# Performance of jet flavor tagging and measurement of $R_b$ using ParticleNet at CEPC

## Libo Liao

Gang Li, Weimin Song, Shudong Wang, and Zhaoling Zhang
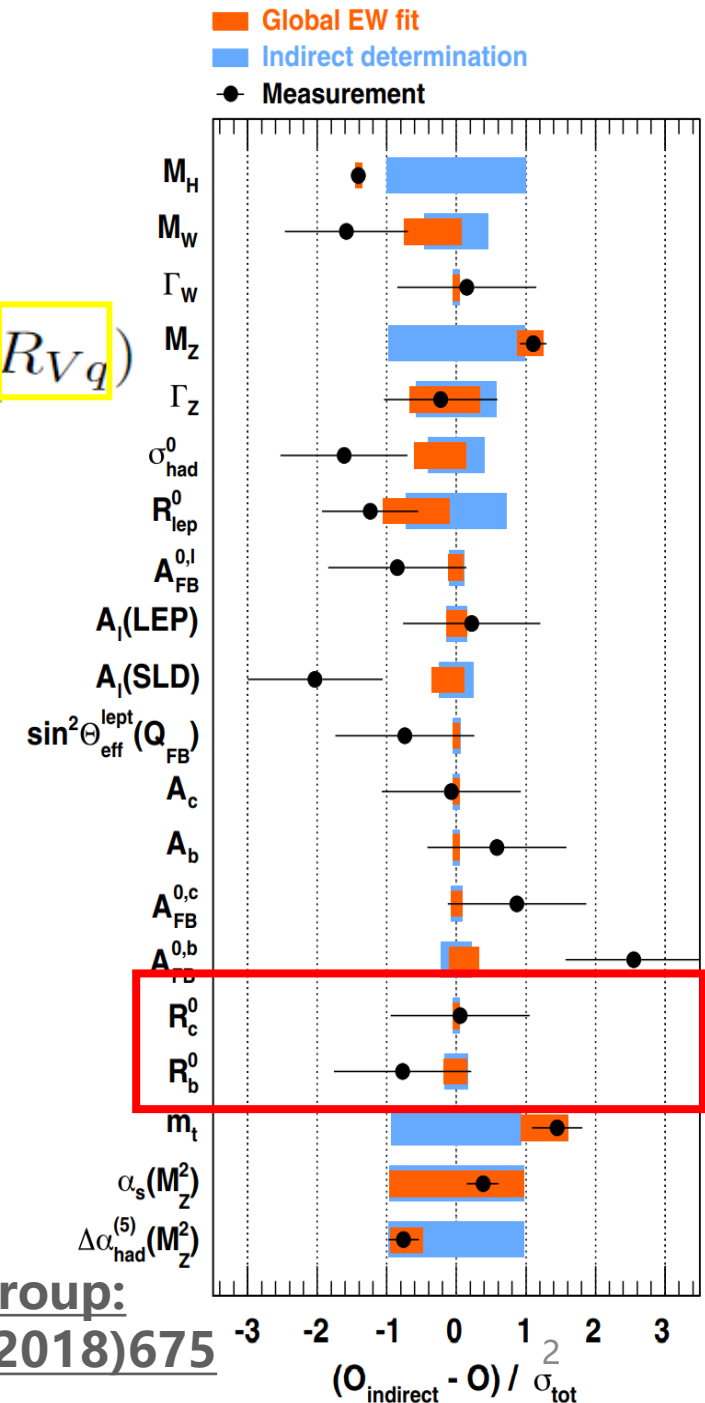
2023.8.17

arXiv: 2208.13503

CEPC味物理-新物理和相关探测技术研讨会

# Introduction

➤ Relative decay width： $R_q = \Gamma_q / \Gamma_h$

- SM testing

$$\Gamma(Z \to q\bar{q}) = 12\Gamma_0 \left( g_{Aq}^2 R_{Aq} + g_{Vq}^2 R_{Vq} \right)$$

- Searching for new physics

- Precision electroweak measurement

➤ Status of $R_b$ and $R_c$ measurements in experiment and theory

- Theoretical << Experimental

|       | Experiment              | Gfitter results          |
|-------|-------------------------|--------------------------|
| $R_b$ | $0.21629 \pm 0.00066$   | $0.21582 \pm 0.00011$    |
| $R_c$ | $0.1721 \pm 0.0030$     | $0.17224 \pm 0.00008$    |



