# Run 3 performance and advances in heavy flavor jet tagging in CMS
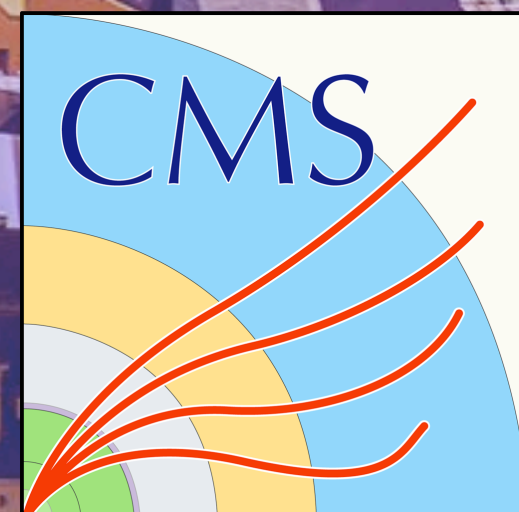
## Uttiya Sarkar
### on behalf of the CMS collaboration
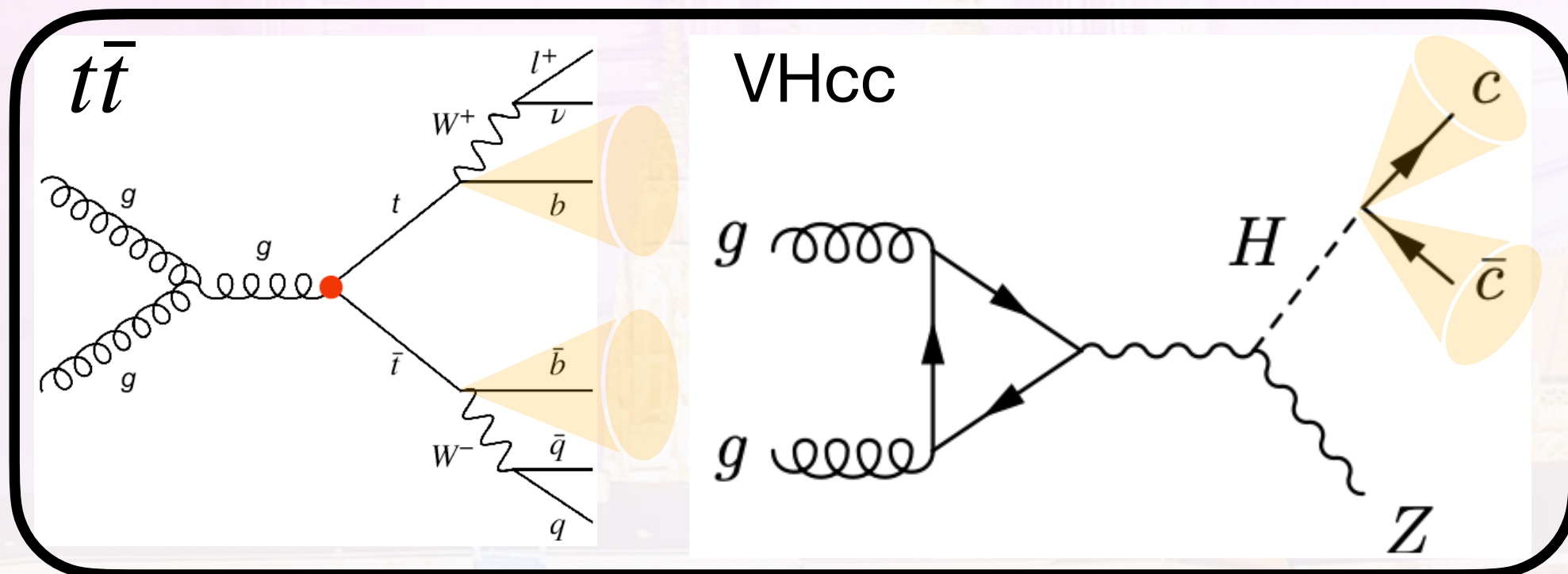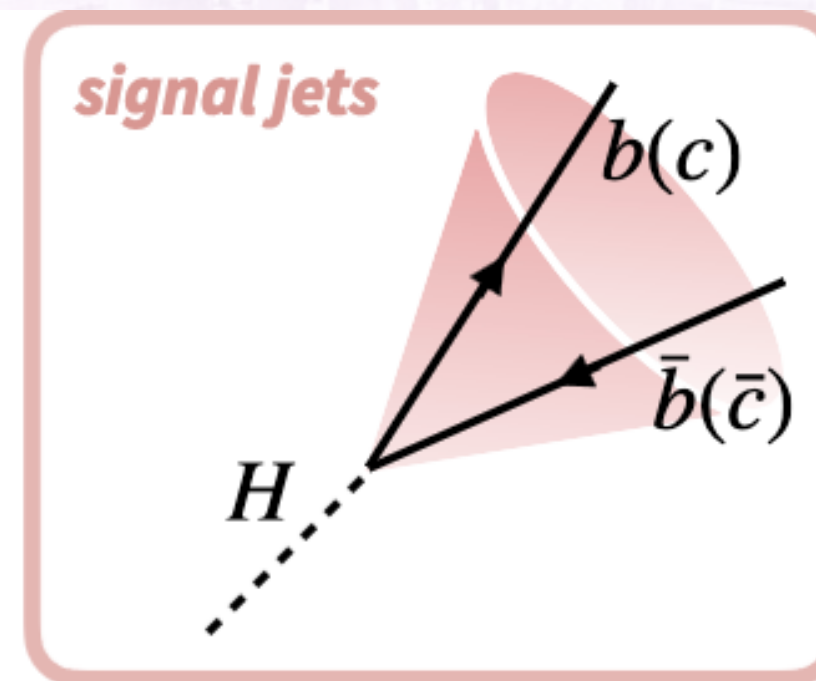
# Overview

- **Heavy-flavor tagging -** jets originating from b (**b jets**) or c (**c jets**) quarks

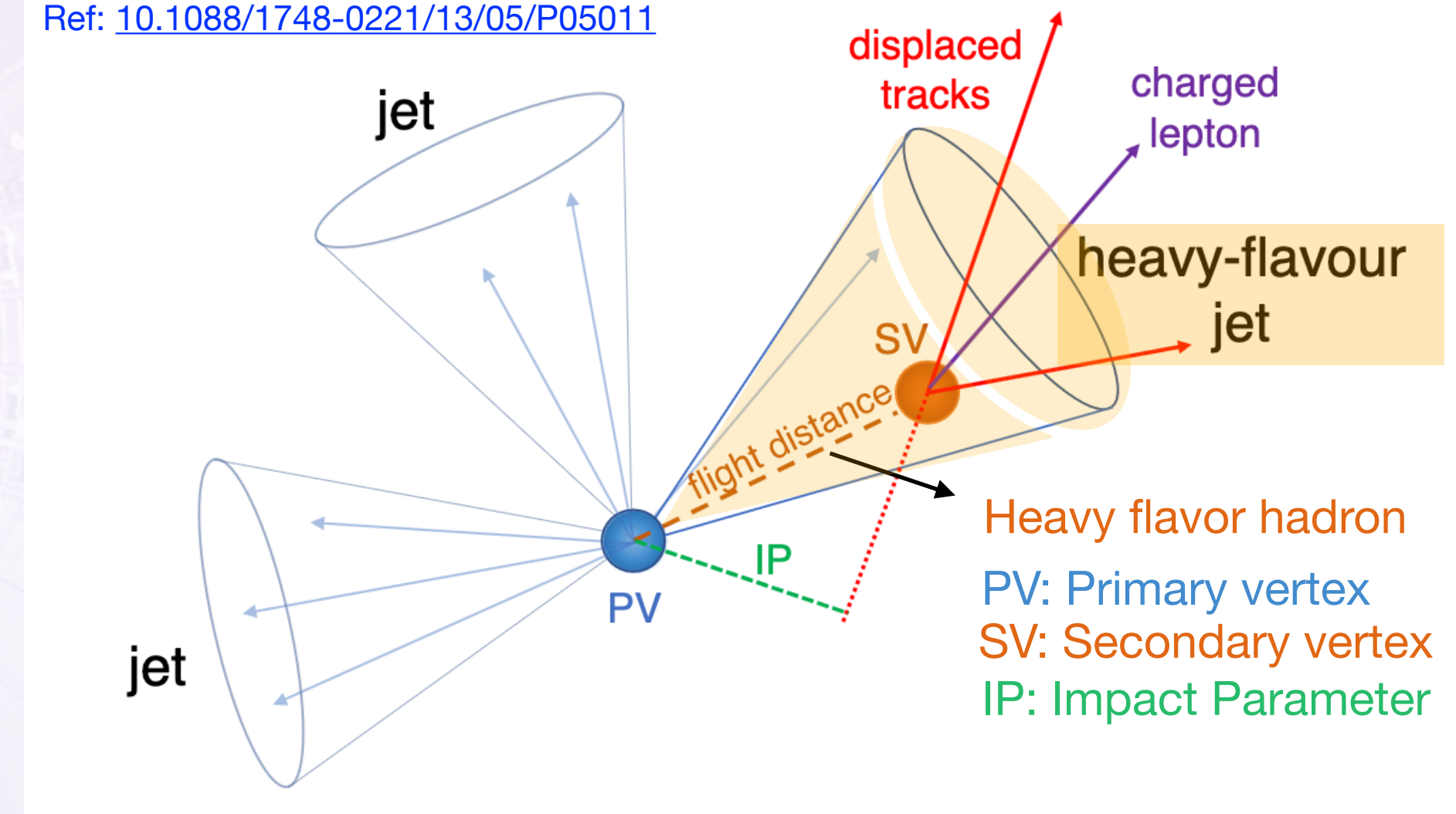- **Physics motivation:** Important in Standard Model (SM), Top, Higgs, BSM and SUSY processes
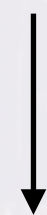
*Resolved jet topology*



*High momentum: merged jet topology*



Ref: 10.1088/1748-0221/13/05/P05011



- Heavy flavor hadron
- PV: Primary vertex
- SV: Secondary vertex
- IP: Impact Parameter

- **Discriminators:**

- Within CMS tracker resolution

- **Significant improvement in machine learning based tagging algorithms**

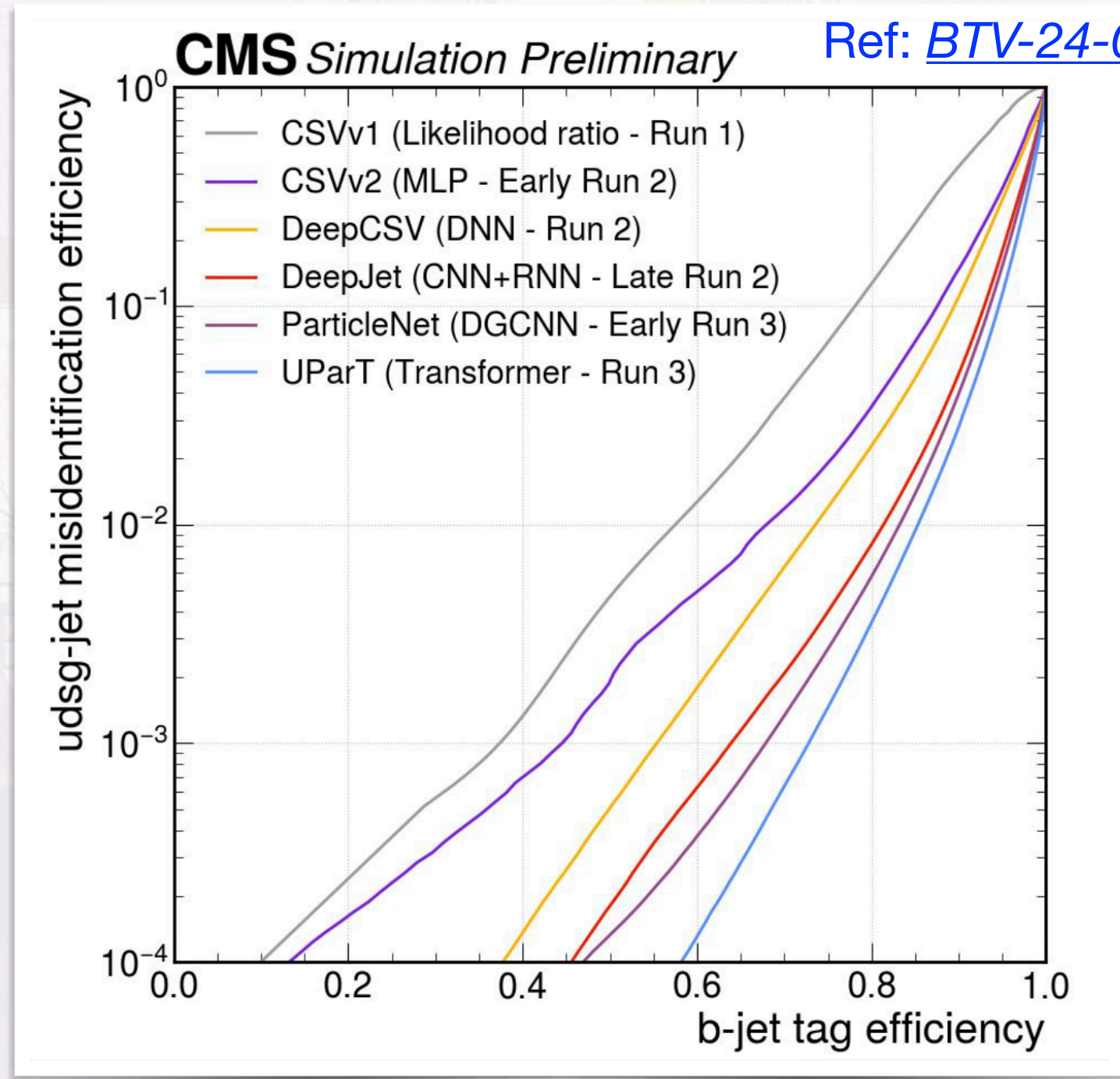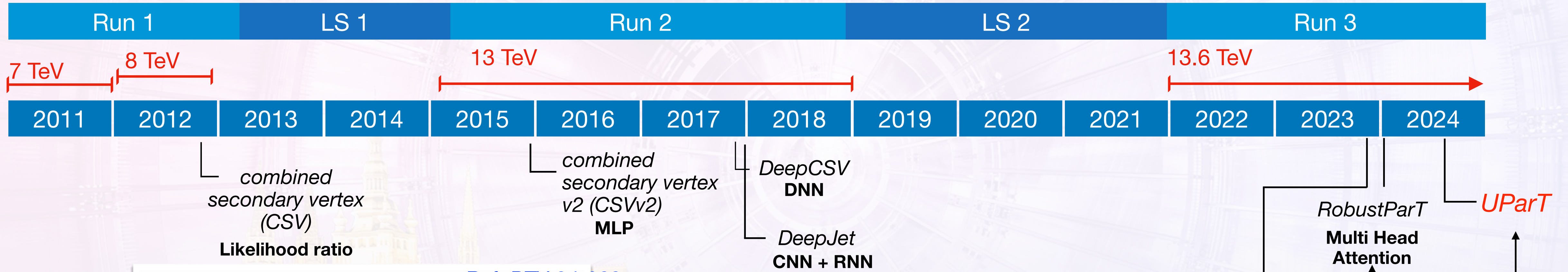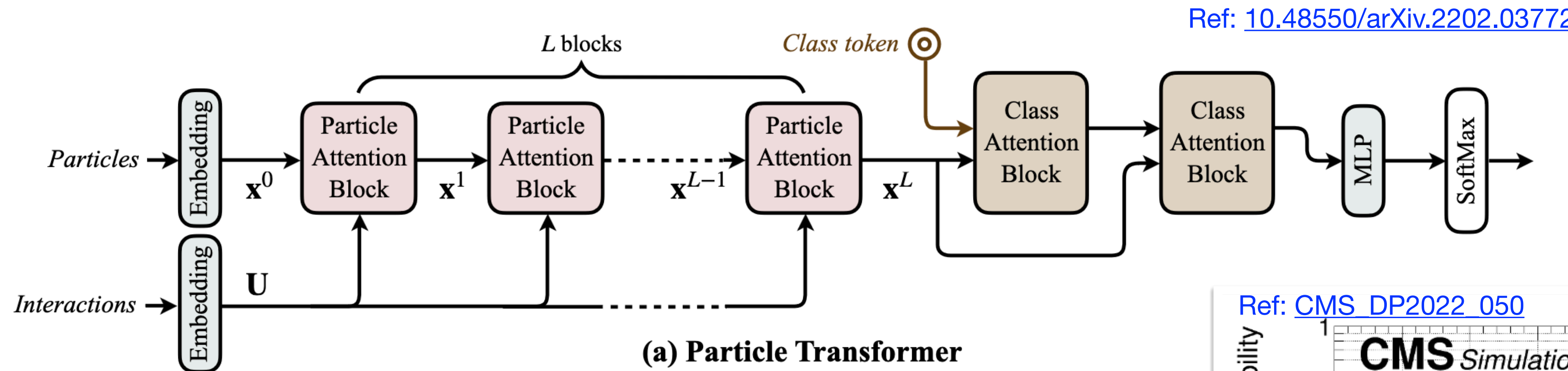*Performs beyond the heavy flavor tagging (can tag τ-jets)*

