



# Systematic Effects in Jet Tagging with the ATLAS Detector

**US LHC User's Association Annual Meeting  
Fermilab, December 14th 2023**

**Kevin Greif, on behalf of the ATLAS collaboration**

# The Jet Tagging Landscape (As of 2019)

[arxiv:1902.09914](https://arxiv.org/abs/1902.09914)

SciPost Physics

Submission

## The Machine Learning Landscape of Top Taggers

G. Kasieczka (ed)<sup>1</sup>, T. Plehn (ed)<sup>2</sup>, A. Butter<sup>2</sup>, K. Cranmer<sup>3</sup>, D. Debnath<sup>4</sup>, B. M. Dillon<sup>5</sup>,  
M. Fairbairn<sup>6</sup>, D. A. Faroughy<sup>5</sup>, W. Fedorko<sup>7</sup>, C. Gay<sup>7</sup>, L. Gouskos<sup>8</sup>, J. F. Kamenik<sup>5,9</sup>,  
P. T. Komiske<sup>10</sup>, S. Leiss<sup>1</sup>, A. Lister<sup>7</sup>, S. Macaluso<sup>3,4</sup>, E. M. Metodiev<sup>10</sup>, L. Moore<sup>11</sup>,  
B. Nachman,<sup>12,13</sup> K. Nordström<sup>14,15</sup>, J. Pearkes<sup>7</sup>, H. Qu<sup>8</sup>, Y. Rath<sup>16</sup>, M. Rieger<sup>16</sup>, D. Shih<sup>4</sup>,  
J. M. Thompson<sup>2</sup>, and S. Varma<sup>6</sup>

“Indeed, we will see that we can consider jet classification based on deep learning at the pure performance level an essentially solved problem.

For a systematic experimental application of these tools **our focus will be on a new set of questions related to training data, benchmarking, calibration, systematics, etc.”**

# The Jet Tagging Landscape (As of 2019)

[arxiv:1902.09914](https://arxiv.org/abs/1902.09914)

SciPost Physics

Submission

## The Machine Learning Landscape of Top Taggers

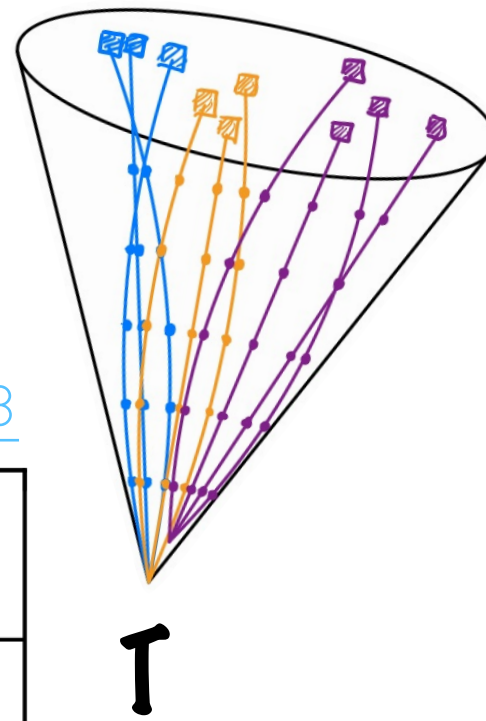
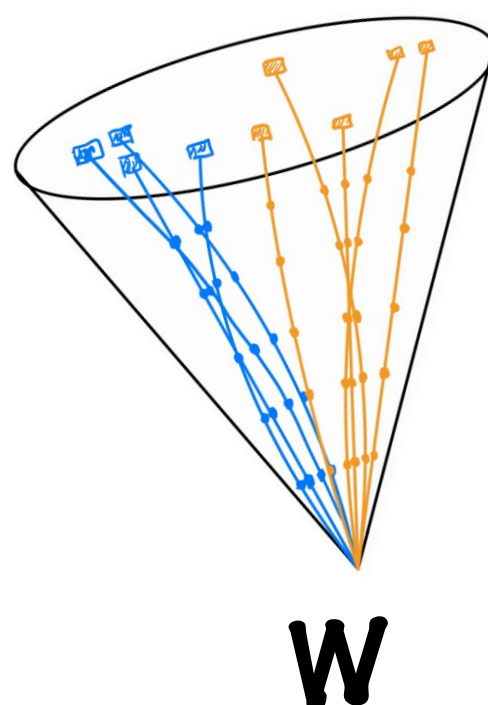
G. Kasieczka (ed)<sup>1</sup>, T. Plehn (ed)<sup>2</sup>, A. Butter<sup>2</sup>, K. Cranmer<sup>3</sup>, D. Debnath<sup>4</sup>, B. M. Dillon<sup>5</sup>,  
M. Fairbairn<sup>6</sup>, D. A. Faroughy<sup>5</sup>, W. Fedorko<sup>7</sup>, C. Gay<sup>7</sup>, L. Gouskos<sup>8</sup>, J. F. Kamenik<sup>5,9</sup>,  
P. T. Komiske<sup>10</sup>, S. Leiss<sup>1</sup>, A. Lister<sup>7</sup>, S. Macaluso<sup>3,4</sup>, E. M. Metodiev<sup>10</sup>, L. Moore<sup>11</sup>,  
B. Nachman,<sup>12,13</sup> K. Nordström<sup>14,15</sup>, J. Pearkes<sup>7</sup>, H. Qu<sup>8</sup>, Y. Rath<sup>16</sup>, M. Rieger<sup>16</sup>, D. Shih<sup>4</sup>,  
J. M. Thompson<sup>2</sup>, and S. Varma<sup>6</sup>

This Talk!

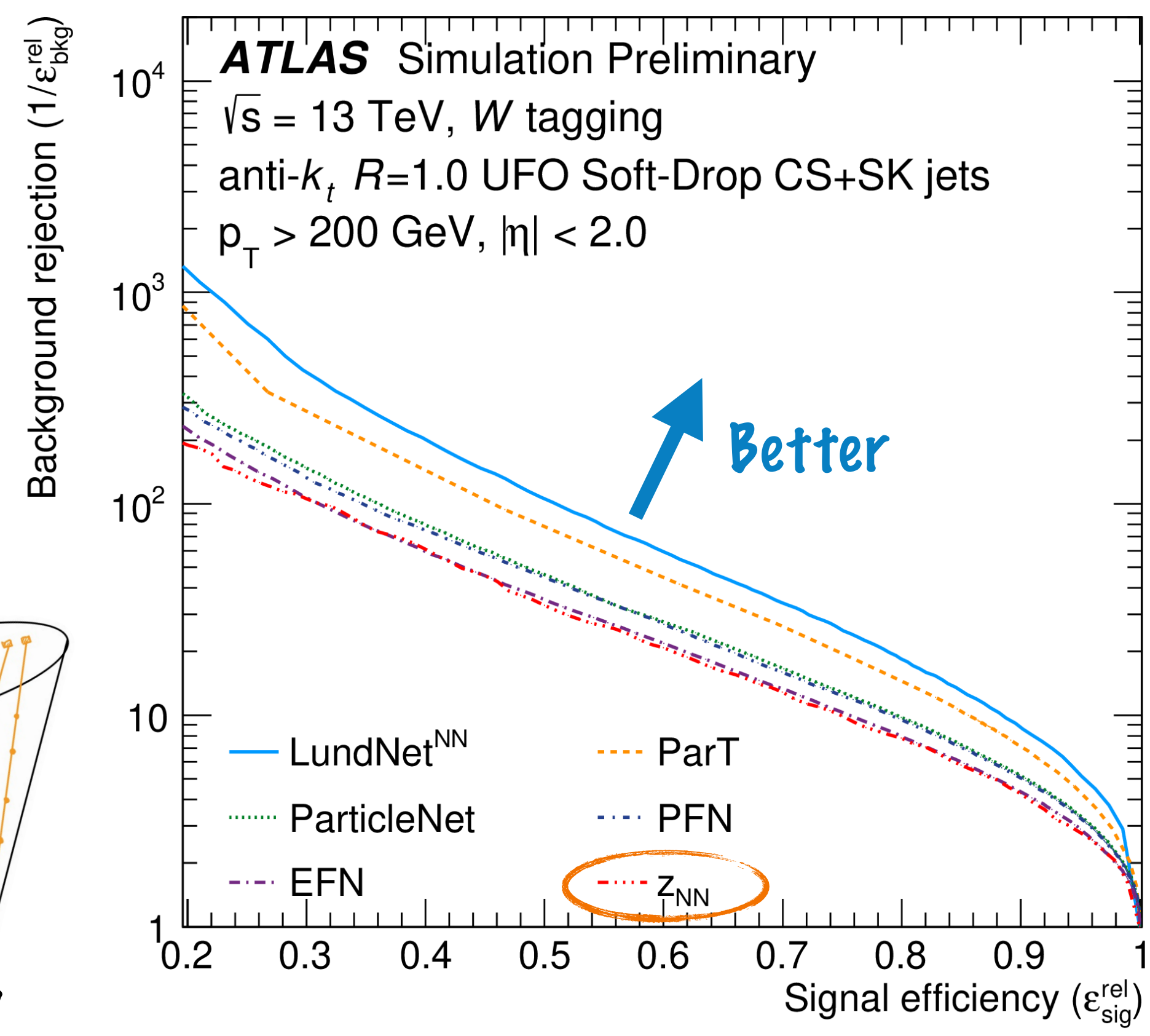
“Indeed, we will see that we can consider jet classification based on deep learning at the pure performance level an essentially solved problem.\*”

For a systematic experimental application of these tools our focus will be on a new set of questions related to training data, benchmarking, calibration, systematics, etc.”

