# How to Unfold Top Decays

Luigi Favaro[1], Roman Kogler[2], Alexander Paasch[3],
Sofia Palacios Schweitzer[1], Tilman Plehn[1], Dennis Schwarz[4]
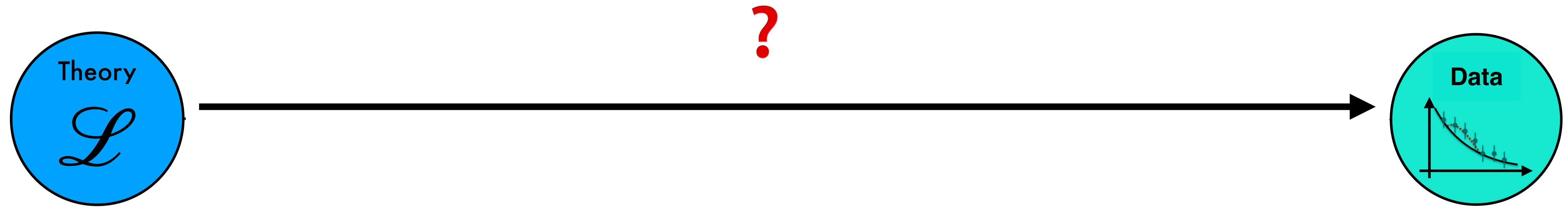
1 - Institut für theoretische Physik, Universität Heidelberg
2 - Deutsches Elektronen-Synchrotron DESY, Hamburg
3 - Institut für Experimentalphysik, Universität Hamburg
4 - Institut für Hochenergiephysik, Österreichische Akademie der Wissenschaft, Wien
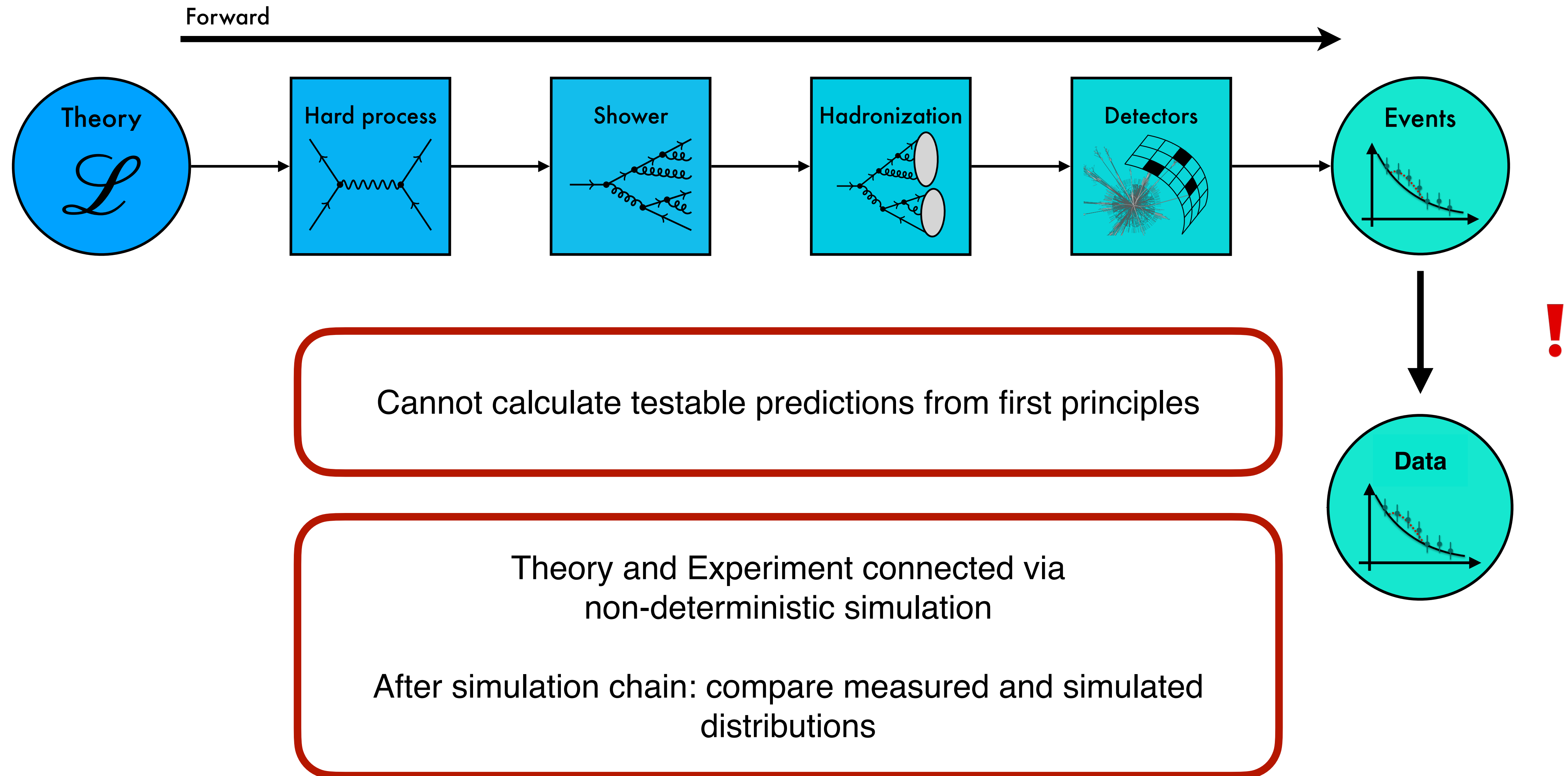
SPONSORED BY THE

