

# Towards an understanding of jet substructure

Mrinal Dasgupta  
University of Manchester

DESY Hamburg, 13 July 2016

With Gavin Salam, Gregory Soyez, Simone Marzani, Andrzej Siodmok, Alessandro Fregoso, Alex Powling Lais Schunk.

# Boosted objects and jet substructure



Boosted regime implies studying particles with

$$P_T \gg M_X.$$

A common situation at the LHC with access to TeV scales in  $P_T$ .

Also relevant for decays of heavy new particles to electroweak scale objects.

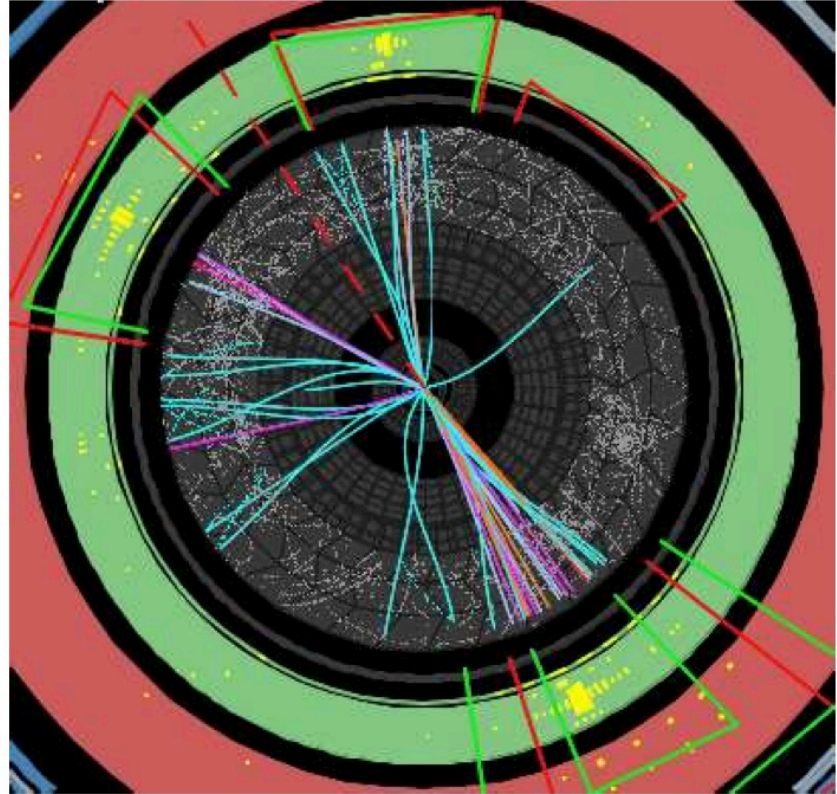
Key observation: Decay products are **collimated**.

$$\theta^2 = \frac{M^2}{p_T^2 z(1-z)}$$

Hadronic two-body decays often reconstructed in single jet.

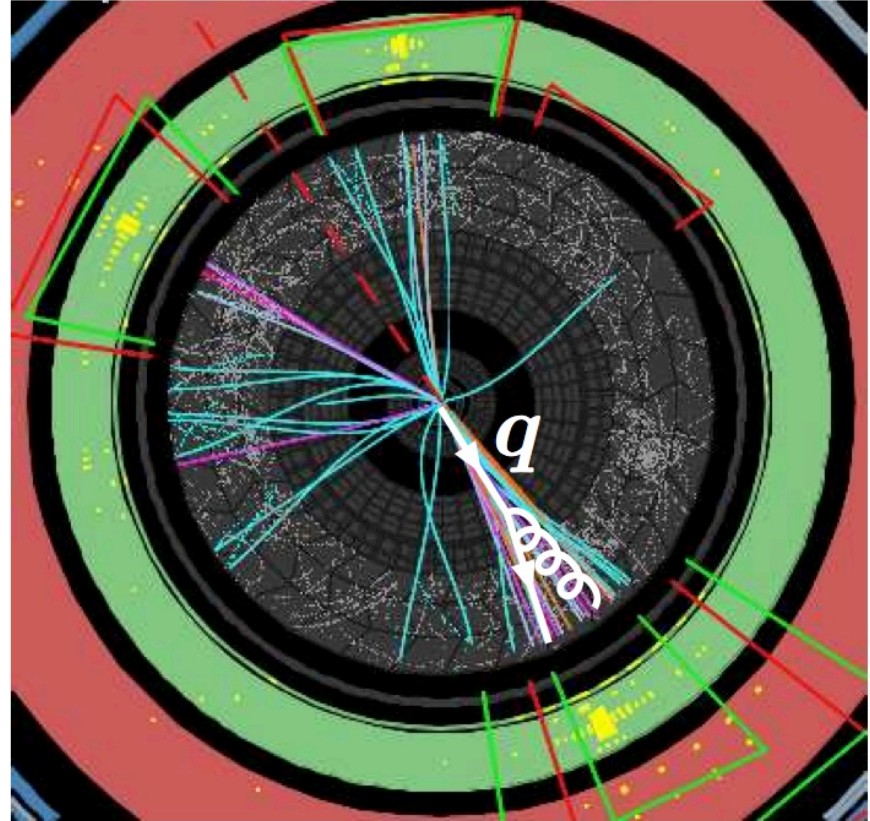
# Jets from QCD vs boosted heavy particles

What jet do we have here?



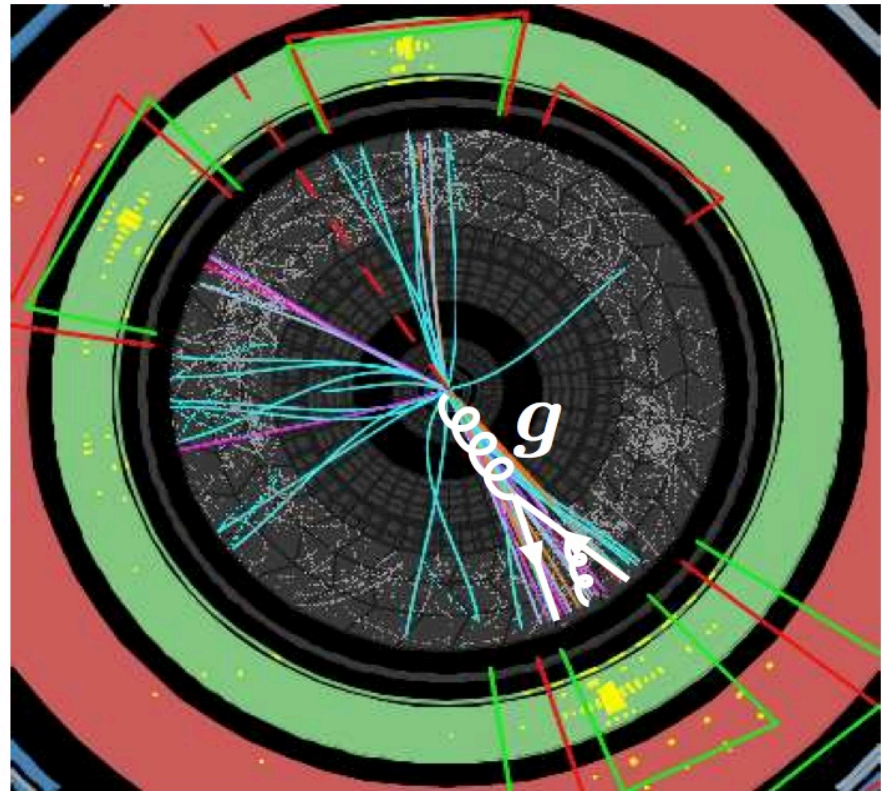
# Jets from QCD vs boosted heavy particles

A quark jet ?



# Jets from QCD vs boosted heavy particles

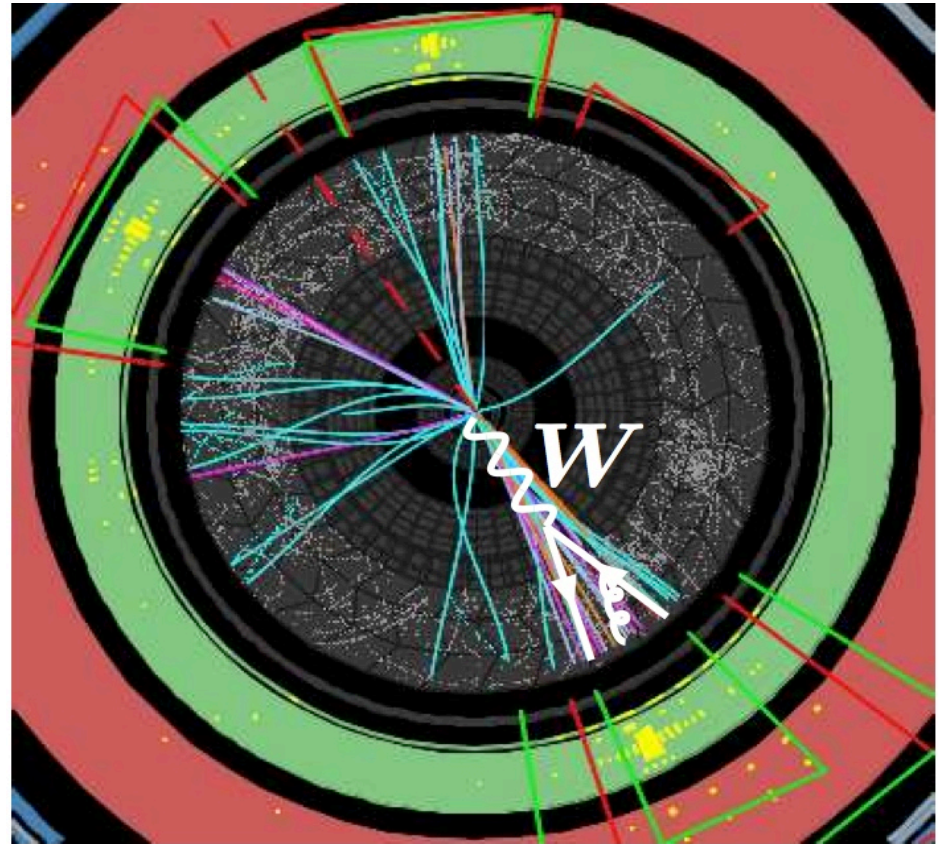
A gluon jet ?