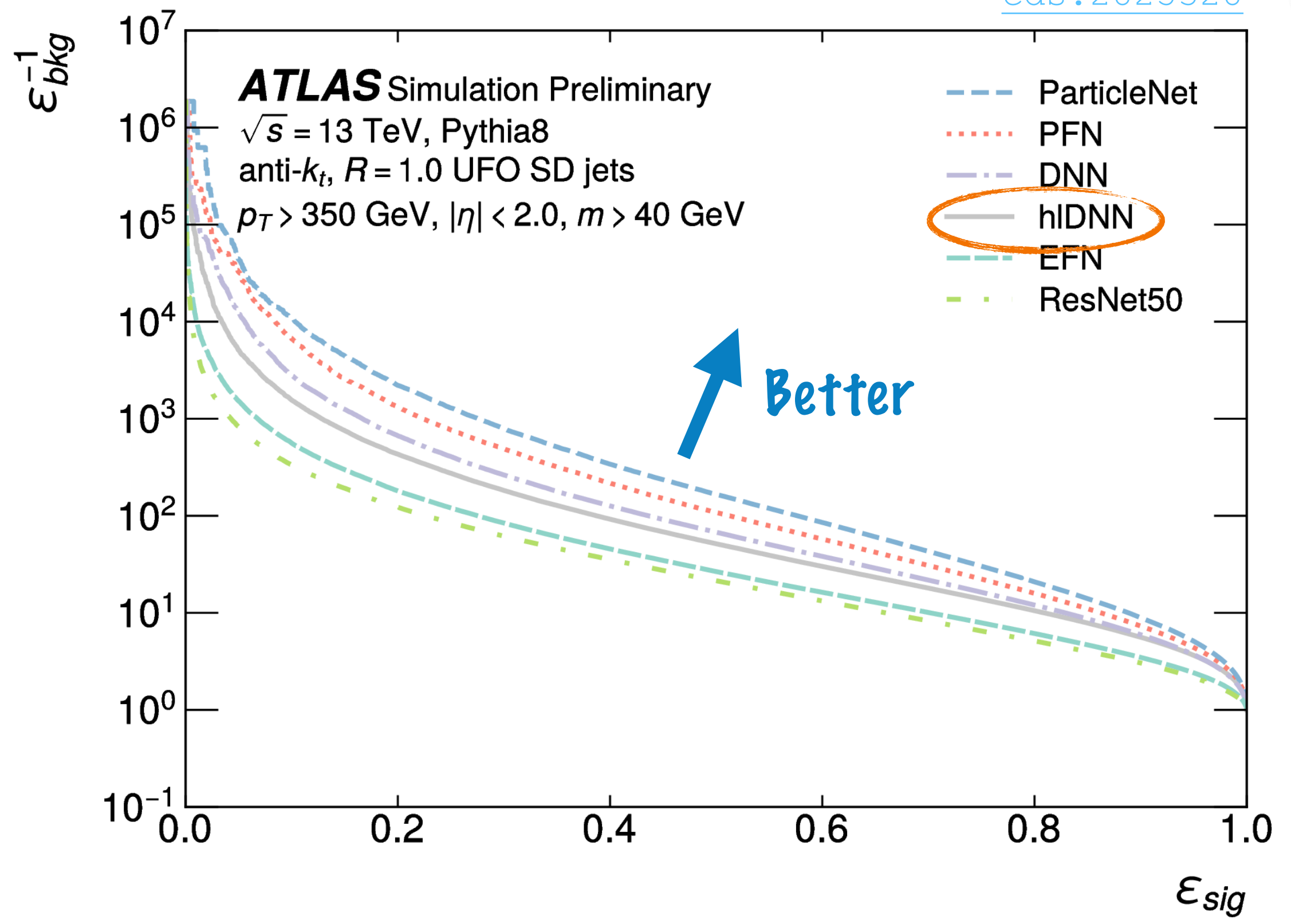
# Today: Point Cloud Taggers in ATLAS



[cds:2864131](#) [cds:2866592](#)



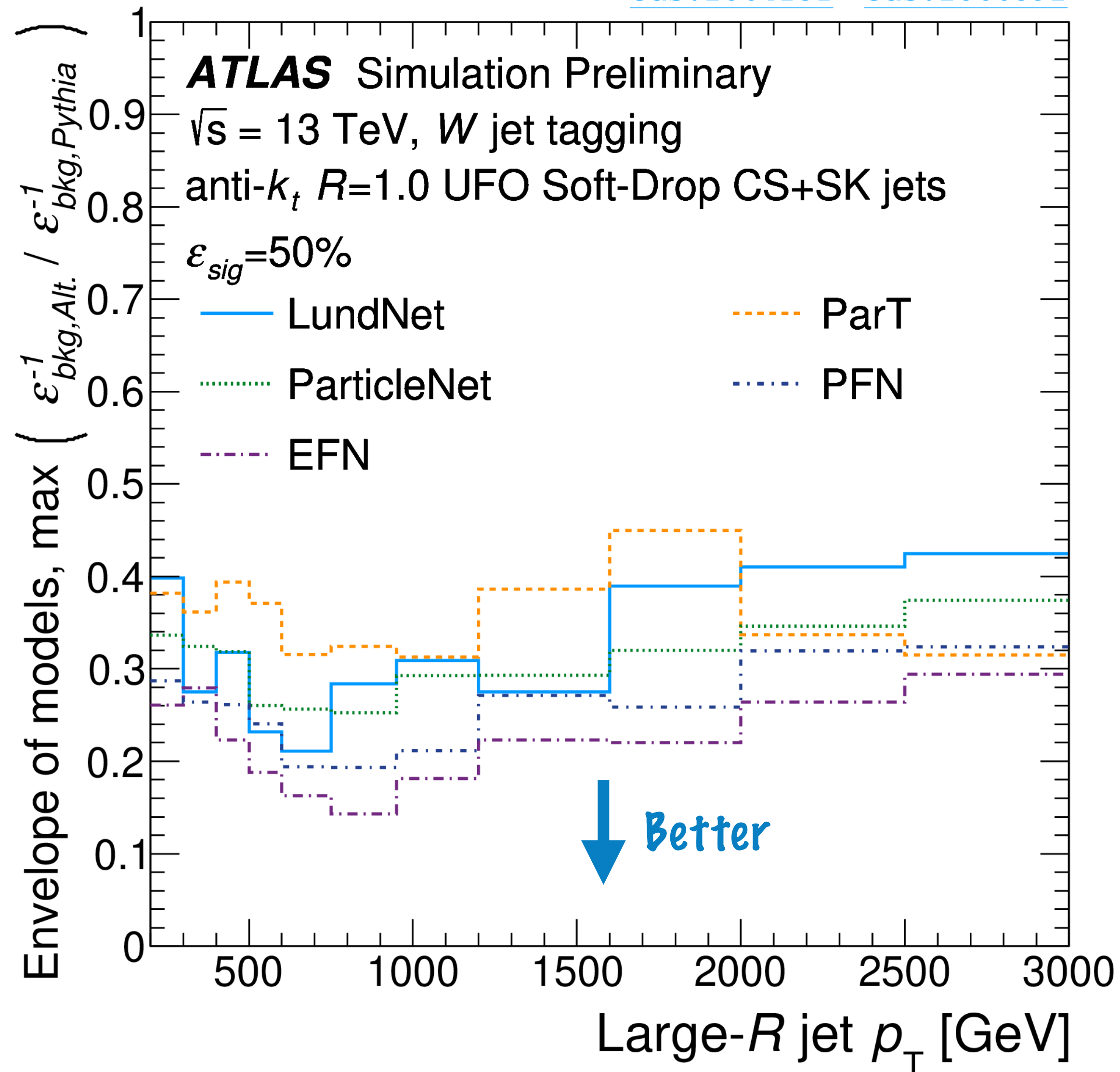
[cds:2825328](#)



Large performance gains for point-cloud taggers over **high-level quantity baselines**

# W Tagger Modeling Dependence

[cds:2864131](#) [cds:2866592](#)



- The most powerful taggers (LundNet, ParT) show variations in performance of up to 40%
- Likely to produce larger scale factor uncertainties

# Beyond Modeling Dependence and Scale Factors

[cds:2724149](#)

Modeling uncertainties were dominant for simple high-level quantity based taggers.

**What about constituent based taggers?**

Systematic Group	$W$ tagger $p_T$ bins [GeV]			
	[200,250]	[250,300]	[300,350]	[350,600]
Statistical	0.01	0.02	0.03	0.04
Theory	< 0.01	< 0.01	< 0.01	< 0.01
$t\bar{t}$ modeling	0.21	0.20	0.15	0.12
Large- $R$ jet	0.01	0.01	< 0.01	< 0.01
Other experimental	< 0.01	< 0.01	< 0.01	< 0.01
$b$ -tagging	< 0.01	< 0.01	< 0.01	< 0.01
<b>Total Uncertainty</b>	<b>0.21</b>	<b>0.20</b>	<b>0.15</b>	<b>0.12</b>

Available on the CERN CDS information server CMS PAS BTV-22-001

CMS Physics Analysis Summary

Contact: cms-pog-conveners-btag@cern.ch 2023/07/29

Performance of heavy-flavour jet identification in boosted topologies in proton-proton collisions at  $\sqrt{s} = 13$  TeV

The CMS Collaboration

Abstract

Physics measurements in the highly Lorentz-boosted regime, including the search for the Higgs boson or beyond standard model particles, are a critical part of the LHC physics program. In the CMS Collaboration, various boosted-jet tagging algorithms, designed to identify hadronic jets originating from a massive particle decaying to  $b\bar{b}$  or  $c\bar{c}$ , have been developed and deployed in a variety of analyses. This note highlights their performance on simulated events, and summarises the novel calibration methods of these algorithms with 2016-2018 data collected in proton-proton collisions at  $\sqrt{s} = 13$  TeV. Three distinct control regions are studied, selected via machine learning techniques or the presence of reconstructed muons from  $g \rightarrow b\bar{b}$  ( $c\bar{c}$ ) decays, as well as regions selected from Z boson decays. The calibration results, derived through a combination of measurements in these three regions, are presented.