*b jet*

*c jet*

*udg jet*

$$\mathrm{BvsX} = \frac{\mathrm{P(B)}}{\mathrm{P(B)} + \mathrm{P(X)}}$$

# Historical evolution of flavor taggers in CMS

| Run 1 | LS 1 | Run 2 | LS 2 | Run 3 |
|---|---|---|---|---|

7 TeV    8 TeV    13 TeV    13.6 TeV

| 2011 | 2012 | 2013 | 2014 | 2015 | 2016 | 2017 | 2018 | 2019 | 2020 | 2021 | 2022 | 2023 | 2024 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

*combined secondary vertex (CSV)*
**Likelihood ratio**

*combined secondary vertex v2 (CSVv2)*
**MLP**

*DeepCSV*
**DNN**

*DeepJet*
**CNN + RNN**

Ref: *BTV-24-066*

*ParticleNet*
**DGCNN -** Treat particles as particle clouds - unordered set of its constituent particle

*RobustParT*
**Multi Head Attention**

*UParT*

**Distortion of input features to enhance the robustness**



**CMS** *Simulation Preliminary*

- CSVv1 (Likelihood ratio - Run 1)
- CSVv2 (MLP - Early Run 2)
- DeepCSV (DNN - Run 2)
- DeepJet (CNN+RNN - Late Run 2)
- ParticleNet (DGCNN - Early Run 3)
- UParT (Transformer - Run 3)

udsg-jet misidentification efficiency

b-jet tag efficiency

*Significant improvement in performance over the last decade!*

# Transformer models: ParticleTransformer

Ref: [10.48550/arXiv.2202.03772](#) ← See more details here



(a) Particle Transformer

Ref: [CMS_DP2022_050](#)



- Based on the "Attention" model designed for particles
- Input embedding:
  - Not only inject single particle information, but also include pair-wise features
- Multi-Head Attention (MHA) Pair-wise feature
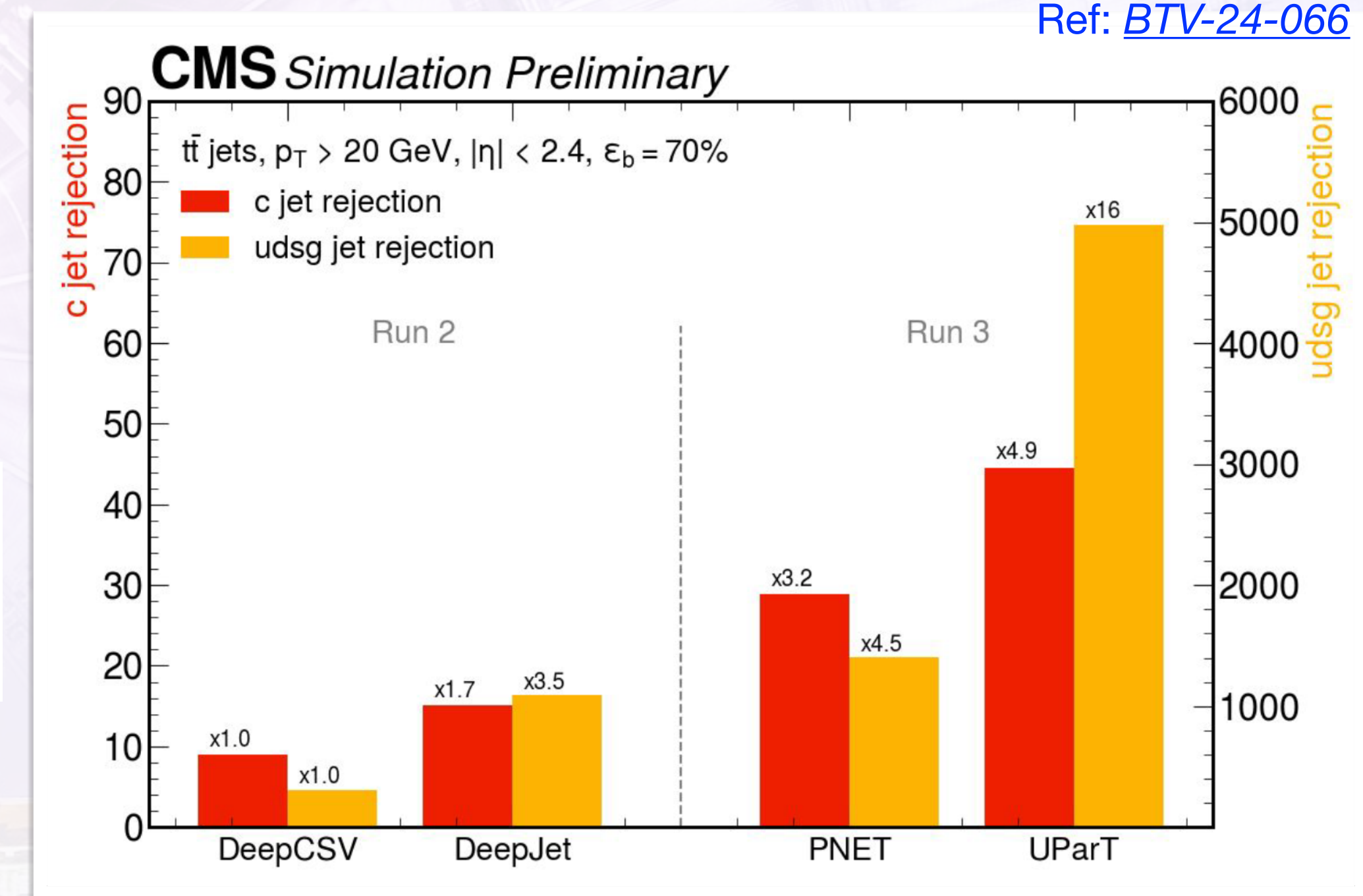
*Significant improvement as compared to DeepJet!*

# Transformer model: UParT

- Extension of ParticleTransformer
  - Extended class: extending from b and c jet identification to s and hadronic tau (one per final state) identification
  - Extended regression: simultaneous flavor aware jet energy and resolution regression

$$L = \underbrace{\text{CatEntropy}(x, x_{\text{truth}})}_{\text{Classification}} + \underbrace{\lambda \times \log(\cosh(y - y_{truth}))}_{\text{Regression}} + \underbrace{\gamma \times [\rho_{0.16}(z - z_{\text{truth}}) + \rho_{0.84}(z - z_{\text{truth}})]}_{\substack{\text{Quantile regression} \\ \text{(resolution estimation)}}}$$

- Input variable distortion:
  - Reduce the observed differences prior to any calibration
  - Improve robustness of the classifier against injected mismodelings
  ➡ Distortions of UParT: Preserving the Particle Cloud representation and the feature importance mapping



**CMS** *Simulation Preliminary*

$t\bar{t}$ jets, $p_T > 20$ GeV, $|\eta| < 2.4$, $\varepsilon_b = 70\%$
- c jet rejection
- udsg jet rejection

*Most performant heavy flavor tagging algorithm so far in CMS!*

# Flavor tagging Performance: UParT

## b-tagging

Ref: *BTV-24-066*



## c-tagging

Ref: *BTV-24-066*



*Significant improvement in b-tagging efficiency!*

*Improvement in c-tagging efficiency and c vs b discrimination*

# Flavor tagging Performance: UParT

**s-tagging**

Ref: *BTV-24-066*



**$\tau$-tagging**

Ref: *BTV-24-066*



*First attempt of s-tagging in CMS!*

*Improvement in $\tau$-tagging performance*

# Data vs. Simulations Mismodeling

- Observed **differences in Data and Simulations**
  - Imperfect modeling of the input variables affects the modeling of the output discriminator distributions
- Prone to changes in the calibration of the detector alignment
- Different estimation in simulations as compared to data

Ref: *DP2024-024*



Data/simulation disagreement due to mismodeling in simulations

**Requires calibrations**

# Calibration and Scale Factors

- Define Scale-Factors
$$SF_f = \epsilon_f^{\text{data}}(p_T, \eta)/\epsilon_f^{\text{sim}}(p_T, \eta)$$

$\epsilon_f^{\text{data}}(p_T, \eta)$ efficiencies for a jet with flavor f in data

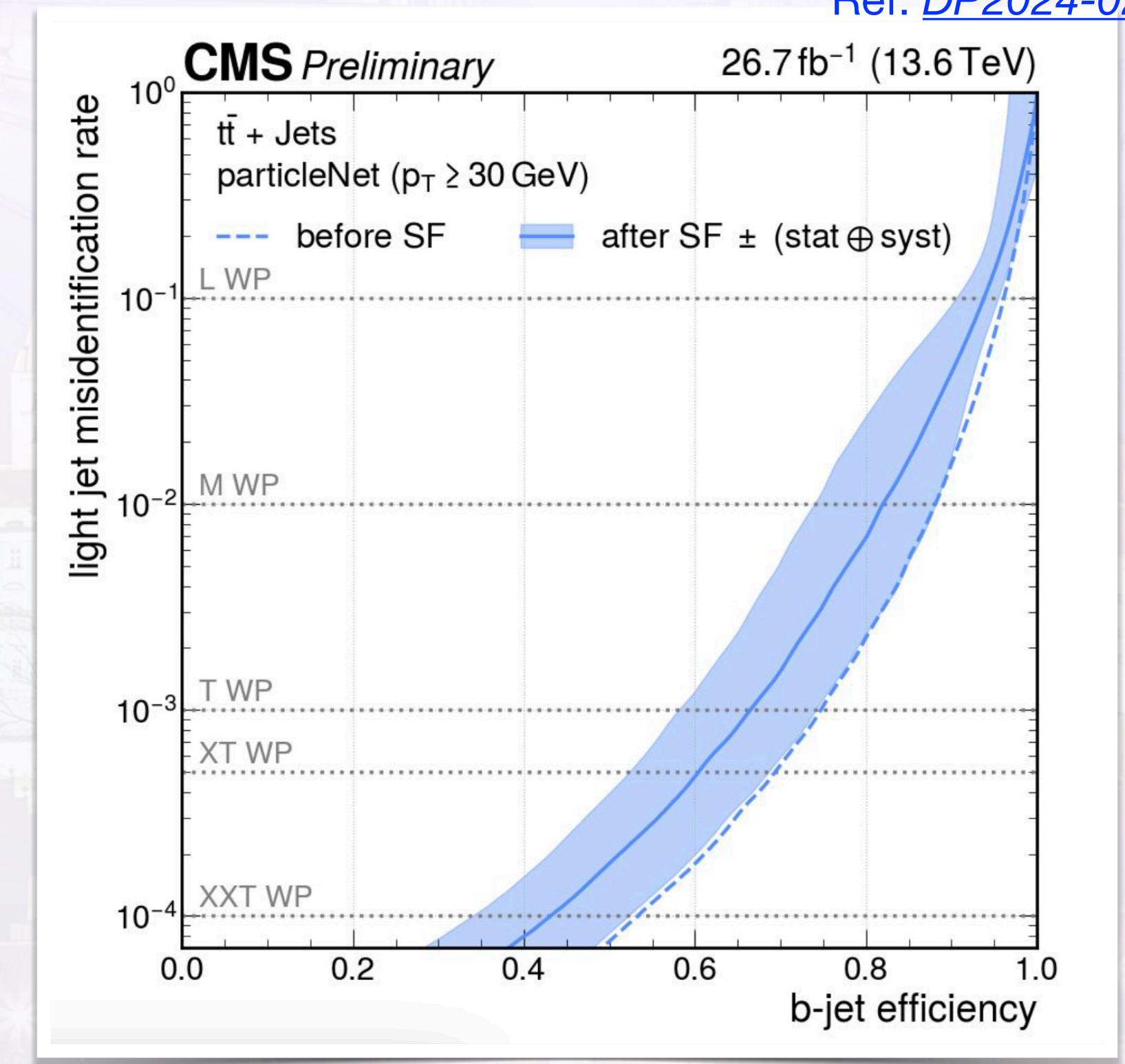$\epsilon_f^{\text{sim}}(p_T, \eta)$ efficiencies for a jet with flavor f in simulation

**Working-point calibration**

- Define heavy flavor enriched data sample to derive the SF
- For c-tagging, calibration of full discriminator shape, simultaneously for CvsL and CvsB  arXiv:2111.03027

Ref: *DP2024-025*

Ref: *DP2024-025*

**Shape calibration**



**First Run-3 SFs derived from 2022 data**



**Good Data vs. Simulations closure after applying the SFs**
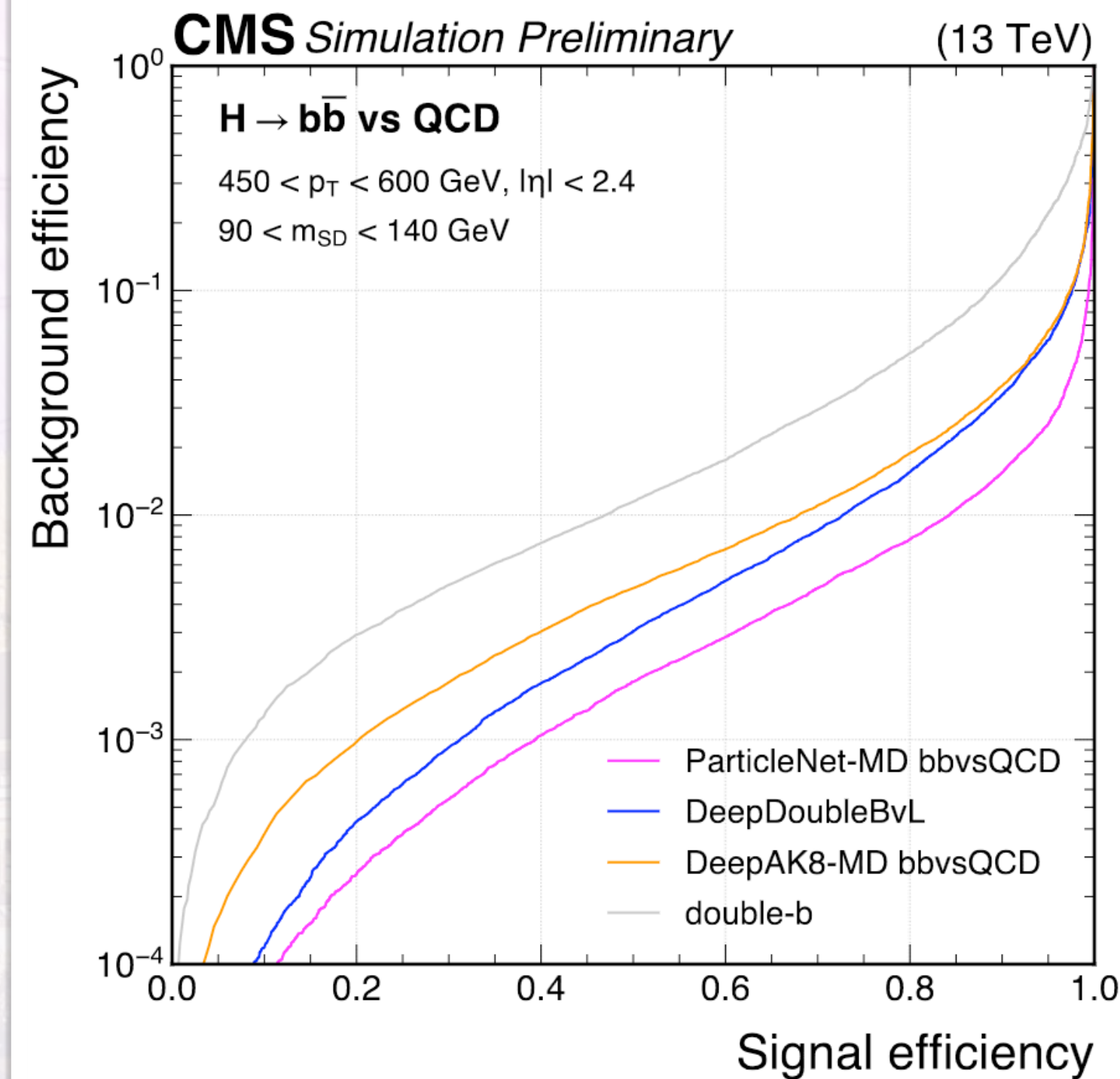
**Change in performance before and after applying the SFs**

# Boosted Object tagging

- Merged jet analyses (H→bb, H→cc) can leverage from the innovative tagging techniques of boosted objects