Federal Ministry
of Education
and Research

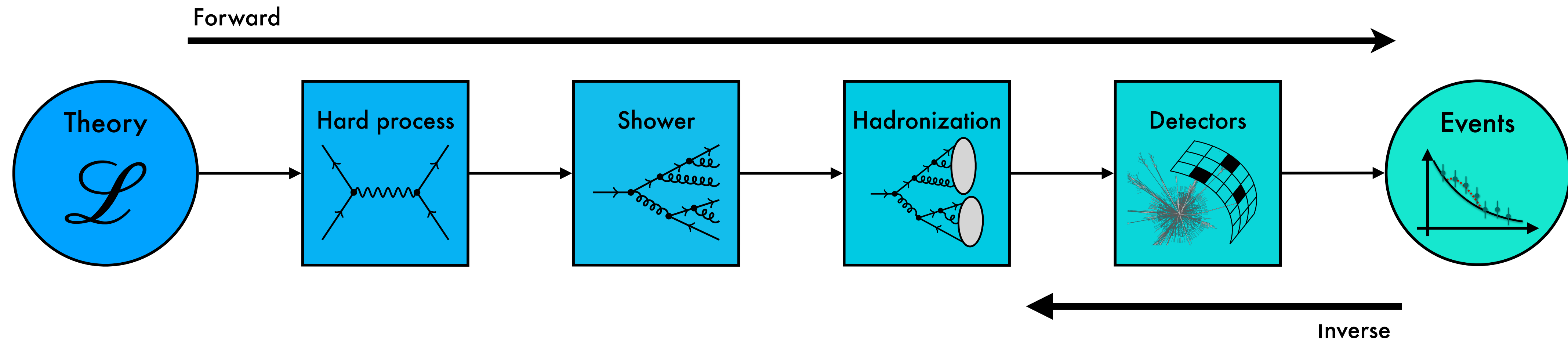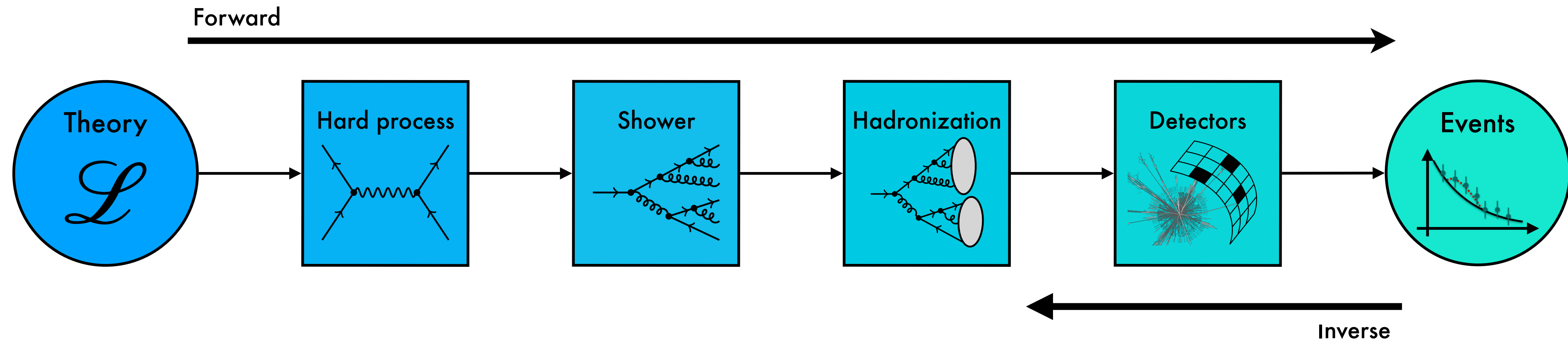Cannot calculate testable predictions from first principles

# Simulation Chain

Forward



Theory → Hard process → Shower → Hadronization → Detectors → Events → Data

Theory

Cannot calculate testable predictions from first principles

Theory and Experiment connected via non-deterministic simulation

After simulation chain: compare measured and simulated distributions

!

