# Demystifying Multivariate Searches
## and the Hunt for the Higgs

Matthew Schwartz

Harvard University

Johns Hopkins University
September 20, 2010

# Part 1:

# Motivation

# WHERE IS THE HIGGS?

| Parameter | Input value |
|-----------|-------------|
| $M_Z$ [GeV] | $91.1875 \pm 0.0021$ |
| $\Gamma_Z$ [GeV] | $2.4952 \pm 0.0023$ |
| $\sigma_{had}^{0}$ [nb] | $41.540 \pm 0.037$ |
| $R_{\ell}^{0}$ | $20.767 \pm 0.025$ |
| $A_{FB}^{0,\ell}$ | $0.0171 \pm 0.0010$ |
| $A_{\ell}$ $^{(\star)}$ | $0.1499 \pm 0.0018$ |
| $A_c$ | $0.670 \pm 0.027$ |
| $A_b$ | $0.923 \pm 0.020$ |
| $A_{FB}^{0,c}$ | $0.0707 \pm 0.0035$ |
| $A_{FB}^{0,b}$ | $0.0992 \pm 0.0016$ |
| $R_c^{0}$ | $0.1721 \pm 0.0030$ |
| $R_b^{0}$ | $0.21629 \pm 0.00066$ |
| $\sin^2\theta_{eff}^{\ell}(Q_{FB})$ | $0.2324 \pm 0.0012$ |

Combine many observables to constrain Higgs mass



Indirect Exclusion (95%)

Theory uncertainty
— Fit including theory errors
···· Fit excluding theory errors

45     80     144

# WHERE IS THE HIGGS?

| Parameter | Input value |
|---|---|
| $M_Z$ [GeV] | $91.1875 \pm 0.0021$ |
| $\Gamma_Z$ [GeV] | $2.4952 \pm 0.0023$ |
| $\sigma_{\text{had}}^0$ [nb] | $41.540 \pm 0.037$ |
| $R_\ell^0$ | $20.767 \pm 0.025$ |
| $A_{\text{FB}}^{0,\ell}$ | $0.0171 \pm 0.0010$ |
| $A_\ell$ $^{(\star)}$ | $0.1499 \pm 0.0018$ |
| $A_c$ | $0.670 \pm 0.027$ |
| $A_b$ | $0.923 \pm 0.020$ |
| $A_{\text{FB}}^{0,c}$ | $0.0707 \pm 0.0035$ |
| $A_{\text{FB}}^{0,b}$ | $0.0992 \pm 0.0016$ |
| $R_c^0$ | $0.1721 \pm 0.0030$ |
| $R_b^0$ | $0.21629 \pm 0.00066$ |
| $\sin^2\theta_{\text{eff}}^\ell(Q_{\text{FB}})$ | $0.2324 \pm 0.0012$ |

Combine many observables to constrain Higgs mass

# WHERE IS THE HIGGS?

| Parameter | Input value |
|-----------|-------------|
| $M_Z$ [GeV] | $91.1875 \pm 0.0021$ |
| $\Gamma_Z$ [GeV] | $2.4952 \pm 0.0023$ |
| $\sigma^0_{\text{had}}$ [nb] | $41.540 \pm 0.037$ |
| $R^0_\ell$ | $20.767 \pm 0.025$ |
| $A^{0,\ell}_{\text{FB}}$ | $0.0171 \pm 0.0010$ |
| $A_\ell$ $^{(\star)}$ | $0.1499 \pm 0.0018$ |
| $A_c$ | $0.670 \pm 0.027$ |
| $A_b$ | $0.923 \pm 0.020$ |
| $A^{0,c}_{\text{FB}}$ | $0.0707 \pm 0.0035$ |
| $A^{0,b}_{\text{FB}}$ | $0.0992 \pm 0.0016$ |
| $R^0_c$ | $0.1721 \pm 0.0030$ |
| $R^0_b$ | $0.21629 \pm 0.00066$ |
| $\sin^2\theta^\ell_{\text{eff}}(Q_{\text{FB}})$ | $0.2324 \pm 0.0012$ |

Combine many observables to constrain Higgs mass

Recent Tevatron exclusion (158-175 GeV)

CDF+D0
arXiv: 1007.4587



If it exists, Higgs is most likely light

# HOW DO WE FIND A LIGHT HIGGS?

## Tevatron

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

LEP Exclusion

Tevatron Exclusion

10

95% CL Limit/SM

······· Expected
——— Observed
■ ±1σ Expected
■ ±2σ Expected

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 $pp \to HZ \to b\bar{b}l^+l^-$

CDF note 10235 (summer 2010)

| | |
|---|---|
| $ZH$ | 0.7 |
| $t\bar{t}$ | 9.9 |
| $WW$ | 0.02 |
| $WZ$ | 0.1 |
| $ZZ$ | 3.6 |
| $Z \to \ell\ell + b\bar{b}$ | 22.1 |
| $Z \to \ell\ell + c\bar{c}$ | 2.4 |
| $Z \to \ell\ell + l.f.$ | 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?

# Part 2:

# The Inputs

# KINEMATIC VARIABLES

Standard Stuff

- $P_T$'s of b's and the leptons
- ⬚ ⬚ ⬚⬚ ⬚ for the b jets and the leptons
- ⬚R of the b's and the leptons
- $P_T$ of the reconstructed Z
- $P_T$ of the reconstructed Higgs
- $m_{bb}$: invariant mass of the b's

Less Standard Stuff

- acoplanarity of the b's: $2\pi - \Delta\theta_{b\bar{b}} - \Sigma\theta_{b\bar{b}}$
- acoplanarity of the leptons
- transverse mass of the bb system $m_T^{b\bar{b}} = m^2 + p_x^2 + p_y^2$
- transverse mass of the lepton system
- invariant mass of 2 leptons and 1 b or 2 b's and 1 lepton

P$_T$ of hardest b-jet

background

signal

⬚ ⬚$_{ll}$

# TWIST

Look at 2D distribution in ☐☐ ☐☐ space:



Higgs Signal

Background Initiated by $gg$

twist

It seems that neither ☐ nor ☐ nor R holds the right information

Introducing twist = polar angle in this plane

Background has pole for zero twist (t-channel singularity)

# TWIST

b-jet twist       lepton twist

Parton level – no cuts

Jet level with detector cuts

Could be more generally useful….

# HELICITY AND AZILICITY ANGLES

Angles in Higgs rest frame
relative to H boost direction

$H$ boost direction

$\theta_h$

$b$

$\phi_a$

beam

$\bar{b}$

$H$ anti-boost

helicity angle

azilicity angle

Parton level
– no cuts

H_Frame_b_coshelicity

H_Frame_b_azilicity

Jet level
with detector cuts

H frame cos($\theta_b^*$) helicity angle

H Frame $\phi_{b1}^*$ azilicity angle

Signal is on-shell
- angles meaningful

Background is not a resonance
- Angles meaningless
- Expect peaked due to collinear singularities

# CHINESE MENU METHOD

- Pick a particle: high-$p_T$ $b$-jet, low-$p_T$ $b$-jet, high-$p_T$ lepton, low-$p_T$ lepton, Higgs, $Z$

- Optionally transform to a boosted frame: Higgs, Z, Center of Mass (CM)

- Optionally rotate the polar axis to point along the initial direction of the particle whose frame you're in (for Helicity and Azilicity Angles).

