# Image and Sequence Based Jet Tagging Applications on Experiments

## Michael Kagan

SLAC

**Model Performance**
Accuracy
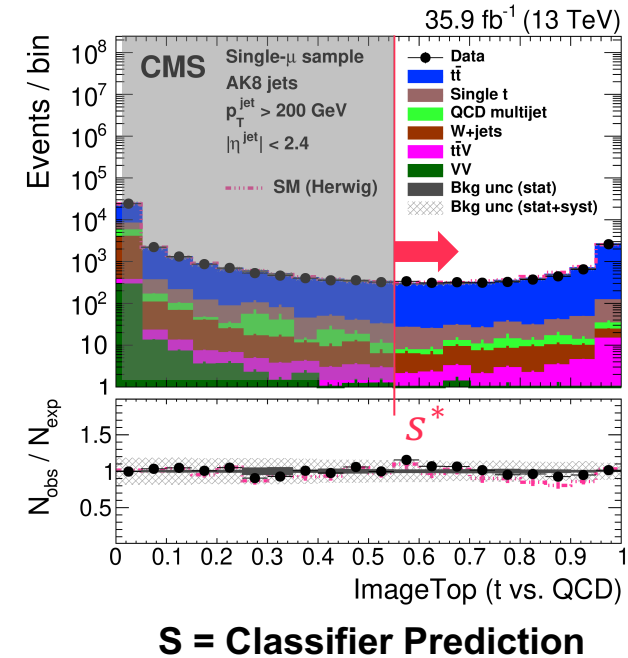
**Interpretability**
Model & Uncertainties

**Calibration**
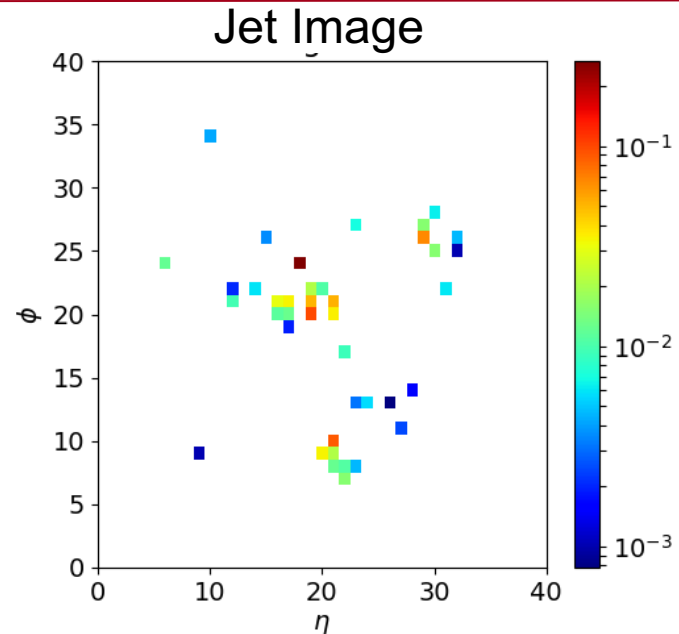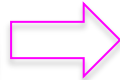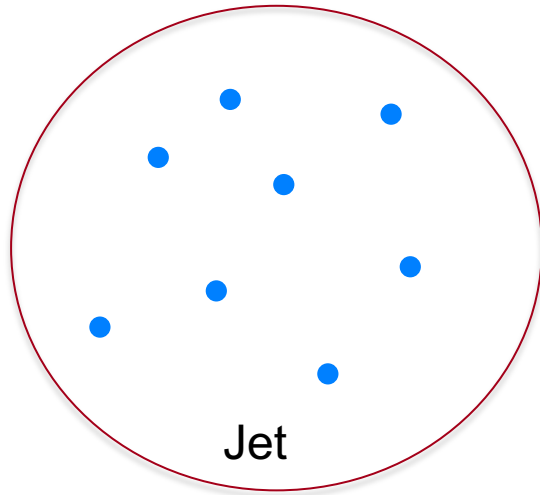Uncertainties

# Calibrating Jet Taggers: Scale Factor

$$SF = \frac{\epsilon_{data}}{\epsilon_{MC}} = \frac{p_{data}(s > s^*)}{p_{MC}(s > s^*)} = \frac{\left(\dfrac{N^{pass} - N^{pass}_{bkg}}{N - N_{bkg}}\right)_{data}}{\left(\dfrac{N^{Pass}_{sig}}{N_{sig}}\right)_{MC}}$$

$$p_{data}(s) = \pi_{sig} p_{sig}(s) + (1 - \pi_{sig}) p_{bkg}(s)$$
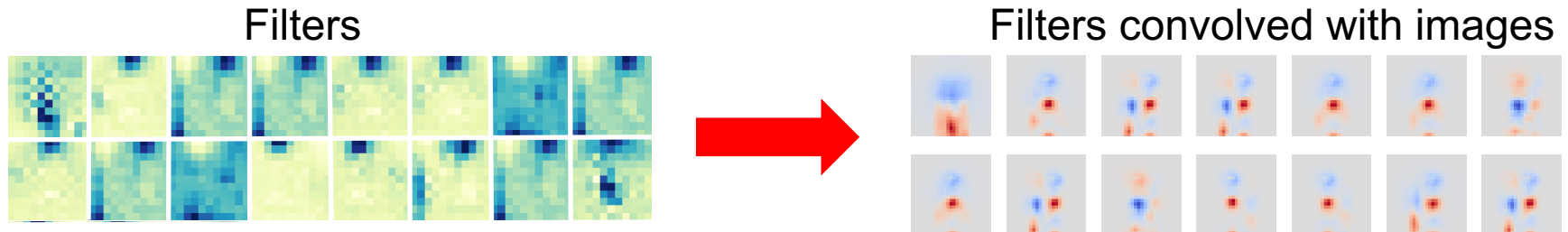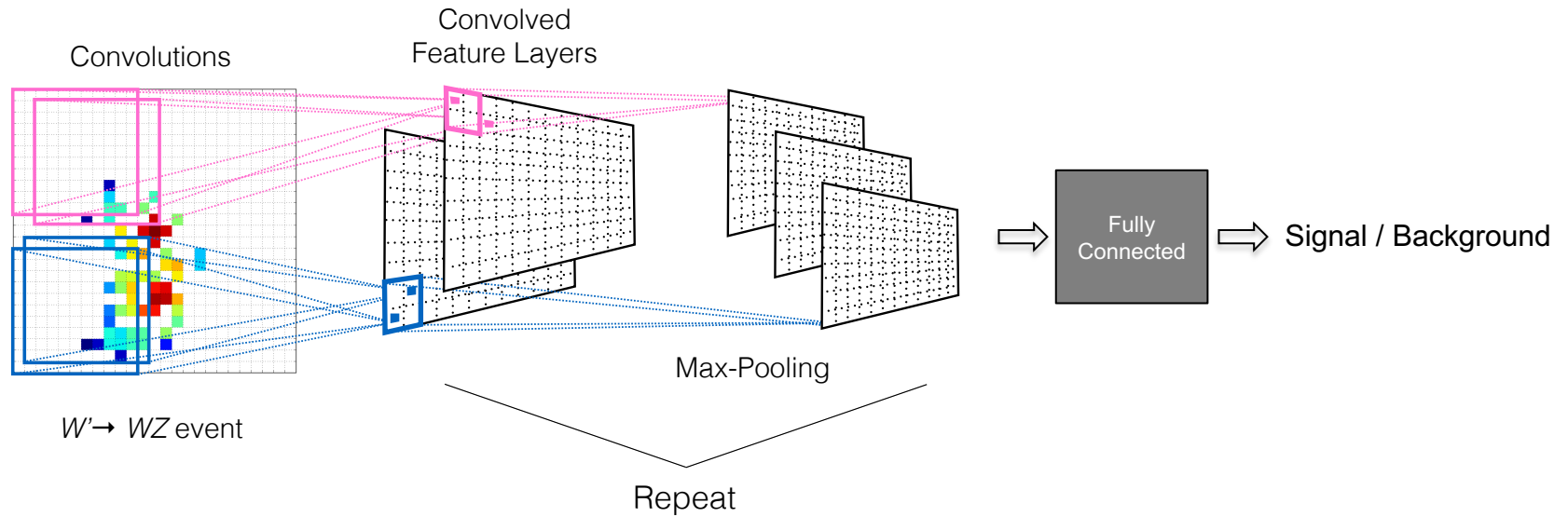


**S = Classifier Prediction**

- Correct the MC efficiency of a cut on classifier output
  - Little insight into "why" a scale factor deviates from unity
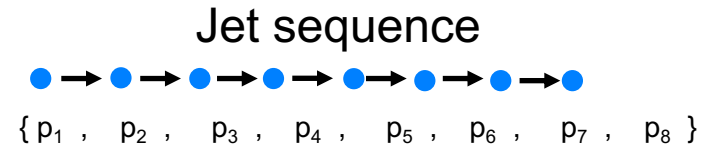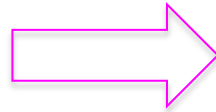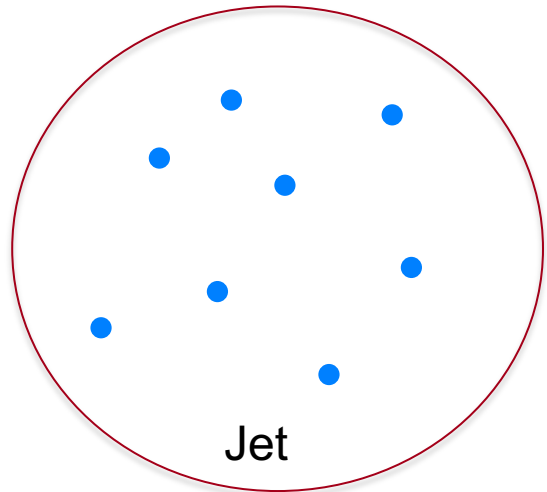
# Computer Vision and Jets



Jet Image

- **Jet-image** $-$2D representation of jet as distribution of energy over $\eta - \phi$

- Multi-channel "color" jet images $-$ separate images for different detectors (calorimeter / track) or particles (charged / neutral hadrons, muon, etc.)