- Boosted-jet tagging algorithms:
  - double-b          **DNN**
  - DeepAK8MD       **CNN**
  - DeepDoubleX     **CNN, RNN**
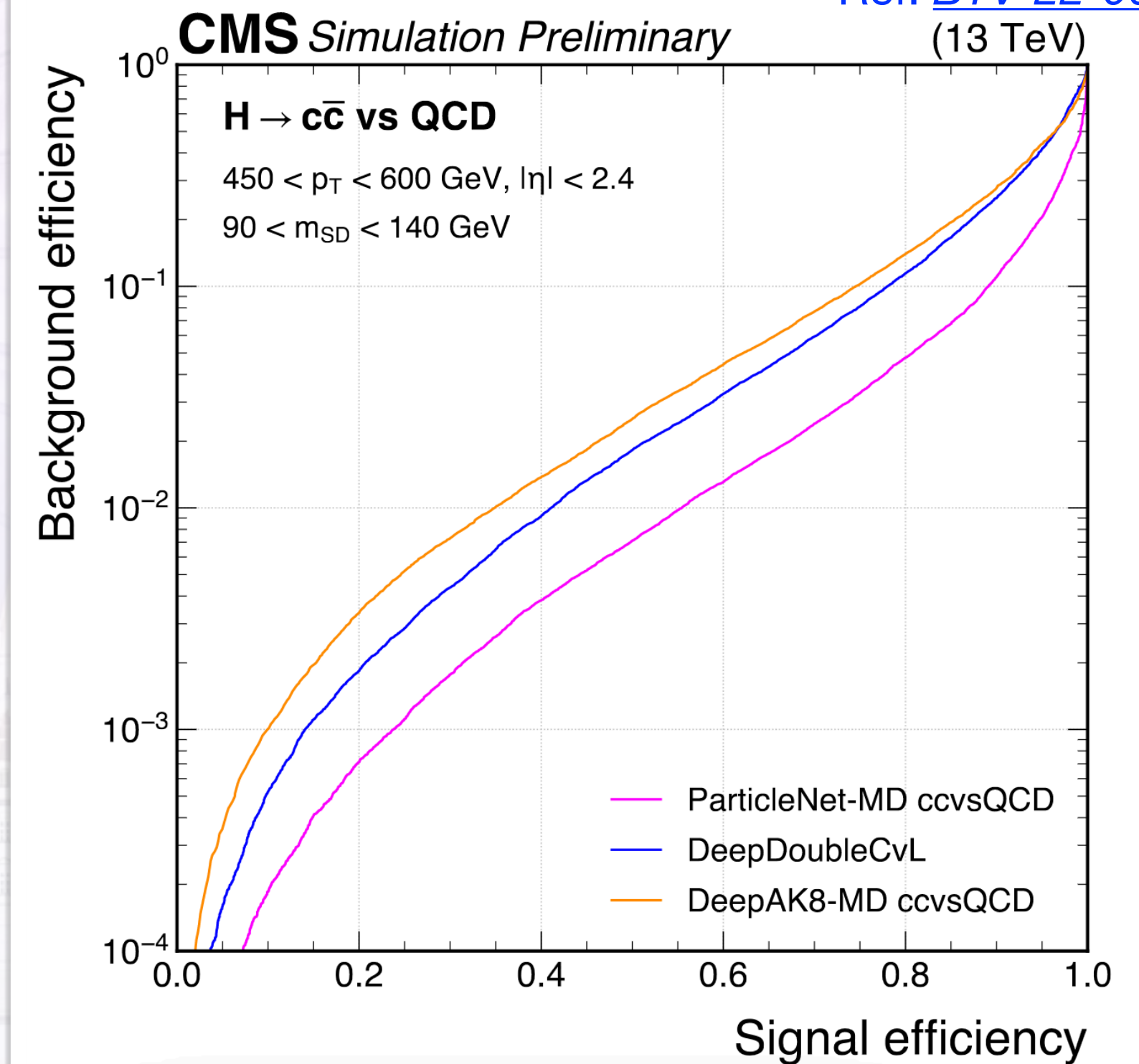  - ParticleNetMD   **DGCNN**

MD = Mass decorrelated

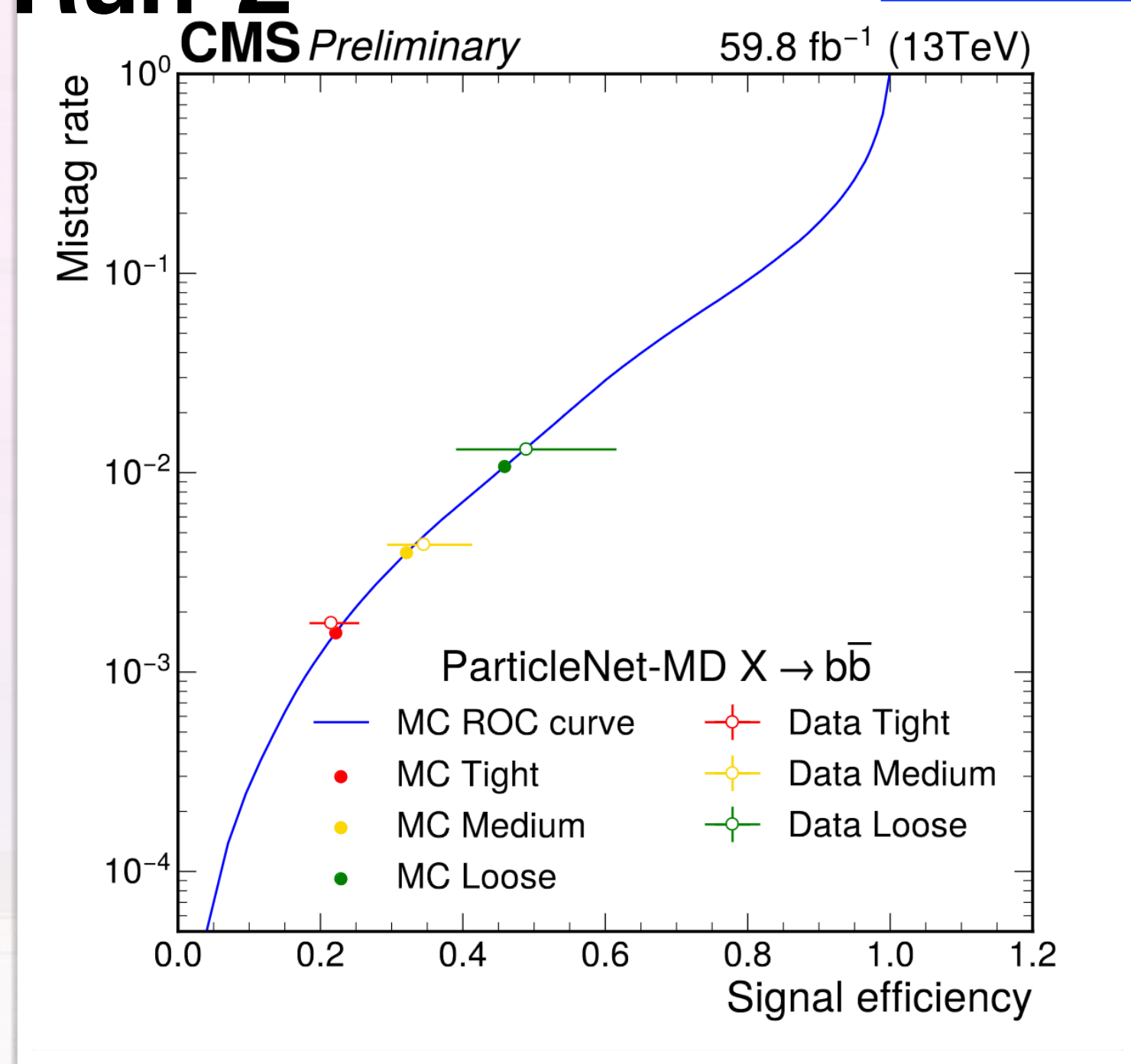**Run-2**          Ref: *BTV-22-001*



**Run-2**          Ref: *BTV-22-001*



*Improved discriminating power with the implementation of newer and better taggers*

# Calibration and Validation

## Run-2

- Calibration performed with Run-2 data
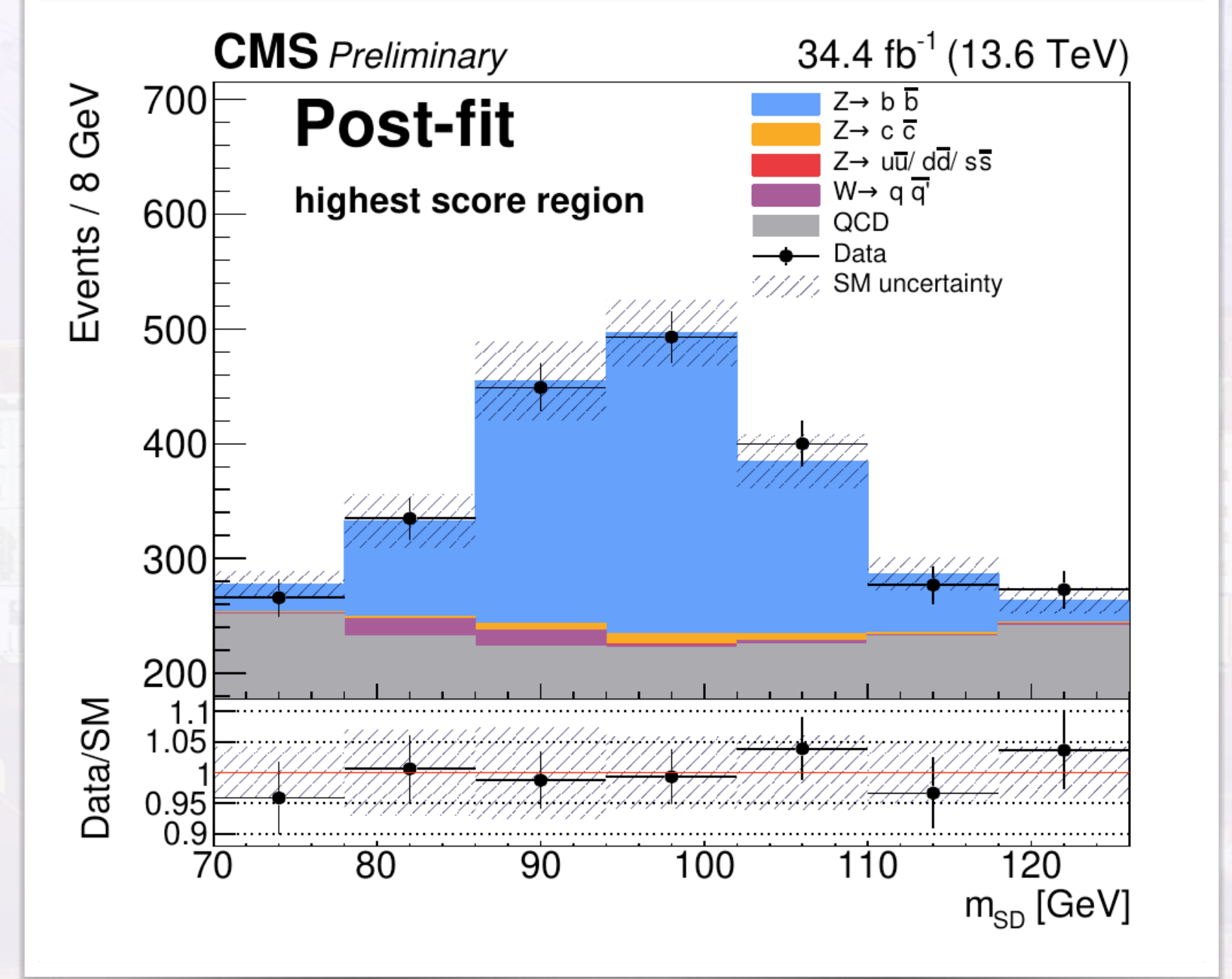- Shows the data vs. simulations ROCs of X→$b\bar{b}$ for different working points

*SFs for 2018 data taking period*

## Run-3

## Validation in Z→$b\bar{b}$ tagging in Run-3

- W→$q\bar{q}$ and QCD as main background
- Result shows the mass distribution in the highest score region (0.988 < ParticleNetMDbbvsQCD ≤ 1) of ParticleNetMD tagger

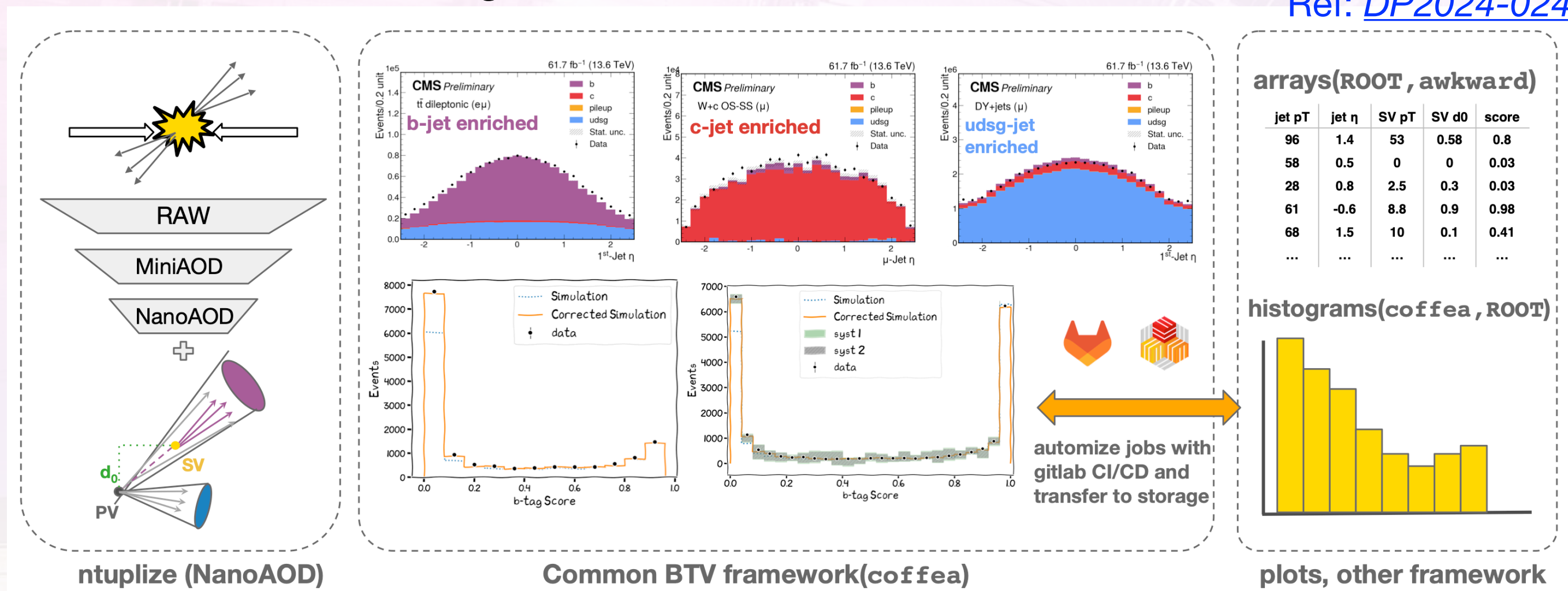Related talk in ICHEP: *"Jet performance and pileup mitigation in Run3 in CMS"*
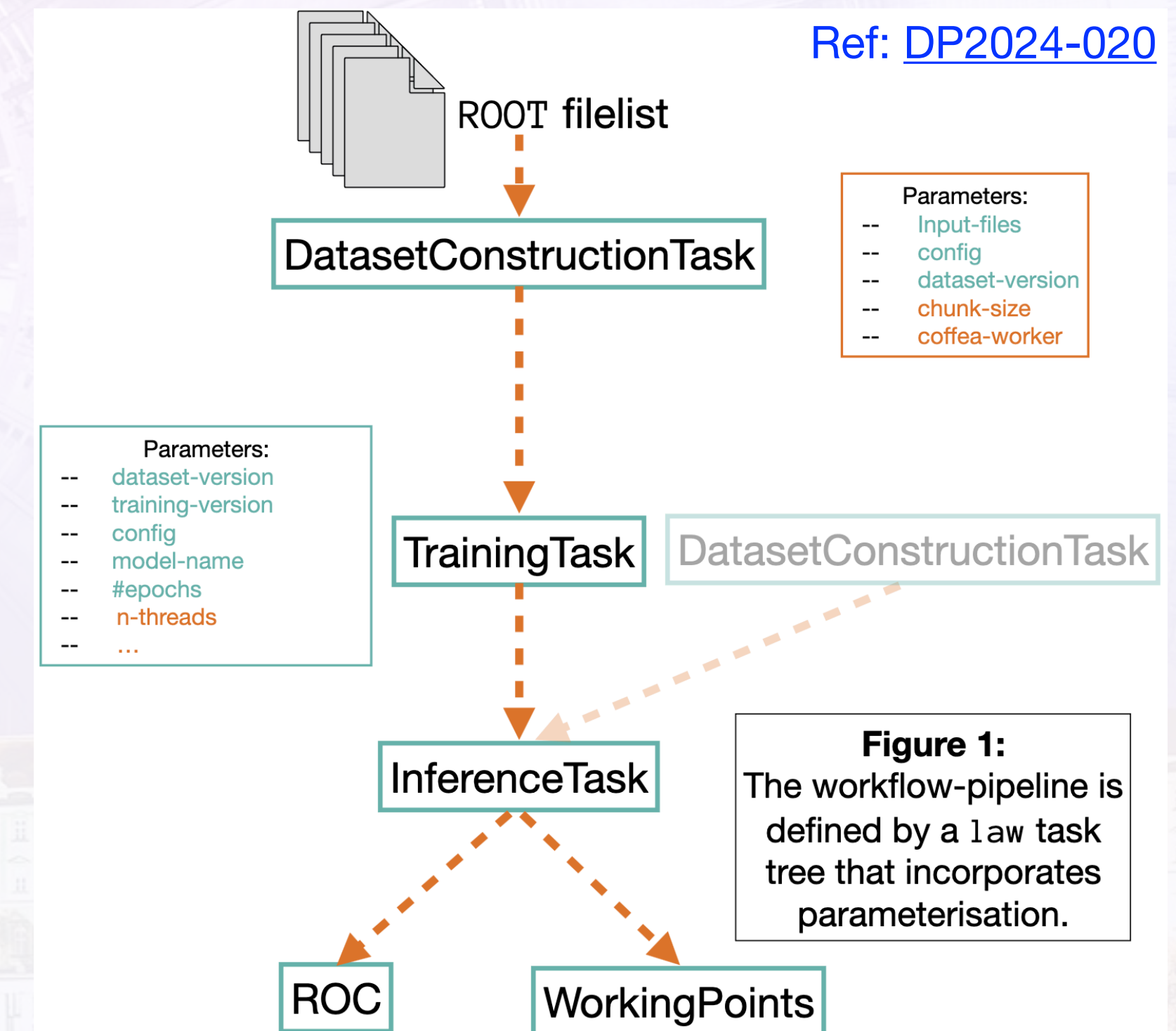
# Frameworks