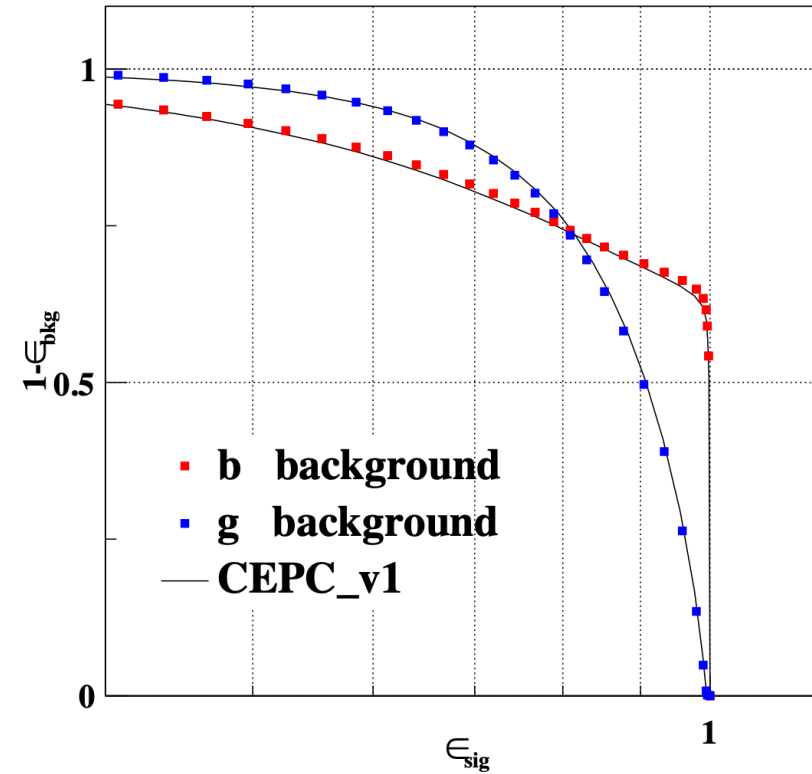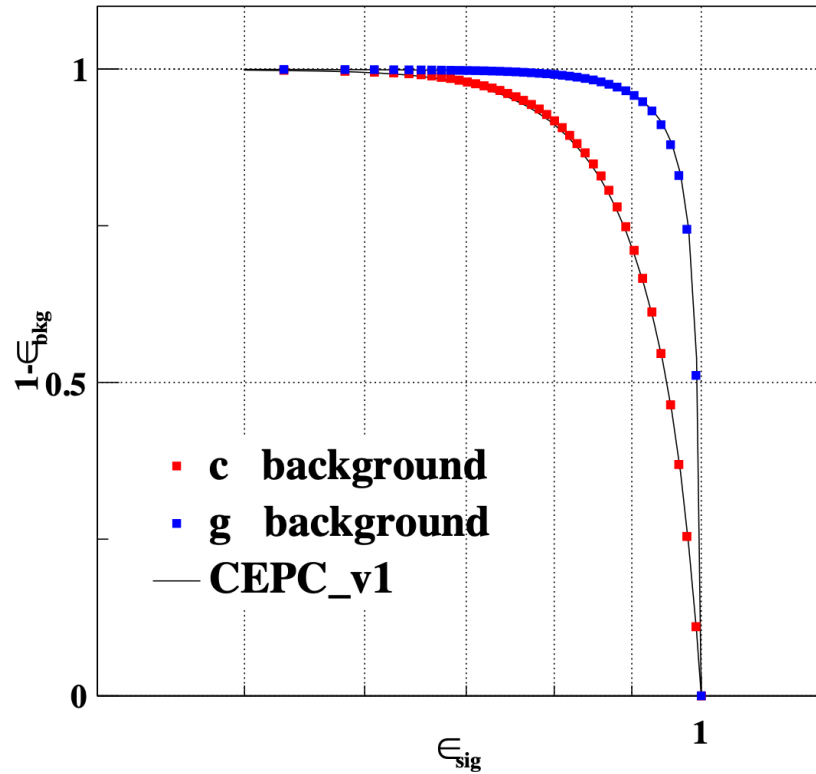**Gfitter Group:**
**EPJC78 (2018)675**

# Introduction

➢ **New colliders to perform precision electroweak studies**

  • CEPC/FCC-ee/ILC/......