# Jets from QCD vs boosted heavy particles

A W/Z/H ?

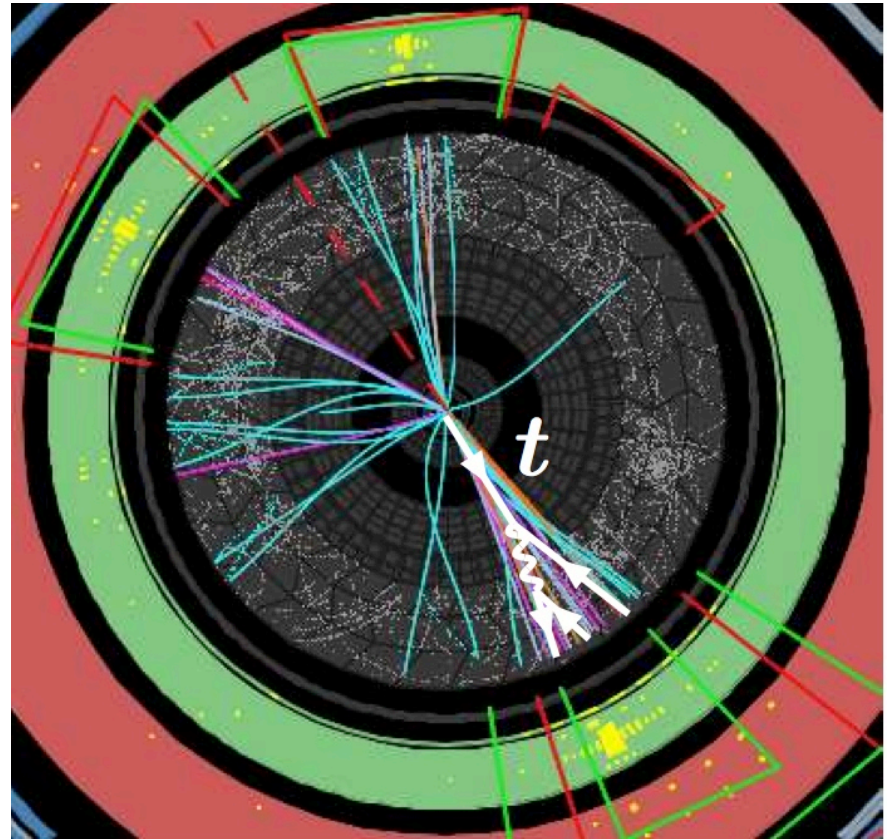


# Jets from QCD vs boosted heavy particles

A top quark?

Source: An ATLAS boosted top candidate

The boosted regime implies a change in paradigm in that jets can be more than quarks and gluons.



# Jet substructure for LHC searches

## Jet substructure as a new Higgs search channel at the LHC

Jonathan M. Butterworth, Adam R. Davison  
*Department of Physics & Astronomy, University College London.*

Mathieu Rubin, Gavin P. Salam  
*LPTHE; UPMC Univ. Paris 6; Univ. Denis Diderot; CNRS UMR 7589; Paris, France.*

It is widely considered that, for Higgs boson searches at the Large Hadron Collider,  $WH$  and  $ZH$  production where the Higgs boson decays to  $b\bar{b}$  are poor search channels due to large backgrounds. We show that at high transverse momenta, employing state-of-the-art jet reconstruction and decomposition techniques, these processes can be recovered as promising search channels for the standard model Higgs boson around 120 GeV in mass.

arXiv:0802.2470v2 [hep-ph] 19 Jun 2008

A key aim of the Large Hadron Collider (LHC) at CERN is to discover the Higgs boson, the particle at the heart of the standard-model (SM) electroweak symmetry breaking mechanism. Current electroweak fits, together with the LEP exclusion limit, favour a light Higgs boson, i.e. one around 120 GeV in mass [1]. This mass region is particularly challenging for the LHC experiments, and any SM Higgs-boson discovery is expected to rely on a combination of several search channels, including gluon fusion  $\rightarrow H \rightarrow \gamma\gamma$ , vector boson fusion, and associated production with  $t\bar{t}$  pairs [2, 3].

Two significant channels that have generally been considered less promising are those of Higgs-boson production in association with a vector boson,  $pp \rightarrow WH, ZH$ , followed by the dominant light Higgs boson decay, to two  $b$ -tagged jets. If there were a way to recover the  $WH$  and  $ZH$  channels it could have a significant impact on Higgs boson searches at the LHC. Furthermore these two channels also provide unique information on the couplings of a light Higgs boson separately to  $W$  and  $Z$  bosons.

Reconstructing  $W$  or  $Z$  associated  $H \rightarrow b\bar{b}$  production would typically involve identifying a leptonically decaying vector boson, plus two jets tagged as containing  $b$ -mesons. Two major difficulties arise in a normal search scenario. The first is related to detector acceptance: leptons and  $b$ -jets can be effectively tagged only if they are reasonably central and of sufficiently high transverse momentum. The relatively low mass of the  $VH$  (i.e.  $WH$  or  $ZH$ ) system means that in practice it can be produced at rapidities somewhat beyond the acceptance, and it is also not unusual for one or more of the decay products to have too small a transverse momentum. The second issue is the presence of large backgrounds with intrinsic

responds to only a small fraction of the total  $VH$  cross section (about 5% for  $p_T > 200$  GeV), but it has several compensating advantages: (i) in terms of acceptance, the larger mass of the  $VH$  system causes it to be central, and the transversely boosted kinematics of the  $V$  and  $H$  ensures that their decay products will have sufficiently large transverse momenta to be tagged; (ii) in terms of backgrounds, it is impossible for example for an event with on-shell top-quarks to produce a high- $p_T$   $b\bar{b}$  system and a compensating leptonically decaying  $W$ , without there also being significant additional jet activity; (iii) the  $HZ$  with  $Z \rightarrow \nu\bar{\nu}$  channel becomes visible because of the large missing transverse energy.

One of the keys to successfully exploiting the boosted  $VH$  channels will lie in the use of jet-finding geared to identifying the characteristic structure of a fast-moving Higgs boson that decays to  $b$  and  $\bar{b}$  in a common neighbourhood in angle. We will therefore start by describing the method we adopt for this, which builds on previous work on heavy Higgs decays to boosted  $W$ 's [4],  $WW$  scattering at high energies [5] and the analysis of SUSY decay chains [6]. We shall then proceed to discuss event generation, our precise cuts and finally show our results.

When a fast-moving Higgs boson decays, it produces a single fat jet containing two  $b$  quarks. A successful identification strategy should flexibly adapt to the fact that the  $b\bar{b}$  angular separation will vary significantly with the Higgs  $p_T$  and decay orientation, roughly

$$R_{b\bar{b}} \simeq \frac{1}{\sqrt{z(1-z)}} \frac{m_H}{p_T}, \quad (p_T \gg m_H), \quad (1)$$

where  $z$ ,  $1-z$  are the momentum fractions of the two quarks. In particular one should capture the  $b, \bar{b}$  and any

First Idea: Seymour  
1993

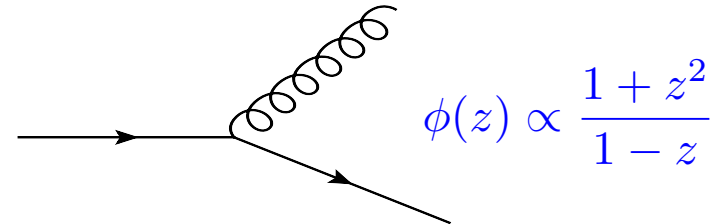
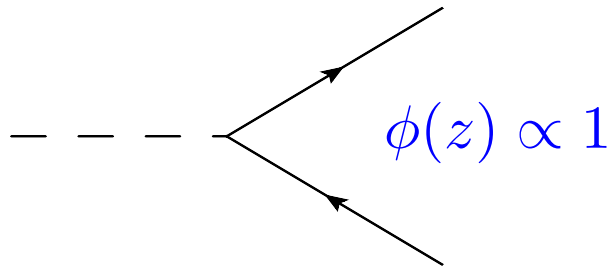
Since 2008 a vibrant  
research field emerged  
based on developing and  
exploiting substructure.

Butterworth, Davison Rubin,  
Salam 2008. Published in PRL.

BDRS paper has over  
600 citations. “Jet  
substructure” title search  
on arXiv gives > 100  
papers post BDRS.



