



# Identification of b-quark jets in the CMS experiment

Sudhir Malik (on behalf of CMS collaboration)

University of Nebraska-Lincoln and LPC (LHC Physics Center), Fermilab, U.S.A.



## What is CMS?

- CMS** – Compact Muon Solenoid detector at
- LHC** – Large Hadron Collider accelerator at
- CERN** – Organisation Européenne pour la Recherche Nucléaire

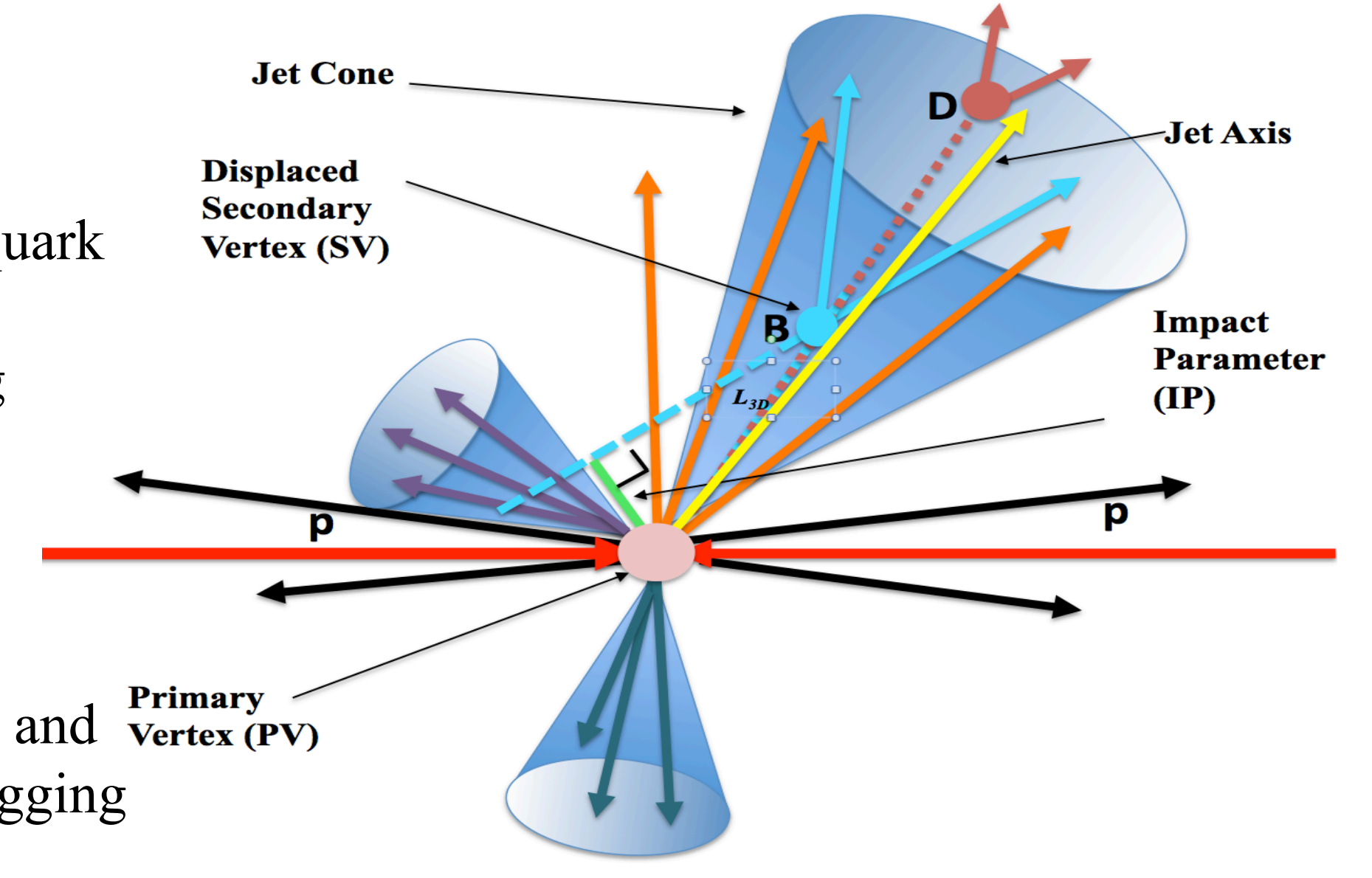
### Highlights of CMS Experiment

- 3.8 Tesla Magnetic Field
- 12500 tonnes
- proton-proton collisions
- Search for Higgs Boson
- Extra dimensions
- Dark Matter
- Discover the Unexpected
- 40 countries
- 200 institutes
- 3500 collaborators
- Collision data taking since over 120 publications



## What is b-tagging and why do we do it?

- Identification (or "tagging") of jets originating from b-quarks is called b-tagging
- Jets from b-quark hadronisation and decay (b jets) present in wide range of processes of interest – top quark decay, Higgs boson, supersymmetric particles
- Identifying b jets is vital to reduce the overwhelming background involving jets from gluons, light-flavour quarks and c-quark fragmentation
- b-quark properties can identify the hadronic jets into which they fragment
- CMS detector has a precise charged particle tracking and robust lepton identification => well matched for b-tagging