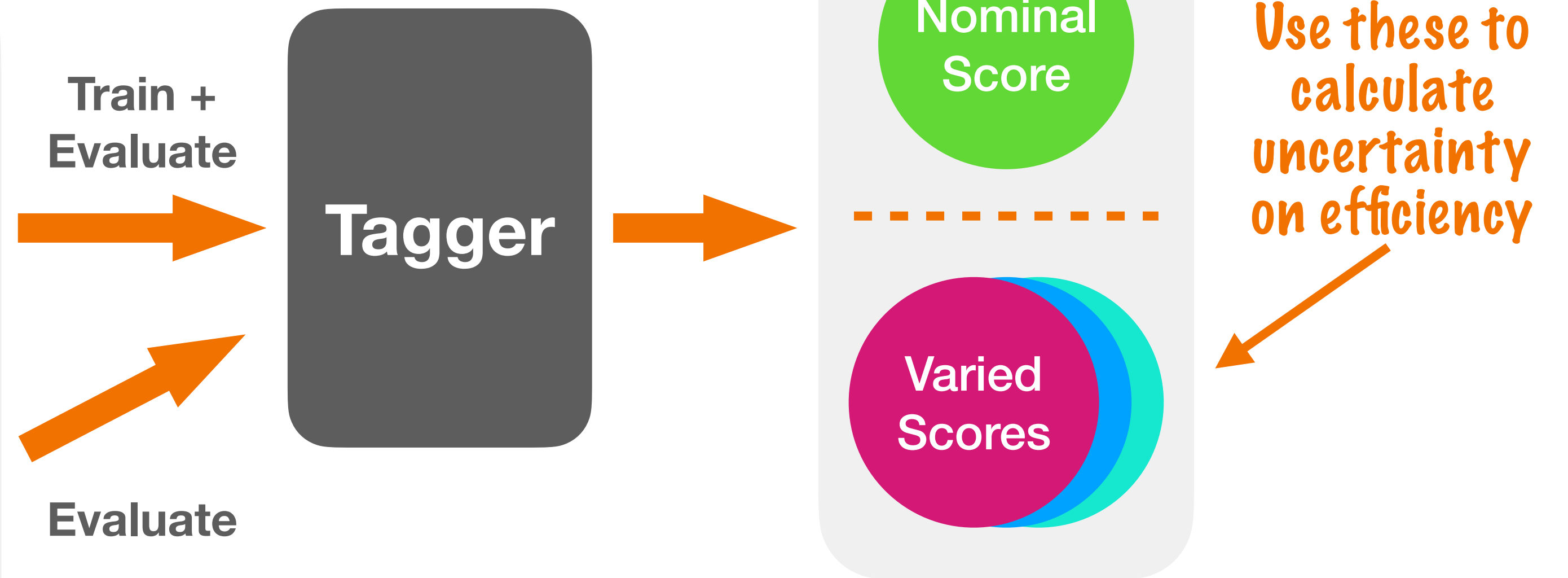
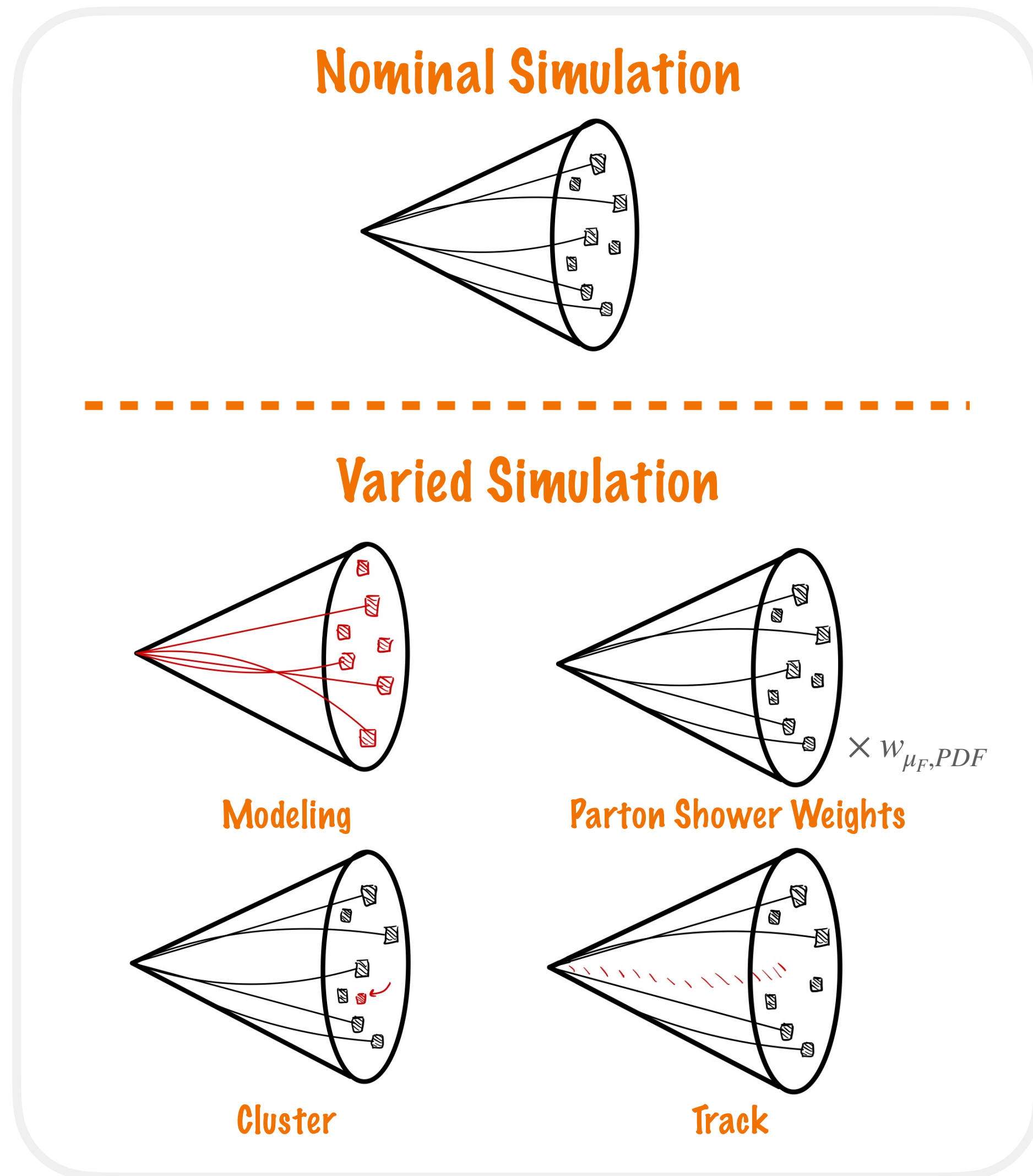
[cds:2866276](#)

© 2023 CERN for the benefit of the CMS Collaboration. CC-BY 4.0 license

Measuring scale factors is difficult, and only possible within collaborations.