➢ **Jet: Key physics object**

  • Vertexing->Clustering->Tagging
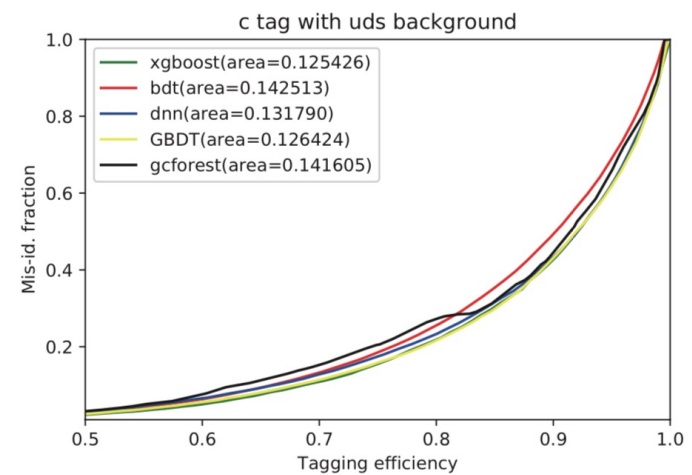
➢ **Tagging methods**
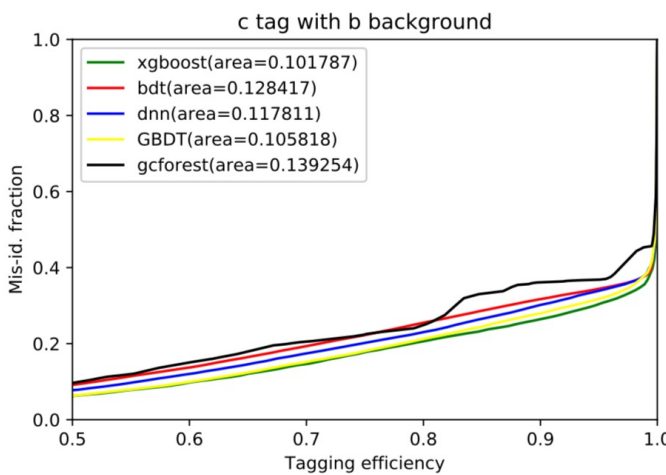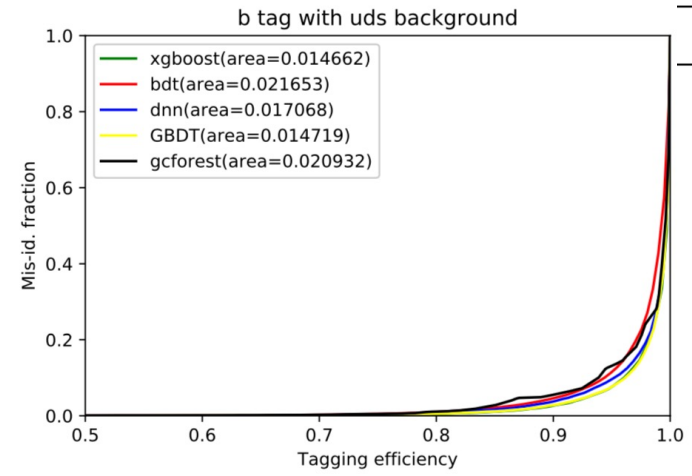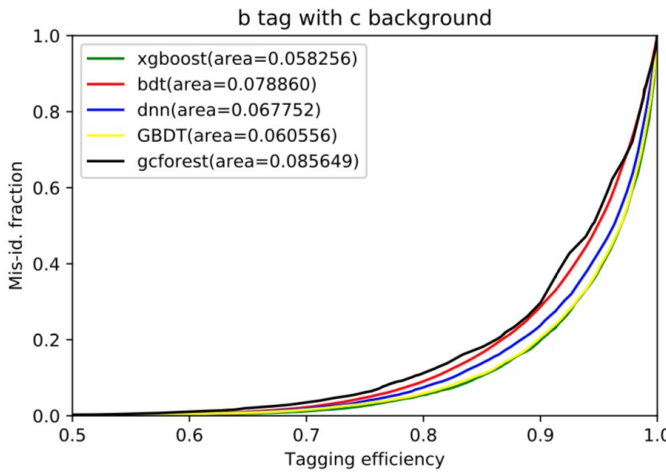
  • Cut-based->TMVA->Deep Learning

# Flavor tagging @ CDR

➤ CEPC baseline detector (TMVA/BDT)

  80% eff. & 90% purity in $b$-tagging

  60% eff. & 60% purity in $c$-tagging

# Flavor tagging @ CDR

Fan Yang: Flavor Tagging Using Machine Learning Algorithms



| Algorithm | DNN | BDT | GBDT | gcforest | xgboost |
|-----------|-----|-----|------|----------|---------|
| Accuracy | 0.788 | 0.776 | 0.794 | 0.785 | 0.801 |

| tag-background | efficiency (%) | Mis-id fraction (%) | | | | |
|----------------|----------------|---------|-----|------|-----|----------|
| | | xgboost | DNN | GBDT | BDT | gcforest |
| b-c | 80 | 5.4 | 7.5 | 5.8 | 9.3 | 10.8 |
| | 90 | 20.1 | 23.7 | 20.6 | 29.2 | 26.3 |
| | 95 | 39.0 | 43.5 | 39.6 | 50.2 | 56.3 |
| b-uds | 80 | 0.5 | 0.7 | 0.5 | 1.0 | 1.1 |
| | 90 | 2.7 | 3.7 | 2.8 | 4.7 | 4.9 |
| | 95 | 7.8 | 9.7 | 7.8 | 11.3 | 13.6 |
| c-b | 80 | 20.8 | 23.1 | 21.5 | 25.6 | 25.1 |
| | 90 | 26.5 | 30.2 | 28.1 | 32.1 | 36.1 |
| | 95 | 30.6 | 33.9 | 31.8 | 34.4 | 36.8 |
| c-uds | 80 | 22.3 | 23.3 | 22.3 | 26.0 | 27.4 |
| | 90 | 43.4 | 43.5 | 43.8 | 51.9 | 43.5 |
| | 95 | 63.6 | 61.7 | 62.1 | 68.8 | 66.1 |