## Algorithms and discriminators for b-jet identification

A variety of objects – tracks, vertices, identified leptons – are used to build observables which discriminate between b and light-flavor jets

### Standard track selections:

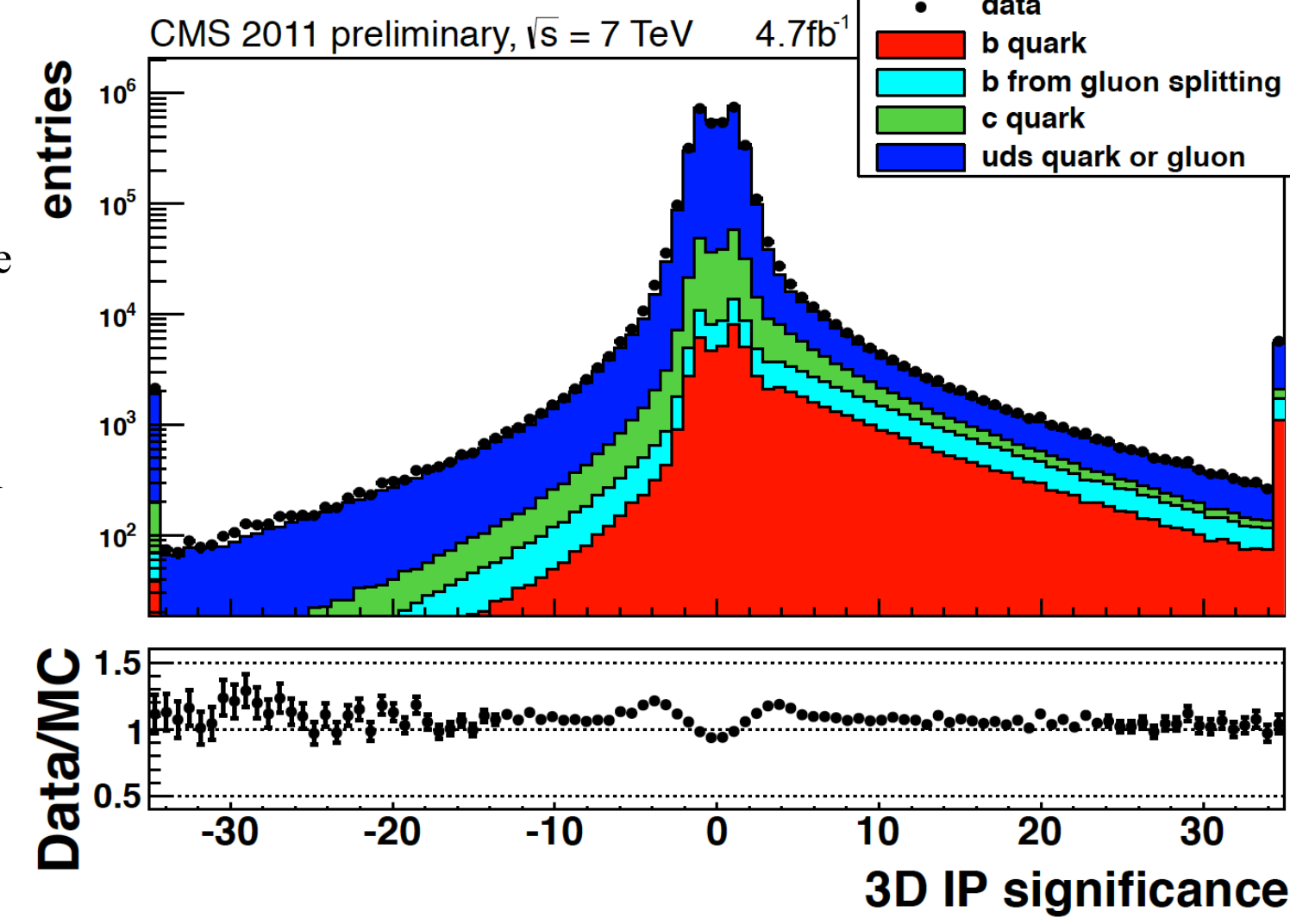
- number of silicon hits  $\geq 8$
- number of pixel hits  $\geq 2$
- transverse impact parameter  $d_{xy} < 0.2\text{cm}$
- longitudinal impact parameter  $d_z < 17\text{cm}$
- $p_T > 1.0\text{GeV}$
- $\chi^2/\text{ndof}$  of the track fit  $< 5.0$
- inside the jet cone:  $\Delta R < 0.5$

### Additional track selections:

- to reduce pile-ups, at the closest point to the jet axis
- the distance to the jet axis  $< 700\mu\text{m}$
- the distance to primary vertex  $< 5\text{cm}$

Impact Parameter (IP) of a track w.r.t. PV can be used to distinguish b hadrons from prompt tracks

Impact Parameter Significance  $S_{IP}$  (defined as ratio between IP and its estimated uncertainty) is used as observable



Algorithms that use IP and  $S_{IP}$  concept are:

**Track Counting (TC)** – sorts tracks in a jet by decreasing values of  $S_{IP}$

Its two versions are:

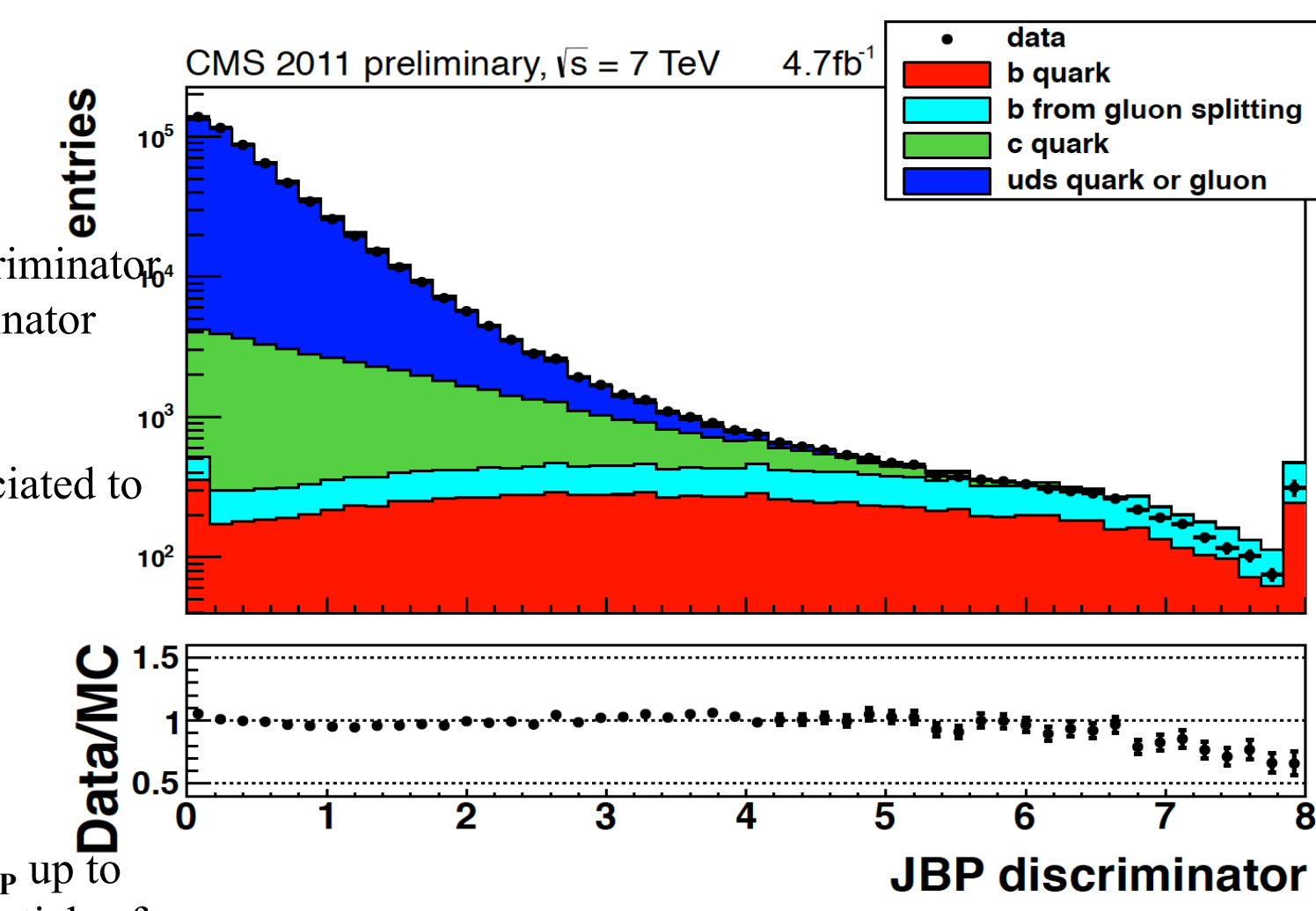
- Track Counting High Efficiency (TCHE)** – uses  $S_{IP}$  of second track discriminator
- Track Counting High Purity (TCHP)** – uses  $S_{IP}$  of third track as discriminator

Those that use info of several tracks in a jet

- Jet Probability (JP)** – uses estimate of the likelihood that all tracks associated to the jet come from PV

$$P_{jet} = \Pi \cdot \sum_{i=0}^{N-1} \frac{(-\ln \Pi)^i}{i!} \quad \text{with} \quad \Pi = \prod_{i=1}^N [\max(P_i, 0.005)]$$

