with quantum mechanics, consider the hydrogen atom. The hydrogen atom is just a proton and an electron. Since the proton is 2000 times heavier than the electron, we treat it as static, generating a potential  $V(r) = \frac{q^2}{4\pi\epsilon_0 r}$ . For the electron, we need to solve the Schrödinger equation in the presence of this potential.

I'm not going to solve it here, since it takes months. I'll just state the result. The general solution can be written as

$$\psi(\vec{x}, t) = \sum_{n=1}^{\infty} e^{i\omega_n t} F_n(\vec{x}) \quad (7)$$

for some frequencies  $\omega_n$  and some functions  $F_n(\vec{x})$  of  $\vec{x}$  only. These functions  $F_n(\vec{x})$  are products of Laguerre polynomials and spherical harmonics, but right now we don't care. What we do care about is that all the time dependence in these solutions has the form of a phase factor. Of course, we can always do a decomposition like this. If we forget about the  $\hbar$ 's, all we have done is find the normal modes of this wave equation.  $\omega_n$  are the normal-mode frequencies. The feature particular to the hydrogen atom is what the frequencies are. For the hydrogen atom,

$$\omega_n = -\frac{m_e c^2 \alpha^2}{2\hbar n^2} = -13.6 \text{ eV} \frac{1}{n^2}, \quad n = 1, 2, 3, \dots \quad (8)$$

where  $\alpha = \frac{e^2}{4\pi\epsilon_0} \approx \frac{1}{137}$  is the fine-structure constant. There are  $n^2$  different solutions for each value of  $\omega_n$ .

Now, we already said that energy is  $\hbar$  times angular frequency, so the hydrogen atom apparently has discrete **energy levels**, with energies  $E_n = -(13.6 \text{ eV}) \frac{1}{n^2}$ . The solution with  $n = 1$  is called the **ground state**, since it has the lowest frequency.  $n > 1$  are called **excited states**. Note that the energies are higher (negative numbers closer to zero) for higher  $n$ .  $E = 0$  corresponds to a free proton and electron. When the two are bound together, the energy is lower (or else they wouldn't bind!). This means that if a photon is to come in and knock one of the electrons out, it needs at least as much energy as the ground state. Thus photons must have a minimal frequency to knock out electrons, which agrees with observations about the photoelectric effect.

If you heat up hydrogen, you can get the electrons into the higher  $n$  energy levels. If a hydrogen atom is in an excited state, it can drop down to a lower energy state and emit a photon. Alternatively, hydrogen can absorb photons to transition between levels. By energy conservation, the photon energy must be the difference in energies of levels of the hydrogen atom. Moreover, since energy is proportional to frequency, we predict that of photons coming out of (or going into) hydrogen should have frequencies given by

$$\nu = \frac{(13.6 \text{ eV})}{h} \left[ \frac{1}{n^2} - \frac{1}{m^2} \right], \quad n = 1, 2, 3, 4, \dots \quad m = 1, 2, 3, 4, \dots \quad (9)$$

The associated wavelength for the  $m \rightarrow n$  transition is

$$\frac{1}{\lambda} = \frac{1}{91.17 \text{ nanometers}} \left[ \frac{1}{n^2} - \frac{1}{m^2} \right] \quad (10)$$

Many of these frequencies are not visible. But some are. The  $n \rightarrow 2$  transitions make up the **Balmer series**. For example,  $3 \rightarrow 2$  has wavelength 656.4 nm,  $4 \rightarrow 2$  has wavelength 486.24 nm, and so on. These agree perfectly with the observed spectrum of hydrogen, as shown in Figure 5.

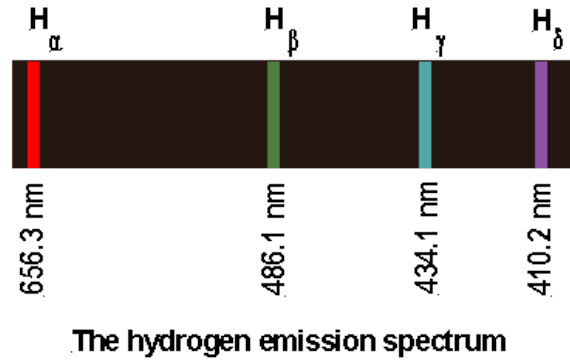


Figure 5. Visible hydrogen emission spectrum.

