

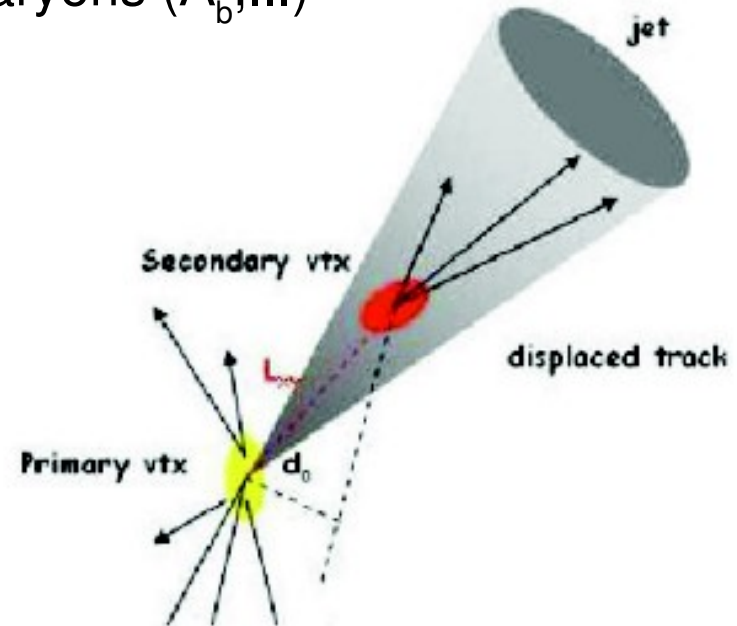
An overview on the b-tagging landscape

algorithms and efficiencies

Joris Maes

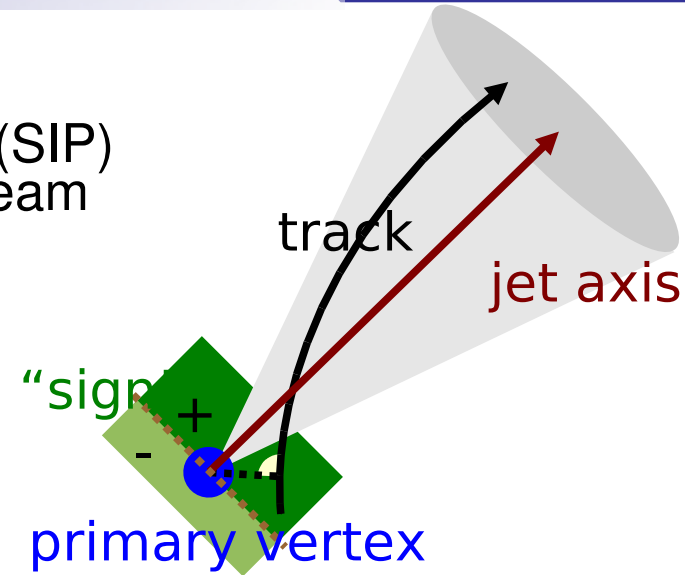
IIHE top quark meeting
24th November 2008

- In different interesting processes at the LHC **b quarks** play an important role
 - $H \rightarrow b\bar{b}$,
 - $t \rightarrow Wb$ for about 100%
- **Identifying jets coming from b quarks** (so called b jets) can help to reduce background in complex final state topologies.
- The **difference** of a jet from a b/c quark w.r.t uds- or gluon-jets can be used to identify them:
 - During fragmentation of heavy flavour quarks (b/c) jets are formed containing heavy mesons (B,D,...) and baryons (Λ_b, \dots)
 - the b hadron has a long lifetime τ of 1.6 ps corresponding to $c\tau$ 490 μm
 - this long lifetime is reflected in the presence of **tracks not compatible with the primary vertex** (displaced secondary vertex)
 - b jets contain (19%) a **non isolated e^\mp/μ^\mp** which has different properties w.r.t. other leptons formed inside jets



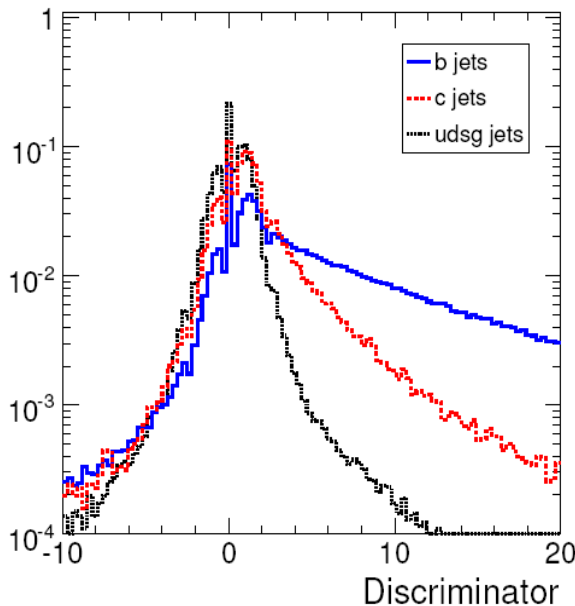
Impact parameter taggers

- Impact parameter IP = smallest distance of a track and the considered primary vertex
- One can assign a **sign** to the impact parameter (SIP) whether the track was produced up- or downstream
- To take into account the experimental resolutions the **significance of the impact parameter** is used = $IP/\sigma(IP)$
- This algorithm needs tracks and primary vertex with some additional quality cuts on the objects

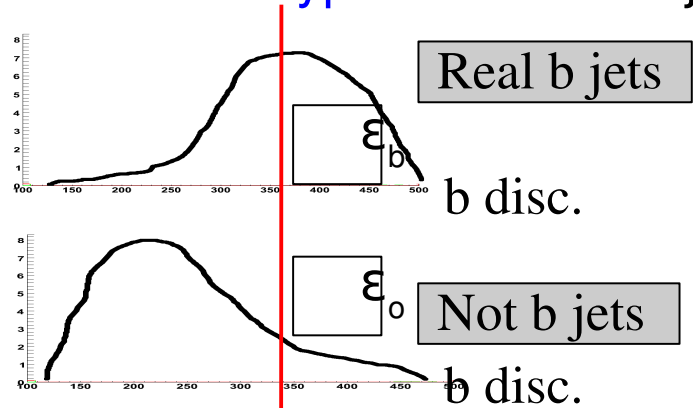


- Algorithms using this info

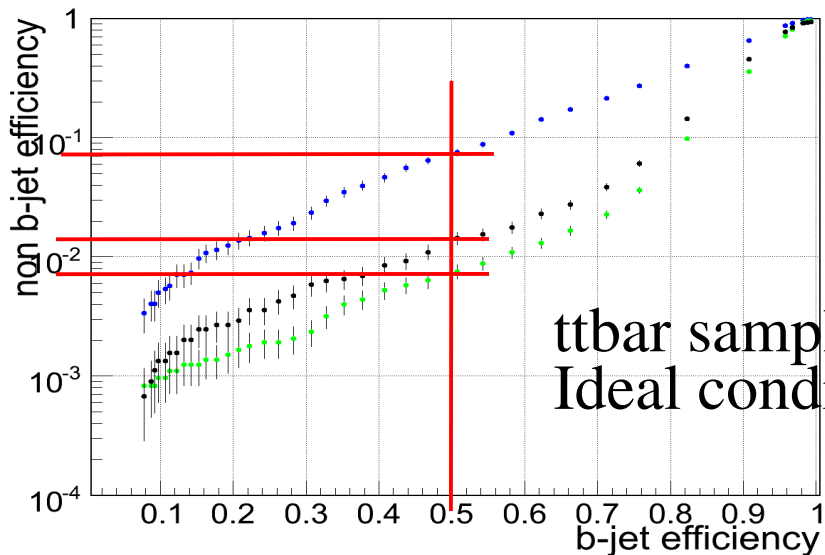
- TrackCountingHighEff: SIP^{3d} of 2nd track (HLT)
- TrackCountingHighPur: SIP^{3d} of 3rd track
- JetProbability: Likelihood that each track comes from a b or light quark, combined for each track in the jet
- JetBProbability: Only use first 4 tracks
- ImpactParameterMVA: Combine all track Ips via MVA