- How to deal with track images?
  - More pixels may improve performance
  - Cost: larger models and more memory needed

Max-Pooling

$W' \to WZ$ event

Repeat

# Boosted Boson Type Tagging

*Jet ETmiss*

Benjamin Nachman and Ariel Schartzman

Filters

Filters convolved with images

*anford Univ*

*ch 26, 201*

# Sequence Modeling



Jet sequence

$\{ p_1 , \quad p_2 , \quad p_3 , \quad p_4 , \quad p_5 , \quad p_6 , \quad p_7 , \quad p_8 \}$

Jet

- Jets are a grouping of a variable number of particles
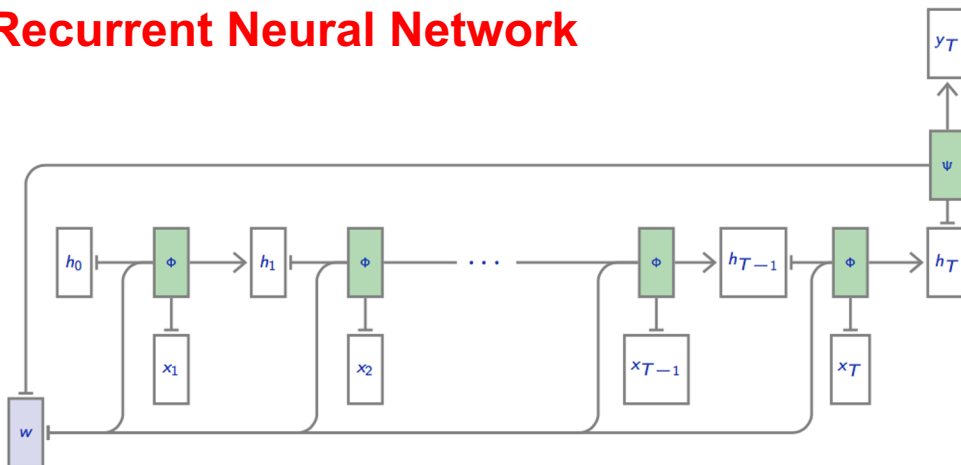- With physically motivated ordering: **jet as a sequence**

**Recurrent Neural Network**

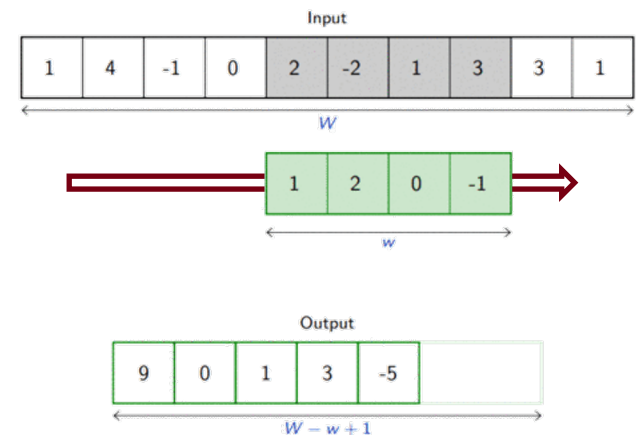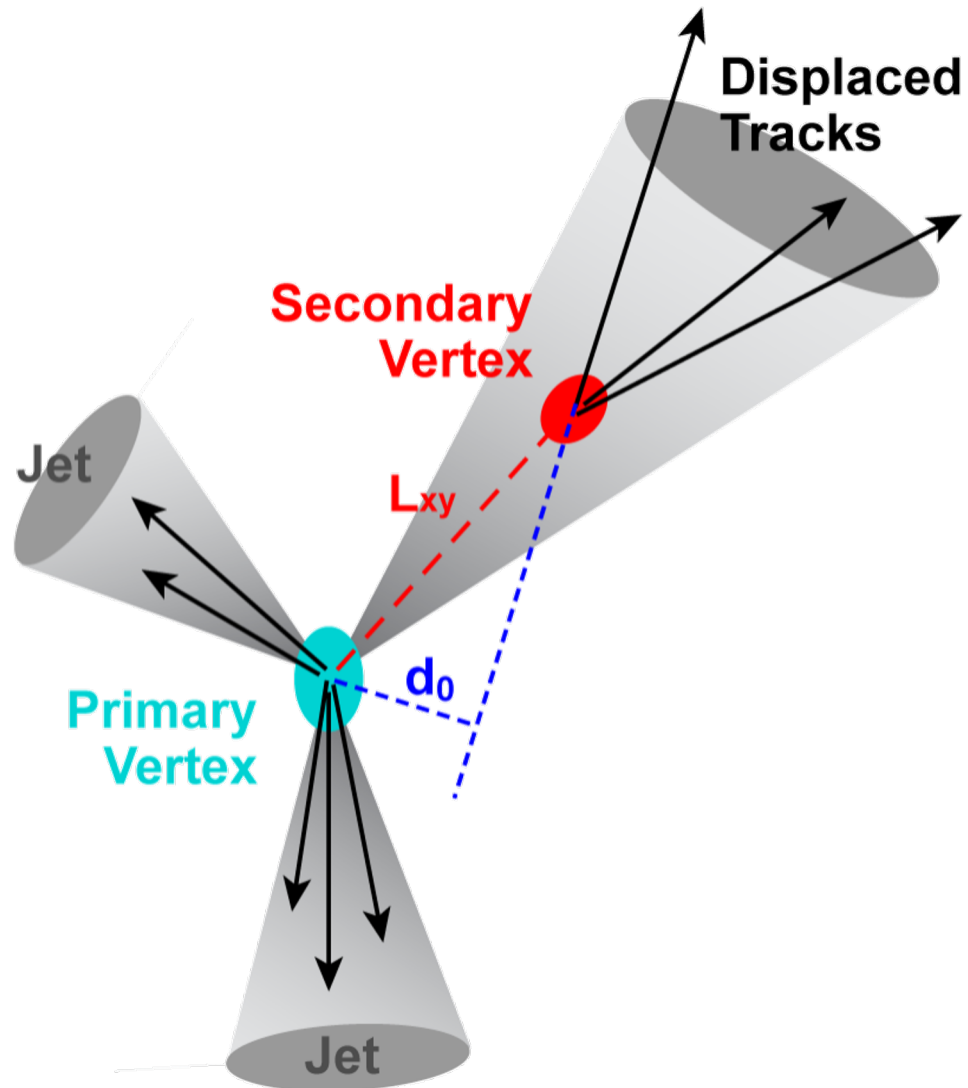**1D Convolutions**



Image Credit: Fleuret, Deep Learning Course
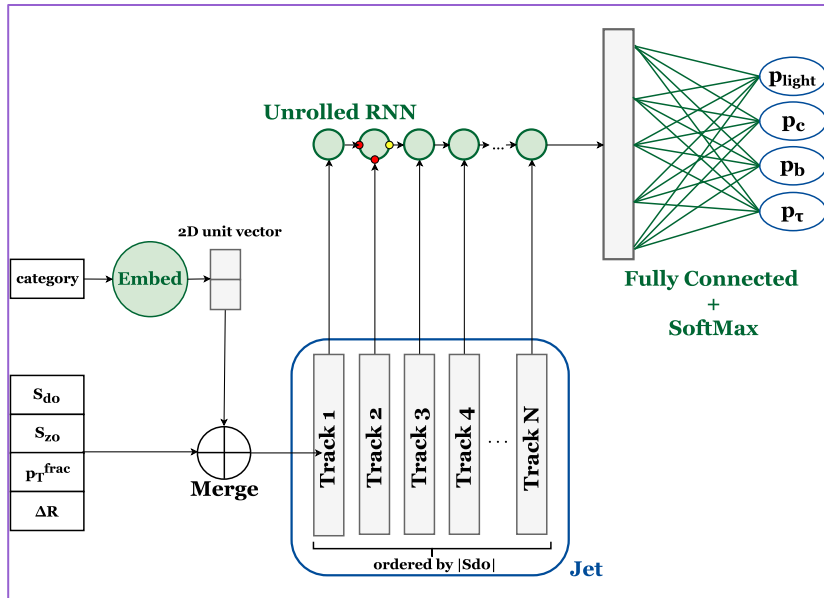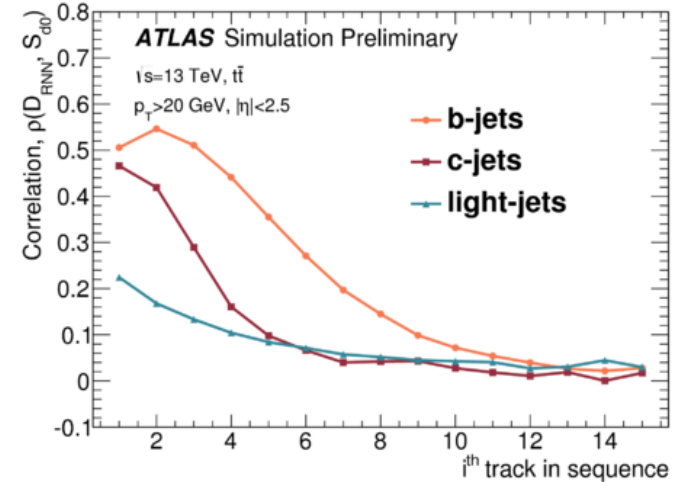
# B-tagging