It is perhaps useful to note that conservation of energy translates to conservation of frequency. This isn't that different from what we have seen classically: a driven oscillator moves at the frequency of the driver. As waves pass across different media, the frequency is fixed, but the wavelength changes.

## 5 Wavepackets

The wavefunction  $\psi(x, t)$  is the amplitude. So  $|\psi(x, t)|^2$  gives the intensity. As we saw in Section 2, however, the intensity for single photons or single electrons does not show up as a very faint smooth pattern, but rather as discrete dots. The reason behind this is one of the hardest parts of quantum mechanics to get your head around. The way it works is that  $|\psi(x, t)|^2$  gives the *probability* of finding an electron at a position  $x$  and  $t$ . This probability interpretation is often taken as a postulate of quantum mechanics, although really it follows from careful consideration of entanglement with the surroundings.

Let's just take the probability interpretation as given. For  $V = 0$ , the Schrödinger equation is simply

$$\left[ i\hbar \frac{\partial}{\partial t} + \frac{\hbar^2}{2m} \vec{\nabla}^2 \right] \psi(x, t) = 0 \quad (11)$$

This is called the **free Schrödinger equation**. It has plane wave solutions

$$\psi(\vec{x}, t) = e^{i\omega t - i\vec{k} \cdot \vec{x}} \quad (12)$$

where  $\omega$  and  $\vec{k}$  satisfy the dispersion relation  $\omega = \frac{\hbar \vec{k}^2}{2m}$ . Where are the electrons in the plane wave? The probability of finding one at a point  $\vec{x}$  at time  $t$  is  $|\psi(\vec{x}, t)|^2 = 1$ . Thus they are everywhere! Clearly plane waves are not useful for describing electrons in any regime where they act particle-like.

There are lots of other solutions to the free Schrödinger equation besides the plane wave ones. For example, we can construct wave-packets. As in Lecture 11, we can write at  $t=0$

$$\psi(\vec{x}) = e^{-\frac{1}{2} \left( \frac{\vec{x} - \vec{x}_0}{\sigma_x} \right)^2} e^{i\vec{k}_0 \cdot \vec{x}} \quad (13)$$

this is a packet centered at  $\vec{x} = \vec{x}_0$  with width  $\sigma_x$  and a carrier wavelength  $\vec{k}_0$ . The Fourier transform is

$$\tilde{\psi}(k) = \frac{\sigma_x}{(2\pi)^{3/2}} e^{-\frac{1}{2} \left( \frac{k - k_0}{\sigma_k} \right)^2} e^{-i x_0 (k - k_0)} \quad (14)$$

with  $\sigma_k = \frac{1}{\sigma_x}$ . This is a wavepacket centered on  $\vec{k}_0$ . But  $\vec{p} = \hbar \vec{k}$ . So the momentum distribution is Gaussian too.

In fact, just as  $|\psi(\vec{x})|^2$  gives the probability that we find the electron at the position  $\vec{x}$  at time  $t$ , it is also true that  $|\tilde{\psi}(\vec{k})|^2$  gives the probability that we find the electron with momentum  $\vec{p} = \hbar\vec{k}$ . Thus the most likely value for momentum is  $\vec{p}_0 = \hbar\vec{k}_0$  and the width of the momentum distribution is  $\sigma_p = \hbar\sigma_k = \frac{\hbar}{\sigma_x}$ . In particular, we find

$$\sigma_x\sigma_p = \hbar \quad (15)$$

This is a special case of the Heisenberg uncertainty principle

$$(\Delta x)(\Delta p) \geq \hbar \quad (16)$$

(sometimes you see this relation with  $\frac{\hbar}{2}$  on the left; the factor of 2 has to do with precisely how  $\Delta x$  and  $\Delta p$  are defined). This says that for any wavefunction, the uncertainty on position  $\Delta x$  times the uncertainty of momentum  $\Delta p$  must be larger than or equal to  $\hbar$ . Gaussian wavepackets saturate this inequality. Note that  $\hbar = 1.04 \times 10^{-34} \text{ J}\cdot\text{s}$  is very small. So we can actually know the position and momentum of the electron quite well at the same time – to less than  $10^{-17}$  each in SI units. Quantum mechanics only shows this weird uncertainty if you are at very short distances.