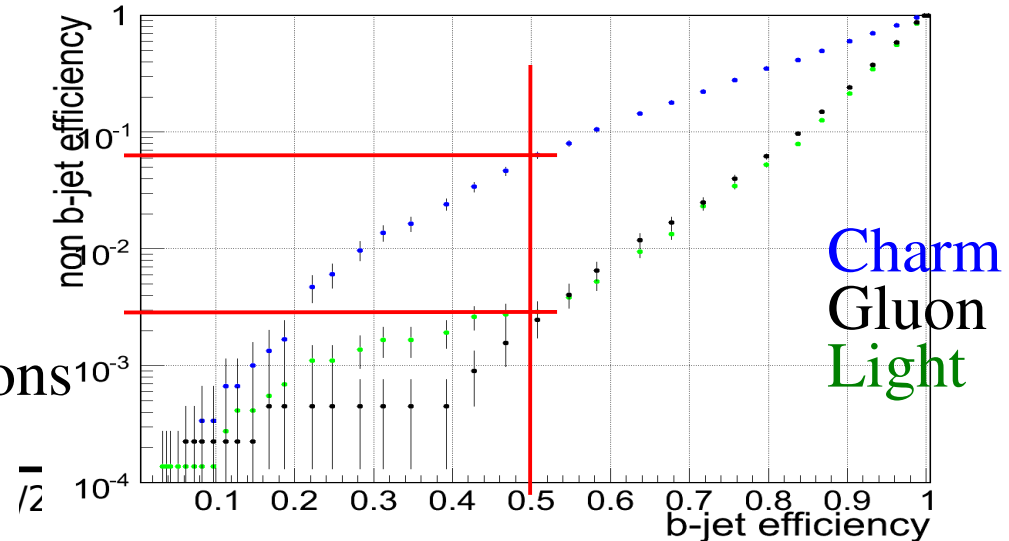
- The algorithms give a b discriminant variable to **test the hypothesis** that a jet is originating of a b quark.
- Two mistakes can be made:
 - accept the hypothesis when it is wrong
 - reject it while it was correct.
- This is reflected in two efficiencies
- ϵ_b : efficiency to tag a b jet as a b jet
- ϵ_o : efficiency to tag an 'other' (udscg) jet as a b jet = **mistag efficiency/rate**
- Interpretation: typically 50% b tag efficiency gives about a 10% c-mistag rate and a 1% udsg mistag rate



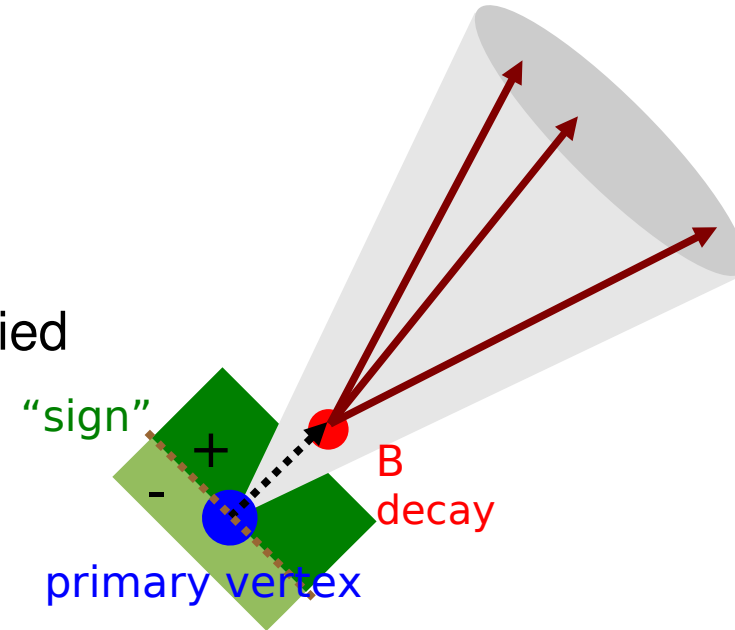
trackCountingHighEff: mistag vs. b tag efficiency



jetProbability: mistag vs. b tag efficiency



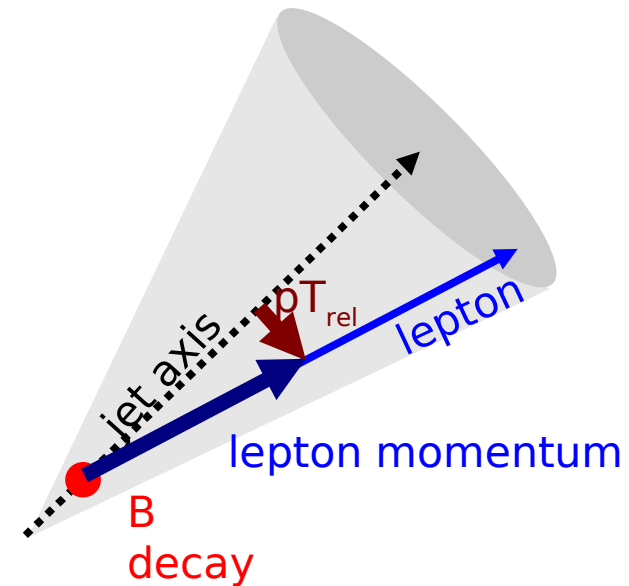
- Besides the primary vertex the **secondary vertex from the B decay** is reconstructed. This one is used as extra input w.r.t impact parameter taggers
- Advantage: in general a lower mistag rate but more input is needed and can only be applied with **enough detector knowledge**



- algorithms

- `SimpleSecondaryVertex`: Discriminate using the (signed) significance 3D flight distance
- `CombinedSecondaryVertex`: likelihood to come from a b, c or light quark. Ratio of combined likelihoods as b-tag discriminator
- `CombinedSecondaryVertexMVA`: neural network instead of Likelihood

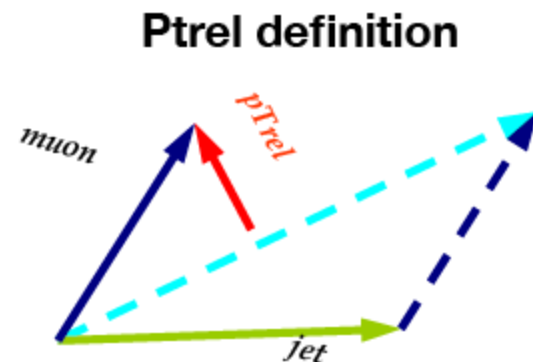
- The branching fraction of direct and **cascade decay of b hadrons** in electrons and muons is about 19% due to $b \rightarrow W^{*-}X$ or $b \rightarrow W^{*-}c \rightarrow W^{*-}W^{+*}X$ where $W^{*-} \rightarrow l\nu$
- This is exploited by looking for a **non isolated lepton in the jet**
- These algorithms are **less correlated** w.r.t track impact and secondary vertex algorithms
- Input used are muons or electrons, primary vertex and tracks to improve jet direction
- Algorithms:
 - (simple) `SoftMuonByPt`: requires a muon with $p_{Trel} > \text{cut}$ (HLT)
 - `SoftMuonByIP3d`: requires a muon with $SIP3d > \text{cut}$
 - `SoftMuon`, `SoftElectron`: Combine lepton kinematic informations via a NN
 - `SoftMuonNoIP`: NN without SIP^{3d}



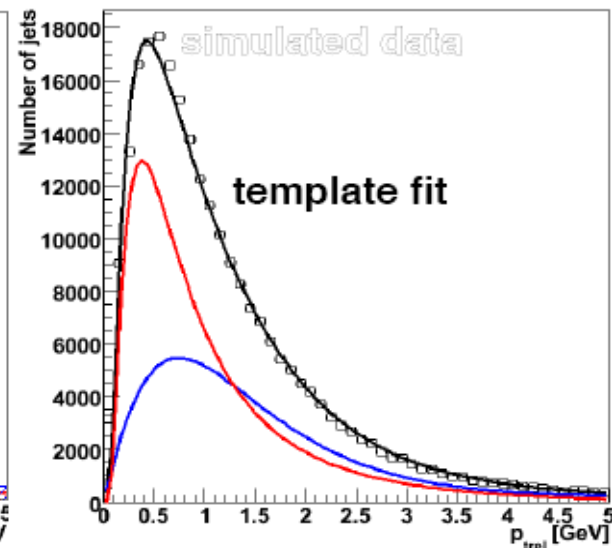
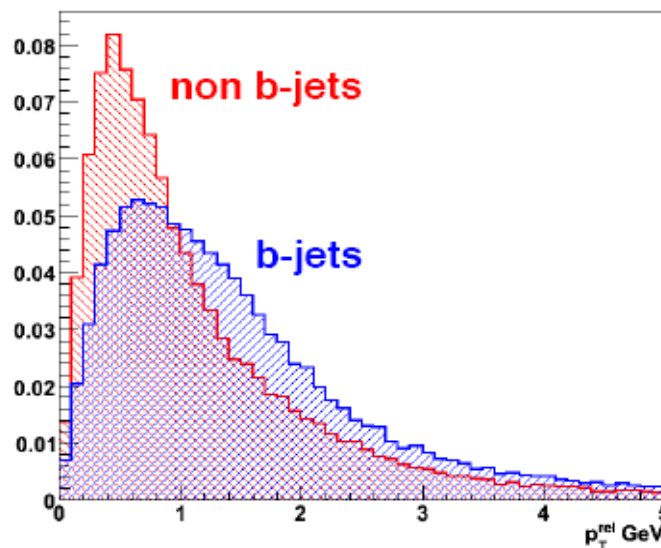
- In CMS different methods are used to measure the b-tagging efficiency from data:

