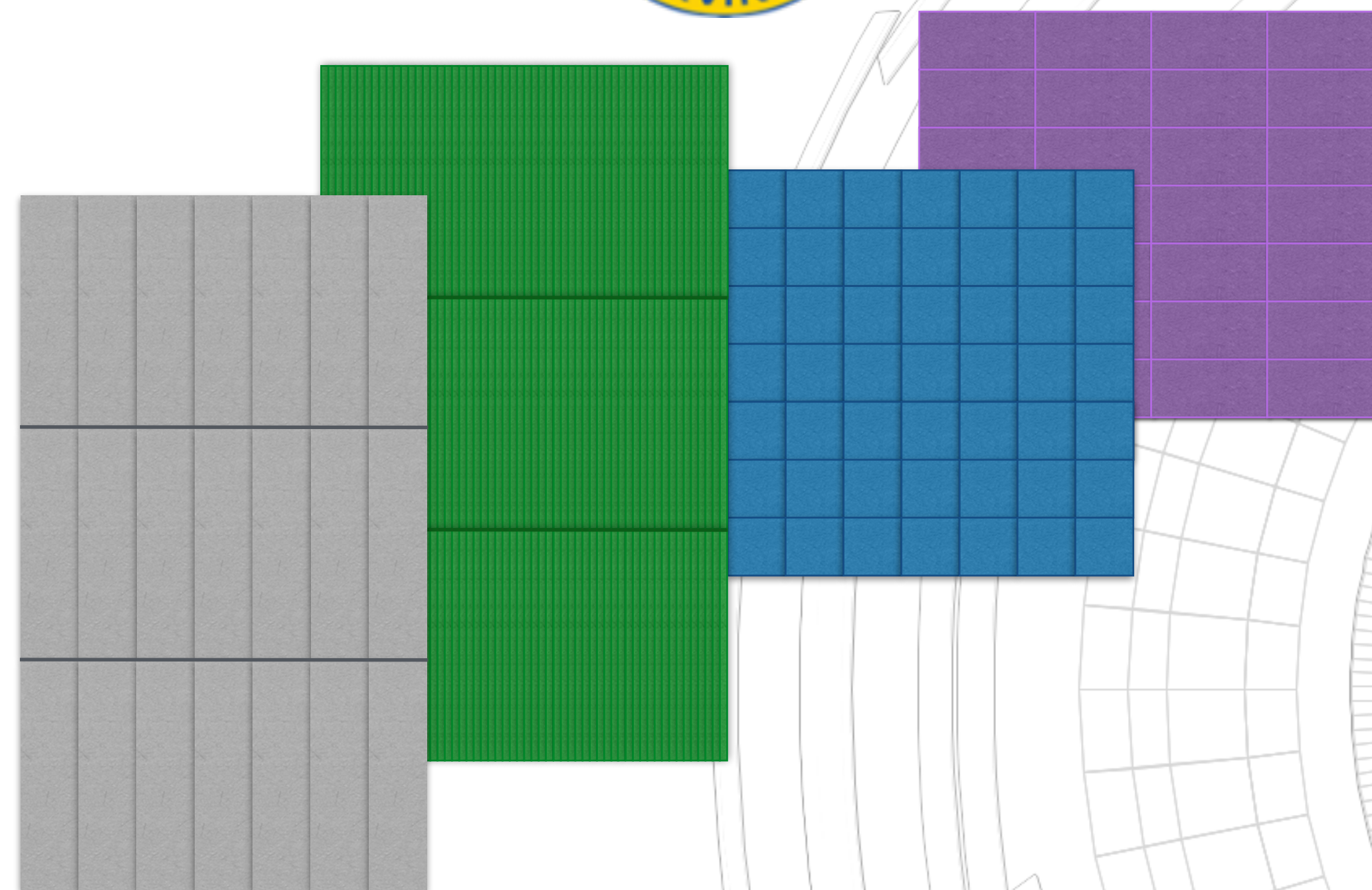# Uncertainties in the era of ML
# (For Particle and Astrophysics)

Aishik Ghosh

Aspen Centre for Physics
28 Mar 2023

# Uncertainties, the bedrock of experimental science

# Uncertainties, the bedrock of experimental science

mH = 125.25 ± 0.17 GeV

# Uncertainties, the bedrock of experimental science

mH = 125.25 ± 0.17 GeV

# Uncertainties, the bedrock of experimental science

$$mH = 125.25 \pm 0.17 \text{ GeV}$$



How sure am I ? How can I reduce my uncertainty ?

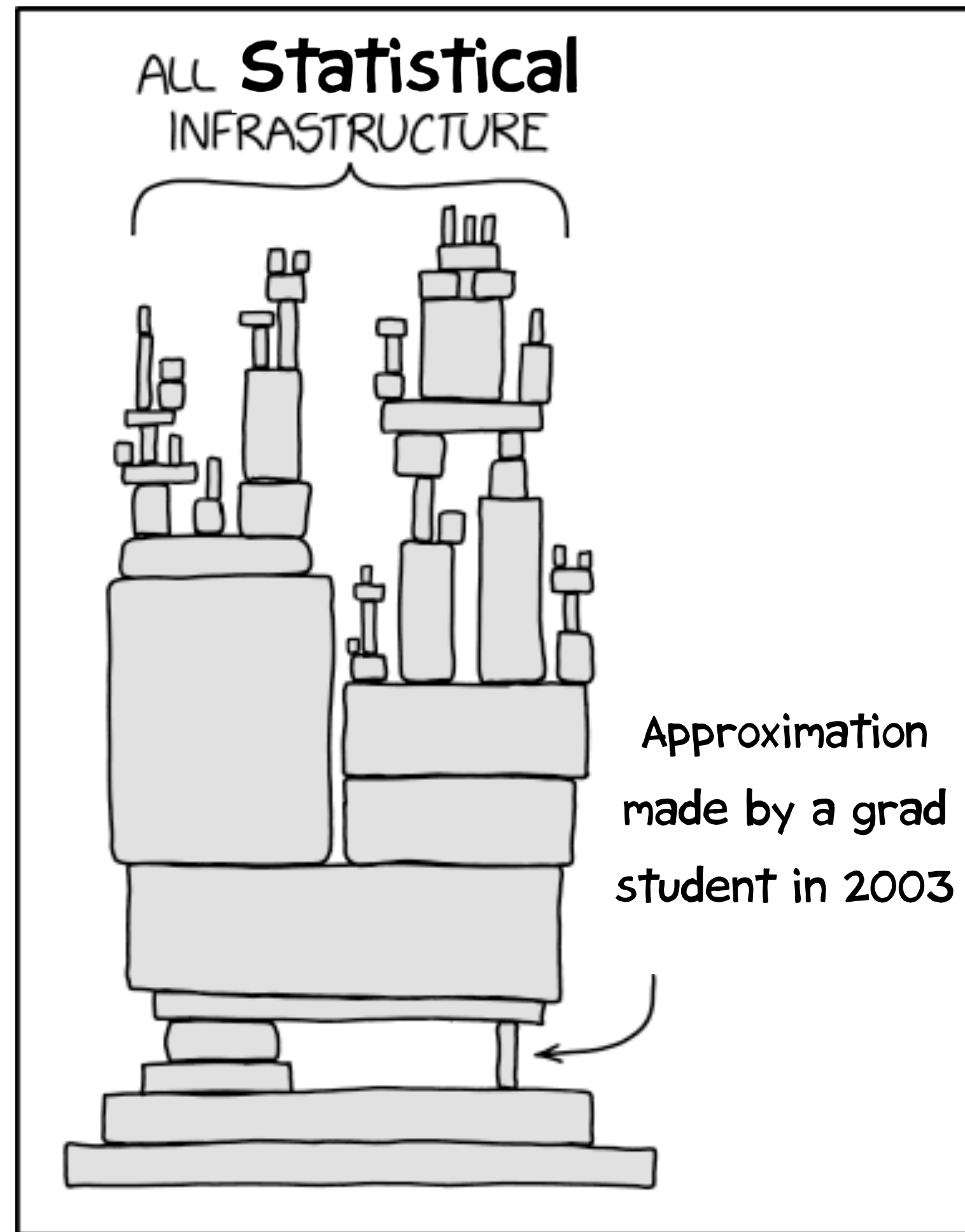# Uncertainties, the bedrock of experimental science

mH = 125.25 ± 0.17 GeV

{statistical, detector systematic, theory systematic, epistemic, ….}



How sure am I ? How can I reduce my uncertainty ?

ALL **Statistical** INFRASTRUCTURE

Approximation made by a grad student in 2003
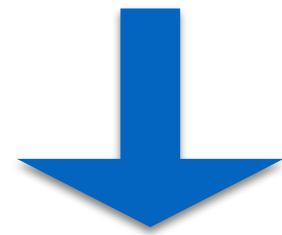
Time to re-examine some of the underlying pieces

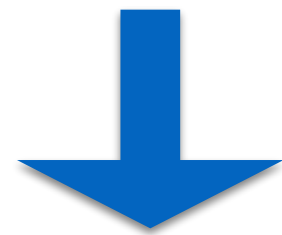Are they up to the task of the precision era?

From Daniel Whiteson
Inspired by XKCD

# Outline: A predictable evolution over ten years

**Fear**: Will ML exacerbate uncertainties in a way human-designed strategies naturally avoid ?

**Solution**: Find ML equivalents of uncertainty mitigation tricks we implicitly use in classical methods. Understand good and bad ways to use ML
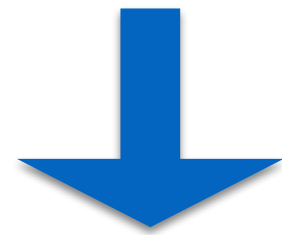
**Opportunity**: ML *for* uncertainty –  Realising that ML unlocks completely new methods to tackle uncertainties in a way classical methods couldn't
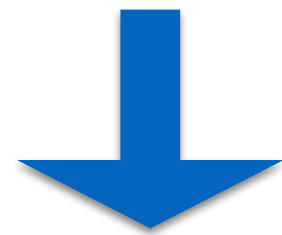
**Revolution**: Novel ML uncertainty quantification & mitigation methods have wider applications, also back-ported to traditional algorithms

# Outline: A predictable evolution over ten years

**Fear**: Will ML exacerbate uncertainties in a way human-designed strategies naturally avoid ?

**Solution**: Find ML equivalents of uncertainty mitigation tricks we implicitly use in classical methods. Understand good and bad ways to use ML

We are
here

**Opportunity**: ML *for* uncertainty –  Realising that ML unlocks completely new methods to tackle uncertainties in a way classical methods couldn't
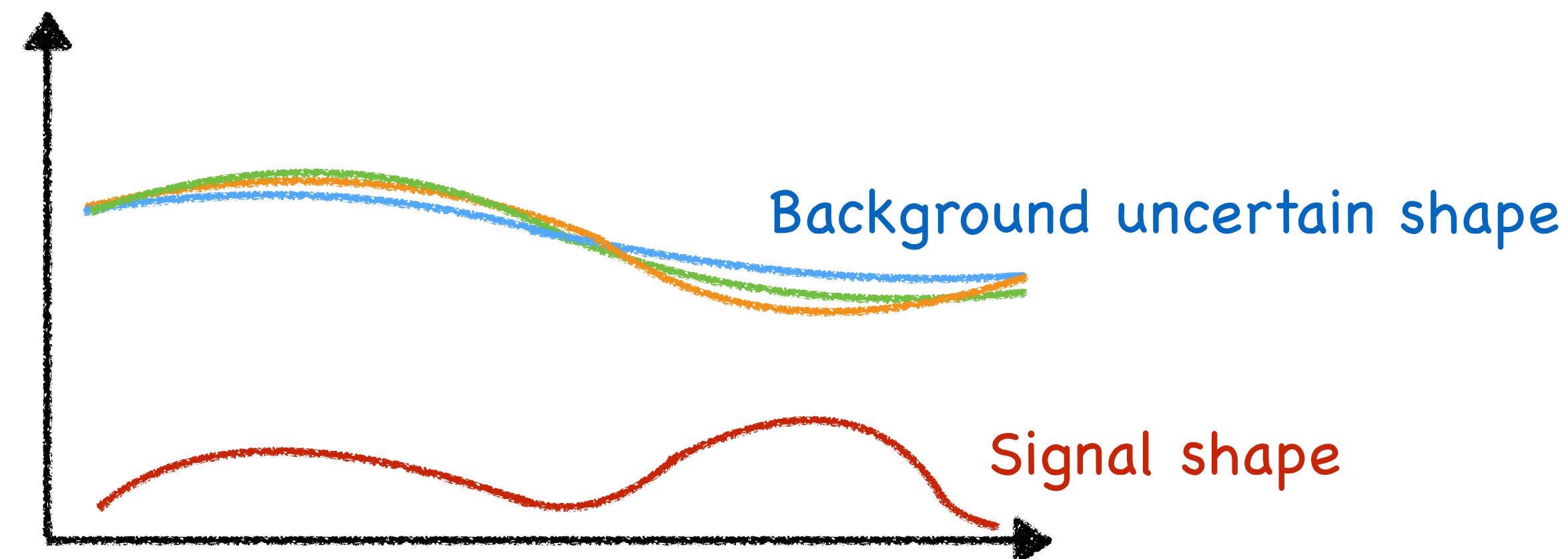
**Revolution**: Novel ML uncertainty quantification & mitigation methods have wider applications, also back-ported to traditional algorithms
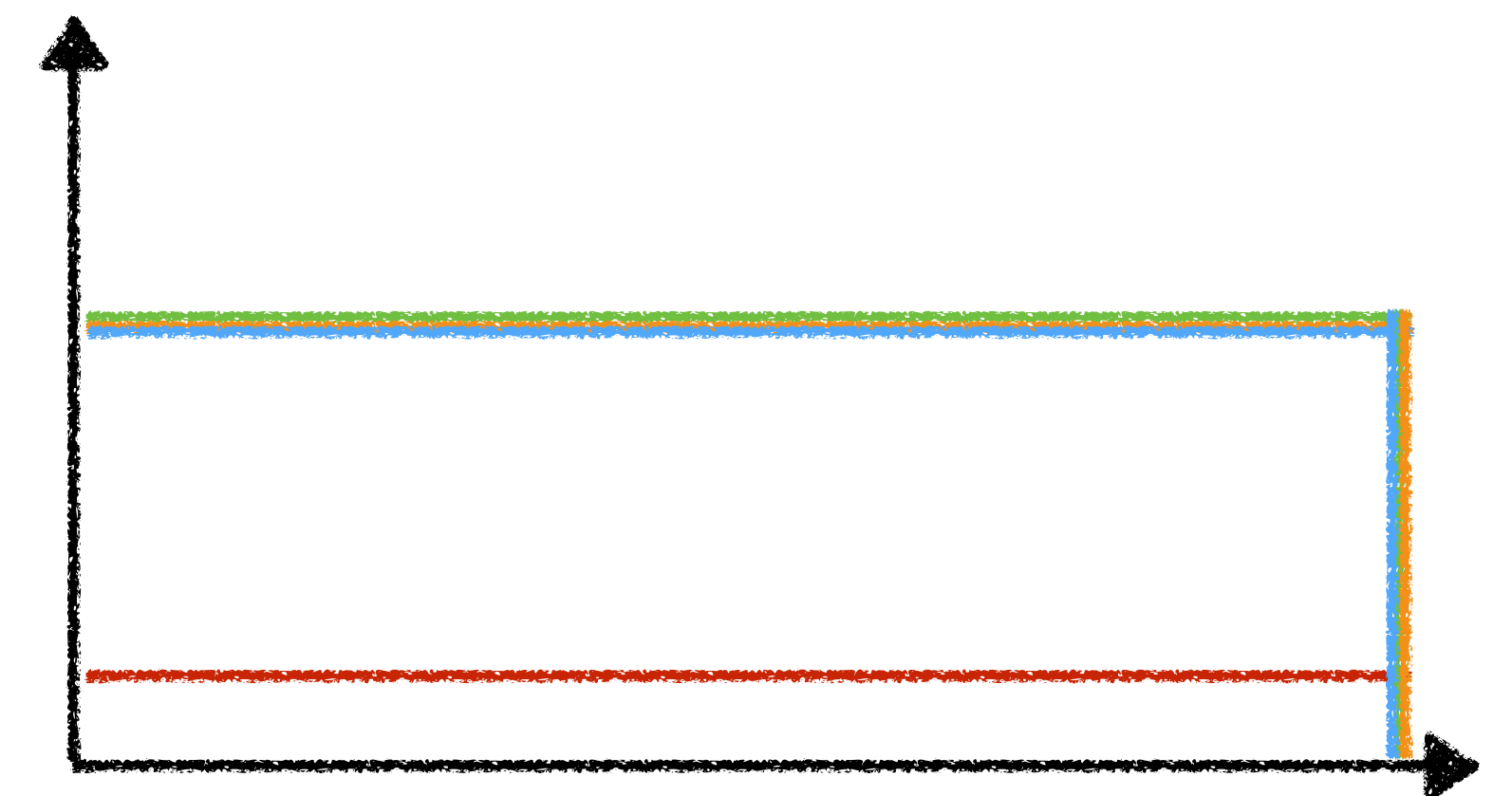
# Observable Sensitive to Nuisance Parameters

Traditionally, we reduce impact of NP by sacrificing something:

- Don't use observable

- Don't use phase space which is badly modelled by simulation

- Reduce sensitivity some other way

Infinite bin analysis, very sensitive to shape uncertainty



Background uncertain shape

Signal shape

Single bin analysis, insensitive to shape uncertainty

# ML equivalent problem: Domain Adaptation

MNIST

SOURCE

TARGET

MNIST-M

Learning to Pivot, Louppe et al.

Similar ideas: Blance et al., Stevens et al., Wunsch at al., Estrade at al. Kasieczka at al.

# Adversarial decorrelation

## S vs B

### Classifier $f$

## Regress NP

### Adversary $r$



NN output

$f(X; \theta_f)$

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

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

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

$\dots$

$Z$

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

$\theta_f$

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

$\theta_r$

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

Learning to Pivot, Louppe et al.

$$L_{Classifier} = L_{Classification} - \lambda \cdot L_{Adversary}$$

7

Learning to Pivot, Louppe et al.

Similar ideas: Blance et al., Stevens et al., Wunsch at al., Estrade at al. Kasieczka at al.
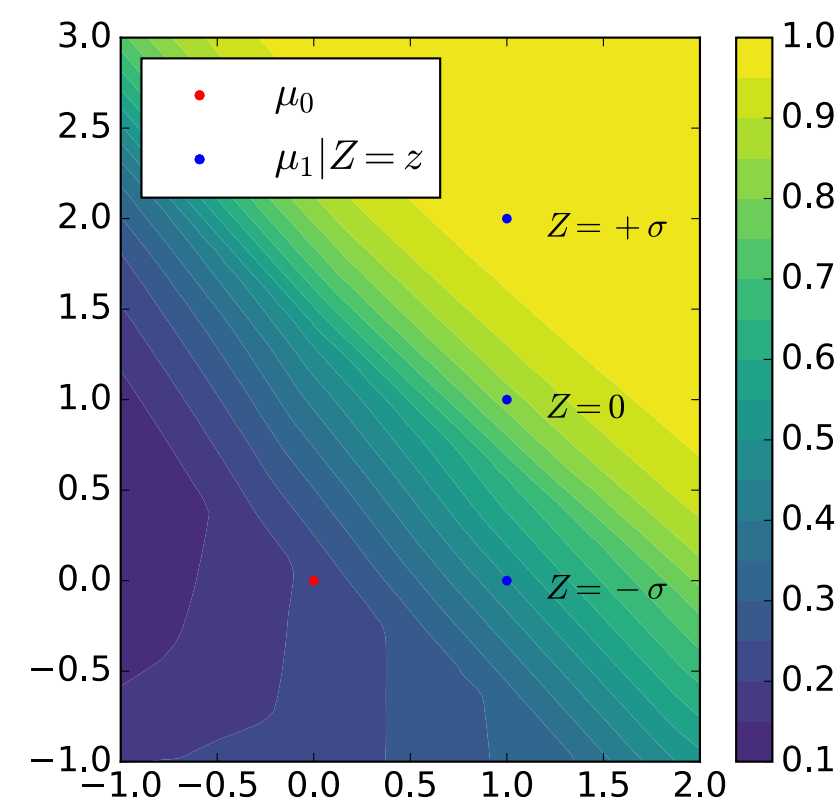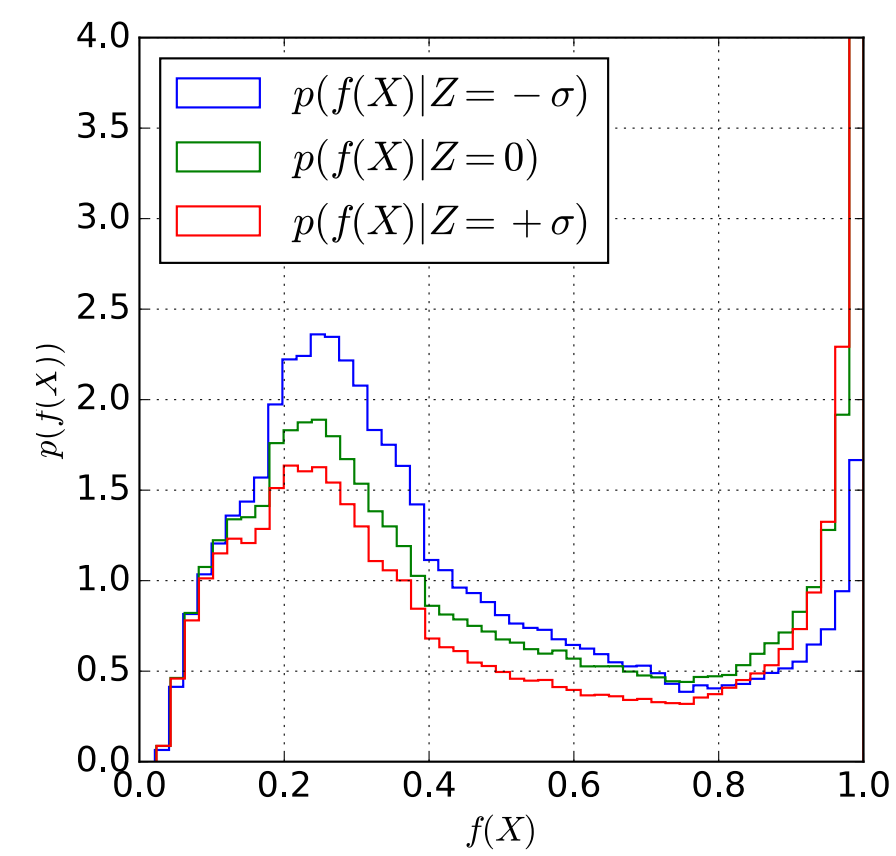
# Adversarial decorrelation

S vs B

Regress NP



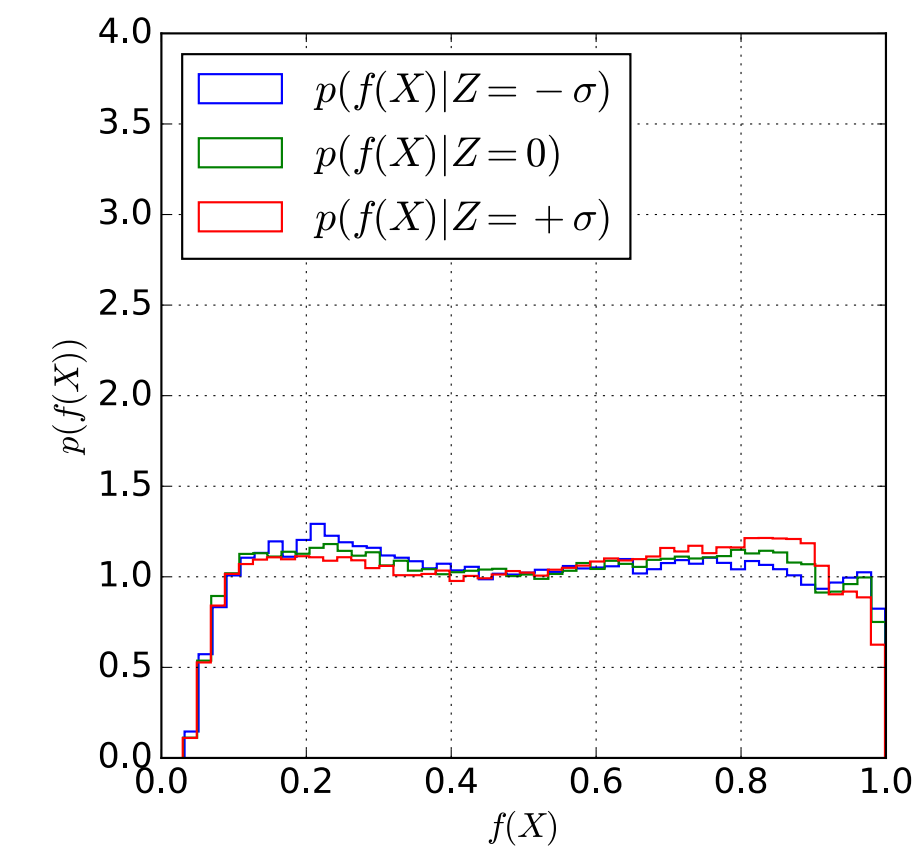To fool the adversary, classifier output should be decorrelated to Z

Learning to Pivot, Louppe et al.