**Can we find something approximate everyone can use?**

# Bottom-up Uncertainties

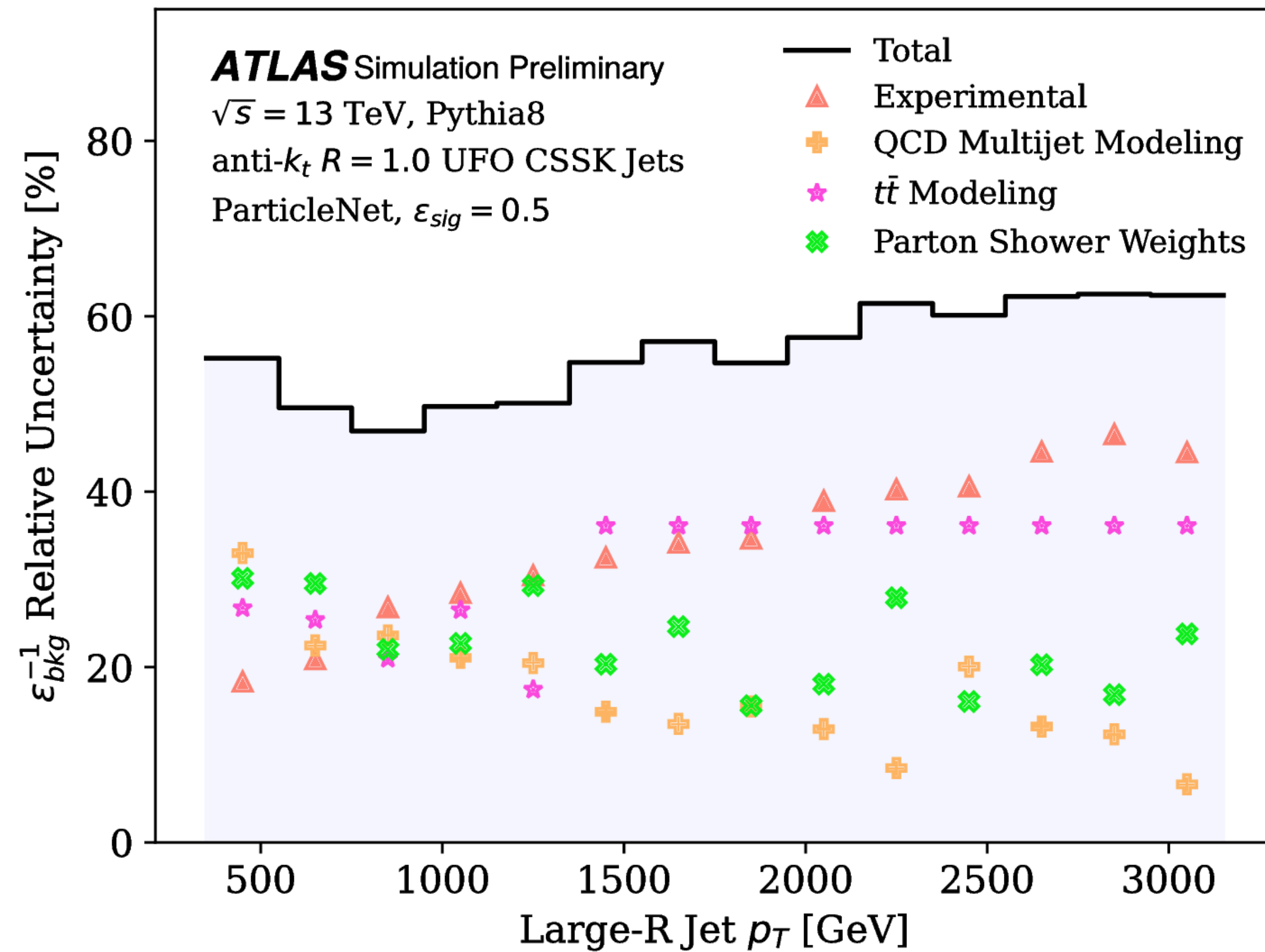


## Benefits

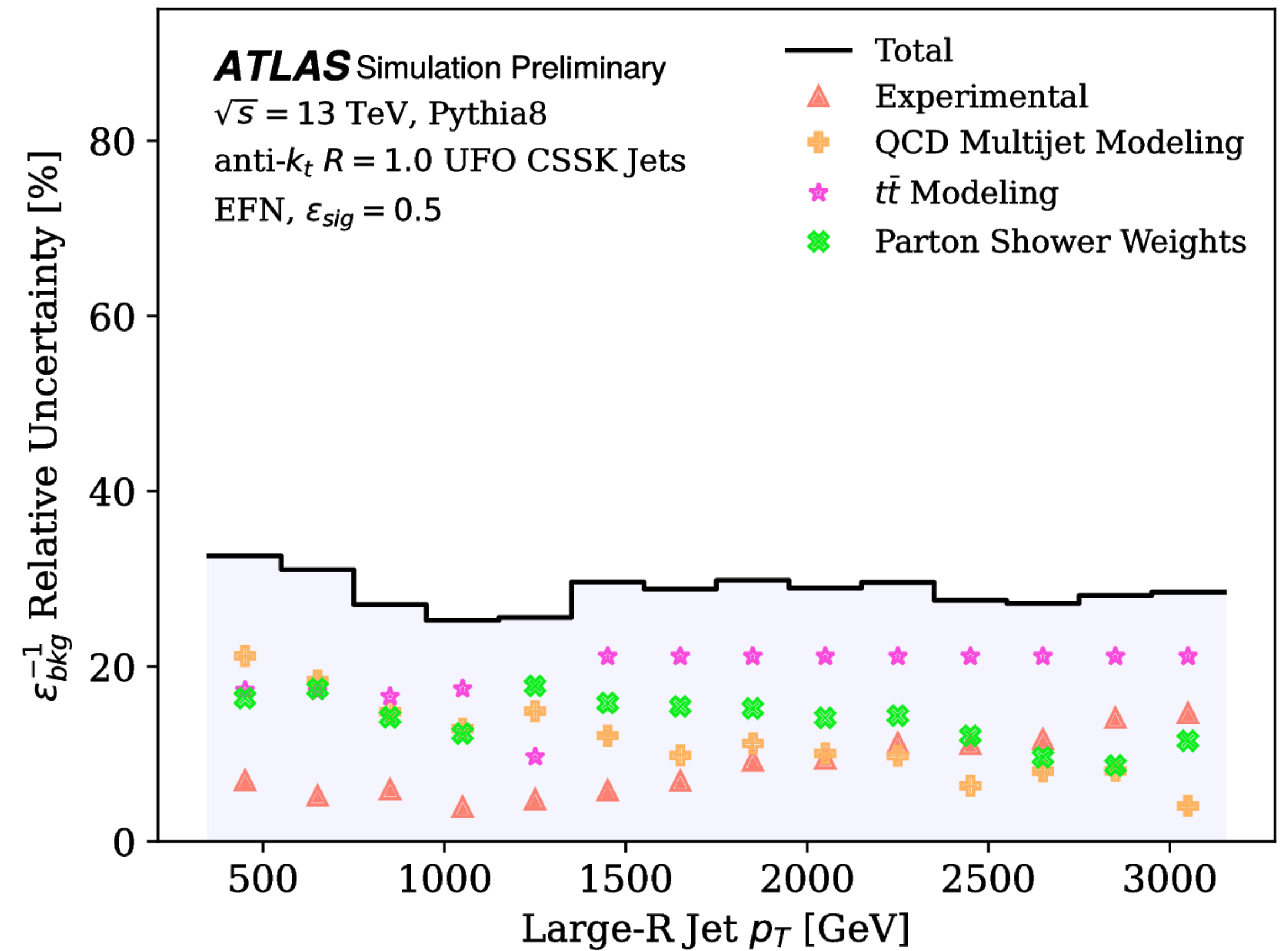
- No data required
- Once varied simulation is generated, can be used for arbitrary tagger
- Can define uncertainties on tagger efficiency with **no signal enriched region in data**

# Top Tagger Uncertainties

Note these **are not** scale factor uncertainties. Expected to be conservative, but relative sensitivity of taggers is important.



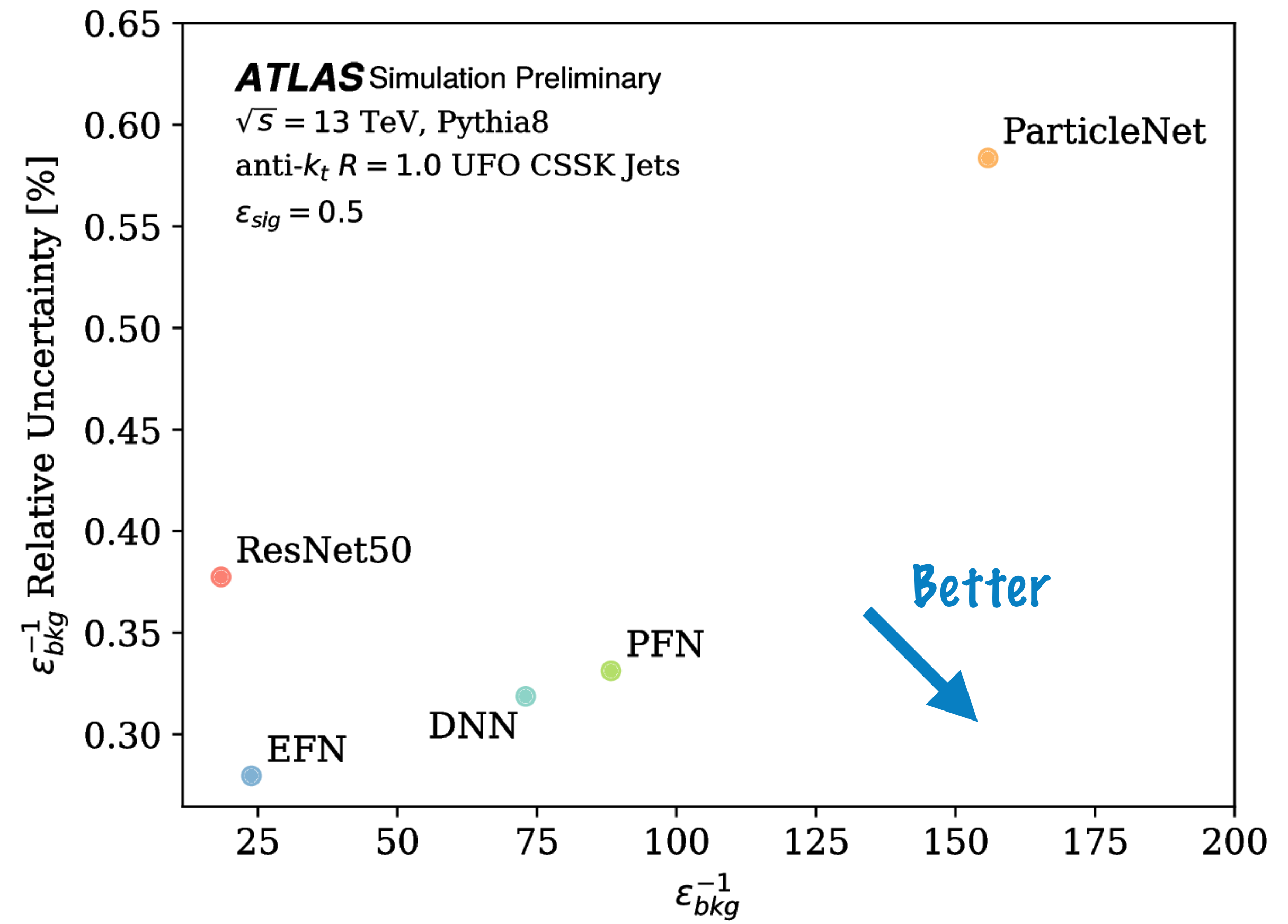
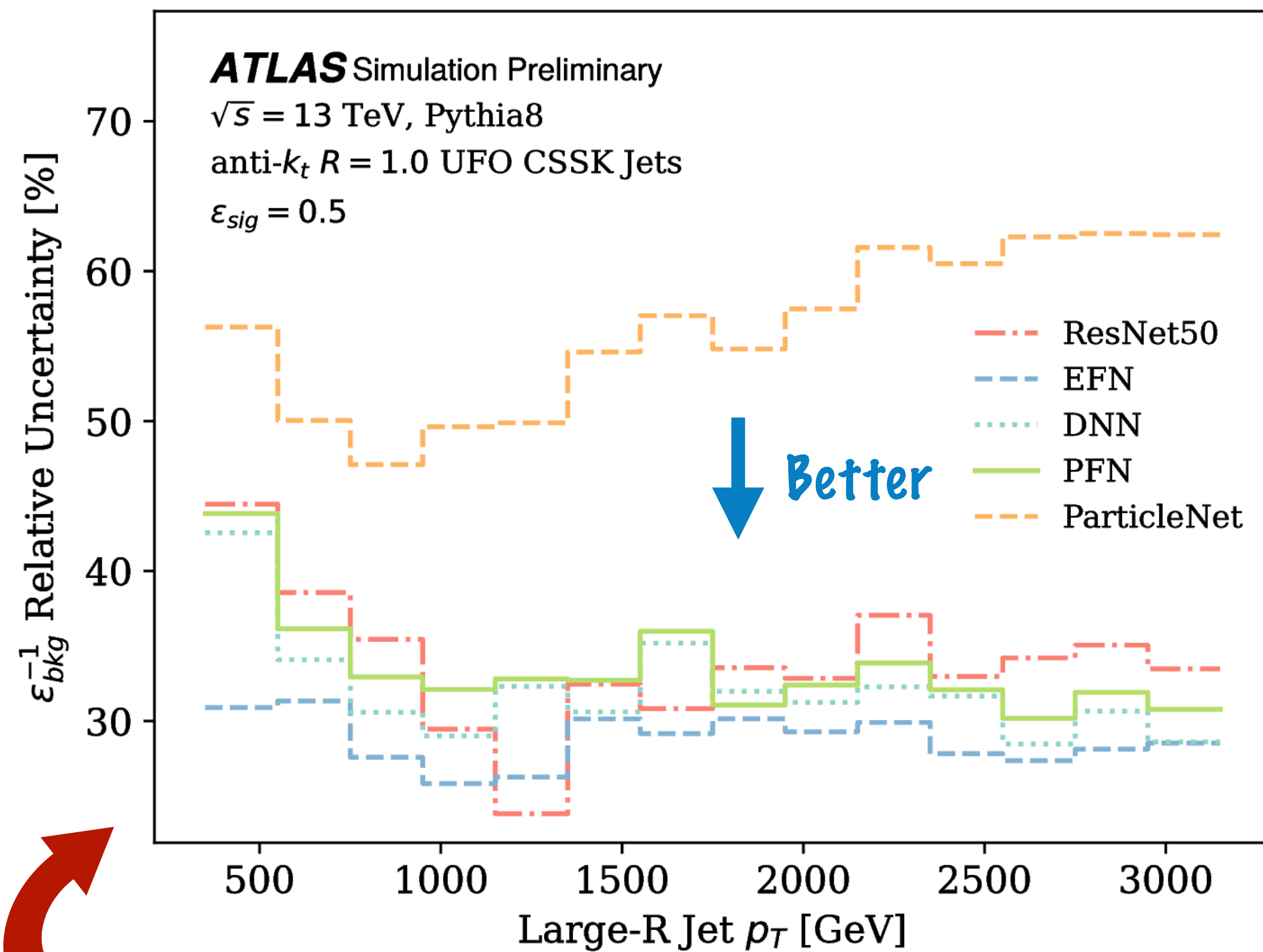
Large and Powerful GNN