P_{Trel} (CMS Note 2007/046)

- Based on events with at least 2 jets and a non-isolated muon ($\Delta R(j, \mu) < 0.4$)
- Method relies directly on fit of p_{Trel} distr.
- P_{Trel} distr. fitted with linear combination of b and c+light jet templates
- Process repeated after tagging
- Btag efficiency calculated as the ratio between b jets after and before tagging



- An extension uses additional info from data: counting method



System 8 (CMS Note 2008/081)

- Based on events with at least 2 jets and a non-isolated muon ($\Delta R(j,\mu) < 0.4$)
- **Two different data samples**
 - Muon-in-jet + away-jet sample (n)
 - Muon-in-jet + tagged-away-jet sample (p)
- **Two independent taggers**
 - Life Time Tag (Track Counting/ Track Prob.)
 - Soft Muon Tag ($p_{Trel} > 0.8$ GeV)

• A system of **8 equations**
with **8 unknowns**.

• Depends minimally on MC:
MC is only used to evaluate
correlation factors
between taggers.

$$\begin{array}{l}
 \text{muon-in-jet+away-jet} \\
 \text{muon-in-jet+tagged-away-jet} \\
 \text{"probe" tagger} \\
 \text{"tag" tagger} \\
 \text{"tag"+"probe" tagger}
 \end{array}
 \left\{
 \begin{array}{l}
 n = n_b + n_{cl} \\
 p = p_b + p_{cl} \\
 n^{\text{tag}} = \epsilon_b^{\text{tag}} n_b + \epsilon_{cl}^{\text{tag}} n_{cl} \\
 p^{\text{tag}} = \beta \epsilon_b^{\text{tag}} p_b + \alpha \epsilon_{cl}^{\text{tag}} p_{cl} \\
 n^{p_{Trel}} = \epsilon_b^{p_{Trel}} n_b + \epsilon_{cl}^{p_{Trel}} n_{cl} \\
 p^{p_{Trel}} = \delta \epsilon_b^{p_{Trel}} p_b + \gamma \epsilon_{cl}^{p_{Trel}} p_{cl} \\
 n^{\text{tag}, p_{Trel}} = \kappa_b \epsilon_b^{\text{tag}} \epsilon_b^{p_{Trel}} n_b + \kappa_{cl} \epsilon_{cl}^{\text{tag}} \epsilon_{cl}^{p_{Trel}} n_{cl} \\
 p^{\text{tag}, p_{Trel}} = \kappa_b \beta \delta \epsilon_b^{\text{tag}} \epsilon_b^{p_{Trel}} p_b + \kappa_{cl} \alpha \gamma \epsilon_{cl}^{\text{tag}} \epsilon_{cl}^{p_{Trel}} p_{cl}
 \end{array}
 \right.$$

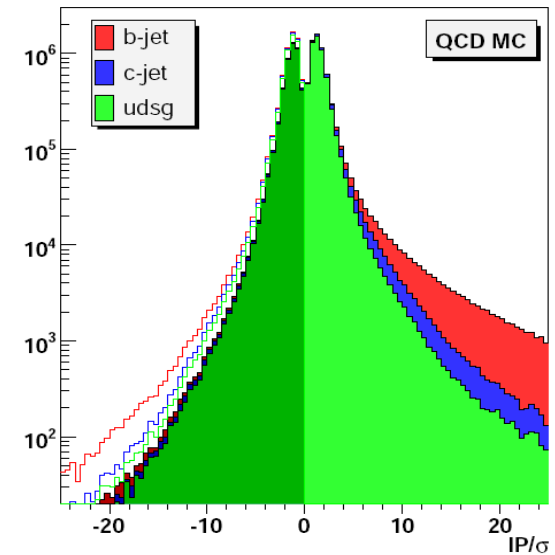
Mistag rate using negative tags CMS Note 2007/048

- Measures the tagging efficiency for uds-g jets (= mistag rates)
- Method uses jets with $p_T > 20$ GeV from QCD samples, this data contains only small fraction of b and c jets
- Then tag the jets by using **tracks with negative IP significance** in order to compute

$$\epsilon_{neg}^{Data}$$

One can then introduce a scale factor R_1

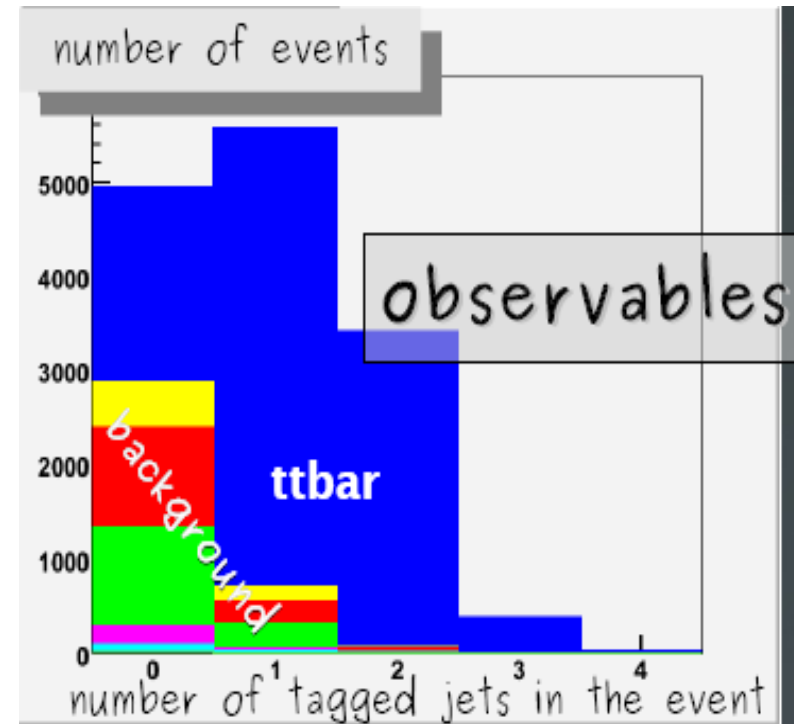
$$\epsilon_{tag}^{Data}(uds-g) = R_1 \times \epsilon_{neg}^{Data}(all)$$