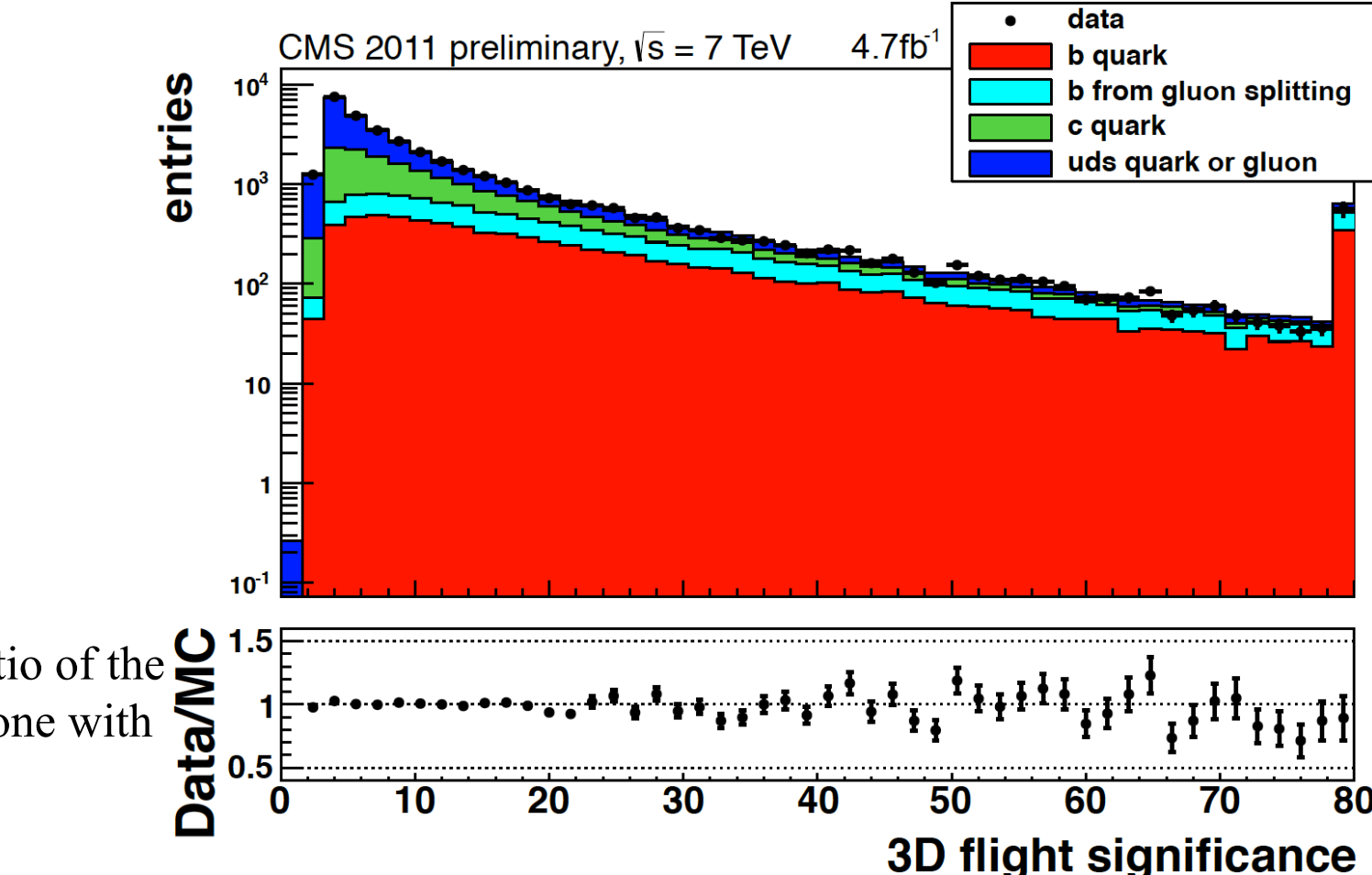
- $P_i$  is the compatibility of track i with PV
- $-\ln P_{jet}$  is plotted using only positive IP tracks



**Jet B probability (JBP)** – gives more weight to tracks with the highest  $S_{IP}$  up to four such tracks, matching the average number of reconstructed charged particles from B-hadron decays

### Secondary vertex candidate selection

- Reject vertices which could be compatible with primary vertex: less than 65% tracks are shared with primary vertex,  $L_{3D} > 3\sigma(L_{3D})$
- Reject interaction vertices and decays of long-lived mesons:  $L_{3D} > 2.5\text{cm}$ ,  $M_{SV} < 6.5\text{GeV}$ ,  $M_{SV} \neq M_{K^*}$
- $\Delta R$  of the flight direction (dotted-blue) and the jet axis (yellow)  $< 0.5$
- Discriminator is  $\log[1 + |L_{3D}|/\sigma(L_{3D})]$



Algorithms using secondary vertices are:

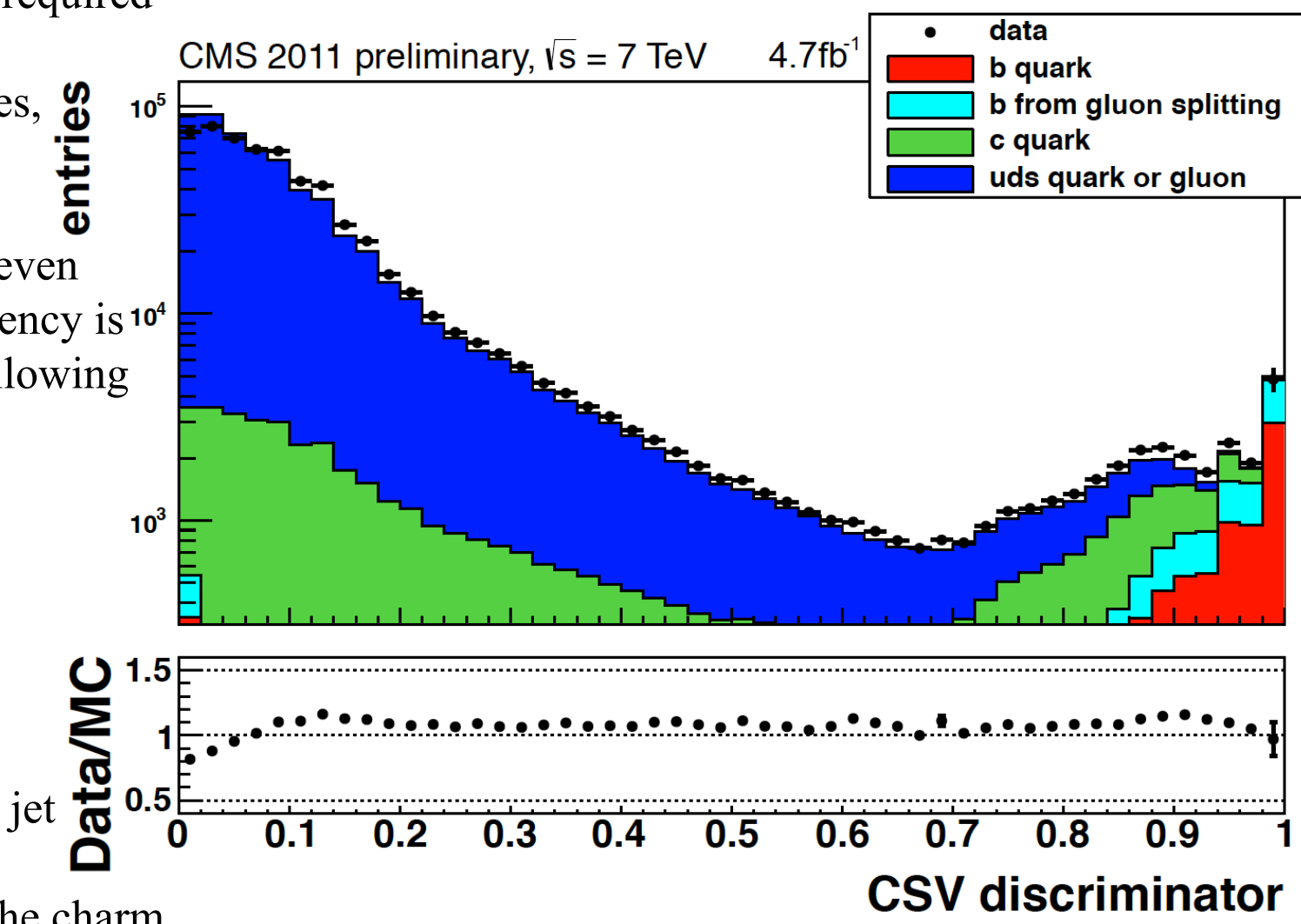
**Simple Secondary Vertex (SSV)** – uses significance of the flight distance (ratio of the flight distance to its estimated uncertainty). If several vertices are present the one with highest significance is used. Its two versions are:

- Simple Secondary Vertex High Efficiency (SSVHE)** – uses vertices with at least two associated tracks
- Simple Secondary Vertex High Purity (SSVHP)** – at least three tracks are required

Another algorithm uses the approach that involves the use of secondary vertices, together with the track based information and is called:

**Combined Secondary Vertex (CSV) algorithm** – provides discrimination even when no secondary vertices are found, so maximum possible b-tagging efficiency is not limited by the secondary vertex reconstruction efficiency. It inputs the following variables into a neural network algorithm. The discriminator is the output.

- the real, "pseudo," or "no" vertex categories
- the 2D flight distance significance
- the vertex mass
- the number of tracks at the vertex
- the number of tracks in the jet
- the 3D signed IP significances for each track in the jet
- the ratio of the energy carried by tracks at the vertex w.r.t. to all tracks in the jet
- the pseudo-rapidity of the tracks at the vertex with respect to the jet axis
- the 2D IP significance of the first track that raises the invariant mass above the charm threshold of 1.5 GeV when subsequently summing up tracks ordered by decreasing IP significance;



## Efficiency measurements

Measure from data using muons jets

### Muon selection:

Global Muon,  $p_T > 5\text{GeV}$ ,  $|\eta| < 2.4$ , number of muon hits  $\geq 1$ , number of muon segments  $\geq 1$ ,  $\chi^2/\text{ndof}$  of the global muon  $< 10$ , number of silicon hits  $\geq 10$ , Number of pixel hits  $\geq 1$ ,  $\chi^2/\text{ndof}$  of the track track  $< 10$ , longitudinal impact parameter  $d_z < 1\text{cm}$