Next, note that the carrier wavenumber  $k_0$  translates to the most likely momentum value. In position space, we see from Eq. (13) that the phase of the wavefunction gives the momentum. Thus a wavefunction at a given time has both position information (in the amplitude) and momentum information (in the phase).

What happens to these electrons with time? Since the dispersion relation  $\omega = \frac{\hbar k^2}{2m}$  is not linear, there is dispersion. That is, wavepackets broaden over time. As we saw in Lecture 11, the second derivative of the frequency with respect to wavenumber determines the broadening rate. In this case,  $\Gamma = \omega''(k_0) = \frac{\hbar}{m}$  is independent of  $k_0$ . In terms of the Compton wavelength in Eq. (4),  $\Gamma = \frac{\lambda_c}{2\pi c}$ . Then the width grows (see Eq. (36) of Lecture 11) as

$$\sigma(t) = \sigma_x \sqrt{1 + \left( \frac{\lambda_c c t}{2\pi \sigma_x^2} \right)^2} \approx \frac{\lambda_c}{2\pi \sigma_x} c t \quad (17)$$

where the last form holds at late times. So no matter how localized the electron is at  $t = 0$ , eventually, we will lose track of it. For example, suppose we start out knowing the electron down to  $\sigma_x = \lambda_c = 10^{-12} \text{ m}$ . Then at late time  $\sigma(t) = \frac{1}{2\pi} c t$ . Thus the electron wavepacket is broadening at the speed of light. Moreover, the smaller  $\sigma_x$  we start out with, the faster the packet broadens. We can take  $\sigma_x$  large to slow the broadening, but for large  $\sigma_x$ , we don't know where the electron is to begin with. Electrons are very hard to pin down indeed!

# Lecture 21: The Doppler effect

## 1 Moving sources

We'd like to understand what happens when waves are produced from a moving source. Let's say we have a source emitting sound with the frequency  $\nu$ . In this case, the maxima of the amplitude of the wave produced occur at intervals of the period  $T = \frac{1}{\nu}$ . If the source is at rest, an observer would receive these maxima spaced by  $T$ . If we draw the waves, the maxima are separated by a wavelength  $\lambda = Tc_s$ , with  $c_s$  the speed of sound.

Now, say the source is moving at velocity  $v_s$ . After the source emits one maximum, it moves a distance  $v_s T$  towards the observer before it emits the next maximum. Thus the two successive maxima will be closer than  $\lambda$  apart. In fact, they will be  $\lambda_{\text{ahead}} = (c_s - v_s)T$  apart. The second maximum will arrive in less than  $T$  from the first blip. It will arrive with period

$$T_{\text{ahead}} = \frac{\lambda_{\text{ahead}}}{c_s} = \left( \frac{c_s - v_s}{c_s} \right) T \quad (1)$$

The frequency of the blips/maxima directly ahead of the siren is thus

$$\nu_{\text{ahead}} = \frac{1}{T_{\text{ahead}}} = \left( \frac{c_s}{c_s - v_s} \right) \frac{1}{T} = \left( \frac{c_s}{c_s - v_s} \right) \nu. \quad (2)$$

In other words, if the source is traveling directly towards us, the frequency we hear is shifted upwards by a factor of  $\frac{c_s}{c_s - v_s}$ .

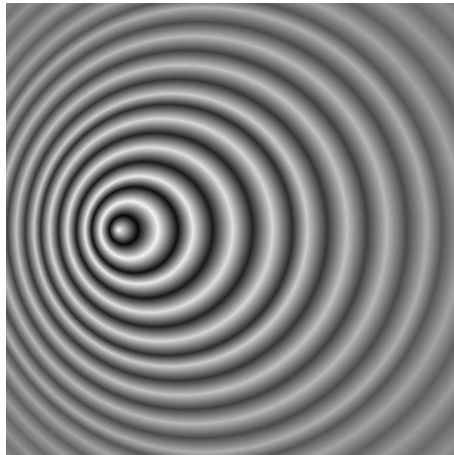
We can do a similar calculation for the case in which the source is traveling directly away from us with velocity  $v$ . In this case, in between pulses, the source travels a distance  $T$  and the old pulse travels outwards by a distance  $c_s T$ . The physical spacing between maxima is therefore  $\lambda_{\text{behind}} = (c_s + v)T$ . The frequency as perceived by an observer behind the siren is thus

$$\nu_{\text{behind}} = \left( \frac{c_s}{c_s + v_s} \right) \nu \quad (3)$$

i.e. lower than for a stationary source.