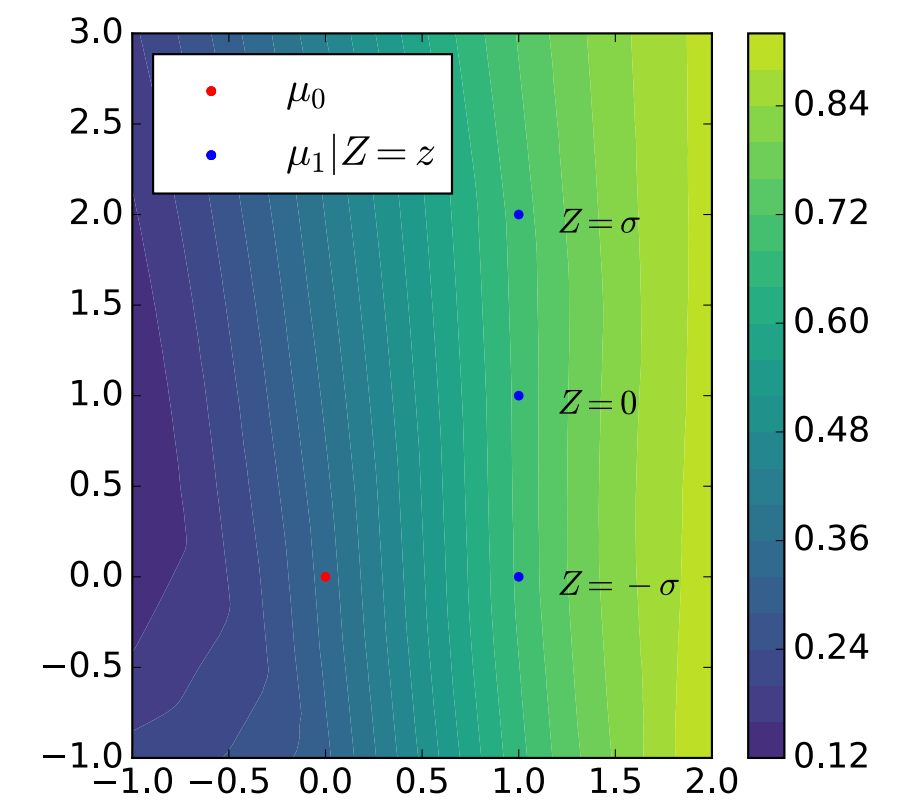$$L_{Classifier} = L_{Classification} - \lambda \cdot L_{Adversary}$$
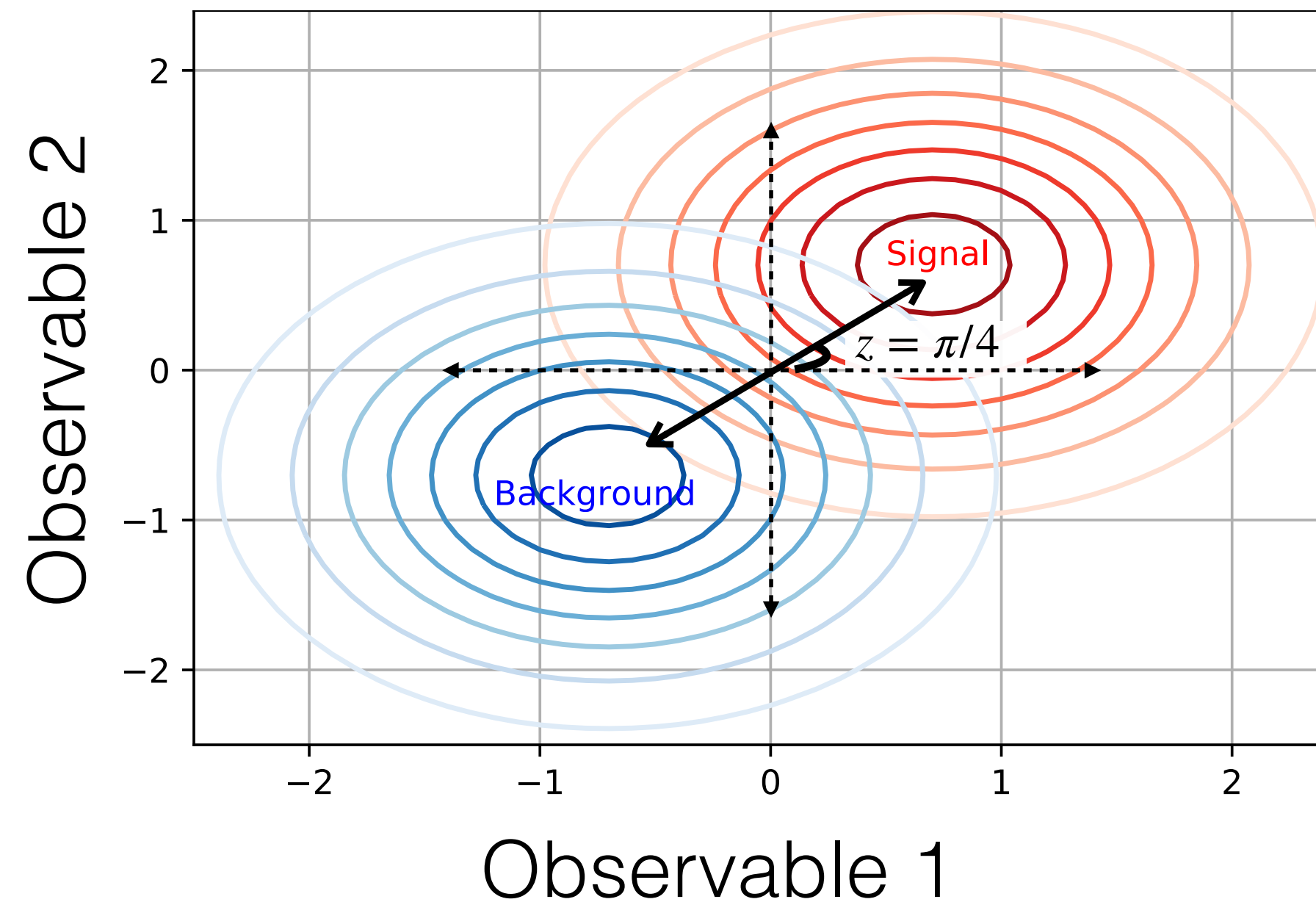
# ML-Decorrelation Methods



Adversarial Decorrelation

Classifier output for various values of Z

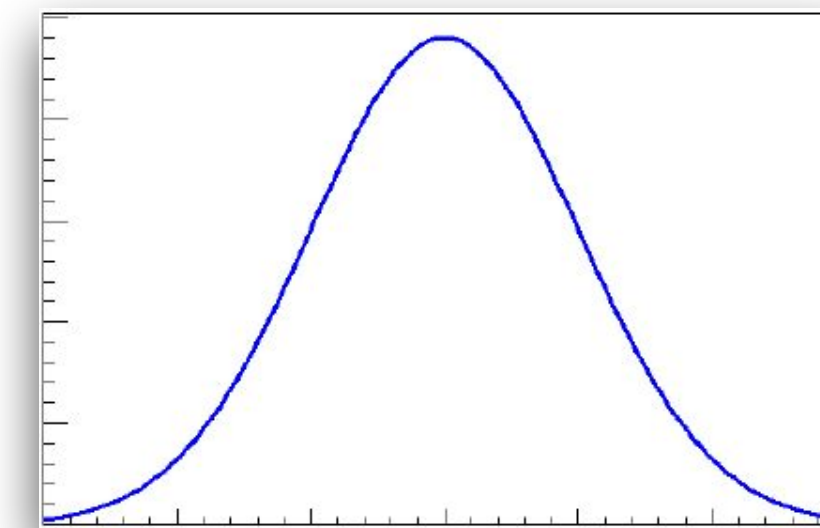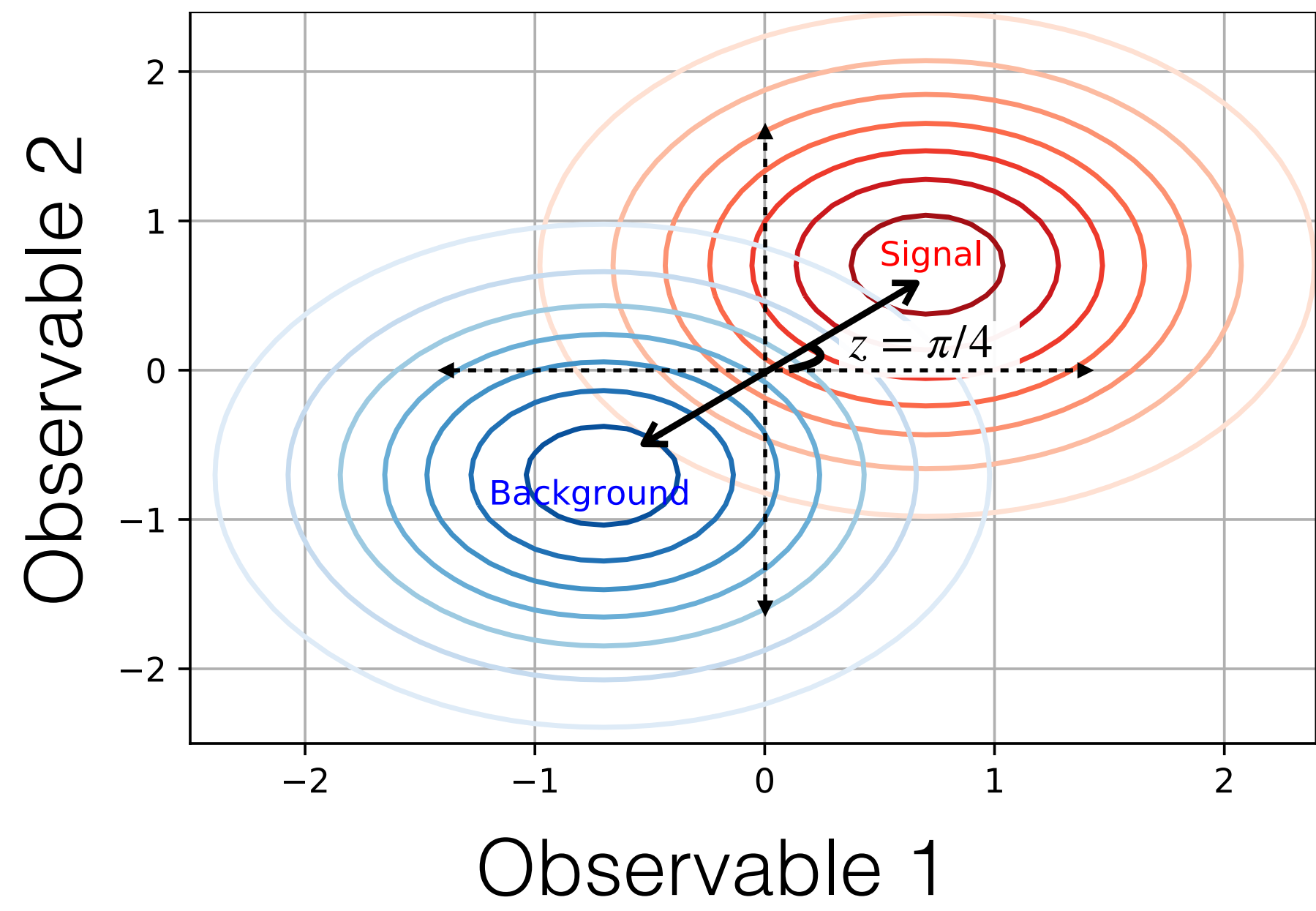Learning to Pivot, Louppe et al.

# What if we could do better ?

# What if we could do better ?



$z$ = Nuisance Parameter

# What if we could do better ?



$$n_i | \mu \cdot S_i(\boldsymbol{\theta}) + B_i(\boldsymbol{\theta})) \times \prod_{j \in syst} \mathcal{G}(\theta_j^0 | \theta_j, \Delta\theta_j)$$

$z$ = Nuisance Parameter
Prior

# What if we could do better ?

16°

$$n_i | \mu \cdot S_i(\boldsymbol{\theta}) + B_i(\boldsymbol{\theta})) \times \prod_{j \in syst} \mathcal{G}(\theta_j^0 | \theta_j, \Delta \theta_j)$$

Signal

$z$ = Nuisance Parameter
Prior

0.16°

16°

Likelihood

Temp

Observable 1

0.16°

.6°

# Opposite of decorrelation: Uncertainty-aware learning

- Propagate uncertainties through the classifier in an "uncertainty aware" way



$f(x_1, x_2, \ldots, z)$

Similar to 1601.07913

Repeat for each hypothesis $z$

Data
with Z = ?

# Opposite of decorrelation: Uncertainty-aware learning

- <u>Propagate uncertainties</u> through the classifier in an "uncertainty aware" way



$f(x_1, x_2, \ldots, z)$

Similar to 1601.07913

Repeat for each hypothesis $z$

Data with Z = ?

- Intuition: Allow the analysis technique to vary with Z
  You always get the best classifier for each value of Z

$$\mathcal{P}(n_i | \mu \cdot S_i(\boldsymbol{\theta}) + B_i(\boldsymbol{\theta})) \times \prod_{j \in syst} \mathcal{G}(\theta_j^0 | \theta_j, \Delta\theta_j) \cdot$$

# Better final measurements!



Narrower ⇒ Smaller [statistical + systematic] uncertainty on measurement

Practical for LHC analysis: Parameterise your main nuisance parameter but no need to train on all 100 NPs

An application in astrophysics

Simulated X-Ray Spectra

Neutron Star Mass and Radius

Equation of State (EOS)

Spectra + NP to (M,R) Regression

(M,R) to EOS Regression

Spectra + NP to EOS Regression

SOTA made a single point estimate + assumed uncorrelated Gaussian uncertainties

Real uncertainties look quite different

# Learn forward process to access the likelihood



Deploy with ONNX Runtime to compute likelihoods on-the-fly

# Forward process step-by-step

Intermediate steps remain interpretable physical quantities

Intermediate steps remain interpretable physical quantities



Learn EOS to M-R

Intermediate steps remain interpretable physical quantities



Learn EOS to M-R

Learn {M,R,NPs} to Spectrum

# Forward process step-by-step

Intermediate steps remain interpretable physical quantities

Learn EOS to M-R

Learn {M,R,NPs} to Spectrum

Nuisance Priors:

M-R likelihoods:

True:

Tight:

Loose:

Back to particle physics

Back to particle physics

We also learnt what not to do …

# ML-decorrelating theory uncertainties

Default

What you want with decorrelation

What you get with decorrelation

Instruction to ML: "Please shrink Pythia vs Herwig difference"

# ML-decorrelating theory uncertainties

Default

What you want with decorrelation

What you get with decorrelation

Instruction to ML: "Please shrink Pythia vs Herwig difference"

Model will learn to fool you !

ML methods don't often generalise the way you would hope

# Theory uncertainties



It's dangerous to use ML methods to mitigate theory uncertainties

But we continue to treat $\Delta_{theory}$ and $\Delta_{exp}$ on same footing in statistical fits

What even is their statistical behaviour?

From Daniel Whiteson

# Scale Uncertainties

Uncertainty of cross-section from truncating QFT series

Sensitivity to scale variation quantifies 'uncertainty'

# Scale Uncertainties

Up: $\mu_+ = 2\,\mu_0$

$$\mu_0 = \frac{H_T}{2} = \frac{1}{2} \sum_{final\ state} \sqrt{m^2 + p_T^2}$$

Uncertainty of cross-section from truncating QFT series

Sensitivity to scale variation quantifies 'uncertainty'

Down: $\mu_- = \frac{1}{2}\,\mu_0$

# Questions

- How accurate are these scale uncertainties ?

- Is 1/2 to 2 a good range ?

**Study pull distribution**

$$t_{scale} = \frac{\sigma_{NLO} - \sigma_{LO}}{\Delta\sigma_{LO\ scale}}$$

# Questions

- How accurate are these scale uncertainties ?

- Is 1/2 to 2 a good range ?

## Madgraph paper

**The automated computation of tree-level and next-to-leading order differential cross sections, and their matching to parton shower simulations**

J. Alwall[a], R. Frederix[b], S. Frixione[b], V. Hirschi[c], F. Maltoni[d], O. Mattelaer[d], H.-S. Shao[e], T. Stelzer[f], P. Torrielli[g], M. Zaro[hi]

**Study pull distribution**

$$t_{scale} = \frac{\sigma_{NLO} - \sigma_{LO}}{\Delta\sigma_{LO\ scale}}$$

| Process | Syntax | Cross section (pb) | | |
|---|---|---|---|---|
| Vector boson +jets | | LO 13 TeV | | NLO 13 TeV |
| a.1 $pp \to W^{\pm}$ | p p > wpm | $1.375 \pm 0.002 \cdot 10^5$ $^{+15.4\%\ +2.0\%}_{-16.6\%\ -1.6\%}$ | | $1.773 \pm 0.007 \cdot 10^5$ $^{+5.2\%\ +1.9\%}_{-9.4\%\ -1.6\%}$ |
| a.2 $pp \to W^{\pm}j$ | p p > wpm j | $2.045 \pm 0.001 \cdot 10^4$ $^{+19.7\%\ +1.4\%}_{-17.2\%\ -1.1\%}$ | | $2.843 \pm 0.010 \cdot 10^4$ $^{+5.9\%\ +1.3\%}_{-8.0\%\ -1.1\%}$ |
| a.3 $pp \to W^{\pm}jj$ | p p > wpm j j | $6.805 \pm 0.015 \cdot 10^3$ $^{+24.5\%\ +0.8\%}_{-18.6\%\ -0.7\%}$ | | $7.786 \pm 0.030 \cdot 10^3$ $^{+2.4\%\ +0.9\%}_{-6.0\%\ -0.8\%}$ |
| a.4 $pp \to W^{\pm}jjj$ | p p > wpm j j j | $1.821 \pm 0.002 \cdot 10^3$ $^{+41.0\%\ +0.5\%}_{-27.1\%\ -0.5\%}$ | | $2.005 \pm 0.008 \cdot 10^3$ $^{+0.9\%\ +0.6\%}_{-6.7\%\ -0.5\%}$ |
| a.5 $pp \to Z$ | p p > z | $4.248 \pm 0.005 \cdot 10^4$ $^{+14.6\%\ +2.0\%}_{-15.8\%\ -1.6\%}$ | | $5.410 \pm 0.022 \cdot 10^4$ $^{+4.6\%\ +1.9\%}_{-8.6\%\ -1.5\%}$ |
| a.6 $pp \to Zj$ | p p > z j | $7.209 \pm 0.005 \cdot 10^3$ $^{+19.3\%\ +1.2\%}_{-17.0\%\ -1.0\%}$ | | $9.742 \pm 0.035 \cdot 10^3$ $^{+5.8\%\ +1.2\%}_{-7.8\%\ -1.0\%}$ |
| a.7 $pp \to Zjj$ | p p > z j j | $2.348 \pm 0.006 \cdot 10^3$ $^{+24.3\%\ +0.6\%}_{-18.5\%\ -0.6\%}$ | | $2.665 \pm 0.010 \cdot 10^3$ $^{+2.5\%\ +0.7\%}_{-6.0\%\ -0.7\%}$ |
| a.8 $pp \to Zjjj$ | p p > z j j j | $6.314 \pm 0.008 \cdot 10^2$ $^{+40.8\%\ +0.5\%}_{-27.0\%\ -0.5\%}$ | | $6.996 \pm 0.028 \cdot 10^2$ $^{+1.1\%\ +0.5\%}_{-6.8\%\ -0.5\%}$ |
| a.9 $pp \to \gamma j$ | p p > a j | $1.964 \pm 0.001 \cdot 10^4$ $^{+31.2\%\ +1.7\%}_{-26.0\%\ -1.8\%}$ | | $5.218 \pm 0.025 \cdot 10^4$ $^{+24.5\%\ +1.4\%}_{-21.4\%\ -1.6\%}$ |
| a.10 $pp \to \gamma jj$ | p p > a j j | $7.815 \pm 0.008 \cdot 10^3$ $^{+32.8\%\ +0.9\%}_{-24.2\%\ -1.2\%}$ | | $1.004 \pm 0.004 \cdot 10^4$ $^{+5.9\%\ +0.8\%}_{-10.9\%\ -1.2\%}$ |

+127 more pp processes from 1405.0301!