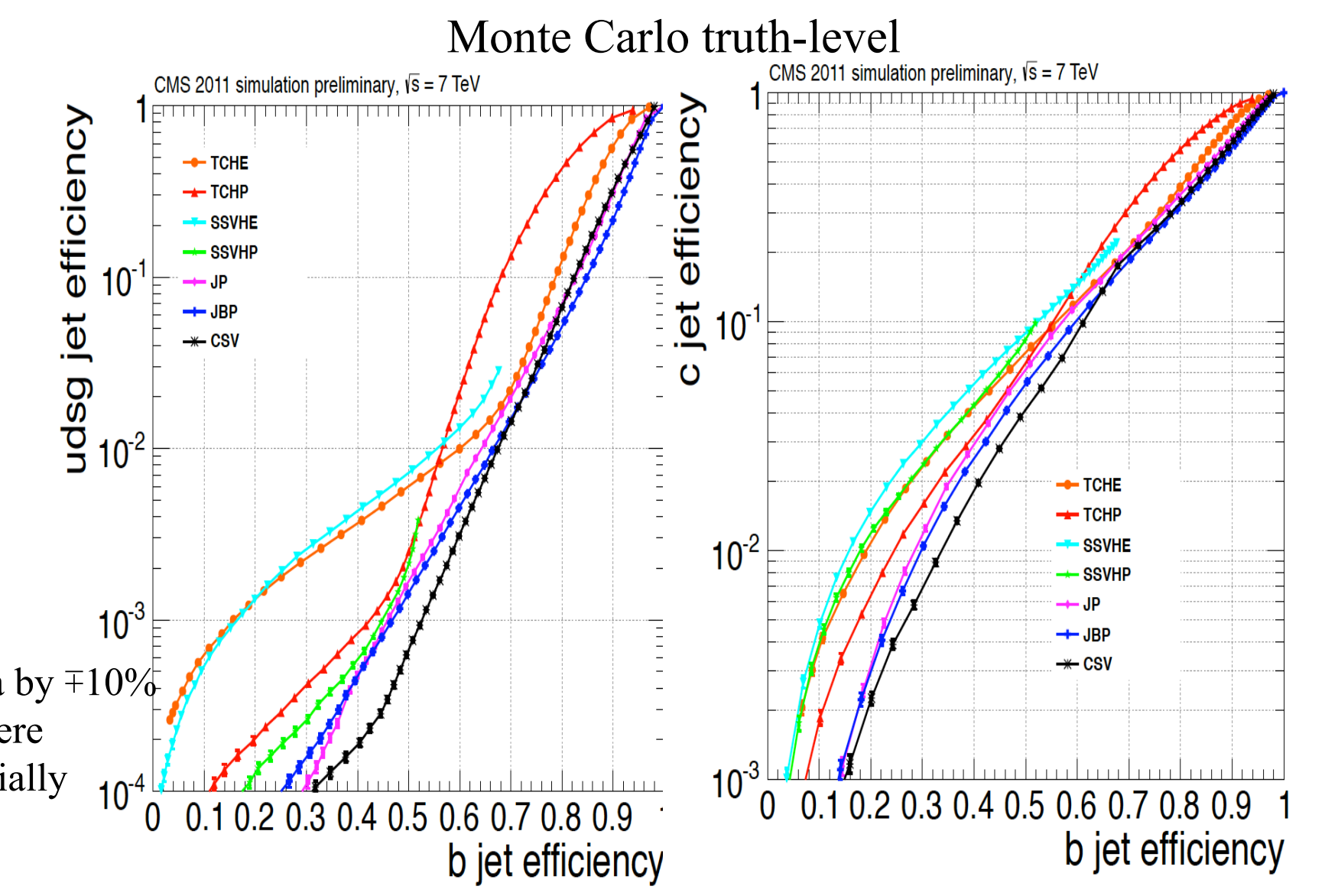
### Muon jet:

the muon is inside the jet cone:  $\Delta R < 0.4$

### Common Systematic Errors Calculation

**pile-up:** vary the average value of the pile-up in data by  $\pm 10\%$   
**gluon splitting:** investigate Monte Carlo sample where the number of events with gluon splitting was artificially changed by 50%

Muon  $p_T$ : vary  $p_T$  cut up to 9 GeV



### Measurement with kinematic properties of the muon jets

Steps:

- the event has another jet fulfilling the Track Counting High Purity (third ranked track) b-tagging criterion at "medium" b-tagging efficiency working point (TCHPM)
- Separate the muon jets into tagged and untagged subsamples by a discriminator working point whose efficiency is to be measured

For the two subsamples separately, fit the spectra of muon jets  $p_{Trel}$  or IP3D using templates of b, c and udsg jets derived from simulation or inclusive jet data. The fraction of b jets is obtained in this step as  $(f_b^{tag} / f_b^{untag})$ . The  $p_{Trel}$  is the transverse momentum of muon associated with the jet direction use the equation to get the result:

$$\epsilon_b^{tag} = \frac{f_b^{tag} \cdot N_{data}^{tag}}{f_b^{tag} \cdot N_{data}^{tag} + f_b^{untag} \cdot N_{data}^{untag}}$$

- Systematic uncertainties:

- use (TCHE, TCHEM, TCHPM) and repeat
- vary the predicted ratio by  $\pm 20\%$  and repeat

### Measurement with the System8 method

Steps:

