# Jet Superstructure
# and Multivariate Studies

## Matthew Schwartz

### Harvard University

LHC Physics Day, CERN
February 4, 2011

# INTRODUCTION

A lot of recent work on jet substructure and some on jet superstructure

- Masses, angularities, filtering/trimming/pruning, subjetiness, planar flow, ...
- Interesting theory questions
  - What is optimal?
  - Can we trust monte carlos?
  - Can we compute them more accurately in QCD?
- Variables are useful, but highly correlated
  - e.g. jet mass and jet $p_T$ are closely related
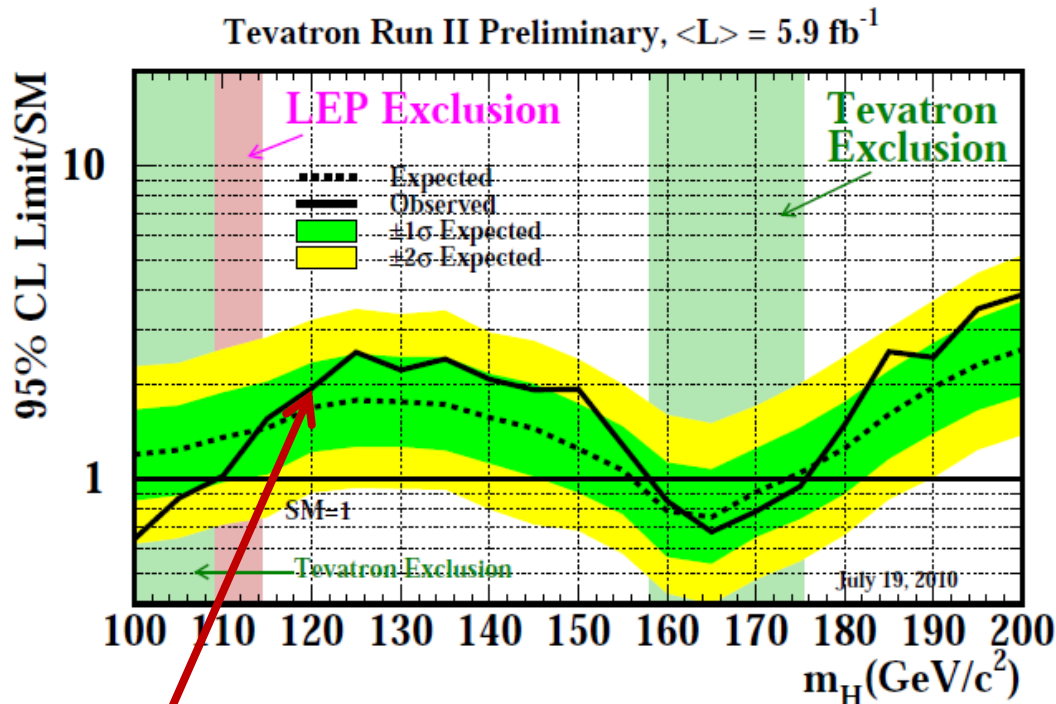
Why do experimentalists use multivariate methods: Neural Networks (NN), Boosted Decision Trees (BDT), etc, but theorists do not?

- Experimentalists want to see things early -- every little bit helps
- To theorists, the difference between 10 fb$^{-1}$ and 100 fb$^{-1}$ is 0
- NNs and BDTs are complicated – theorists are scared of black boxes

To properly appreciate jets, we must get used to studying variables *and* their correlations

# HOW DO WE FIND A LIGHT HIGGS?

## Tevatron



Tevatron Run II Preliminary, $\langle L \rangle = 5.9$ fb$^{-1}$

LEP Exclusion

Tevatron Exclusion

Expected
Observed
$\pm 1\sigma$ Expected
$\pm 2\sigma$ Expected

95% CL Limit/SM

10

1

SM=1

Tevatron Exclusion

July 19, 2010

100 110 120 130 140 150 160 170 180 190 200
$m_H$(GeV/c$^2$)

- Need a factor of 2 improvement in significance for $m_H$=120
- Double statistics gives √2
- Where will the other √2 come from?

## LHC

- Important search channel is
$$pp \rightarrow W/Z + H$$
$$H \rightarrow bb$$

- Abandoned by ATLAS and CMS
too much background

- Recently high $P_T$ W/Z + H revived,
  - Requires $P_T$ > 200
  - Lose 95% of signal

How good can we do
in W/Z + (H → bb)?

# FOCUS ON $pp \rightarrow HZ \rightarrow b\bar{b}l^+l^-$

CDF note 10235 (summer 2010)

| | |
|---|---|
| $ZH$ | 0.7 |
| $t\bar{t}$ | 9.9 |
| $WW$ | 0.02 |
| $WZ$ | 0.1 |
| $ZZ$ | 3.6 |
| $Z \rightarrow \ell\ell + b\bar{b}$ | 22.1 |
| $Z \rightarrow \ell\ell + c\bar{c}$ | 2.4 |
| $Z \rightarrow \ell\ell + l.f.$ | 1.2 |
| fakes | 0.9 |
| Total Bkg | 40.3 |

Dominant background
is the irreducible one

CDF employs multivariate approach

Inputs to the neural net are
- Missing transverse energy
- Dijet mass
- tt matrix element output ⎤ Parton-level
- ZH matrix element output ⎦ kinematics
- Sum of leading jet Pt's
- number of jets

Questions:
- Are there smarter more comprehensive inputs?
- Can we trust the multivariate approach?