- Pick a kinematic property: $p_T$, $\eta$, $\phi$, $\cos(\theta)$, etc.

- Optionally pick a second particle to form a sum or difference, sometimes with a coordinate transformation as in $\Delta R$ and twist $\tau$, and sometimes with a more complicated combination as in invariant-mass.

- For vector quantities optionally take the magnitude of vector sums, $|\vec{p}_1 \pm \vec{p}_2|$ or scalar sums, $|\vec{p}_1| \pm |\vec{p}_2|$.

  - $\Delta y_{H,b1}$ and $\Delta y_{H,b2}$: Difference in rapidity between Higgs and higher-$p_T$ or lower-$p_T$ $b$-jet

  - $\cos(\theta_{b2}^*)$: Center of Mass frame $\cos(\theta)$ of the lower-$p_T$ $b$-jet. Same for higher-$p_T$ $b$-jet.

  - $\Delta p_T^{Z,l1}$: Difference in $p_T$ between the reconstructed Z and the higher-$p_T$ lepton

  - $\Delta p_T^{b1,l2}$: Difference in $p_T$ between the higher-$p_T$ b-jet and the lower-$p_T$ lepton

  - $\Sigma p_T^{b1,l2}$: Sum of $p_T$'s of the higher-$p_T$ b-jet and the lower-$p_T$ lepton

  - $\Delta \eta_{b1,l2}$: Difference in $\eta$ between the higher-$p_T$ b-jet and the lower-$p_T$ lepton

# EVENT SHAPE VARIABLES

Nothing to do with the particular signal or background

- $H_T$ = Scalar sum of all $E_T$

- $\Sigma p_T$ = Scalar sum of all $p_T$ (which differs from $H_T$ for massive jets)

- $H_z$ = Boost of the center-of-mass system along the beam

- $E_{vis}$ = Scalar sum of all visible energy

- $\hat{s}$ = CM energy for hard collision, or invariant mass of the reconstructed Higgs and Z.

- Centrality

- Aplanarity and Aplanority

- Sphericity and Spherocity

- DShape and Yvariable (related to the eigenvalues that go into defining the above)

- Fox-Wolfram Moments

# JET SUPERSTRUCTURE

What is not in the parton-level kinematics?

- **Global information**
    - Event shapes

- **Color:**
    - Color **charge**: Quark vs. Gluon jets
    - Color **connections**

$$\text{Tr}[T^A T^B] \propto \delta^{AB}$$

Signal

Background

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

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

# 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:

# HOW CAN WE USE IT?

Higgs:

$q\bar{q}$

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

accumulated all Pt weight - notZ

| | |
|---|---|
| Entries | 23715 |
| Mean x | 0.0137 |
| Mean y | 3.139 |
| RMS x | 1.8 |
| RMS y | 1.822 |

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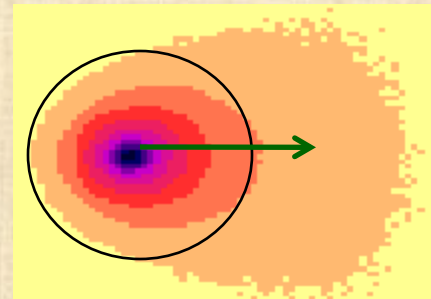- Can we find a *simpler* or more *universal* discriminant?

# PULL



- Find **jets** (e.g. anti-$k_T$)
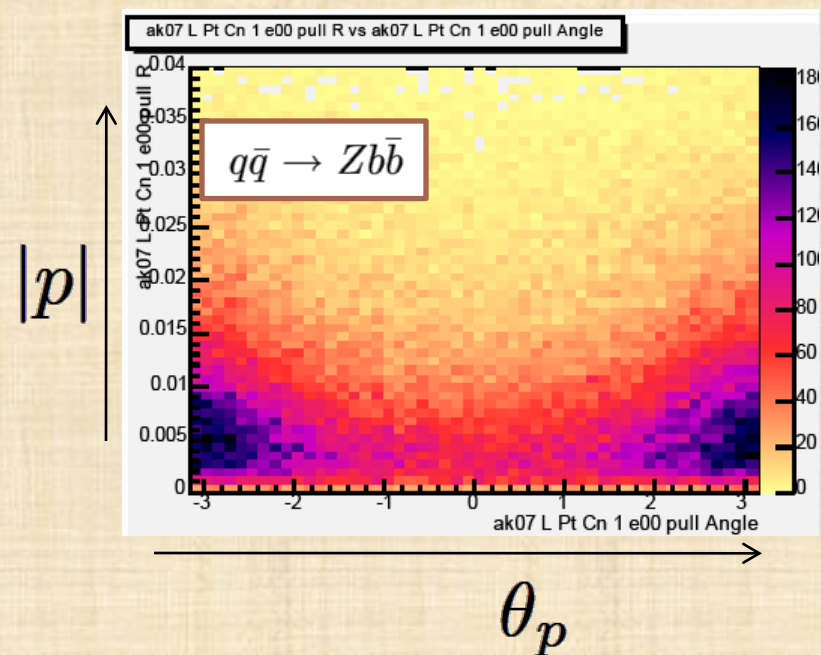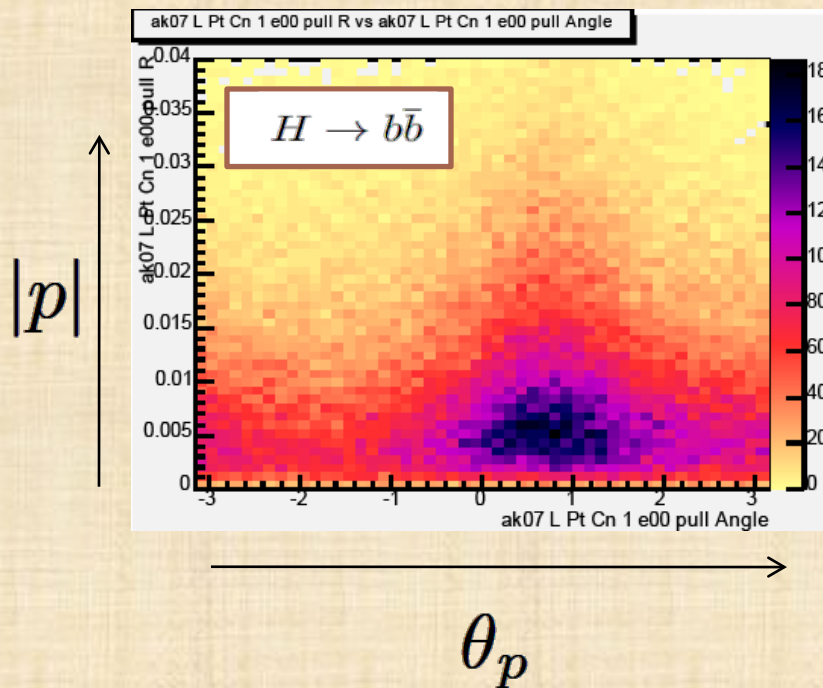- Construct <span style="color:green">pull vector</span> (~ dipole moment) on radiation in **jet**