- take all muon jets as sample 1 (n)
- take the muon jets in those events where another jet is tagged by the Track Counting High Purity b-tagging criterion at "low" working point (TCHPL) as sample 2 (p)
- tag the muon jet by  $p_{Trel} > 0.8\text{GeV}$  or a discriminator working point whose efficiency is to be measured.
- solve the following eight equations and get the efficiency

$$n = n_b + n_{cl}$$

$$p = p_b + p_{cl}$$

$$n^{tag} = \epsilon_b^{tag} n_b + \epsilon_{cl}^{tag} n_{cl}$$

$$p^{tag} = \beta_b^{tag} \epsilon_b^{tag} p_b + \alpha_{cl}^{tag} \epsilon_{cl}^{tag} p_{cl}$$

$$n^{P_{Trel}} = \epsilon_b^{P_{Trel}} n_b + \epsilon_{cl}^{P_{Trel}} n_{cl}$$

$$p^{P_{Trel}} = \beta_b^{P_{Trel}} \epsilon_b^{P_{Trel}} p_b + \alpha_{cl}^{P_{Trel}} \epsilon_{cl}^{P_{Trel}} p_{cl}$$

$$n^{tag, P_{Trel}} = \beta^n \epsilon_b^{tag} \epsilon_b^{P_{Trel}} n_b + \alpha^n \epsilon_{cl}^{tag} \epsilon_{cl}^{P_{Trel}} n_{cl}$$

$$p^{tag, P_{Trel}} = \beta^p \epsilon_b^{tag} \epsilon_b^{P_{Trel}} p_b + \alpha^p \epsilon_{cl}^{tag} \epsilon_{cl}^{P_{Trel}} p_{cl}$$

$\alpha$ 's and  $\beta$ 's are:

correlation factors between the sample1 and sample2 which are obtained from MC simulation

variables to solve:

n: number of events of sample1

p: number of events of sample2

$\epsilon$ : efficiency

subscripts:

b: b jet; cl: non-b jet

superscripts:

tag: tagged by the discriminator working point to be measured

$P_{Trel}$ : tagged by  $p_{Trel} > 0.8\text{GeV}$

Systematic uncertainties:

- use (TCHE, TCHEM, TCHPM) and repeat
- vary from 0.5 to 1.2GeV and take the largest discrepancy

### Measurement with JP as a reference tagger

It can be applied on any jet (not only muon jet) since most jets have JP information

Steps:

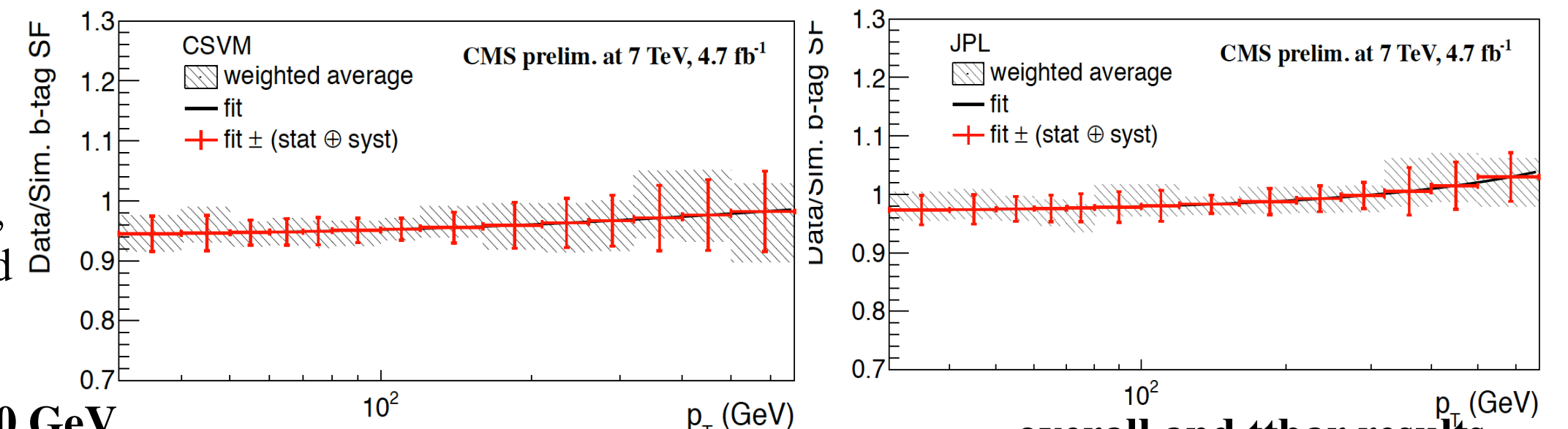
- acquire  $C_b$ , the fraction of b jets with JP information (before tagging) in the Monte Carlo.
- fit  $f_b^{beforetag}$  distribution by templates for b, c and udsg jets to get the fraction of b jets
- tag with a discriminator working point whose efficiency is to be measured
- repeat step 2 on step 3 tagged sample and get the fraction of b jets  $f_b^{tag}$
- the efficiency is the ratio of  $f_b^{tag} \cdot N_{data}^{tag}$  and  $(f_b^{beforetag} \cdot N_{data}^{beforetag}) / C_b$

Systematic uncertainties

- the difference between muon jets and inclusive jets is considered as a systematic error
- estimate the systematic uncertainty of the half residual correction
- instead of JP, use CSV as a reference tagger and then compare

## Results

The result is shown as the ratio of the efficiencies measured in data and the Monte Carlo truth efficiency, SFb. The errors are presented as statistical+systematic.