# ONE THING THEY IGNORE: COLOR

Signal

Background

$$H \rightarrow b\bar{b}$$

$$q\bar{q} \rightarrow Zb\bar{b}$$

$$gg \rightarrow Zb\bar{b}$$

# HOW DO THEY SHOW UP?

Monte Carlo simulation
- • Color coherence (angular ordering, e.g. Herwig)
- • Color string showers in its rest frame (pt ordering, e.g. Pythia)
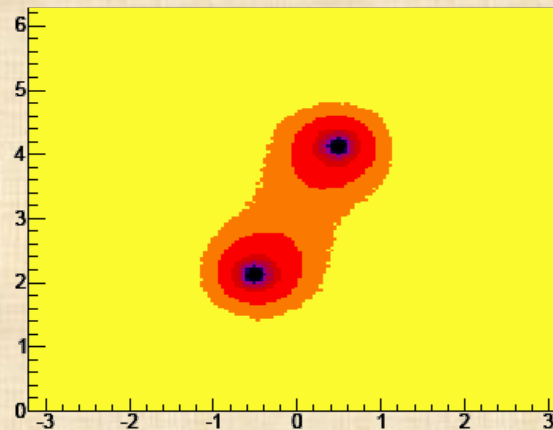    - • Boost → string showers in string-momentum direction

Shower same event
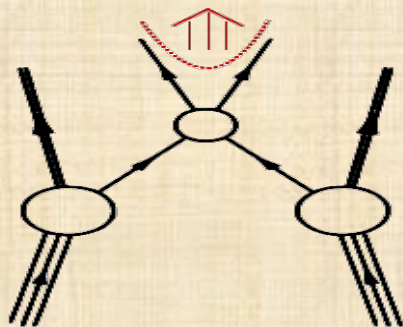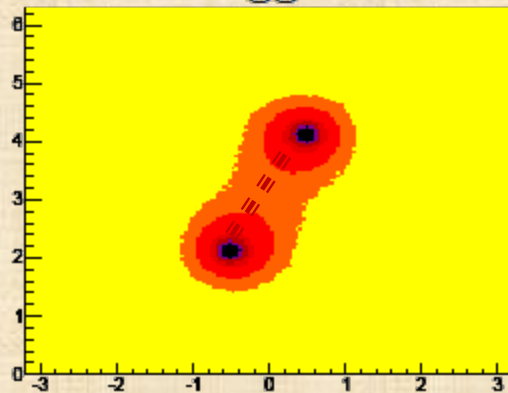*millions* of times

Higgs:

$$\Delta\eta_{b\bar{b}} = 1$$
$$\Delta\phi_{b\bar{b}} = 2$$

Add up $E_T$ in
each cell:

# SIGNAL VS BACKGROUND
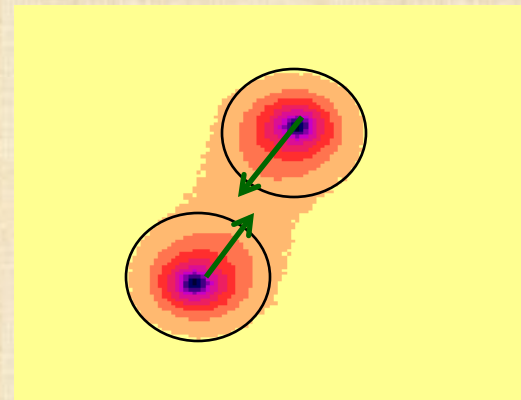
# HOW CAN WE USE IT?

Higgs:



$q\bar{q}$



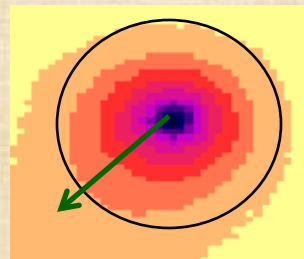Baysean **probability** that
each bit of radiation is **signal**



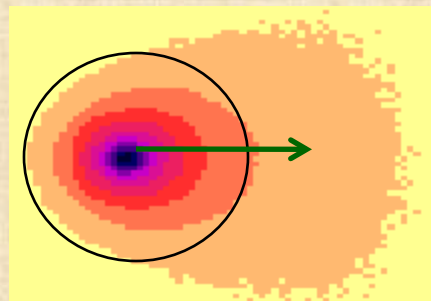- Most useful radiation is
  **R = 0.5 – 1.5** away
- Pattern depends strongly on kinematics
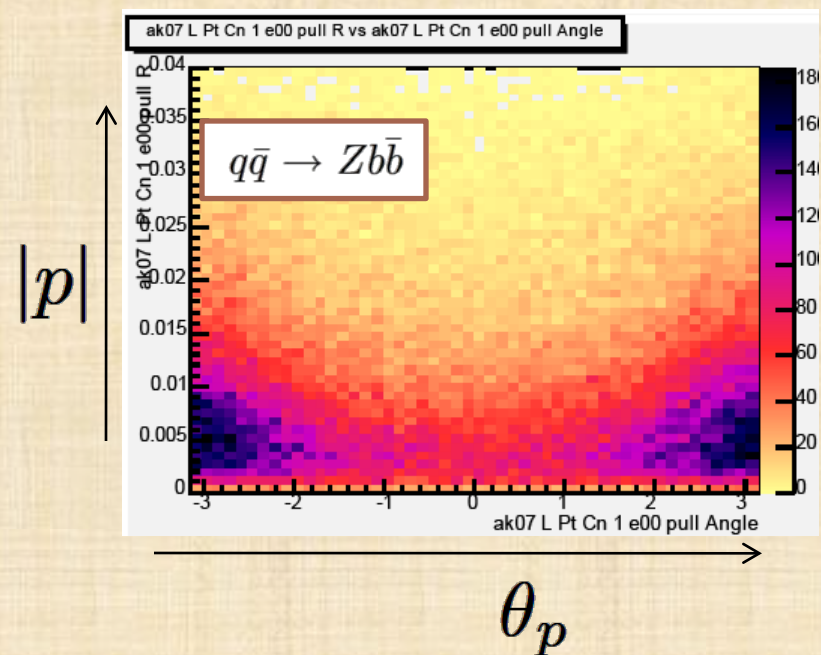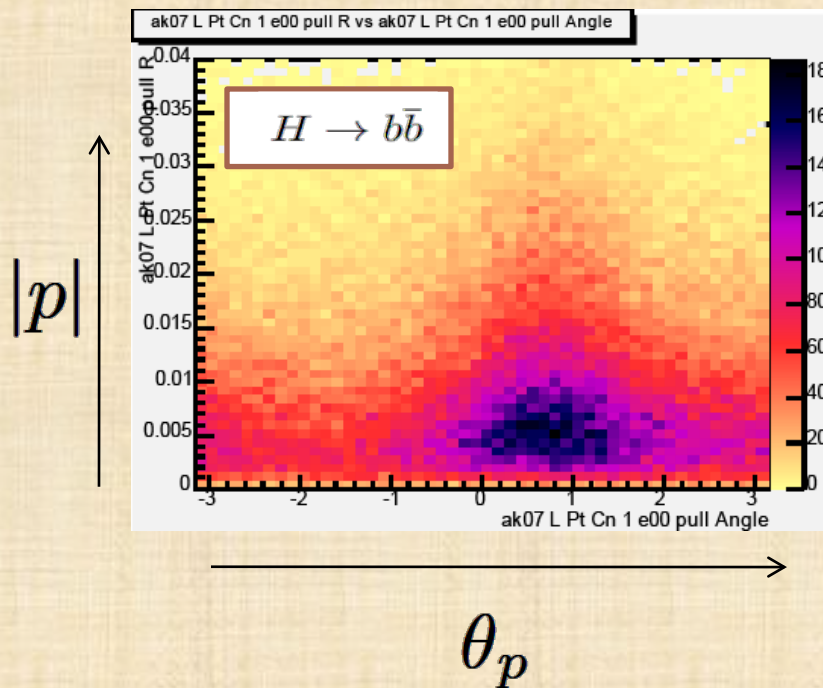- Can we find a *simpler* or more *universal* discriminant?