## Commissioning Workflows

### BTVNanoCommissioning

Ref: *DP2024-024*



ntuplize (NanoAOD)     Common BTV framework(coffea)     plots, other framework

- Fast and efficient
- Pythonic array based manipulation instead of loops
- Automatized using Gitlab Continuous integration - monitor performance in regular intervals
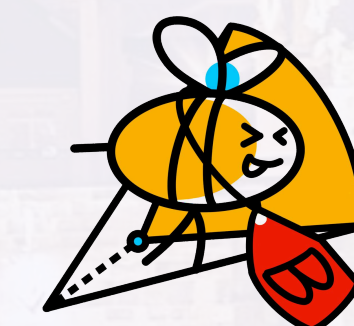
## Training
### B-Hive

Ref: *DP2024-020*



**Figure 1:** The workflow-pipeline is defined by a `law` task tree that incorporates parameterisation.

- Framework dedicated for training
- Easily customizable to introduce your own model

# Summary

- **Rapid Development of Machine Learning Architectures (2017-2023)**:
  - Evolution from **CSVv1** to **UParT -** increasing complexity and capability of models
- **Innovative Calibration Techniques in Run-2**

- **Improvement in Jet Tagging Performance:**
  - ParticleNet algorithm deployed in the online High Level Trigger (Ref: DP2023-021 )
  - Notable advancements from Run-1 to the present
  - Demonstrates significant gains in accuracy and efficiency

- **Ongoing Enhancement in Tagging and Calibration Methods:**
  - Continued development expected in Run-3
  - Aim to further refine and optimize performance
  - New taggers and many more, stay tuned!
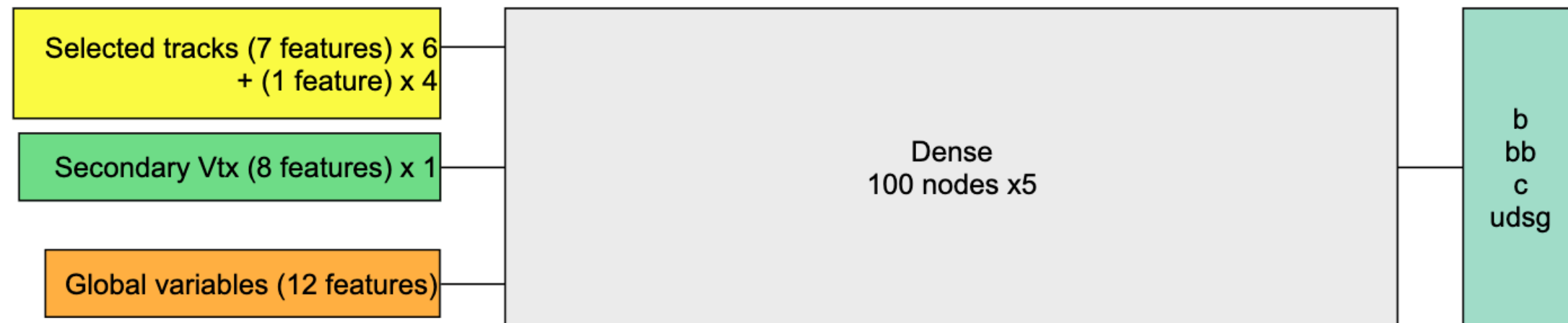
**Thank you for listening!**

Related poster from Donato Troiano: *Identification of Lorentz-boosted jets in the CMS experiment*
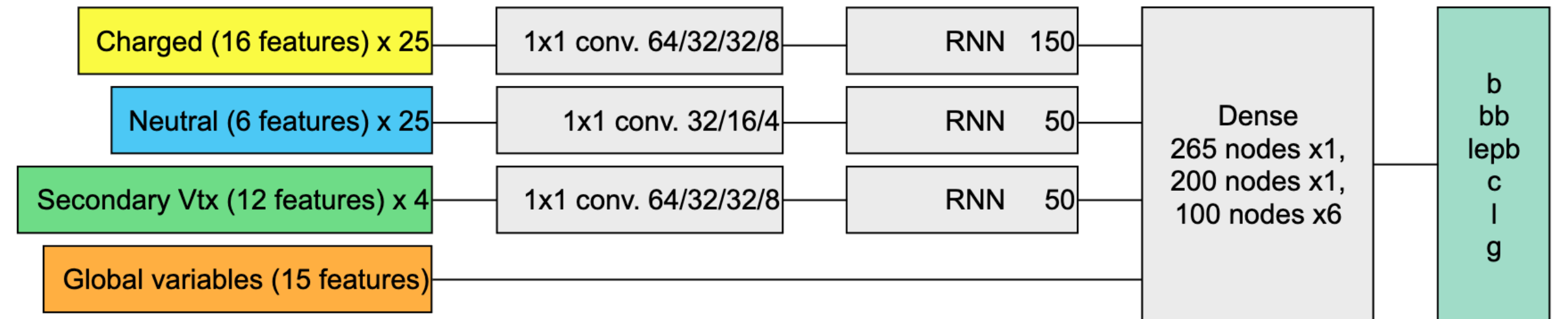
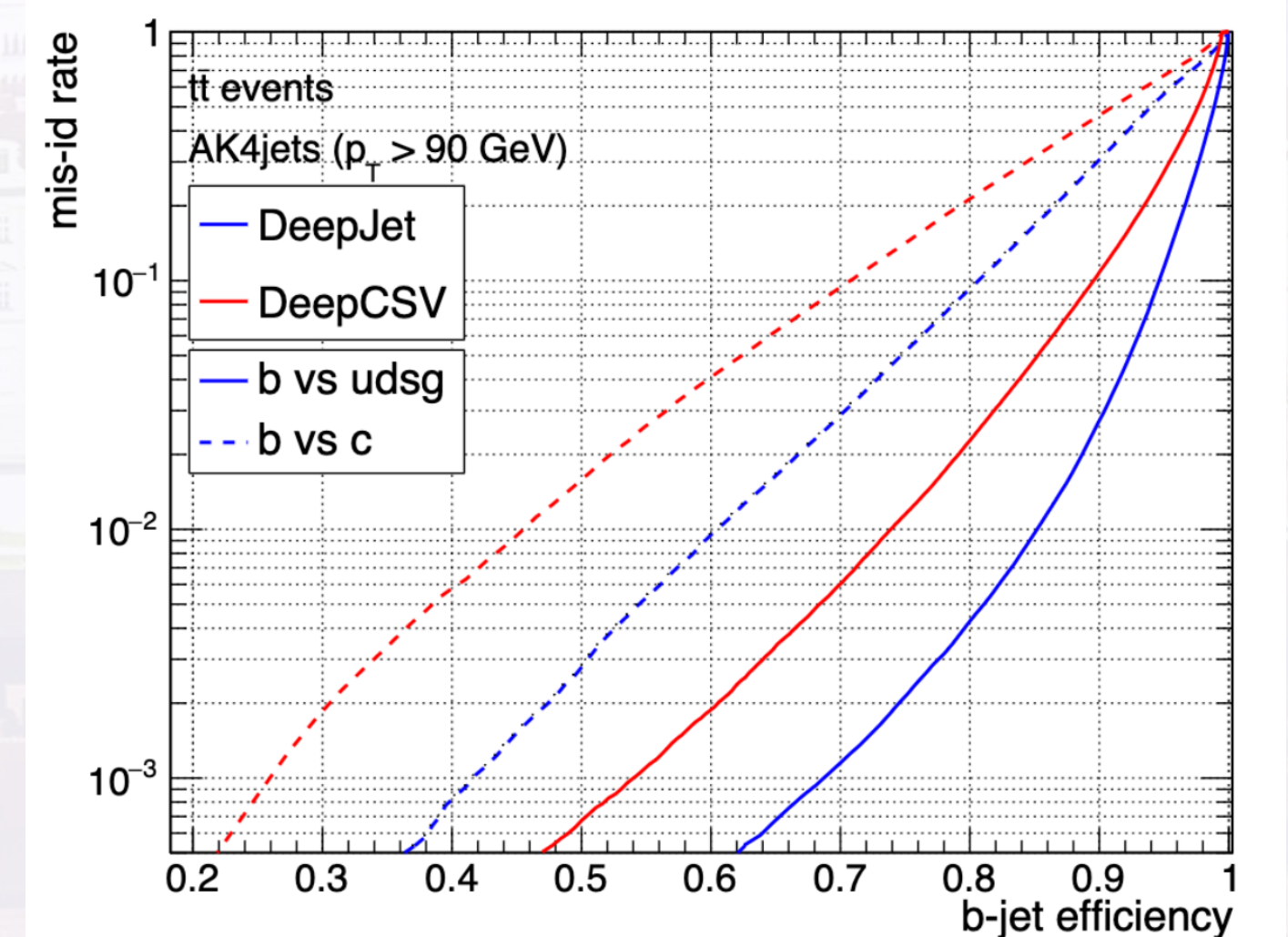# BACKUP

# DeepCSV and DeepJet

**DeepCSV[1]**

| Selected tracks (7 features) x 6 + (1 feature) x 4 | | |
|---|---|---|
| Secondary Vtx (8 features) x 1 | Dense 100 nodes x5 | b bb c udsg |
| Global variables (12 features) | | |

**DeepJet[2]**

| Charged (16 features) x 25 | 1x1 conv. 64/32/32/8 | RNN 150 | | |
|---|---|---|---|---|
| Neutral (6 features) x 25 | 1x1 conv. 32/16/4 | RNN 50 | Dense 265 nodes x1, 200 nodes x1, 100 nodes x6 | b bb lepb c l g |
| Secondary Vtx (12 features) x 4 | 1x1 conv. 64/32/32/8 | RNN 50 | | |
| Global variables (15 features) | | | | |

- **Fully connected neural network (dense)**
- Combines properties from **selected** tracks, secondary vertices and global variables directly (66 features)
- Only a small subset of the charged jet constituents pass stringent quality criteria
  - clean and simple environment for the classifier
  - **information loss - potential performance degradation**