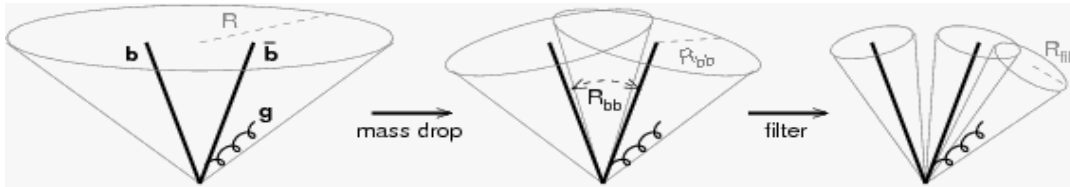
# Signal vs background



BDRS studied the process  $pp \rightarrow VH, H \rightarrow b\bar{b}$

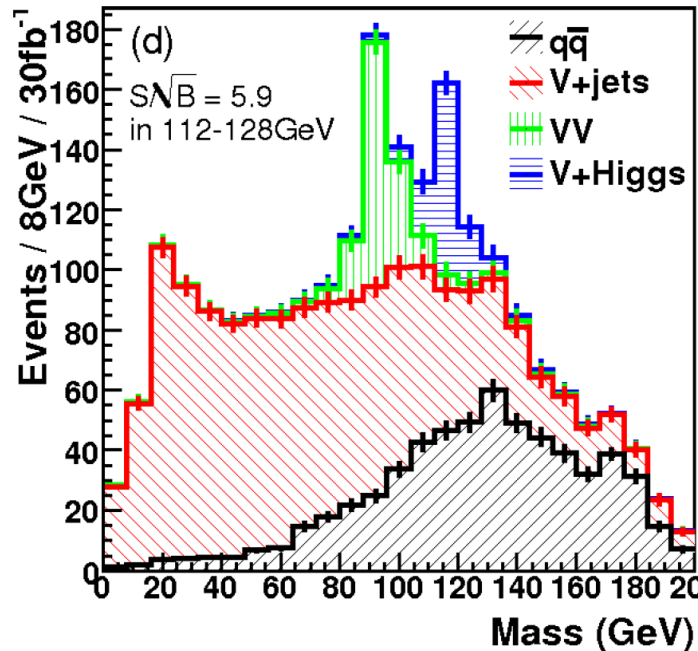
- This was considered an unpromising channel for Higgs discovery due to large QCD backgrounds.
- In boosted limit Higgs decay products are reconstructed in a single fat jet and need to distinguish a signal jet from a plain QCD jet.
- One key is that QCD branchings have **soft enhancements**. Asymmetric sharing of energy compared to Higgs case.

# BDRS mass drop+filtering



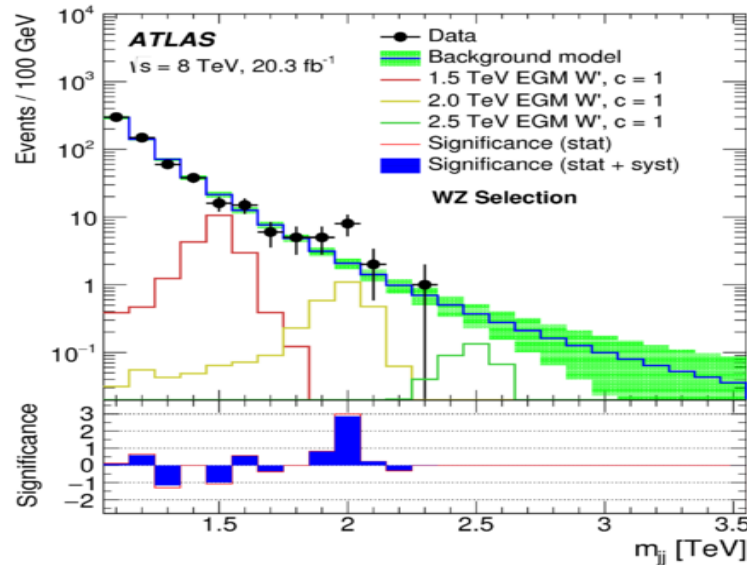
- Break the jet into two subjets  $j_1$  and  $j_2$  such that  $m_{j_1} > m_{j_2}$ .
- If there is a mass drop  $m_{j_1} < \mu m_j$  and the splitting is not too asymmetric  $y = \min(p_{tj_1}^2, p_{tj_2}^2) \Delta R_{j_1 j_2}^2 / m_j^2 > y_{cut}$  then deem the jet tagged or if not discard  $j_2$  and continue.
- Also called the “mass drop” tagger (MDT). more about this later.....
- Filtering method designed to clean the jet of contamination from the Underlying event (grooming).

# BDRS method results



Signal significance of  $4.5\sigma$  was demonstrated in MC studies for a Higgs boson of 115 GeV. Turned this unpromising channel into one of the best discovery channels for light Higgs.

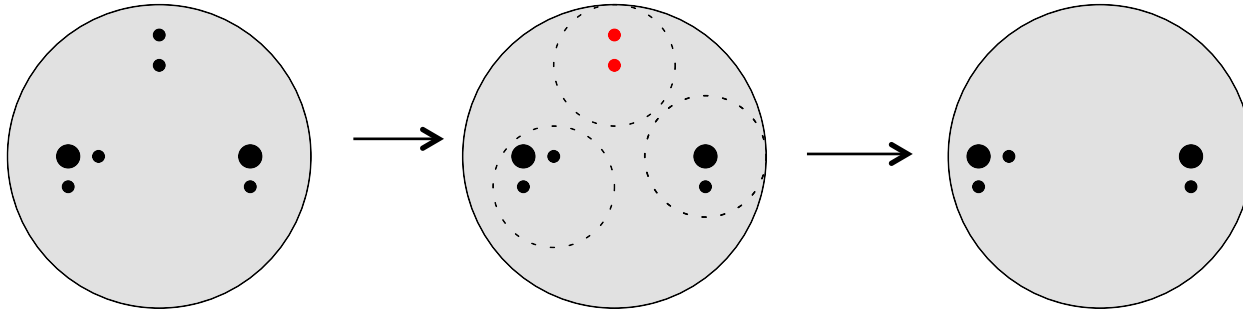
# Jet substructure and LHC searches



- Several methods being used in experimental searches for new physics at LHC.
- Example was recent Run-1 2 TeV **diboson anomaly** observed by ATLAS in hadronic channel. Search for resonances decaying to WZ studied invariant mass of dijets with each jet tagged as a boson jet. Used MDT analysis.



# Several other methods exist



**Trimming** re-clusters jet with smaller radius  $R_{\text{trim}}$ .

Discards subjets with  $p_{t,\text{subjet}} < f_{\text{cut}} p_{t,\text{jet}}$ .

Krohn, Thaler, Wang 2010

**Pruning** is similar but uses a dynamical radius  $R_{\text{prune}} \sim m_j/p_t$ .

Ellis, Walsh, Vermillion 2009

Many other methods: Y-splitter, Atlas top tagger, HEP top tagger, CMS top tagger, JH top tagger, Template Overlap, Planar Flow, Shower Deconstruction, Qjets, N-subjettiness, ECF's etc.

Shall give them collective name of “taggers” for this talk.

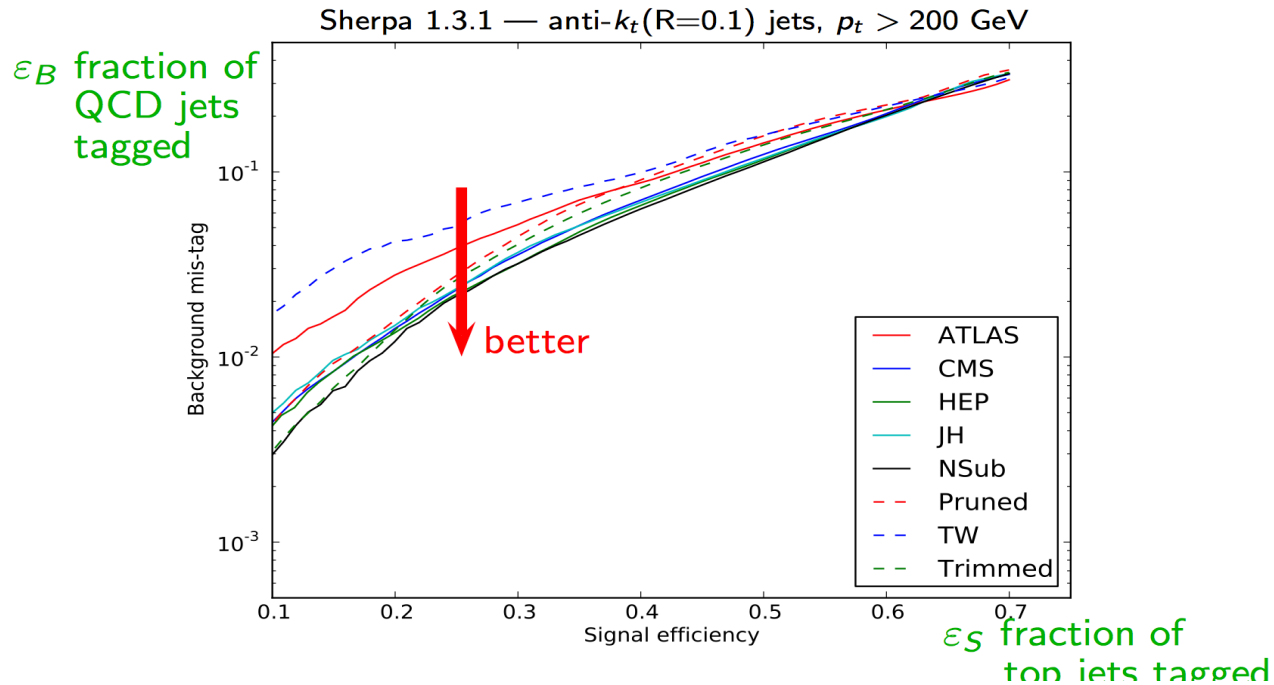
# Some open questions

- Why so many methods?
- Are they really different?
- How to compare methods: number of parameters, vast kinematic range?
- Are tools robust? What is the connection to QCD predictions?

Monte Carlo studies **alone** are insufficient to provide detailed answers to these and other questions.