# Deep learning architectures
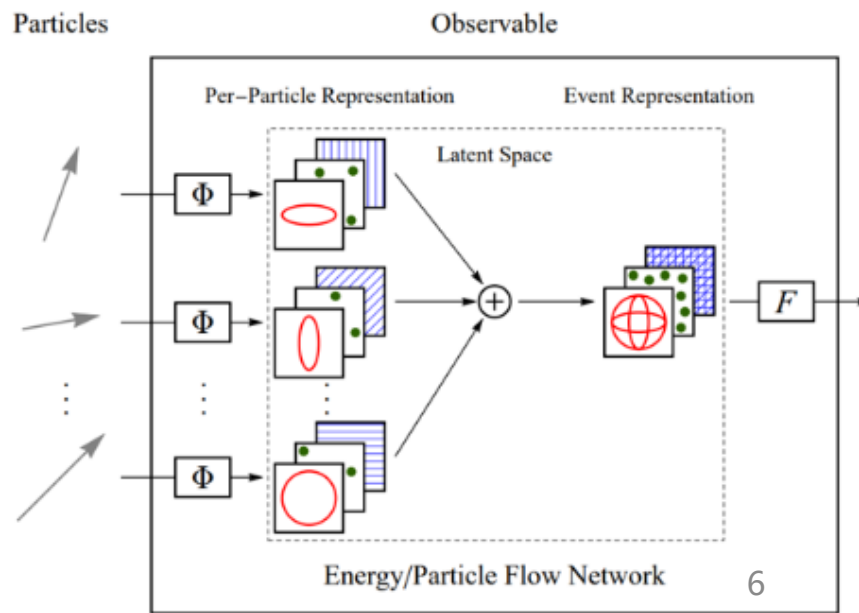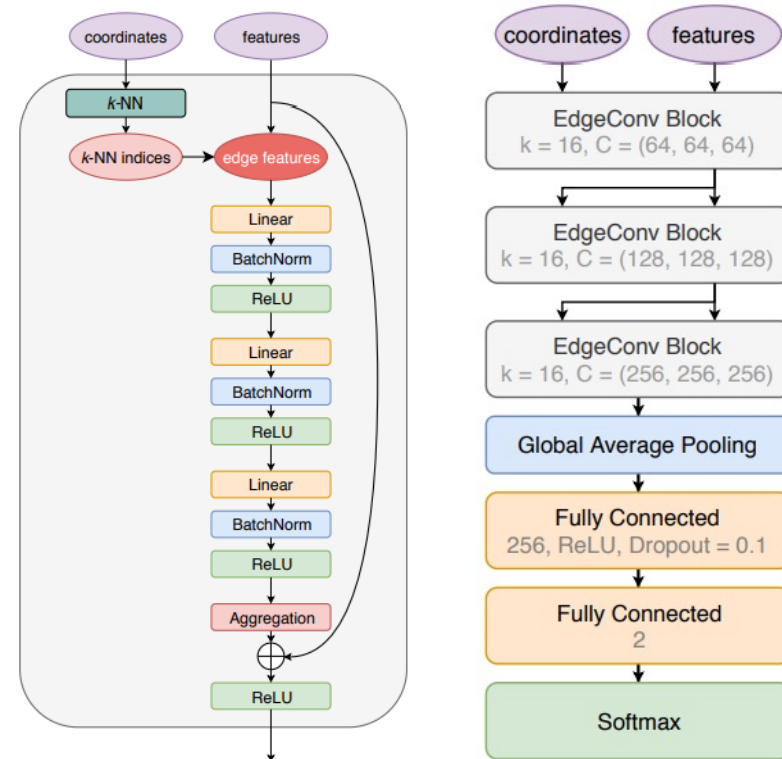
[Jet tagging via particle clouds]

➢ **ParticleNet**

- Treating a jet as an unordered set of particles in space
- Using permutation-invariant graph neural networks

[Energy flow networks: deep sets for particle jets]

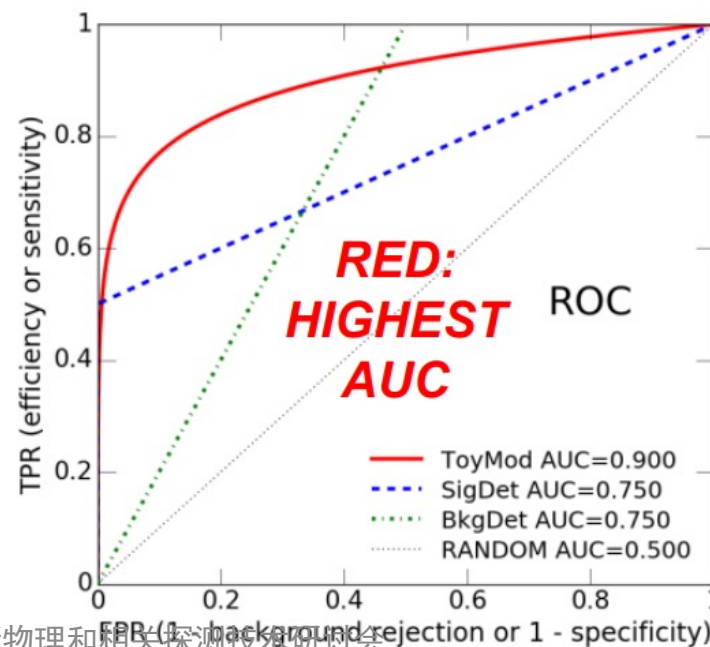➢ **Particle Flow Network (PFN)**

- Based on "point clouds"
- As a cross check





Energy/Particle Flow Network

# Evaluation metrics

➢ Efficiency $\epsilon_s$ = TP/(TP+FN)