*(Not a random sampling)*

Alwall et al.

# Plot the pulls

$$t_{scale} = \frac{\sigma_{NLO} - \sigma_{LO}}{\Delta\sigma_{LO\ scale}}$$

# Which of these distributions do you expect?



$$t_{scale} = \frac{\sigma_{NLO} - \sigma_{LO}}{\Delta\sigma_{LO\ scale}}$$

# Pull distribution

# Pull distribution

# What processes populate the tail ?

| Process | $n_{\text{part}}$ | $\Delta\sigma/\sigma_0$ | $\frac{\sigma_{\text{NLO}} - \sigma_0}{\Delta\sigma}$ |
|---|---|---|---|
| p p > wpm | 1 | $1.54 \times 10^{-1}$ | 1.84 |
| p p > wpm j | 2 | $1.97 \times 10^{-1}$ | 1.96 |
| p p > wpm j j | 3 | $2.45 \times 10^{-1}$ | 0.59 |
| p p > wpm j j j | 4 | $4.10 \times 10^{-1}$ | 0.25 |
| p p > z | 1 | $1.46 \times 10^{-1}$ | 1.87 |
| p p > z j | 2 | $1.93 \times 10^{-1}$ | 1.82 |
| p p > z j j | 3 | $2.43 \times 10^{-1}$ | 0.56 |
| p p > z j j j | 4 | $4.08 \times 10^{-1}$ | 0.27 |
| p p > a j | 2 | $3.12 \times 10^{-1}$ | 5.33 |
| p p > a j j | 3 | $3.28 \times 10^{-1}$ | 0.85 |
| p p > w+ w- wpm | 3 | $1.00 \times 10^{-3}$ | 610.69 |
| p p > z w+ w- | 3 | $8.00 \times 10^{-3}$ | 92.39 |
| p p > z z wpm | 3 | $1.00 \times 10^{-2}$ | 85.00 |
| p p > z z z | 3 | $1.00 \times 10^{-3}$ | 302.75 |
| p p > a w+ w- | 3 | $1.90 \times 10^{-2}$ | 42.33 |
| p p > a a wpm | 3 | $4.40 \times 10^{-2}$ | 47.24 |
| p p > a z wpm | 3 | $1.00 \times 10^{-3}$ | 1244.49 |
| p p > a z z | 3 | $2.00 \times 10^{-2}$ | 17.24 |

| Process | $\dfrac{\Delta\sigma}{\sigma_0}$ | | $n$ | $\dfrac{\Delta\sigma}{n\,\sigma_0}$ | |
|---|---|---|---|---|---|
| p p > j j | $+2.49 \times 10^{-1}$ | $-1.88 \times 10^{-1}$ | 2 | $+1.24 \times 10^{-1}$ | $-9.40 \times 10^{-2}$ |
| p p > b b | $+2.52 \times 10^{-1}$ | $-1.89 \times 10^{-1}$ | 2 | $+1.26 \times 10^{-1}$ | $-9.45 \times 10^{-2}$ |
| p p > t t | $+2.90 \times 10^{-1}$ | $-2.11 \times 10^{-1}$ | 2 | $+1.45 \times 10^{-1}$ | $-1.06 \times 10^{-1}$ |
| p p > j j j | $+4.38 \times 10^{-1}$ | $-2.84 \times 10^{-1}$ | 3 | $+1.46 \times 10^{-1}$ | $-9.47 \times 10^{-2}$ |
| p p > b b j | $+4.41 \times 10^{-1}$ | $-2.85 \times 10^{-1}$ | 3 | $+1.47 \times 10^{-1}$ | $-9.50 \times 10^{-2}$ |
| p p > t t j | $+4.51 \times 10^{-1}$ | $-2.90 \times 10^{-1}$ | 3 | $+1.50 \times 10^{-1}$ | $-9.67 \times 10^{-2}$ |
| p p > b b j j | $+6.18 \times 10^{-1}$ | $-3.56 \times 10^{-1}$ | 4 | $+1.54 \times 10^{-1}$ | $-8.90 \times 10^{-2}$ |
| p p > b b b b | $+6.17 \times 10^{-1}$ | $-3.56 \times 10^{-1}$ | 4 | $+1.54 \times 10^{-1}$ | $-8.90 \times 10^{-2}$ |
| p p > t t j j | $+6.14 \times 10^{-1}$ | $-3.56 \times 10^{-1}$ | 4 | $+1.53 \times 10^{-1}$ | $-8.90 \times 10^{-2}$ |
| p p > t t t t | $+6.38 \times 10^{-1}$ | $-3.65 \times 10^{-1}$ | 4 | $+1.60 \times 10^{-1}$ | $-9.12 \times 10^{-2}$ |
| p p > t t b b | $+6.21 \times 10^{-1}$ | $-3.57 \times 10^{-1}$ | 4 | $+1.55 \times 10^{-1}$ | $-8.93 \times 10^{-2}$ |
| average | | | | $+1.47 \times 10^{-1}$ | $-9.34 \times 10^{-2}$ |

Table 1: Scale dependence for LHC processes with only QCD particles in the final state. For each process, we report the relative scale uncertainty, the number of final state particles, and the per-particle relative scale uncertainty.

→Tilman Plehn's 'reference process' method

$$\frac{\Delta\sigma_{\mathrm{ref}}}{\sigma_0} = n \times \left\langle \frac{\Delta\sigma}{n\sigma_0} \right\rangle_{\mathrm{QCD}} .$$

# Make correction in UQ for EW processes



Much reduced tails

Tilman Plehn's 'reference process' method

# Make correction in UQ for EW processes



Much reduced tails

Tilman Plehn's 'reference process' method

# Leaves us wanting more …

- Would be even more interesting to repeat study for NLO → NNLO, differential distributions

- Can we use ML to automatically find patterns of failure ?

- Why did we find a Gaussian-ish core ?

- Impact: A new method for cross-checks within experiment collaborations

# Snowmass Whitepaper: Recommendations for the future

- Common language for uncertainty between ML and Physics communities

- Funding to test ML UQ methods for physics

- Create benchmark datasets for uncertainty tests

- Develop and study interpretability methods

# Snowmass Whitepaper: Recommendations for the future



Snowmass 2021: Summary of past work and future roadmap

- Common language for uncertainty between ML and Physics communities

- Funding to test ML UQ methods for physics

- Create benchmark datasets for uncertainty tests

- Develop and study interpretability methods

# Conclusion

- ML more sensitive to simulation artefacts → building better uncertainty quantification tools

- ML lets us better propagate experimental uncertainties and build analyses optimised for all possibilities: HEP, Astro

- Solutions have wider use cases

  - Tractable likelihoods

  - Optimise true objective with differentiable programming [Inferno, NEOS]

  - Uncertainty quantification of ML-simulators? [Performance metrics, Bayesian networks]

  - Learn physics from machine: Mapping ML into a human-readable space [CNN to EFPs]

And more cool solutions to come !

Thank you!

# Known unknowns

## Simulation using Standard Model of particle physics



Train ML models on simulation, apply on data

## Unlabelled data from LHC



Detector state in sim: Z=1

Detector state in data: Z = ?

Systematic differences lead to systematic uncertainties

# Make correction in UQ for EW processes

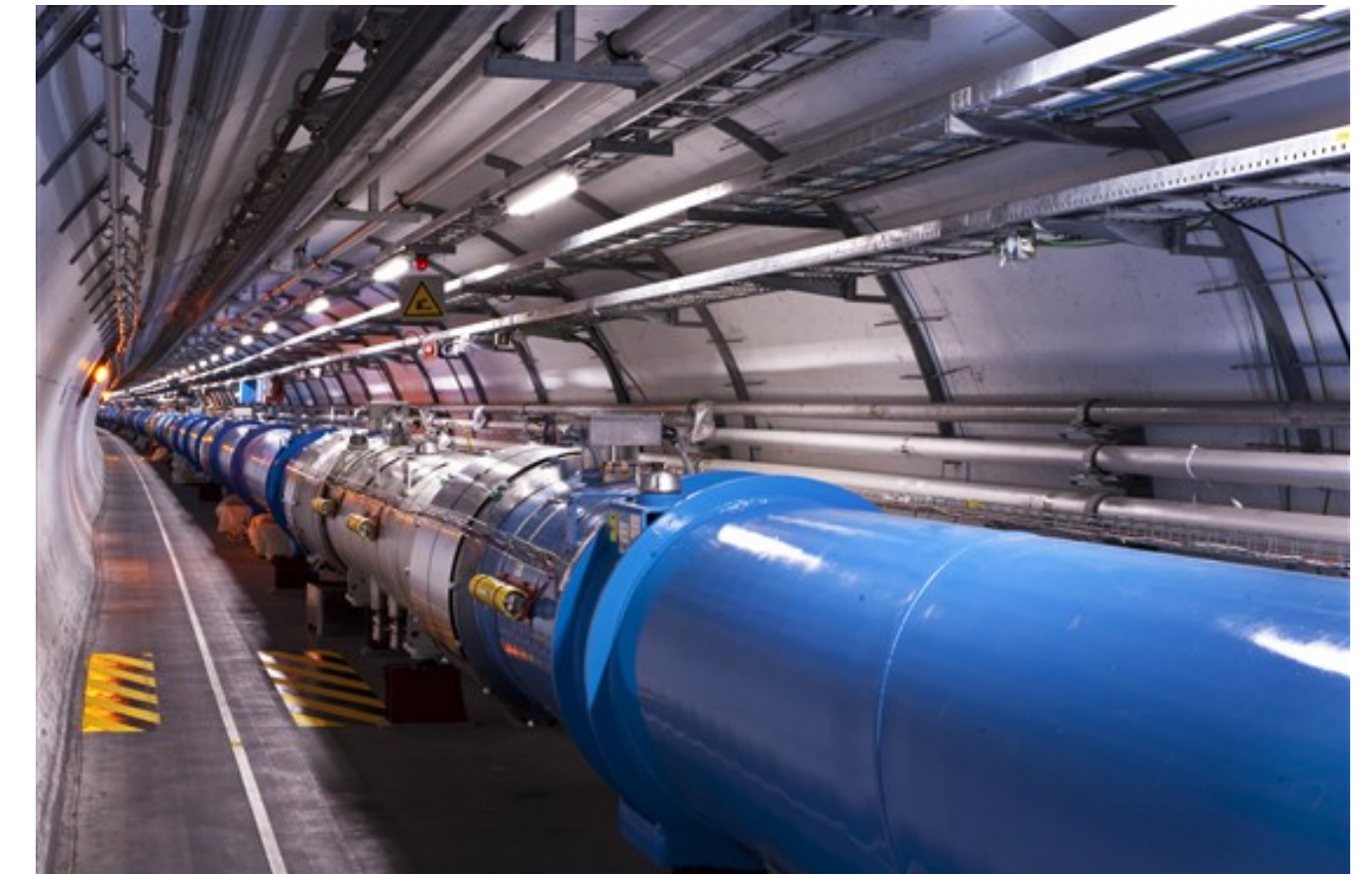| Process | $n_{\text{part}}$ | $\Delta\sigma/\sigma_0$ | $\frac{\sigma_{\text{NLO}}-\sigma_0}{\Delta\sigma}$ | $\Delta\sigma_{\text{ref}}/\sigma_0$ | $\frac{\sigma_{\text{NLO}}-\sigma_0}{\Delta\sigma_{\text{ref}}}$ |
|---|---|---|---|---|---|
| p p > wpm | 1 | $1.54 \times 10^{-1}$ | 1.84 | $1.47 \times 10^{-1}$ | 1.92 |
| p p > wpm j | 2 | $1.97 \times 10^{-1}$ | 1.96 | $2.94 \times 10^{-1}$ | 1.31 |
| p p > wpm j j | 3 | $2.45 \times 10^{-1}$ | 0.59 | $4.41 \times 10^{-1}$ | 0.33 |
| p p > wpm j j j | 4 | $4.10 \times 10^{-1}$ | 0.25 | $5.88 \times 10^{-1}$ | 0.18 |
| p p > z | 1 | $1.46 \times 10^{-1}$ | 1.87 | $1.47 \times 10^{-1}$ | 1.86 |
| p p > z j | 2 | $1.93 \times 10^{-1}$ | 1.82 | $2.94 \times 10^{-1}$ | 1.19 |
| p p > z j j | 3 | $2.43 \times 10^{-1}$ | 0.56 | $4.41 \times 10^{-1}$ | 0.31 |
| p p > z j j j | 4 | $4.08 \times 10^{-1}$ | 0.27 | $5.88 \times 10^{-1}$ | 0.19 |
| p p > a j | 2 | $3.12 \times 10^{-1}$ | 5.33 | $2.94 \times 10^{-1}$ | 5.66 |
| p p > a j j | 3 | $3.28 \times 10^{-1}$ | 0.85 | $4.41 \times 10^{-1}$ | 0.63 |
| p p > w+ w- wpm | 3 | $1.00 \times 10^{-3}$ | 610.69 | $4.41 \times 10^{-1}$ | 1.39 |
| p p > z w+ w- | 3 | $8.00 \times 10^{-3}$ | 92.39 | $4.41 \times 10^{-1}$ | 1.68 |
| p p > z z wpm | 3 | $1.00 \times 10^{-2}$ | 85.00 | $4.41 \times 10^{-1}$ | 1.93 |
| p p > z z z | 3 | $1.00 \times 10^{-3}$ | 302.75 | $4.41 \times 10^{-1}$ | 0.69 |
| p p > a w+ w- | 3 | $1.90 \times 10^{-2}$ | 42.33 | $4.41 \times 10^{-1}$ | 1.82 |
| p p > a a wpm | 3 | $4.40 \times 10^{-2}$ | 47.24 | $4.41 \times 10^{-1}$ | 4.72 |
| p p > a z wpm | 3 | $1.00 \times 10^{-3}$ | 1244.49 | $4.41 \times 10^{-1}$ | 2.82 |
| p p > a z z | 3 | $2.00 \times 10^{-2}$ | 17.24 | $4.41 \times 10^{-1}$ | 0.78 |

Tilman Plehn's 'reference process' method

# Surviving tails

| Process | $n_{\text{part}}$ | $\Delta\sigma/\sigma_0$ | $\frac{\sigma_{\text{NLO}}-\sigma_0}{\Delta\sigma}$ | $\Delta\sigma_{\text{ref}}/\sigma_0$ | $\frac{\sigma_{\text{NLO}}-\sigma_0}{\Delta\sigma_{\text{ref}}}$ |
|---------|-------------------|-------------------------|-----------------------------------------------------|--------------------------------------|------------------------------------------------------------------|
| p p > h | 1 | $3.48 \times 10^{-1}$ | 3.02 | $1.47 \times 10^{-1}$ | 7.15 |

Large corrections loop-induced  2->1 process

# Mapping machine-learned physics into a human-readable space



Signal/Background Pairs

Black-Box Guided Search

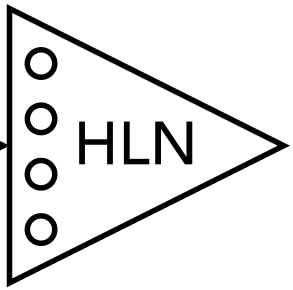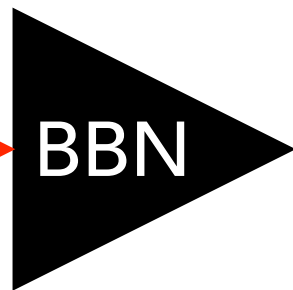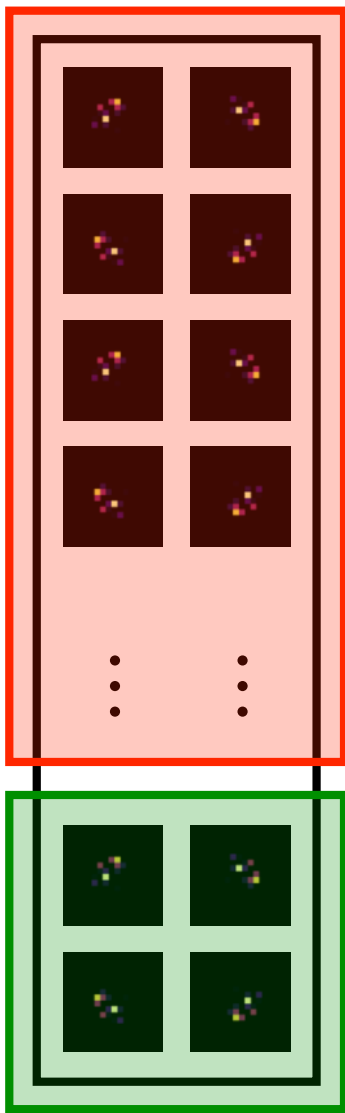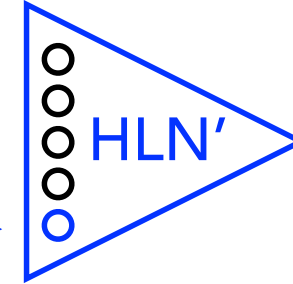| Rank | EFP | $\kappa$ | $\beta$ | Chrom # | $ADO[EFP, CNN]_{x_6}$ | $AUC[EFP]$ | $ADO[6HL + EFP, CNN]_{x_{all}}$ | $AUC[6HL + EFP]$ |
|---|---|---|---|---|---|---|---|---|
| 1 | | 2 | $\frac{1}{2}$ | 3 | 0.6207 | 0.8031 | 0.9714 | $0.9528 \pm 0.0003$ |
| 2 | | 2 | $\frac{1}{2}$ | 3 | 0.6205 | 0.8203 | 0.9714 | 0.9524 |
| 3 | | 0 | − | 1 | 0.6205 | 0.6737 | 0.9715 | 0.9525 |
| 4 | | 2 | $\frac{1}{2}$ | 3 | 0.6199 | 0.8301 | 0.9715 | 0.9527 |
| 5 | | 2 | $\frac{1}{2}$ | 3 | 0.6197 | 0.8290 | 0.9714 | 0.9527 |
| 6 | | 2 | $\frac{1}{2}$ | 3 | 0.6196 | 0.8251 | 0.9715 | 0.9522 |
| 7 | | 0 | $\frac{1}{2}$ | 2 | 0.6187 | 0.7511 | 0.9715 | 0.9526 |
| 8 | | 2 | $\frac{1}{2}$ | 3 | 0.6184 | 0.8257 | 0.9712 | 0.9527 |
| 9 | | 2 | $\frac{1}{2}$ | 3 | 0.6182 | 0.8090 | 0.9714 | 0.9527 |
| 10 | | 2 | $\frac{1}{2}$ | 3 | 0.6180 | 0.8314 | 0.9714 | 0.9526 |
| 60 | | 0 | 1 | 2 | 0.6163 | 0.7194 | 0.9715 | 0.9525 |
| 341 | | −1 | $\frac{1}{2}$ | 4 | 0.6142 | 0.6286 | 0.9714 | 0.9509 |
| 589 | | 0 | 2 | 2 | 0.6109 | 0.7579 | 0.9714 | 0.9523 |
| 3106 | | −1 | − | 1 | 0.5891 | 0.5882 | 0.9714 | 0.9510 |
| 3519 | | $\frac{1}{2}$ | $\frac{1}{2}$ | 2 | 0.5664 | 0.7698 | 0.9715 | 0.9524 |
| 3521 | | $\frac{1}{2}$ | − | 1 | 0.5663 | 0.7093 | 0.9714 | 0.9522 |
| 5531 | | 1 | 2 | 1 | 0.5290 | 0.7454 | 0.9714 | 0.9507 |
| 5554 | | 1 | $\frac{1}{2}$ | 2 | 0.5279 | 0.8210 | 0.9713 | 0.9505 |
| 5610 | | 2 | − | 1 | 0.5245 | 0.7117 | 0.9714 | 0.9507 |
| 5657 | | 1 | 1 | 3 | 0.5224 | 0.8257 | 0.9712 | 0.9506 |
| 5793 | | 1 | 1 | 2 | 0.5191 | 0.8640 | 0.9714 | 0.9505 |
| 6052 | | 1 | 2 | 3 | 0.5153 | 0.8500 | 0.9716 | 0.9504 |
| 7438 | | 1 | 2 | 2 | 0.5011 | 0.8835 | 0.9716 | 0.9506 |

Faucett et al.