# Simulation Chain — Inversion



Forward

Theory $\mathcal{L}$ → Hard process → Shower → Hadronization → Detectors → Events

Inverse

Figure from A. Butter et al.: arXiv:2203.07460, R. Winterhalder

# Why unfolding?

Forward

Theory $\mathcal{L}$ → Hard process → Shower → Hadronization → Detectors → Events

Inverse

Theory analyses don't care about detectors

Comparing data from different experiments (Global Analysis)
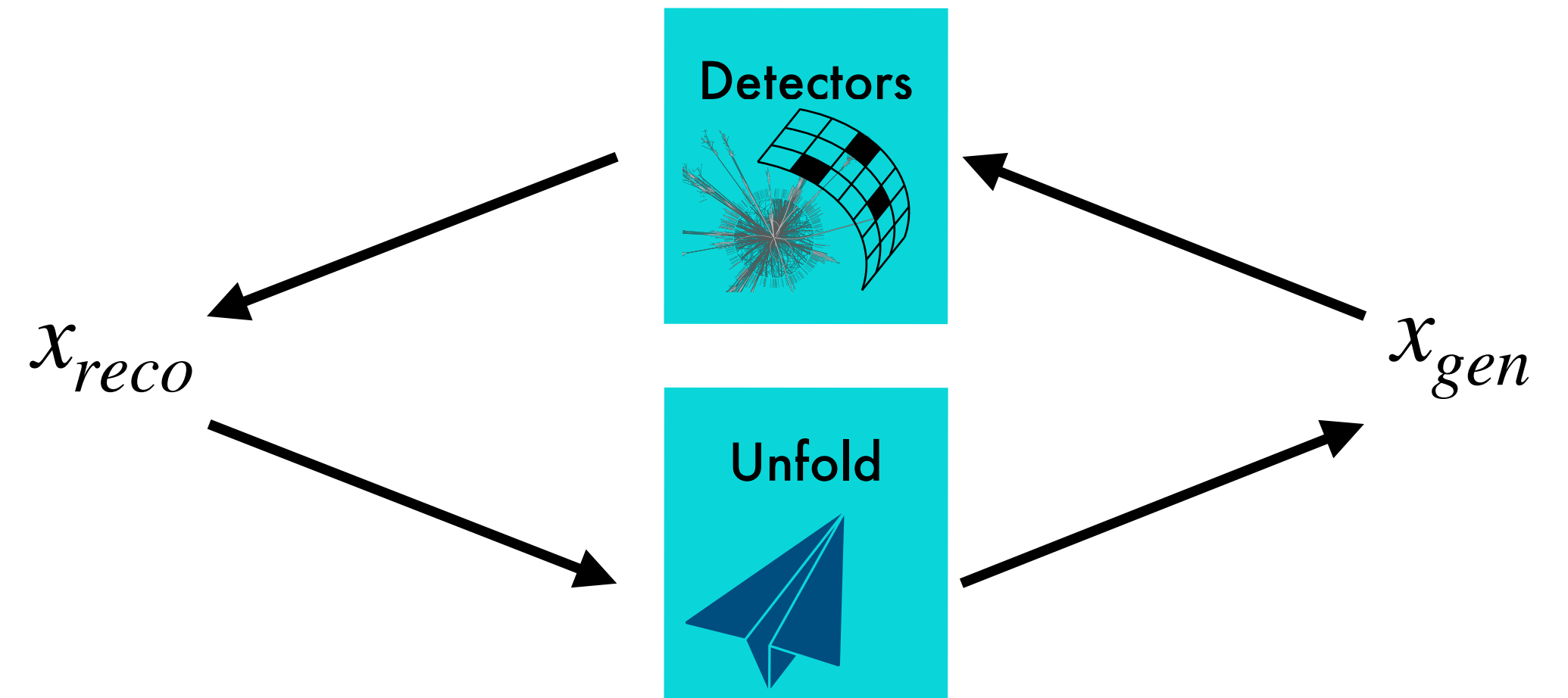
For some analysis direct access to theory parameters