In other words, the frequency goes up when the source is approaching us, and goes down when it is traveling away from us.

The sound waves produced by the moving source are depicted in Fig. 1.



**Figure 1.** Sound waves emitted by a source moving to the left. The peaks in pressure ahead of the source are more closely spaced than the peaks behind the source.

We can summarize these two results by saying that for a stationary observer directly ahead or behind the moving source,

$$\nu' = \left( \frac{c_s}{c_s + v_s} \right) \nu \quad (4)$$

where  $v_s$  is positive if the source is moving away from the observer, and negative if the source is moving towards the observer.

It is not hard to include also the case when the observer is in motion. Say the source is moving away from the observer. Then the spacing between adjacent peaks are spread out as before  $\lambda_{\text{behind}} = (c_s + v_s)T = \frac{c_s + v_s}{\nu}$ . If the observer is moving towards the source with velocity  $v_r$ , then she passes these peaks faster, at the rate

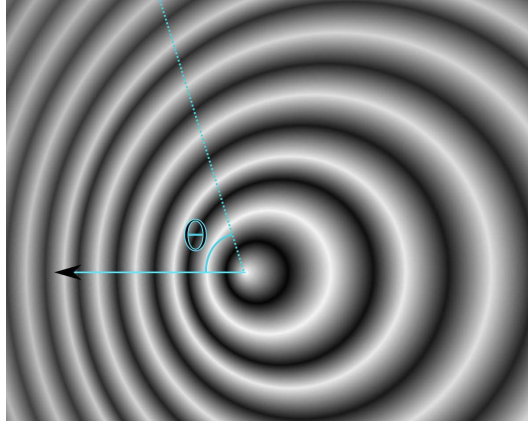
$$\nu' = \frac{c_s + v_r}{\lambda_{\text{ahead}}} = \left( \frac{c_s + v_r}{c_s + v_s} \right) \nu \quad (5)$$

In the cases where the direction of either the source or observer is flipped, the sign of  $v_s$  or  $v_r$  flips, but this equation still holds. Note that if the observer and source are moving at the same speed in the same direction, no frequency change is detected.

This type of change in frequency due to motion is called the **Doppler effect**.

## 2 Motion at an angle

What happens if the source is not moving directly towards or away from the receiver? Say the source is moving at an angle  $\theta$  with respect to a stationary receiver, as shown in Fig. 2.



**Figure 2.** Source moving at angle  $\theta$  relative to the axis connecting the source and the receiver. In this picture, the observer is along the angled blue line, say on the top of the image.

It's easiest to see what the observed frequency is in this case by looking at the picture. Now the maxima are spaced  $\lambda_\theta = (c_s - v_s \cos \theta)$ . Check that for  $\theta = 0$ , this reduces to the ahead case, for  $\theta = \pi$  it reduces to the behind case and for  $\theta = \frac{\pi}{2}$ , where the observer is orthogonal to the direction, there is no change. Following the same logic as before, the frequency of the sound the receiver hears is given by

$$\nu' = \left( \frac{c_s}{c_s - v_s \cos \theta} \right) \nu \quad (6)$$

This angular dependence explains why the siren of a police car or ambulance sounds the way it does when it passes you. While approaching at a distance, the car is basically going towards you and the frequency is increased. When it's going away, the frequency is lowered. As the car passes us, the angle transitions pretty quickly, and the sound transitions from high to low as the car goes through the intermediate angles.

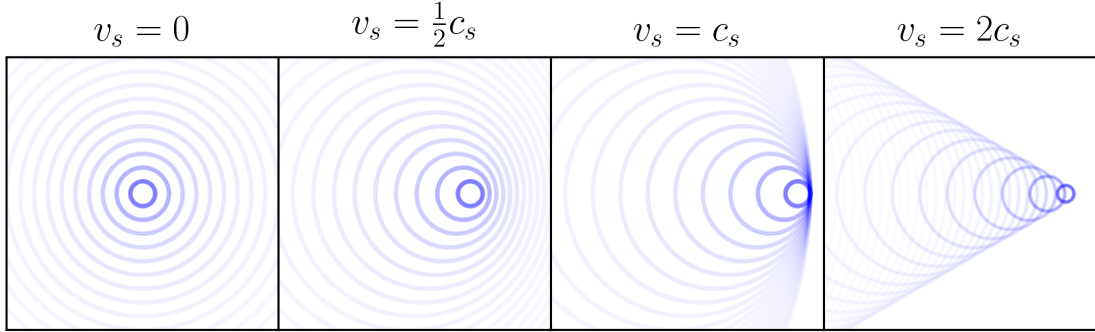
Note that the angle  $\theta$  denotes the angle to the velocity vector of the police car when the source was emitted, when the sound is received. To see this, imagine that we are very far away from the source, and it took a couple of days for the sound to get to us. Clearly what the source did during those days, such as the position it ended up at, is irrelevant to the frequency we hear.