# Monte Carlo studies

[Boost 2011 proceedings]

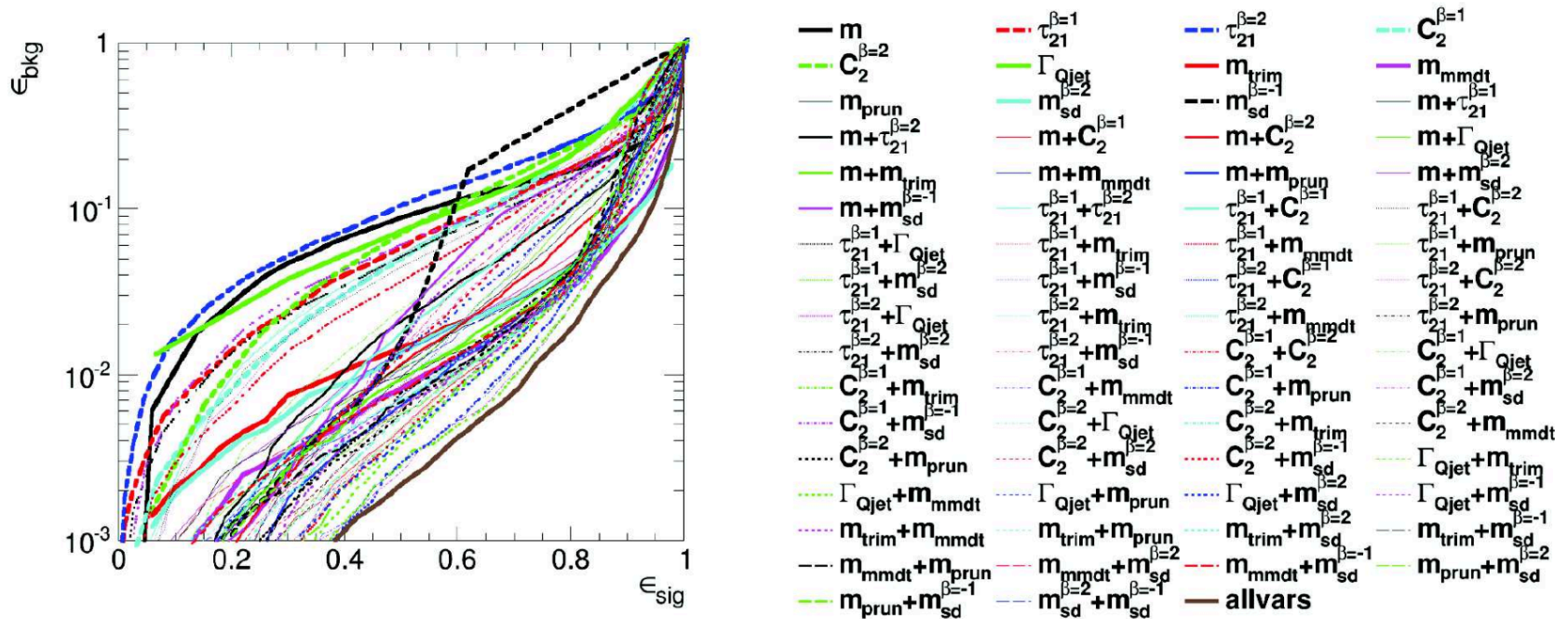


Studies are for fixed parameter settings. No idea about why something works better or if picture changes with parameters.

# More games with Monte Carlo

[Boost 2013 WG]

$W$  v.  $q$  jets: combination of “2-core finder” + “radiation constraint”



Combinations help but details far from obvious.



# An analytical approach?

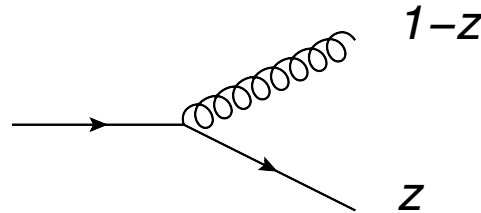


Schwartz, Boost 2012

→ Precision QCD

- Prior to 2013 widely believed that MC studies were only option.
- Analytics was thought impossible due to complexity of taggers and number of scales and parameters involved.
- The tools and precision QCD were largely thought to be incompatible.

# What to compute?



First step in tagging is always cut on jet mass so jet mass distributions of QCD jets before and after grooming are of interest. But plain jet mass distributions at hadron colliders are already **very hard to compute precisely**.

Natural to calculate the distributions in  $\rho = \frac{m^2}{p_T^2 R^2}$  which is invariant under boosts along jet axis.

At LO in soft+collinear limit:

$$\rho \frac{d\sigma}{d\rho} = \frac{C_F \alpha_s}{2\pi} \int dz \frac{d\theta^2}{\theta^2} \left( \frac{1+z^2}{1-z} \right) \delta(\rho - z(1-z)\theta^2)$$

# Plain jet mass

LO result in soft-collinear limit:  $\rho \frac{d\sigma}{d\rho} = \frac{C_F \alpha_s}{\pi} \left( \ln \frac{1}{\rho} - \frac{3}{4} \right)$

Integrated distribution  $\int_0^\rho \frac{d\sigma}{d\rho'} d\rho'$  contains up to double logarithms  $\alpha_s^n L^{2n}$

Logarithmic enhancements spoil convergence of perturbation series  
so **fixed-order is inadequate at small  $\rho$  which is the boosted limit.**

Need to **resum large logarithmic** terms to all orders in perturbation theory. For phenomenology one needs to control also single logarithmic terms  $\alpha_s^n L^n$

# Issues with jet mass

Resummed formula looks like

$$\frac{1}{\sigma} \frac{d\sigma}{d\rho} = \exp [Lg_1(\alpha_s L) + g_2(\alpha_s L) + \dots]$$

Double  
logarithms  
and running  
coupling

Single logarithms from hard  
collinear, soft large-angle and **non-  
global logs**. **Very complicated and  
only possible numerically in large  
 $N_C$  limit.** Dasgupta and Salam 2002

Inspite of complications and large NP corrections  
resummation gives the basic features of jet mass  
distributions. **Can we do the same for taggers?**

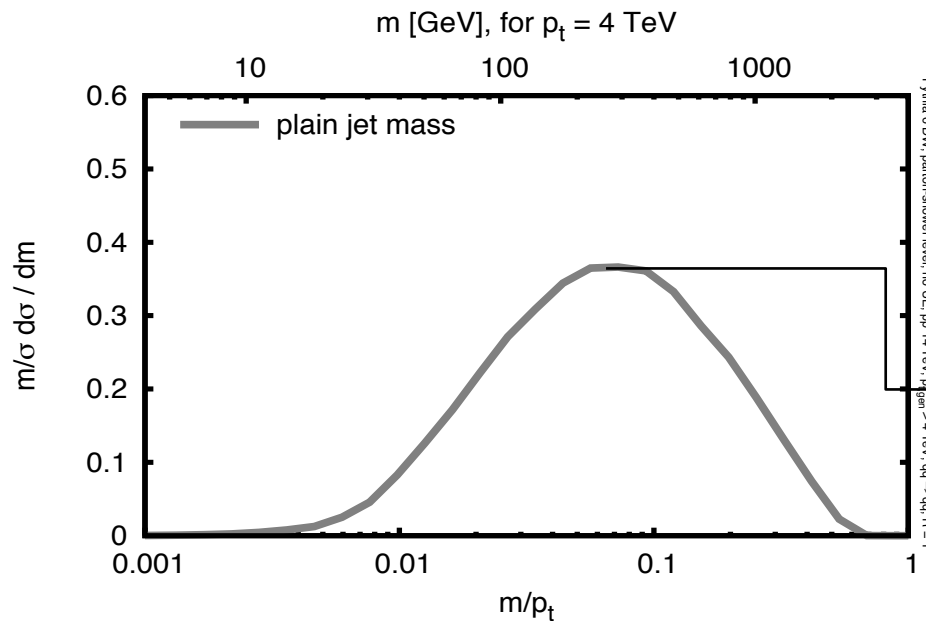


# Current understanding

Analytical studies have paved the way for a sophisticated understanding of this field.

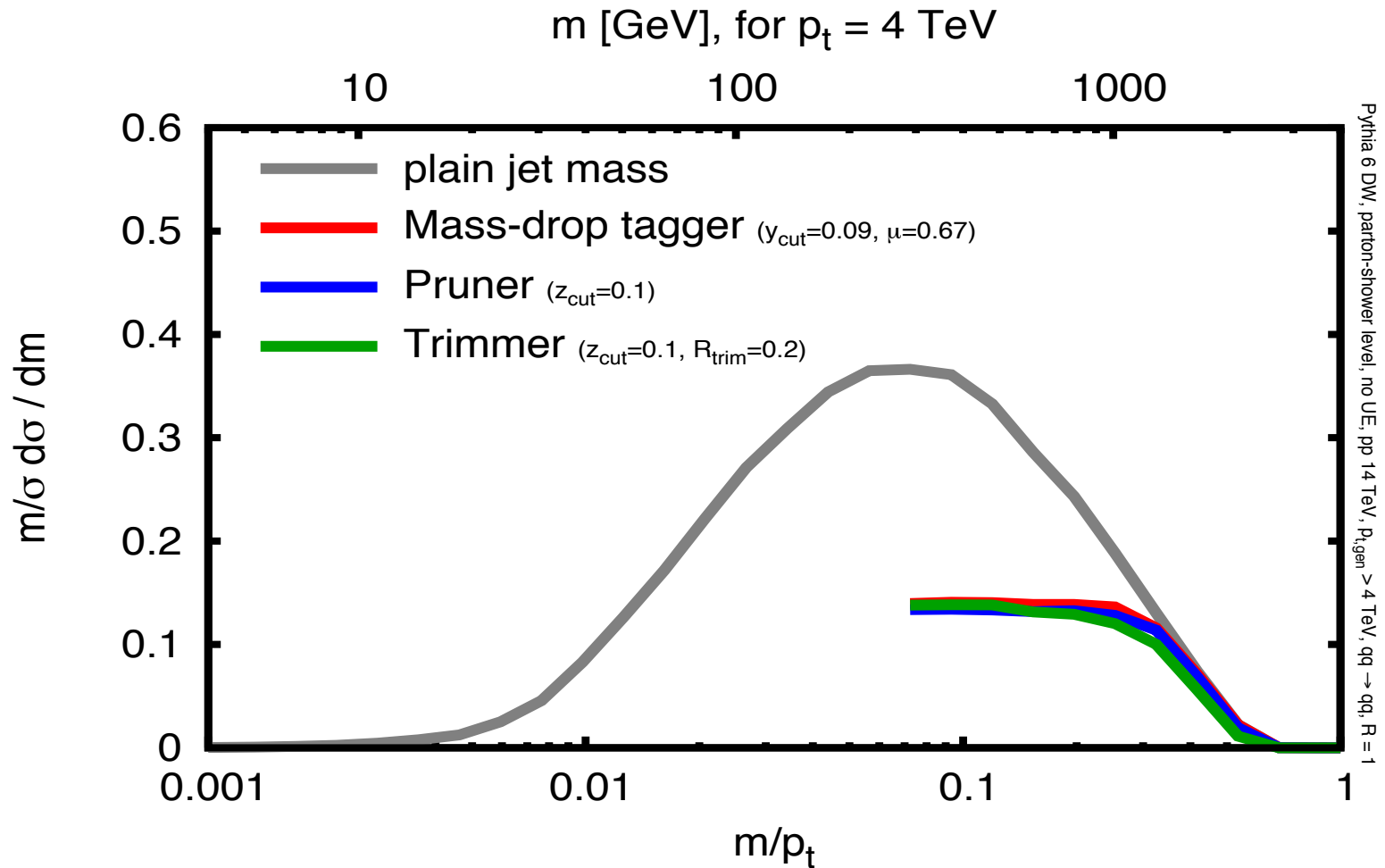
Dasgupta, Fregoso, Marzani, Salam  
2013