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

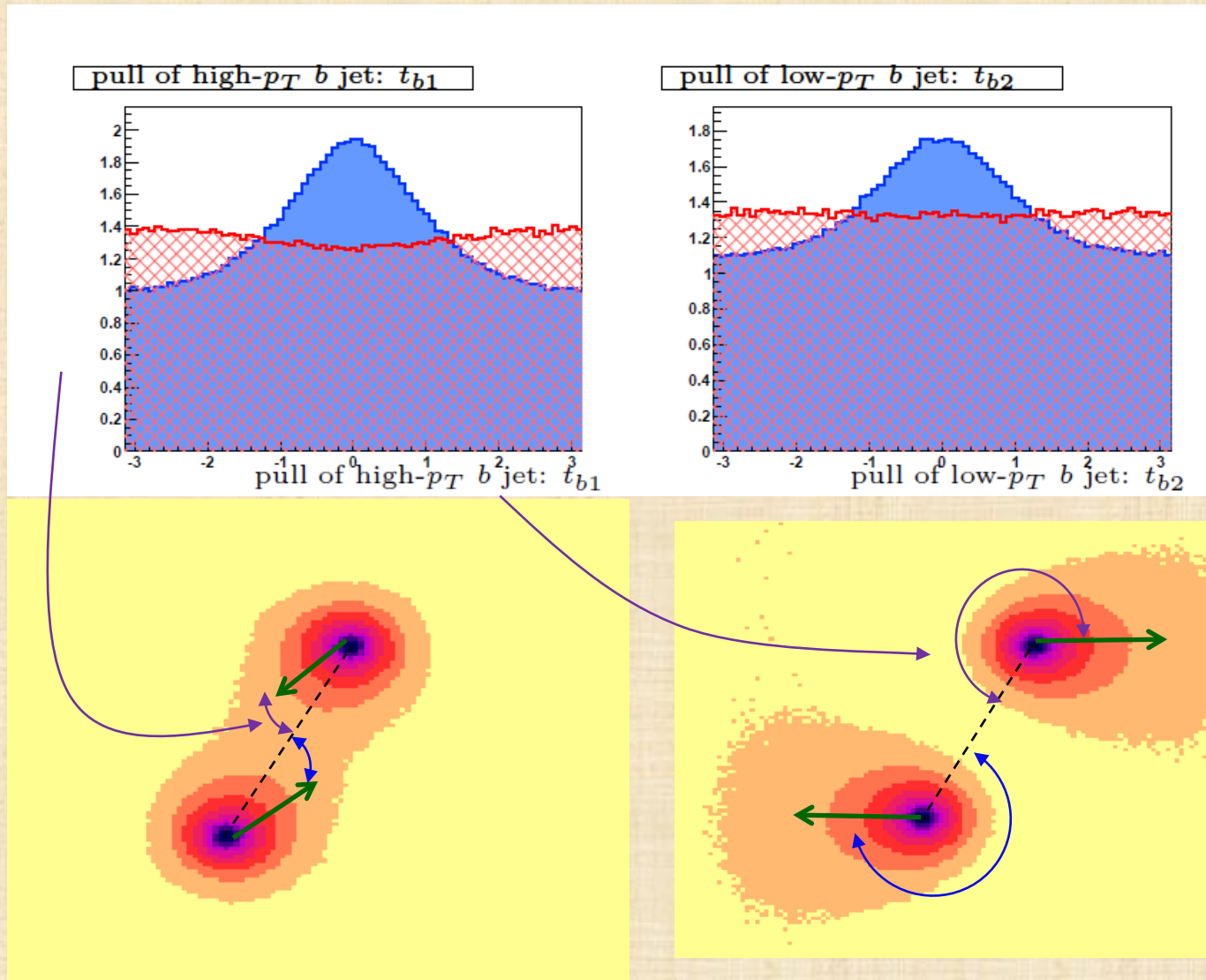- Can use bigger jets for pull, but R = 0.7 seems optimal

# 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
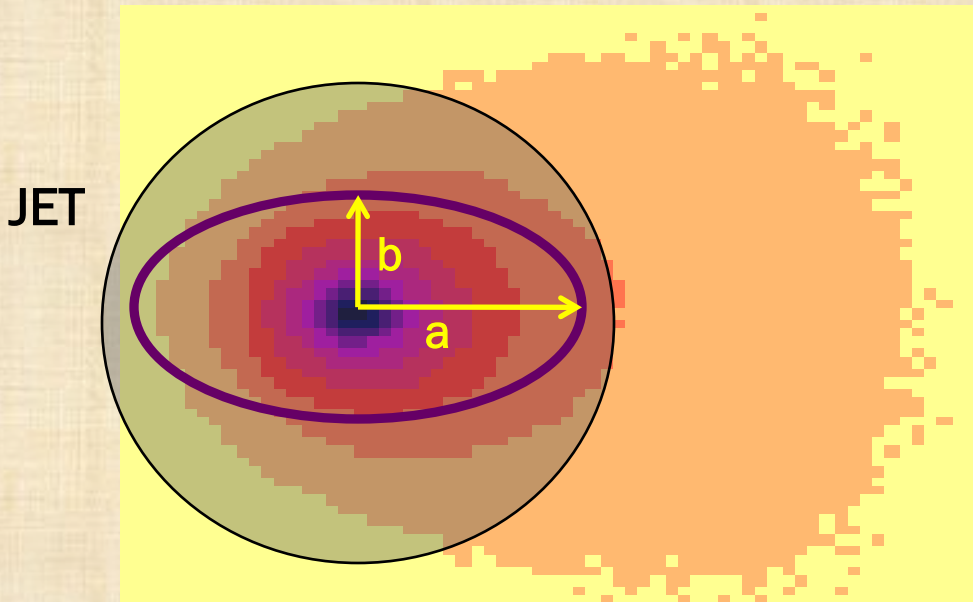- Look at radial pull angle (like for twist)

# SECOND MOMENTS

What about higher moments?

$$\mathbf{I} = \sum_i \frac{E_T^i |r_i|}{E_T^{jet}} \begin{pmatrix} \Delta\phi_i^2 & -\Delta\phi_i\,\Delta\eta_i \\ -\Delta\eta_i\,\Delta\phi_i & \Delta\eta_i^2 \end{pmatrix} \longrightarrow$$ Eigenvalues
a and b