Resolution

Data preservation

$$p(x_{reco}) = \int p(x_{gen})R(x_{reco}, x_{gen}) \, dx_{gen}$$

$$p(x_{gen}) = \int p(x_{reco})p(x_{gen}|x_{reco}) \, dx_{reco}$$
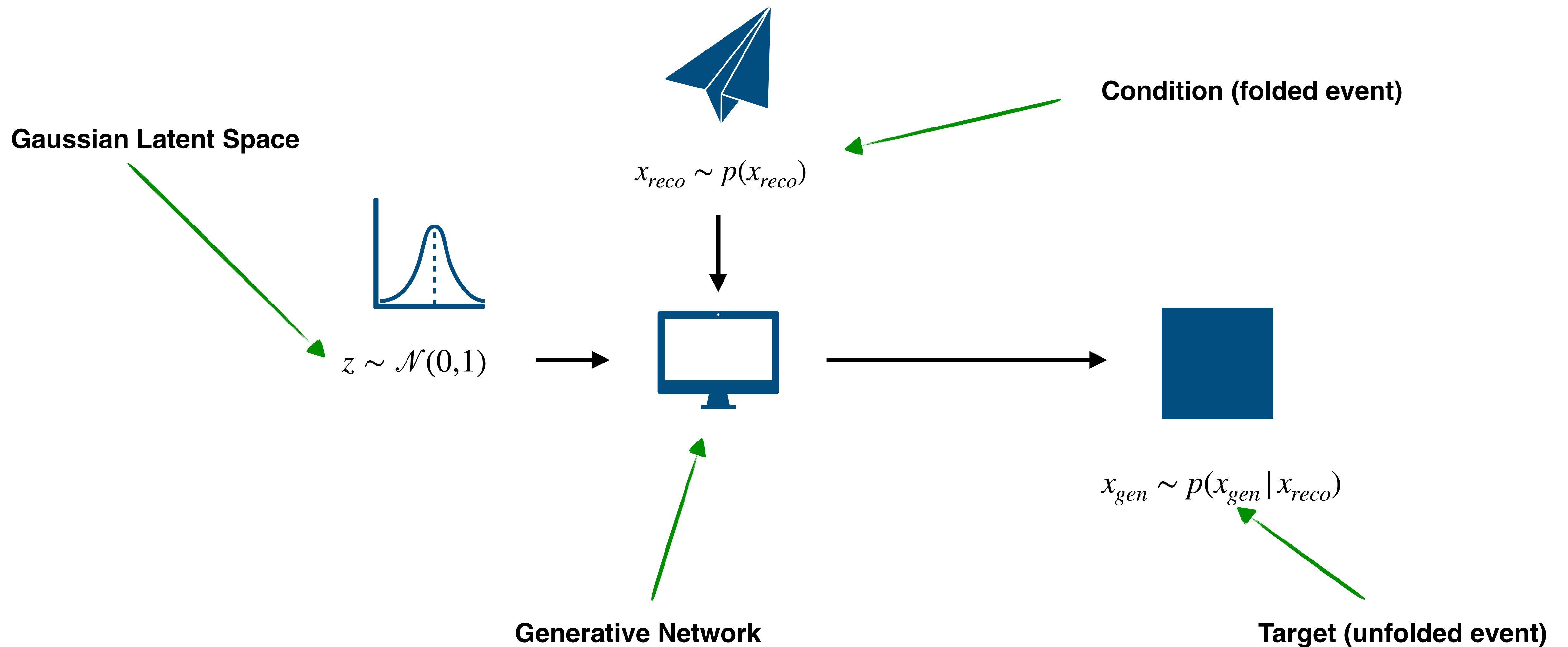
$$p(x_{reco}) = \int p(x_{gen}) \underbrace{R(x_{reco}, x_{gen})}_{\mathbf{p(x_{reco} \mid x_{gen})}} \, dx_{gen}$$

$$p(x_{gen}) = \int p(x_{reco}) p(x_{gen} \mid x_{reco}) \, dx_{reco}$$

Classical methods are restricted to binned, one-dimensional distributions

We would like to learn high-dimensional, unbinned unfolding probability

7

$$p(x_{reco}) = \int p(x_{gen}) \underbrace{R(x_{reco}, x_{gen})}_{\mathbf{p(x_{reco} | x_{gen})}} \, dx_{gen}$$

$$p(x_{gen}) = \int p(x_{reco}) \underbrace{p(x_{gen} | x_{reco})}_{\textbf{target probability}} \, dx_{reco}$$

Classical methods are restricted to binned, one-dimensional distributions

We would like to learn high-dimensional, unbinned unfolding probability

**Condition (folded event)**

**Gaussian Latent Space**

$$x_{reco} \sim p(x_{reco})$$

$$z \sim \mathcal{N}(0,1)$$

$$x_{gen} \sim p(x_{gen} | x_{reco})$$

**Generative Network**

**Target (unfolded event)**

9

**Goal:** learn transformation latent → gen phase space conditioned on reco event

During training, use paired events of forward simulation

After training, repeated sampling from latent space with constant condition allows probabilistic single event unfolding

$x_{rec} \sim p(x_{reco})$

$z \sim \mathcal{N}(0,1)$

$x_{gen} \sim p(x_{gen}|x_{reco})$

**Goal:** learn transformation latent → gen phase space conditioned on reco event

During training, use paired events of forward simulation

After training, repeated sampling from latent space with constant condition allows probabilistic single event unfolding

$x_{rec} \sim p(x_{reco})$

$z \sim \mathcal{N}(0,1)$

$x_{gen} \sim p(x_{gen}|x_{reco})$

Bellagente et al. **1912.00477**
Bellagente et al. **2006.06685**
Backes et al. **2212.08674**
Huetsch et al. **2404.18807**