# PULL



- Find **jets** (e.g. anti-$k_T$)
- Construct pull vector (~ dipole moment) on radiation in **jet**

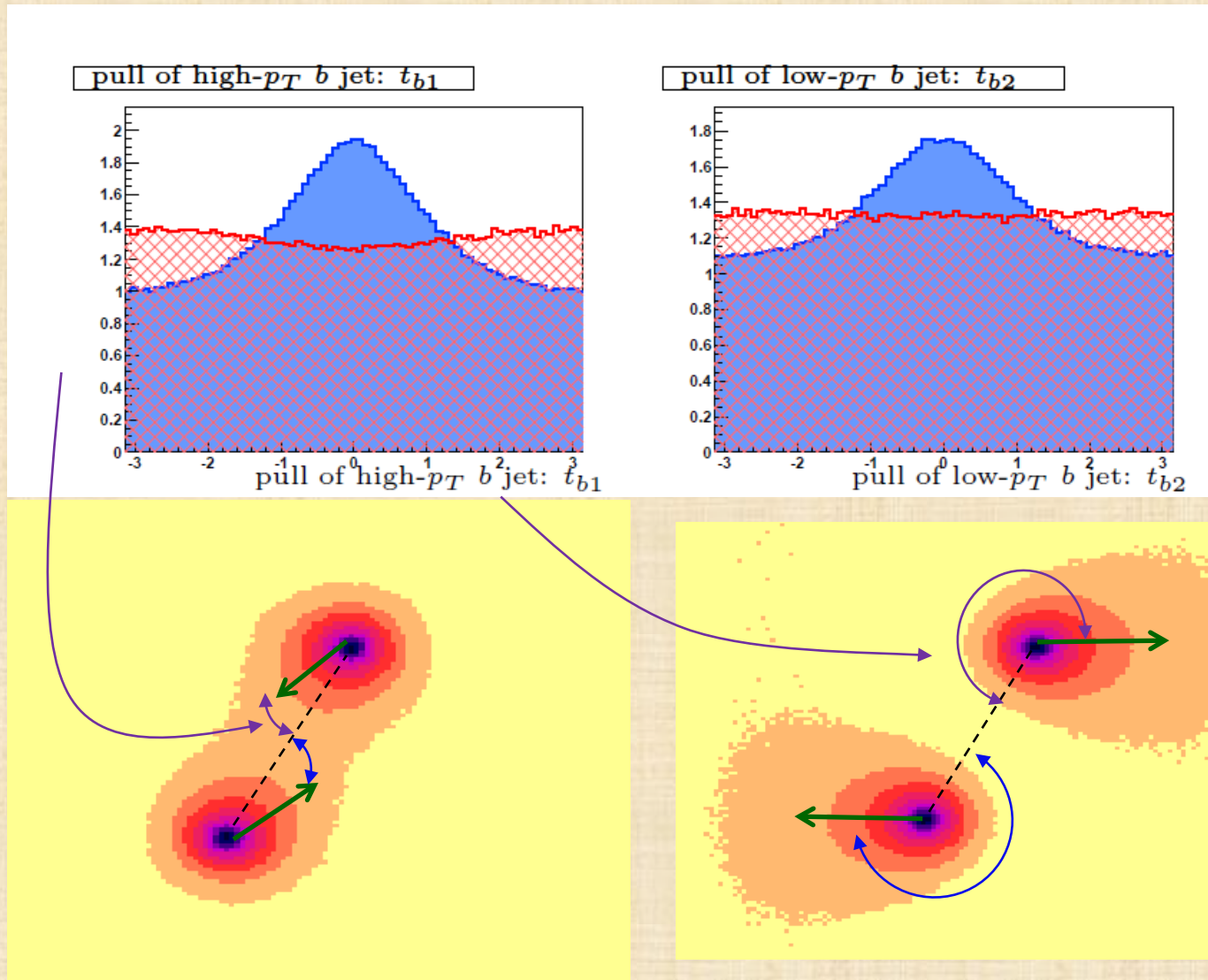$$\vec{p} = \sum_i \frac{E_T^i \, |r_i|}{E_T^{jet}} \, \vec{r}_i$$