JET

Eccentricity

$$e = \sqrt{\frac{a^2 - b^2}{a^2}}$$

Girth

$$g = \sqrt{a^2 + b^2}$$

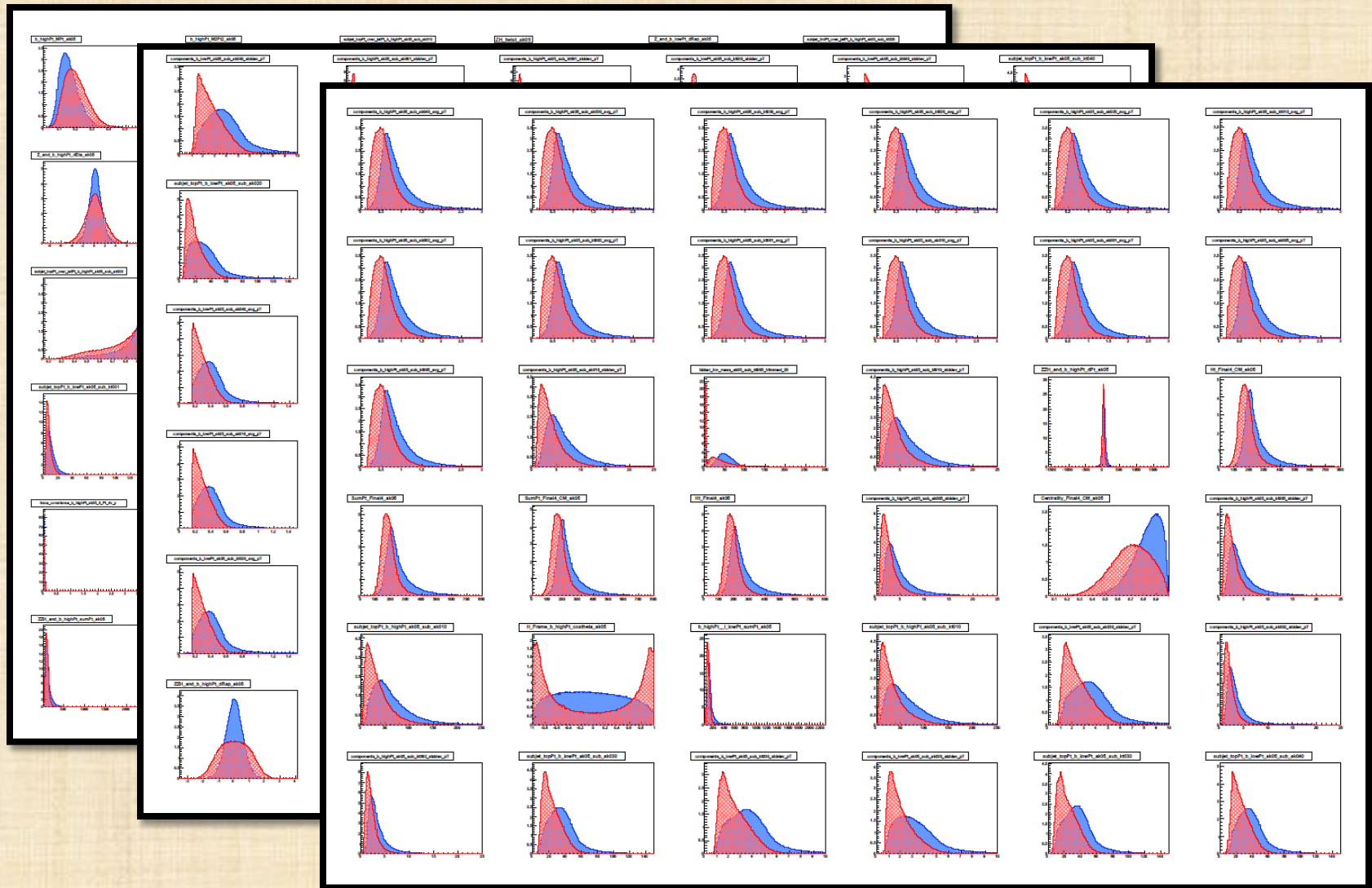# OTHER "SHOWERED" VARIABLES

Many variables vanish at the parton level
- Do not enter the matrix element method
- Complimentary and uncorrelated with kinematic variables

- Mass of each b-jet and the jet mass to $p_T$ ratio

- Rapidity $y$ in addition to pseudorapidity $\eta$ of each massive b-jet

- Subjet multiplicity for each b-jet

- Average $p_T$ of the small subjets within each b-jet

- $p_T$ of hardest, 2nd hardest, and 3rd hardest subjets within each b-jet

- Radial moments ("girth") of each b-jet: $g = \sum_i \frac{p_T^i |r_i|}{p_T^{jet}}$

- Angularity: $\tau_a = \frac{1}{m_{jet}} \sum_i E_i \sin^a \left( \frac{\pi\theta_i}{2R} \right) \left[ 1 - \cos \left( \frac{\pi\theta_i}{2R} \right) \right]^{1-a}$ for $-1 < a < 1$
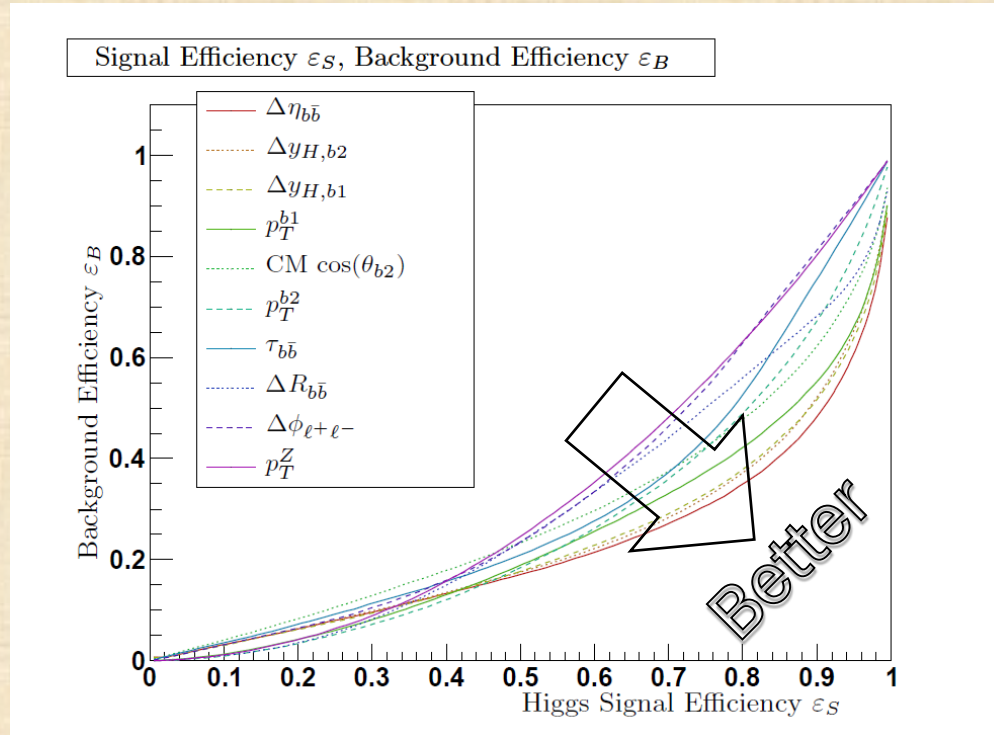
# SUMMARY

- We looked at ~ 900 discriminents!

# Part 3:

# The Output

# EFFICIENCIES

ROC curve:

Background efficiency as a function of signal efficiency

Receiver
Operator
Characteristic



Signal Efficiency $\varepsilon_S$, Background Efficiency $\varepsilon_B$

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

Better

Which variable is best?