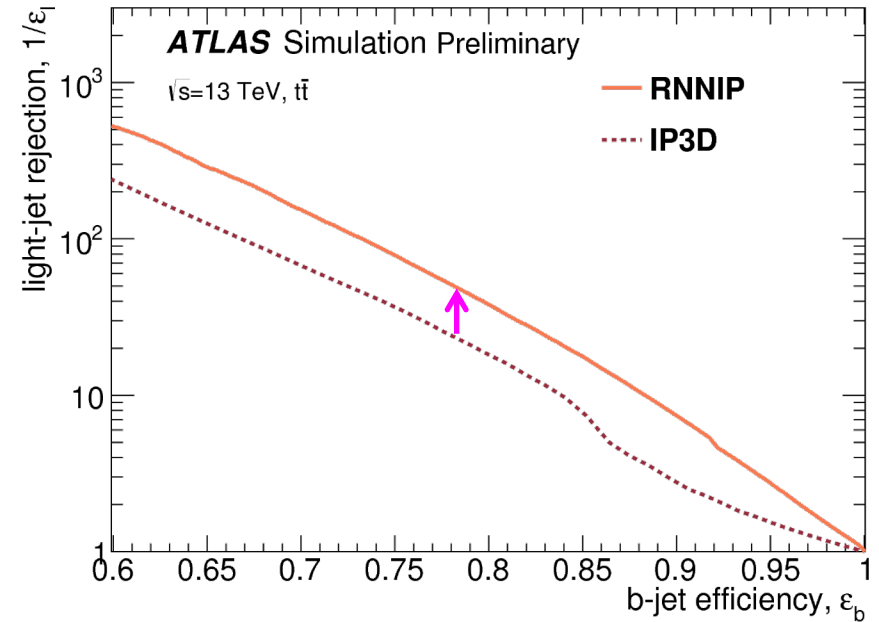Displaced Tracks

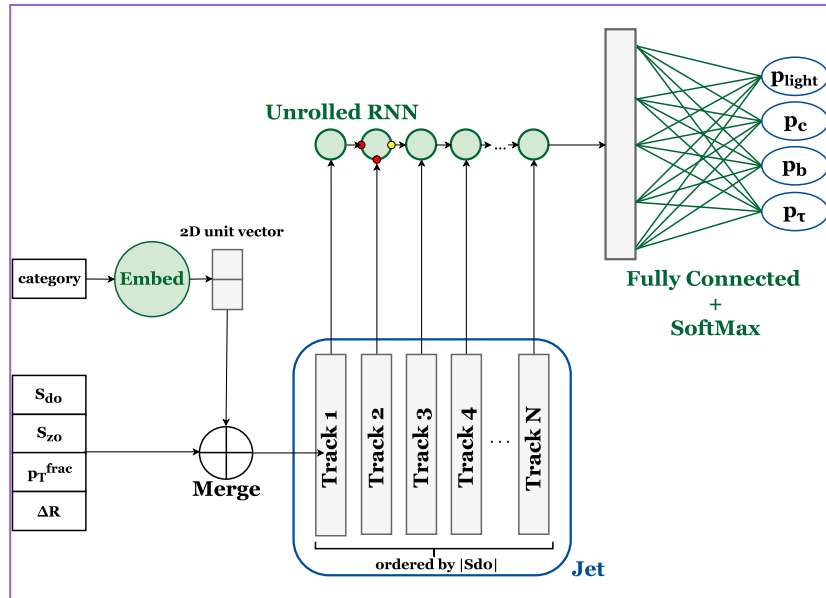Secondary Vertex

Jet

$L_{xy}$

$d_0$
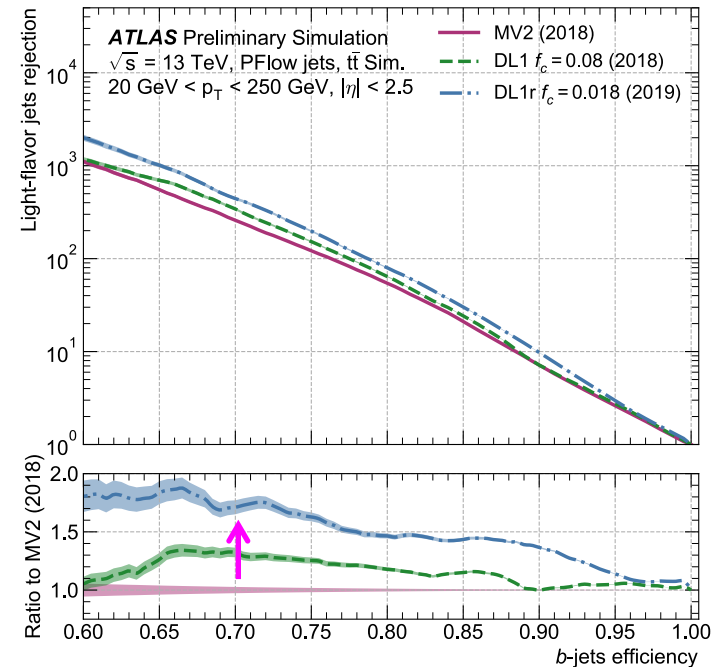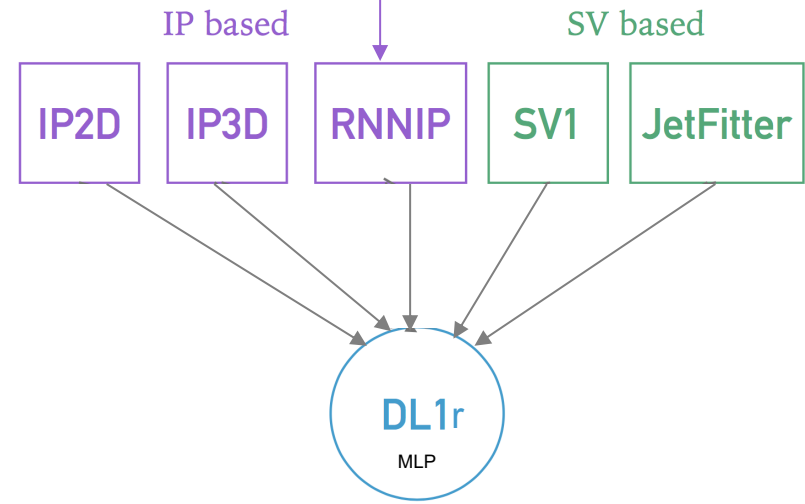
Primary Vertex

Jet

# Recurrent Neural Net b-tagging – RNNIP

ATL-PHYS-PUB-2017-003

# Recurrent Neural Net b-tagging – RNNIP

# DL1r

# DL1r Calibration

# DL1r Calibration Uncertainties

- Theory modeling among largest uncertainties

- Significant recent SF reduction, independent of tagger, e.g. from
  - Improving charge deposition modeling in Silicon
  - Better method to estimate sample flavour composition

- Still difficult to separate "Model is learning bad correlations" from imperfect calibration methods

# DeepJet

**Number of particles/SV**

**Feature extractor** convolution performed on each particle / SV [1x1]

**LSTM layers** Builds a summary of the information extracted in each set of features

**Correlations and classification**

# Boosted Jet Tagging



q/g

W/Z→qq

h→bb

t→Wb→qqb

Image Credit: arXiv:1909.12285

# ML Boosted Jet Taggers on CMS

## Deep AK8



| Category | Label |
|----------|-------|
| Higgs | H (bb) |
| | H (cc) |
| | H (VV*→qqqq) |
| Top | top (bcq) |
| | top (bqq) |
| | top (bc) |
| | top (bq) |
| W | W (cq) |
| | W (qq) |
| Z | Z (bb) |
| | Z (cc) |
| | Z (qq) |
| QCD | QCD (bb) |
| | QCD (cc) |
| | QCD (b) |
| | QCD (c) |
| | QCD (others) |

## ImageTop

- 6 channels
  - All PF candidates
  - Charged hadron
  - Neutral hadron
  - Photon
  - Electron
  - Muon

# Tagger Performance

# Top Tagging Scale Factors

# Top Misidentification Scale Factors



35.9 fb$^{-1}$ (13 TeV)

**CMS**

Dijet sample
Top quark tagging
QCD multijet: MG+P8

Legend:
- $m_{SD}+\tau_{32}$
- $m_{SD}+\tau_{32}+b$
- HOTVR
- $N_3$-BDT (CA15)
- BEST
- ImageTop
- ImageTop-MD
- DeepAK8
- DeepAK8-MD

jet p$_T$ [GeV]

# Quark versus Gluon with Jet Images

**ATLAS** Simulation Preliminary

convolutional filters

Max-pooling

dense layer

quark jet

gluon jet

3x

Truth

Tracks

Topo Clusters

Towers

# Quark versus Gluon with Jet Images

# Sensitivity to Generators → Representation Learning



**ATLAS** Simulation Preliminary
$\sqrt{s}$ = 13 TeV
Anti-$k_t$ EM+JES R=0.4
$|\eta| < 2.1$, 150 GeV $< p_T <$ 200 GeV

- CNN Truth, Trained on Pythia, Test on Pythia
- CNN Truth, Trained on Herwig, Test on Herwig
- CNN Truth, Trained on Pythia, Test on Herwig
- CNN Truth, Trained on Herwig, Test on Pythia



JHEP 01 (2017) 110

- Pythia CNN on Pythia Color Images
- Herwig CNN on Pythia Color Images
- Pythia CNN on Herwig Color Images
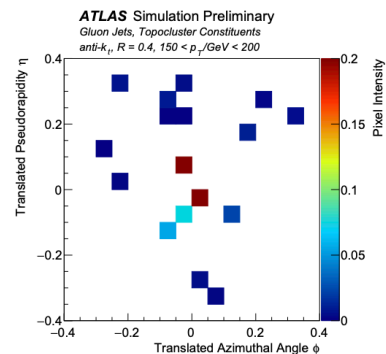- Herwig CNN on Herwig Color Images



Deep Neural Network Performance

- PYTHIA
- PYTHIA+VINCIA
- Herwig (angular)
- Herwig (dipole)
- Sherpa

- PYTHIA $\tau_{21}$
- PYTHIA mass

$50 < m < 110$ GeV, $250 < p_T < 300$ GeV

- What you train on seems to have smaller impact than what you test on

- Robustness of the learned representations?