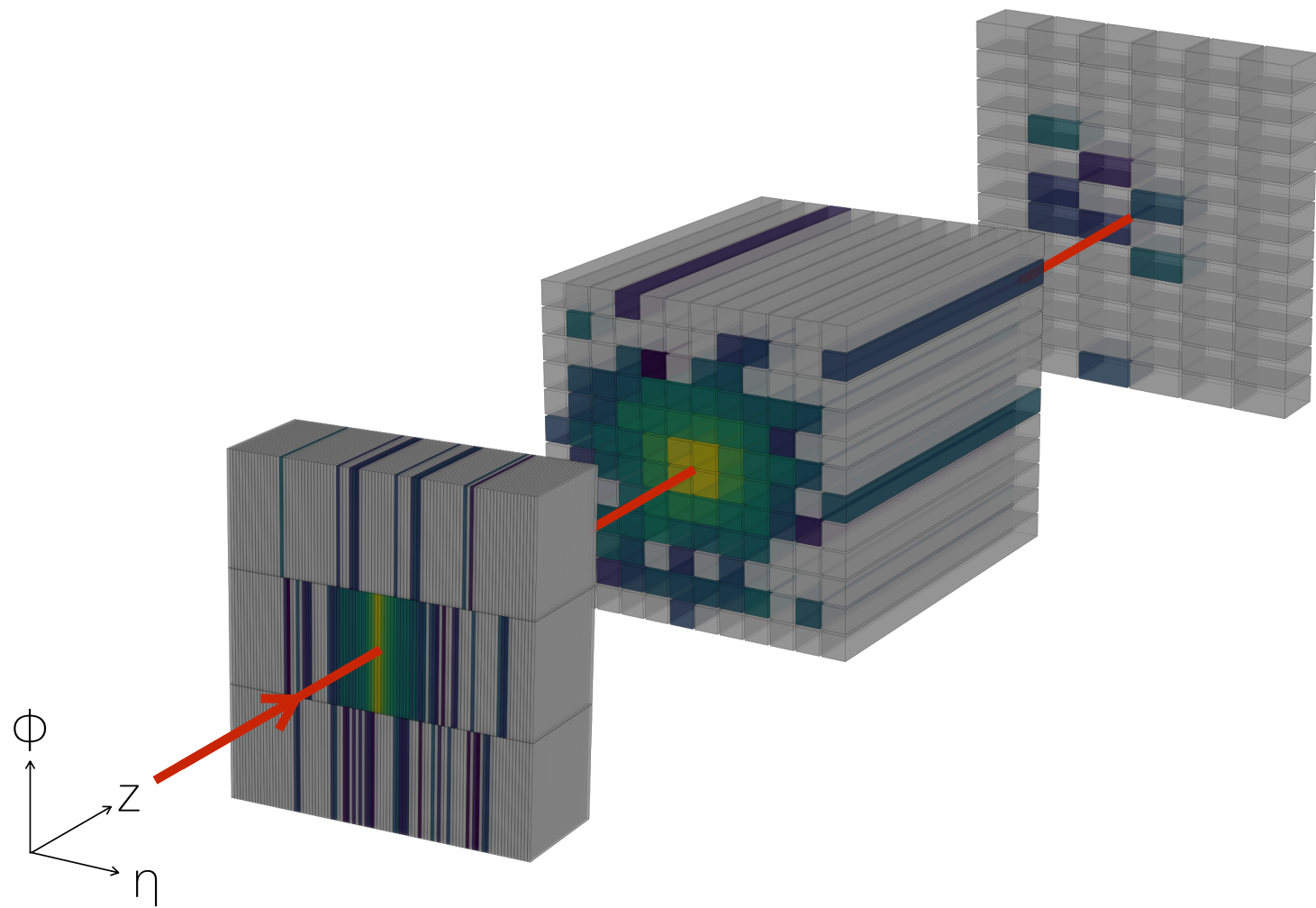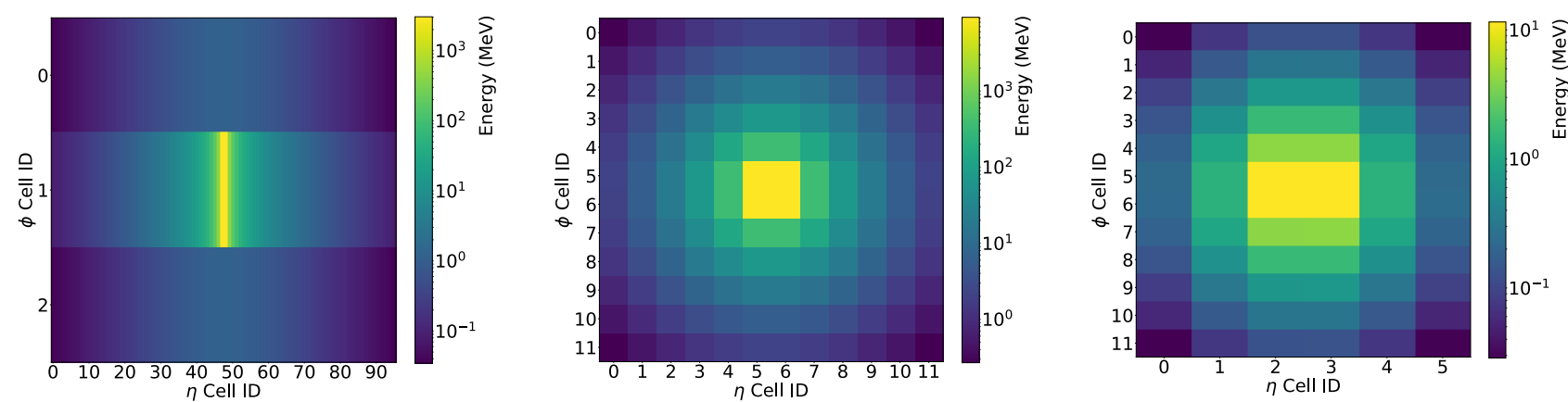# Performance Metrics for Generative Models

# Generative Models for Simulation

Ghosh, ATLAS Collaboration, 2019

Paganini et al.

# Evaluating Fast Calo Simulators

## 5.4 Classifier metrics

In much of the GAN literature (see e.g. [8]), a common metric is to train classifiers to distinguish between different categories of data (e.g. $e^+$ vs. $\pi^+$), and to see if there is any difference in classifier performance when real data and generated data are interchanged. For example, one might train a classifier on $e^+$ vs. $\pi^+$ GEANT4 images, and compare this to a classifier trained on $e^+$ vs. $\pi^+$ GAN images. If the classifier trained on real images performs similarly to the classifier trained on generated images, then this is evidence that the generated images are approximating the real images well. One can repeat this test for different combinations of real and generated data.

The ultimate test of whether $p_{\text{generated}}(x) = p_{\text{data}}(x)$ would be a direct binary classifier between real and generated images of the *same* type. If the generated and true probability

# Can we automise the evaluation ?

Krause and Shih, 2021

## 5.4 Classifier metrics

In much of the GAN literature (see e.g. [8]), a common metric is to train classifiers to distinguish between different categories of data (e.g. $e^+$ vs. $\pi^+$), and to see if there is any difference in classifier performance when real data and generated data are interchanged. For example, one might train a classifier on $e^+$ vs. $\pi^+$ GEANT4 images, and compare this to a classifier trained on $e^+$ vs. $\pi^+$ GAN images. If the classifier trained on real images performs similarly to the classifier trained on generated images, then this is evidence that the generated images are approximating the real images well. One can repeat this test for different combinations of real and generated data.

The ultimate test of whether $p_{\text{generated}}(x) = p_{\text{data}}(x)$ would be a direct binary classifier between real and generated images of the *same* type. If the generated and true probability
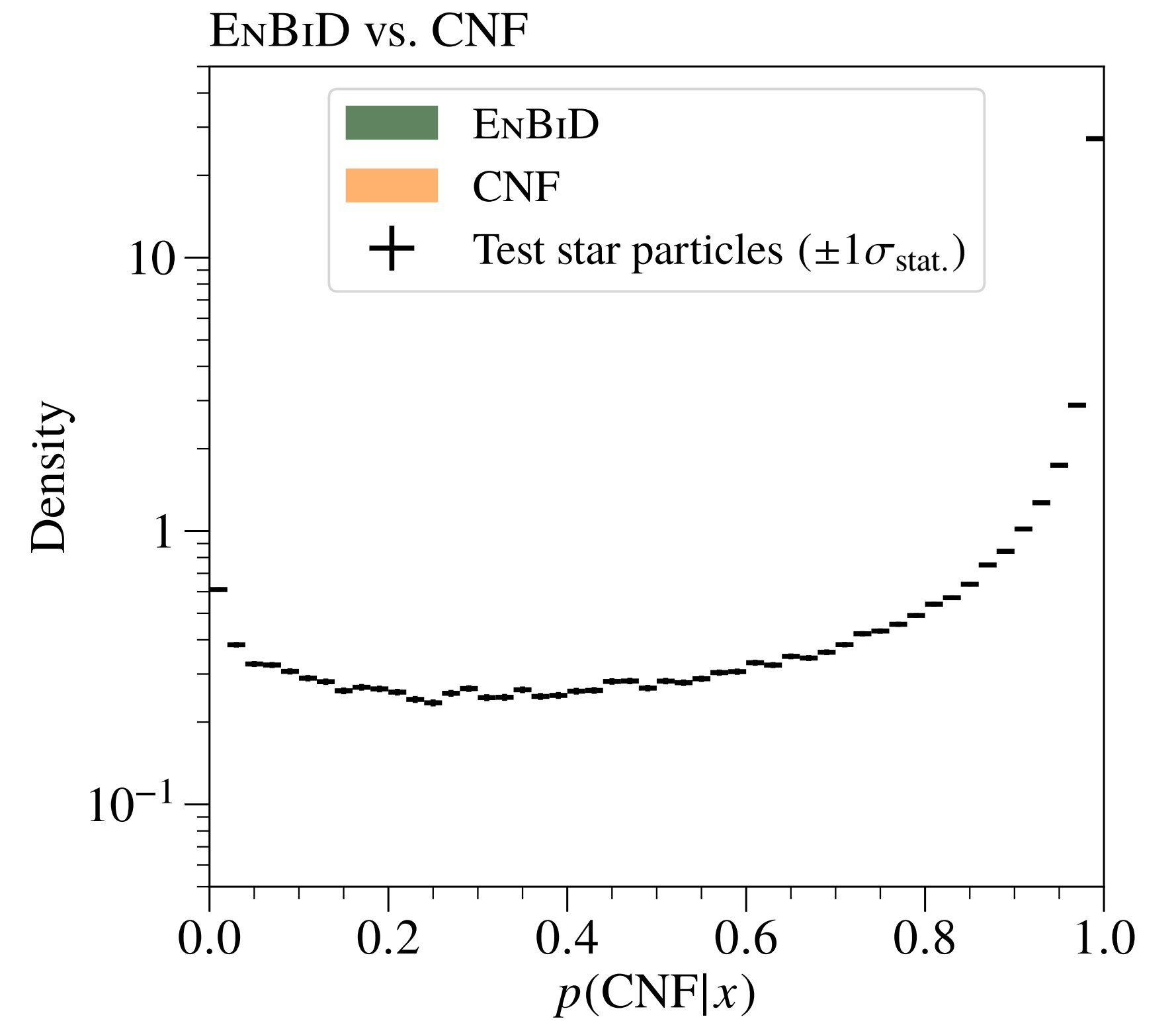
Classify Geant4 vs generated and use AUC as single metric

# Another classifier test

Compare two generative models:

Classify generative model1 vs model2, check if test dataset agrees better with one or the other



ᴇɴʙɪᴅ vs. CNF

# A comparison of metrics

**On the Evaluation of Generative Models in High Energy Physics**

Raghav Kansal,* Anni Li, and Javier Duarte
*University of California San Diego, La Jolla, CA 92093, USA*

Nadezda Chernyavskaya, Maurizio Pierini
*European Center for Nuclear Research (CERN), 1211 Geneva 23, Switzerland*

Breno Orzari, Thiago Tomei
*Universidade Estadual Paulista, São Paulo/SP, CEP 01049-010, Brazil*

(Dated: November 21, 2022)

Detailed comparison on Gaussian toys where you have full control

Application on jet dataset with hand designed distortions

# Study



- $FGD_\infty$, MMD unbiased
- W too expensive for large N

$FGD_\infty$ most promising

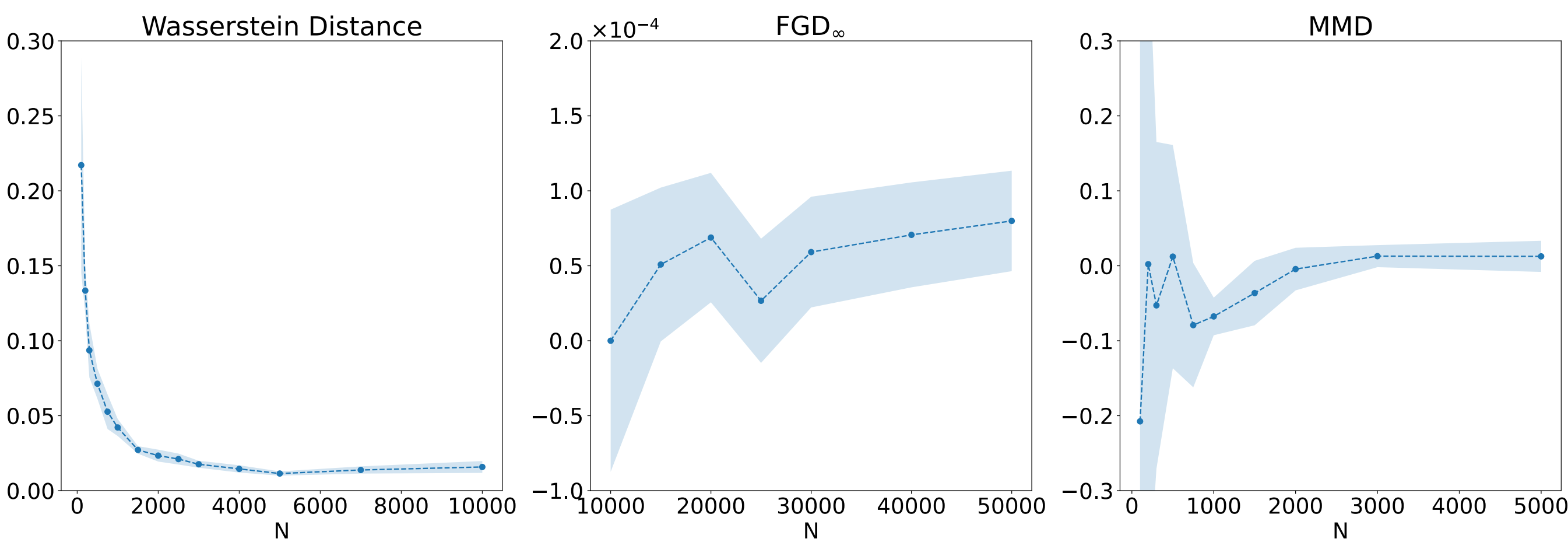(but no sensitivity to higher moments, requires extrapolation)

| Metric | Truth | Shift $\mu_x$ by $1\sigma$ | Shift $\mu_x$ by $0.1\sigma$ | Zero covariance | Multiply (co)variances by 10 | Divide (co)variances by 10 | Mixture of Two Gaussians 1 | Mixture of Two Gaussians 2 |
|---|---|---|---|---|---|---|---|---|
| Wasserstein | $0.016 \pm 0.004$ | $1.14 \pm 0.02$ | $0.043 \pm 0.008$ | $0.077 \pm 0.006$ | $9.8 \pm 0.1$ | $0.97 \pm 0.01$ | $\mathbf{0.036 \pm 0.003}$ | $\mathbf{0.191 \pm 0.005}$ |
| $FGD_\infty \times 10^3$ | $0.08 \pm 0.03$ | $\mathbf{1011 \pm 1}$ | $\mathbf{11.0 \pm 0.1}$ | $\mathbf{32.3 \pm 0.2}$ | $9400 \pm 8$ | $\mathbf{935.1 \pm 0.7}$ | $0.07 \pm 0.03$ | $0.03 \pm 0.03$ |
| MMD | $0.01 \pm 0.02$ | $16.4 \pm 0.9$ | $0.07 \pm 0.04$ | $0.40 \pm 0.08$ | $\mathbf{19k \pm 1k}$ | $4.3 \pm 0.1$ | $0.06 \pm 0.02$ | $0.35 \pm 0.03$ |
| Precision | $0.972 \pm 0.005$ | $0.91 \pm 0.01$ | $0.976 \pm 0.004$ | $0.969 \pm 0.006$ | $0.34 \pm 0.01$ | $1.0 \pm 0.0$ | $0.975 \pm 0.003$ | $0.9976 \pm 0.0007$ |
| Recall | $0.997 \pm 0.001$ | $0.992 \pm 0.003$ | $0.997 \pm 0.001$ | $0.9976 \pm 0.0006$ | $0.998 \pm 0.001$ | $0.58 \pm 0.02$ | $0.996 \pm 0.001$ | $0.9970 \pm 0.0009$ |
| Density | $3.23 \pm 0.06$ | $2.48 \pm 0.08$ | $3.19 \pm 0.07$ | $3.1 \pm 0.1$ | $0.60 \pm 0.02$ | $5.7 \pm 0.3$ | $2.99 \pm 0.09$ | $0.989 \pm 0.009$ |
| Coverage | $0.876 \pm 0.002$ | $0.780 \pm 0.006$ | $0.872 \pm 0.005$ | $0.872 \pm 0.004$ | $0.60 \pm 0.01$ | $0.406 \pm 0.008$ | $0.871 \pm 0.002$ | $0.956 \pm 0.006$ |

0.0    0.000  0.025  0.050  0.075  0.100  0.125  0.150  0.175  0.200
Jet $m/p_T$

0.0    0.000  0.025  0.050  0.075  0.100  0.125  0.150  0.175  0.200
Jet $m/p_T$

# Jet Study

| Metric | Truth | Smeared | Shifted | Removing tail | Particle features smeared | Particle $\eta^{\rm rel}$ smeared | Particle $p_{\rm T}^{\rm rel}$ smeared | Particle $p_{\rm T}^{\rm rel}$ shifted |
|---|---|---|---|---|---|---|---|---|
| $W_1^M \times 10^3$ | $0.28 \pm 0.05$ | $2.1 \pm 0.2$ | $6.0 \pm 0.3$ | $0.6 \pm 0.2$ | $1.7 \pm 0.2$ | $0.9 \pm 0.3$ | $0.5 \pm 0.2$ | $5.8 \pm 0.2$ |
| Wasserstein EFP | $0.02 \pm 0.01$ | $0.09 \pm 0.05$ | $0.10 \pm 0.02$ | $0.016 \pm 0.007$ | $0.19 \pm 0.08$ | $0.03 \pm 0.01$ | $0.03 \pm 0.02$ | $0.06 \pm 0.02$ |
| FGD$_\infty$ EFP $\times 10^3$ | $0.01 \pm 0.02$ | $\mathbf{21.5 \pm 0.3}$ | $\mathbf{26.8 \pm 0.3}$ | $\mathbf{2.31 \pm 0.07}$ | $23.4 \pm 0.3$ | $\mathbf{3.59 \pm 0.09}$ | $2.29 \pm 0.05$ | $28.9 \pm 0.2$ |
| MMD EFP $\times 10^3$ | $-0.006 \pm 0.005$ | $0.17 \pm 0.06$ | $0.9 \pm 0.1$ | $0.03 \pm 0.02$ | $0.35 \pm 0.09$ | $0.08 \pm 0.05$ | $0.01 \pm 0.02$ | $1.8 \pm 0.1$ |
| Precision EFP | $0.9 \pm 0.1$ | $0.94 \pm 0.04$ | $0.978 \pm 0.005$ | $0.88 \pm 0.08$ | $0.7 \pm 0.1$ | $0.94 \pm 0.06$ | $0.7 \pm 0.1$ | $0.79 \pm 0.09$ |
| Recall EFP | $0.9 \pm 0.1$ | $0.88 \pm 0.07$ | $0.97 \pm 0.01$ | $0.92 \pm 0.06$ | $0.83 \pm 0.05$ | $0.92 \pm 0.07$ | $0.8 \pm 0.1$ | $0.8 \pm 0.1$ |
| Wasserstein PN | $1.65 \pm 0.06$ | $1.7 \pm 0.1$ | $2.4 \pm 0.4$ | $1.71 \pm 0.08$ | $4.5 \pm 0.1$ | $1.79 \pm 0.05$ | $4.0 \pm 0.4$ | $7.6 \pm 0.2$ |
| FGD$_\infty$ PN $\times 10^3$ | $0.8 \pm 0.7$ | $40 \pm 2$ | $193 \pm 9$ | $5.0 \pm 0.9$ | $\mathbf{1250 \pm 10}$ | $20 \pm 1$ | $\mathbf{1230 \pm 10}$ | $\mathbf{3640 \pm 10}$ |
| MMD PN $\times 10^3$ | $-2 \pm 2$ | $4 \pm 8$ | $80 \pm 10$ | $-1 \pm 4$ | $500 \pm 100$ | $3 \pm 2$ | $560 \pm 60$ | $1100 \pm 40$ |
| Precision PN | $0.68 \pm 0.07$ | $0.64 \pm 0.04$ | $0.71 \pm 0.06$ | $0.73 \pm 0.03$ | $0.09 \pm 0.04$ | $0.75 \pm 0.08$ | $0.08 \pm 0.04$ | $0.39 \pm 0.08$ |
| Recall PN | $0.70 \pm 0.05$ | $0.61 \pm 0.04$ | $0.61 \pm 0.08$ | $0.73 \pm 0.06$ | $0.014 \pm 0.009$ | $0.7 \pm 0.1$ | $0.01 \pm 0.01$ | $0.57 \pm 0.09$ |
| Classifier LLF AUC | $0.50$ | $0.52$ | $0.54$ | $0.50$ | $0.97$ | $0.81$ | $0.93$ | $0.99$ |
| Classifier HLF AUC | $0.50$ | $0.53$ | $0.55$ | $0.50$ | $0.84$ | $0.64$ | $0.74$ | $0.92$ |

- $FGD_\infty$ on EFPs does quite well in these tests

- Would be interesting to see it used and stress tested !

# Bayesian Generative Models
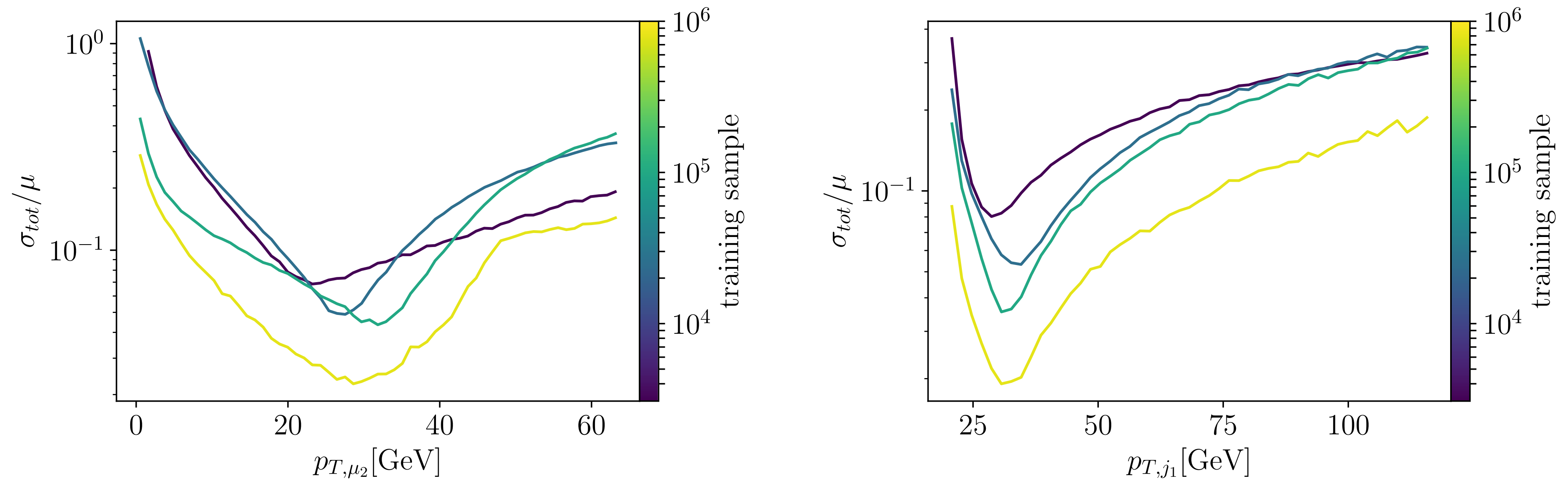
# Bayesian Generative Models

Figure 12: Relative uncertainty from the BINN for the $Z + 1$ jet sample, as a function of the size of the training sample.