# OTHER VISUALIZATIONS

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

$$r_{\frac{S}{B}} \equiv \frac{\varepsilon_S}{\varepsilon_B}$$

$$p_T^Z > 200 \text{ GeV}$$

$\varepsilon_S = 1/20$

$\varepsilon_B = 1/360$

$$\frac{\varepsilon_S}{\varepsilon_B} = 18$$

**S**ignificance
**I**mprovement
**C**haracteristic

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

$$r_\sigma \equiv \frac{\varepsilon_S}{\sqrt{\varepsilon_B}}$$



$$\frac{\varepsilon_S}{\sqrt{\varepsilon_B}} = 0.94$$

- Has maximum
- Maximum $r_\sigma$ can rank variables
- Effective visualization
  - Contains lots of information

# TOP VARIABLES
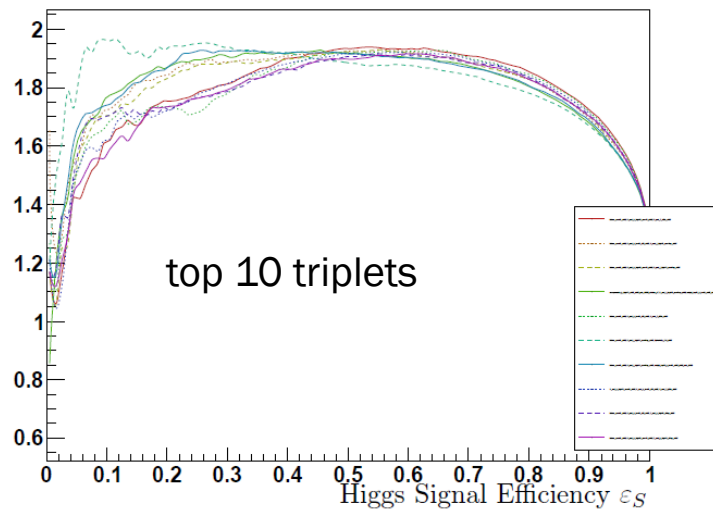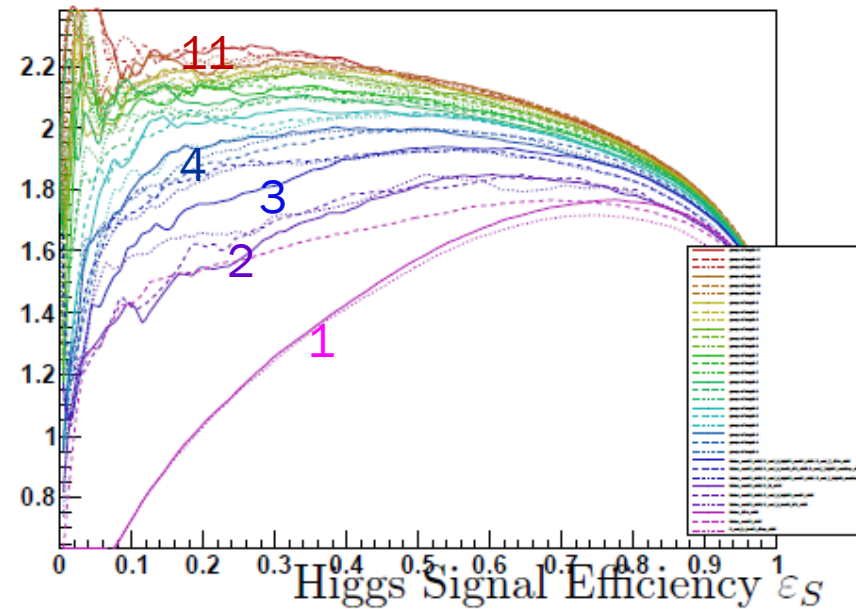
top 10 variables





top 10 pairs, with boosted decision trees

# ADDING MORE

Sequential variable addition
- Take top 3 sets of n variables
- Add any of original 900
- Take top 3 sets of n+1 variables



top 10 triplets

# OBSERVATIONS
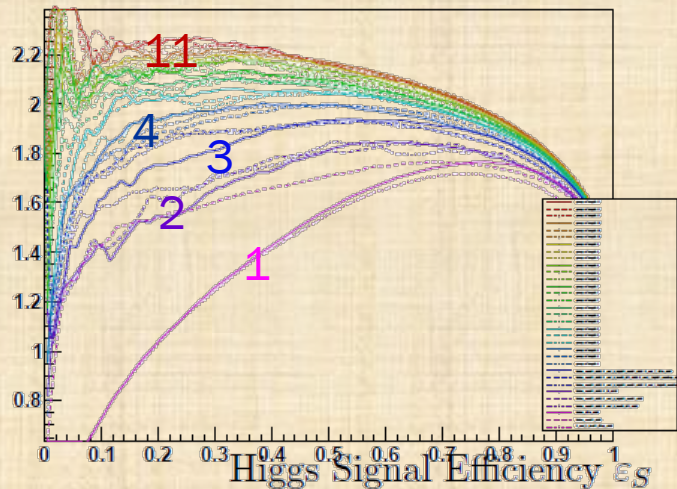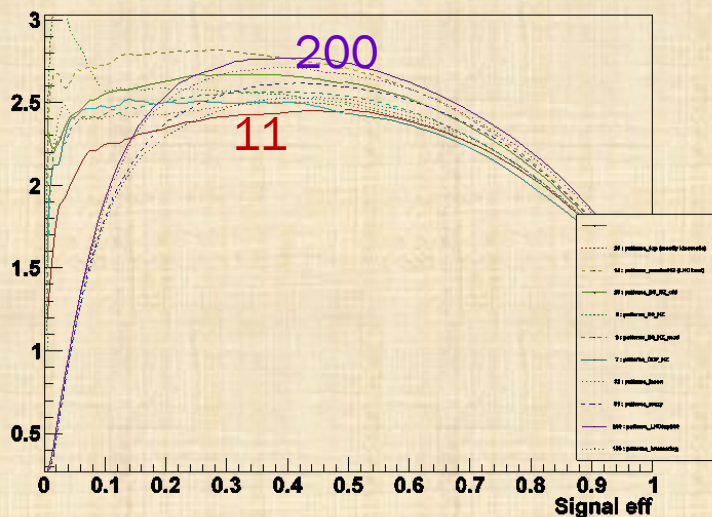


Improvement in $\sigma$



LHC HZ : Significance

- Converges slowly
- Sensitivity to statistics apparent
  - $r_{\Box} = 2$ $\Box_S = 0.05$ gives $\Box_B = 1/1600$
  - 1 million events down to 600