Thus the only possibly relevant piece of information is where the source was and how it was moving when it emitted the sound. So the fastest frequency change is not when the car is moving perpendicular to your line of sight, but slightly before the car gets to that point.

### 3 Source moving at or faster than sound

What happens if the source approaches the speed of sound, or surpasses it? Again, the result is easiest to understand with pictures. Here are some wavefronts as the speed of sound is approached and surpassed.



**Figure 3.** Pulses given off by a source moving to the right at various speeds. The older pulses are depicted as dimmer. All the circles are growing with time.

The first two panels we’ve already discussed: stationary source and a source moving **subsonically** (less than the speed of sound).

As the source approaches the speed of sound, you can see from the picture that the wavefronts in the forward direction start bunching up. When the speed of sound is hit, they all come together. Remember for sound these are wavefronts of pressure. The large increase in pressure from the accumulation of maxima is followed by a large decrease in pressure from the minima. This rapid change in pressure is very loud. It is called a **sonic boom**. You heard it as the cracking of a whip.

When the source goes **supersonically** (faster than the speed of sound), the sonic boom goes from a straight wavefront perpendicular to the motion of the source to being bent backwards in a cone. The sonic boom is now at an angle to the motion of the source, but it is still there. So, to be clear, the sonic boom is a sign that the speed of sound has been surpassed; it is not something that happens only exactly when  $v = c_s$ .

### 4 Redshift

The picture with the wavefronts works just as well for light as for sound, at least for the case  $v_s \ll c$ . Of course, the source can never go faster than the speed of light, due to special relativity, so these equations need some modification.

The nice thing about relativity is that you can pick whatever inertial reference frame you want. So let’s work in the rest frame of the source, and call  $v$  the velocity of the observer. Since the source is at rest, the wave crests are spaced  $\lambda = \frac{c}{\nu}$  apart. If the moving observer is heading away from the source, she passes the crests at a rate  $\nu_{\text{move}} = \frac{c-v}{\lambda} = \frac{c-v}{c}\nu$ , as in Eq. (5). However, since the observer is moving very fast there is also a time dilation effect. Time is slower for her, so she really sees successive crests at the lower frequency

$$\nu' = \frac{\nu_{\text{move}}}{\sqrt{1 - \frac{v^2}{c^2}}} = \sqrt{\frac{1 - \frac{v}{c}}{1 + \frac{v}{c}}} \nu \quad (7)$$



More generally,  $v$  is the relative velocity of the source and the observer:  $v = v_s - v_r$  in the notation from before. If  $v > 0$  the two are moving away from each other and  $v < 0$  if they are moving towards each other.

Note that for small  $v \ll c$ , we can Taylor expand Eq. (7) giving  $\nu' = \nu(1 + \frac{v_r - v_s}{c} + \dots)$ . Taylor expanding Eq. (5) gives  $\nu' = \nu(1 + \frac{v_r - v_s}{c_s} + \dots)$ . So in the small velocity limit, this relativistic analysis reduces to our previous results.

In astronomy, it is useful to define the **redshift** of a signal as

$$z \equiv \frac{\Delta \lambda}{\lambda} = -\frac{\Delta \nu}{\nu} = \frac{\nu - \nu'}{\nu} = 1 - \frac{\nu'}{\nu}. \quad (8)$$

Negative redshift is referred to as **blueshift**. These names are used because a signal which is redshifted is shifted to longer wavelengths, and red is longest-wavelength light visible to humans. Blueshifted signals are shifted to shorter wavelengths, and blue is the shortest-wavelength light visible to humans. Astronomers commonly use “red” as a synonym for long-wavelength, and “blue” as a synonym for short-wavelength. For objects at low velocity compared to the speed of light, plugging Eq. (8) into the above formula for redshift yields

$$z \approx \frac{v}{c} \quad (9)$$

where  $v$  is positive if the source is moving away from the receiver. Sources receding from the observer thus appear redshifted, while sources moving towards the observer appear blueshifted. This is an incredibly useful fact in astrophysics, as it allows us to measure the velocity with which distant sources of light are receding from or approaching us.