Theory Motivated IRC Safe Tagger

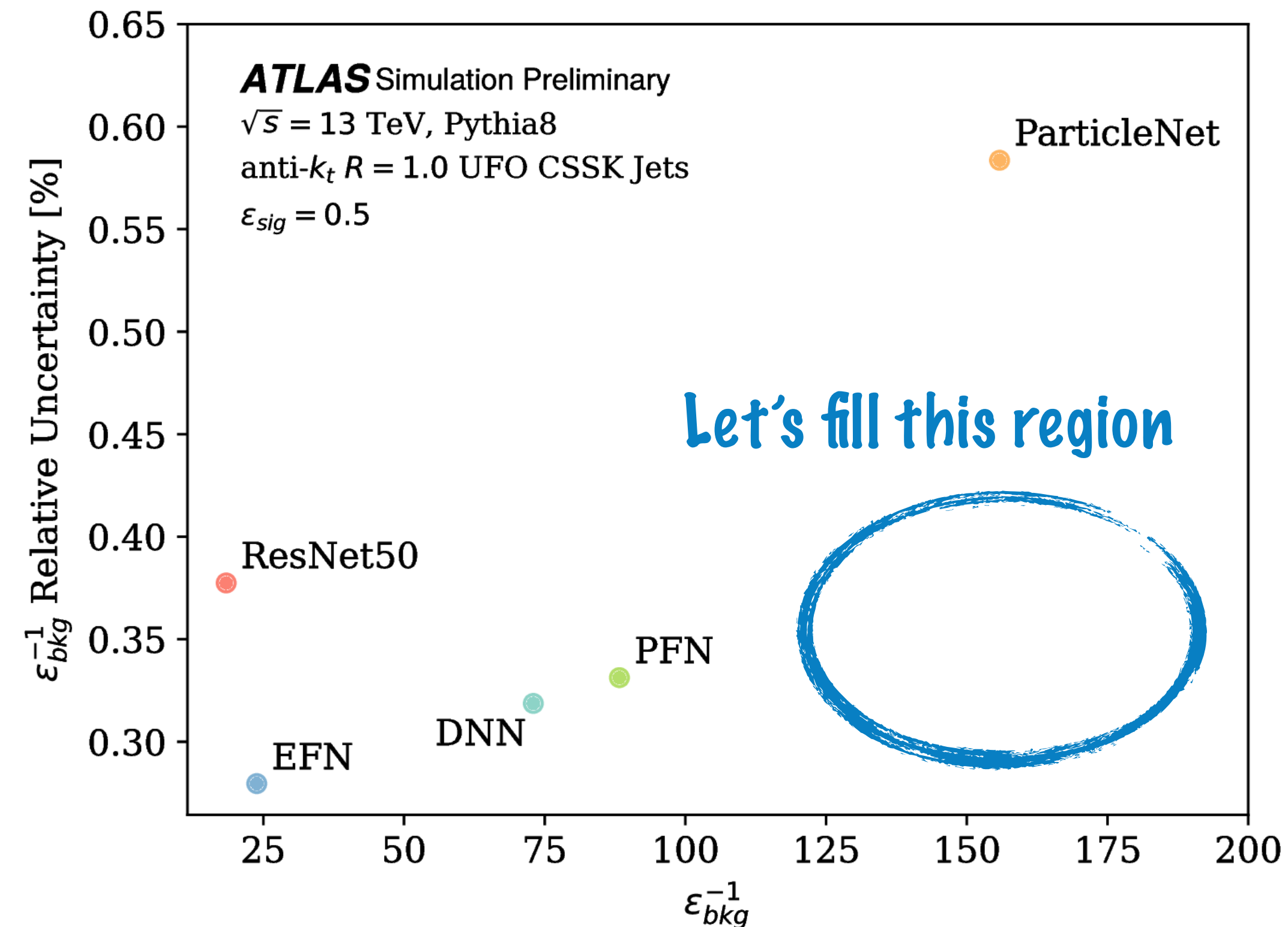


# Uncertainty Comparison



Larger uncertainties here are expected to produce larger SF uncertainties

# Conclusions



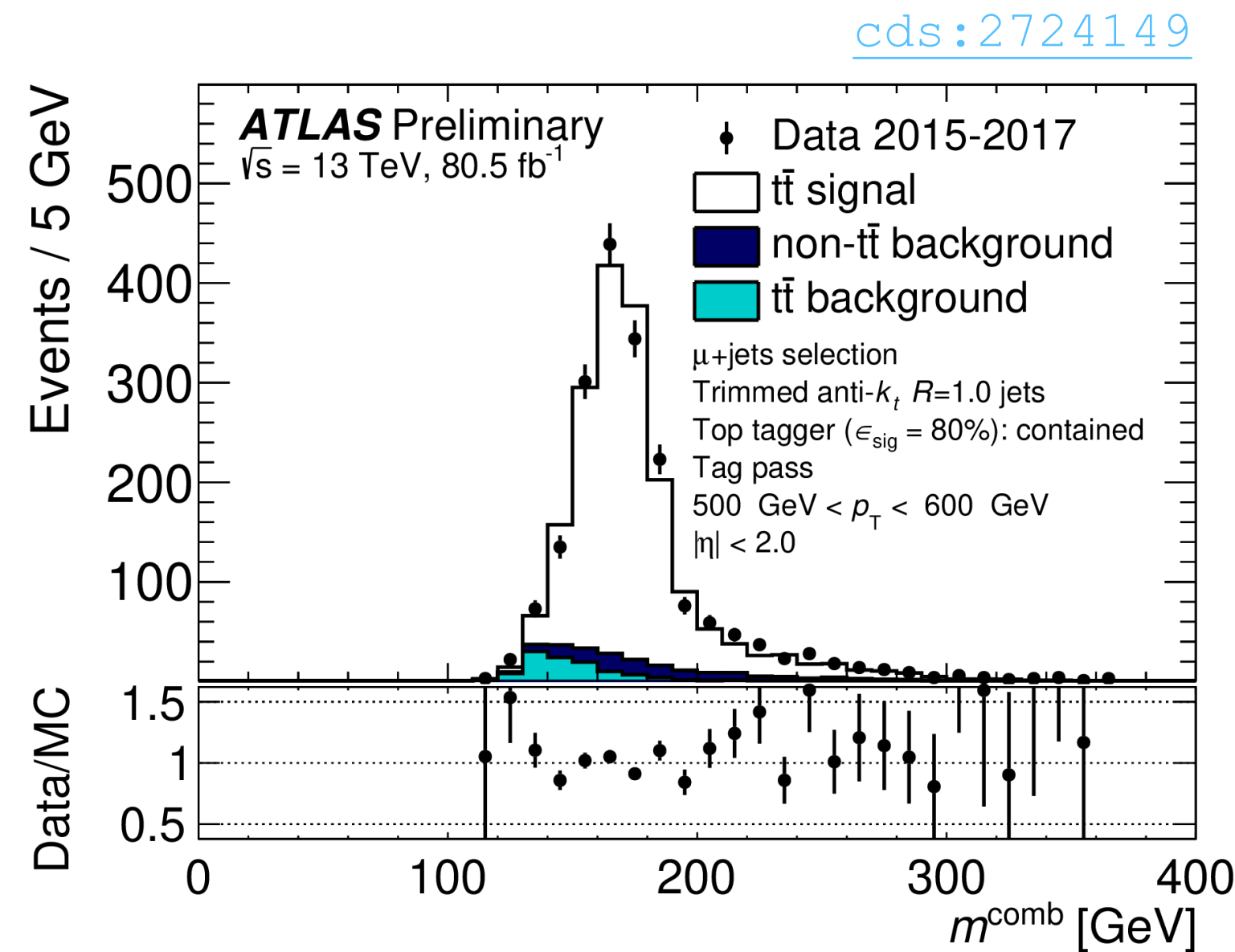
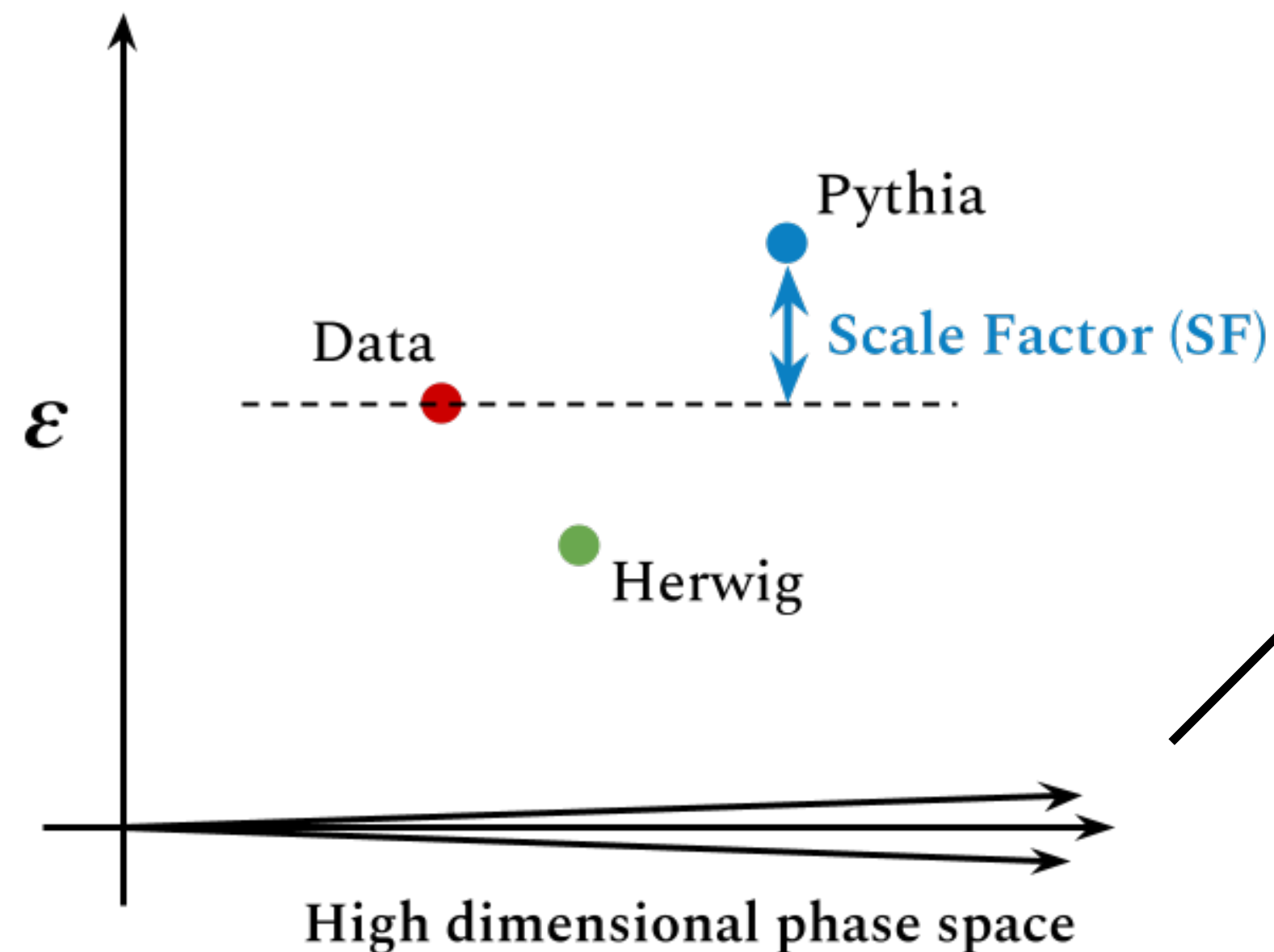
- Powerful ML based jet tagging is deployed and producing physics!
- However, the more powerful the tagger, the larger the uncertainties
  - Could be limiting for some analyses

**The new frontier is high performance and low uncertainties**

**Backup**

# A Brief Aside on Scale Factors

- Both ATLAS and CMS train taggers on MC, but need to know efficiency in data