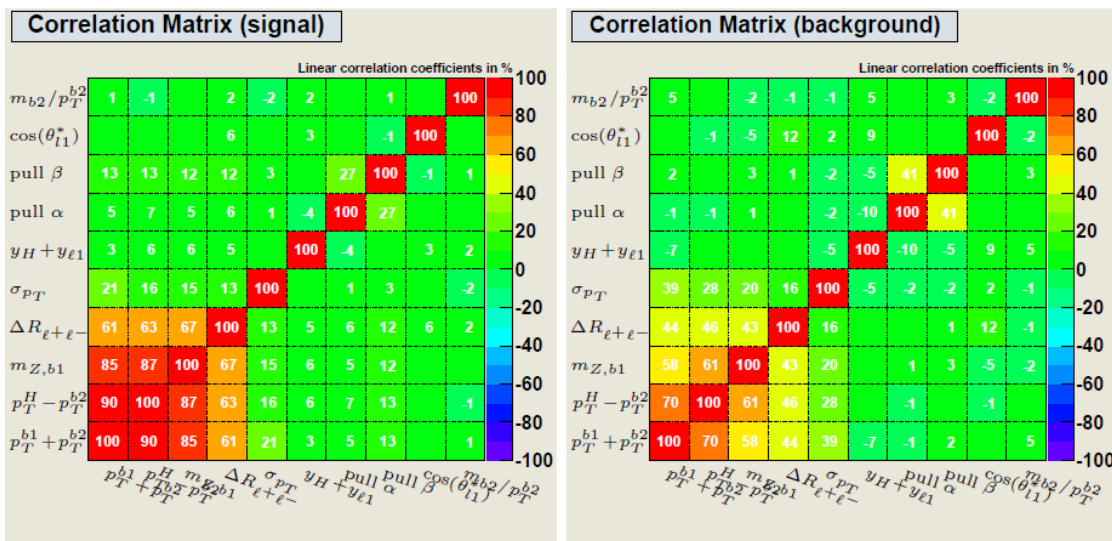- Some variables very poor by themselves, but show up as 5th or 6th variable

Top 10 include

- Higgs $p_T$
- $\Box R_{ZH}$
- Pull
- Twist y (twist with y not $\Box$)
- Event shape D
- Determinant of covariance matrix for radiation in low $p_T$ b jet
- Scalar sum of the b jet $p_T$s

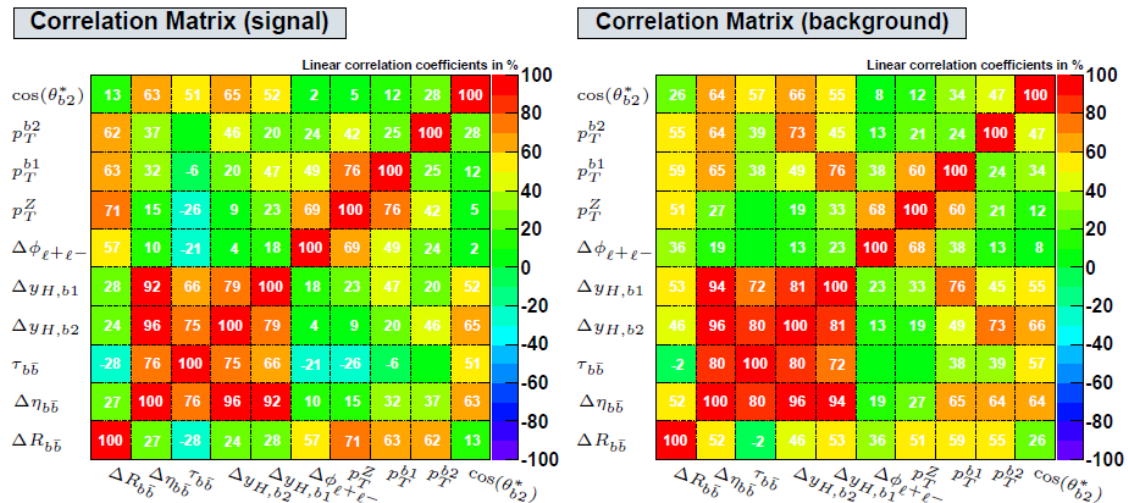# CORRELATIONS OF GOOD 10 COMBO

Best 10 combo
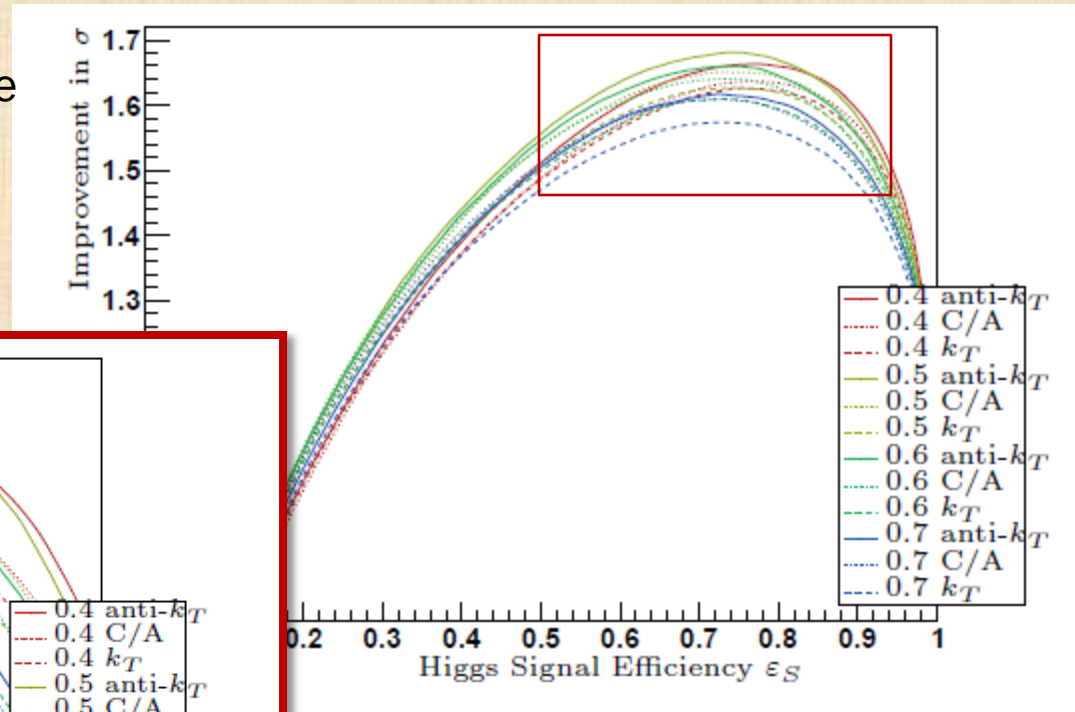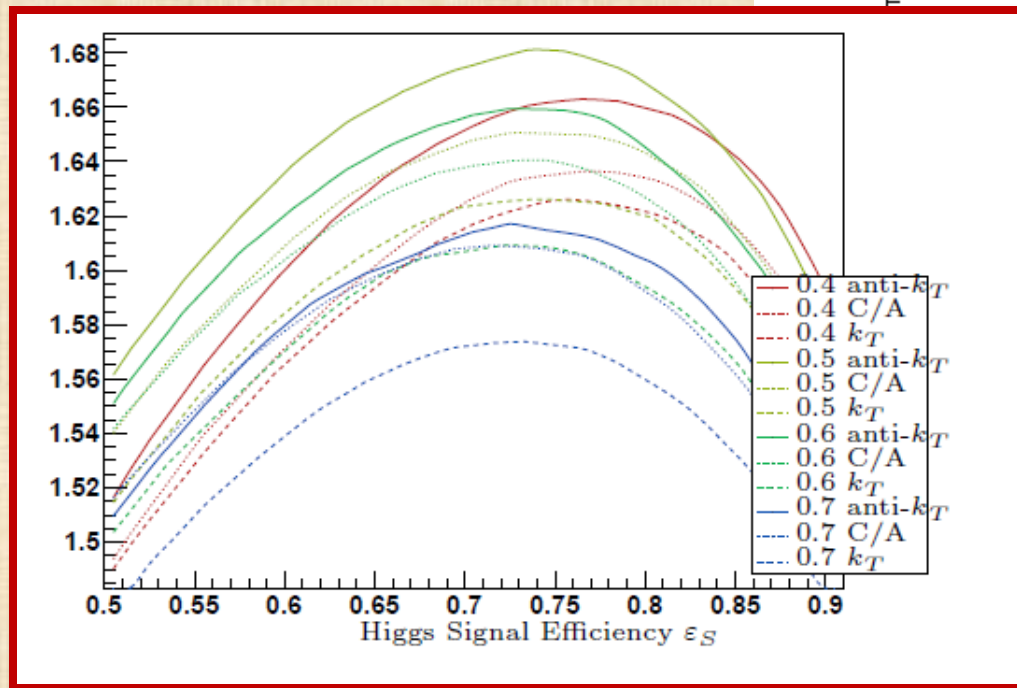
Best 10 individuals