So what do we find? Everywhere we look, the objects are redshifted. We can sometimes measure the distance to an object by how bright it is, or using parallax. When we do this we find an essentially linear relationship between distance and redshift.

$$z = \frac{1}{c} H_0 r \quad (10)$$

where  $H_0$  is a constant, called **Hubble’s constant**, named after Edwin Hubble who first observed this relation. The value is  $H_0 = 20 \frac{\text{km}}{\text{s}} / 10^6 \text{ light years} = (13.8 \text{ billion years})^{-1}$ . This means that stars, galaxies, and everything else is moving away from us, and the farther stuff is moving faster. That’s exactly what would happen in a big explosion: the stuff farther away is moving faster (that’s how it got to be farther away), and distance is proportional to velocity. So the proportionality of  $z$  to distance gives direct evidence that there was once a big bang. Extrapolating the velocities back to when they all met, it seems that the big bang was around 13.8 billion years ago. That’s not really very long, considering the earth is only 5 billion years old. It’s also pretty neat that when we look farther away in the universe (at larger redshift, we are looking back in time).

## 5 Cosmic microwave background

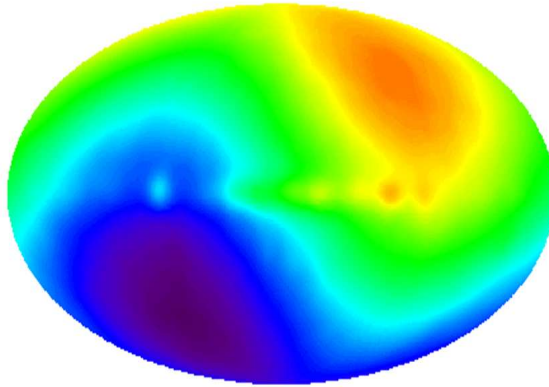
When you go closer to the big bang, you have to fit all the energy of the universe in a smaller volume, so it has to be hotter. Working out the relevant physics, the result is that the temperature  $T$  in Kelvin is related to time  $t$  in seconds by

$$T \approx \frac{10^{10}}{\sqrt{t}} \quad (11)$$

At an early enough time, the universe was so hot that the thermal energy was enough to ionize atoms. At this time, the universe was essentially a plasma. In a plasma, photons cannot get very far without hitting a free electron or proton, so the universe is opaque. The ionization energy of atoms is around the Rydberg constant 13 eV, which corresponds to a temperature of  $T = 160,000 \text{ K}$ . By the time the universe has cooled to around 3000 K, the rate for electrons to be captured by protons exceeds the rate for photons to knock them out. This temperature is called the **reionization temperature**. Using Eq. (11) this happened at a time  $t = 380,000 \text{ yrs}$ . This time corresponds to a redshift of  $z = 1100$ , which is called the **surface of last scattering**. We can’t see anything farther than this, since the universe was opaque.

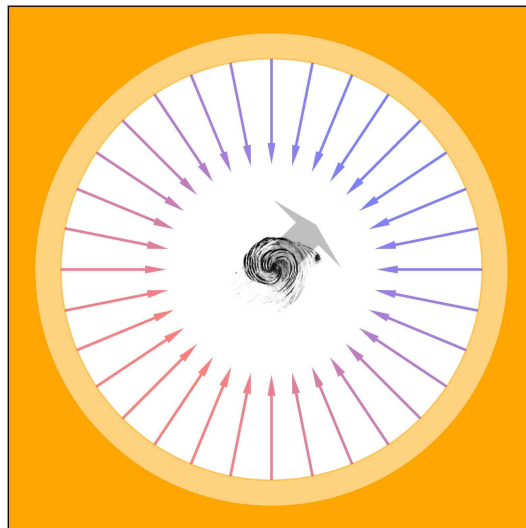
What does the surface of last scattering look like? Well, the last photons that scattered had a blackbody spectrum at  $T \approx 3000K$ . But these photons all get redshifted since they are very far away. The effect is that the spectrum looks like a blackbody at  $\frac{T}{z} \approx 3K$ . People sometimes say that the photons have cooled as the universe expanded, but I think it's more accurate to say that they redshifted. The net result is that there is a thermal bath of photons at  $3K$ . This bath was first observed by 1964 by Penzias and Wilson (somewhat by accident – they thought it was reducible noise in another experiment they were doing, but couldn't get rid of it!). The peak photon wavelength at  $3K$  is in the microwave region, and this is called the **cosmic microwave background (CMB)**.