➢ Purity $\rho_s$ = TP/(TP+FP)

➢ Accuracy = (TP+TN)/ALL

➢ ROC/AUC

➢ $\epsilon_s \times \rho_s$: between 0 and 1

  • The higher, the better

  • Proportional to 1/error2

$$(\Delta R_i)^2 \propto \frac{1}{\epsilon_i \rho_i}$$



| TP ($S_{sel}$) | FP ($B_{sel}$) | | TP ($S_{sel}$) | FP ($B_{sel}$) | | TP ($S_{sel}$) | FP ($B_{sel}$) |
|---|---|---|---|---|---|---|---|
| FN ($S_{rej}$) | TN ($B_{rej}$) | | FN ($S_{rej}$) | TN ($B_{rej}$) | | FN ($S_{rej}$) | TN ($B_{rej}$) |

| $\mathrm{TPR} = \dfrac{\mathrm{TP}}{\mathrm{TP}+\mathrm{FN}}$ | $\mathrm{PPV} = \dfrac{\mathrm{TP}}{\mathrm{TP}+\mathrm{FP}}$ | $\mathrm{TNR} = \dfrac{\mathrm{TN}}{\mathrm{TN}+\mathrm{FP}} = 1-\mathrm{FPR}$ |
|---|---|---|
| HEP: "efficiency" | HEP: "purity" | HEP: "background rejection" |
| $\epsilon_s = \dfrac{S_{sel}}{S_{tot}}$ | $\rho = \dfrac{S_{sel}}{S_{sel}+B_{sel}}$ | $1-\epsilon_b = 1-\dfrac{B_{sel}}{B_{tot}}$ |



**RED: HIGHEST AUC**  ROC

ToyMod AUC=0.900
SigDet AUC=0.750
BkgDet AUC=0.750
RANDOM AUC=0.500

**Andrea Valassi : ROC's, AUC's and alternatives in HEP and other domains**

# Datasets

➢ Full simulation with CEPC baseline detector at $Z$-pole

➢ PID used as a feature by matching reconstruction and MC truth

➢ 900k jets for each flavor($b, c, o = uds$)

➢ Clustered by $ee - kt$ into 2 jets

| Variable | Definition |
|---|---|
| $\cos\theta$ | cosine of polar angle of particle |
| $\phi\sin\theta$ | azimuth angle times sine of polar theta of particle |
| $\Delta R$ | $\sqrt{\delta\theta^2 + \delta\phi^2}$, angular separation between the particle and the jet axis |
| PID | particle ID |
| $E$ | energy of a particle |
| $Q$ | electric charge of a particle |
| $\log E$ | logarithm of the particle's energy |
| $\log P$ | logarithm of the particle's momentum |
| $D_0$ | impact parameter of a track in the r-$\phi$ plane |
| $Z_0$ | impact parameter of a track along the $z$ axis |
| $D_0/\sigma_{D_0}$ | significance of the impact parameter in the r-$\phi$ plane |
| $Z_0/\sigma_{Z_0}$ | significance of the impact parameter along the $z$ axis |
| prob | the probability for a certain Chi-squared and number of degrees of freedom |

# Jet features
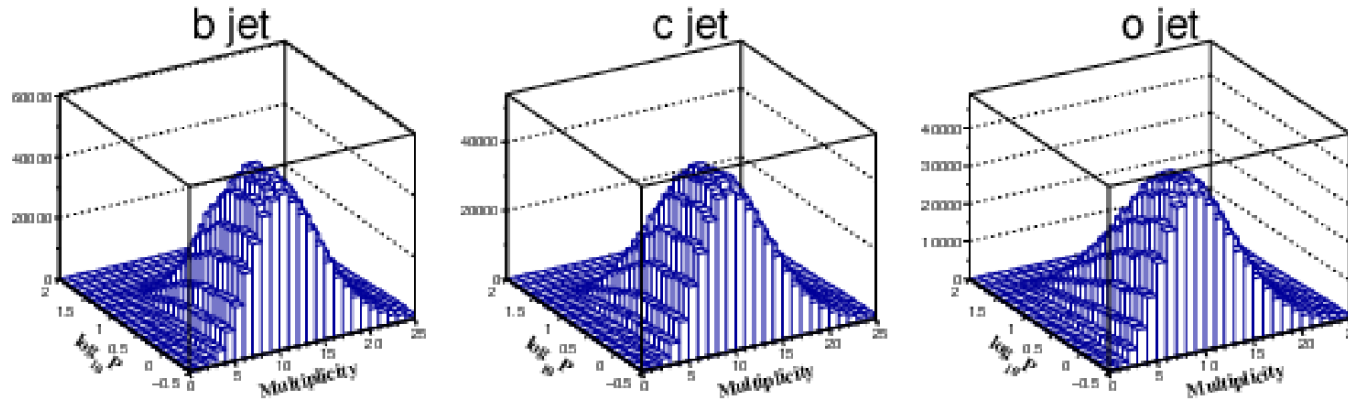
- (a) Multiplicity
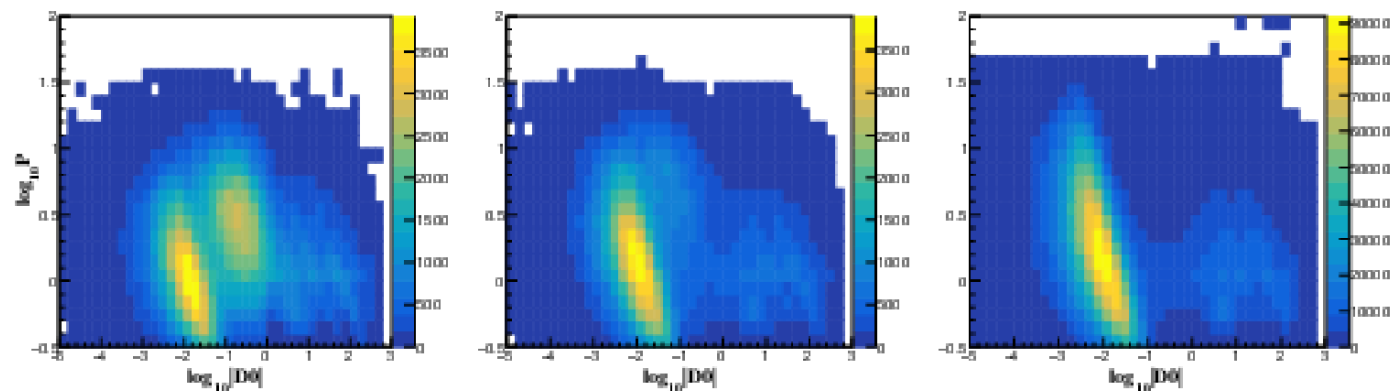  - $b > c > o$