# JET ALGORITHMS

- Main observable is $m_{bb}$
- Look at jet algorithm dependence



- The winner is …
  anti-kT with R = 0.5
- Optimal mass window

$$90 \text{ GeV} < m_{b\bar{b}} < 124 \text{ GeV}$$

# TRIMMING

Krohn, Thaler, Wang

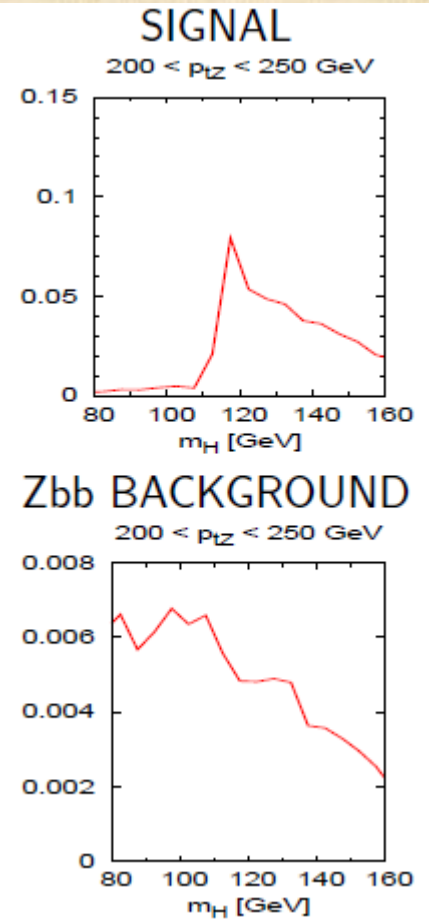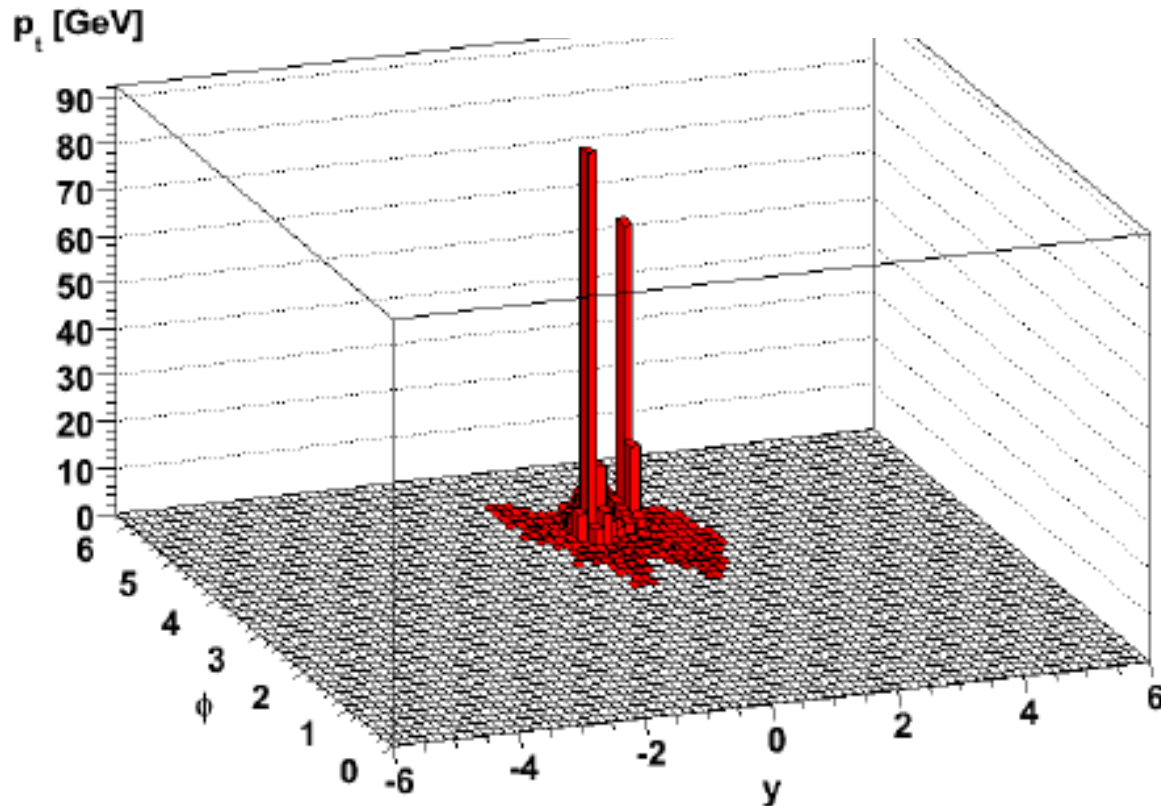1. **Recluster** jet constituents into *very* thin jets