$p_{T,J} > 400$ GeV

Reconstruct triple jet mass
$M_{jjj}$ to measure $m_t$

Tag side

$p_{T,J} > 400$ GeV



Reconstruct triple jet mass $M_{jjj}$ to measure $m_t$

Tag side

Previously done in CMS with TUnfold (classical binned unfolding algorithm)



CMS **2211.01456**

$p_{T,J} > 400$ GeV

Reconstruct triple jet mass
$M_{jjj}$ to measure $m_t$

Tag side

Previously done in CMS with TUnfold
(classical binned unfolding algorithm)



CMS **2211.01456**

<u>BUT</u> leading uncertainty: choice of $m_t$ in simulation + no access to full phase space

→ **Could generative unfolding help?**

**1. Mulitresonant phase space**

## 1. Mulitresonant phase space
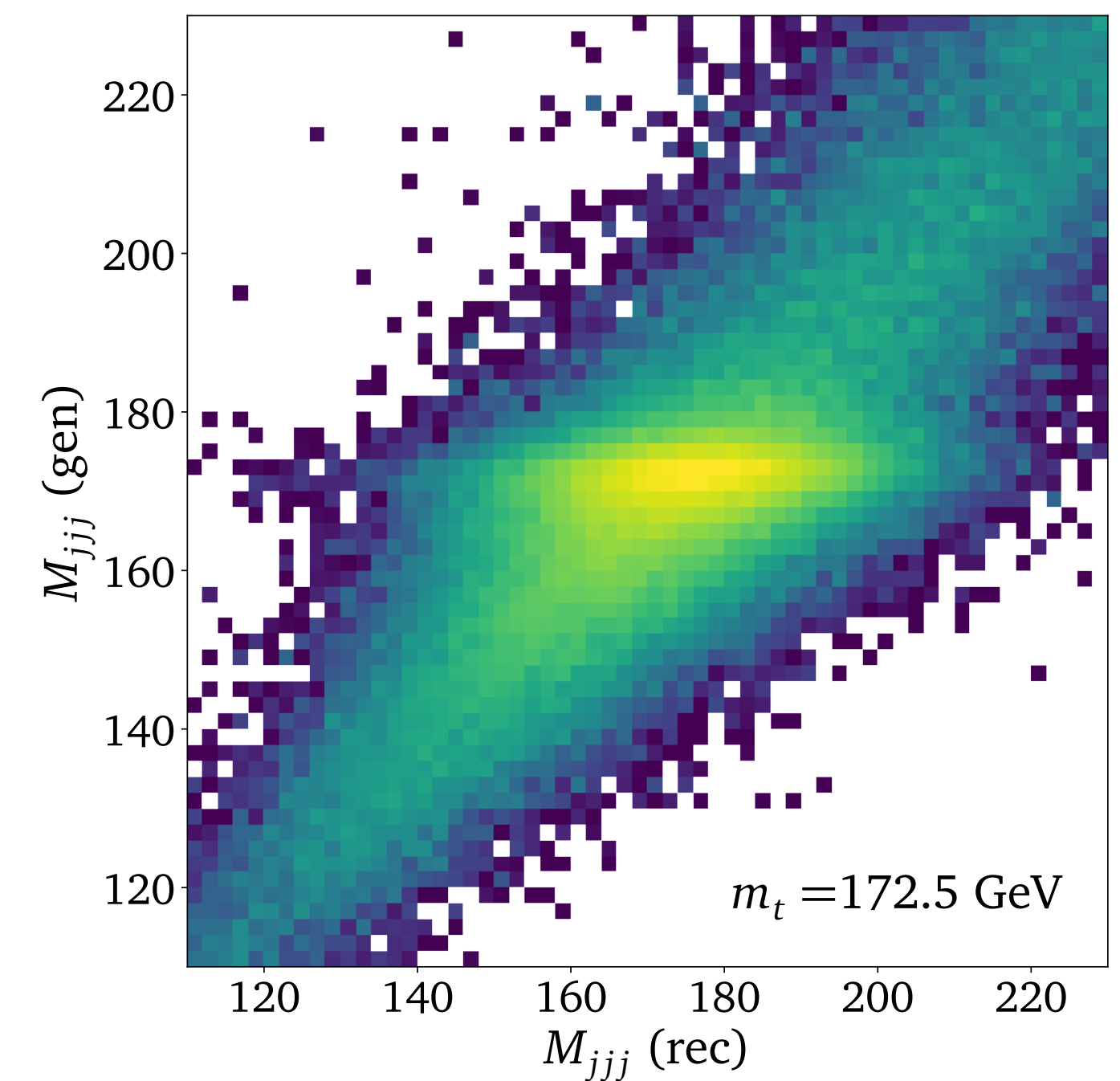