In 1995, the Cobe satellite found that the CMB wasn't perfectly thermal. It has temperature fluctuations around the blackbody spectrum at around one part in  $10^5$ . Here is a picture of their observations



**Figure 4.** Temperature anisotropies in the Cosmic Microwave Background, as seen by NASA's COBE mission. This oval represents all the directions in the universe, much like a map of the earth can be projected on a similar ellipse. Notice that one end of the sky is hotter than the opposite end.

One thing we see from this is that the universe is slightly hotter in one direction than the other. This means that we are moving with respect to the CMB, as explained in this figure



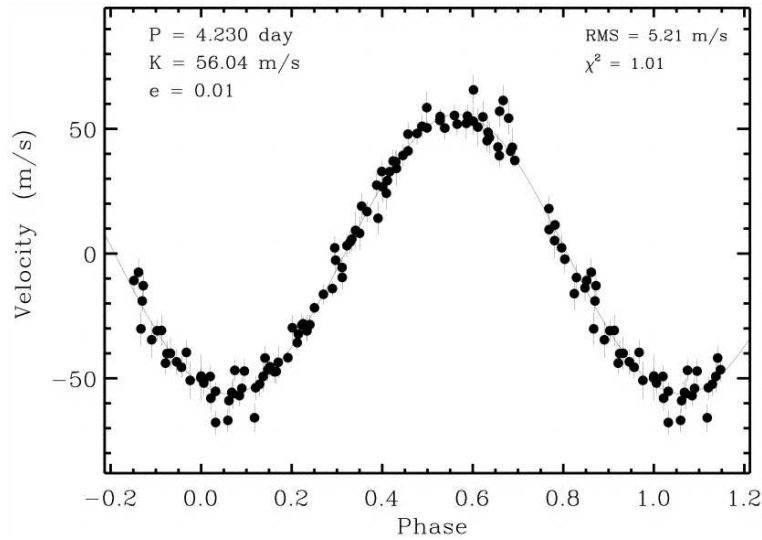
**Figure 5.** A schematic representation of the dipole anisotropy in the Cosmic Microwave Background. CMB photons emitted from the surface of last scattering approach the Milky Way. The relative velocity of the Milky Way to the surface of last scattering determines the redshift of the photons an observer in the Galaxy perceives. Photons traveling inwards along the Milky Way's velocity vector appear blueshifted, while photons traveling inwards from the opposite point on the sky appear redshifted.

More precision studies of the CMB have told us (and continue to tell us) a lot about our universe: the amount of dark matter, dark energy, the total number of particles.

## 6 Doppler spectroscopy

Gasses in stars emit and absorb light at specific frequencies, corresponding to transition energies between electron orbitals. The spectrum of a star therefore contains many narrow peaks and troughs, called spectral lines. If the star is moving relative to our telescope, then its entire spectrum is either redshifted or blueshifted. The redshift of a star can therefore be determined very precisely by measuring the offset of certain spectral lines from their standard frequencies. This method, called **Doppler spectroscopy**, allows a precise measurement of the velocity of a star along the line of sight. This velocity is called the radial velocity of the star.

Here is an example of the motion of a star (51 Pegasi) determined from Doppler spectroscopy:



**Figure 6.** Radial velocity of the star 51 Pegasi, inferred by Doppler shifting of emission and absorption lines in the star's spectrum. Detected in 1995, 51 Pegasi b was the first planet found orbiting another star. This plot is from Marcy *et al.* (1997).

We see that there is a periodic oscillation of the velocity. This can be explained if there is a planet orbiting the star. For example, imagine some alien is looking at the earth from far away. The average velocity of the earth is the average velocity of our solar system. But when the earth is on one side of the sun, it moves away from the alien and on the other side, it moves towards the alien. Thus the earth's velocity as viewed from far away has wiggles due to the annual orbit. If the earth were much bigger, the sun would also have detectable wiggles. So what we are seeing in the 51 Pegasi spectrum above are wiggles due to a very large planet in a tight orbit with that star.