Post analytics it is easy to do the right MC studies

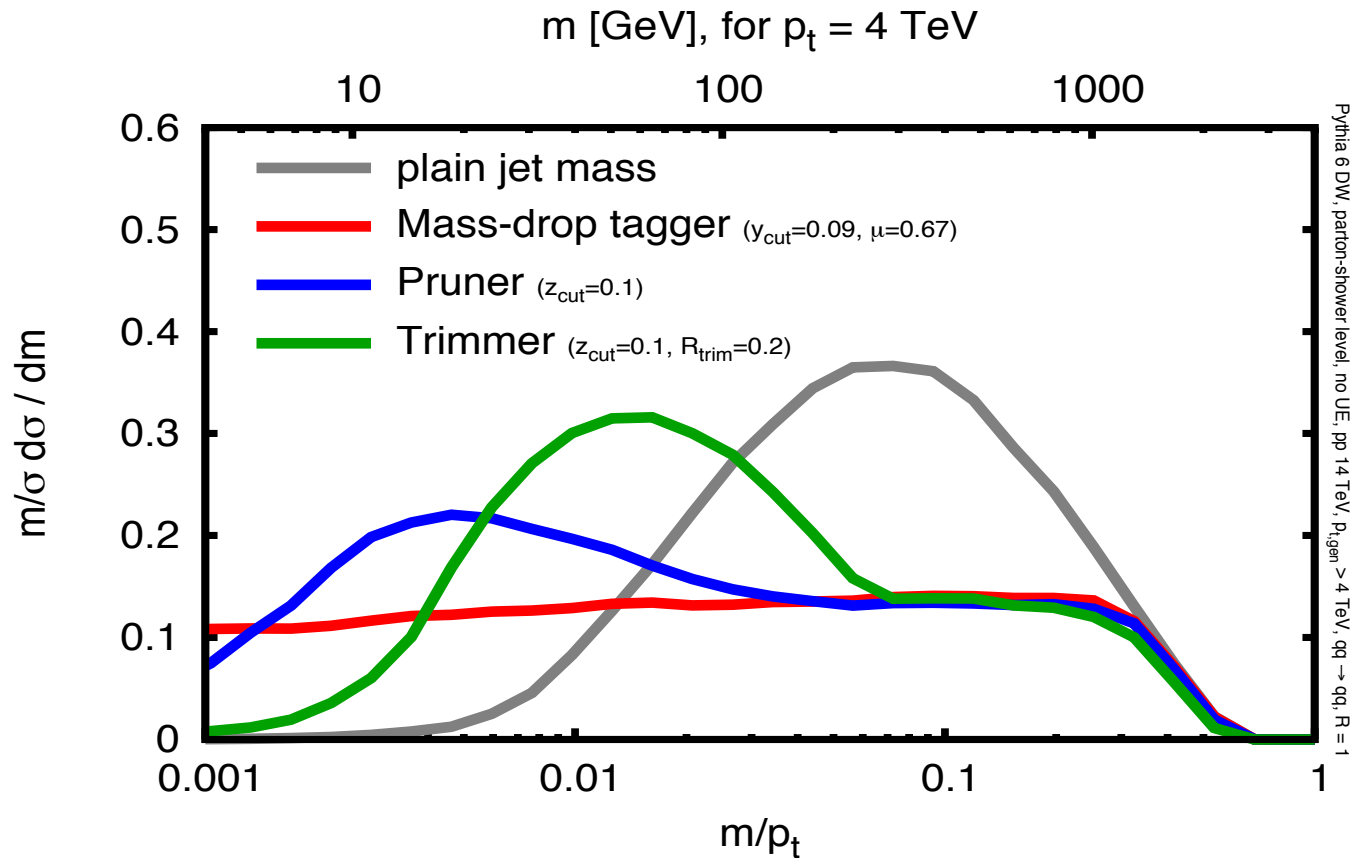


Note Sudakov peak  
in  $\sim 300$  GeV region

# Taggers look similar

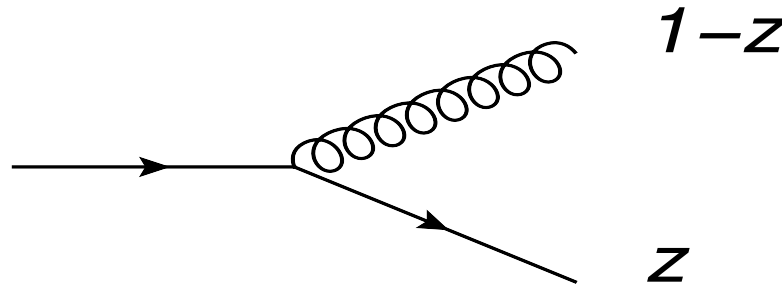


# But only over limited mass range



How do we understand what we are seeing? Positions of kinks, peaks etc.  
Needs analysis and calculation.

# Mass drop at leading order



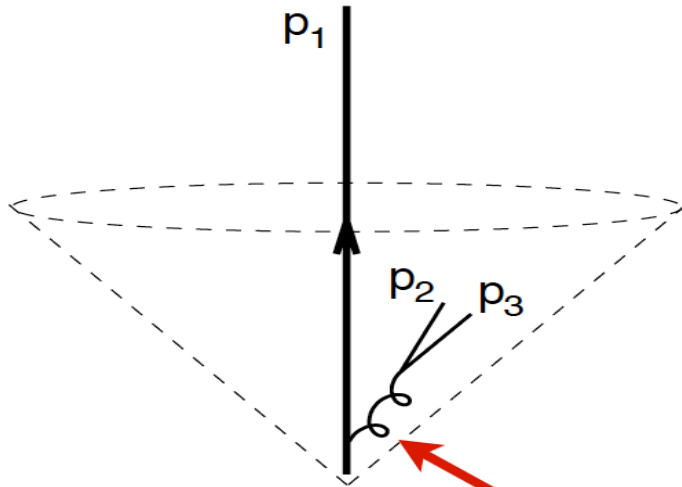
$$\frac{d\sigma}{d\rho} = \frac{C_F \alpha_s}{2\pi} \int dz \frac{d\theta^2}{\theta^2} \left( \frac{1+z^2}{1-z} \right) \delta(\rho - z(1-z)\theta^2) \Theta(z - y_{\text{cut}}) \Theta(1 - y_{\text{cut}} - z)$$

$$\sim \frac{1}{\rho} \frac{C_F \alpha_s}{\pi} \left( \ln \frac{1}{y_{\text{cut}}} - \frac{3}{4} \right) \Theta(y_{\text{cut}} - \rho) + \Theta(\rho - y_{\text{cut}}) \frac{1}{\rho} \frac{C_F \alpha_s}{\pi} \left( \ln \frac{1}{\rho} - \frac{3}{4} \right)$$

- Transition point at  $y_{\text{cut}}$
- Only single logarithmic behaviour for small jet mass/ $p_T$
- Logs have simple origin in pure collinear physics i.e. are of DGLAP type. No soft enhancements!
- We neglected terms of order  $y_{\text{cut}}$



# Beyond LO and a flaw in MDT



## What MDT does wrong beyond LO:

Follows a soft branch ( $p_2 + p_3 < y_{\text{cut}} p_{\text{jet}}$ ) with “accidental” small mass, when the “right” answer was that the (massless) hard branch had no substructure

**Subjet is soft, but has more substructure than hard subjet**

MDT's leading logs (LL, in  $\Sigma$ ) are:

$$\alpha_s L, \alpha_s^2 L^3, \dots \text{ I.e. } \alpha_s^n L^{2n-1}$$

quite complicated to evaluate

# Modified mass drop and all-orders

Modified mass drop tagger to follow **harder** rather than more massive branch. Small phenomenological effect but drastic simplification to logarithmic structure.

We performed an all-orders resummation of the jet mass distribution with mMDT.