## 2. Combinatorics

# Challenging aspects of top - unfolding

**1. Mulitresonant phase space**



**2. Combinatorics**



**3. Detector Smearing**

# Choosing the right parametrization

**1. The naive**

Reco and gen level difference not significantly visible, only in correlations

$$p_1 = (E_1, \overrightarrow{p}_1)$$

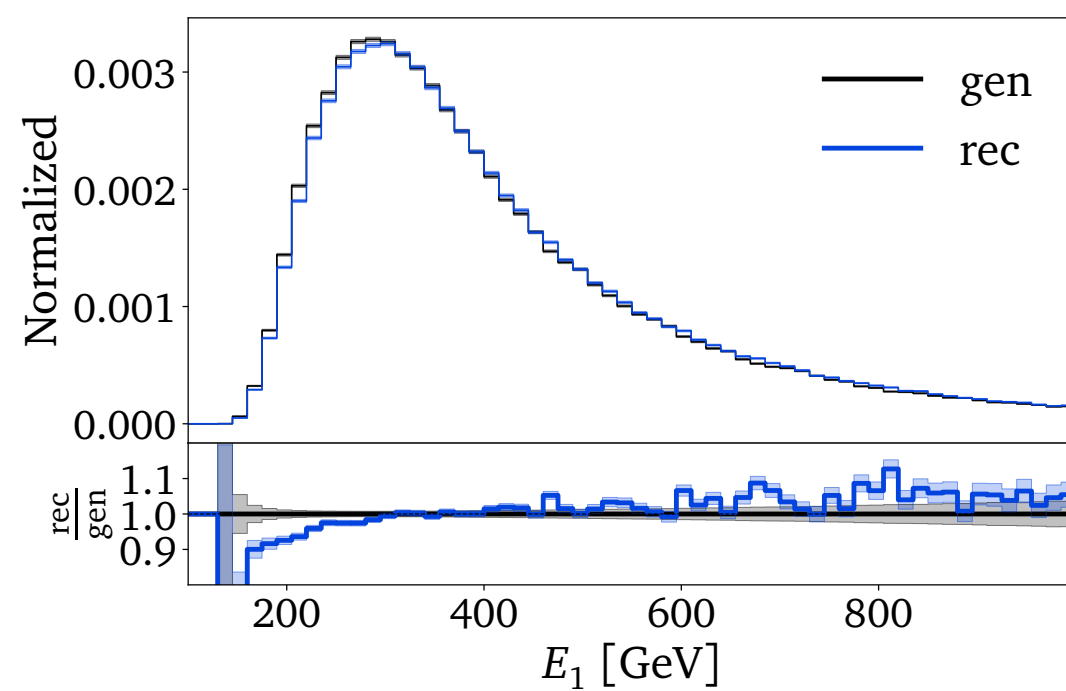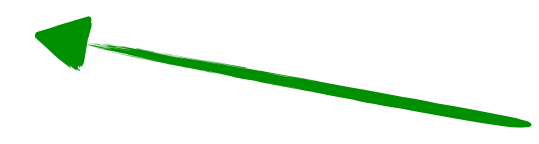$$p_2 = (E_2, \overrightarrow{p}_2)$$

$$p_3 = (E_3, \overrightarrow{p}_3)$$

$$M_{jjj}(p_1, p_2, p_3)$$

$$M_{ij}(p_i, p_j)$$

12 dimensional correlation

8 dimensional correlation + combinatorics difficult

**2. The less naive**

$$p_1 = (p_{T,1}, \phi_1, \eta_1, m_1)$$

$$p_2 = (p_{T,2}, \phi_2, \eta_2, m_2)$$

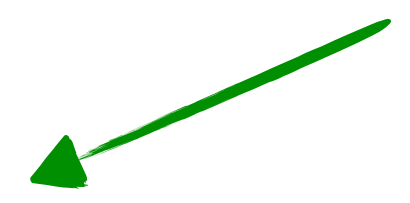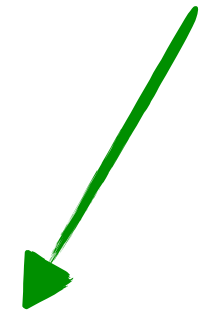$$p_3 = (p_{T,3}, \phi_3, \eta_3, m_3)$$