# TRIMMING

Krohn, Thaler, Wang

1. **Recluster** jet constituents into *very* thin jets

2. **Throw away** thin jets that are too soft

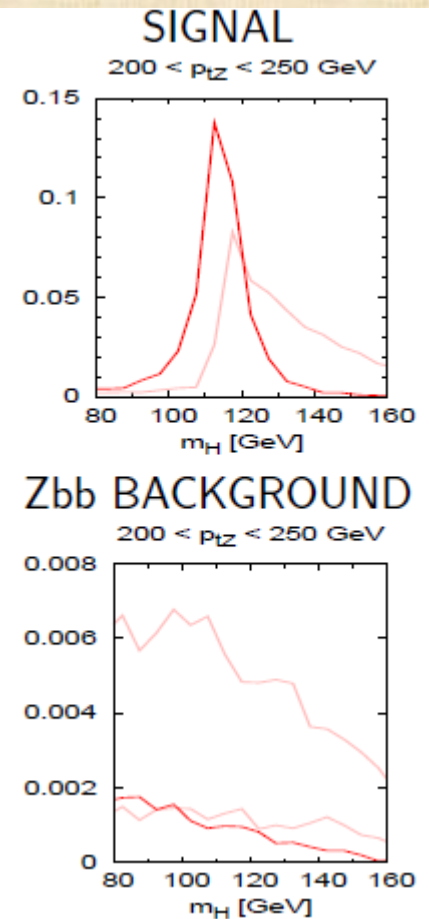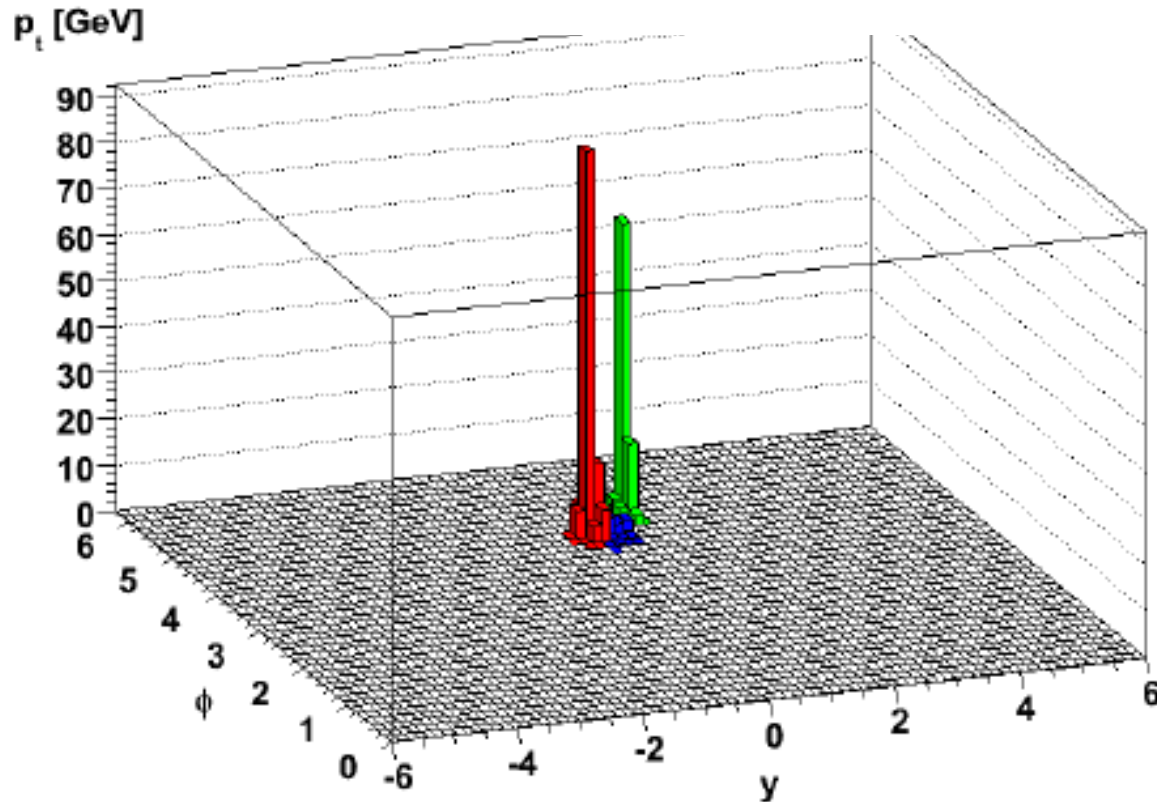# TRIMMING     Boosted H → bb

# TRIMMING

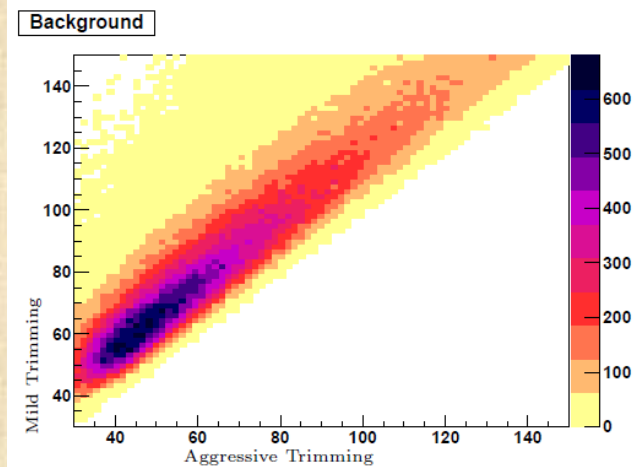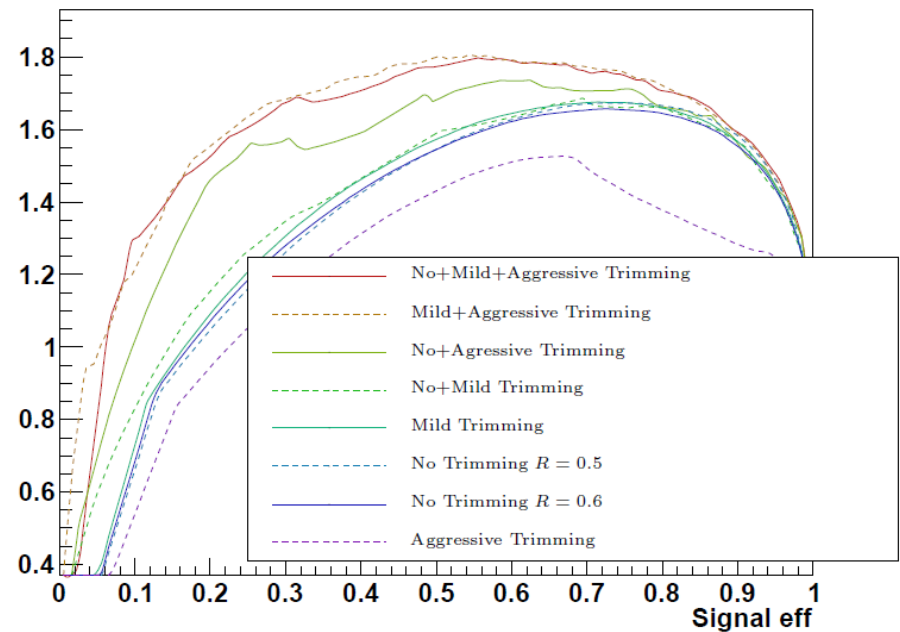Boosted H → bb

# MULTIPLE TRIMMINGS

Trimming does not seem to help much in our case...



Multiple trimmings do help!



(inspired by Soper and Spannowsky)

# CONCLUSIONS

- Final efficiencies still under construction
- Looks like we can help the Tevatron searches
    - around 10% with variables (relative to the ones they use)
    - around 10% with masses (assuming they can trim)
- W/Z + H is totally feasible at the LHC
    - Do not need large $p_T$
    - Discovery potential with 30 fb$^{-1}$

General Observations
- SIC curves provide a useful visualization
    - demonstrate instabilities
    - show covergance
    - visually compare variables' performance
- Uncorrelated variables helpful after kinematics exhausted
- Multiple mass measures useful

Future
- Compare boosted decision trees, random forest, neural networks, etc.
- Compare different generators (Herwig/Pythia)
- Study reducible backgrounds