### 80 GeV < jet pT < 120 GeV

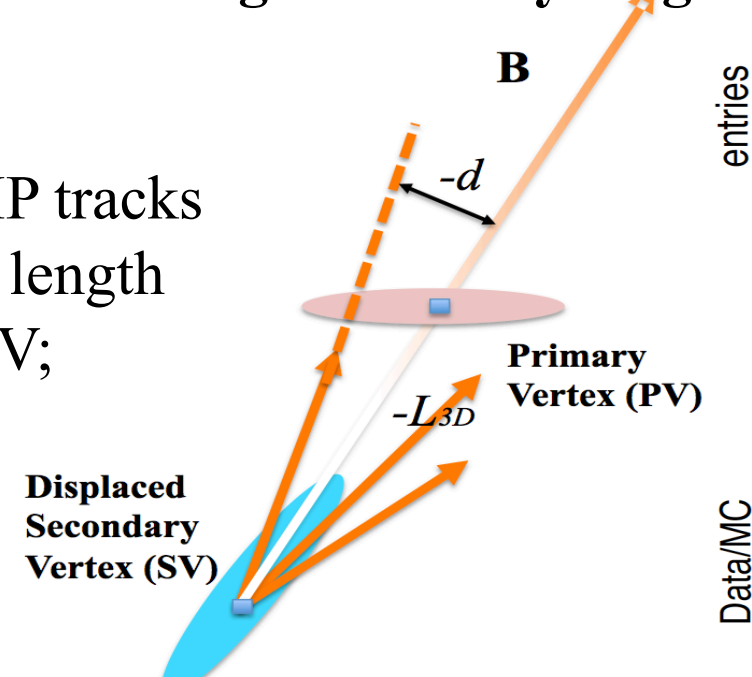
b tagger	SFb (PtRel)	SFb (System8)	SFb (JP)	SFb (comb.)	b tagger	SFb in multijet events	SFb in tt events
JPM	0.90 ± 0.01 ± 0.03	0.91 ± 0.03 ± 0.06	0.99 ± 0.02 ± 0.05	0.92 ± 0.03	JPM	0.92 ± 0.03	0.95 ± 0.03
JBPM	0.92 ± 0.01 ± 0.02	0.94 ± 0.03 ± 0.08	0.99 ± 0.02 ± 0.05	0.92 ± 0.03	JBPM	0.92 ± 0.03	0.93 ± 0.04
TCHEM	0.94 ± 0.01 ± 0.03	0.97 ± 0.03 ± 0.07	0.98 ± 0.01 ± 0.03	0.95 ± 0.02	TCHEM	0.95 ± 0.03	0.96 ± 0.04
TCHPM	0.95 ± 0.01 ± 0.03	0.93 ± 0.02 ± 0.09	0.97 ± 0.01 ± 0.02	0.96 ± 0.02	TCHPM	0.94 ± 0.03	0.93 ± 0.04
SSVHEM	0.93 ± 0.01 ± 0.02	0.91 ± 0.03 ± 0.05	0.97 ± 0.01 ± 0.02	0.95 ± 0.02	SSVHEM	0.95 ± 0.03	0.96 ± 0.04
CSV	0.93 ± 0.01 ± 0.02	0.95 ± 0.03 ± 0.06	0.97 ± 0.01 ± 0.03	0.95 ± 0.02	CSV	0.95 ± 0.03	0.97 ± 0.04

## Mistag Rate Measurement

### The definition of negative discriminator:

- negative IP: when  $\theta$  (see the top plot)  $> 90^\circ$
- TC has the same sign with IP
- JP negative is formed from only negative IP tracks
- SV is negative when it has negative decay length
- CSV: when SV is found, same sign with SV; otherwise same sign with TC.

### negative decay length



### Get the efficiency directly by:

$$\epsilon_{data}^{mistag} = \epsilon_{data}^{-} \cdot \epsilon_{MC}^{mistag} / \epsilon_{MC}^{-}$$

Result: It is shown as the ratio of the efficiencies measured in data and the Monte Carlo truth efficiency

b tagger	mistag rate ( $\pm \text{stat}$ )	scale factor ( $\pm \text{stat} \pm \text{syst}$ )
JPM	0.0109 ± 0.0002	1.02 ± 0.02 ± 0.16
JBPM	0.0112 ± 0.0001	0.94 ± 0.01 ± 0.11
TCHEM	0.0286 ± 0.0003	1.20 ± 0.01 ± 0.14
TCHPM	0.0306 ± 0.0003	1.24 ± 0.01 ± 0.12
SSVHEM	0.0209 ± 0.0002	0.93 ± 0.01 ± 0.08
CSV	0.0152 ± 0.0002	1.10 ± 0.01 ± 0.11

Systematic uncertainties: b and c fractions, gluon fraction, long lived  $K_S^0$  and  $\Lambda$  decays, photon conversion and nuclear interactions, mismeasured tracks, the ratio of negative over positive tagged jets, pile-up, and event sample.