- (b) Impact parameters
  - Larger impact parameters, more energetic tracks in $b$

- (c) The weighted fractions of particles
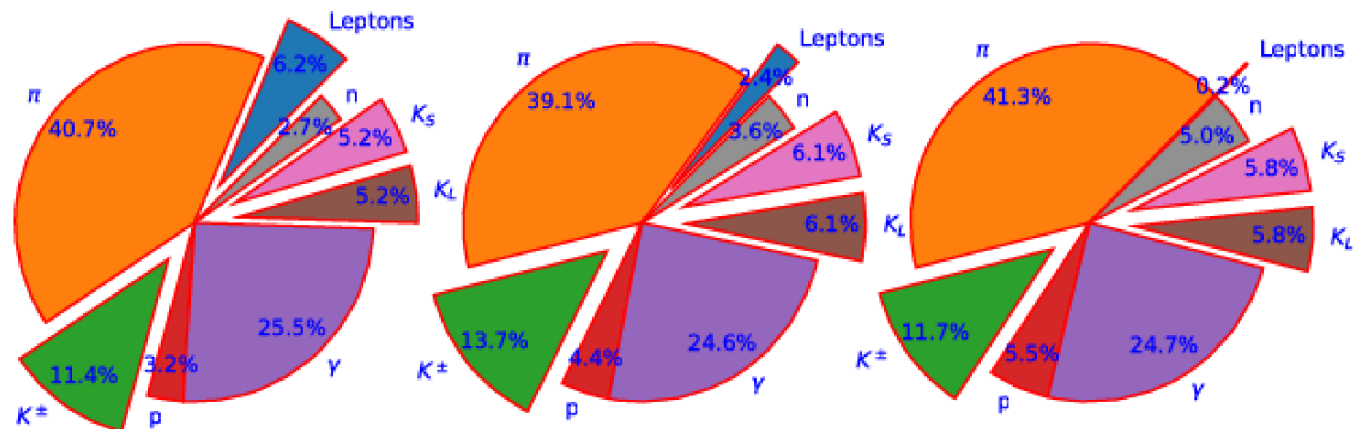  - Far more energetic leptons in $b$
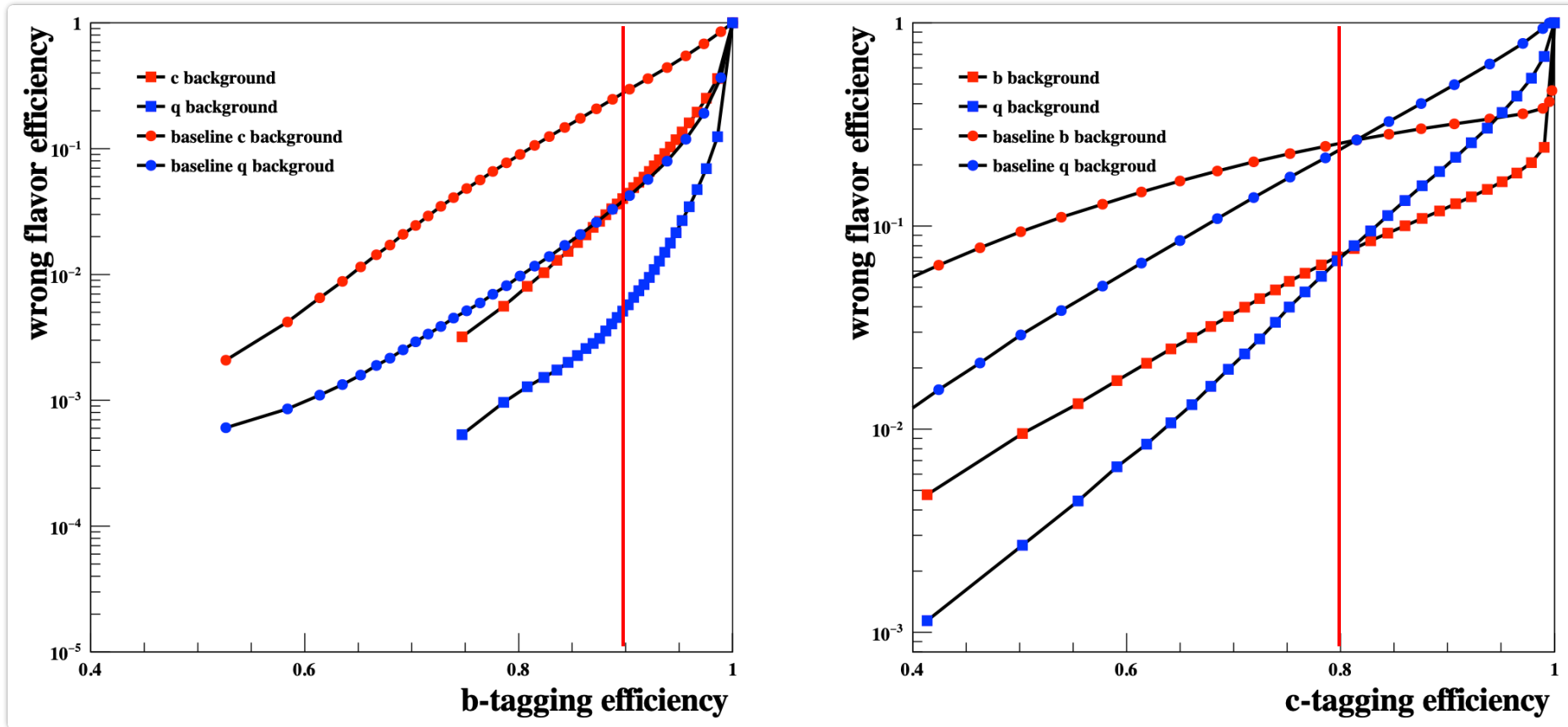  - Slightly more energetic $K$ in $c$