Other uncertainty methods

# Differentiable Programming: Optimise your final objective directly
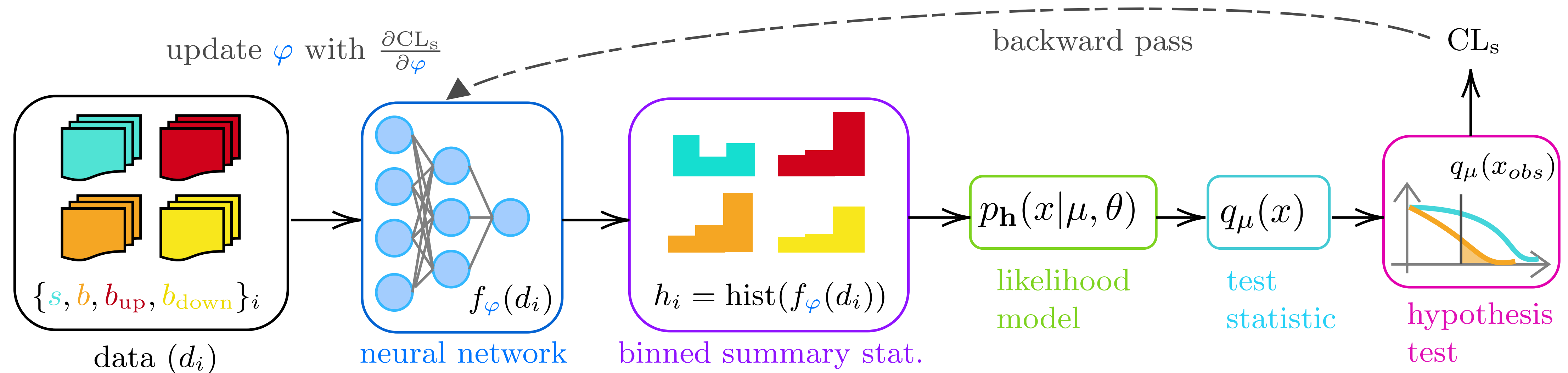
Following Inferno [de Castro et al.]



**Figure 1.** The pipeline for `neos`. The dashed line indicating the backward pass involves updating the weights $\varphi$ of the neural network via gradient descent.

# Unfolding with nuisance parameters

FIG. 6. Higgs boson cross section: the nominal detector-level spectra $m_{\gamma\gamma}$ (left) and $p_{\gamma\gamma}^{\mathrm{T}}$ (right) with $\epsilon_\gamma = 1$ reweighted by the trained $w_1$ conditioned at $\epsilon_\gamma = 1.2$ and compared to the spectra with $\epsilon_\gamma = 1.2$.

# More on uncertainty-aware networks

# Baseline Approach to Uncertainty Quantification

Train AI classifier on nominal data (assume detector state Z=1) and estimate uncertainties using alternate simulations

# Baseline Approach to Uncertainty Quantification

Train AI classifier on nominal data (assume detector state Z=1) and estimate uncertainties using alternate simulations



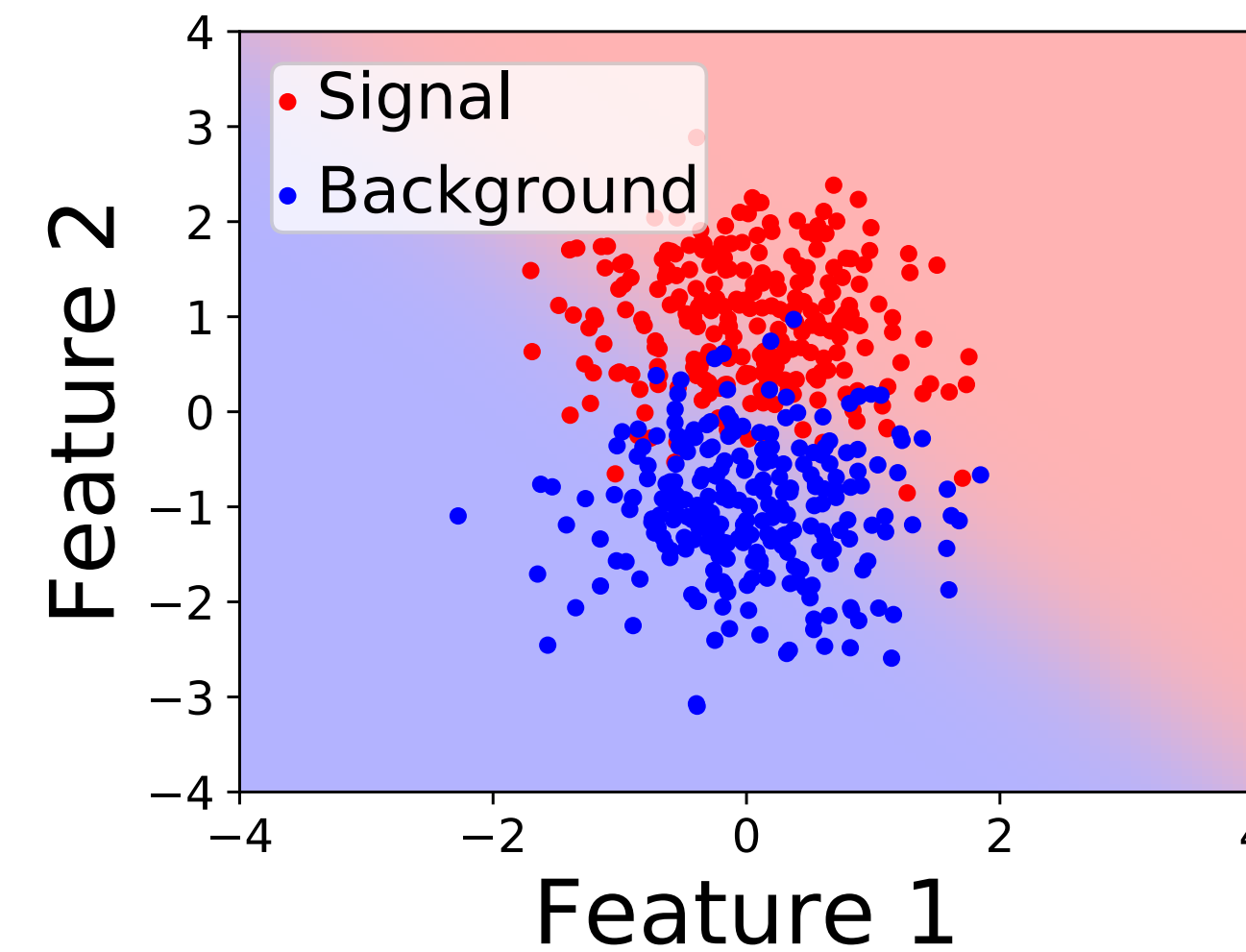Full statistical treatment → Expensive 'Profile Likelihood'

# Nominal and Systematic Up Examples

## Baseline Classifier



Nominal "Data"

AUC=0.978

Optimal

## Uncertainty-Aware Classifier



AUC=0.978

Optimal

SystUp "Data"

# Nominal and Systematic Up Examples

# Nominal and Systematic Up Examples



Baseline Classifier
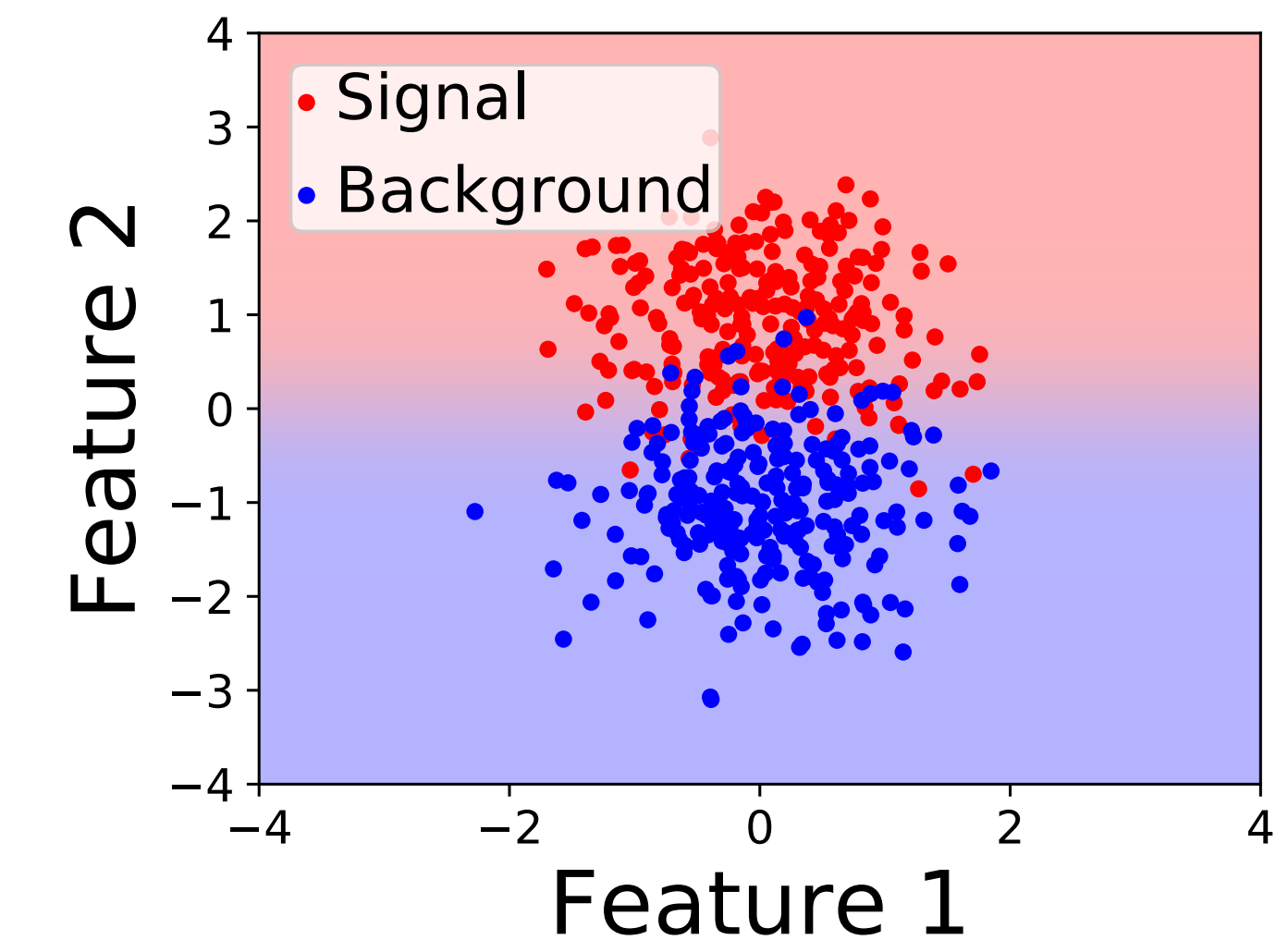
Uncertainty-Aware Classifier

Nominal "Data"

AUC=0.978
Optimal

AUC=0.978
Optimal
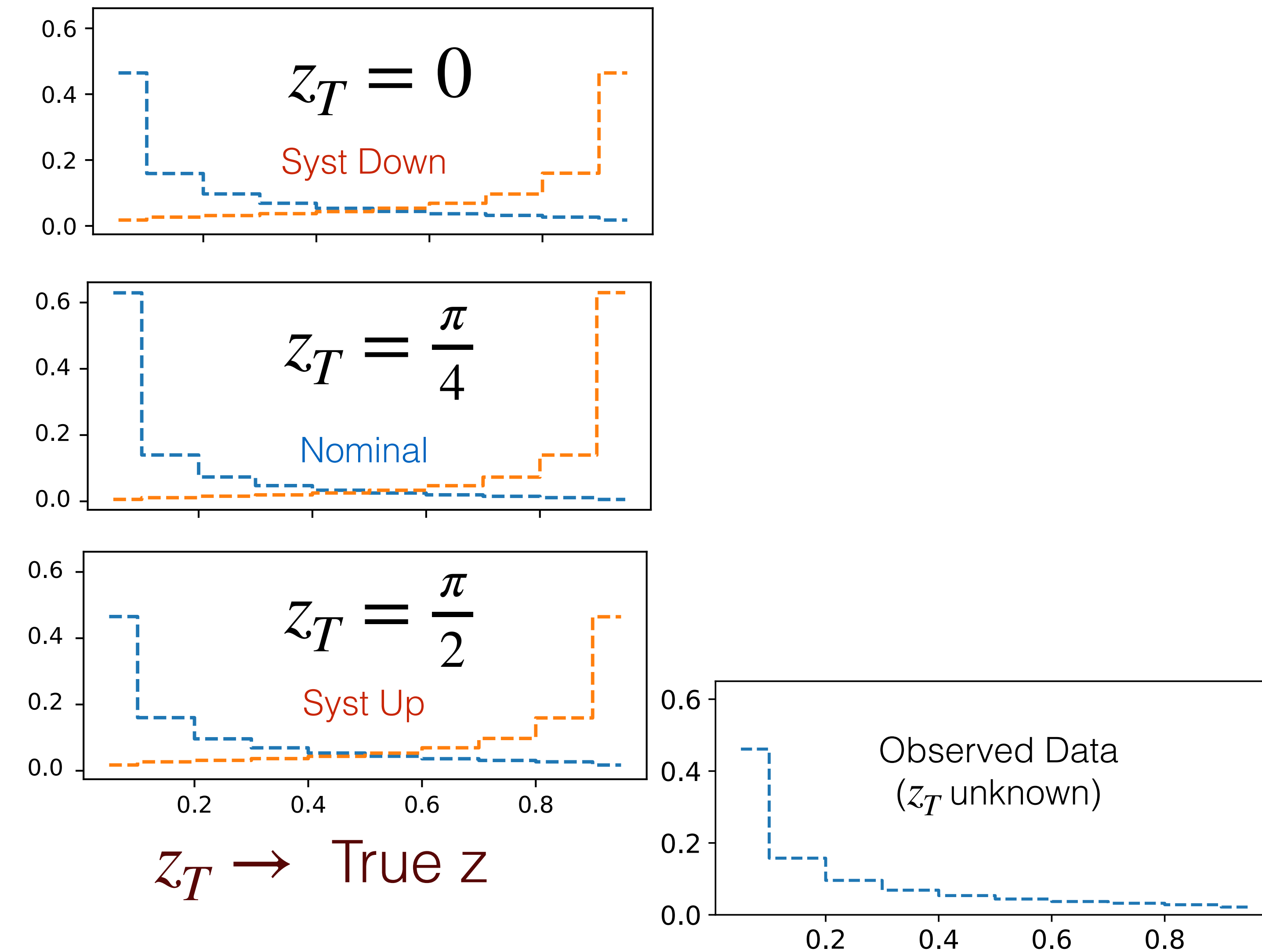
SystUp "Data"

AUC=0.924
**Sub-Optimal**

AUC=0.978
**Optimal**

Uncer-Aware Classifier is able to rotate its decision function based on Z while the Baseline Classifier decision function remains frozen[49]
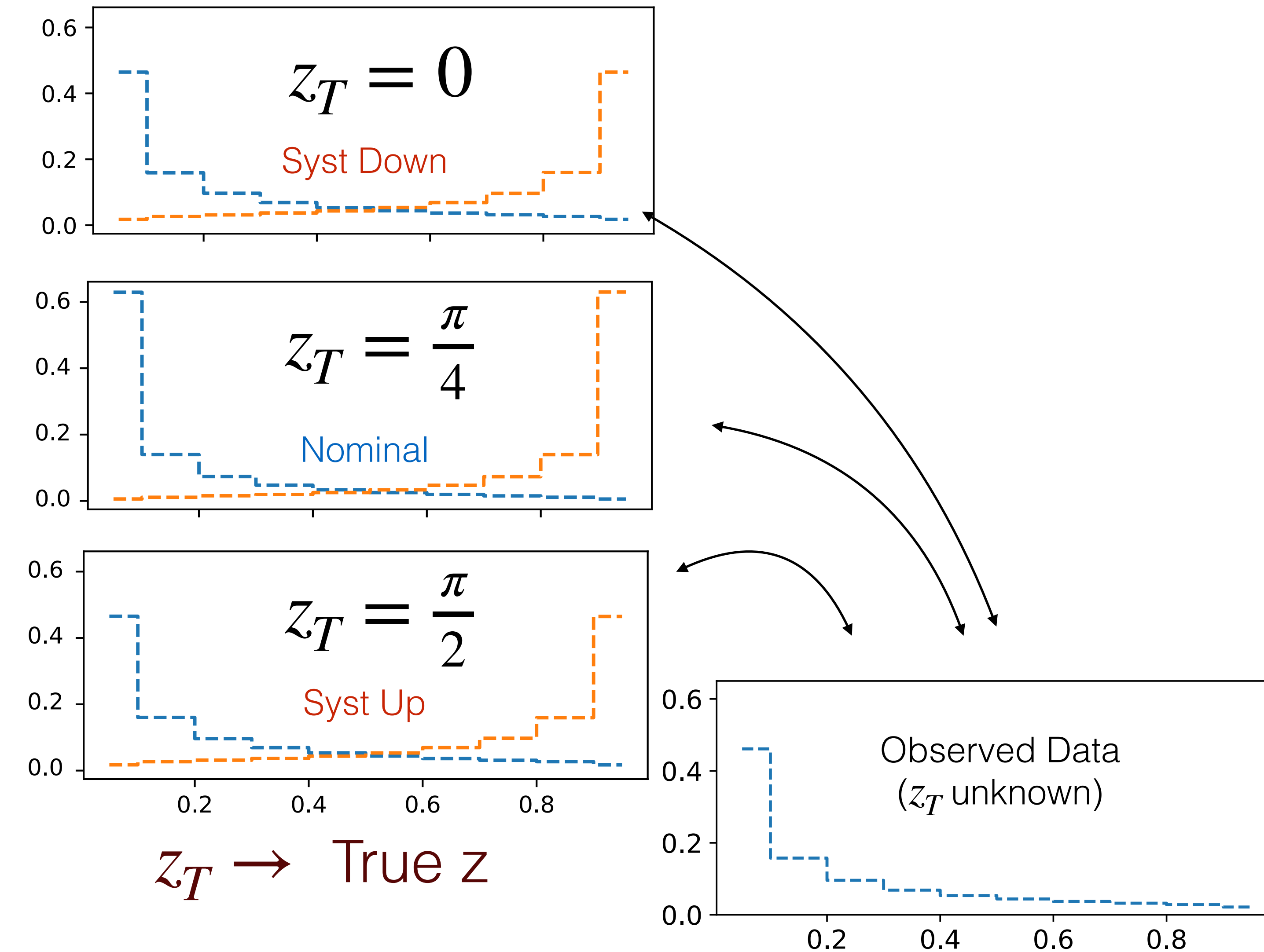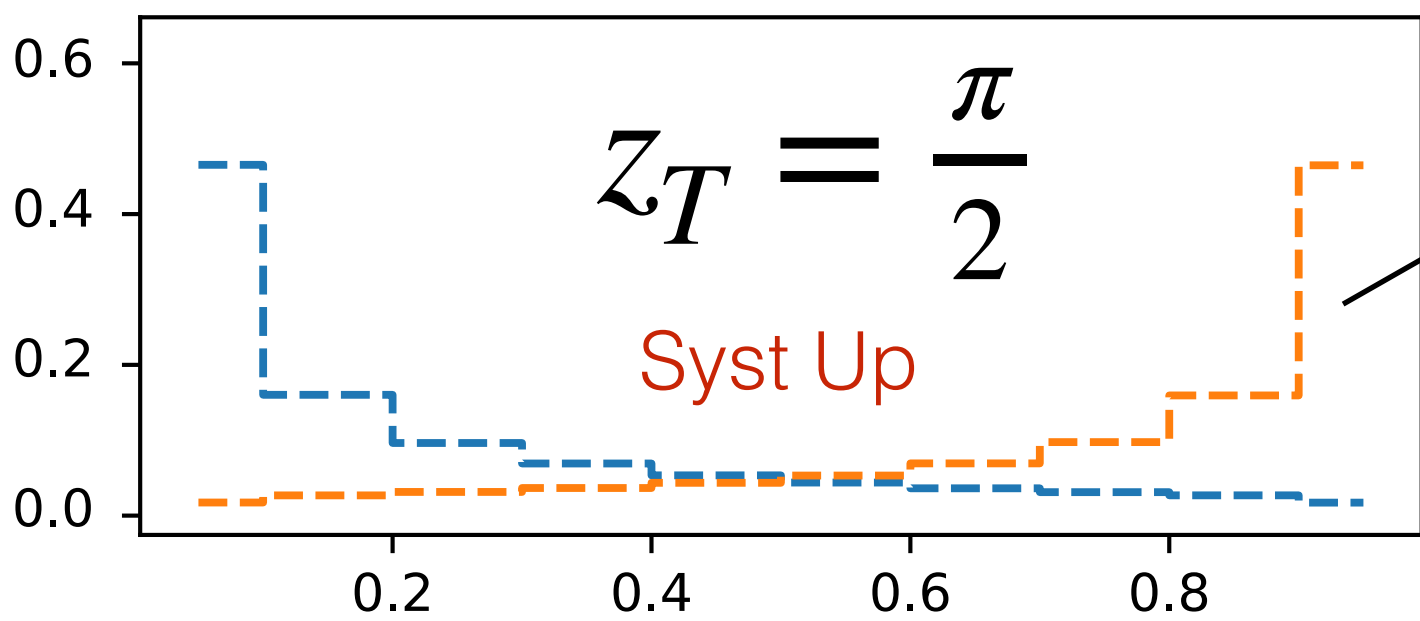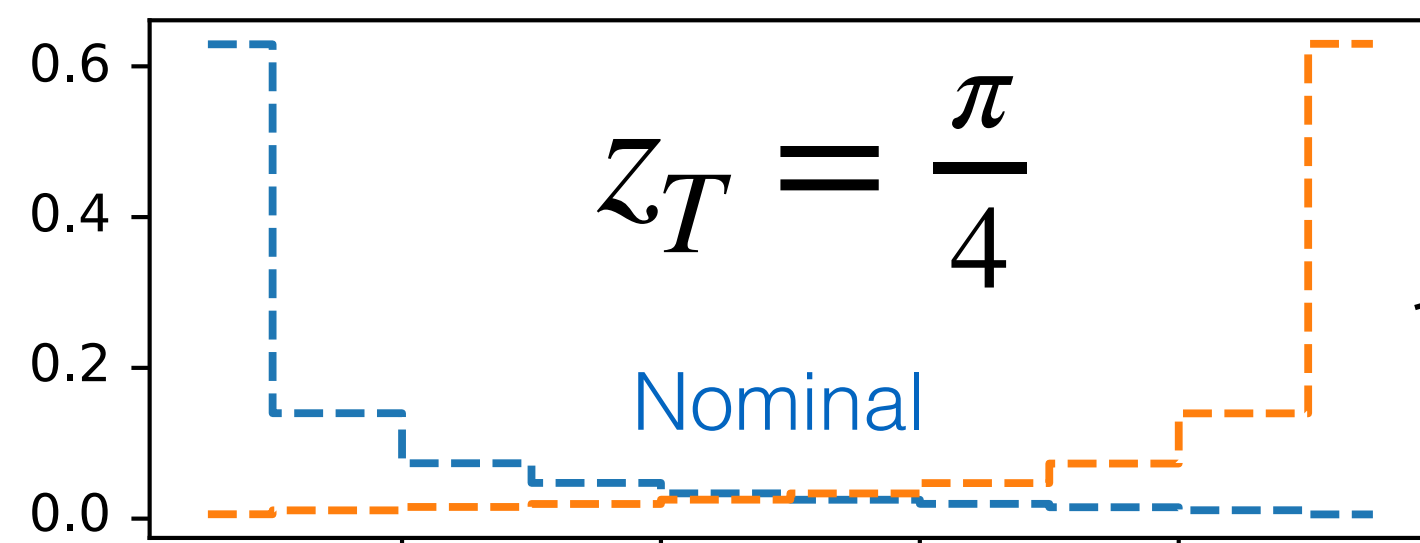
0.