Phys. Rev. D 95, 014018 (2017

# Mitigating Dependencies

# Mitigating Dependencies

- With flexibility comes complexity:

  - Hard to control how models learn and utilize information
  - Potentially unwanted sensitivity to poorly modeled aspects of simulation
  - Potentially unwanted sculpting of key physics distributions like mass

- Decorrelation methods

  - Reweighting training distributions
  - DDT: Designing decorrelated taggers JHEP 05 (2016) 156
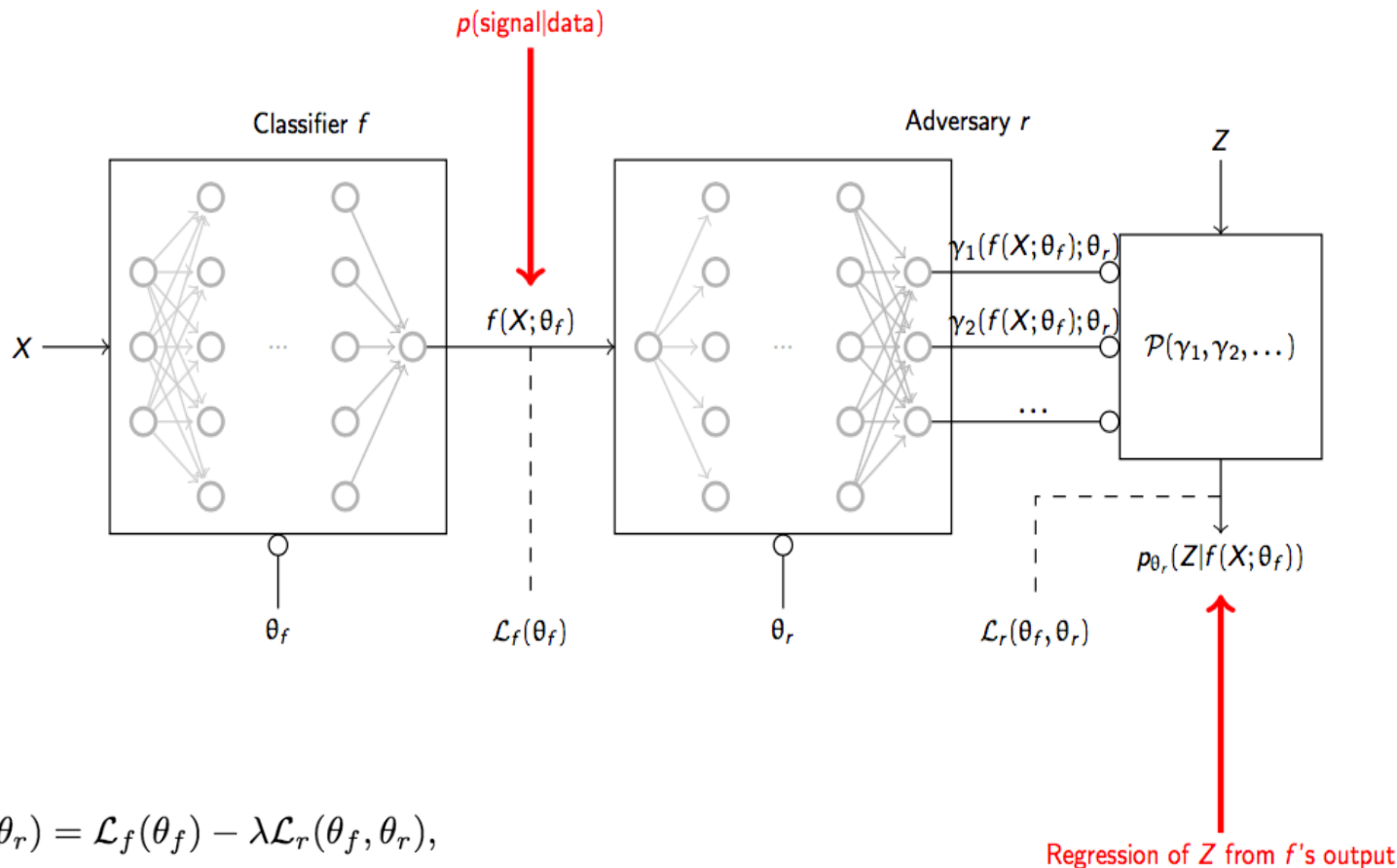  - DisCo: Distance Correlation regularization Phys. Rev. Lett. 125, 122001 (2020)
  - Adversarial Learning NeurIPS 2017, 981-990, Phys. Rev. D 96, 074034 (2017)

# Adversarial Learning
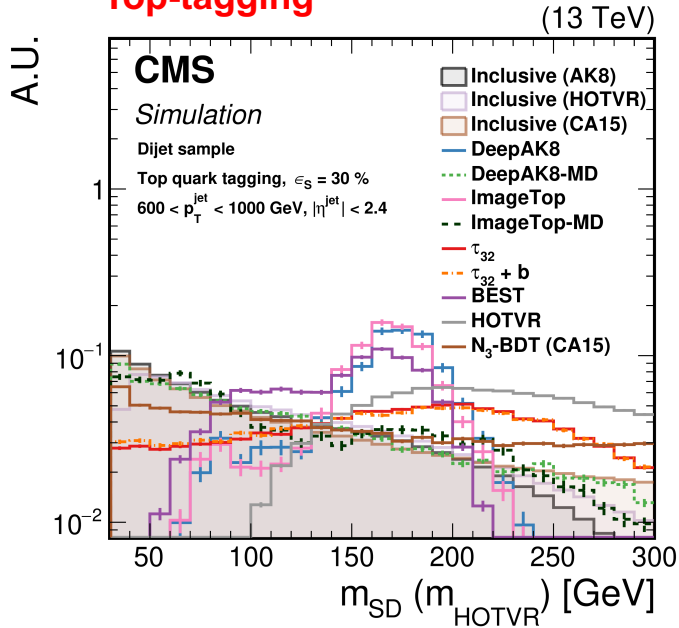


$p(\text{signal}|\text{data})$

Classifier $f$

Adversary $r$

$Z$

$X$

$f(X; \theta_f)$

$\gamma_1(f(X; \theta_f); \theta_r)$

$\gamma_2(f(X; \theta_f); \theta_r)$

$\mathcal{P}(\gamma_1, \gamma_2, \dots)$

$\theta_f$

$\mathcal{L}_f(\theta_f)$

$\theta_r$

$\mathcal{L}_r(\theta_f, \theta_r)$

$p_{\theta_r}(Z|f(X; \theta_f))$

Regression of $Z$ from $f$'s output

$$E_\lambda(\theta_f, \theta_r) = \mathcal{L}_f(\theta_f) - \lambda \mathcal{L}_r(\theta_f, \theta_r),$$