- **Convolution, RNN,** and **Dense** layers
- **Does not rely on a selection** of the jet constituents
  - better purity, more number of inputs
- Full information of all **jet constituents, charged and neutral particles, secondary vertices, and global event variables simultaneously**
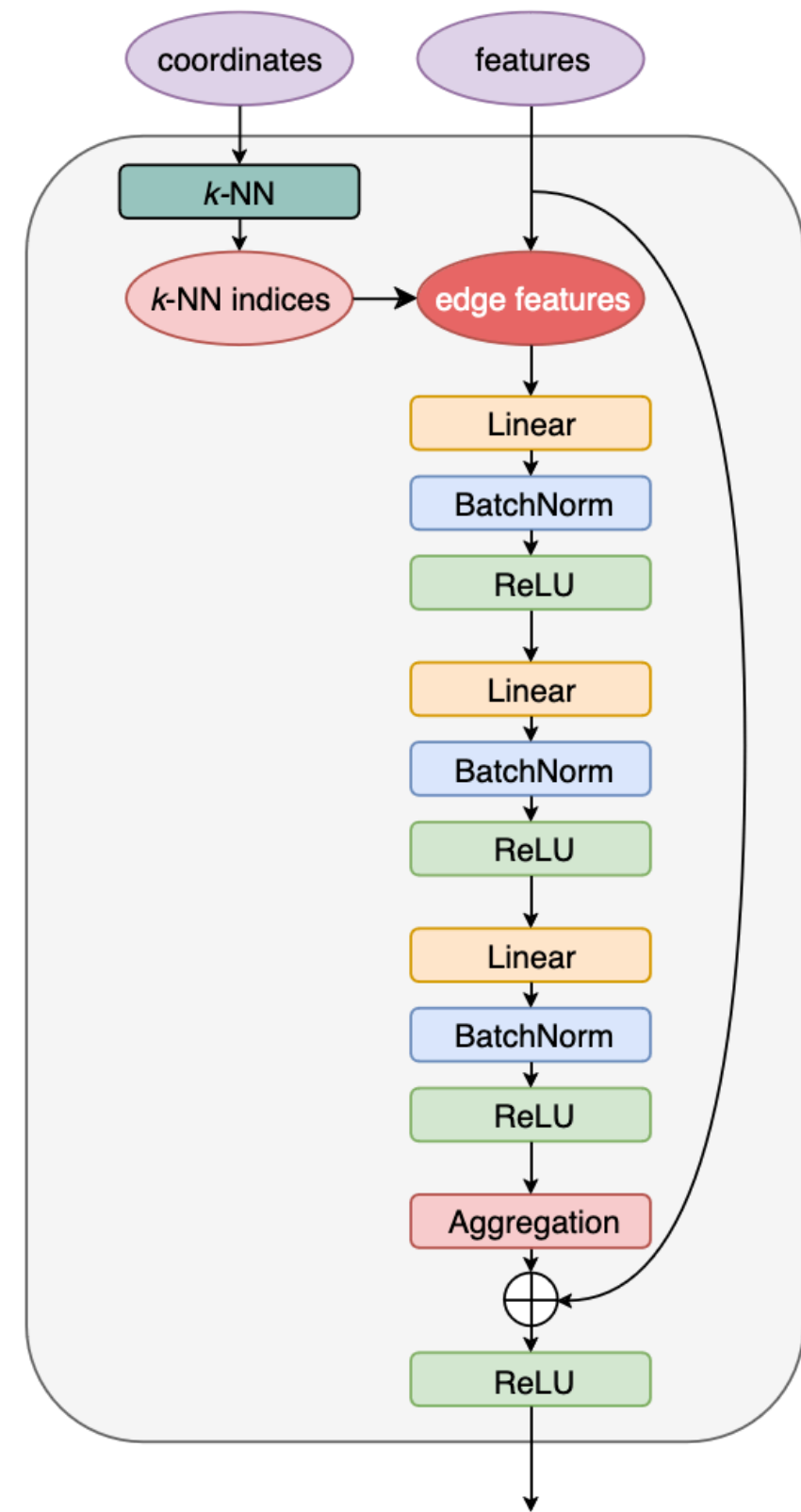
~20% gain in efficiency at $10^{-3}$ misidentification probability[2]
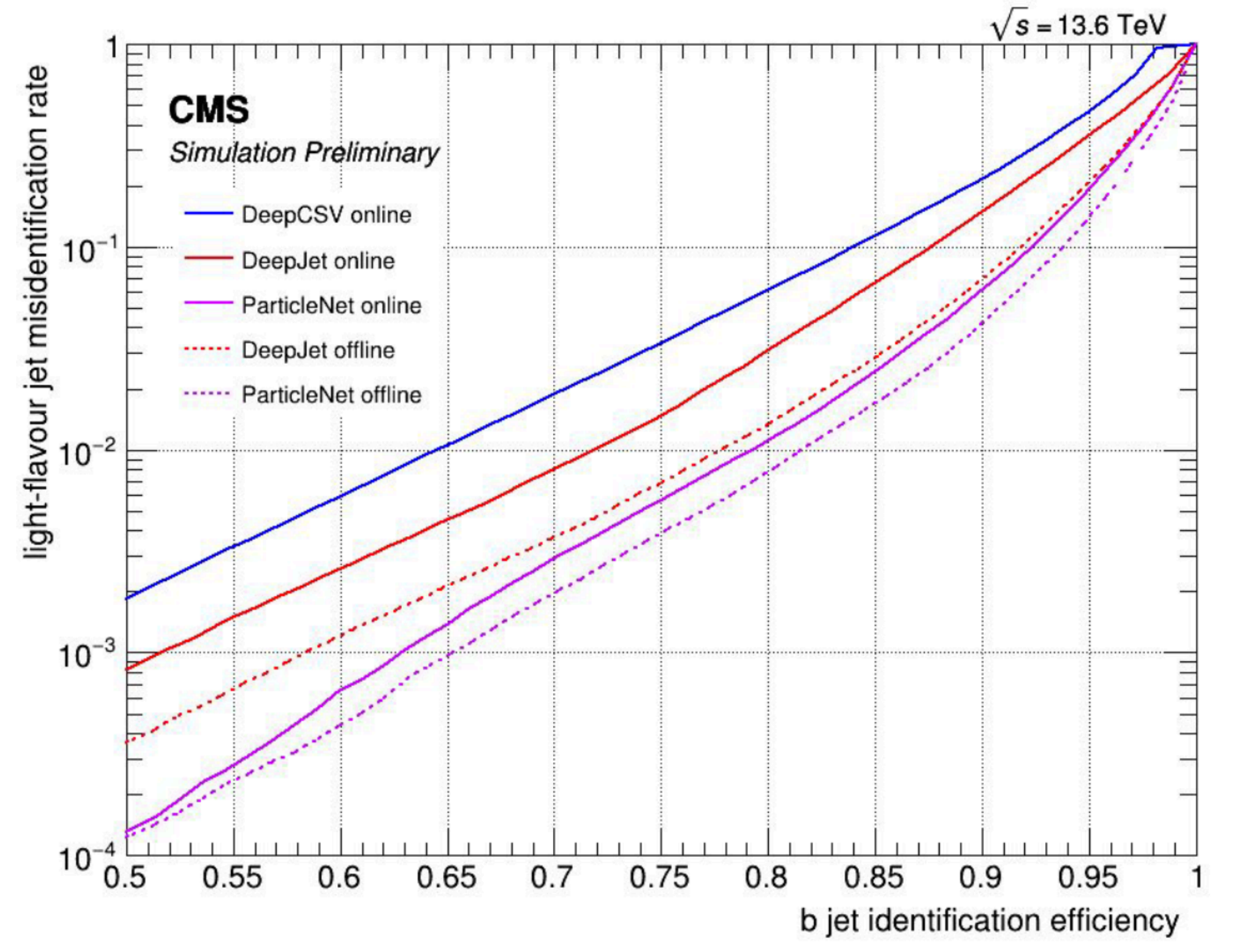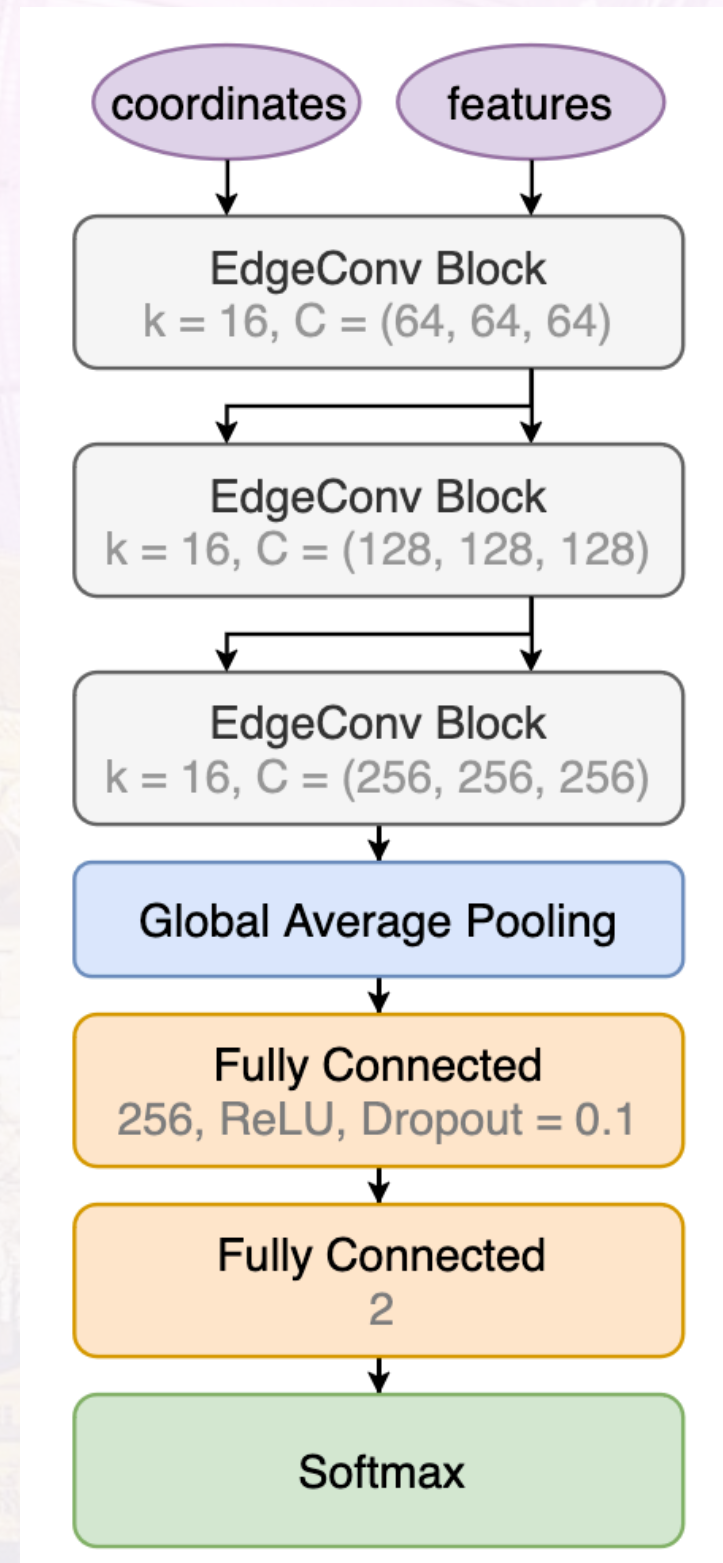
[1] 10.1103/PhysRevD.94.112002
[2] 10.1088/1748-0221/15/12/P12012

# ParticleNet [1]

Treat a jet as unordered sets of constituent particles

Perform Edge-convolution and Dynamic Graph Convolutional Neural Network (DGCNN)
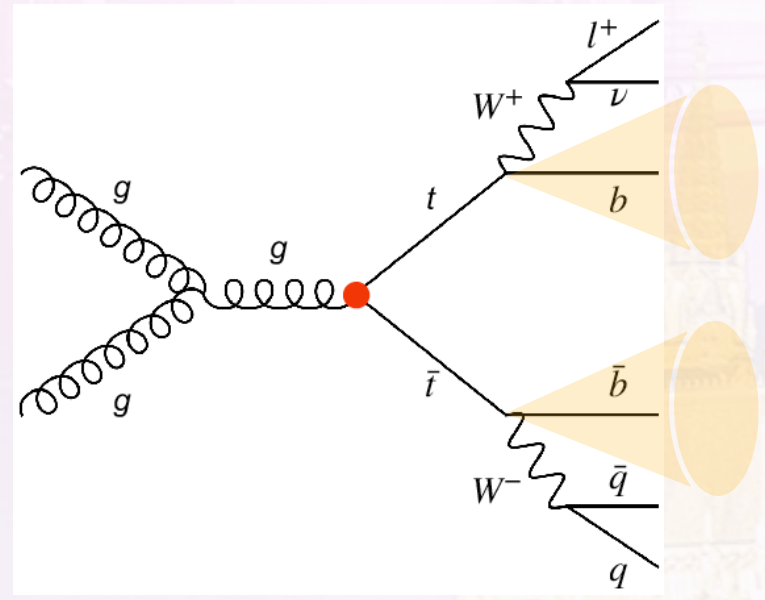


One edge-convolution block





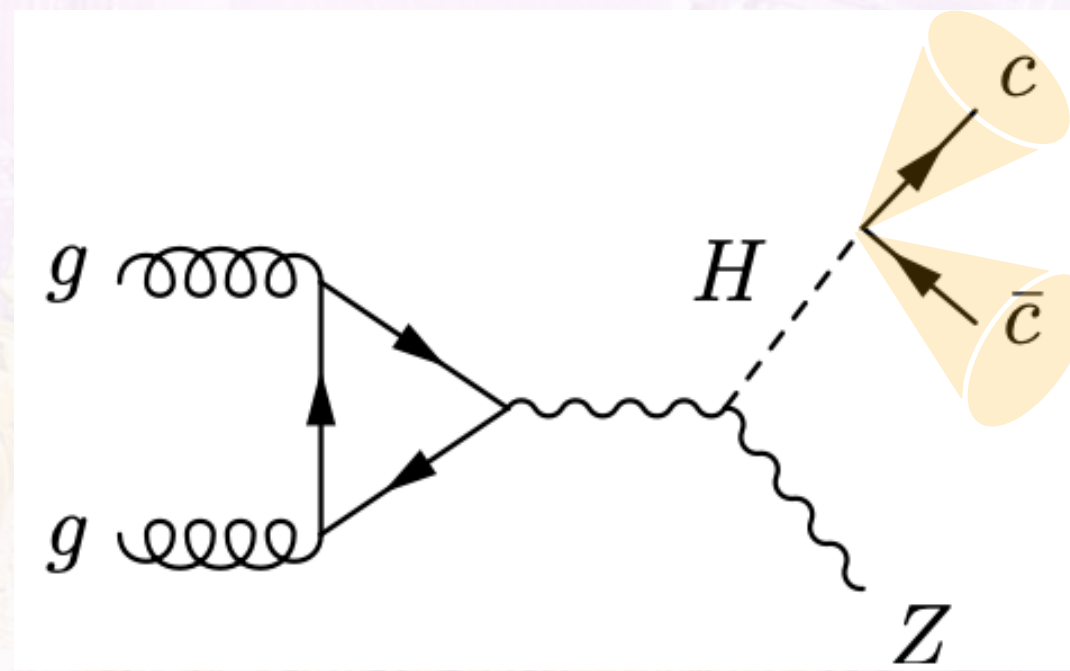~10% gain in efficiency at $10^{-3}$ misidentification probability[2]
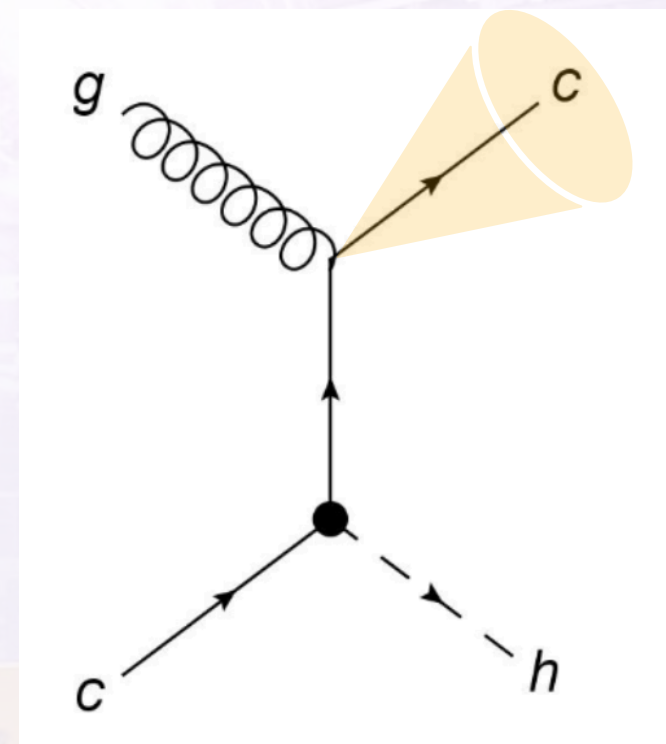
[1] arxiv1902.08570
[2] CMS_DP2023_021

- **Heavy flavor jets =** jets originating from b (*b jets*) or c (*c jets*) quarks arising from the process of hadronization
- Important in Standard Model (SM), Top, Higgs(H->bb,cc), BSM and SUSY processes
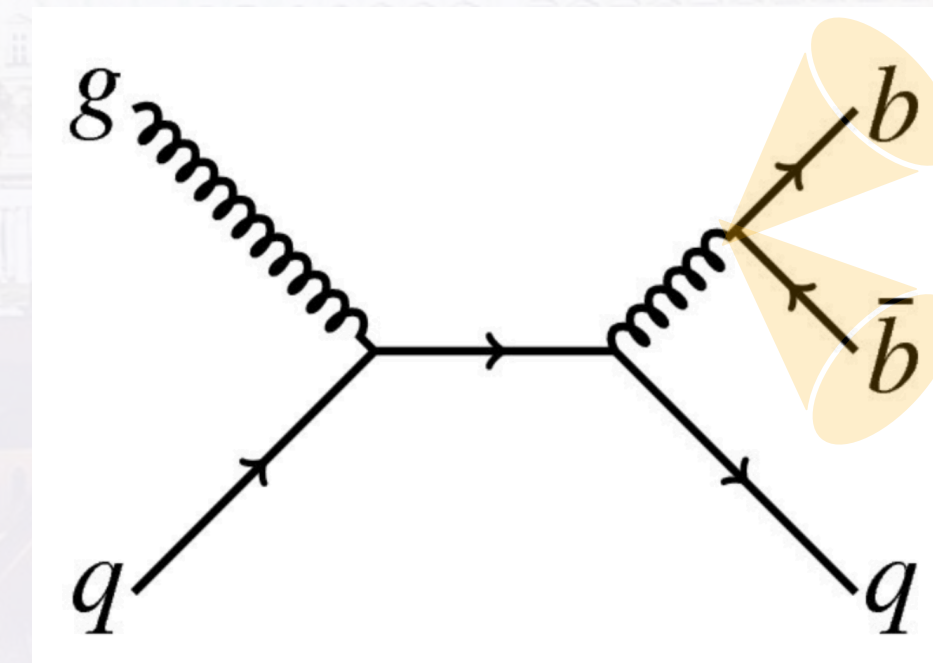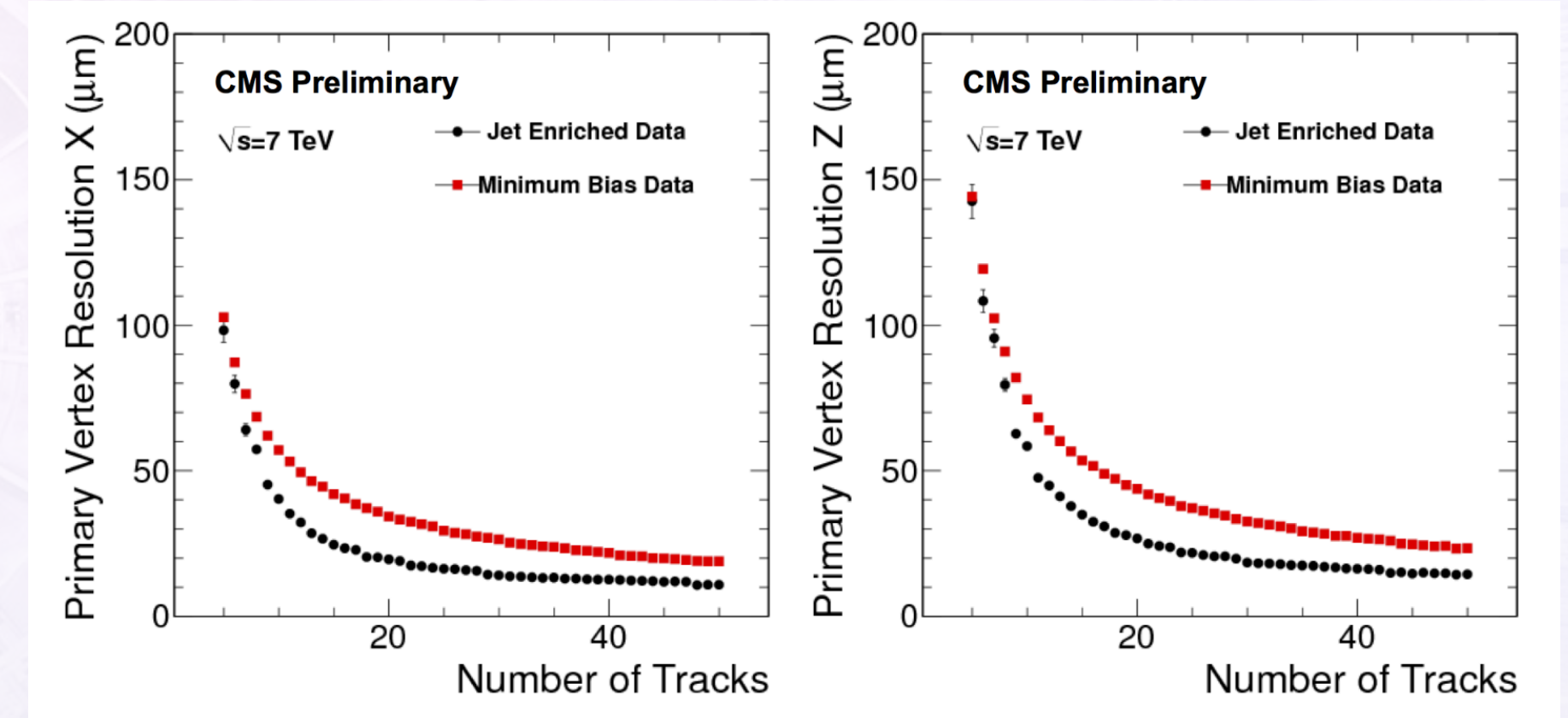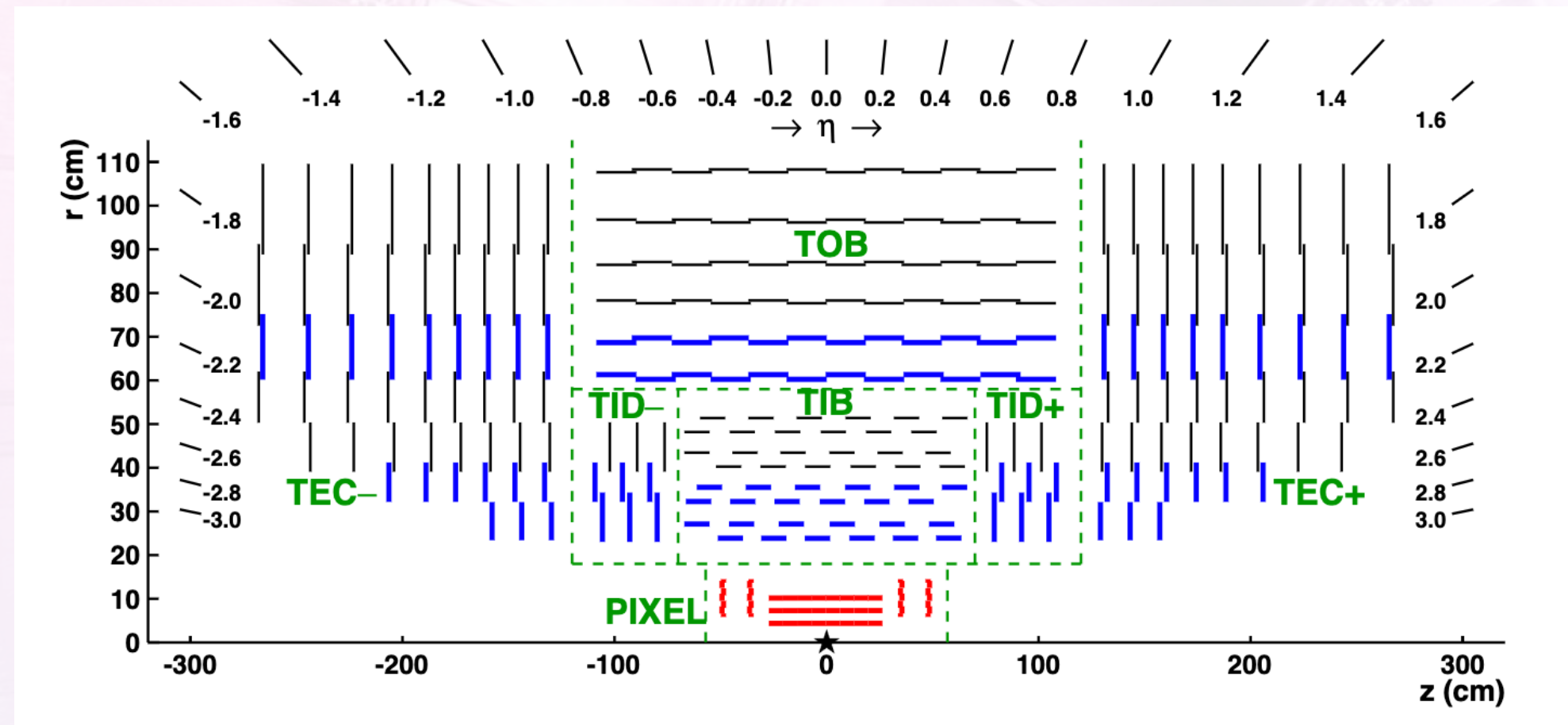