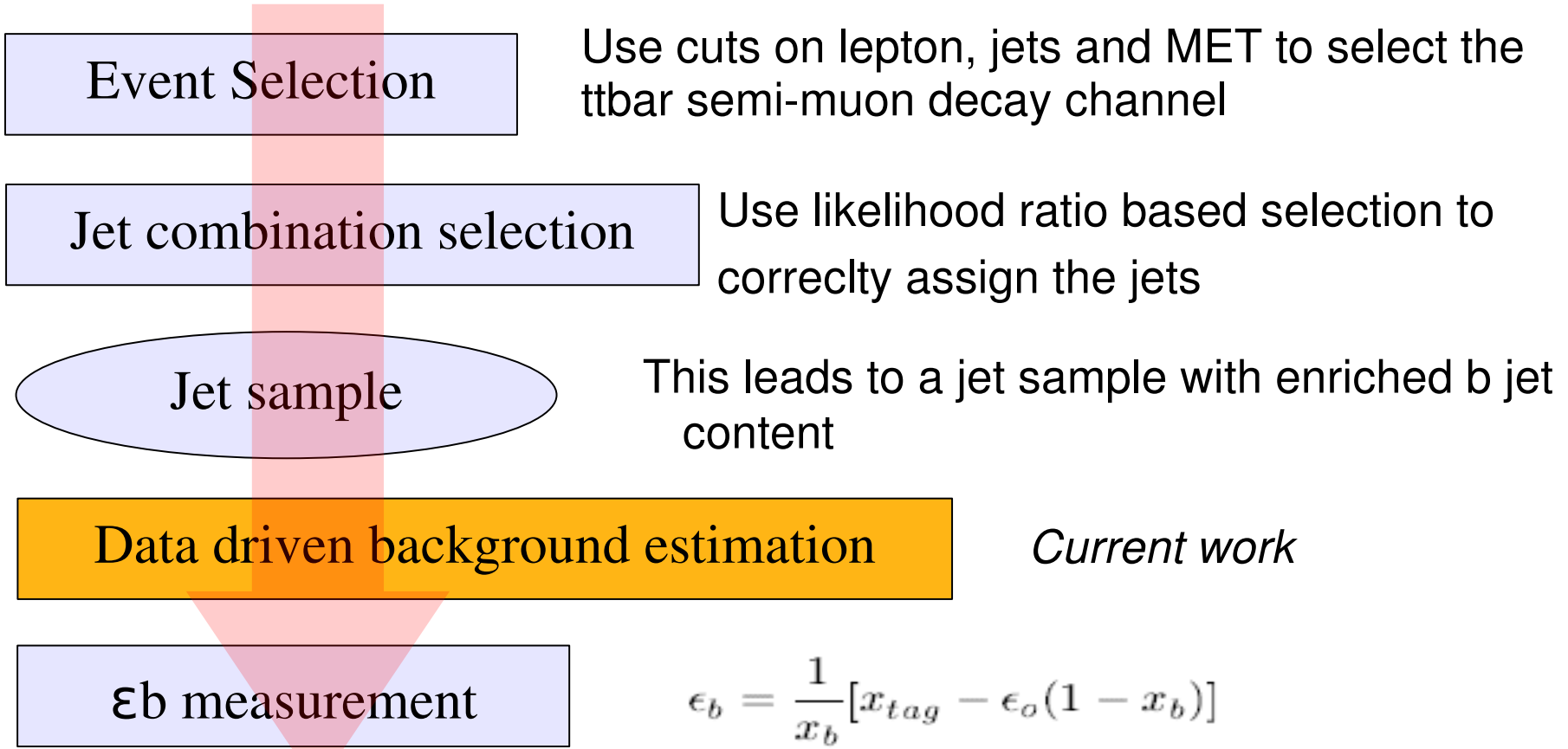
$$R_1 = \frac{\epsilon_{tag}^{MCqcd}(uds-g)}{\epsilon_{neg}^{MCqcd}(all)}$$

B-tag efficiency using $t\bar{t}$ samples tag consistency method

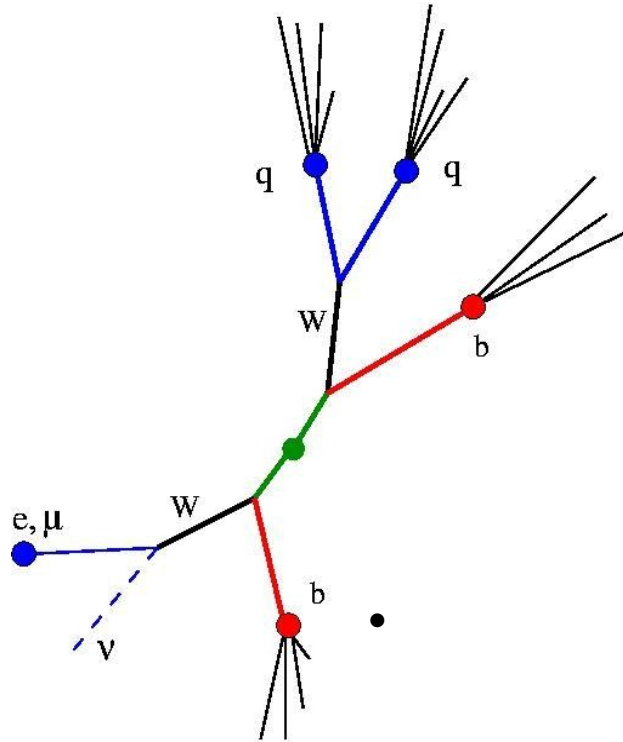
- Currently implemented in the semi-muonic decay channel, dilepton channel is coming
- The number of events (F_{ijk}) with i,j,k , of b,c and light jets is extracted from MC (= jet flavor composition of sample)
- These fractions are used together with the $t\bar{t}$ cross section, the acceptance and the $\epsilon_b \epsilon_c \epsilon_l$ to compare the expected and observed number of tagged jets in an event
- The comparison is done using a maximum likelihood fit to extract the b and c tagging efficiency



Implemented in CMSSW_1_6_12 with TQAF running on CSA07 samples, for the [semi-muonic](#) and the dileptonic decay channels



$$\epsilon_b = \frac{1}{x_b} [x_{tag} - \epsilon_o(1 - x_b)]$$



Event Selection

- **Muon** with $|\eta| < 2.1$, $p_T > 30$ GeV and isolation criteria (trackIso < 3.0 GeV, calIso < 6.0 GeV)
- At least 4 **jets** with $|\eta| < 2.4$ & $p_T > 40$ GeV
- **Missing $E_T > 30$ GeV**

Jet combination selection

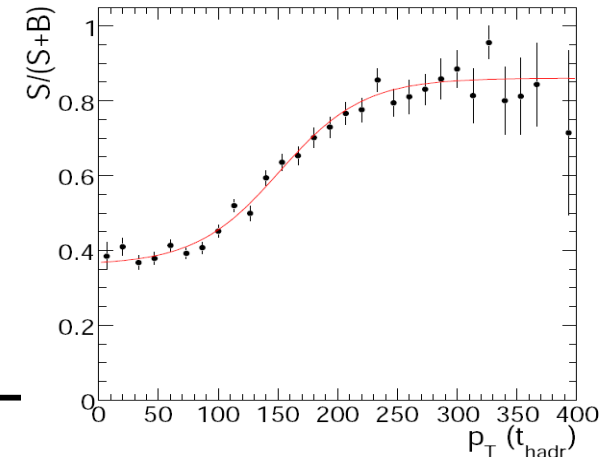
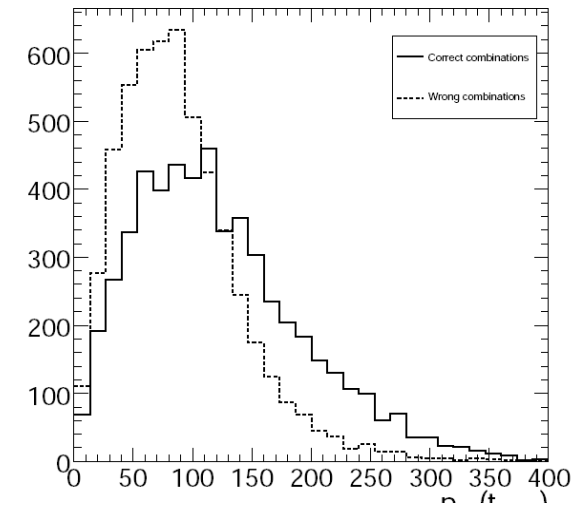
Construct a combined likelihood ratio to **choose the jet combination** among the 12 possible solutions

- For chosen jet combination
 - Require the b coming from the hadronic decaying top to be **loosely btagged**: trackCountingHighEffJetTags > 2 (reduce W+jets background)
 - Require only solutions with **converged kinematic fit** $\chi^2 > 0$ for m_{thad} , m_{tlep} , m_{Whad} , m_{Wlep} constraints.
- **The jet from the leptonic decaying top is used to form**

Jet sample

Likelihood ratio

- After the basic event selection to select the semi-muonic $t\bar{t}$ events the **correct assignment of jets** should be found
- I use a Likelihood ratio based method to find the **most likely solution** among the 12 possible jet combination. (Based on the Neyman-Pearson hypothesis)
- As input 11 observables are used which have a different shape for good and bad assignments (solutions)
- **Full line** shows p_T of hadronic top for the good combination
- **Dashed line** for the bad combinations (norm.)
- From these distributions in a **bin-by-bin** way the **ratio $S/(S+B)$** is calculated
- Via a fit you obtain a function which gives for each value of the had. top p_T **the probability to have signal** (correct jet combination) at this value
- These probabilities are **multiplied** for several observables to obtain a total probability



