

b-Tagging (at multi-purpose LHC experiments)





Christophe Saout

CERN, Karlsruhe Institute of Technology (KIT)

- Introduction
- Track based algorithms
- Lepton based algorithms
- Vertex based algorithms
- Combined algorithms
- Mistag/Efficiency measurements
- Tracker misalignment
- Conclusions





Why b-Tagging?

- b physics
 - b production (cross-sections, ...)
 - b decays (mixing, oscillations, CP violation, CKM, ...)
- Searches: b-tagging is of crucial importance for physics analyses that need to identify jets from heavy quark flavours (b, top), e.g. top, SUSY, Higgs
 - \rightarrow background suppression!

Where can b-Tagging be applied?

- Offline (what we are focussing on)
 - Full detector information
 - CPU-intense and precise reconstruction
 - Possibility to use very sophisticated algorithms
- Trigger
 - Limited information
 - Timing constraints \rightarrow needs simple and fast algorithms





An exemplary "big picture": pp $\rightarrow t\bar{t}$



collision \rightarrow heavy resonance production (here: via QCD)







heavy resonance decay \rightarrow b quark final states







Parton Shower: QCD radiation, gluon splitting, ...







Hadronization: *jets of particles*







What we see in the detector: jets of hadrons

CMS





heavy hadron decays: this is what allows us to do b-tagging



b Decays





b-hadron decays are weak decays

- only allowed through family change
- Strongly suppressed by CKM matrix \rightarrow long lifetime!
 - \rightarrow b \rightarrow c cascade decays preferred
- O(20%) of the decays occur semi-
 - \rightarrow lepton "in jet"
- ... but so are c-hadron decays...

heavy hadron decays: this is what allows us to do b-tagging



b Decays



- **b**-quarks significantly differ from light flavour quarks by:
 - mass: m = 4.2 GeV
 - lifetime: $\tau \approx 1.5 \text{ ps} \rightarrow \sim 1.8 \text{mm}$ (at 20 GeV) before decay
 - decay: weak, mostly into c-quarks ($\rightarrow 3^{rd}$ decay) $\rightarrow 20\%$ into leptons
- Signatures:
 - tracks: high decay multiplicity, significant displacement (~5 charged tracks on average)
 - secondary vertices (SV): tracks intersecting at a common vertex
 - b fragmentation: large portion of jet energy carried by b-hadron
 - \rightarrow need good tracking resolution at impact point! (\rightarrow alignment!)





Limitations & Backgrounds



Track reconstruction

 \rightarrow key role in b-tagging!

- impact parameter resolution
- measurement errors
- tracking in jets difficult
- Efficiency / minimum p_T limitations
- can only reconstruct charged hadrons
- bad / fake tracks (pattern recognition / combinatorics)
- Real lifetime:
 - c-jets
 - Long-lived decays: κ_s^{0} , Λ^{0} , hyperons
- Leptons
 - Muon identification (punch-through)
 - Real muons (e.g. muons from in flight π decays)
 - electron ID in jets difficult (no calorimeter isolation)
- Other
 - Pile-up \rightarrow more than one collision \rightarrow need z measurement!

 \rightarrow 2D pixel detector

- "jet overlap" \rightarrow neighboring tracks isolation criteria
- QCD: $g \rightarrow b\overline{b} / g \rightarrow c\overline{c}$ splitting



Tracking Detector







Impact Parameters





Impact Parameter Resolution





(multiple scattering, nuclear interactions, ...)

But also: *Tracking detector alignment!*

CMS



Impact Parameters





Simple I.P. based Algorithm



- Compute Impact Parameters for all tracks in jet
- Sort tracks by descending Signed IP Significances (2D or 3D)
- Select nth track (hence called "Track Counting" in CMS)
 - 2^{rd} track \rightarrow "high efficiency" tag
 - 3^{rd} track \rightarrow "high purity" tag

CMS

- Use IP significance as discriminator
- Simple, fast \rightarrow suitable for online trigger



- Fake tracks
- "V" decays









Choosing the optimal "n"-th track. → trade-off between background reduction b-tagging efficiency sacrifice

CMS

Simple I.P. Algo Performance



We can interprete the discriminator distributions with MC truth to determine working points and performance:

- We choose a working point, e.g. discr cut = 2.5
- We call a jet, which has a discr > 2.5 tagged

CMS

- The percentage of b-jets we correctly tag, is our efficiency
- The percentage of non-b jets we tag, is our mistag rate
- We can scan all working points and collect efficiencies and mistag rates



Track Probabilities



Method originally from LEP (ALEPH):

• For each *track*, determine the *probability* that it comes from the PV

 $\rightarrow P(S') = \int_{S'} pdf(S) dS$ (pdf being the normalized IP Sig. distr.)