$t\bar{t}$



VHcc[1]



Higgs+c[2]



SUSY stop-> SM + MET

- **QCD:** understanding heavy-parton effects



[1] https://www.dpg-verhandlungen.de/year/2024/conference/karlsruhe/part/t/session/20/contribution/6
[2] https://www.dpg-verhandlungen.de/year/2024/conference/karlsruhe/part/t/session/71/contribution/7

# CMS Tracking





| Subdetector | Radius [cm] | Sensor size [$\mu$m] | Resolution [$\mu$m] | <Hits on track> |
|---|---|---|---|---|
| Pixel | 4.4-10.2 | 100×150 | R$\phi$:10 z:20 | 3 |
| Strip tracker | 25.5-110 | $\sim 100$ | $\sim 15- \sim 45$ | 13 |

https://pos.sissa.it/190/041/pdf

Ideal to observe in CMS, though challenging!

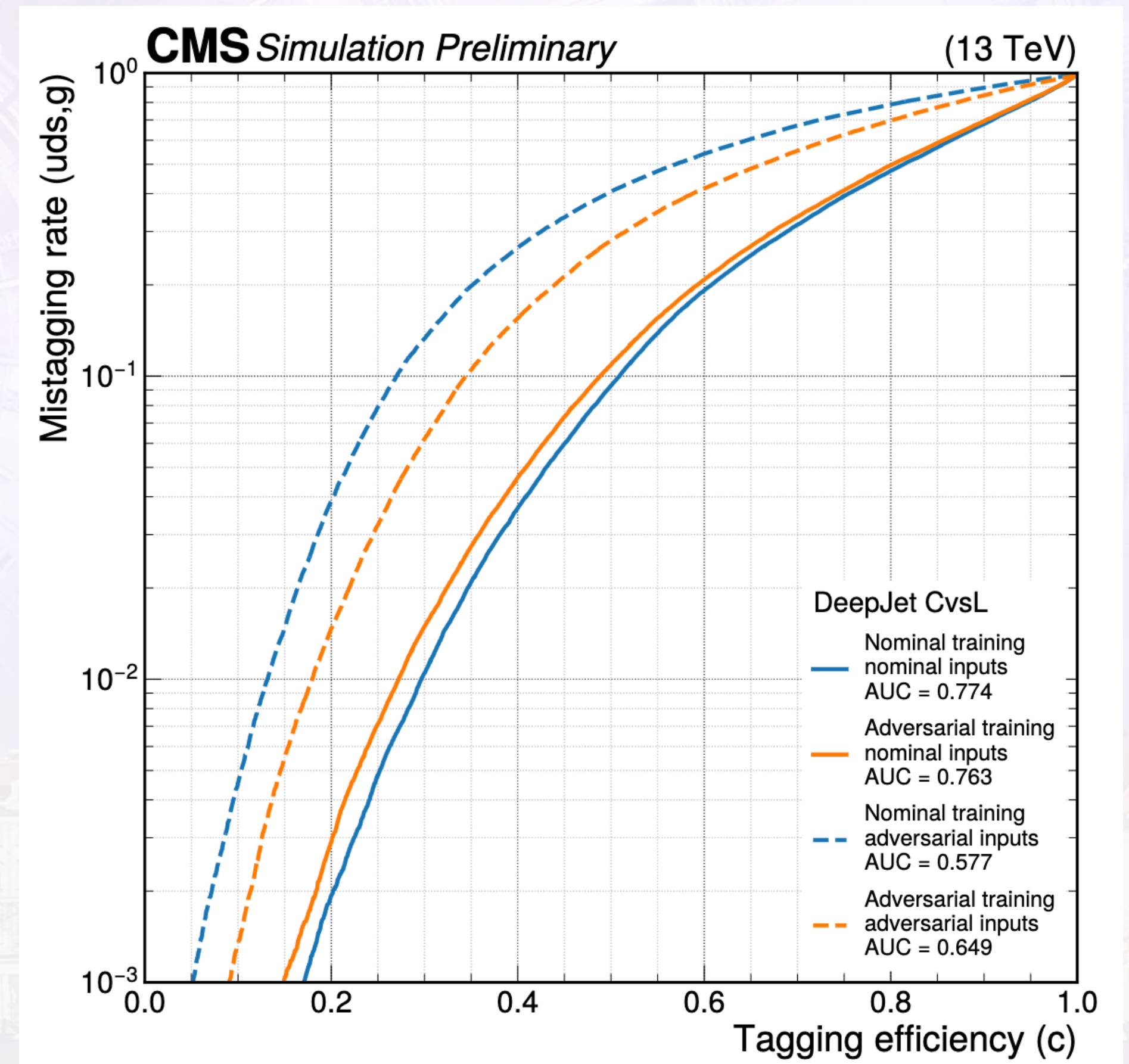[1] http://cds.cern.ch/record/1279383/files/TRK-10-005-pas.pdf

# Adversarial Attacks

- Adversarial training strategy:
  - Reduce the observed differences prior to any calibration,
  - Improve robustness of the classifier against injected mismodelings