Likelihood ratio

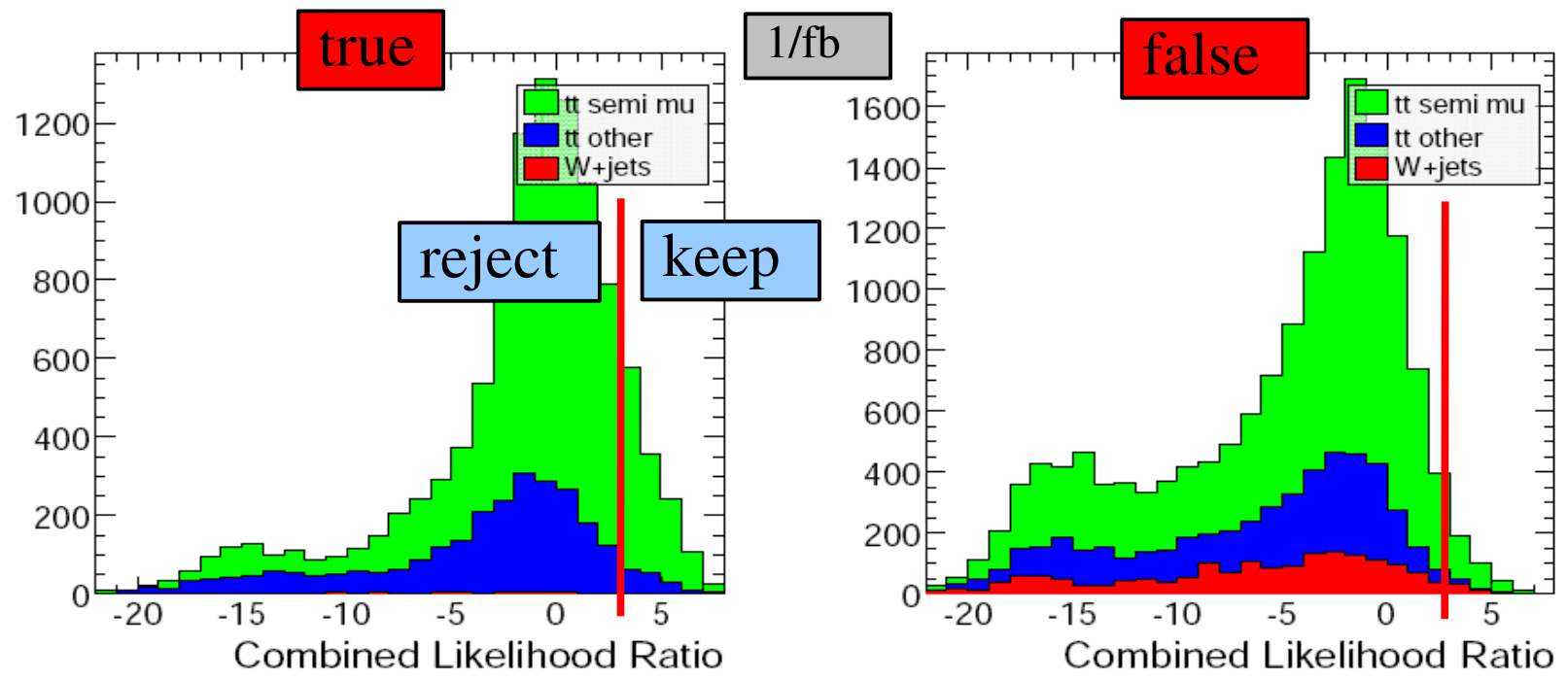
- For each observable the $S/(S+B)$ are calculated and **parameterized** with $f_i(x_i)$, from this the Likelihood Ratio is calculated

$$\mathcal{L}_i(x_i) = \frac{f_i(x_i)}{1 - f_i(x_i)}$$

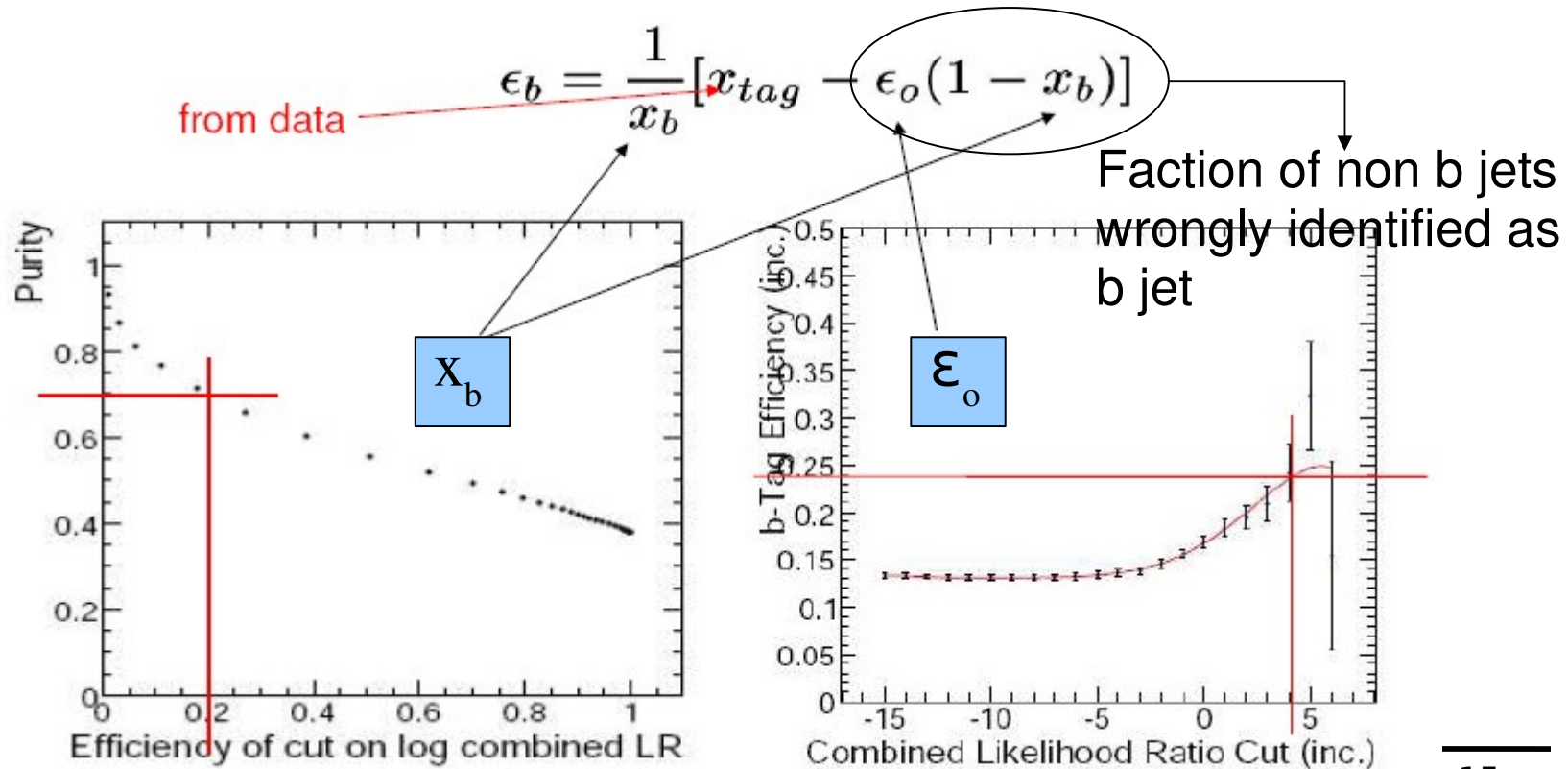
- The Likelihood Ratio functions are **combined** into one observable:

$$\mathcal{L} = \prod_i \mathcal{L}_i(x_i)$$

- Solution with highest value is most probably the correct one
- Leptonic b jet has b flavour** ($\Delta(b_{\text{quark}}, b_{\text{jet}}) < 0.3$) in most likely jet combination?