$$M_{jjj}(p_1, p_2, p_3)$$

$$M_{ij}(p_i, p_j)$$

12 dimensional correlation

8 dimensional correlation + combinatorics difficult
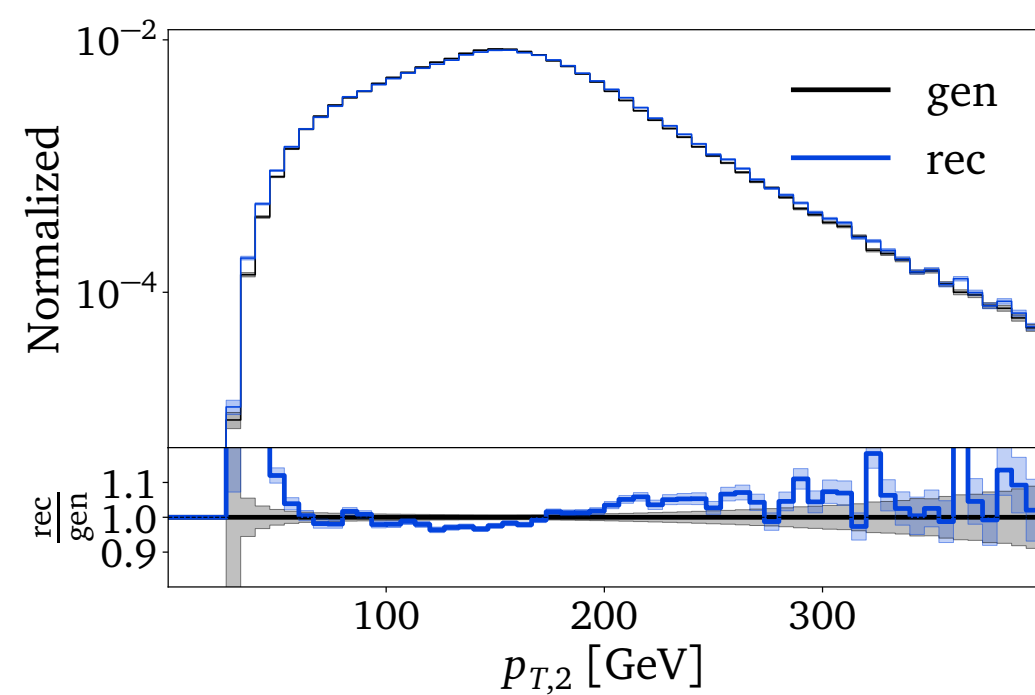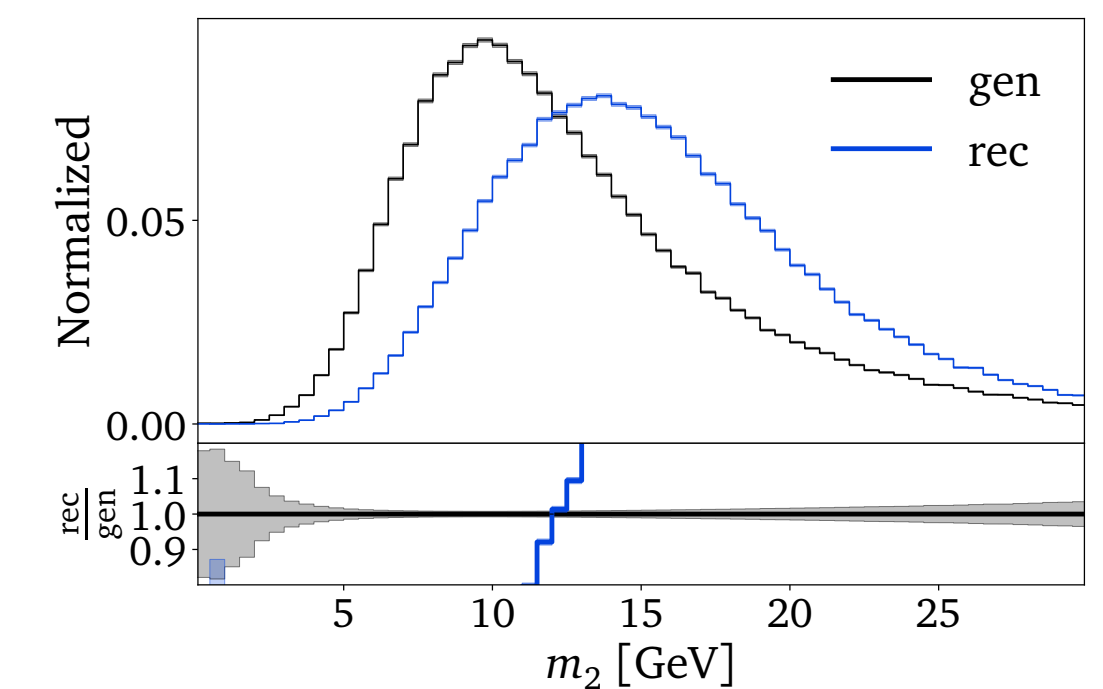