# Scan the 2D Likelihood space in $Z$ vs $\mu$

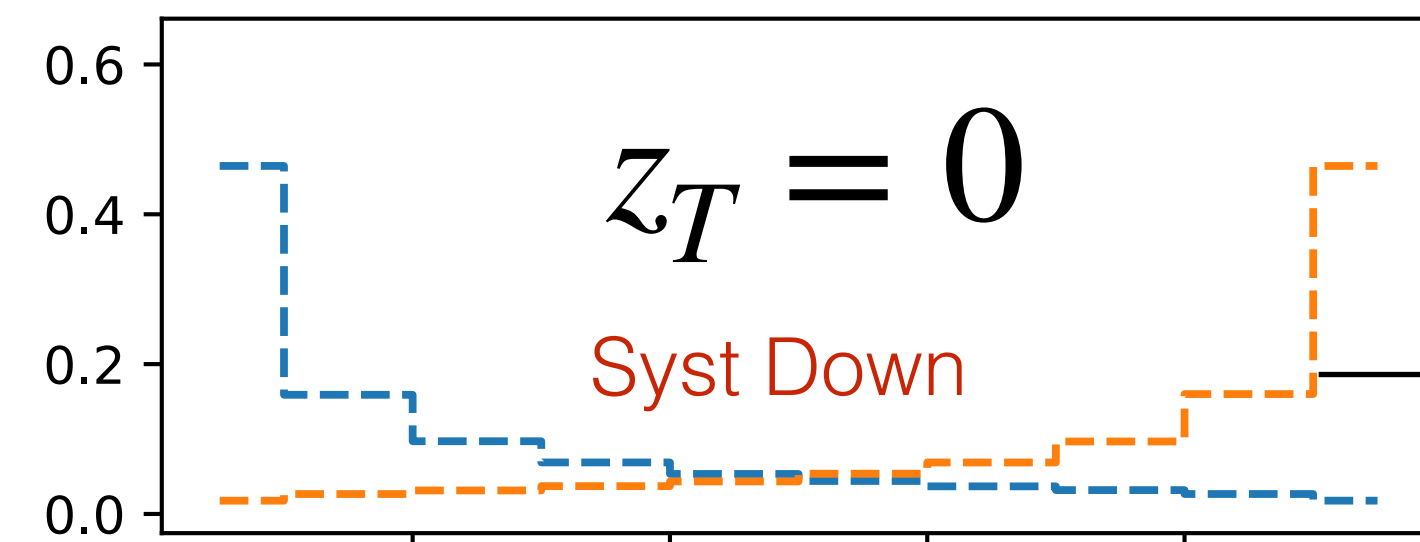Template **Baseline Classifier** Score Histograms for various Z

Template Scores for Clf for different Angles

Syst Down

0.08°

Template Scores for Clf for different Angles

0.31°

Nominal

Number of events normalized to unity

0.08°

Syst Up

0.16°

Observed Data
($z_T$ unknown)

0.55°

$z_T \rightarrow$ True z   0.31°

Score  0.39°

0.79°

Scan the 2D Likelihood space in $Z$ vs $\mu$

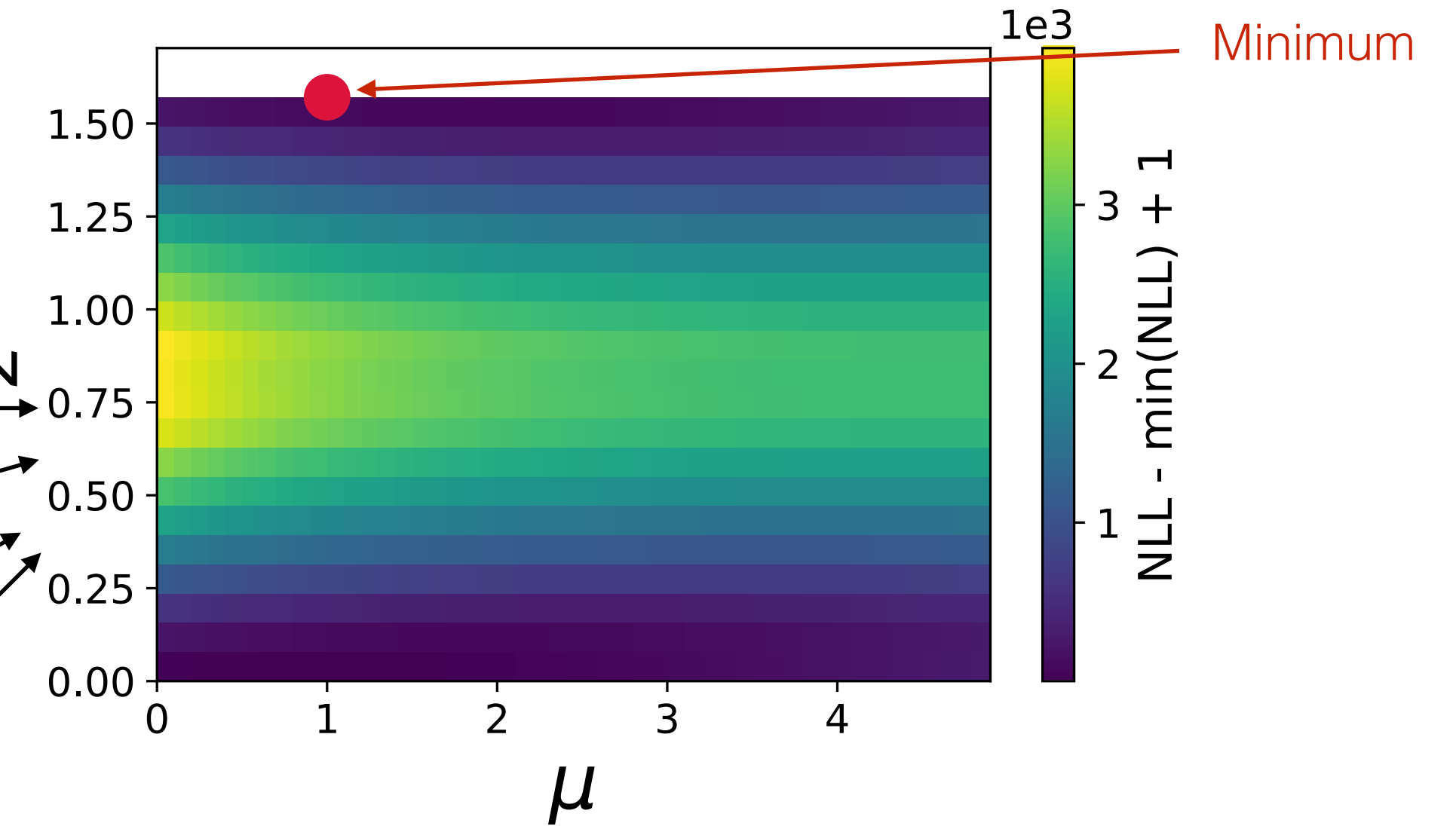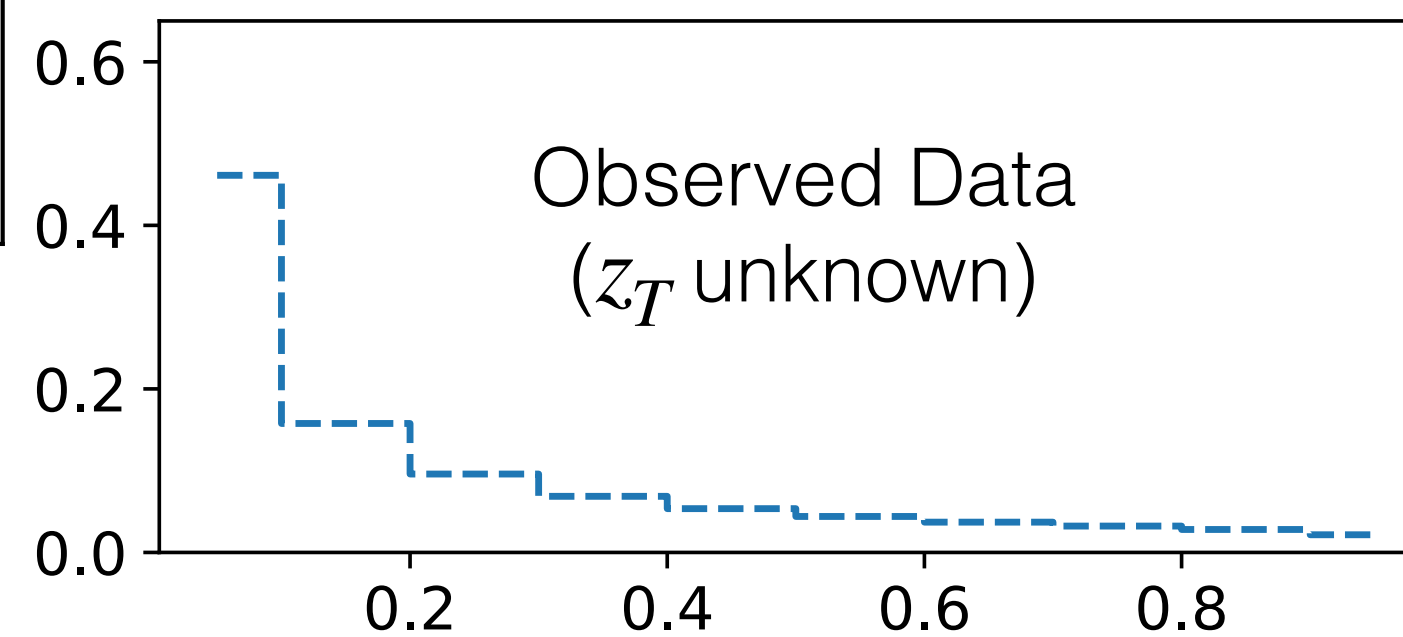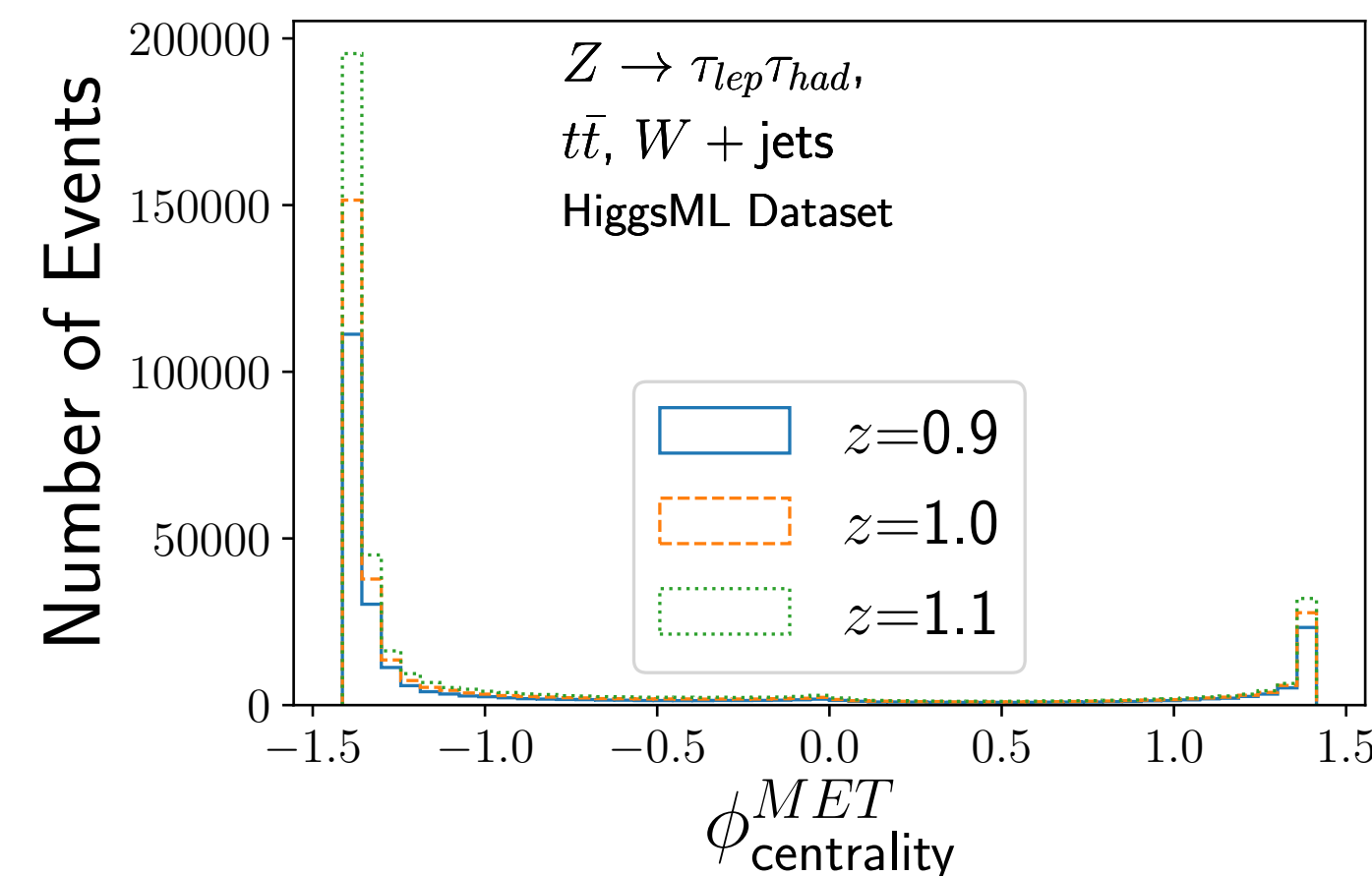Template **Baseline Classifier** Score Histograms for various Z
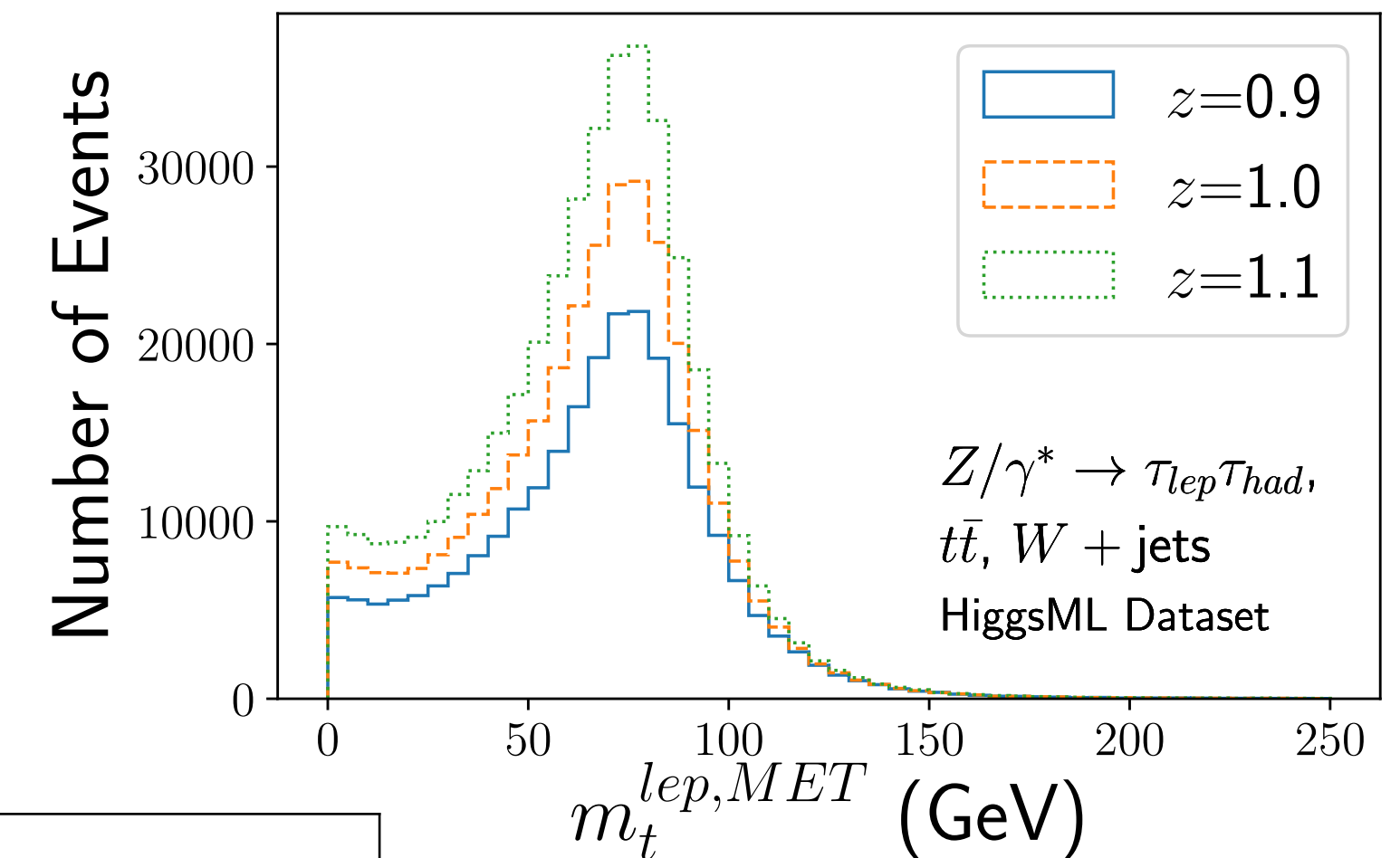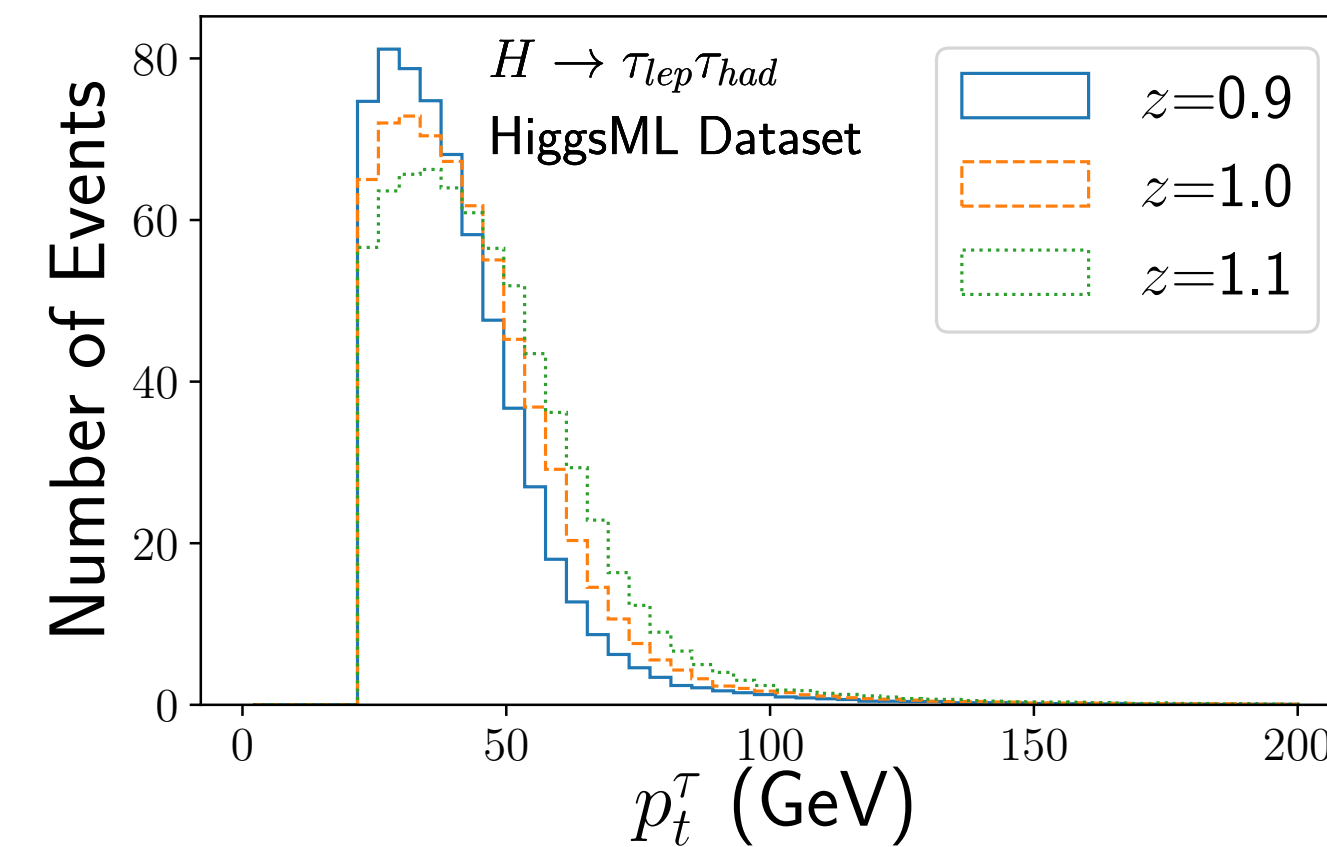
Template Scores for Clf for different Angles

Syst Down

Template Scores for Clf for different Angles

Nominal

Number of events normalized to unity

0.08°

0.31°

0.08°

0.16°

0.55°

Syst Up

Observed Data
($z_T$ unknown)

$z_T \rightarrow$ True z   0.31°

Score   0.39°

0.79°

50

Template Scores for Awe for different Angles

0.63°

0.16°

16°

different Angles

Minimum

$z$

$1.50$
$1.25$
$1.00$
$0.75$
$0.50$
$0.25$
$0.00$

$1e3$

NLL - min(NLL) + 1

$0$    $1$    $2$    $3$    $4$

$\mu$

0.08°

0.31°

0.39°

$z_T \xrightarrow{0.16°}$ True $z$

0.55°

0.63°

Parameter of Interest is Higgs signal strength μ, and TES is the nuisance parameter Z

Uncertainty-Aware coincides with classifier trained on true Z
⇒ Can't get much better than that!