(a) Log(P) vs. Multiplicity

(b) Log(P) vs. log(D0)

CEPC味物理-

# Jet tagging

tagging efficiencies vs. the corresponding wrong flavour efficiencies

# Jet tagging

Fan Yang: Flavor Tagging Using Machine Learning

| Algorithm | ParticleNet | PFN | DNN | BDT | GBDT | gcforest | XGBoost |
|-----------|-------------|-------|-------|-------|-------|----------|---------|
| Accuracy | 0.876 | 0.850 | 0.788 | 0.776 | 0.794 | 0.785 | 0.801 |

➢ At least 9% improvement in ParticleNet at global accuracy

- Richer information
- Strong inductive bias

➢ The performance of $b$-tagging and $o$-tagging are much better than $c$-tagging

➢ ParticleNet is better than the PFN

- Consistent with the study
  [Jet tagging via particle clouds]





| tag | ParticleNet | | PFN | |
|-----|-------------|-----|-------------|-----|
| | Efficiency | AUC | Efficiency | AUC |
| $b$ | 0.908 | 0.986 | 0.870 | 0.979 |
| $c$ | 0.798 | 0.951 | 0.765 | 0.930 |
| $o$ | 0.923 | 0.974 | 0.911 | 0.966 |

# Physics impacts of jet tagging

➢ Take LCFIPlus & XGBoost(CDR baseline) as reference

- ● ParticleNet & PFN are better than the baseline, especially in $c$-tagging

➢ Statistical uncertainty can be improved

- ● roughly 30%(sqrt(0.597/0.345)) in counting $c$ jets

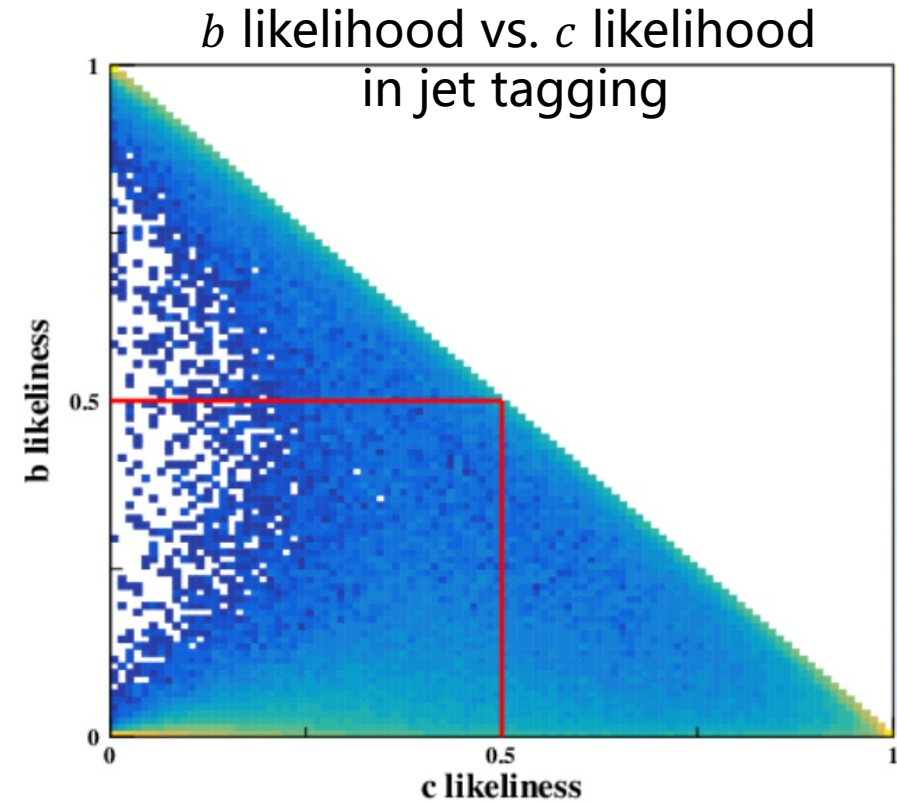| tag | $\epsilon_S(\%)$ | $\epsilon \times \rho$ | | | |
|-----|------|----------|---------|-------------|------|
| | | LCFIPlus | XGBoost | ParticleNet | PFN |
| b | 80 | – | 0.747 | 0.786 | 0.763 |
| | 90 | 0.72 | 0.713 | 0.821 | 0.752 |
| c | 60 | 0.36 | – | 0.554 | 0.485 |
| | 70 | – | – | 0.605 | 0.497 |
| | 80 | – | 0.345 | 0.597 | 0.467 |
| | 90 | – | 0.292 | 0.532 | 0.402 |

Applied in $R_q$ measurement

# $R_b$ & $R_c$ measurement

$$N_s^{i,\mathrm{obs}} = 2N^{h,\mathrm{pro}} \cdot (R_b\varepsilon_{ib} + R_c\varepsilon_{ic} + R_o\varepsilon_{io}) ,$$

$$N_d^{i,\mathrm{obs}} = N^{h,\mathrm{pro}} \cdot [R_b\varepsilon_{ib}^2(1 + C_{ib}) + R_c\varepsilon_{ic}^2(1 + C_{ic})$$
$$+ R_o\varepsilon_{io}^2(1 + C_{io})] ,$$

➤ Double tagging:

- Neglect the correlation of jets

- Choose the working point

- Solved 6 equations by the least square method

➤ References

- LEP+SLC: Limited by statistics & flavor tagging