- Build loss that encodes performance of classifier and an adversary
- Classifier penalized when adversary does well predicting $Z$
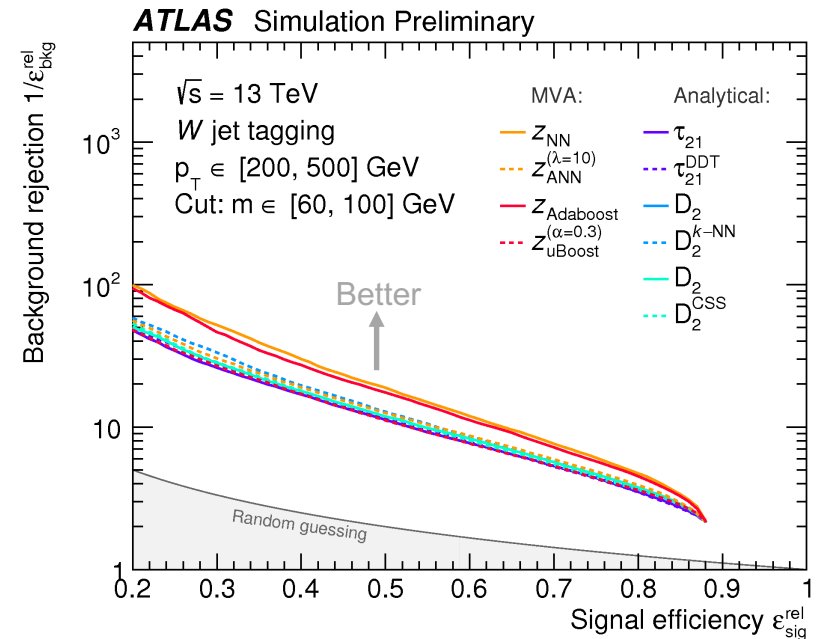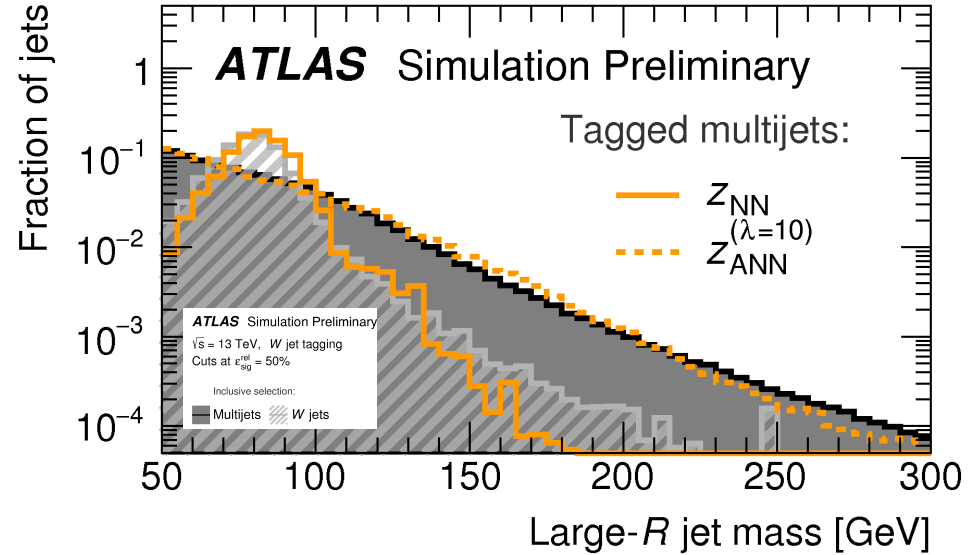- Training is a min-max game targeting saddle point solution
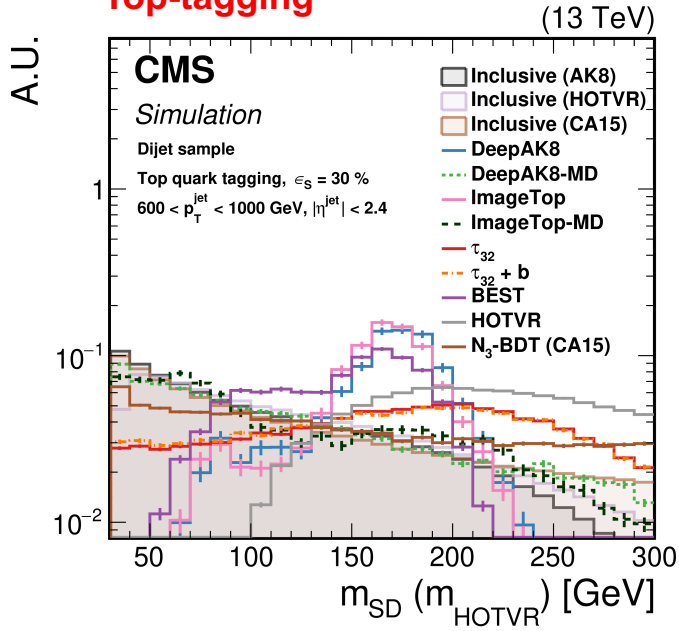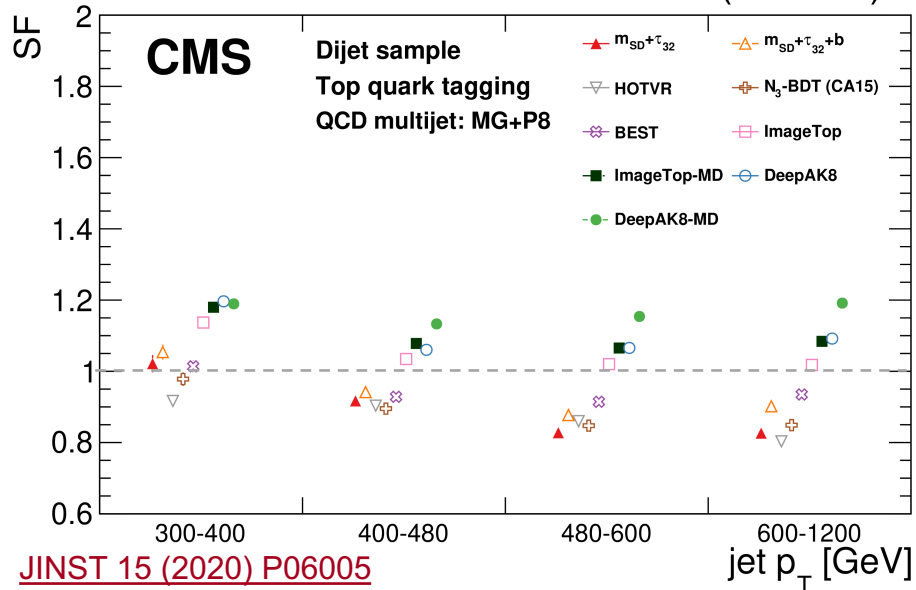
# Mass Decorrelation

**Top-tagging**

(13 TeV)



**W-tagging**

# Mass Decorrelation

**Top-tagging**

# Mitigating Data / MC Differences in LLP Jet Tagging
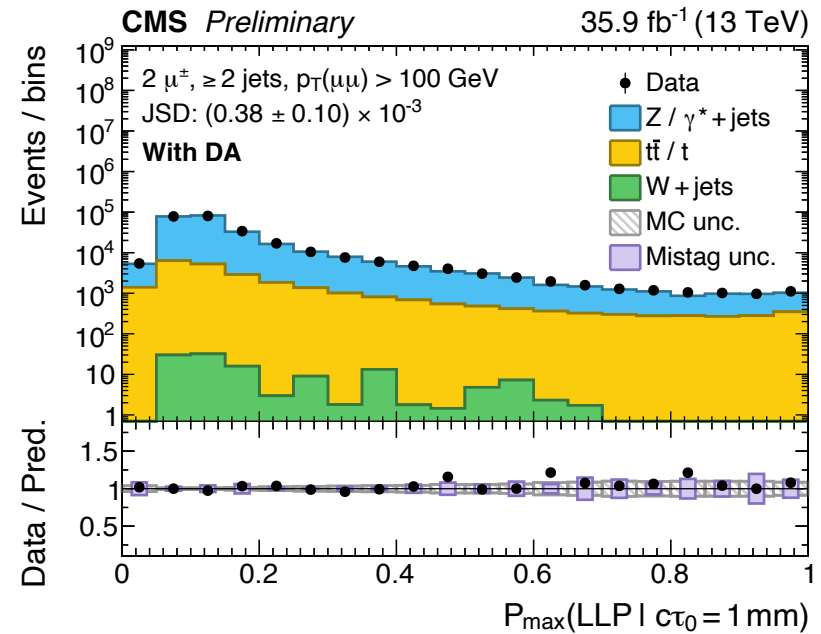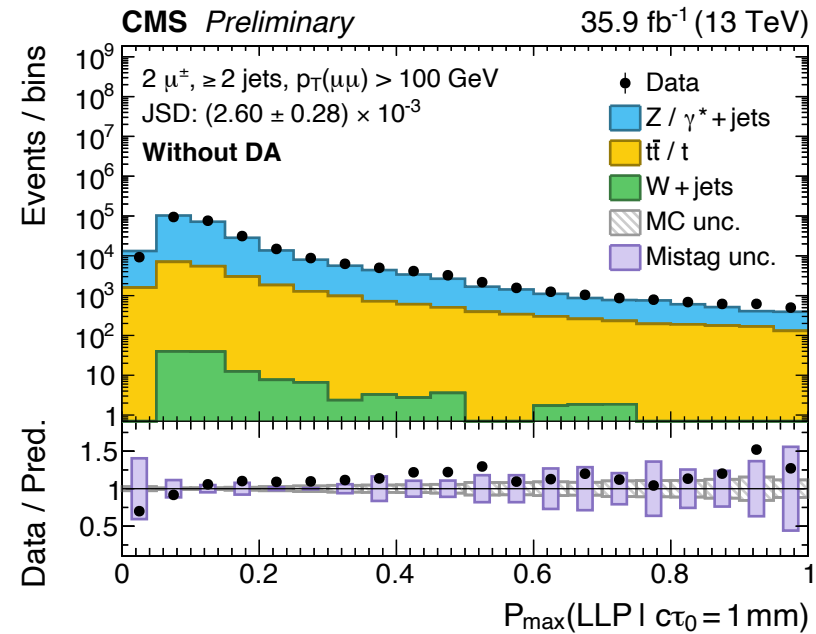
- Modified version of DeepJet