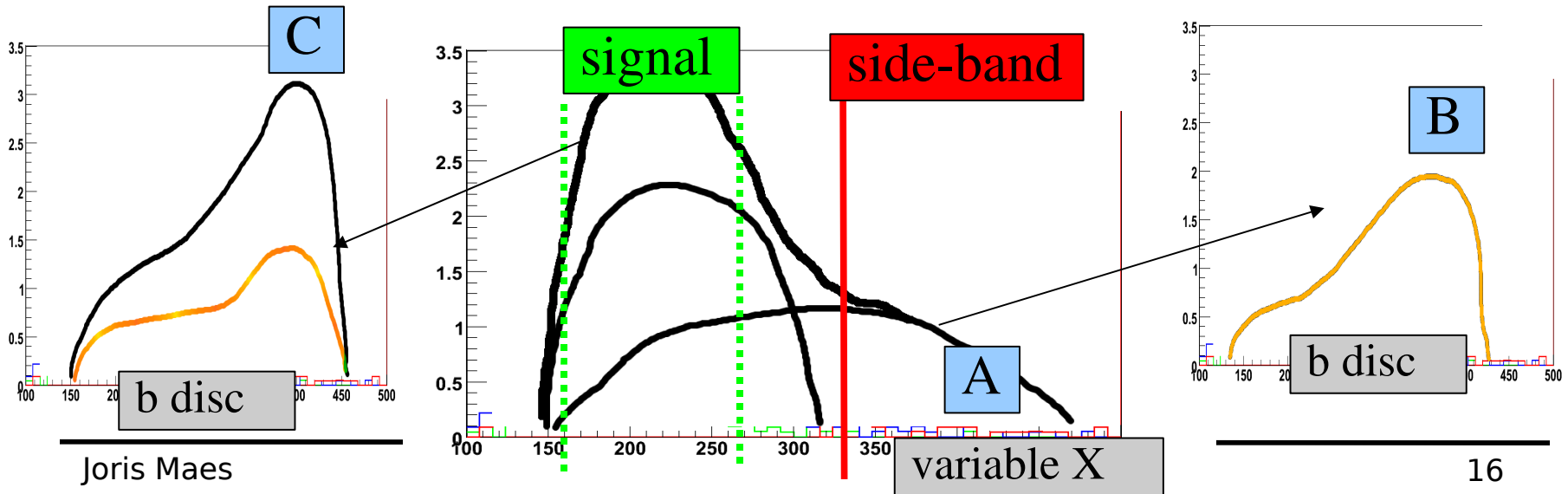


- By cutting harder on the combined likelihood ratio value we tend to have more correct event solutions so the **purity** of the jet sample, formed from the leptonic b jets **increases** → b-enriched jet sample
- a cut at 4 results in a purity x_b of 75% resulting in a jet sample with about 1100 jets for 100/pb
- For these jets **we test if they are tagged by a b-tagging algorithm**
 $x_{tag} = \text{\#positive identified} / \text{\#total number of jets in sample}$

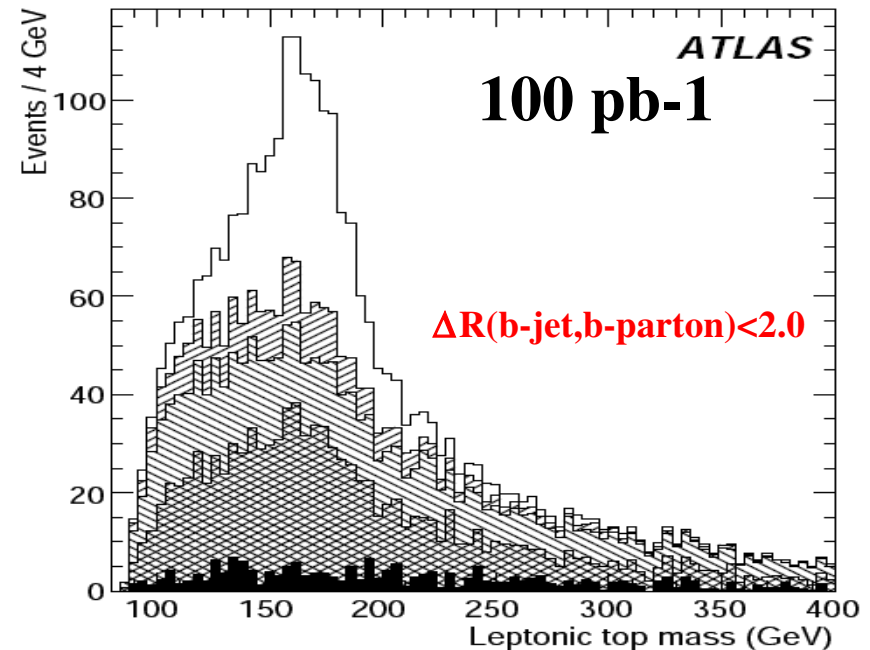
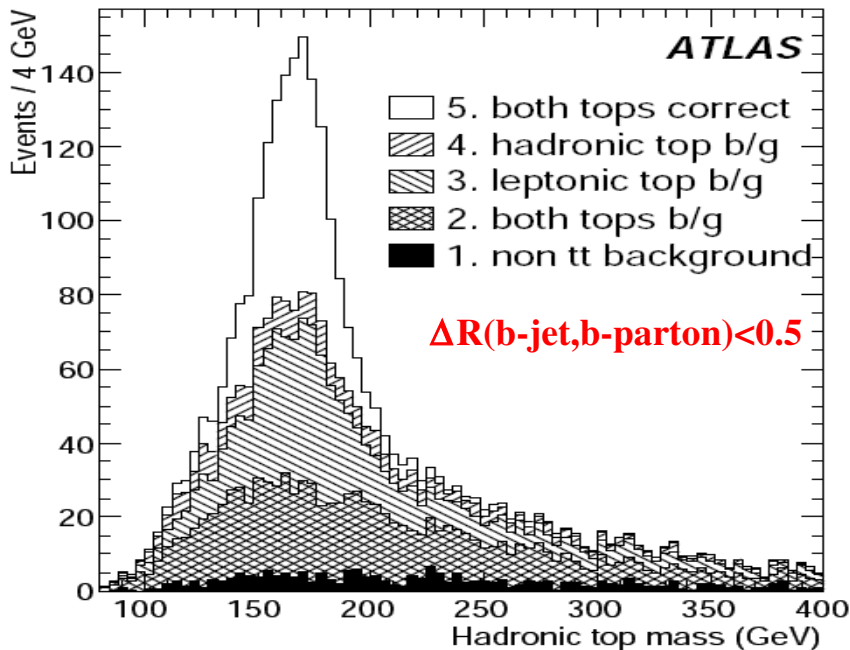


Data driven background control

- The method uses the **flavor composition from Monte Carlo** and is therefore strongly dependent on the model we use.
- To overcome this we would like to estimate the b-tagging efficiency in a slightly different way **using data**
 - A. Find a variable with a clear background dominated region (**side-band**)
 - B. Obtain the **b-tag discriminator distribution** in this region (data)
 - C. **Correct** the b-tag discriminator value in the signal region (data) with the sideband distribution normalized with a **factor from data**
 - The side-band region distribution of the b-tag discriminator should be the same as in the signal region (to be checked)



- **Hadronic top:**
 - Di-jet combinations with $60 < m_{jj} < 100 \text{ GeV}/c^2$ (j_1, j_2 are anti-b-tagged)
 - Combine this di-jet system with a b-tagged jet
 - $E_T(j_1 \text{ or } j_2) \text{ \& } E_T(b) > 40 \text{ GeV}$
- **Leptonic top:**
 - W formed with lepton and neutrino (p_z from W mass constraint)
 - One of the remaining jets is used to reconstruct the leptonic top
- **Combination with largest scalar sum of p_T of the 2 top quarks is chosen**



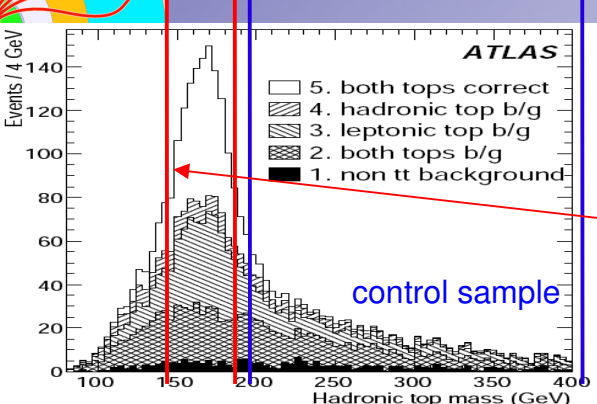
- **Goal: calculate b-tag weight for jets in the b-enriched sample, measure ϵ_b → background has to be well-estimated (next slide) and subtracted (next-to-next slide)**



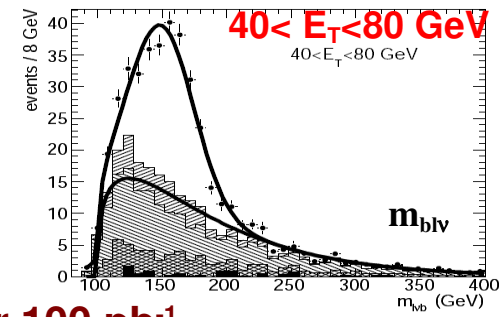
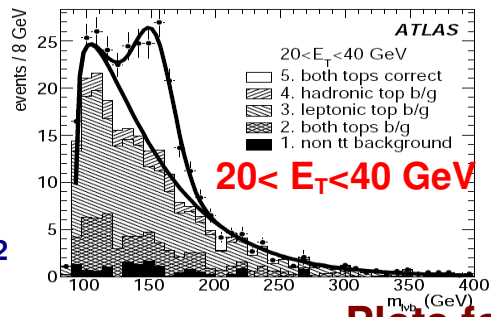
ATLAS