## APPROXIMATE SQUARED MATRIX ELEMENT

$$\sum_n \frac{1}{n!} \prod_i^n \frac{d\theta_i^2}{\theta_i^2} \frac{dz_i}{z_i} \frac{\alpha_s(\theta_i z_i p_t^{\text{jet}}) C_F}{\pi}$$

**can use QED-like independent emissions, as if gluons don't split**

+ virtual corrections, essentially from unitarity

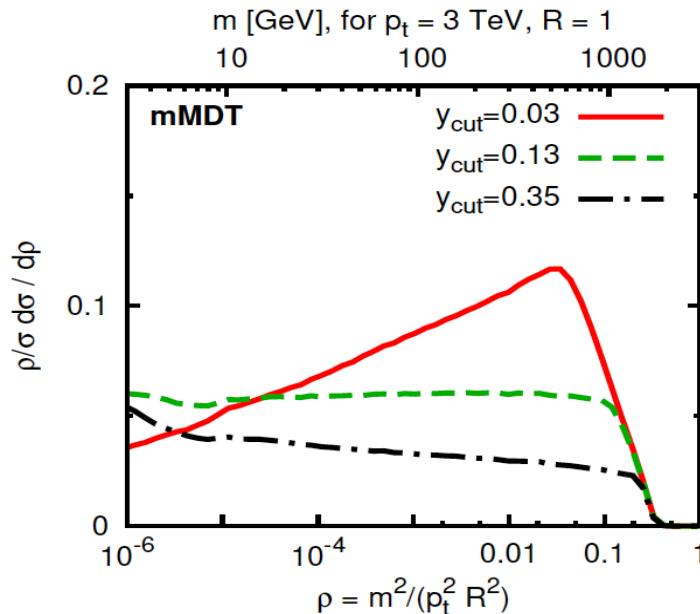
# All orders results

$$\left(\rho \frac{d\sigma}{d\rho}\right)^{\text{fixed-coupling}} = \rho \frac{\partial}{\partial \rho} \exp \left[ -C_F \frac{\alpha_s}{\pi} \left( \ln \frac{1}{y_{\text{cut}}} - \frac{3}{4} \right) \ln \frac{1}{\rho} \right] \quad \rho < y_{\text{cut}}$$

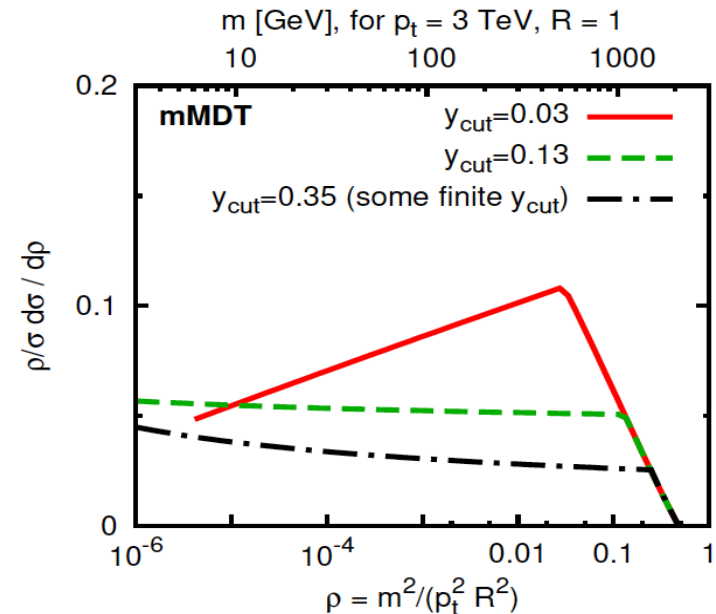
- One finds a pure collinear single logarithmic structure that exponentiates straightforwardly.
- Transition to plain jet mass at large masses.
- No soft logs or non-global logs unlike jet masses. Possible to compute this with high precision. We have computed only the leading collinear logs  $(\alpha_s L)^n$ .
- First time a jet observable of this type was ever seen.
- **No dependence on mass-drop** cut but only on asymmetry parameter.

# Comparison to MC

## Monte Carlo

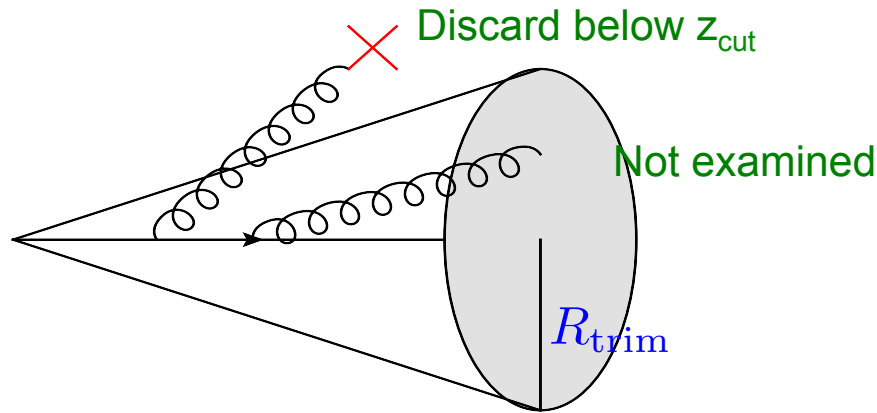


## Analytic



Excellent agreement of analytic and MC results indicate we have captured the relevant physics with our simple formulae.

# Trimming



$$\frac{\rho}{\sigma} \frac{d\sigma^{\text{trim,LO}}}{d\rho} = \frac{C_F \alpha_s}{\pi} \left[ \Theta(\rho - z_{\text{cut}}) \ln \frac{1}{\rho} + \Theta(z_{\text{cut}} - \rho) \ln \frac{1}{z_{\text{cut}}} - \frac{3}{4} + \Theta(z_{\text{cut}} r^2 - \rho) \ln \frac{z_{\text{cut}} r^2}{\rho} \right]$$

Plain jet mass  
result

Single log  
mMDT like  
region

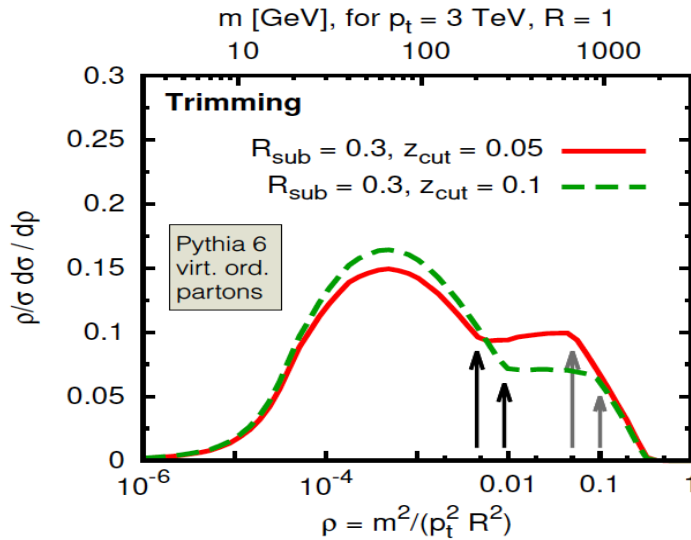
Double  
logarithmic  
behaviour

$$r = \frac{R_{\text{trim}}}{R}$$

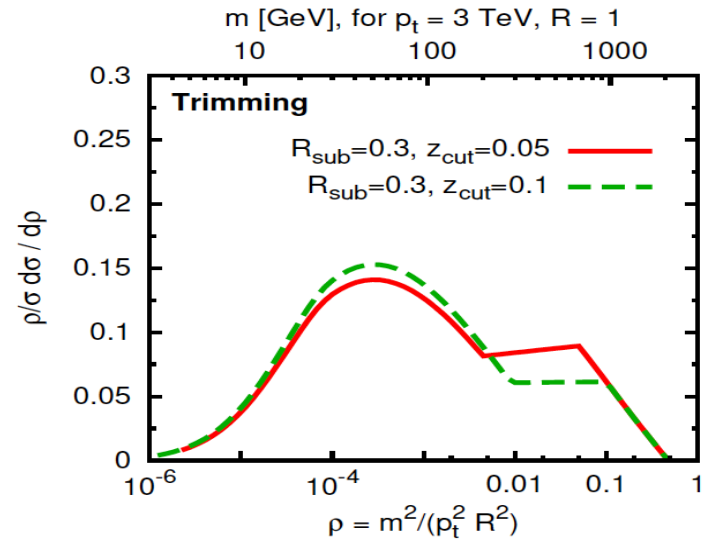
Three transition points seen

# All order result and MC comparison

## Monte Carlo



## Analytic



**Non-trivial agreement!**  
(also for dependence on parameters)

$$\frac{d\sigma^{\text{trim,resum}}}{d\rho} = \frac{d\sigma^{\text{trim,LO}}}{d\rho} \exp \left[ - \int_{\rho}^1 d\rho' \frac{1}{\sigma} \frac{d\sigma^{\text{trim,LO}}}{d\rho'} \right]$$

# Pruning results

Recall that pruning is like trimming but with a dynamical radius  $R_{\text{prune}} \sim m/p_t$ .

LO result is single logarithmic like (m)MDT.

$$\rho \frac{d\sigma}{d\rho} \sim \alpha_s \ln \frac{1}{z_{\text{cut}}}$$

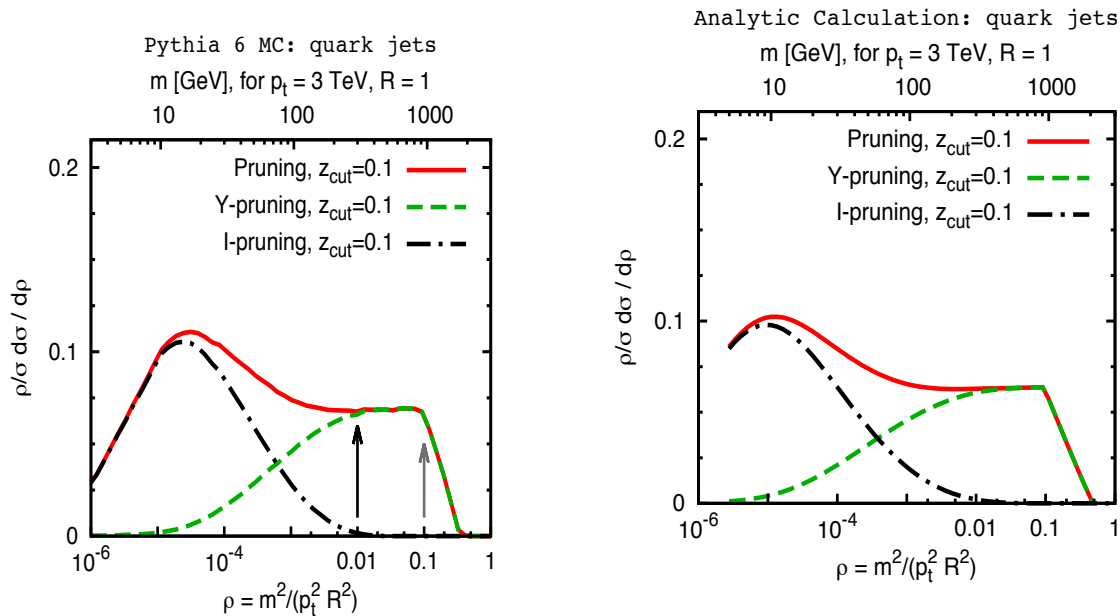
However at NLO one encounters terms as singular as the plain jet mass i.e. double logarithms.

$$\rho \frac{d\sigma}{d\rho} \sim \alpha_s^2 \ln^3 \frac{1}{\rho}$$

It turns out that pruning is a sum of two components only one of which is sane. We initially called the other component “anomalous”.