Doppler spectroscopy can be extremely powerful when the spectral features one is interacting with are very narrow (or high  $Q$ ). One of the most high  $Q$  spectroscopic systems are Mossbauer transitions in radioactive nuclei ( $^{56}\text{Fe}$  for example). A Mossbauer transition has a  $Q > 10^{11}$ . Mossbauer transitions are such high  $Q$  due to an unusual effect where nuclei locked in a crystal must all recoil together in response to the emission of a photon. Additionally, nuclear energy levels are highly protected from perturbing electric fields, mostly because they are well shielded by all of the orbiting electrons. Thus, every nucleus in a solid block of iron will have a nearly identical resonant frequency.

Because the  $Q$  of a Mossbauer transition is so high, they can be used for extremely high precision Doppler spectroscopy measurements. In general, these experiments work as follows. One Mossbauer emitter is placed at rest and a second absorbing piece of Mossbauer material on a moving stage. By measuring the absorption of photons emitted by the first stationary piece in the second moving piece, one can measure velocities as low as  $10^{-6} \frac{m}{s}$ .

One of the most clever experiments ever performed using this effect was the measurement of the gravitational redshift. A gravitational redshift happens from a combination of quantum mechanics and relativity. Photons have energies  $E = h\nu$ . Since  $E = mc^2$ , this energy acts just like mass from the point of view of gravity. Thus if we shine a photon towards the center of the earth, it must gain energy, just like dropped mass would. This energy can only be stored in the photon's frequency. So this frequency must change. In 1959 Pound and Rebka placed a iron 57 Mossbauer emitter in the basement of Jefferson labs and another on the fourth floor. Working out the numbers, the expected energy shift is just  $\frac{dE}{E} < 10^{-14}$  of the photon. For the 14 keV photon, this corresponds to a Doppler shift of the second absorber moving at  $10^{-6} \frac{m}{s}$ . They moved the second iron piece at the speed required to have perfect absorption (on resonance) of the emitted photons. With this carefully designed experiment, they were eventually able to observe the gravitational redshift to within 1% of the predicted value.

In gases, the  $Q$  values for spectral lines are not nearly as high as in Mossbauer materials. The reason that emission peaks, say from Hydrogen, are not infinitely narrow is due to **Doppler broadening**. If the gas is at a temperature  $T$ , the probability  $M(v)dv$  of finding an atom with velocity between  $v$  and  $v + dv$  is given by the Maxwell-Boltzmann distribution

$$M(v) = \sqrt{\frac{m}{2\pi k_B T}} e^{-\frac{mv^2}{2k_B T}} \quad (12)$$

where  $m$  is the mass of the Hydrogen molecules  $k_B$  is the Boltzmann constant. At rest, a Hydrogen atom would have an emission at  $\lambda_0 = 656\text{nm}$  or equivalently at a frequency of  $\nu_0 = 4.5 \times 10^{17}\text{Hz}$ . If the molecules are moving towards the observer, the frequency will go up due to the Doppler effect, and if they are moving away, it will go down. The probability of finding the emission at  $\nu$  is then  $P(\nu)d\nu$  where

$$P(\nu) = \frac{c}{\nu_0} M\left(c\left(\frac{\nu}{\nu_0} - 1\right)\right) = \sqrt{\frac{mc^2}{2\pi k_B T \nu_0^2}} \exp\left[-\frac{mc^2(\nu - \nu_0)^2}{2k_B T \nu_0^2}\right] \quad (13)$$

This is a Gaussian distribution with width

$$\frac{\Delta\nu}{\nu_0} = \frac{\Delta\lambda}{\lambda} = \sqrt{\frac{8k_B T}{mc^2}} = \sqrt{\frac{8}{\gamma}} \frac{c_s}{c} \quad (14)$$

where  $c_s = \sqrt{\gamma \frac{k_B T}{m}}$  is the speed of sound. For example, at room temperature, where  $c_s = 343 \frac{m}{s}$  and  $\gamma = \frac{7}{5}$  for Hydrogen gas,  $\frac{\Delta\lambda}{\lambda} = 2 \times 10^{-6}$  and the emission line at  $\lambda = 656\text{nm}$  will have a width of  $\Delta\lambda = 0.002\text{nm}$ .