- Adversarial training to penalize differences in performance on MC vs Data

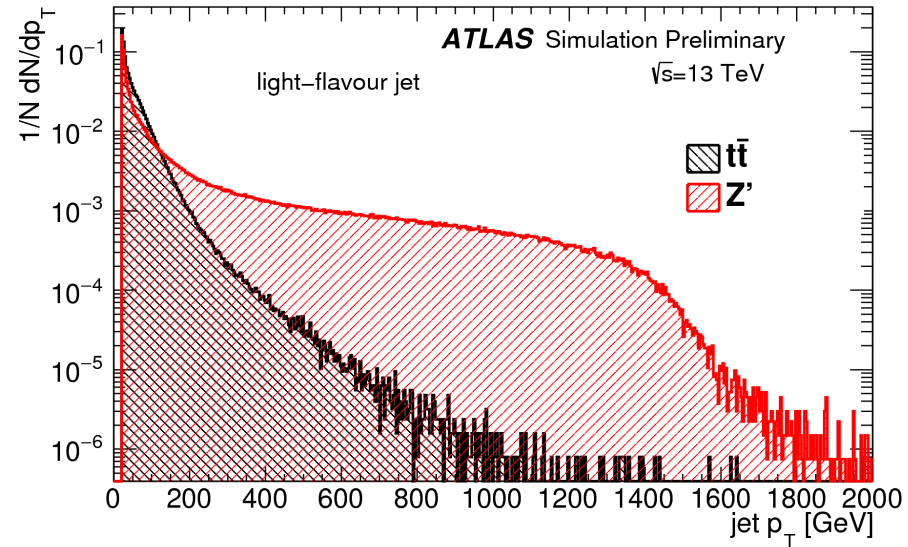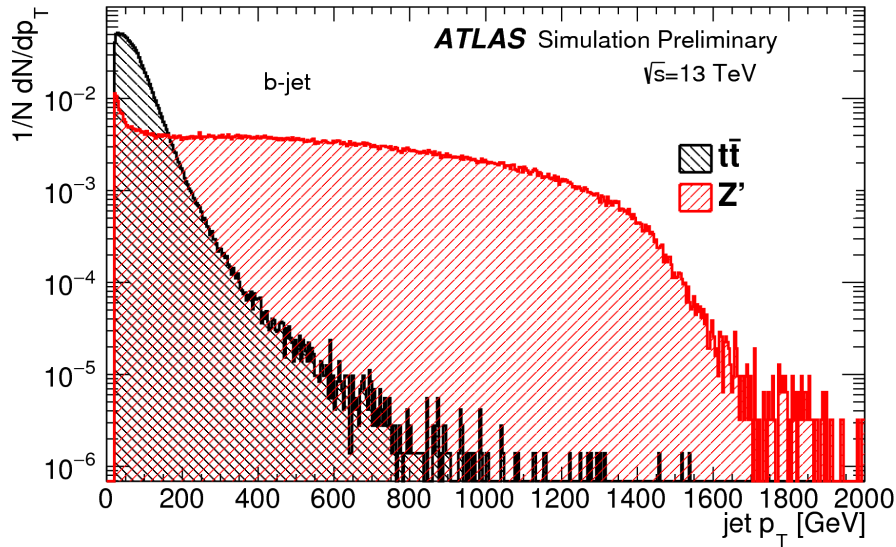# Mitigating Data / MC Differences in LLP Jet Tagging

# Conclusion

- Image and Sequence taggers deployed for Boosted jet tagging and b-tagging → show many of the expected performance gains

- Scale factors are reasonable
  - Mismodeling is not out of control
  - Interesting potential for mitigating Data /MC differences

- SF uncertainties worse in samples with more background / flavour fraction uncertainty
  - Must separate uncertainties from calibration method and from learning mismodeled features

- Intriguing questions open about learned representations and how they are expressed

# Backup

# Mitigating Sample Kinematic Bias in Training



- Want tagger to understand how features change with kinematics

- Don't want to be sensitive to training distribution of kinematics
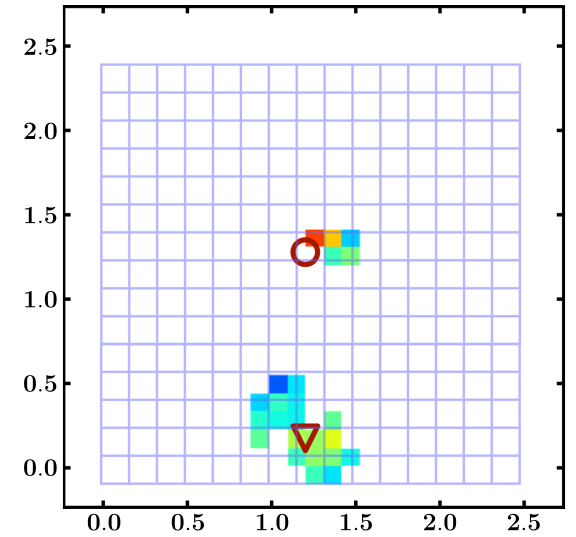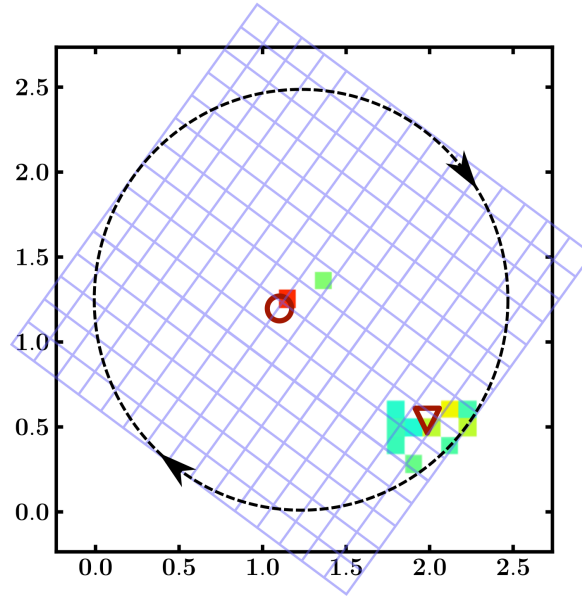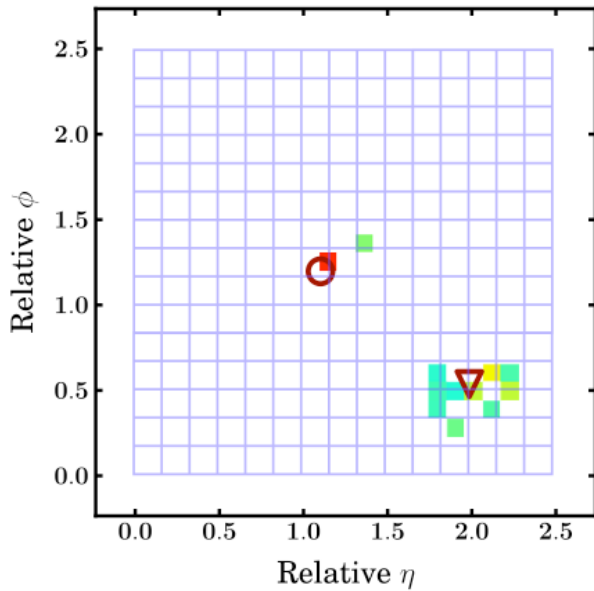  - $p_T$ is a pretty good discriminant! But distribution changes in analysis!

- Match key kinematic distributions between Signal / Background
  - Reweighting
  - Down sampling → ATLAS b-tag found this more stable for training

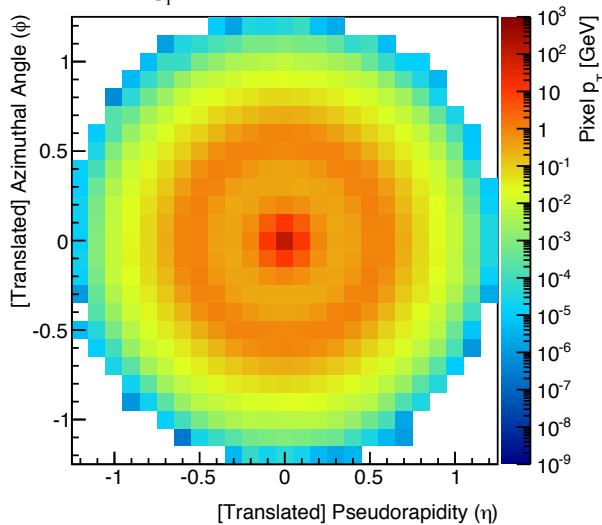# Combining Substructure Variables



arXiv:1808.07858

- Wide array of physics insight has gone into developing jet substructure observables

- Direct application of ML for combining power of multiple partially correlated substructure features

- First calibrations look quite reasonable!