# PULL VECTOR IN RADIAL COORDS

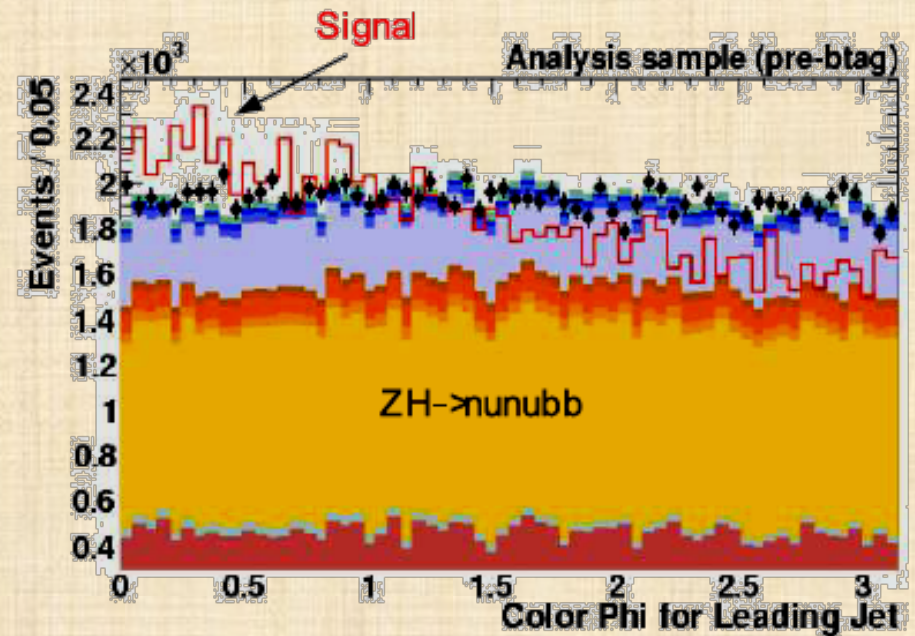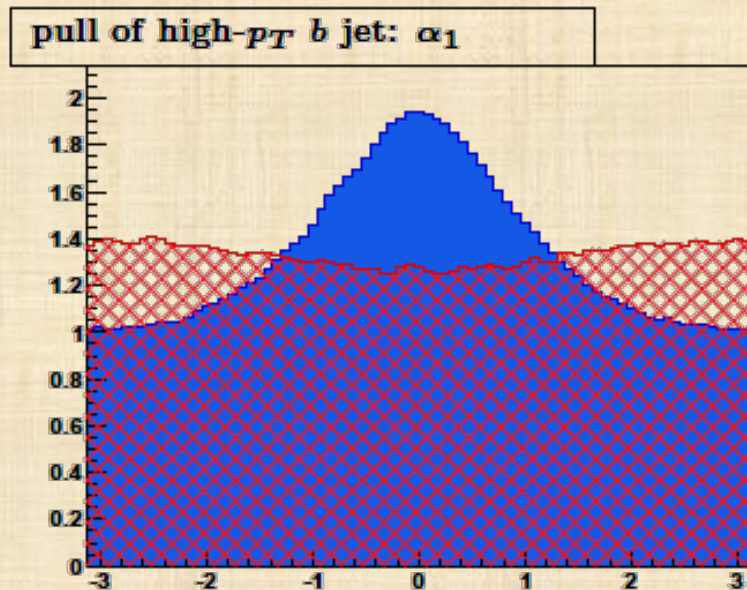$$\vec{p} = \sum_i \frac{E_T^i \, |r_i|}{E_T^{jet}} \, \vec{r_i}$$


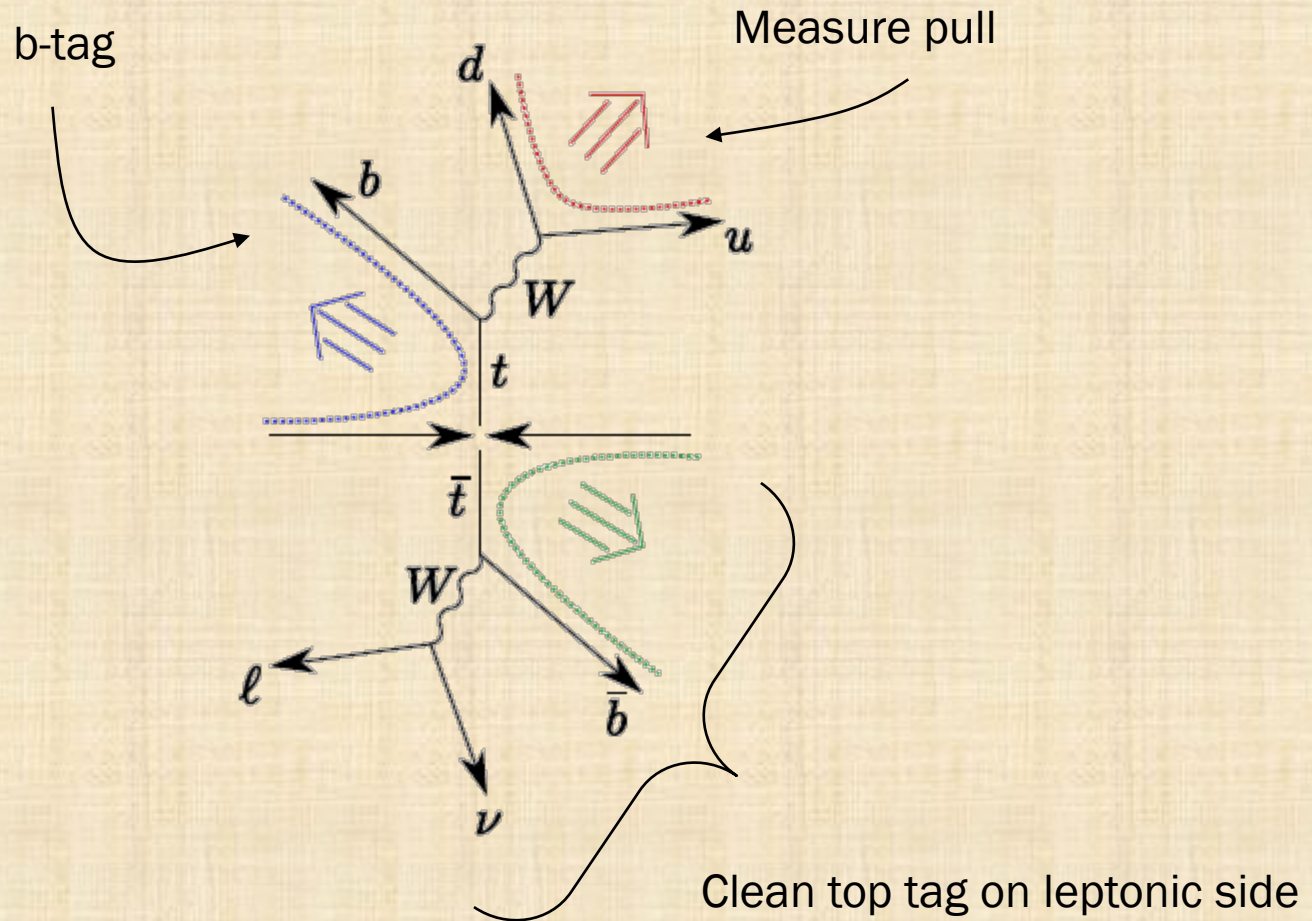
- Angle much more important than length
- Look at relative pull angles

# PULL HAS BEEN MEASURED!



Note 6087-CONF Aug 2010, Andy Haas: $ZH \to b\bar{b}\nu\bar{\nu}$
(consistent with flat background)