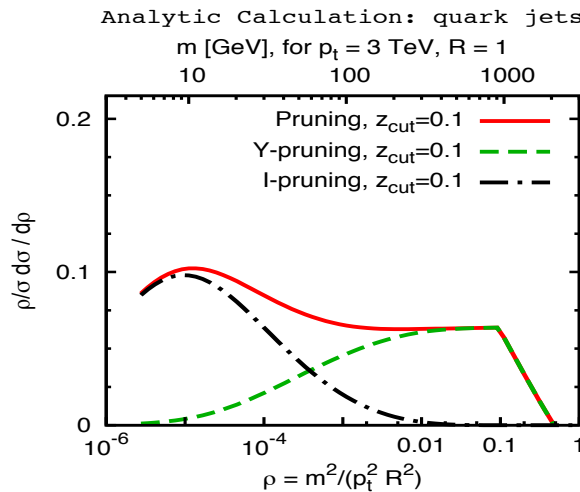
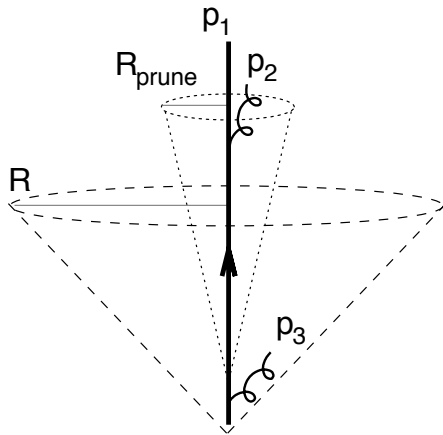


# Pruning MC v analytics



The black line denotes the anomalous component (I-pruning). The green line is the sane component (Y-pruning).

# A new tagger – Y pruning



$$\frac{\rho}{\sigma} \frac{d\sigma^{\text{Y-prune}}}{d\rho} \sim \frac{C_F \alpha_s}{\pi} \left( \ln \frac{1}{z_{\text{cut}}} \right) \exp \left( -\frac{C_F \alpha_s}{2\pi} \ln^2 \rho \right)$$

I pruning eliminated by demanding that at least one emission passes pruning.

# Non-perturbative effects

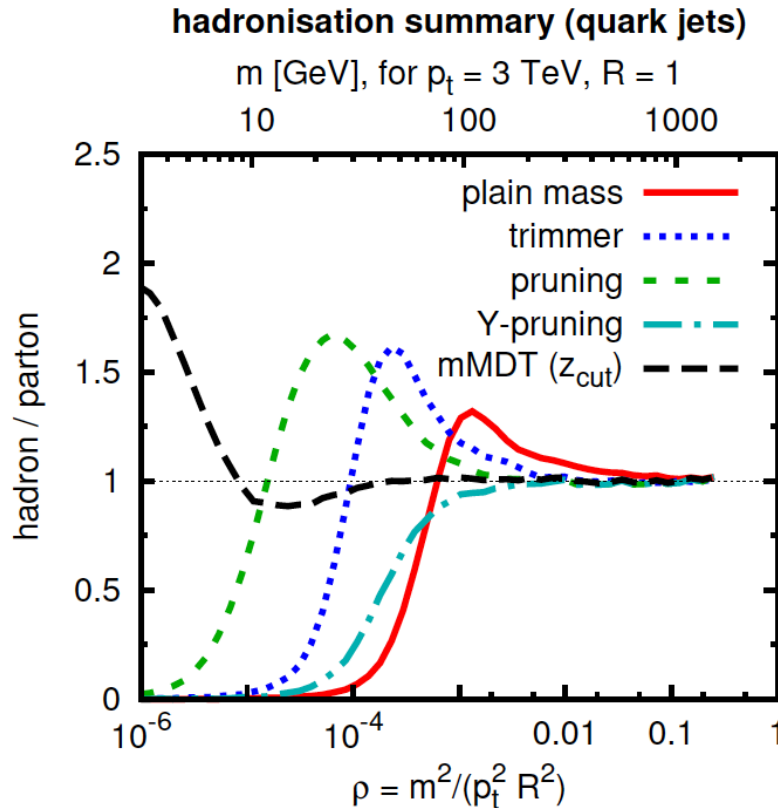
Do we really need to worry about these on TeV scale jets?

Consider the fact that a 1 GeV hadron can produce a squared jet mass  $M_j^2 = 1\text{GeV} \times R^2 p_T$  which for a 3 TeV jet leads to a mass of 55 GeV quite close to the electroweak scale!

Need to worry about both hadronisation and the Underlying Event (UE, radiation uncorrelated with the hard process)

The most common way of studying these is via Monte Carlo though analytical models for hadronisation are common and successful.

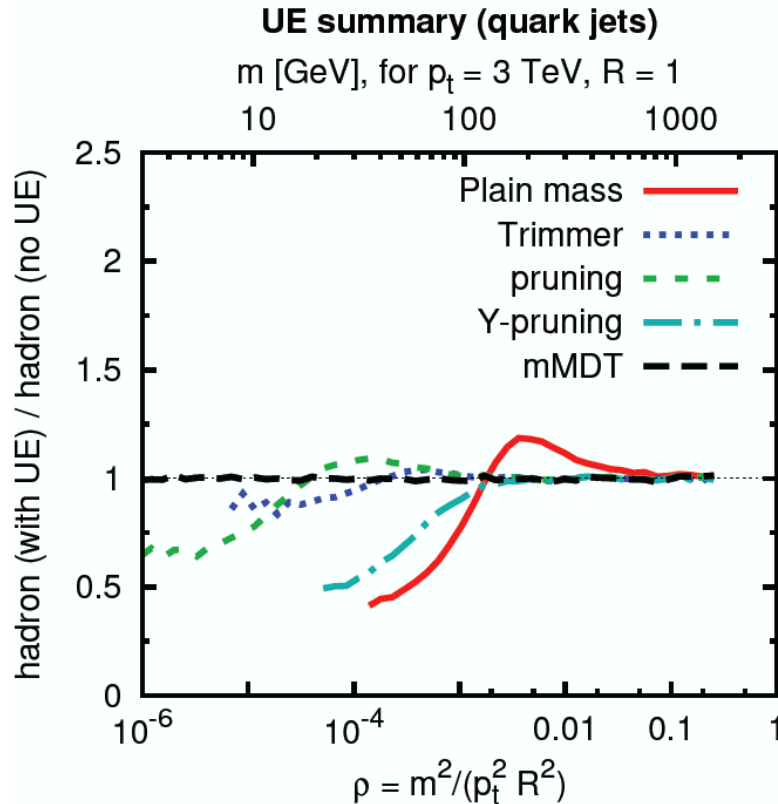
# Hadronisation



Nearly all taggers have  
large hadronisation  
effects:

15 – 60%  
for  $m = 30 - 100$  GeV

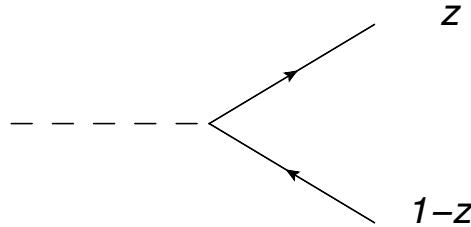
# Underlying event



Underlying event impact  
much reduced relative to  
jet mass

Almost zero for mMDT  
(this depends on jet  $p_t$ )

# Signal jets



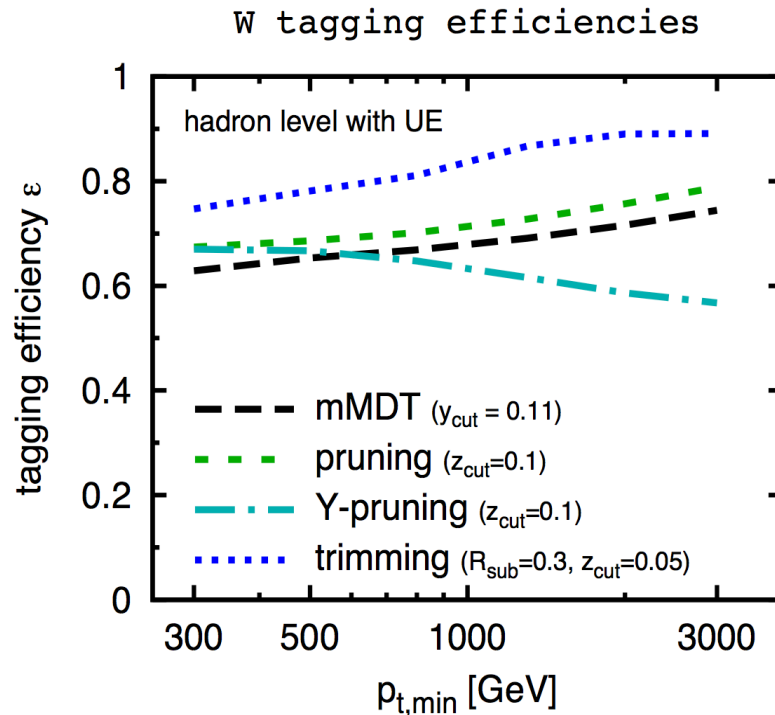
The action of taggers on signal jets reveals less surprises than the case of QCD backgrounds. Basic LO result for mMDT and pruning for Higgs decays:

$$\epsilon_s = \int_{y_{\text{cut}}}^{1-y_{\text{cut}}} dz = 1 - 2y_{\text{cut}}$$

Trimming has a more involved structure even at LO

$$(1 - 2y)\Theta(1 - 2y) + \sqrt{1 - \frac{4\Delta}{r_{\text{trim}}^2}}\Theta\left(\frac{1}{4} - \frac{\Delta}{r_{\text{trim}}^2}\right)\Theta\left(y - \frac{1}{2}\right) + \\ \left(2y - 1 + \sqrt{1 - \frac{4\Delta}{r_{\text{trim}}^2}}\right)\Theta\left(\frac{1}{4} - \frac{\Delta}{r_{\text{trim}}^2}\right)\Theta\left(\frac{1}{2} - y\right)\Theta\left(y - \frac{1}{2}\sqrt{1 - \frac{4\Delta}{r_{\text{trim}}^2}}\right)$$

# Signal efficiencies with taggers



Tree level is a good approximation with small effects from ISR and FSR effects.