# Jet Image Pre-Processing



Pythia 8, W' → WZ, √s = 13 TeV

$240 < p_T/GeV < 260$ GeV, $65 < mass/GeV < 95$

**Translate → Rotate → Flip**

[Translated] Pseudorapidity (η)

[Translated] Azimuthal Angle (φ)

Pixel $p_T$ [GeV]

Relative η

Relative φ

$Q_1$

$Q_2$

# W Tagging

CMS-DP-2020-002

EPJC 79 (2019) 375



NOTE: different $p_T$ ranges

# DeepAK8 variations



(13 TeV)

**CMS**

*Simulation*

**Top quark vs. QCD multijet**

$1000 < p_T^{gen} < 1500$ GeV, $|\eta^{gen}| < 2.4$

$105 < m_{SD}^{AK8} < 210$ GeV

Background efficiency

Signal efficiency

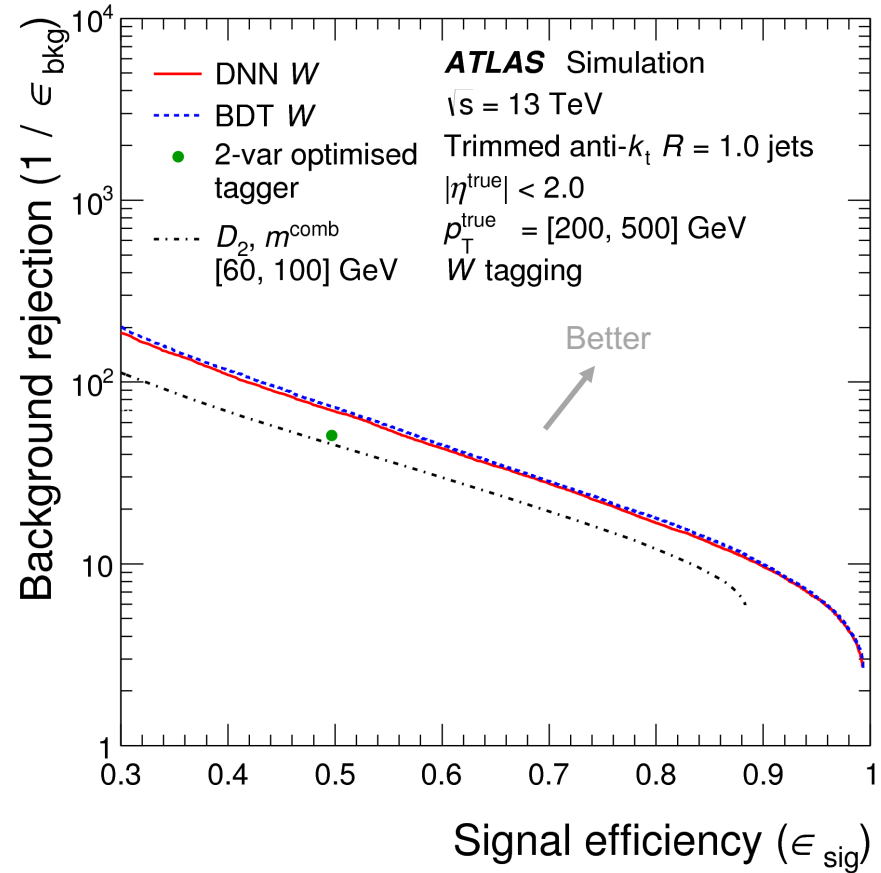- - - Particle (Kinematics)
— · — Particle (w/o Flavor)
——— Particle + SV (DeepAK8)

# DeepJet

**Number of particles/SV**

**LSTM layers**
**Builds a summary of the information extracted in each set of features**

**Feature extractor**
**convolution performed on each particle / SV [1x1]**

**Correlations and classification**

# DeepJet and Training Size

# Energy deposition in Si modeling

# Top Misidentification Scale Factors

# W-tagging Scale Factors

# Top Misidentification Scale Factors

# Adversarial Networks

$$\hat{\theta}_f, \hat{\theta}_r = \arg \min_{\theta_f} \max_{\theta_r} E(\theta_f, \theta_r).$$

$$E_\lambda(\theta_f, \theta_r) = \mathcal{L}_f(\theta_f) - \lambda \mathcal{L}_r(\theta_f, \theta_r),$$

- Loss encodes performance of classifier and adversary

  - Classifier penalized when adversary does well at predicting Z

- Hyper-parameter λ controls trade-off

  - Large λ enforces f(…) to be pivotal, e.g. robust to nuisance
  - Small λ allows f(…) to be more optimal

# Learning to Pivot: Toy Example

2D example

$$x \sim \mathcal{N}\left((0,0), \begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix}\right) \quad \text{when } Y = 0,$$

$$x \sim \mathcal{N}\left((1, 1+Z), \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}\right) \quad \text{when } Y = 1.$$



- **Without adversary (top) large variations in network output with nuisance parameter**

- **With adversary (bottom) performance is independent!**

# Modeling Comparisons

$C_\mu \to qq$ vs $H \to gg$

$C_\mu \to qq$ vs $H \to qq$

# Next Generation Taggers

# Next Generation: DIPS – Deep Impact Parameter Sets

- **Challenges of RNN Tagging**
  - Must choose sequence ordering, not inherent, which is best?
  - Requires iteration over tracks, can't be parallelized

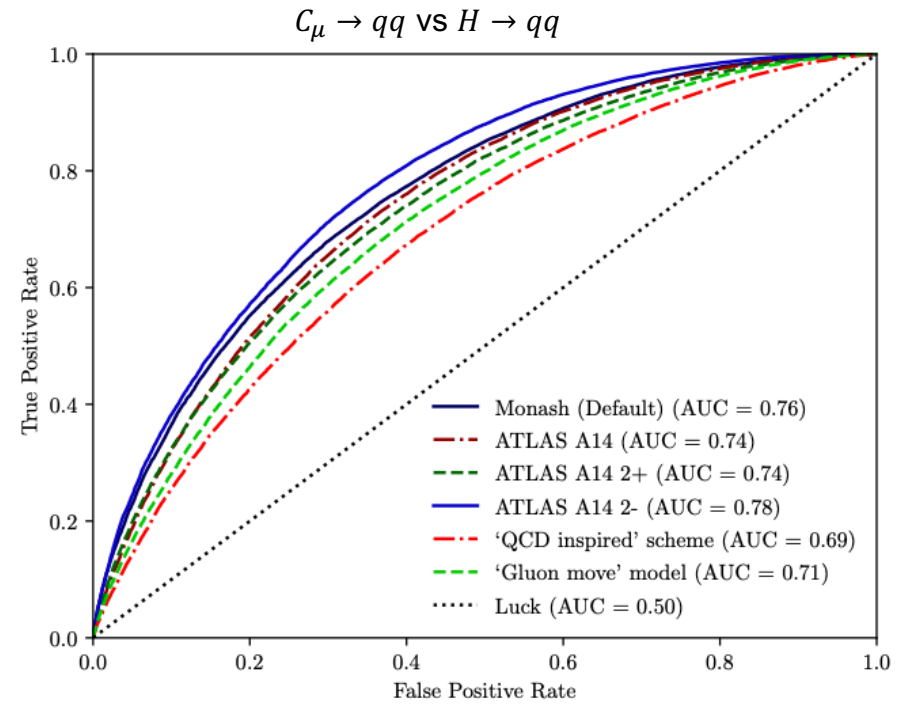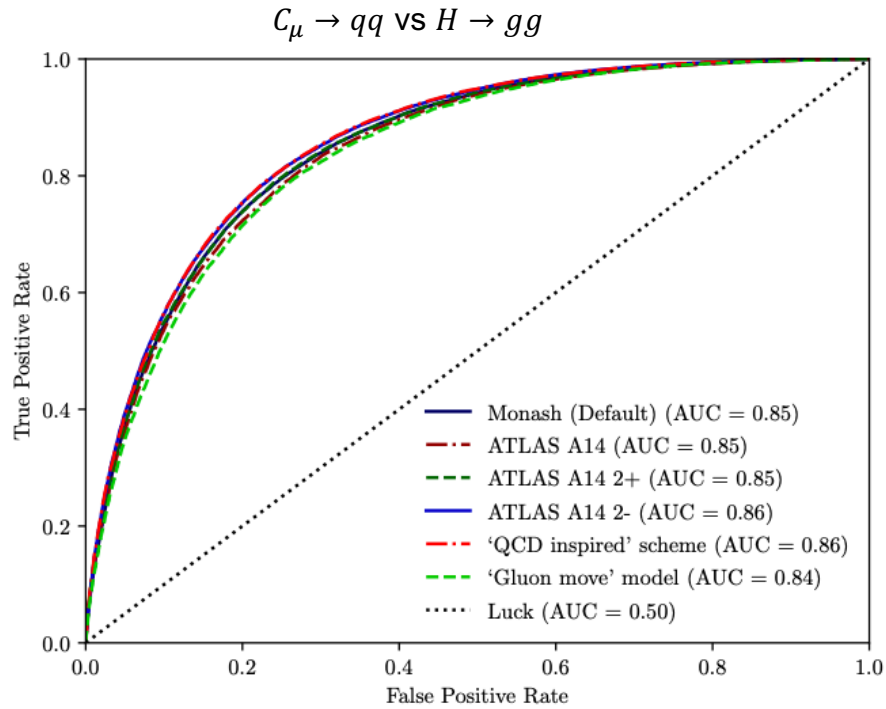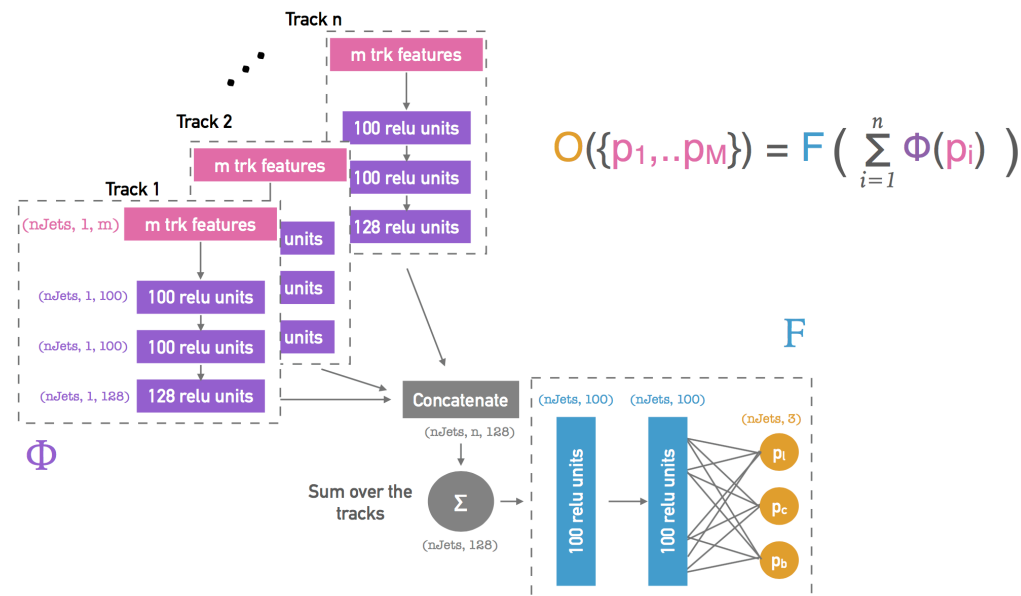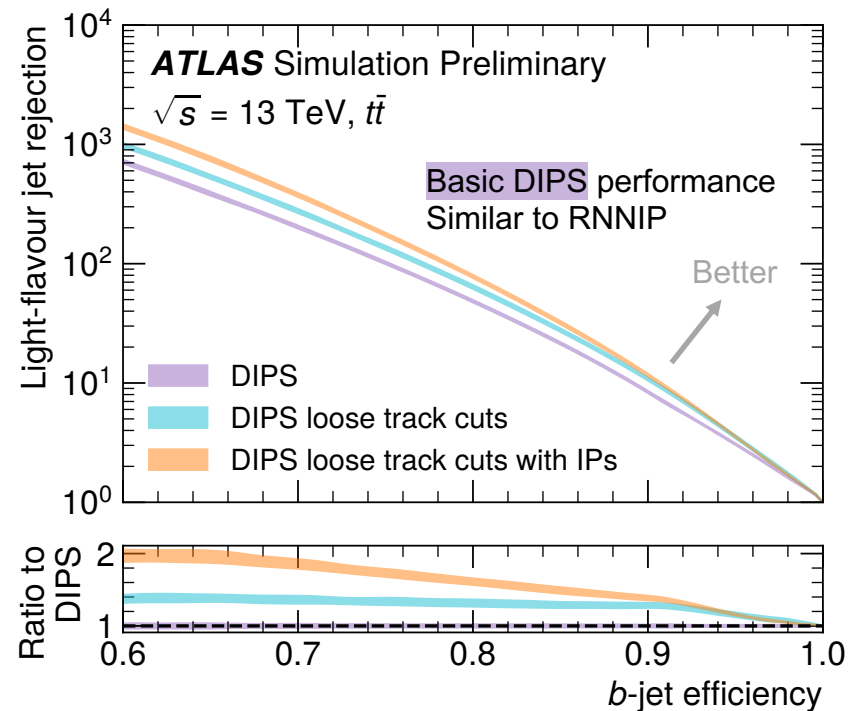- **Deep Sets: permutation invariant and parallelizable model**