# MANY PULL ANGLES

# MEASURED BY DZERO

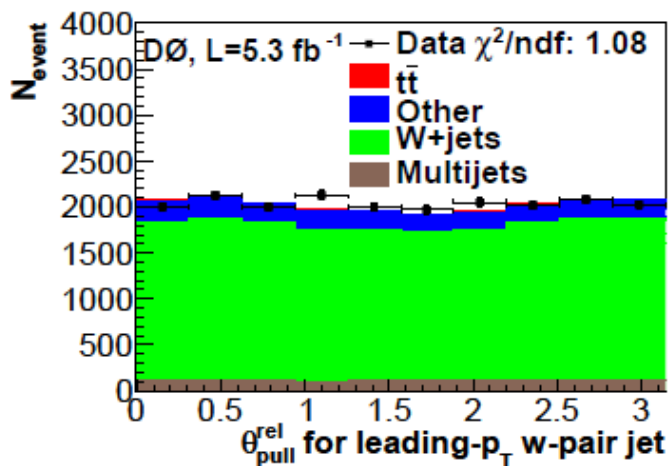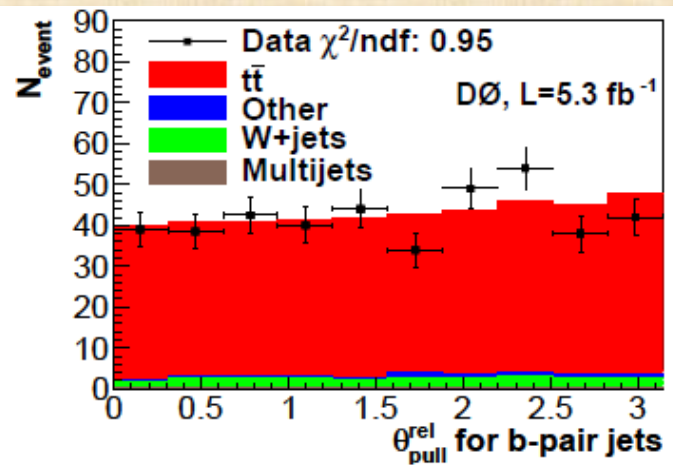Andy Haas and Yvonne Peters, hep-ex:1101.0648
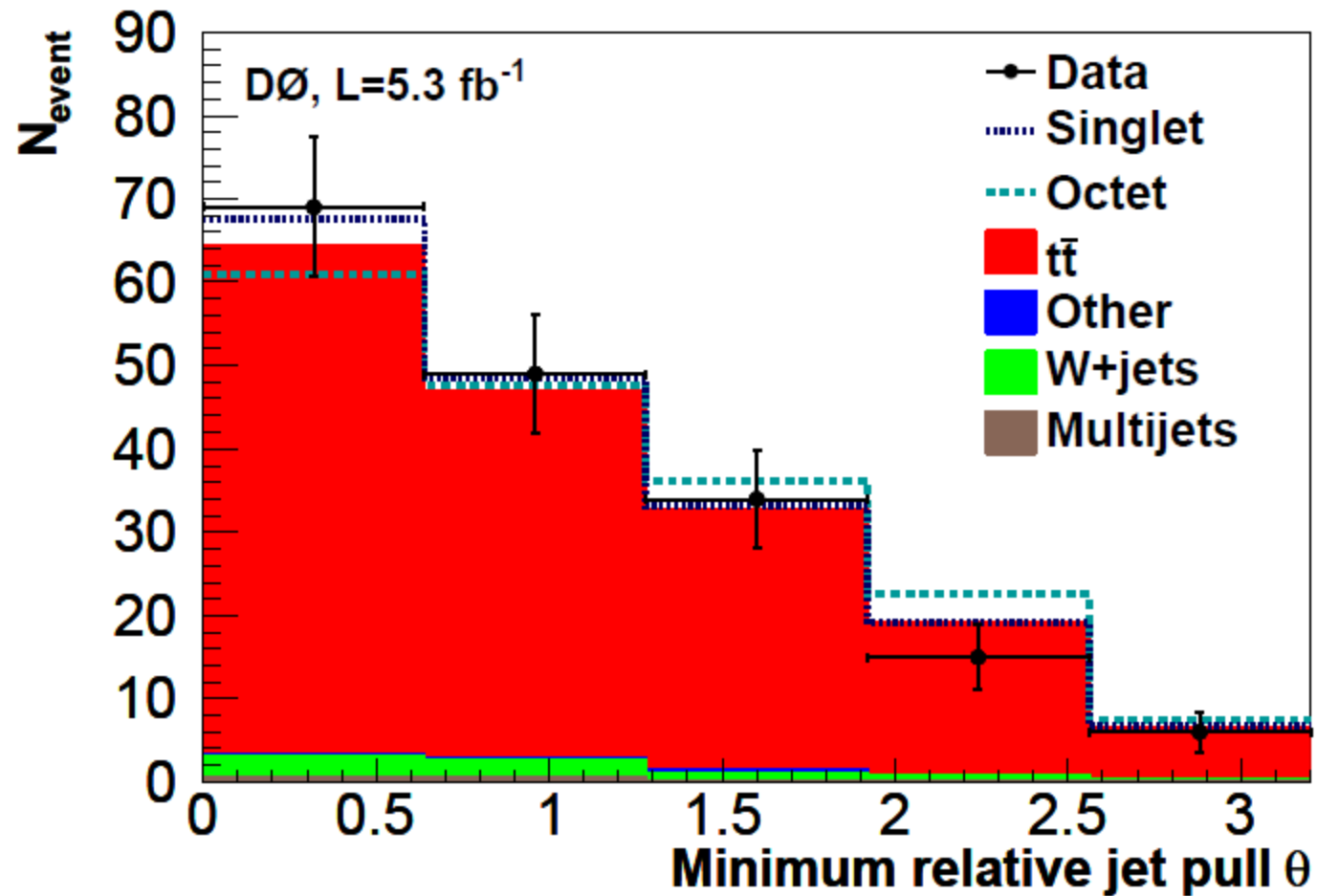
# NON-FLAT PULL
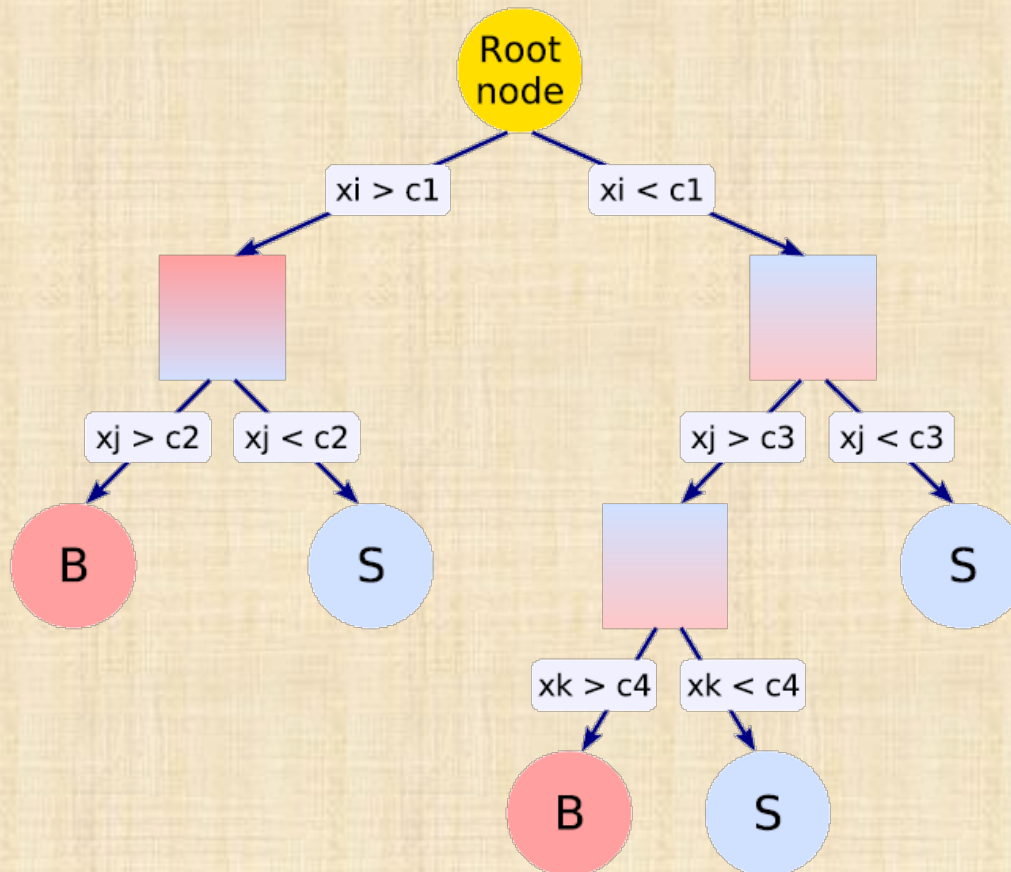
# RULED OUT COLOR OCTET W