# Test performance for "observed" data at nominal and above nominal Z



In every case the Aware Classifier is as good as the optimal one, no other technique matches its performance everywhere

# Idea fascinating also to ML researchers !

# Idea fascinating also to ML researchers !

- ML researchers assume i.i.d

- This technique exploits correlations between samples – a different paradigm

- Interesting applications outside of physics

# Idea fascinating also to ML researchers !

- ML researchers assume i.i.d

- This technique exploits correlations between samples – a different paradigm

- Interesting applications outside of physics

# Idea fascinating also to ML researchers !

- ML researchers assume i.i.d

- This technique exploits correlations between samples – a different paradigm

- Interesting applications outside of physics

For my handwriting this is '2', for yours it might be 'a'
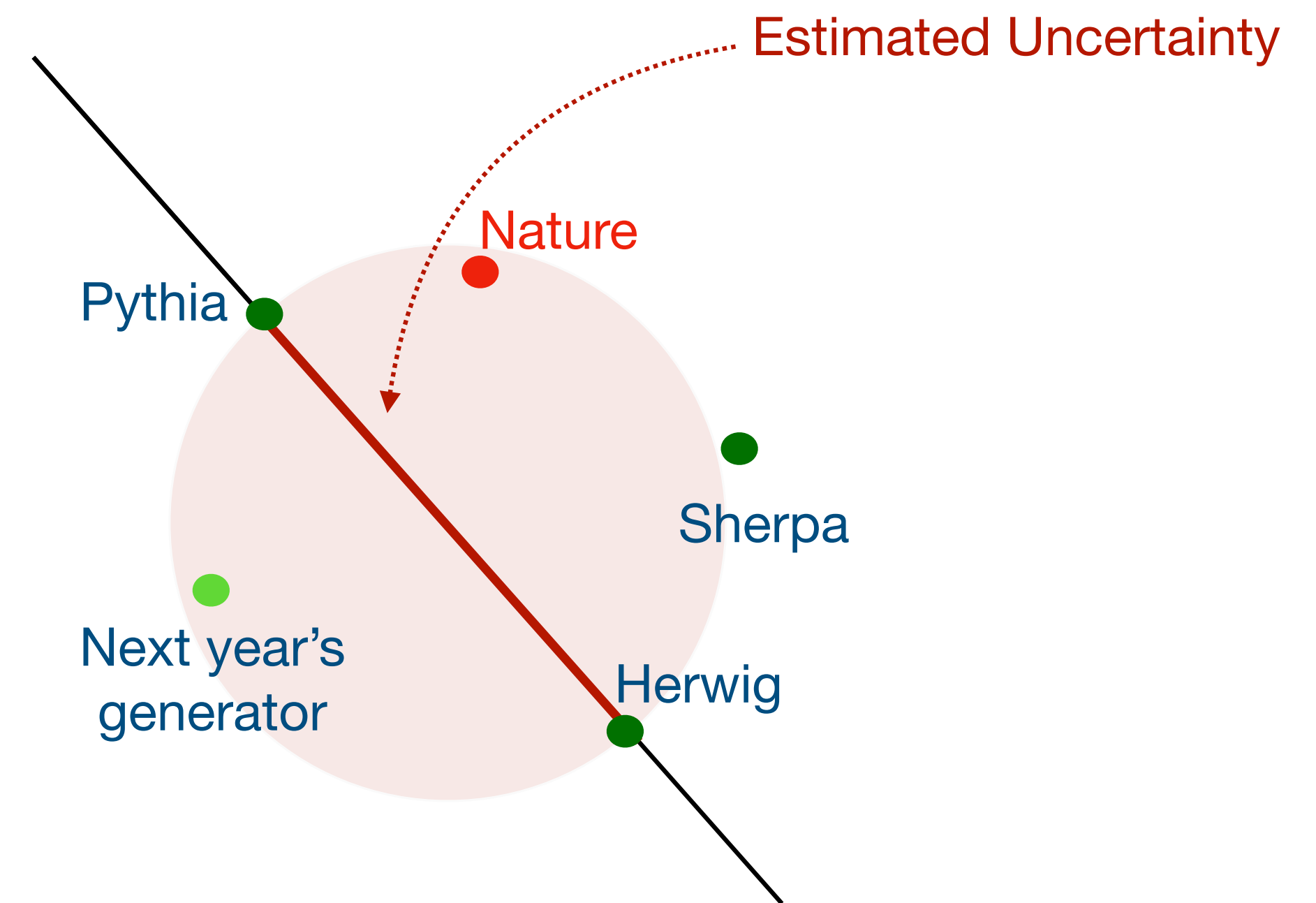ARM: Adapt to the individual + classify

ERM → 2
ARM → a

# Theory Uncertainties

# What are they ?

Theory uncertainties often describe our <u>lack of understanding / ability to calculate</u>

No statistical origin for them (such as auxiliary measurement)

Estimated Uncertainty

Nature

Pythia

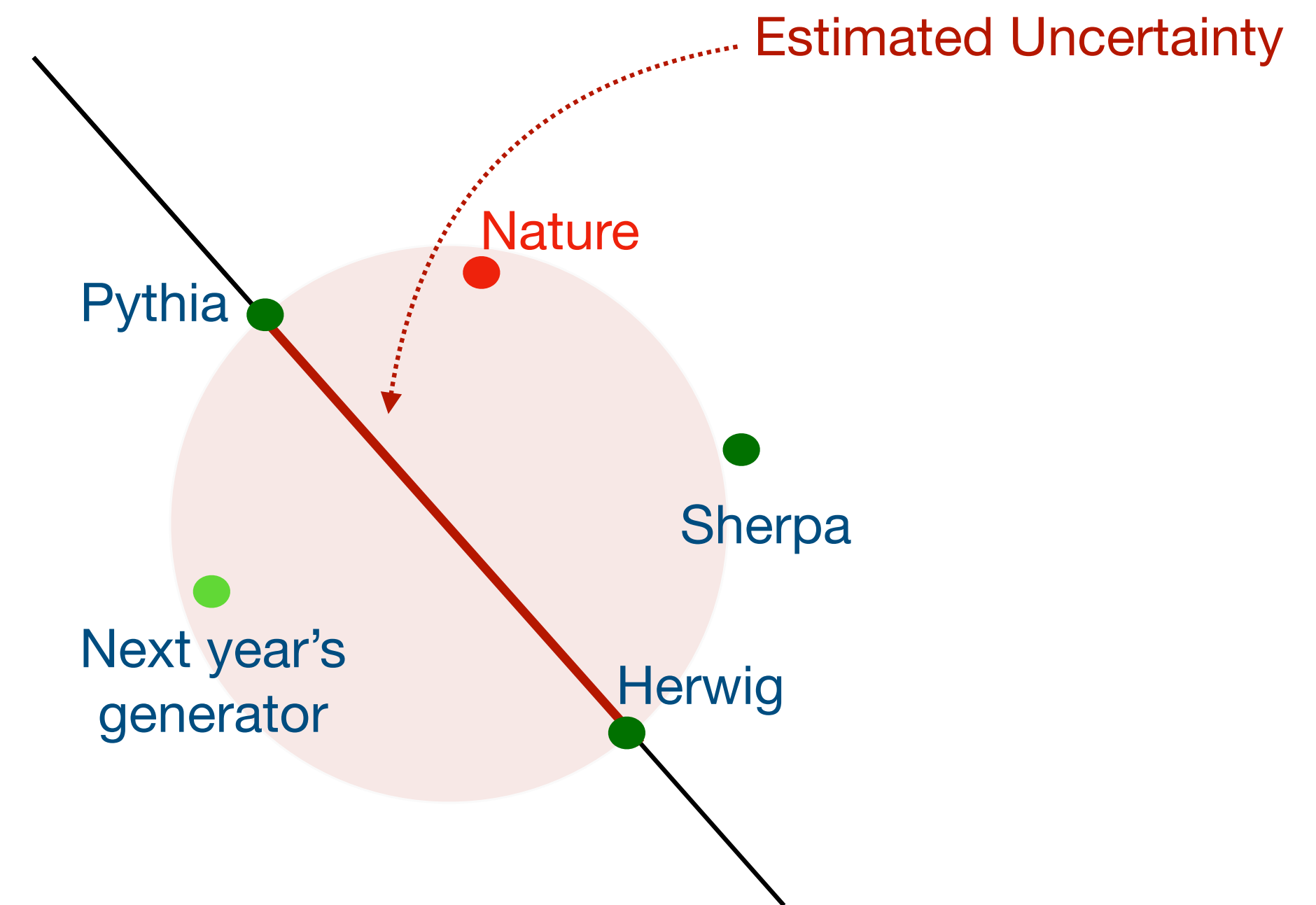Sherpa

Next year's generator

Herwig

# What are they ?

Theory uncertainties often describe our <u>lack of understanding / ability to calculate</u>

No statistical origin for them (such as auxiliary measurement)

Eg. <u>Hadronisation</u>:

- Few different packages to simulate it

- None are correct!

- Use difference in performance of your data analysis algorithm on **Pythia simulator** vs **Herwig simulator** ad-hoc estimate of uncertainty

Estimated Uncertainty

Nature

Pythia

Sherpa

Next year's generator

Herwig

# Goodhart's Law

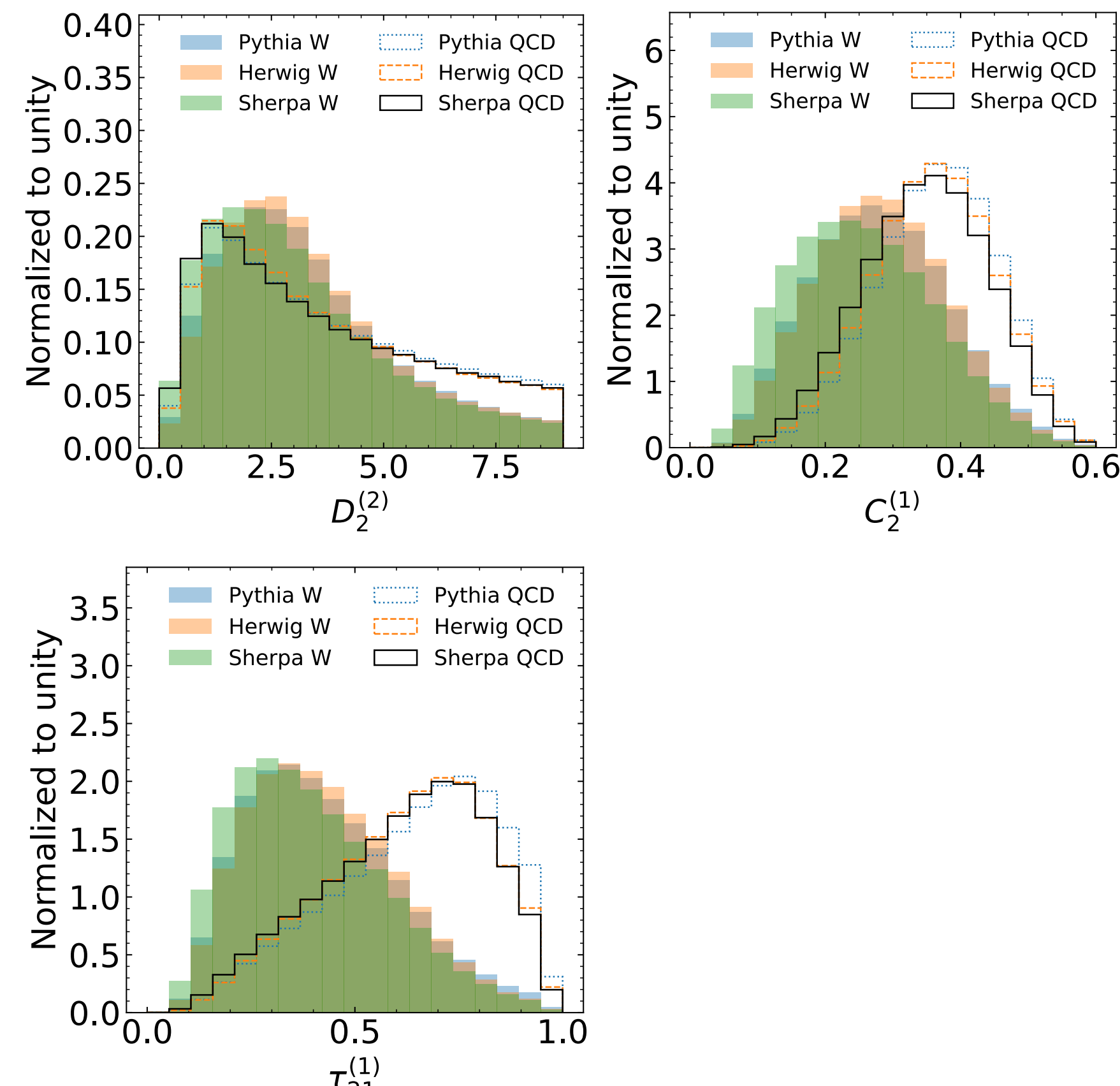When a measure becomes a target, it ceases to be a good measure

=> Dangerous to optimise proxy metrics of uncertainty

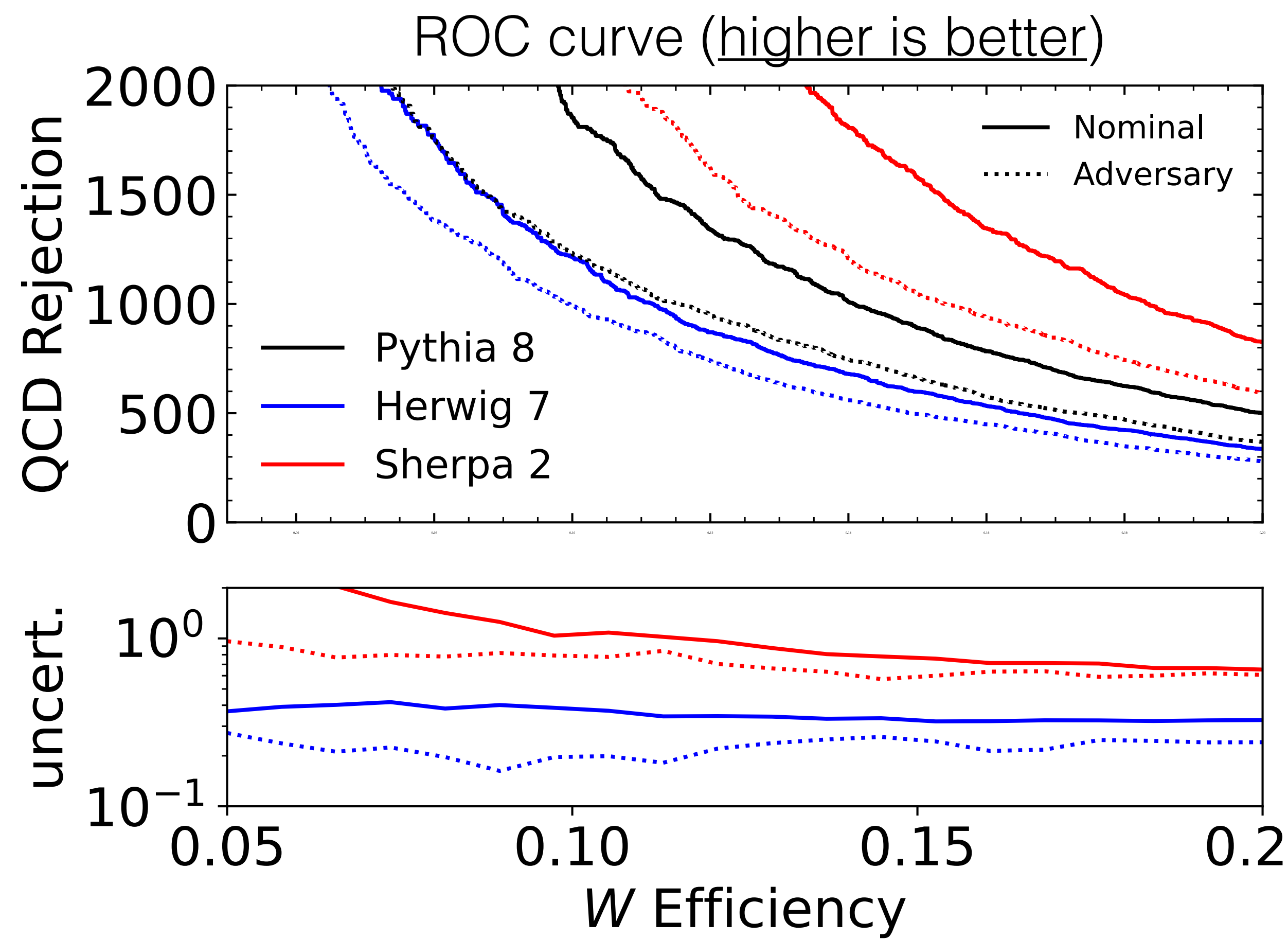# Case Study 1: Two-point uncertainty (fragmentation modelling)

Goal: W jets vs QCD jets
Decorrelation: Reduce difference in performance on Herwig vs Pythia
Cross-check: Test uncertainty estimate from {Herwig vs Pythia} using Sherpa

# Case Study 1: Two-point uncertainty - Result

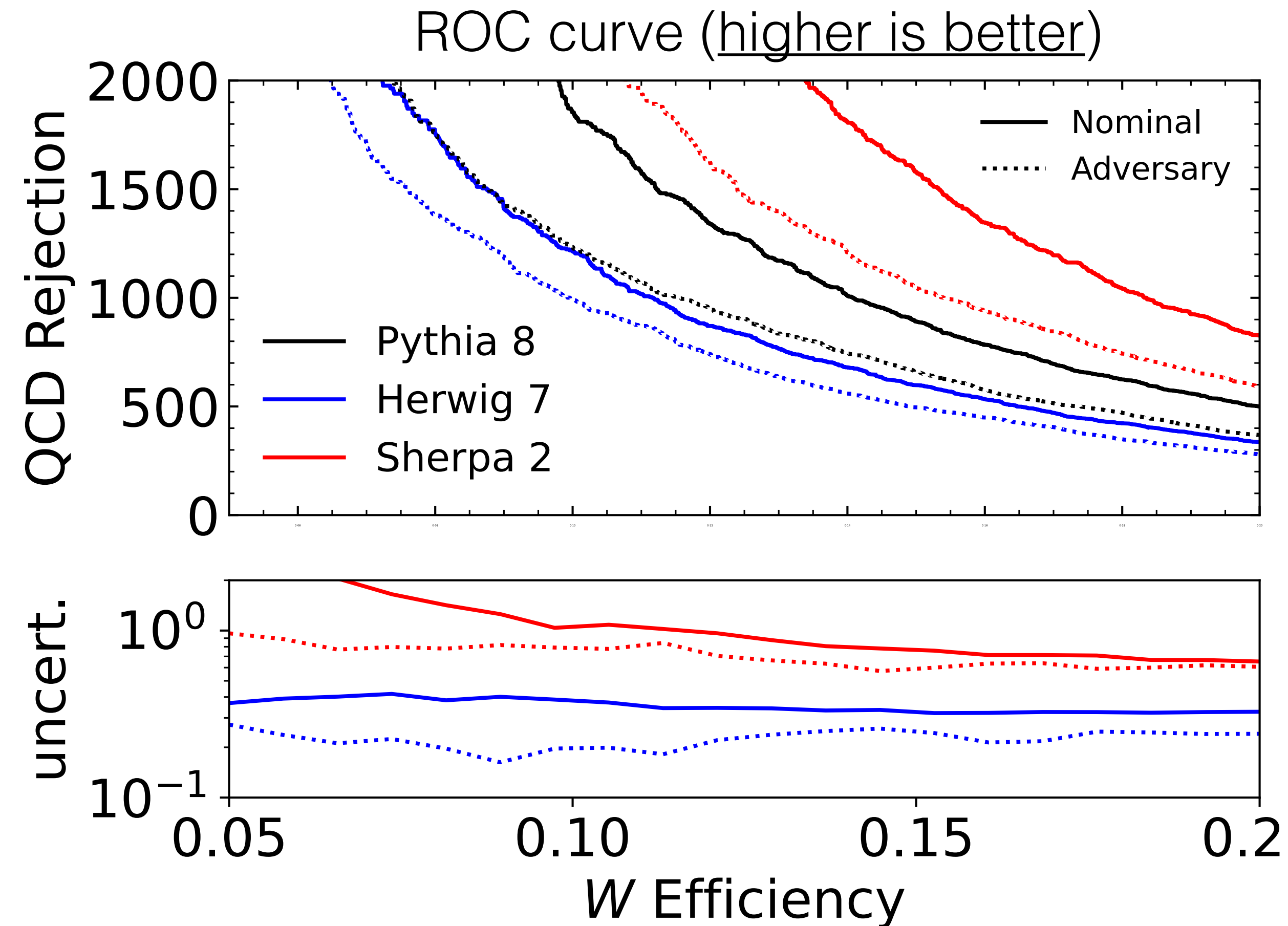

ROC curve (<u>higher is better</u>)

# Case Study 1: Two-point uncertainty - Result

Adversary successfully <u>sacrifices separation power</u> in order to reduce difference in performance between <span style="color:blue">Herwig</span> and `Pythia`

Cross-check with <span style="color:red">Sherpa</span> reveals <u>uncertainty severely underestimated</u> by usual <span style="color:blue">Herwig</span> vs `Pythia` comparison

In an typical LHC analysis, a cross-check with third generator rarely performed, similar to prior work suggesting decorrelation for theory uncertainties

ROC curve (<u>higher is better</u>)

Adversary successfully <u>sacrifices separation power</u> in order to reduce difference in performance between <span style="color:blue">Herwig</span> and Pythia

Cross-check with <span style="color:red">Sherpa</span> reveals <u>uncertainty severely underestimated</u> by usual <span style="color:blue">Herwig</span> vs Pythia comparison