- Measure **scale factor** to correct MC efficiency to data efficiency



Fit normalizations (N) of MC distributions to data

$$\epsilon_{\text{data}}(p_T) = \frac{N_{\text{fitted signal}}^{\text{tagged}}(p_T)}{N_{\text{fitted signal}}^{\text{tagged}}(p_T) + N_{\text{fitted signal}}^{\text{not tagged}}(p_T)}$$

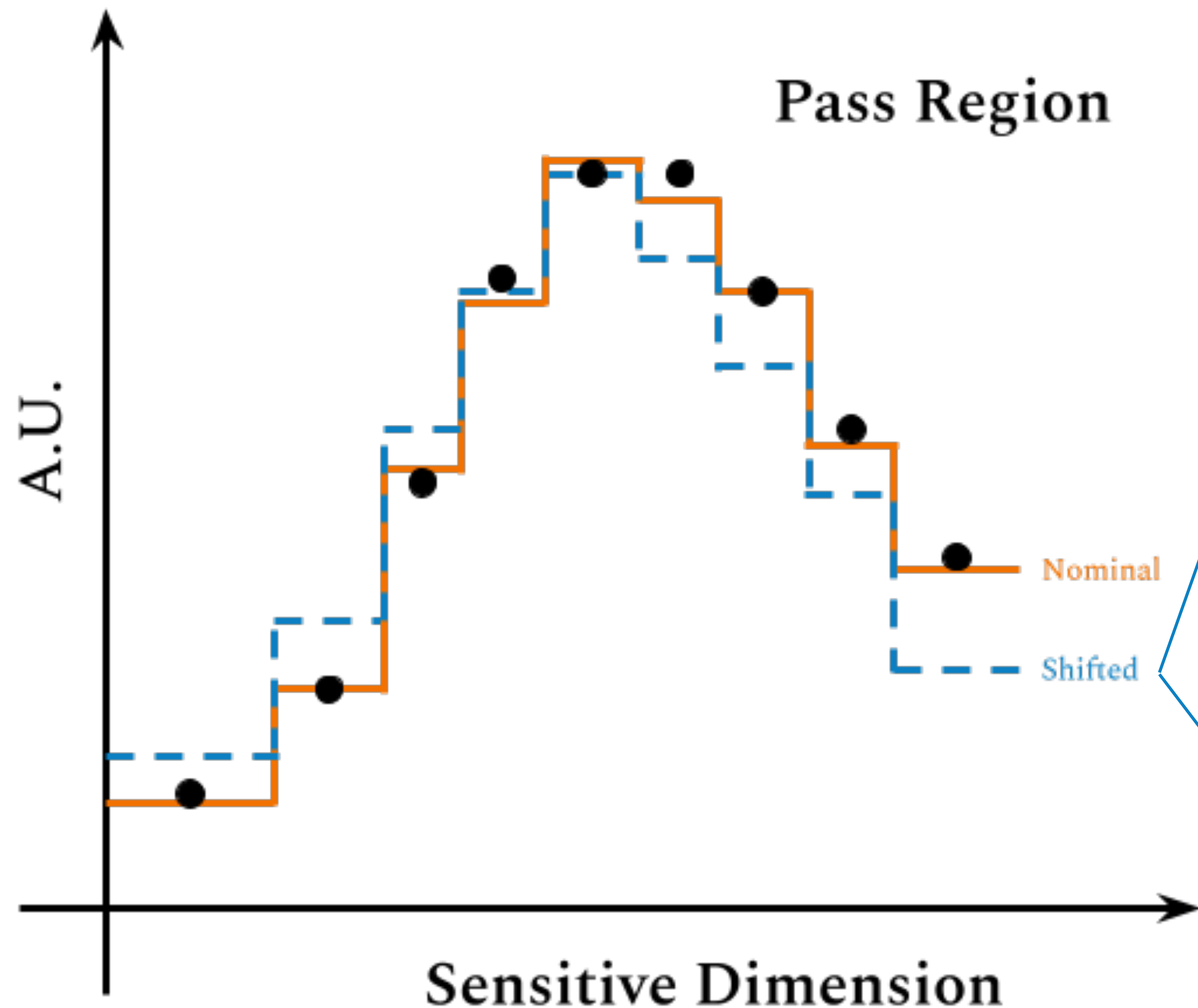
$$\text{SF}(p_T) = \frac{\epsilon_{\text{data}}(p_T)}{\epsilon_{\text{MC}}(p_T)}$$

Project to one dimension sensitive to efficiency

# Scale Factor Uncertainties

$$SF(p_T) = \frac{\epsilon_{\text{data}}(p_T)}{\epsilon_{\text{MC}}(p_T)}$$

Like any measurement SFs have **uncertainties**:



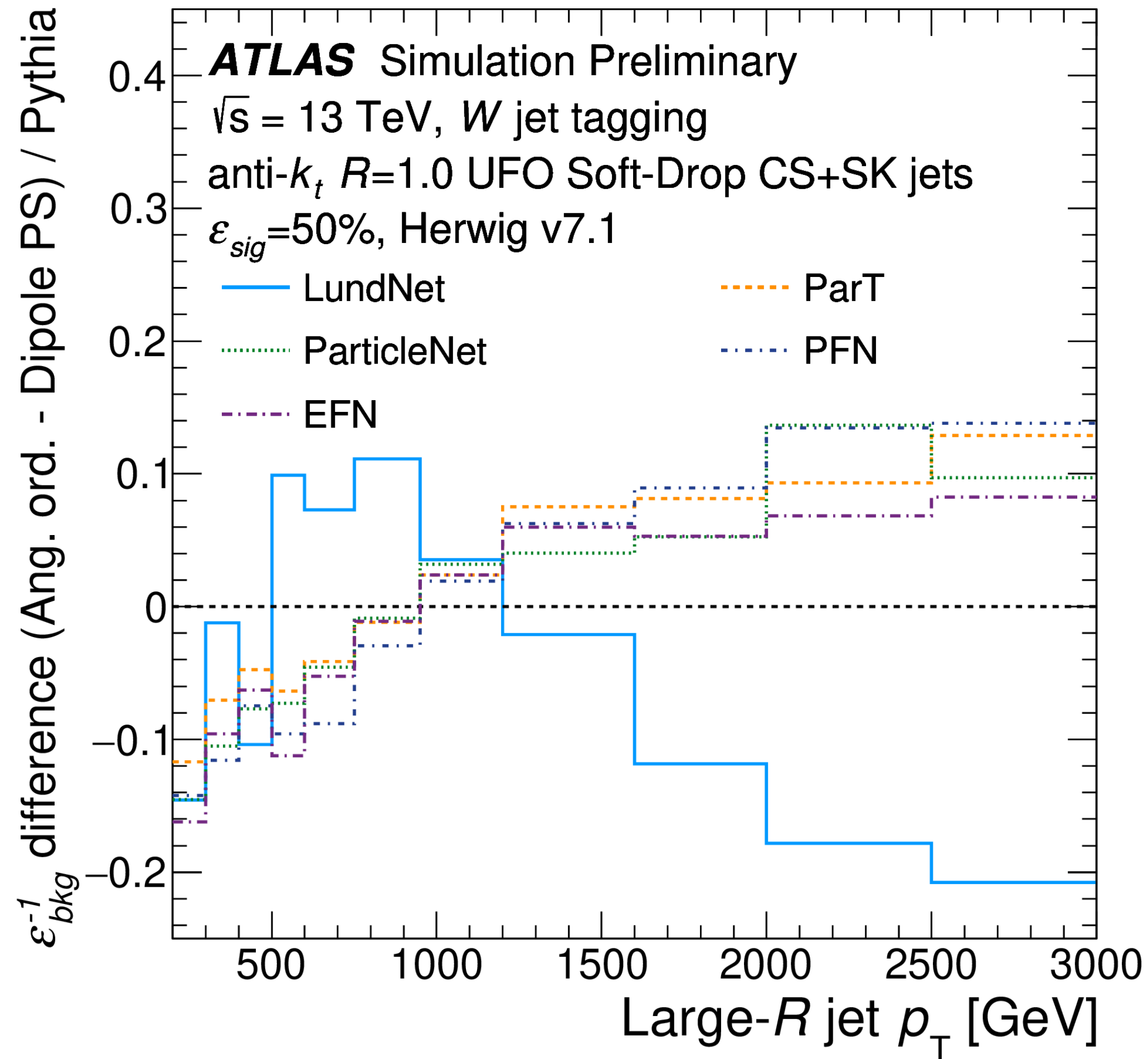
## Theoretical

- Parton shower + hadronization modeling
- Renormalization scale
- Cross sections
- PDFs

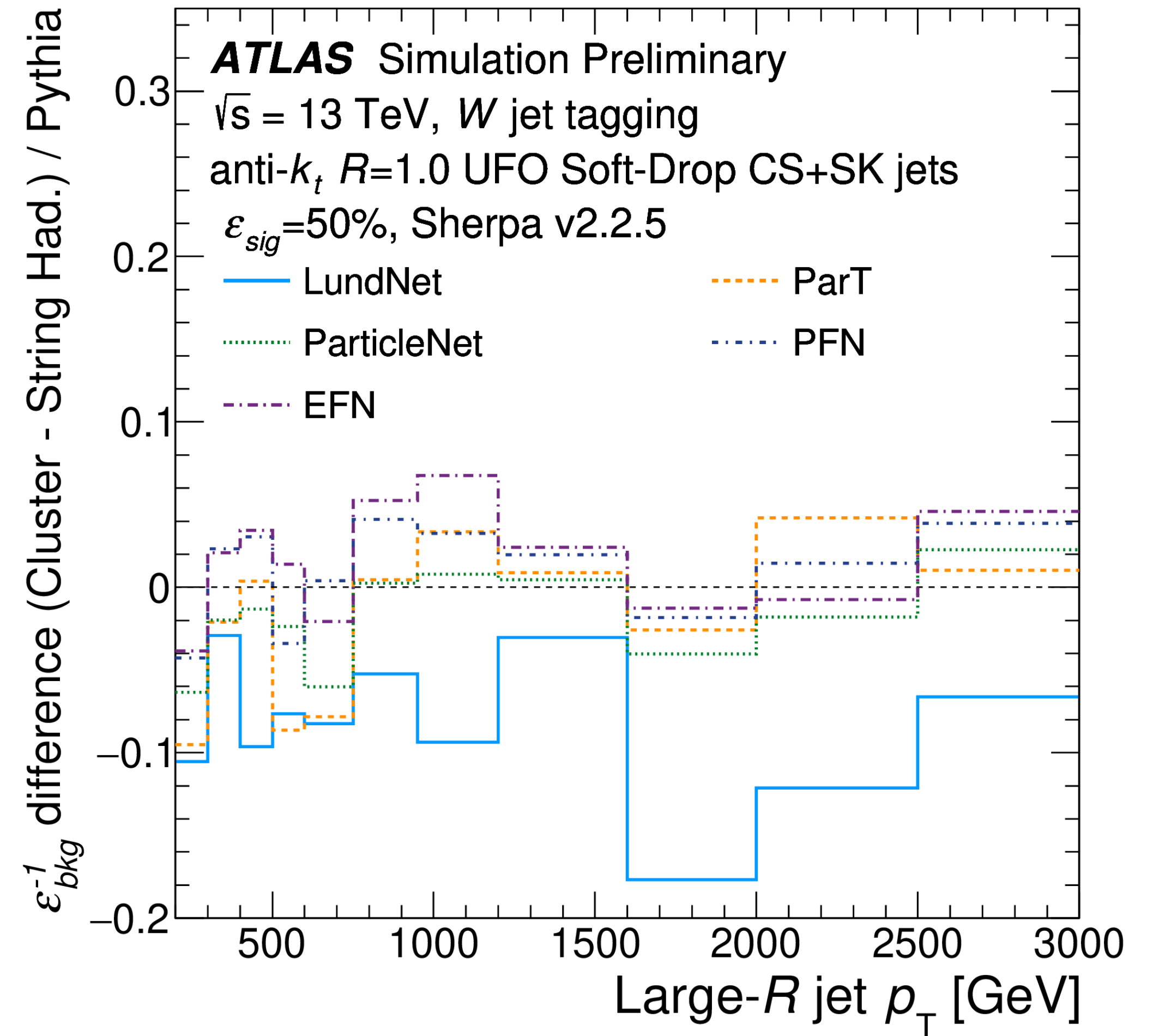
## Experimental

- Jet mass scale / resolution or similar
- Statistical
- Other SFs (e.g. b-tagging)
- Luminosity

# W Tagger Modeling Dependence



Parton Shower



Hadronization

# Top Tagging Systematic Variations

Modify nominal  
Alternative samples  
Pythia shower weights

## Experimental

- Calorimeter Clusters<sup>1</sup>
  - Energy Scale (Up / Down)
  - Energy Resolution
  - Position resolution
- Tracks
  - Fake rate
  - Efficiency
  - Sagitta bias

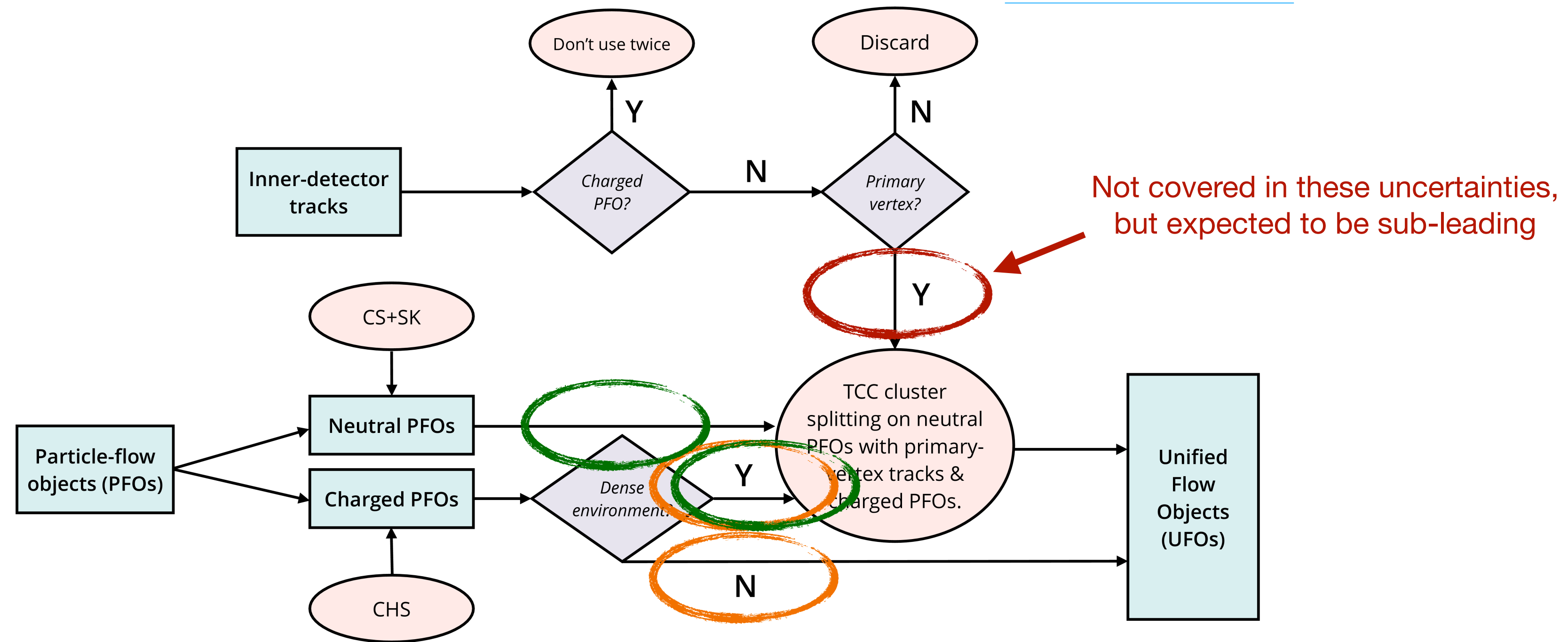
<sup>1</sup> - [arxiv: 1912.0983](https://arxiv.org/abs/1912.0983), [arxiv:1903.02942](https://arxiv.org/abs/1903.02942), [arxiv:2108.09043](https://arxiv.org/abs/2108.09043)