# PYTHIA VS HERWIG



Seems robust.

Can we calculate pull??? Good theory question…

# HOW DOES PULL HELP

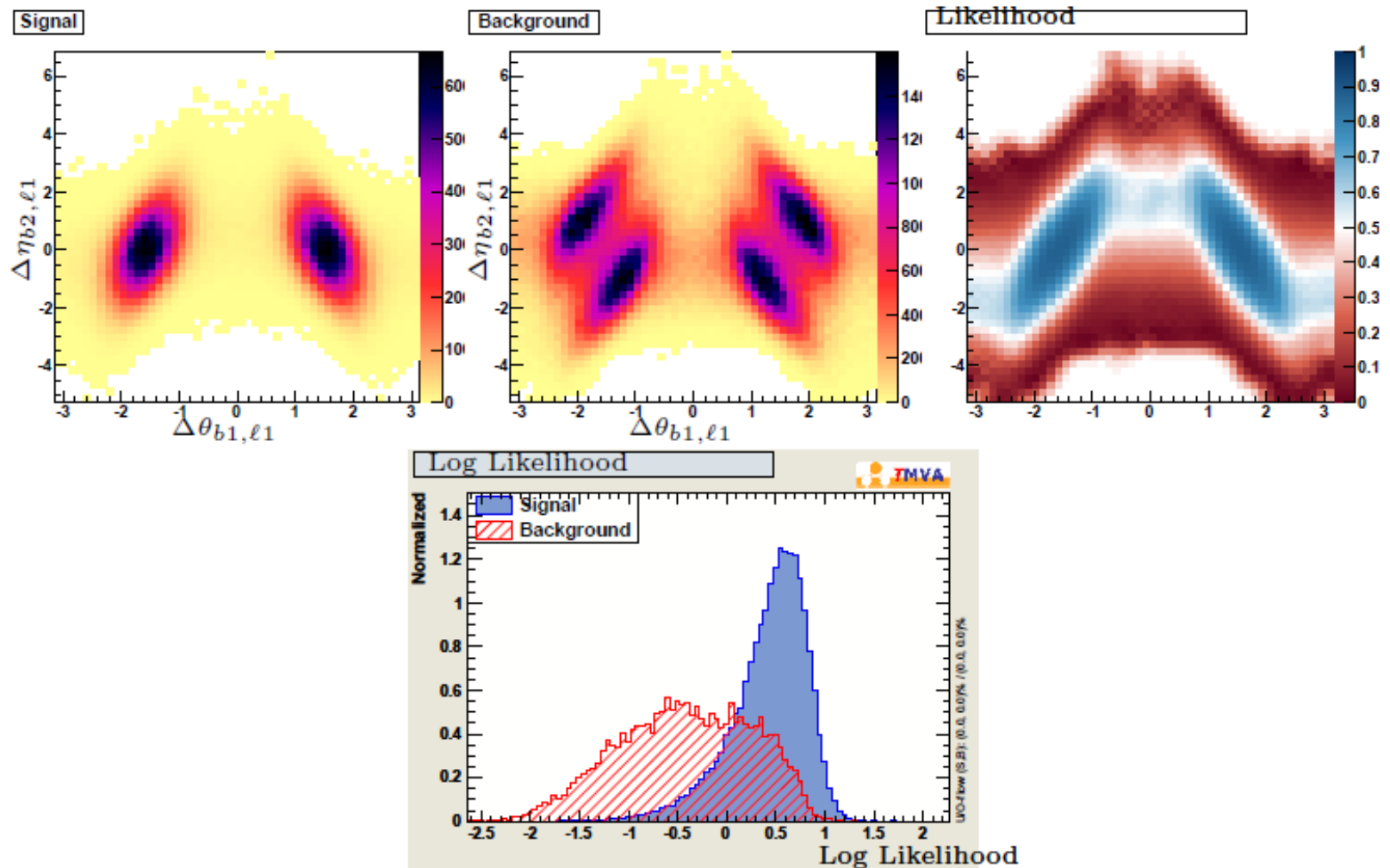- According to Dzero it gives around 5% improvement in pp -> ZH -> νν bb

- How does that work?  **Boosted Decision Trees**
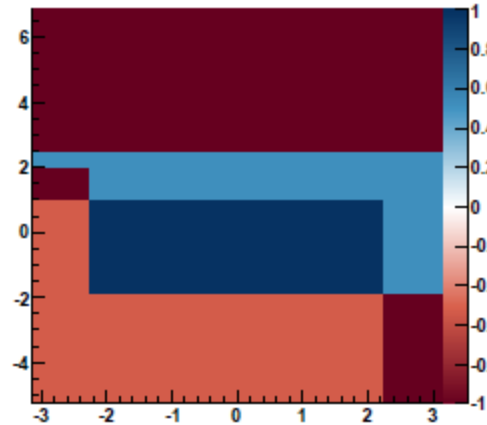


- Train multiple trees and have them vote

- Approximates exact solution.