$$O(\{p_1, .. p_M\}) = F \left( \sum_{i=1}^{n} \Phi(p_i) \right)$$

| Model | Parameters | Training time [min] | Time / epoch [s] |
|-------|-----------|---------------------|------------------|
| RNNIP | 47k | 86 ± 13 | 241 ± 14 |
| DIPS | 49k | 44 ± 4 | 78 ± 4 |

| Model | Parameters | GPU Evaluation time [s] | CPU evaluation time [s] |
|-------|-----------|-------------------------|-------------------------|
| RNNIP | 47k | 170 ± 2 | 685 ± 84 |
| DIPS | 49k | 46 ± 2 | 206 ± 98 |



**ATLAS** Simulation Preliminary
$\sqrt{s}$ = 13 TeV, $t\bar{t}$

Basic DIPS performance Similar to RNNIP

Better

- DIPS
- DIPS loose track cuts
- DIPS loose track cuts with IPs

Light-flavour jet rejection

Ratio to DIPS

$b$-jet efficiency

ATL-PHYS-PUB-2020-014     JHEP 01 (2019) 121
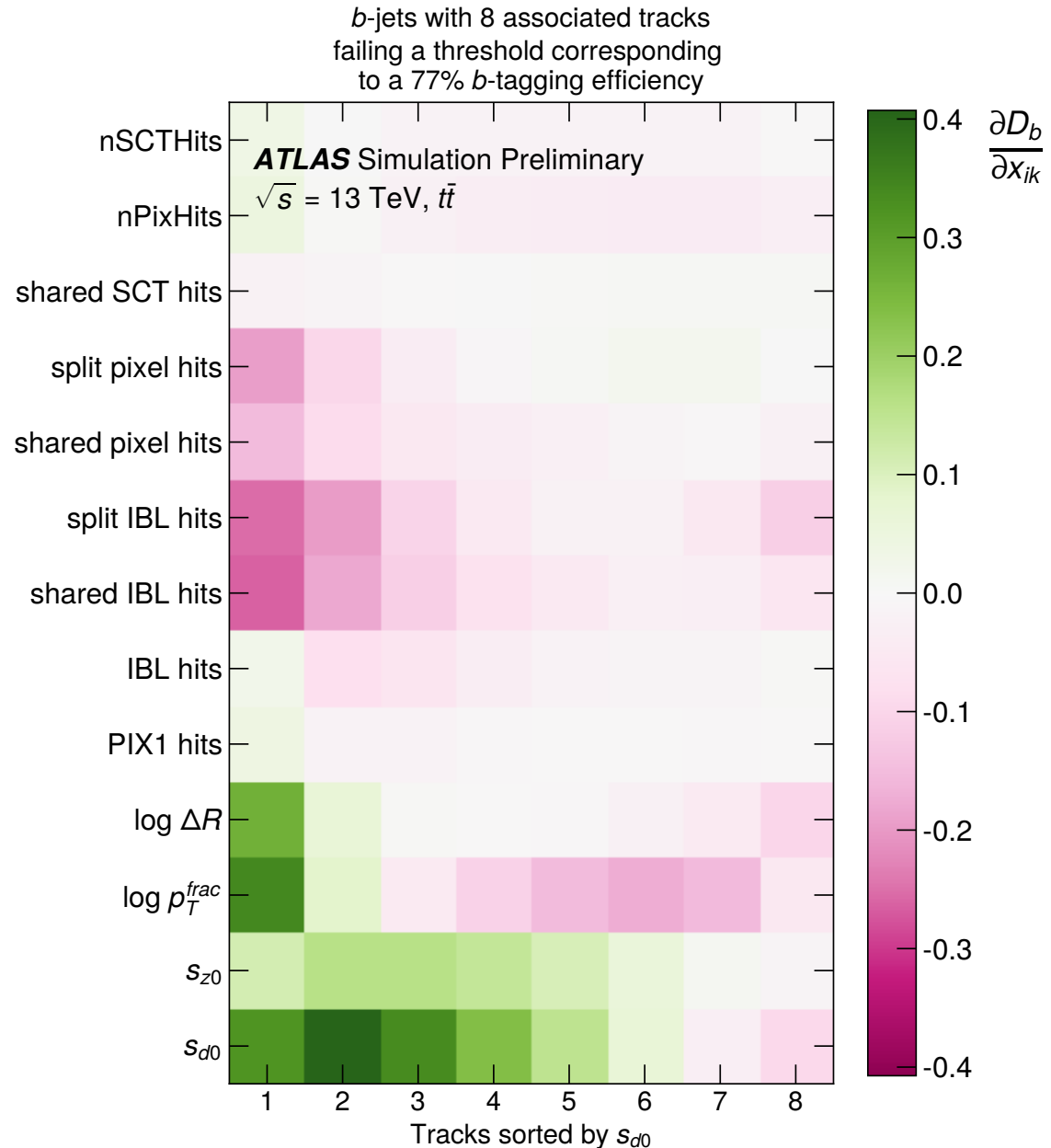
# Sanity Checks

Salience Map

$$S_{jk} = \sum_{i=1}^{N_{jets}} \frac{dD_b(x^{(i)})}{dx_{jk}^{(i)}}$$
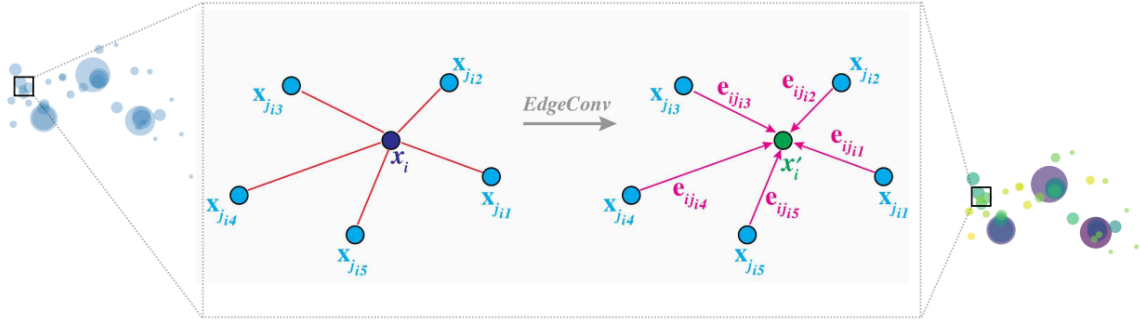
$x^{(i)}$ = all tracks/features of $i^{th}$ jet
i = jet
j = feature
k = track



b-jets with 8 associated tracks failing a threshold corresponding to a 77% b-tagging efficiency

ATL-PHYS-PUB-2020-014

# Next Generation: ParticleNet with Graph Neural Networks



Phys. Rev. D 101, 056019 (2020)
CMS-DP-2020-002