## Theoretical

- $t\bar{t}$  modeling
  - Compare Pythia to Herwig in SM  $t\bar{t}$  samples
- QCD multijet modeling
  - Compare Herwig angular ordered to dipole parton shower
  - Compare Sherpa cluster to string based hadronization model
- Renormalization scale
  - Vary scale up/down by factors of 2
- PDFs
  - Vary PDFs up/down

# Experimental Uncertainties

[arxiv:2009.04986](https://arxiv.org/abs/2009.04986)



## Tracks

- Apply to charged and “merged” UFOs
- Track fake rate and efficiency
- Track bias

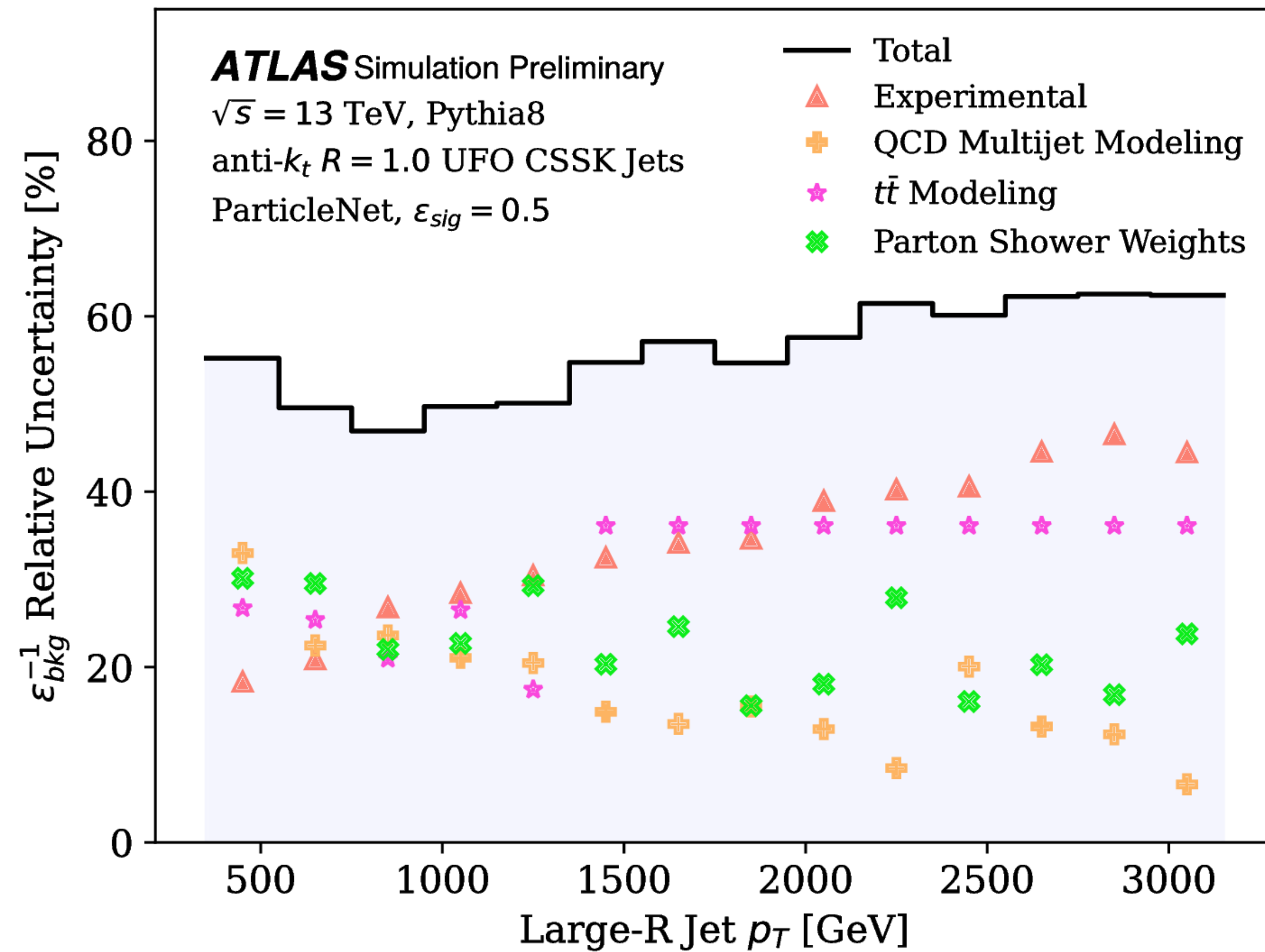
## Calorimeter Clusters

- Apply to neutral and “merged” UFOs
- Cluster energy scale and resolution
- Cluster position resolution

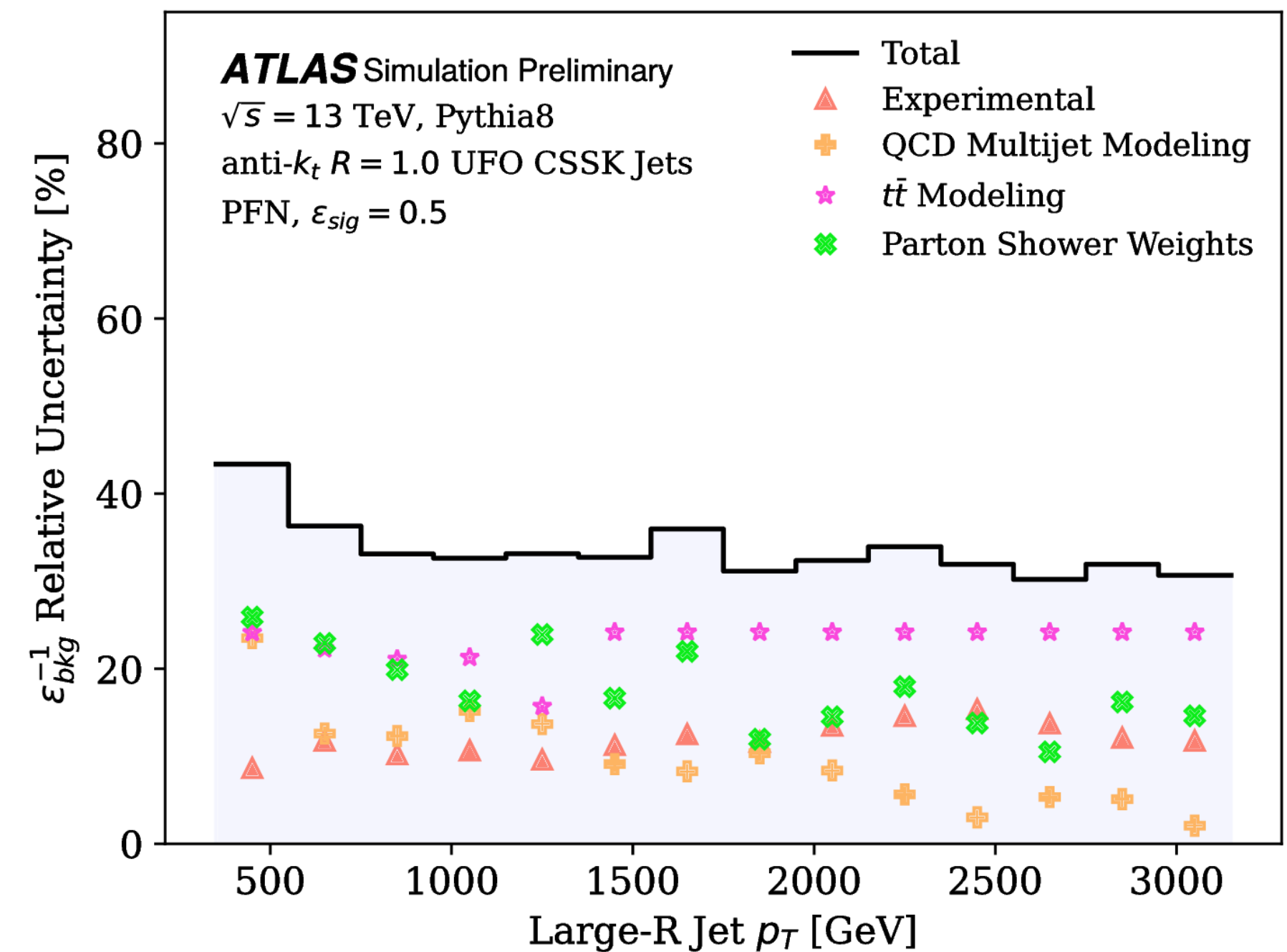


# Top Tagger Uncertainties

Note these **are not** scale factor uncertainties. Expected to be conservative, but relative sensitivity of taggers is important.



Large and Powerful GNN



Deep Set (PFN)

- ParticleNet more sensitive to modeling
- Dramatically more sensitive to experimental variations

# Top Tagger Uncertainties

Note these **are not** scale factor uncertainties. Expected to be conservative, but relative sensitivity of taggers is important.