The Fast Gradient Sign Method (FGSM) is used to systematically distort inputs
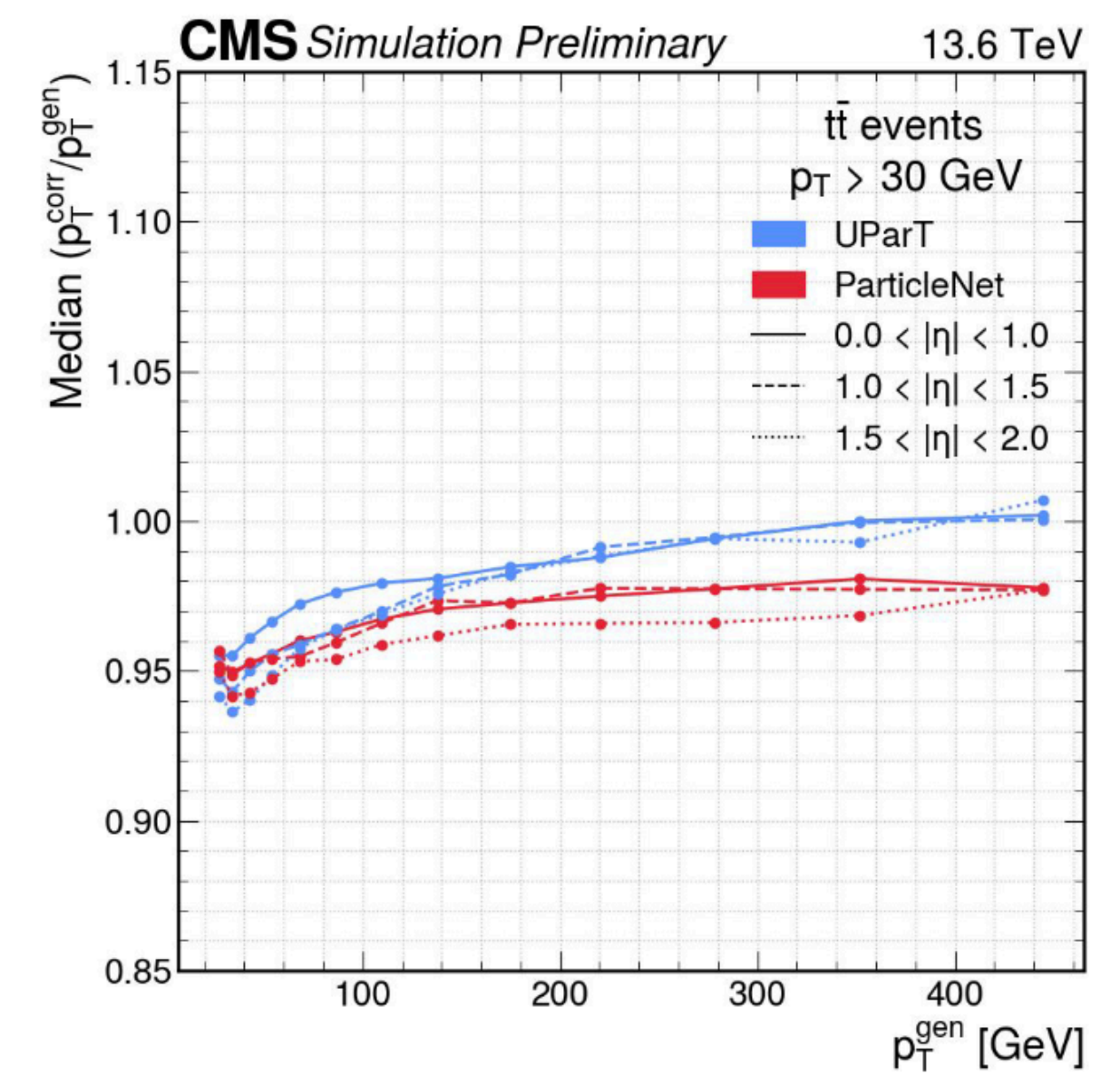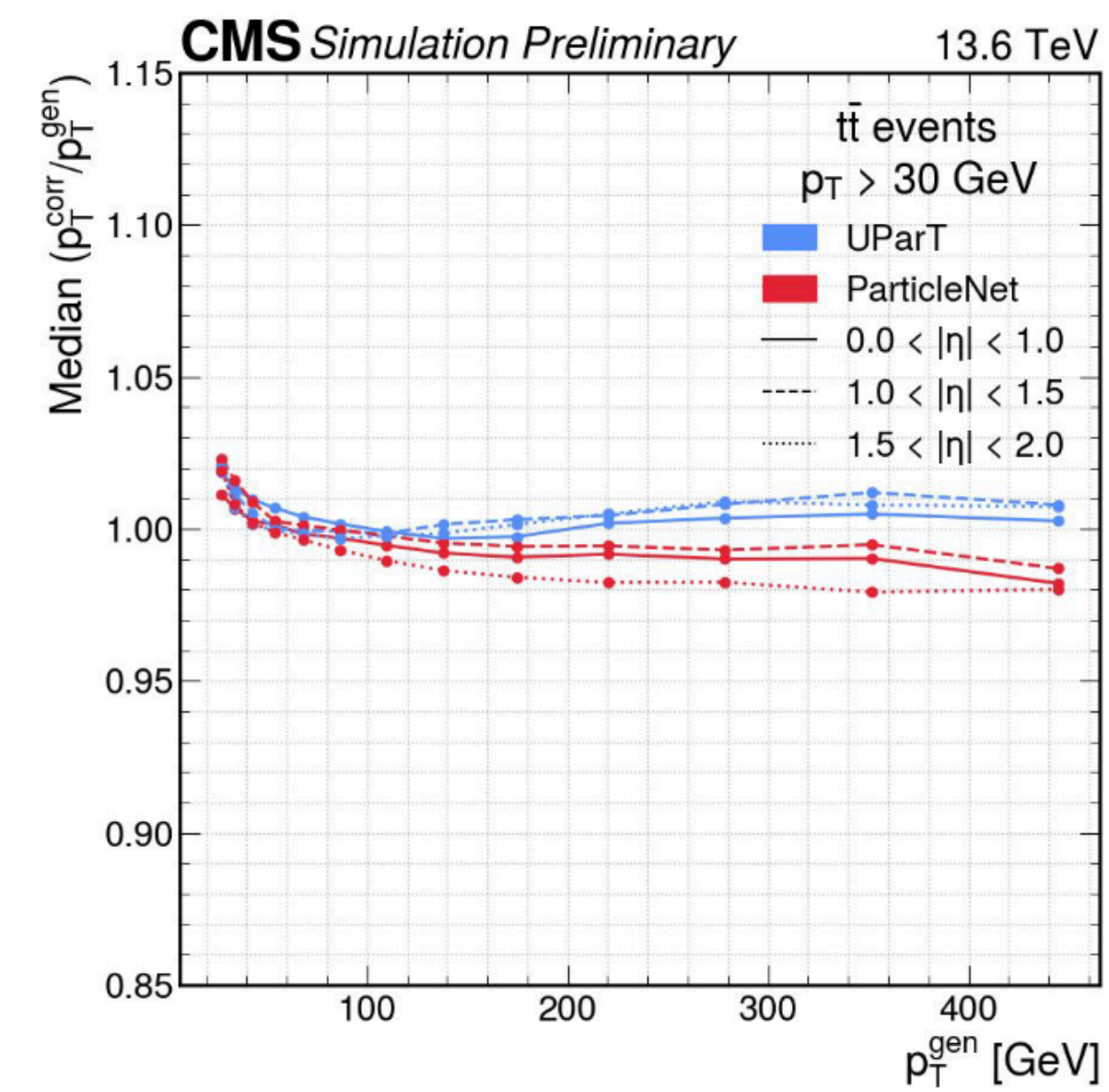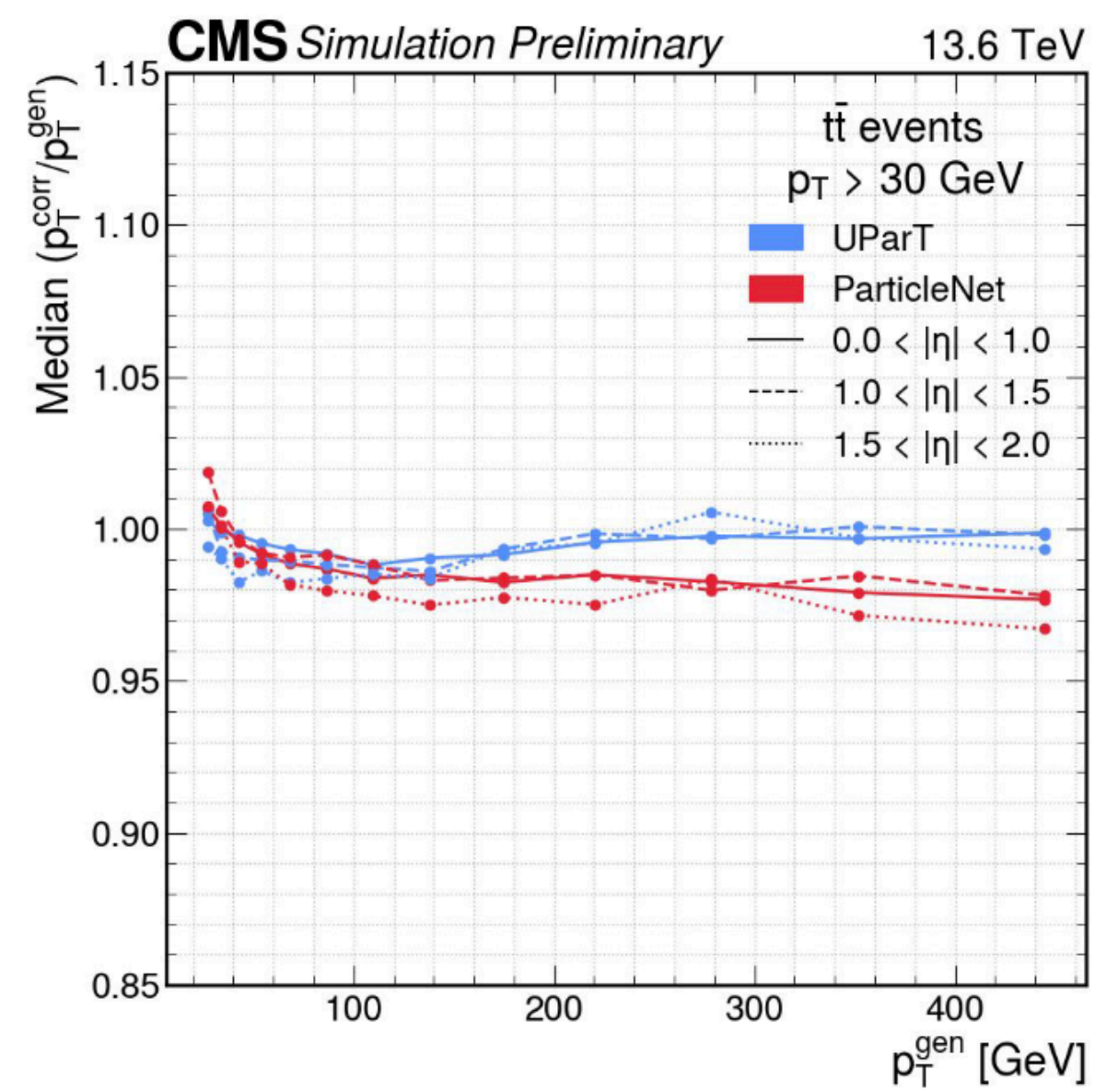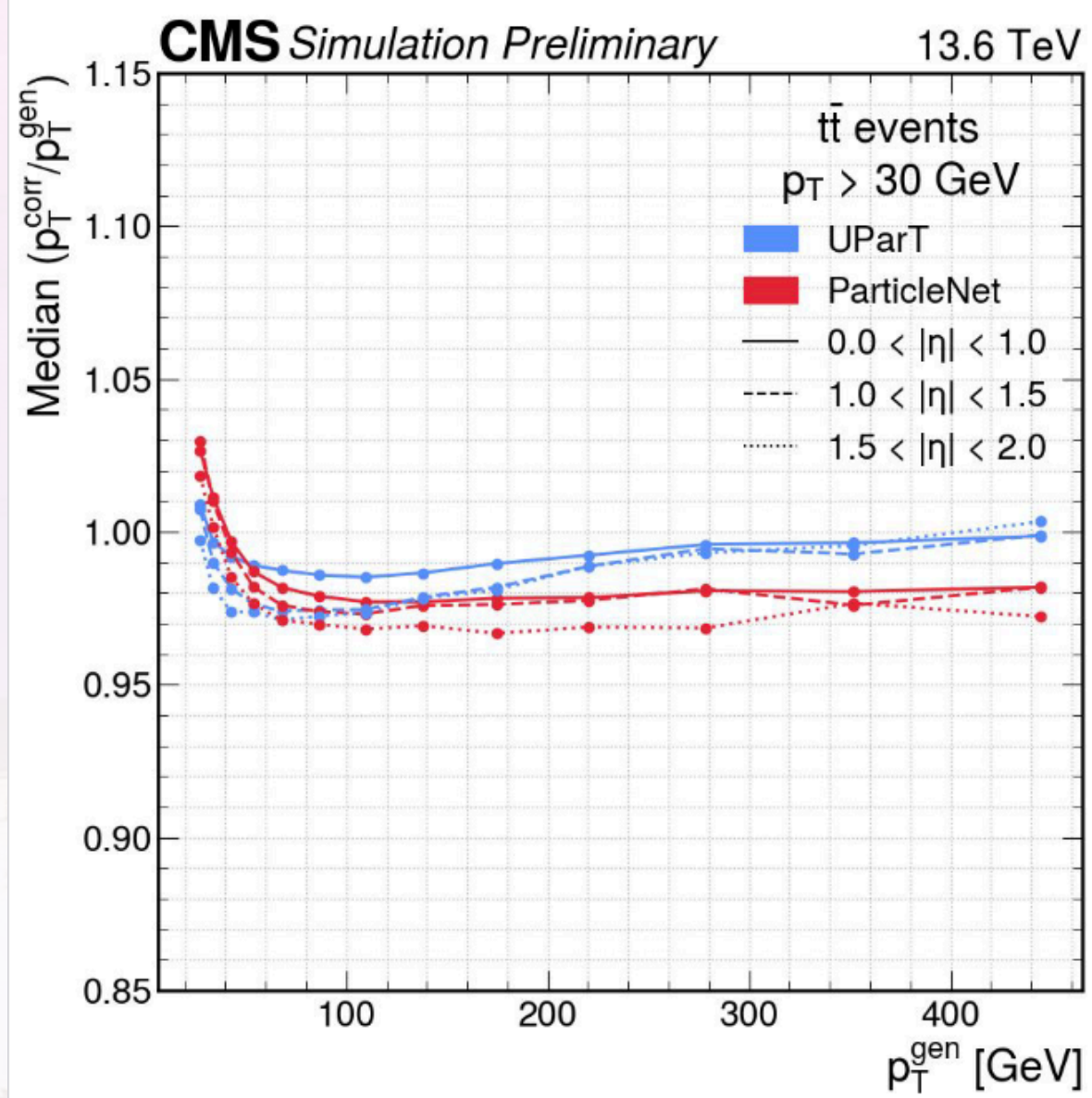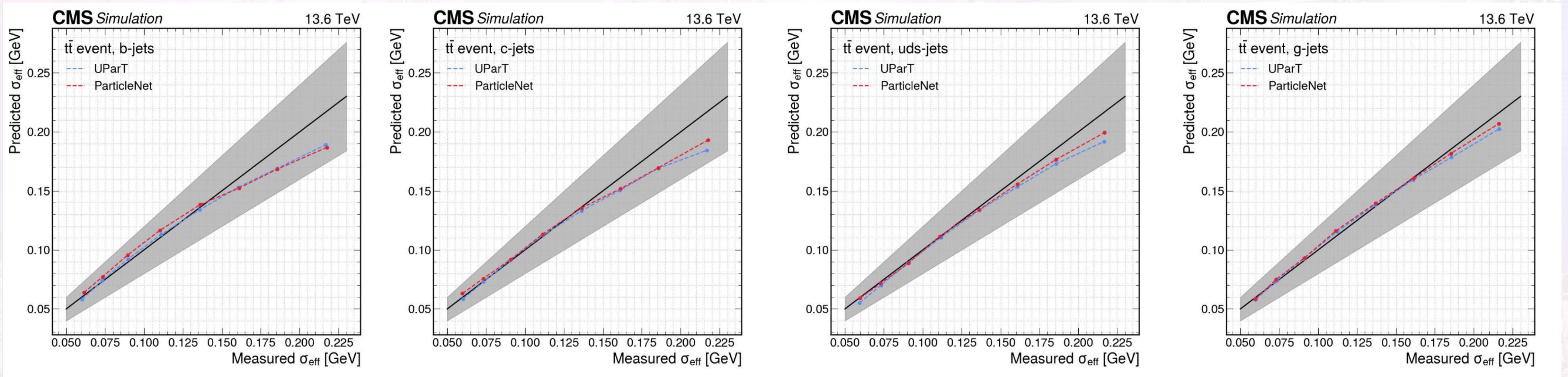


DeepJet architecture

# Jet energy regression: UParT

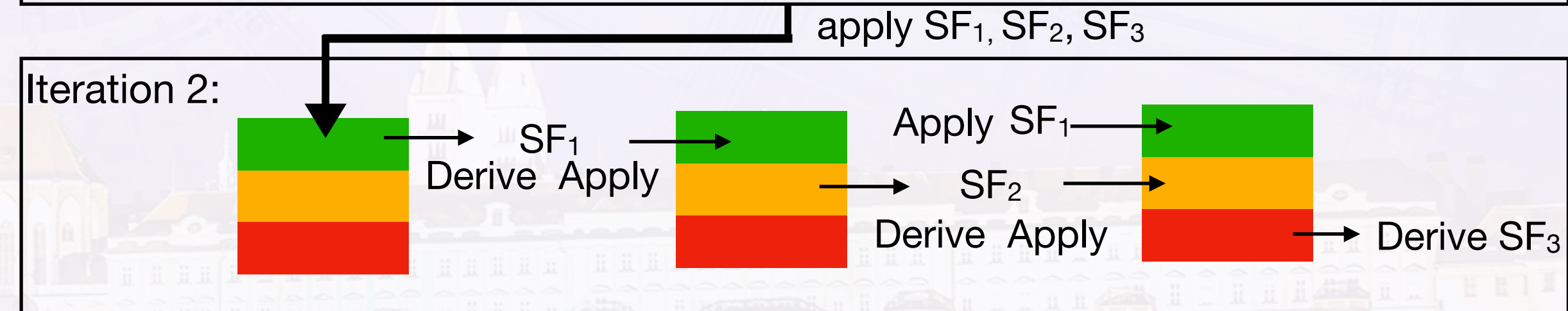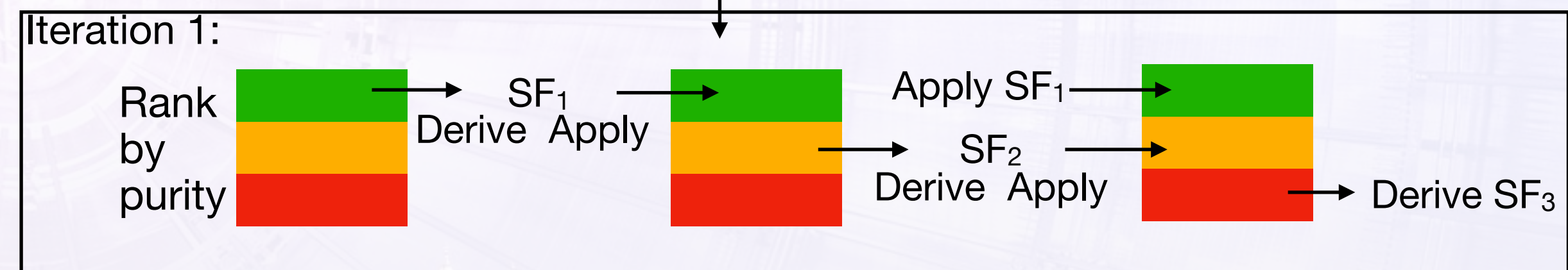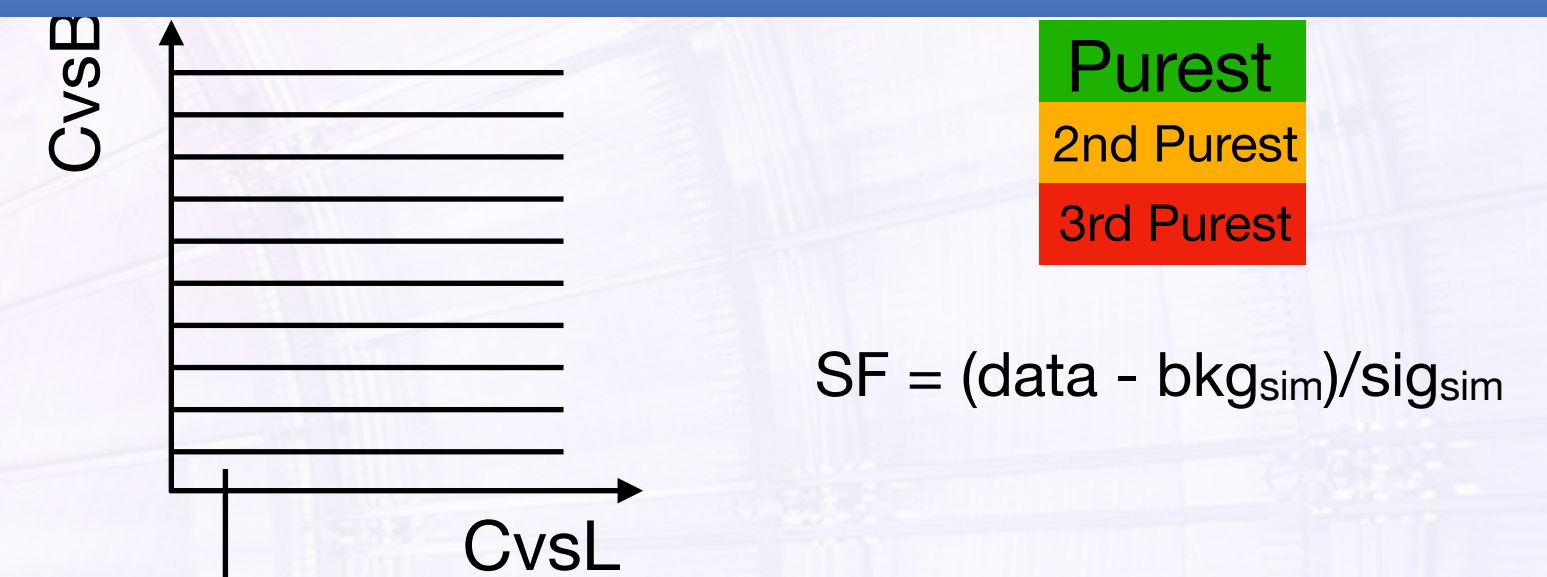# Jet energy resolution: UParT

Ref: *BTV-24-XXX*

# Iterative Fit: c tagging

- Fit iteratively in the 2D plane of CvsB-CvsL over three selections allowing each flavour component to vary independently

  - Divide CvsB-CvsL plane into 10 "bin slices" along the CvsB axis
  - Variable binning along CvsL axis - perform fit and derive SF in each bin