In an typical LHC analysis, a cross-check with third generator rarely performed, similar to prior work suggesting decorrelation for theory uncertainties
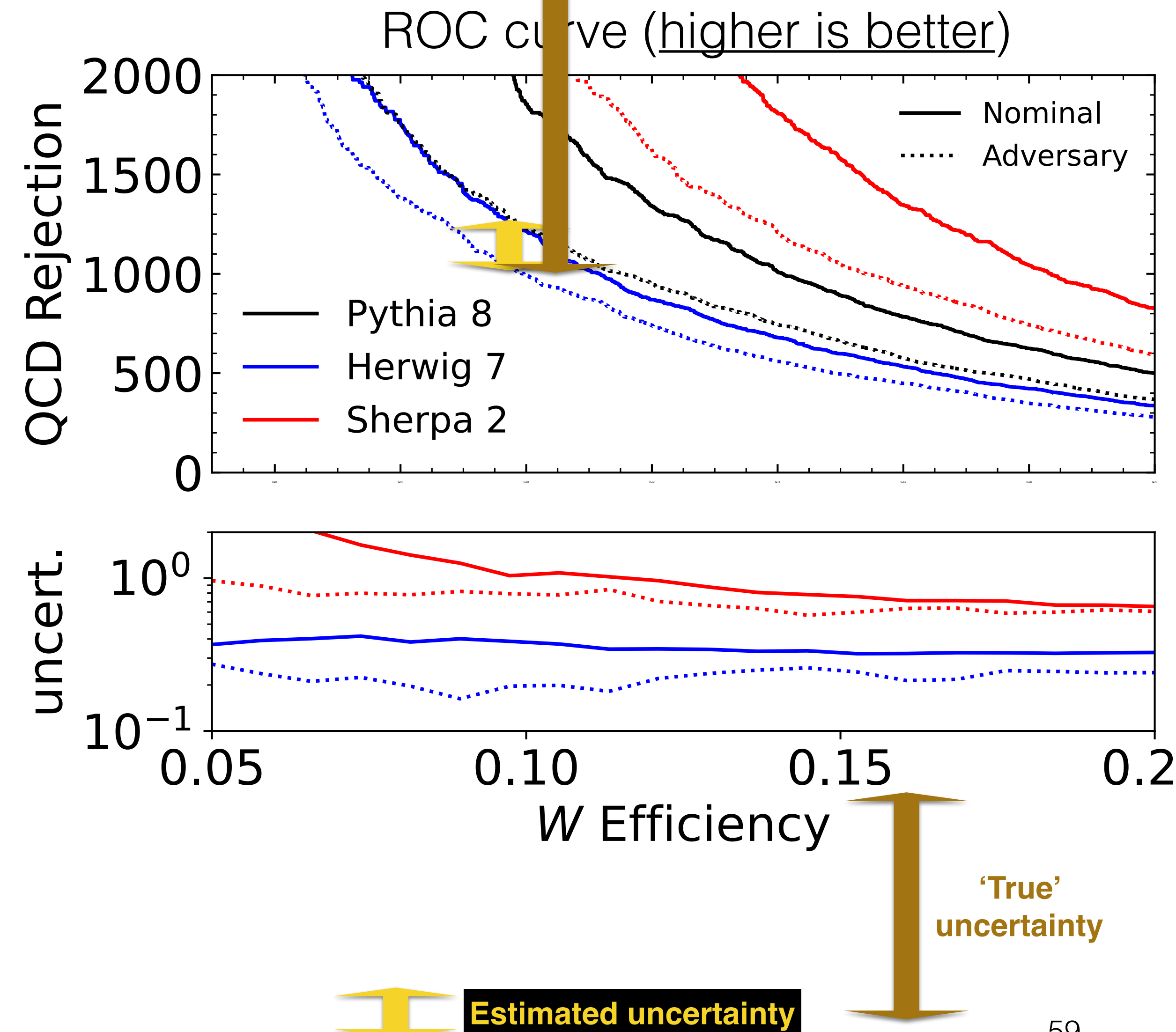
ROC curve (<u>higher is better</u>)
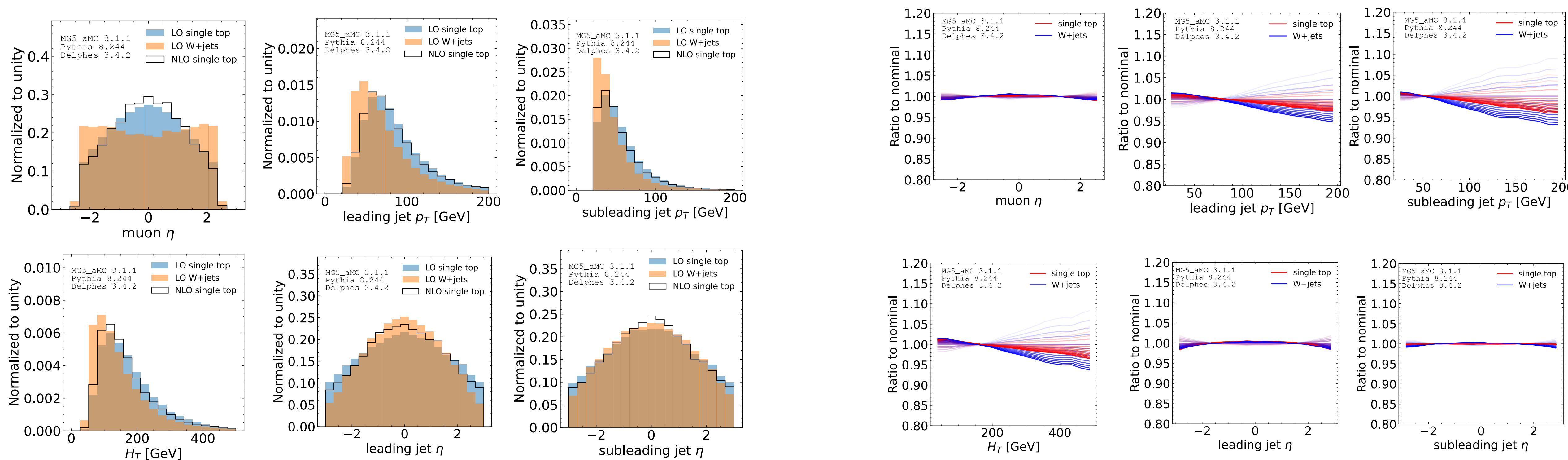
# Case Study 2: Higher-order corrections

- We can't calculate QFT to infinite order

- Artefact of truncation of series: Varying certain unphysical scales changes predictions

- Uncertainty quantification: Vary scales (renormalization scale, factorisation scale) between 1/2 to 2 in MC, see change in prediction

Goal: Single top vs W+Jets
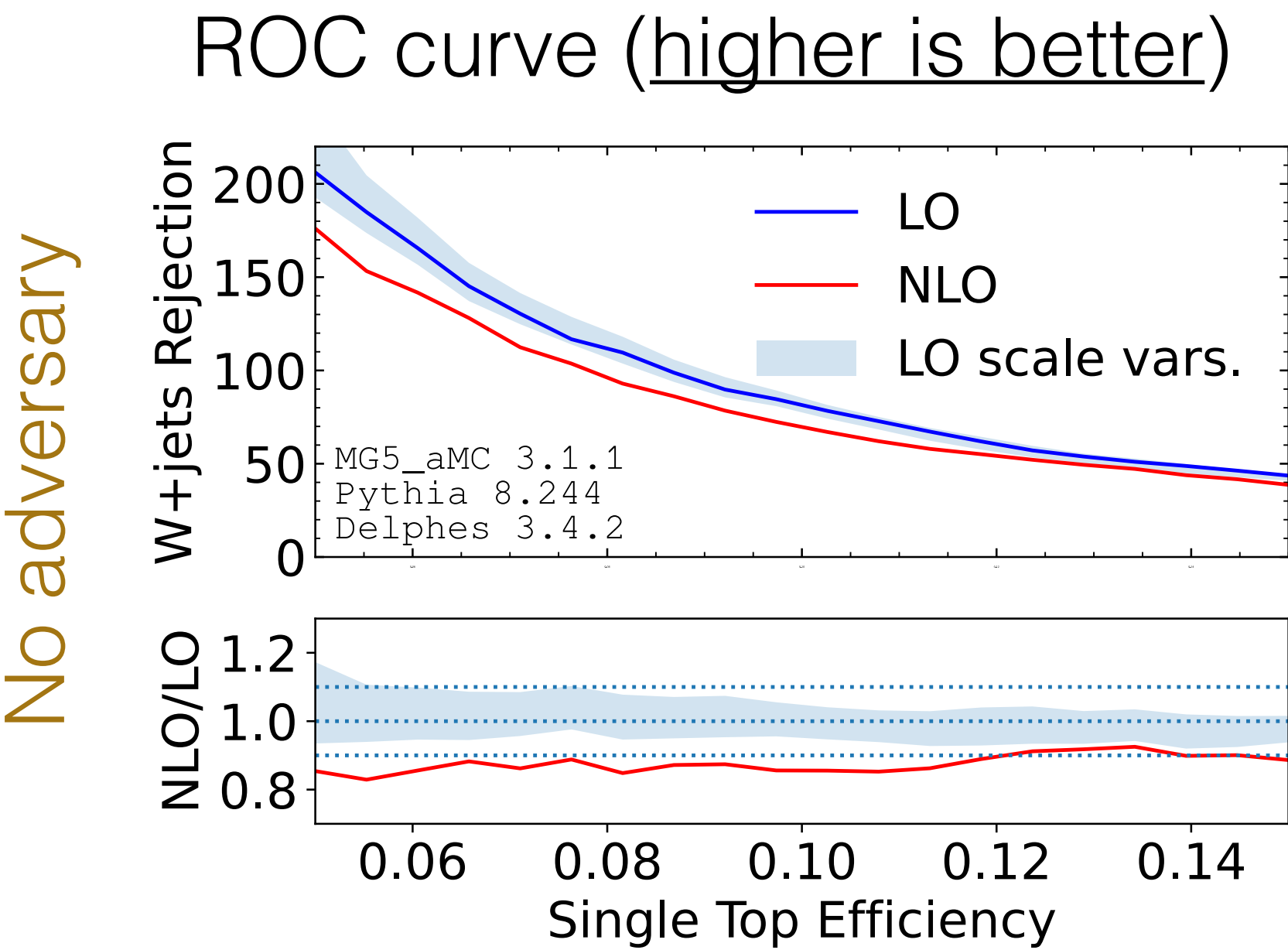Decorrelation: Reduce difference in performance on scale variations at LO
Cross-check: Test uncertainty estimate from {scale variations at LO} using NLO
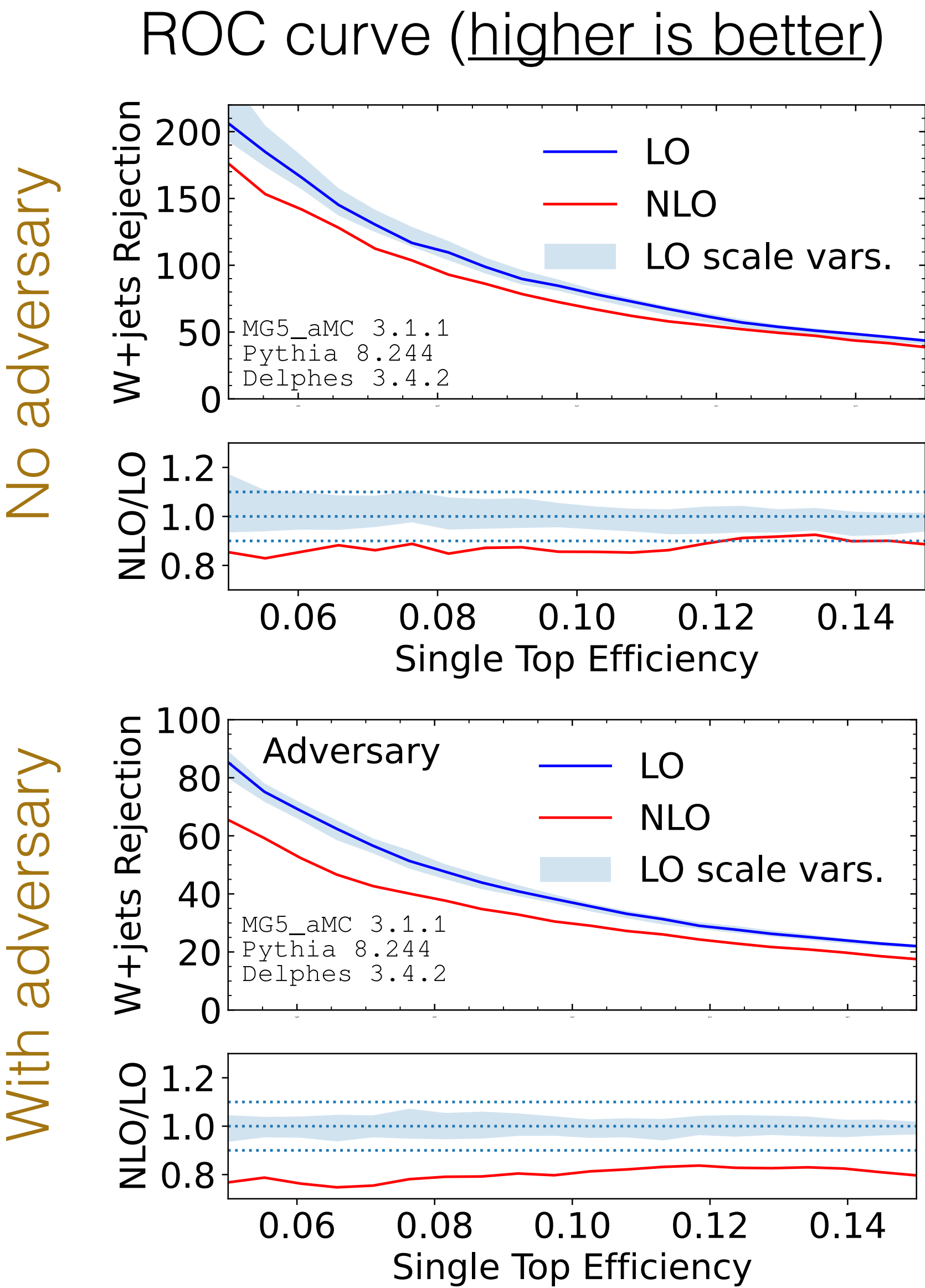


NLO vs LO

Factorisation scale variations going
from 1/2 to 2

ROC curve (<u>higher is better</u>)

No adversary



LO

NLO

LO scale vars.

MG5_aMC 3.1.1
Pythia 8.244
Delphes 3.4.2

W+jets Rejection

NLO/LO

Single Top Efficiency

## ROC curve (<u>higher is better</u>)

**No adversary**



W+jets Rejection

LO
NLO
LO scale vars.

MG5_aMC 3.1.1
Pythia 8.244
Delphes 3.4.2

NLO/LO

Single Top Efficiency

**With adversary**



W+jets Rejection

Adversary

LO
NLO
LO scale vars.

MG5_aMC 3.1.1
Pythia 8.244
Delphes 3.4.2

NLO/LO

Single Top Efficiency

ROC curve (<u>higher is better</u>)

Decorrelation:
Only the error bars
shrink, not the actual
distance to NLO

# Case Study 2: Continuous uncertainty - Result

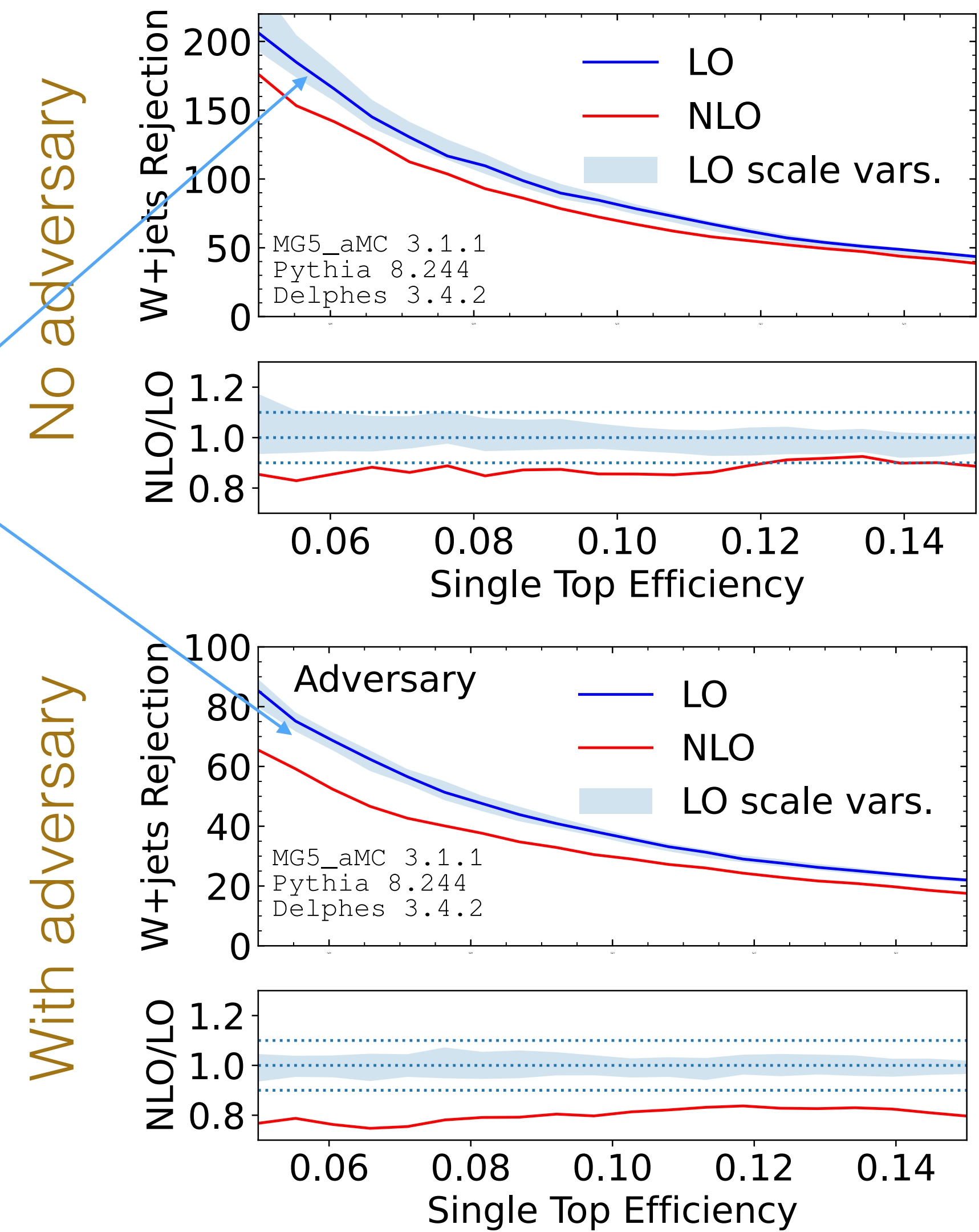Adversary successfully **sacrifices separation power** in order to reduce difference in performance between scale variations

Cross-check with NLO reveals **uncertainty severely underestimated** by decorrelation approach

In an typical LHC analysis, a cross-check with higher-order usually unavailable

Decorrelation:
Only the error bars shrink, not the actual distance to NLO
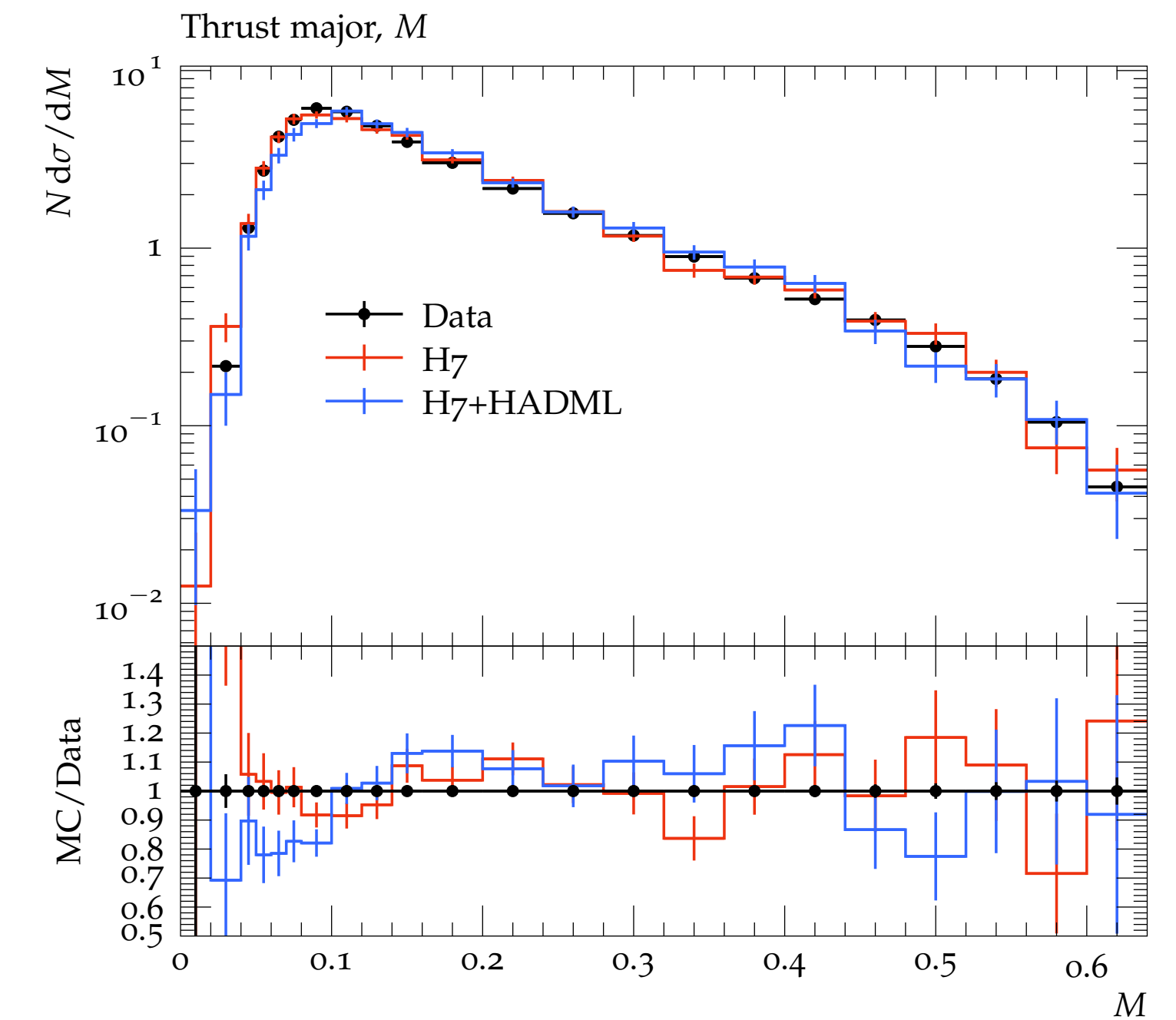
ROC curve (<u>higher is better</u>)

# Universe is a perfect simulator

# Universe is a perfect simulator

Bypass theory, can we learn hadronization directly from data ?

# Overconstraining NP

## Our modelling of NPs might be over-simplified

- If you assume one NP – chances are that your physics Likelihood will exploit this oversimplified JES model to overconstrain JES for high $p_T$ jets!



*Jet Energy Scale miscalibration*

5%

$\alpha_{JES}$

i.e. JES miscalibration is coherent for all jets
→ You can calibrate high $p_T$ jets with a low $p_T$ jet sample

*Jet $p_T$*



*Jet Energy Scale miscalibration*

5%   5%   5%   5%   5%

$\alpha_{JES1}$   $\alpha_{JES2}$   $\alpha_{JES3}$   $\alpha_{JES4}$   $\alpha_{JES5}$

i.e. JES miscalibration is not coherent across $p_T$ but still has 5% uncertainty for each $p_T$ bin

*Jet $p_T$*