# VISUALIZE IMPROVEMENT
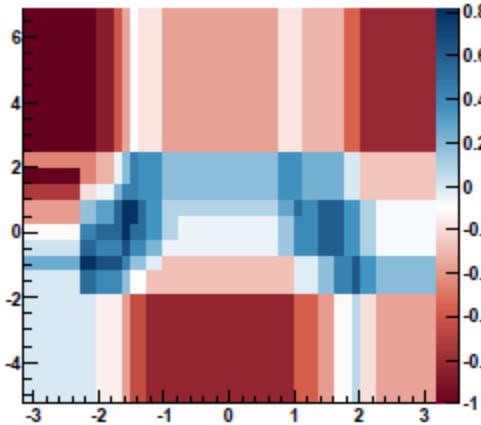
Receiver Operator Characteristic (ROC)

# TAKE RATIOS

$$\frac{S}{B} \quad \xrightarrow{\text{cut}} \quad \frac{\varepsilon_S S}{\varepsilon_S B} = \left(\frac{\varepsilon_S}{\varepsilon_B}\right) \frac{S}{B}.$$



LHC HZ : Signal over Background

- $\Delta\eta_{b\bar{b}}$
- $\Delta y_{H,b2}$
- $\Delta y_{H,b1}$
- $p_T^{b1}$
- $p_T^{b2}$
- CM $\cos(\theta_{b2})$
- $\tau_{b\bar{b}}$
- $\Delta R_{b\bar{b}}$
- $\Delta\phi_{\ell^+\ell^-}$
- $p_T^Z$

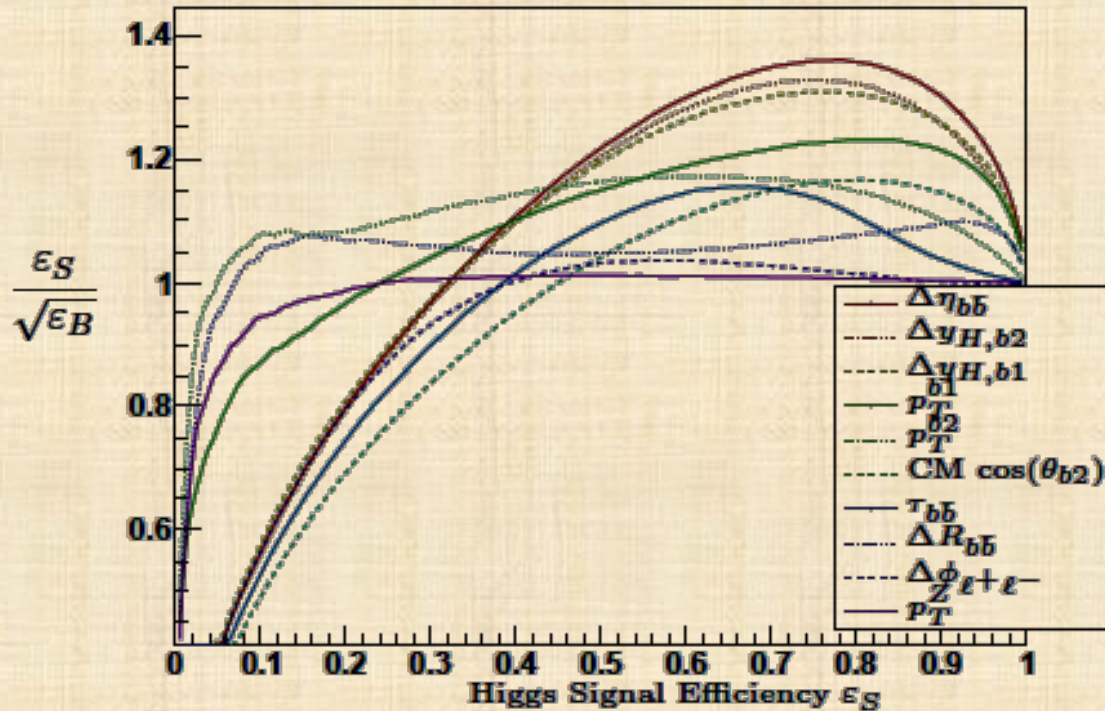$\frac{\varepsilon_S}{\varepsilon_B}$ vs Higgs Signal Efficiency $\varepsilon_S$

- Hard cuts make arbitrary large improvement in S/B
- S/B important, but misleading to optimize or compare variables

# SIC CURVES

$$\sigma \equiv \frac{S}{\sqrt{B}} \quad \xrightarrow{\text{cut}} \quad \frac{\varepsilon_S S}{\sqrt{\varepsilon_B B}} = \left(\frac{\varepsilon_S}{\sqrt{\varepsilon_B}}\right)\sigma$$