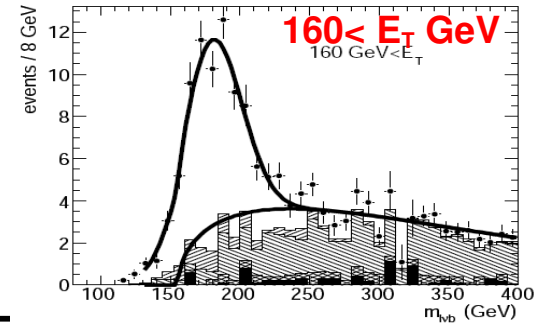
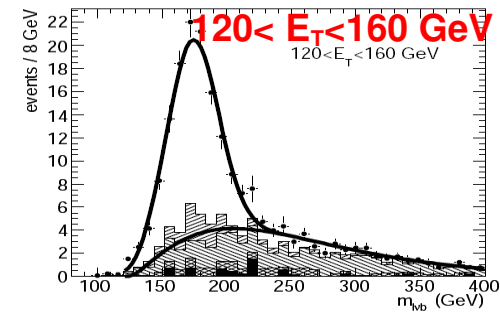
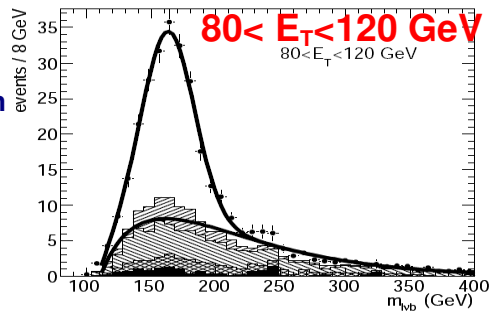
Background estimation



- Function to describe $m_{bl\nu}$
 - Signal sample: $140 < m_{bjj} < 190 \text{ GeV}/c^2$
5 subsamples for leptonic b-jet ET
 - Control sample: $200 < m_{bjj} < 400 \text{ GeV}$,
leptonic b-jet anti-b-tagged \rightarrow
describes background-shape in $m_{bl\nu}$
 - Fit simultaneously background shape and normalization to $m_{bl\nu}$ distributions of signal and control sample (fully data-driven)



Plots for 100 pb⁻¹



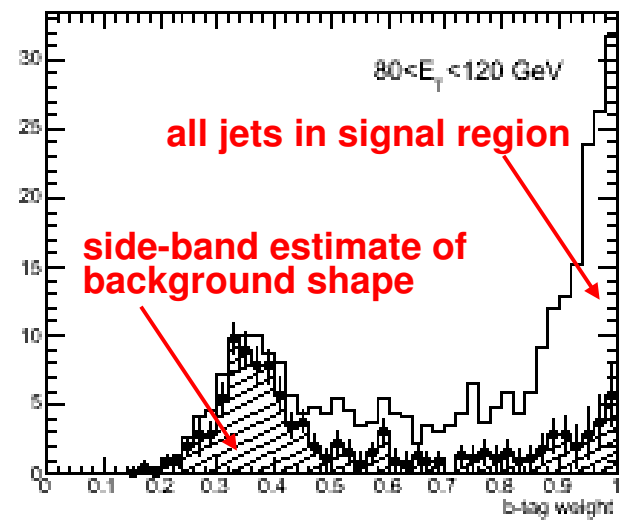
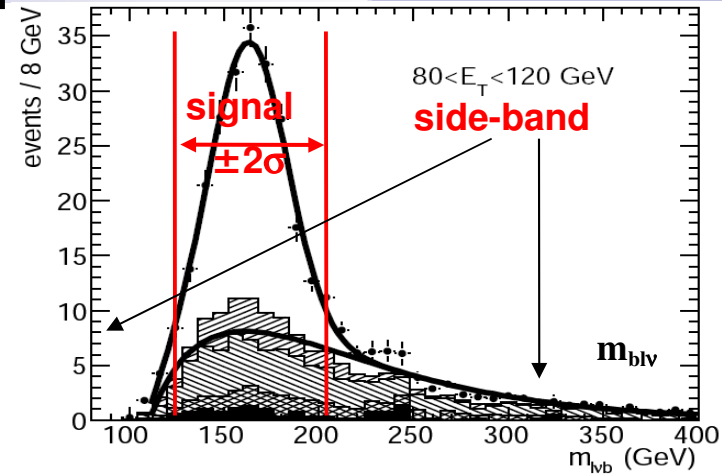
$$F_s(m_{bl\nu}) = c_4 F_b(m_{bl\nu}) + c_5 G((m_{bl\nu} - c_6)/c_7)$$

$$F_b(m_{bl\nu}) = E((m_{bl\nu} - c_0)/400), \quad E(x) = \begin{cases} c_1 x^2 \exp(-c_3 x) & \text{if } x > 0 \\ 0 & \text{if } x < 0 \end{cases}$$

ATLAS

Background subtraction

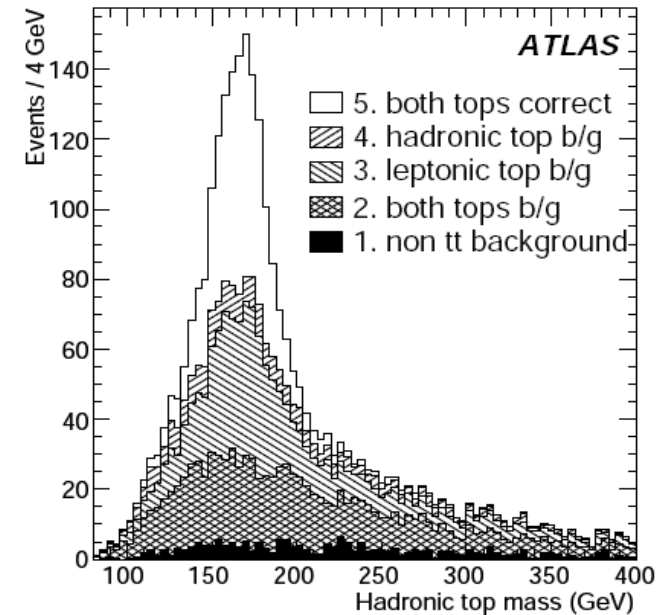
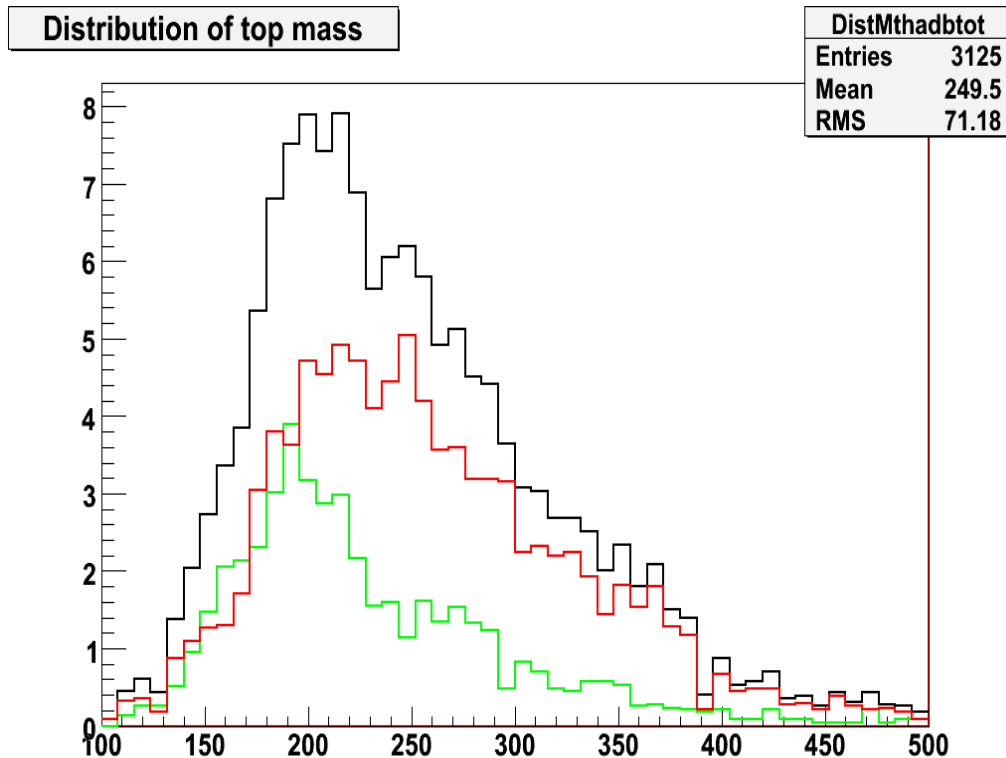
- Obtain b weight distribution in signal region
 - Jet sample from **signal region**:
b-jets + other jets (from comb. background)
 - Other jets have same b-tag weight distribution as jets from events in the side-band
- **Subtract scaled side-band b-tag weight distribution from signal b-tag weight distribution**



$$S = \frac{\int_{s_1}^{s_2} F_b(m_{blv}) dm_{blv}}{\int_{b_1}^{s_1} F_b(m_{blv}) dm_{blv} + \int_{s_2}^{b_2} F_b(m_{blv}) dm_{blv}}$$