- Template fit: Much larger statistics & better flavor tagging in CEPC baseline

➤ Our work:  *Int.J.Mod.Phys.A* 36 (2021) 27, 2150207

  ➤ Statistic of $10^{11}$ $Z$ bosons, same as template fit

  ➤ Comparable with template fit in $R_b$

  ➤ Improved more than 60% in $R_c$ measurement



*b* likelihood vs. *c* likelihood in jet tagging

| | $\sigma_{R_b}$ | $\sigma_{R_c}$ | $\sigma_{R_q}$ |
|---|---|---|---|
| LEP+SLC | 659 | 3015 | - |
| Template fit | 1.2 | 2.3 | 2.1 |
| Double tag | 1.3 | 1.4 | - |

All results in 10^-6

# Conclusion

➢ Two novel deep learning methods are used to enhance the performance of jet flavor tagging

- Significant improvement in jet tagging, especially for $c$ tagging
- Maximize the usage of information in a jet
- Strong inductive bias

➢ $R_q$ measurement is taken to demonstrate the physics impacts

- Statistical uncertainty improved 60+% in $R_c$ measurement
- Systematic uncertainties pose a significant challenge and require careful investigation

# Thank you!

# Comparison of the performance of ParticleNet with three alternative models

TABLE II: Performance comparison on the top tagging benchmark dataset. The ParticleNet, ParticleNet-Lite, P-CNN and ResNeXt-50 models are trained on the top tagging dataset starting from randomly initialized weights. For each model, the training is repeated for 9 times using different randomly initialized weights. The table shows the result from the median-accuracy training, and the standard deviation of the 9 trainings is quoted as the uncertainty to assess the stability to random weight initialization. Uncertainty on the accuracy and AUC are negligible and therefore omitted. The performance of PFN on this dataset is reported in Ref. [52], and the uncertainty corresponds to the spread in 10 trainings.

| | Accuracy | AUC | $1/\varepsilon_b$ at $\varepsilon_s = 50\%$ | $1/\varepsilon_b$ at $\varepsilon_s = 30\%$ |
|---|---|---|---|---|
| ResNeXt-50 | 0.936 | 0.9837 | $302 \pm 5$ | $1147 \pm 58$ |
| P-CNN | 0.930 | 0.9803 | $201 \pm 4$ | $759 \pm 24$ |
| PFN | - | 0.9819 | $247 \pm 3$ | $888 \pm 17$ |
| ParticleNet-Lite | 0.937 | 0.9844 | $325 \pm 5$ | $1262 \pm 49$ |
| **ParticleNet** | **0.940** | **0.9858** | $\mathbf{397 \pm 7}$ | $\mathbf{1615 \pm 93}$ |