- Can use the measured negative signed IP distribution for calibration!
- Use information from all tracks by combining the probabilities

$$P_{N} = \pi \cdot \sum_{j=1}^{N-1} \frac{-\log \pi}{j!} \quad \text{with} \quad \pi = \prod_{i=1}^{N} P(S^{i}) \quad \text{and using} \quad discr = -\log P_{N}^{+}$$

Alternative: Likelihood Ratio (using both b and light flavour PDFs) (ATLAS IPxD)





Soft Lepton b-Tagging



- As seen earlier, the branching ratio $b \rightarrow l$ ($l = e, \mu$) is of the order of 20%
- This allows to use reconstructed leptons in jets as a complementary way of doing b-tagging
- This is particularly useful for b-tagging efficiency measurement methods
- Or can, combined with lifetime-based methods, improve efficiency
- Leptons from b-decays are *soft* (as apposed to leptons from hard process)
- They are within the jet and non-isolated (typically $\Delta R < 0.4$)
- The b-quark mass gives them a significant transverse momentum with respect to the jet direction (called " p_T^{rel} ")
- As for other tracks from the b-hadron, the tracks are displaced
- Reconstruction of electrons in jets is particularly tricky (Bremsstrahlung, calorimeter-track matching, calorimeter isolation, ...)
- Background: photon conversion, decays in flight



Soft Lepton b-Tagging



More complex alternatives:

- Combine multiple variables using MVA techniques:
 - Pt rel
 - IP significance
 - Other kinematic variables (angle, momentum ratios)
 - or lepton ID discriminator

• ...



For the two "simple" CMS soft muon taggers





- Requiring the reconstruction of secondary vertices is an excellent indication for a heavy particle decay:
 - \rightarrow not only requiring displaced tracks, but also
 - \rightarrow displaced tracks intersecting at a common vertex
- Different approaches to vertex finding
 - Repeated fitting with outlier removal
 - Finding track pairs that form secondary vertices
 - Seeding vertices from compatible with a b-hadron direction hypothesis





Secondary Vertices



- The reconstruction efficiency is strongly depending on the jet kinematics
 - low energy tracks \rightarrow strongly reduced impact parameter resolution
 - low energy jet \rightarrow jet is wieder \rightarrow jet / track association efficiency loss
 - high energy jet \rightarrow high track density \rightarrow track reconstruction problems \rightarrow bad / fake tracks



calorimeter cone jets with $\Delta R \le 0.5$

tracks associated to jet if $\Delta R \le 0.3$





- The vertex position can be used to probe the lifetime:
 - The flight distance: $D^{\mathfrak{D}} = |\overrightarrow{SV} \overrightarrow{PV}|$ (in 2D or 3D)
 - Similarly to the impact parameters: flight distance significance
 - (Can also define a flight distance sign and use vertices behind the primary vertex to measure light flavour fake vertex distributions)





Note that the error on the vertex position longitudinally along the jet is usually very large





- In addition to the existence one or more secondary vertices, one gets an additional amount of observables: One...
 - can obtain the charged decay multiplicity
 - can compute an invariant vertex mass (from charged tracks only)
 - can probe the fragmentation function: energy fraction carried by vertex





Secondary Vertex b-Tagging



- The fact that a secondary vertex has been reconstructed can already be used as a b-tag (e.g. CDF at Tevatron)
- One can define different working points by tightening requirements
 - minimum nuber of tracks
 - track qualitites
 - vertex fit quality
 - vertex flight distance significance \rightarrow e.g. as discriminator to cut on:





ATLAS "Jet Fitter"



First presented by SLD under the acronym "Ghost Track" b-tagging:

 Instead of finding vertices directly, first attempt to reconstruct the b-hadron flight line (takes into account b-c decay chain)



• This is done by first fitting a straight line ("the ghost track") that minimizes the distance to all track:



- Then tracks can be clustered along the "ghost track" to form vertices
 - \rightarrow b-c decay chains can be reconstruced, including "1-track vertices" decay topology as additional observable \rightarrow higher efficiency/purity



Combined Algorithms



- Instead of using individual observables \rightarrow combine variables
- "track probability" such an example (uses background hypothesis only)
- More common: Growing family of "Multivariate Analysis Techniques"
- e.g. Likelihood Ratio or more advanced: Neural Networks, Boosted Decision Trees, ...
 - → machine-learning (from signal & background samples)



In 3 categories: IP only, IP + PseudoVertex, IP + SV \rightarrow whole b-tag eff. range covered

Mistag Measurement (on Data)



Efficiency/Mistag measurements on data are much more complicated Also, unfortunately, there are no simple "tag & probe" type methods

CMS

Light flavour jet observables approximately symmetric around zero \rightarrow sign of b-tagging observables (signed IP, signed distance) \rightarrow mistag measurement using negative side

→ Idea: flip the sign and you get a b-tagging algorithm that will behave as if the jet was always a light flavour jet (in principle works for IP + SV observables and even combined taggers)





Efficiency Measurement



muon-jet

(tag jet)





Efficiency Measurement



You can overcome these efficiencies by solving for both tag & probe taggers simulatenously

- \rightarrow System of 8 equations
- Only slight MC dependency: "probe" / "tag" tagger correlations

 \rightarrow use uncorrelated taggers (e.g. soft muon & lifetime)





System 8 method

(developed at D0)

Solutions:
$$(\epsilon_{b}^{tag}), \epsilon_{cl}^{tag}, \epsilon_{b}^{\mu}, \epsilon_{cl}^{\mu}, n_{b}, n_{cl}, p_{b}, p_{cl}$$

SLT efficiency







Distributions for the 2^{rd} -highest signed impact parameter in jets: (\rightarrow "track counting" b-tag algorithm)

error on the IP measurement dominated by pixel hit resolution, extrapolation from innermost hits







Effect on Impact Parameters



The errors are chosen such that the pull distributions have a width of 1!

With the old "Startup 2008" scenario b-tagging would have been useless!







Effect on Secondary Vertex Reconstruction

Since essentially, candidate tracks for SV fit are those incompatible with the PV *(related to significance)*, the SV finding efficiency also decreases with the tracker misalignment:

	Secondary vertex fraction [%]		
Misalignment scenario	b-jets	c-jets	udsg-jets
No Misalignment	62.6	22.0	2.7
100 pb ⁻¹ Misalignment	62.1	19.6	2.4
10 pb ⁻¹ Misalignment	53.0	12.7	2.9
10 pb^{-1} Pixel L1 Off Misal.	39.2	7.6	1.9
Startup Misalignment	37.8	7.7	3.5





2rd -highest IP significance

non b-jet efficiency Misalignment scenario: 10⁻³ None ---- 100 pb-1 10 pb-1 Startup 10 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 b-jet efficiency "Jet Probability": Combination of all "Track Probabilities" (per-track IP sig. pdf's)

udsg iets. JetProbability

CMS Preliminary

data points represent light flavour mistag vs. b-tagging efficiency for different working points (discriminator cuts)





General observation: The more complicated an algorithm (and the more efficient), the more sensitive it is to misalignment



The simple "Track Counting" and "Simple SV" algorithm are presumably easiest To get under control with early data. They do not need any "training" on MC. \rightarrow with 10pb⁻¹ of data, b-tagging will already be usable



Conclusions



- Many exciting discovery channels, as well as well-known SM channels have b-jets in the final state $(tt \rightarrow bWbW, H \rightarrow bb, SUSY, ...)$
- Both CMS and ATLAS have high-resolution state-of the art tracking detectors, with pixels allowing for impact parameter measurements in the r- ϕ and z plane
- Both CMS and ATLAS have an exhaustive set of b-tagging algorithms in place, both suitable for the startup phase and the long-term discovery phase, as well as for online triggers
- Methods and infrastructure to measure tagging performance and mistag rates on data are in place
- Track reconstruction is key ingredient
 → will be challenging to keep up with alignment progress!
- Given the current state of alignment from cosmics b-tagging should be usable very early
- A lot of work and exciting first years ahead of us!