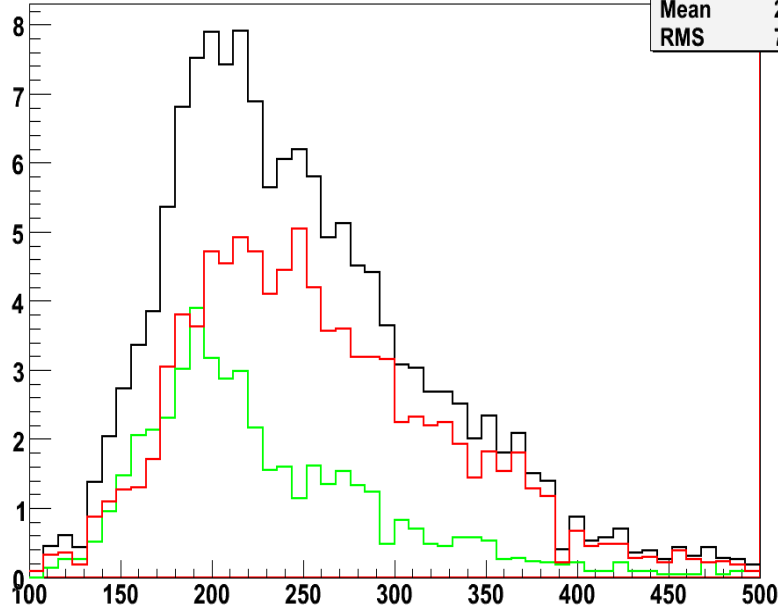
- Do this for each E_T sub-sample
- The side-band estimate does subtract the background in the signal region properly
- Apply b-tagging criterion (cut on b-tag weight)
- Calculate ϵ_b from the remaining jets above this cut value

- We look at the distribution for **calibrated hadronic top mass**, 2 categories
 - Signal: leptonic top matches with generated parton $\Delta R < 0.3$ (green)
 - Background: leptonic top doesn't match with generated parton (red)
- Sideband is contaminated with signal events**
- Mass is **not as peaked** as in ATLAS case

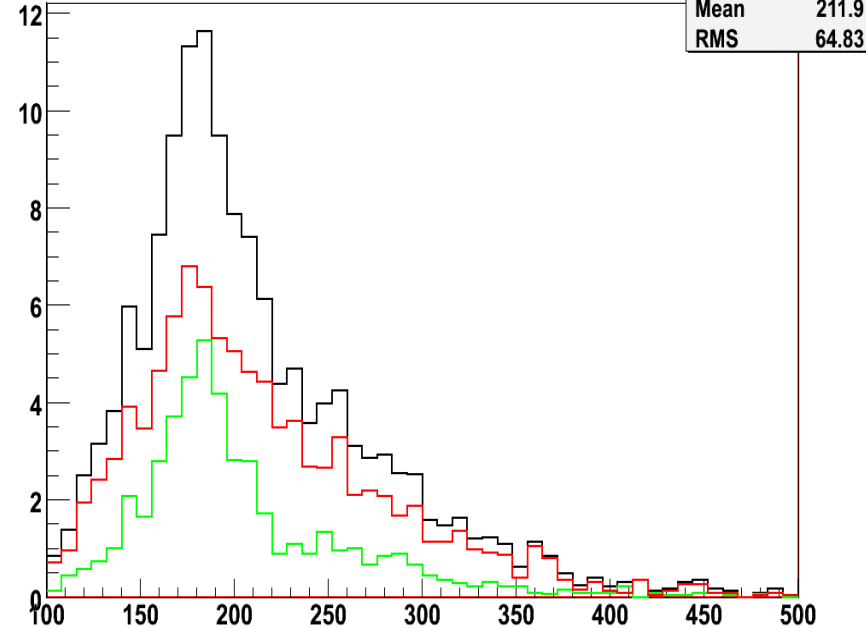


- In stead of the calibrated hadronic top mass we look at the **fitted hadronic top mass**. The 3 jets serving as input for the top mass calculated are adjusted by using the **kinematic fit with $m_{W_{had}}$ and $m_{W_{lep}}$ constraint**.
- The mass is more peaked
- **Background region still contains too much b jets**
- Other observables (e.g. $\Delta\Phi(\text{top}, \text{antitop})$) show similar problem

Distribution of top mass



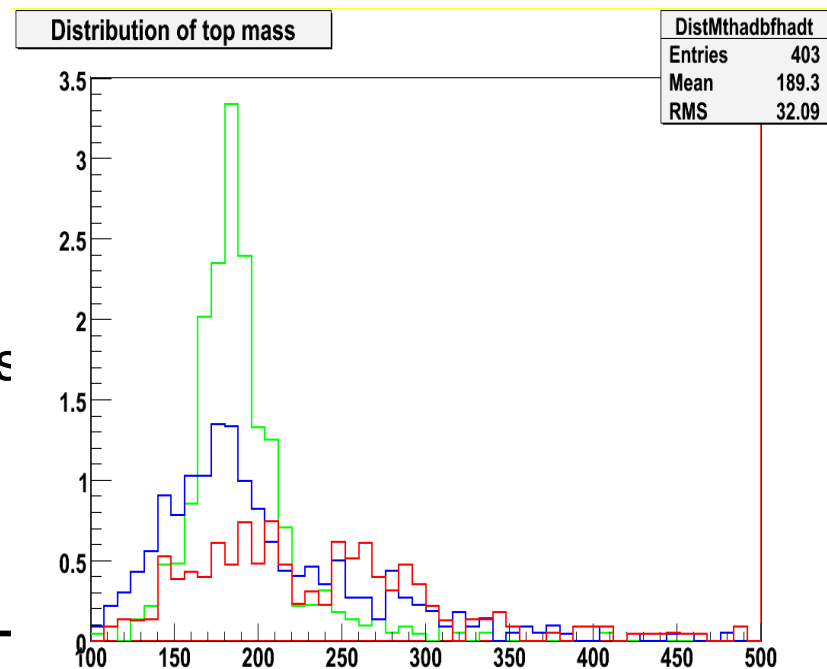
Distribution of fitted top mass



- A **breakdown of the top mass distribution** (green) where the leptonic b jet has really b flavor, this shows that there is an irreducible part.
 - Green: top mass reconstructed with leptonic b and hadronic b jet matched with generated b jets
 - Red: top mass reconstructed with two b jets interchanged
 - Blue: hadronic b jet has not b flavour
 - can be that a light jet is interchanged with hadronic b jet
 - A wrong jet is picked up in the event
 - ...

- This can be understood by the fact that the **likelihood ratio is biasing our sample**.

The likelihood ratio tends to select events which are compatible with the $t\bar{t}b$ topology, so distinguishing inside this sample might be impossible



- Since it is very hard to find an observable to isolate background after applying the likelihood ratio we should look for a **different way to select the correct jet combination**
- Study for which event selection which observables suits best for isolating a background sample

- Should have **significant** amount of events

$$\text{Sign} = \frac{S}{\sqrt{S+B}}$$

- Should have a good **separation** of the signal vs. background

$$\text{Sep} = \sum_{\text{bin}} (S(i) * B(i))$$

- See if the b discriminant distribution can be used to subtract the light/c jets in the signal region

- Port code to CMSSW_2_1_X, use PAT
- Run on the Summer08 samples