  Apply $SF_c$ to the same distributions they were derived from

  - Propagate uncertainties (from data):
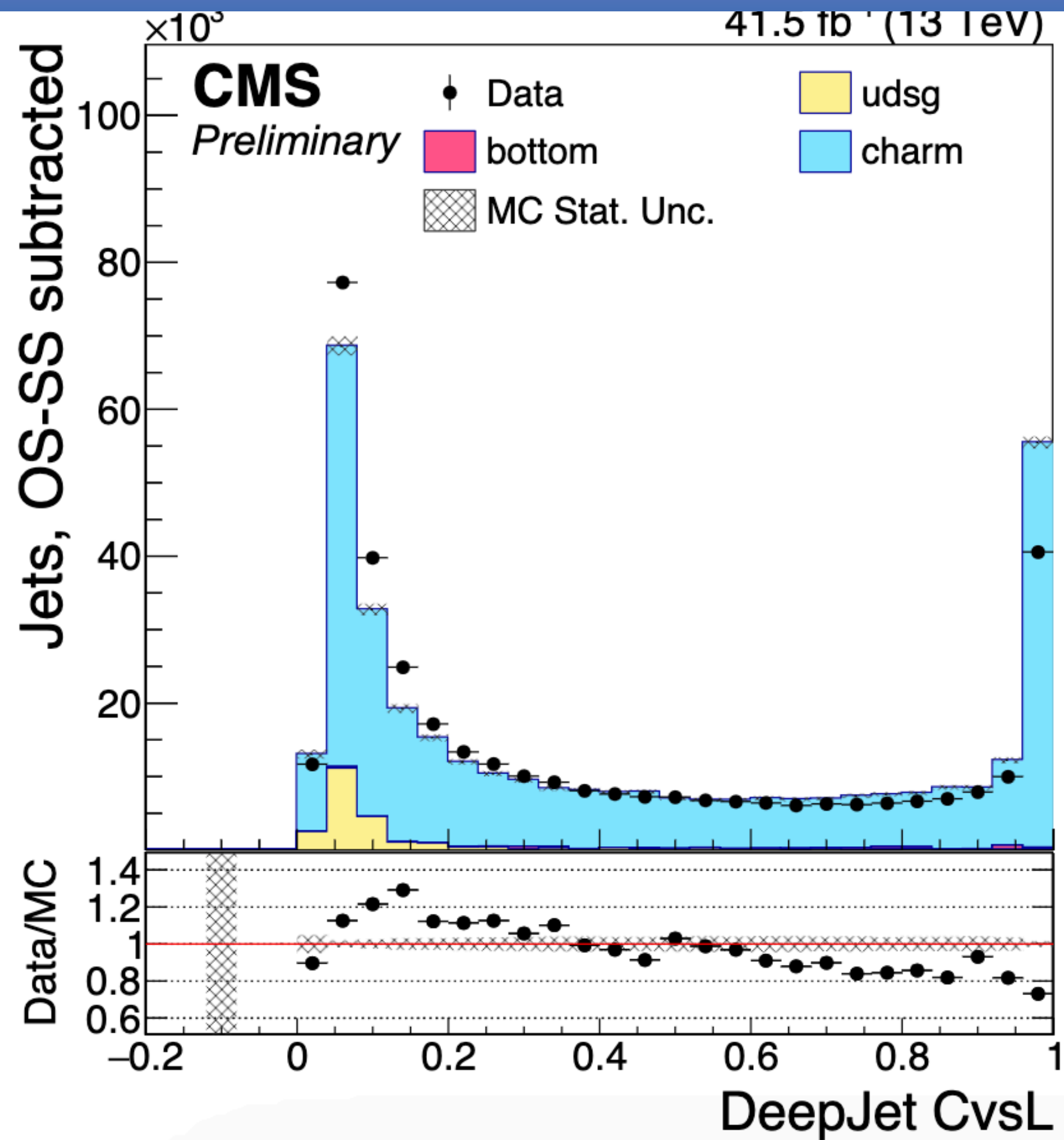    - For each systematic source, redo iterative fit with selections shifted by 1σ on each side



Purest
2nd Purest
3rd Purest

$SF = (data - bkg_{sim})/sig_{sim}$

CvsB

CvsL

Iteration 1:

Rank by purity | $SF_1$ Derive Apply | Apply $SF_1$ | $SF_2$ Derive Apply | Derive $SF_3$

apply $SF_1, SF_2, SF_3$

Iteration 2:

$SF_1$ Derive Apply | Apply $SF_1$ | $SF_2$ Derive Apply | Derive $SF_3$

Continue to iteration N until a Global $\chi^2$ minimization is achieved

$$\chi^2_{sel} = \sum_{i=1}^{nBins} \frac{(N_{sim,i} - N_{data,i})^2}{\sigma^2_{sim,i} + \sigma^2_{data,i}}$$
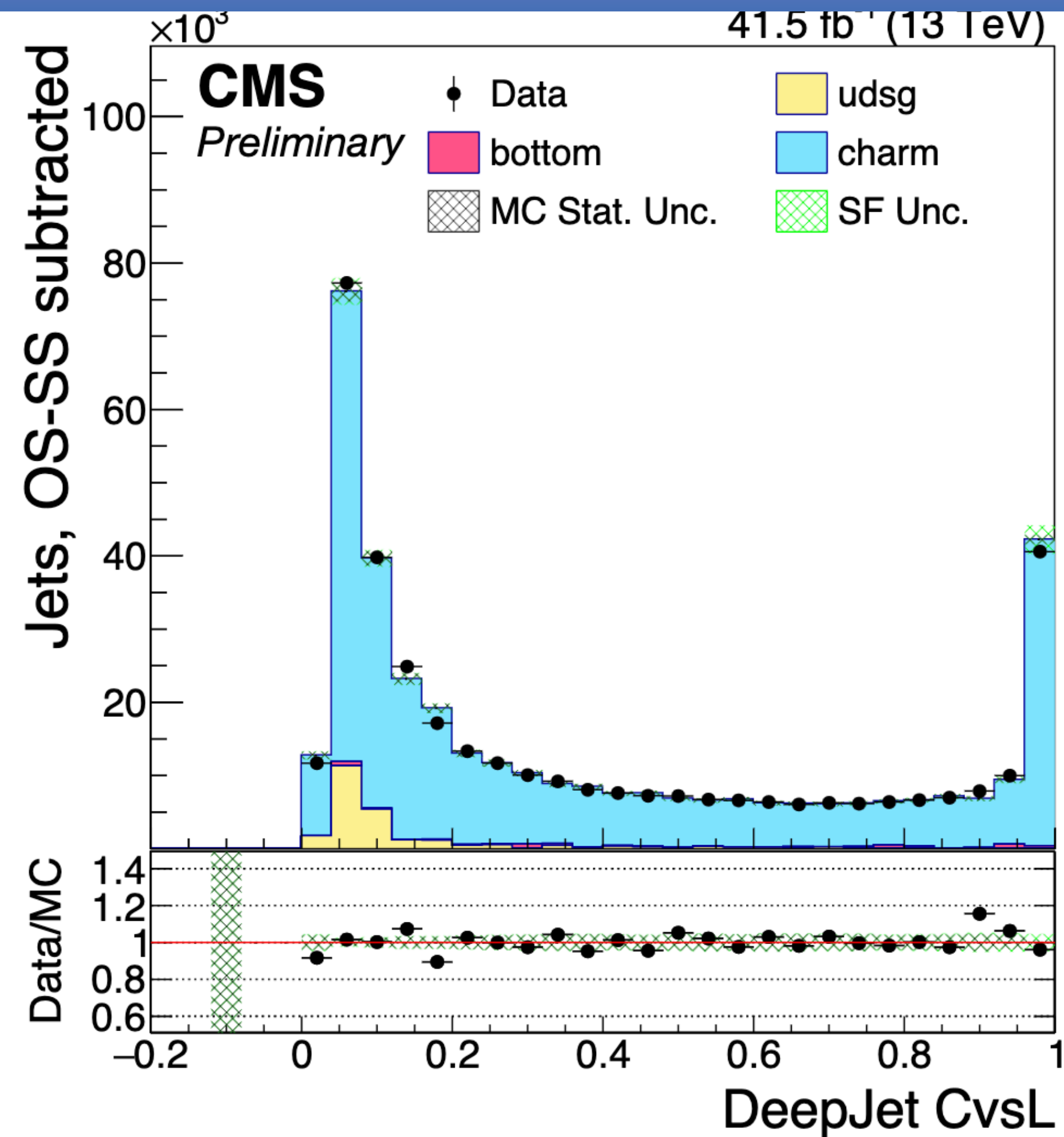
Global $\chi^2 = \chi^2_{Wc} + \chi^2_{t\bar{t}} + \chi^2_{DY}$

# Performance: c tagging



Before applying SF$_c$

Apply SF and
Uncertainties

After applying SF$_c$

*Good agreements post calibration!*
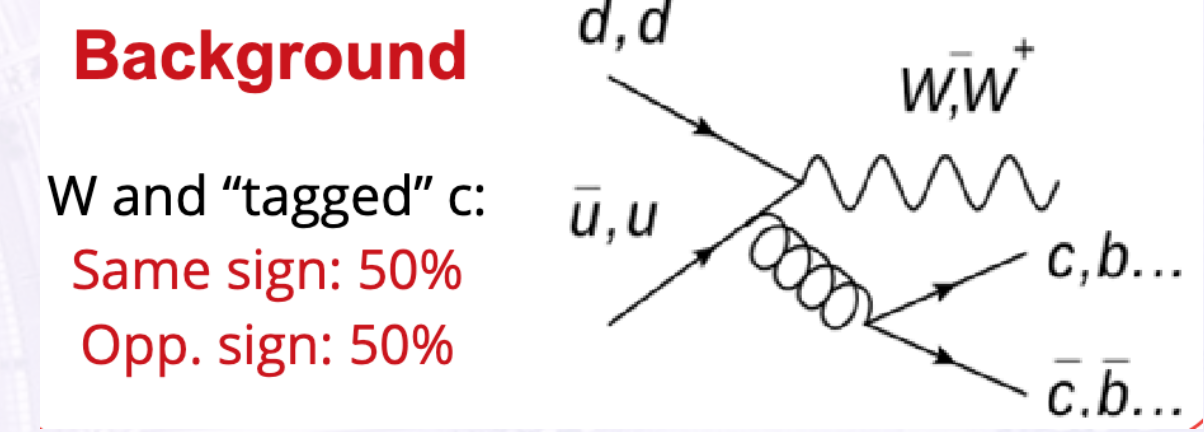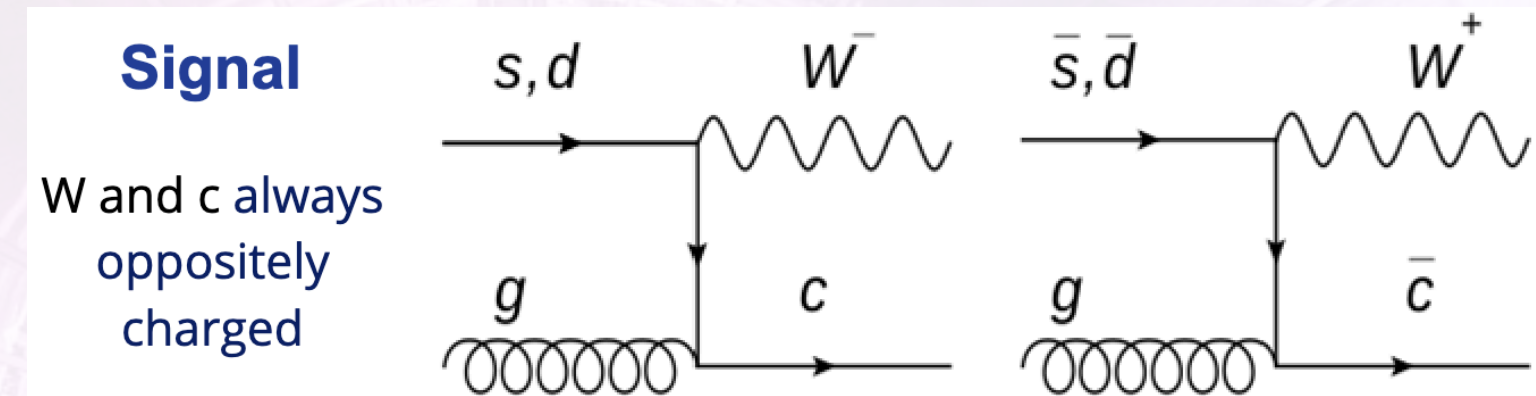
# Calibration

**Variable binning along CvsL axis**

- Binning for the SFs for each flavour is determined based on:

  - **Stats available:** Set the width to the minimum width required to reach a target stat uncertainty ($\epsilon_{max}$) per bin. Practically, we start with a minimum width, $b_{min}$, and increase in steps of $b_{min}$.

  - **An upper limit:** If stat unc < $\epsilon_{max}$ is not satisfied even with width >= $b_{max}$, set bin width to $b_{max}$.

- Parameters to be preset: $\epsilon_{max}$, $b_{min}$, $b_{max}$.

  - Final choice of parameters: $\epsilon_{max} = 2\%$, $b_{min} = 0.02$, $b_{max} = 0.10$.
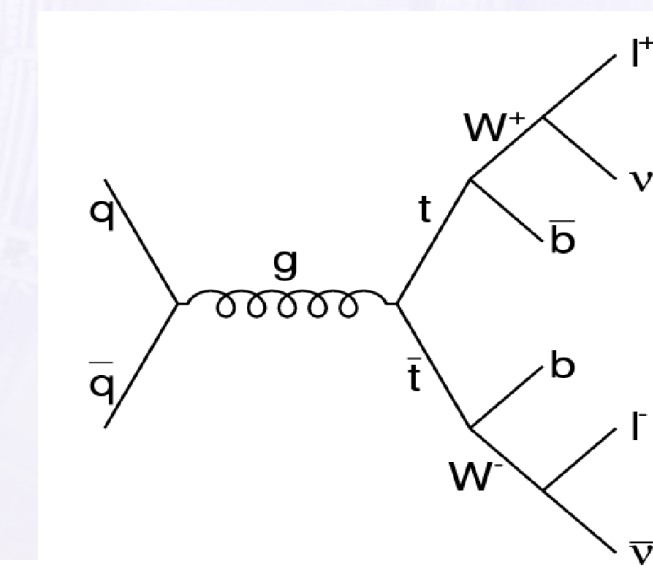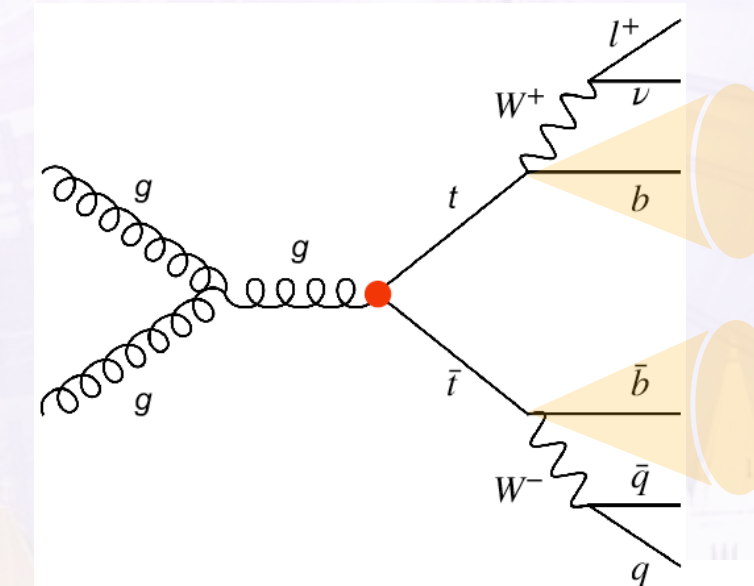
# Calibration: Event selection for c SF

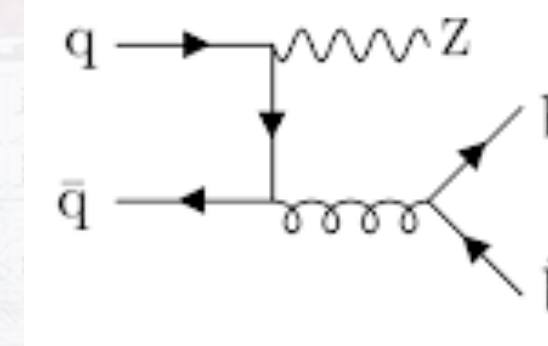- **Define three regions:**

  - **c enriched:** $W(l\nu)$ and at least one jet with a soft, non-isolated mu,

    - OS – SS selection gives pure W + c distributions

  - **b enriched:** Semileptonic and dileptonic $t\bar{t}$

  - **Light-flavor jet enriched:** DY(ll)+jets(udsg)

    *No heavy-flavor tagging in the selections*

**Signal**

W and c always oppositely charged

**Background**

W and "tagged" c:
Same sign: 50%
Opp. sign: 50%

| Selection | Jet yield | c % | b % | udsg % |
|-----------|-----------|------|------|--------|
| W+c | 360000 | **93.0** | 1.0 | 6.0 |
| $t\bar{t}$ | 380000 | 12.0 | **81.1** | 6.9 |
| DY+jet | 8509000 | 8.9 | 5.1 | **86.0** |

*Reasonably pure regions enriched in respective flavors*

# Uncertainties: c-taggging

- **Statistical uncertainties:** Evaluated as the statistical uncertainty in the ratio of Data and MC in each bin

- **Systematics:**

  We consider the
  following sources:

|  |  |
| --- | --- |
| ✓ Electron ID | ✓ Factorisation scale |
| ✓ Muon ID | ✓ Renormalisation scale |
| ✓ PileUp | ✓ Parton shower: ISR |
| ✓ JES (total) | ✓ Parton shower: FSR |
| ✓ JER |  |

**Cross section:**

- ✓ W+Jets
- ✓ DY+Jets
- ✓ $t\bar{t}$
- ✓ Single Top

[1] https://cds.cern.ch/record/2866276