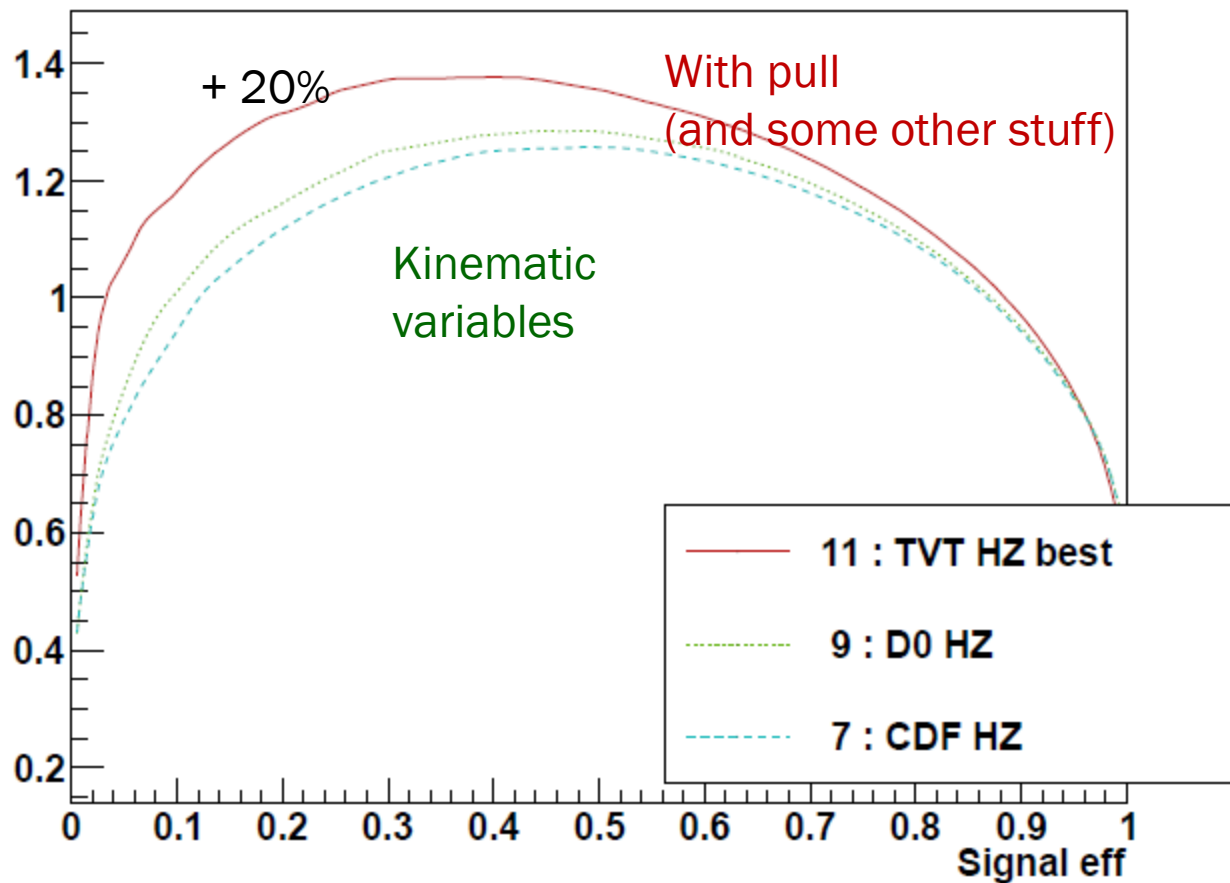
$$\text{SIC} \equiv \frac{\varepsilon_S}{\sqrt{\varepsilon_B}}$$

- Nice visualization
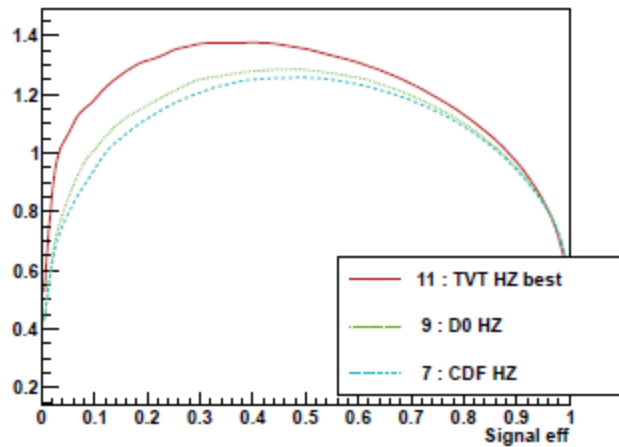- Has maximum at interesting place
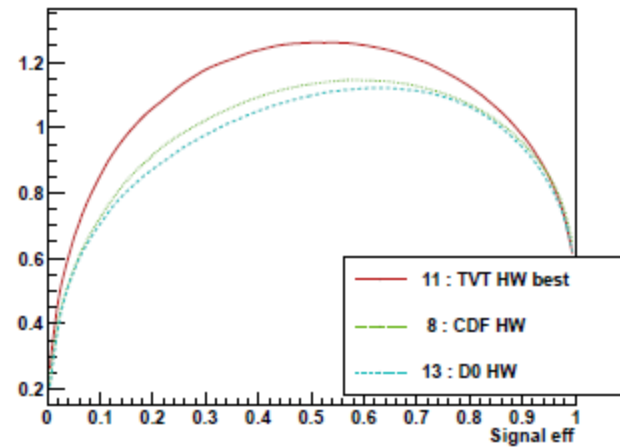- Well defined way to compare variables

# ADDING PULL HELPS

# ALSO AT THE LHC

# OPTIMIZE W TAGGING



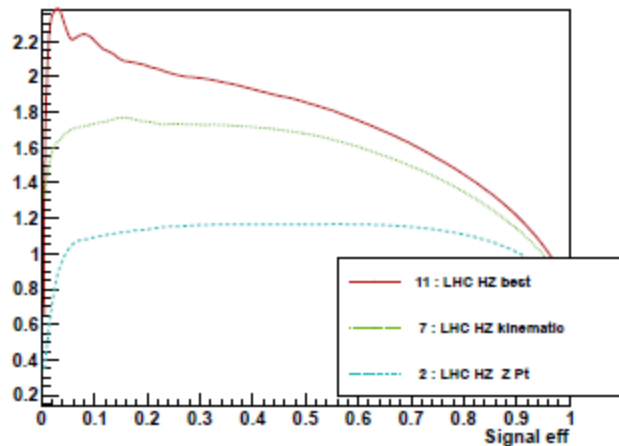Huge improvement in significance with multivariate approach

# CONCLUSIONS

## New machine needs new tricks

- Jet substructure extremely useful
    - Jet mass, Jet shapes, R-cores, Splitting scales
    - Filtering, Trimming, Pruning

- Complimentary information in superstructure
    - Sensitive to other jets – global information
    - Measures color flow

- Correlations are subtle
    - Multivariate techniques are essential
    - Boosted Decision Trees work well if used carefully
    - Proper visualization makes comparisons much easier
        - e.g. SIC curves

$$\mathrm{SIC} \equiv \frac{\varepsilon_S}{\sqrt{\varepsilon_B}}$$