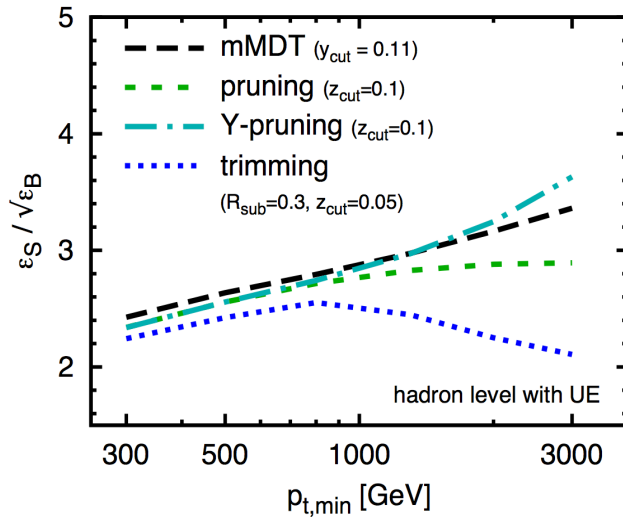
Y-pruning suffers a loss of efficiency at high  $p_t$

All this also understood analytically and with MC studies.

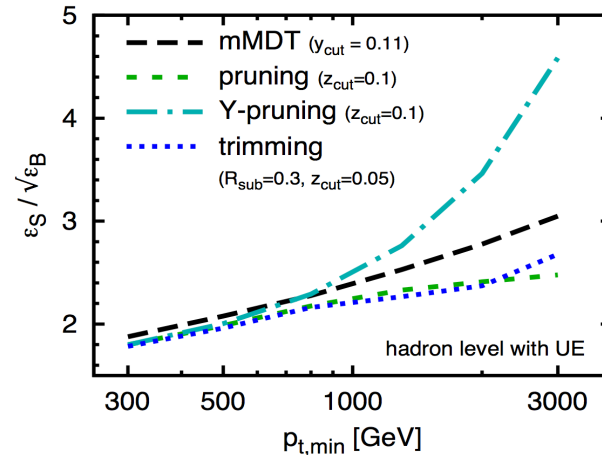


# So which is the best tagger?

signal significance with quark bkgds



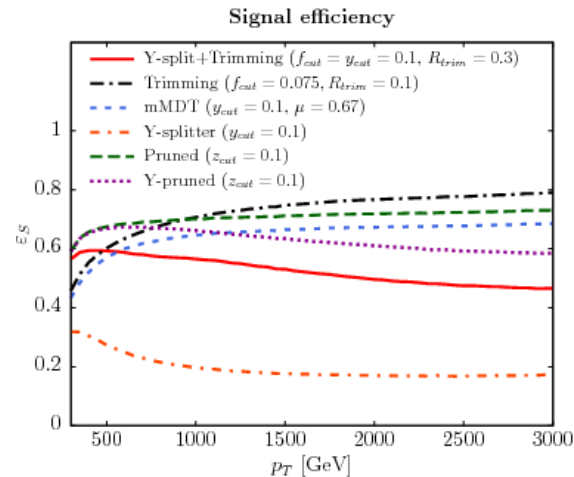
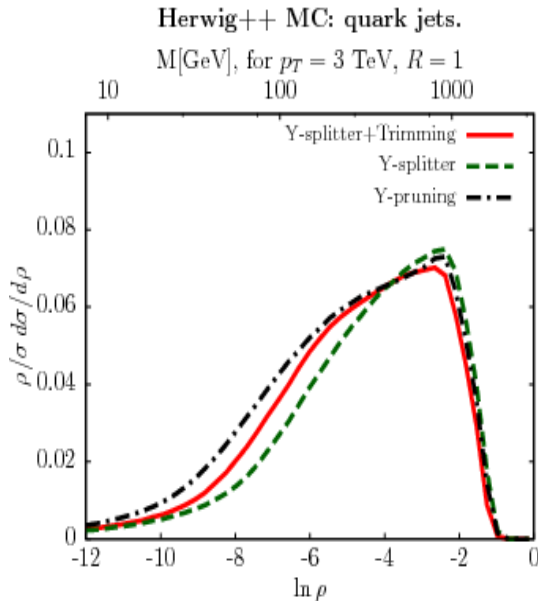
signal significance with gluon bkgds



mMDT has some nice features. Simple analytical structure so precision calculations and phenomenology possible. Closest to being a scale invariant tagger. Only one transition point etc. Good for QCD phenomenology and robust for data driven background estimates in searches.

However the **Sudakov suppression** of background in Y-pruning gives it the best signal to square root of background ratio amongst the tools studied here.

# Can we do better?



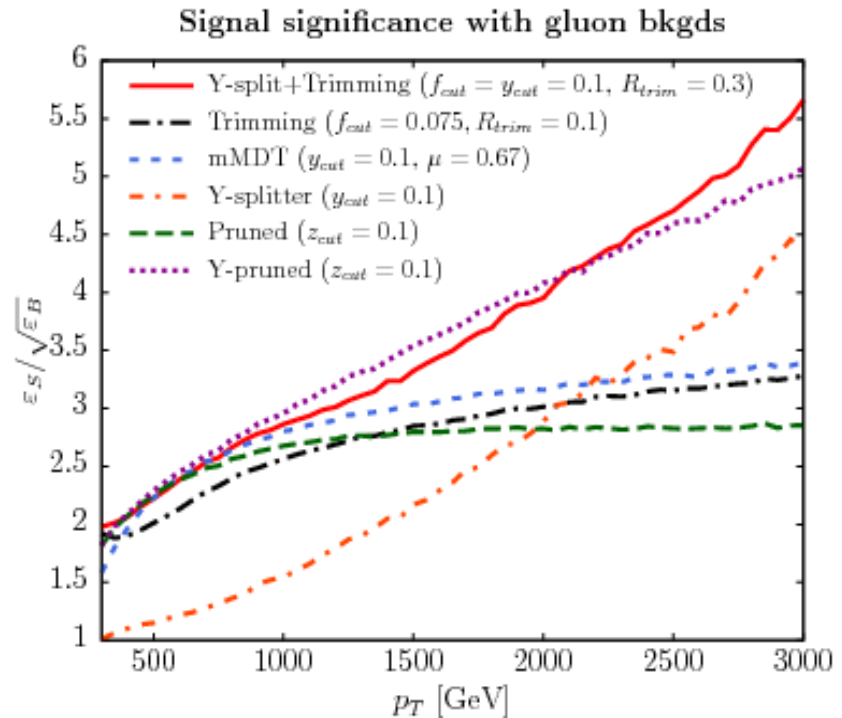
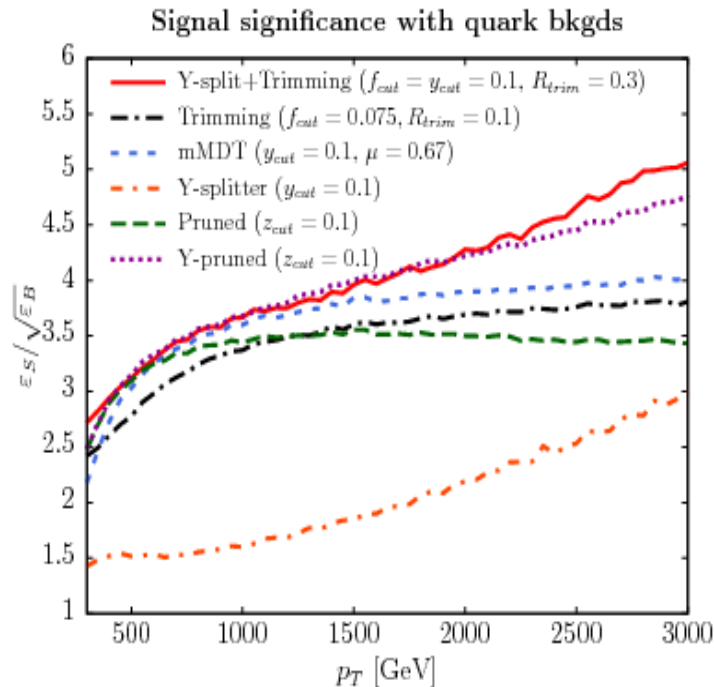
A systematic understanding of substructure can be used to create efficient, robust and high-performance tools.

We applied such an understanding to Y-splitter which was not commonly used. Discovered it has **excellent background rejection but poor signal response**.

Butterworth, Cox and Forshaw 1995

We considered its combination with a groomer such as trimming and found that it improves signal while not modifying background rejection much.

# Y-splitter+trimming



The combination of Y-splitter with trimming outperforms other taggers at high  $p_t$ . Also understood analytically.

Dasgupta, Powling, Siodmok, Soye, Sarem-Schunk in progress

# Summary

- An analytical understanding of jet substructure is possible and significant progress has been made.
- There are several tools developed so far and we have only examined some of them
- Radiation constraining jet shapes such as N-subjettiness and energy correlation functions have also been studied analytically with some success. Dasgupta, Soyez, Sarem-Schunk 2015

The quest for the best tools for LHC run-2 and beyond continues. The hope is that analytical understanding has put this field on much firmer ground than before.