Reco and gen level difference visible

**3. The least naive**

$$p_1 = (p_{T,1}, M_{12}, \eta_1, m_1)$$

$$p_2 = (p_{T,2}, M_{23}, \eta_2, m_2)$$

$$p_3 = (p_{T,3}, M_{13}, \eta_3, m_3)$$

$$M_{jjj}^2 = \sum_{ij,i>j} M_{ij}^2 - \sum_{i} m_i^2$$

6 dimensional correlation

Direct input + combinatorics simple



Reco and gen level difference visible

**3. The least naive**

$$|\Delta\phi_{ij}|$$
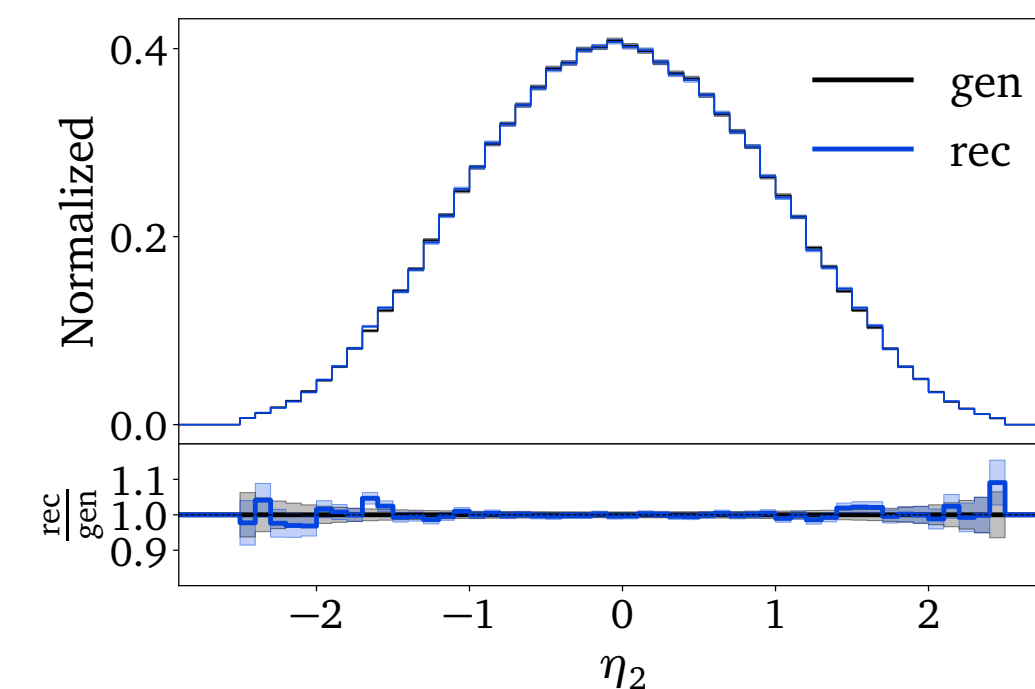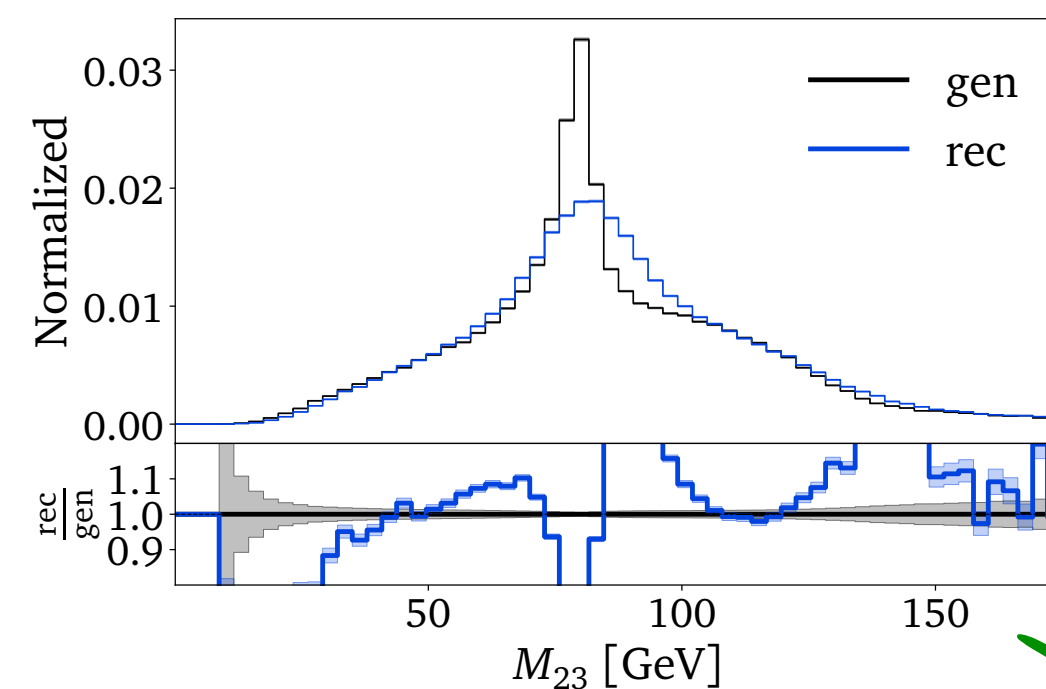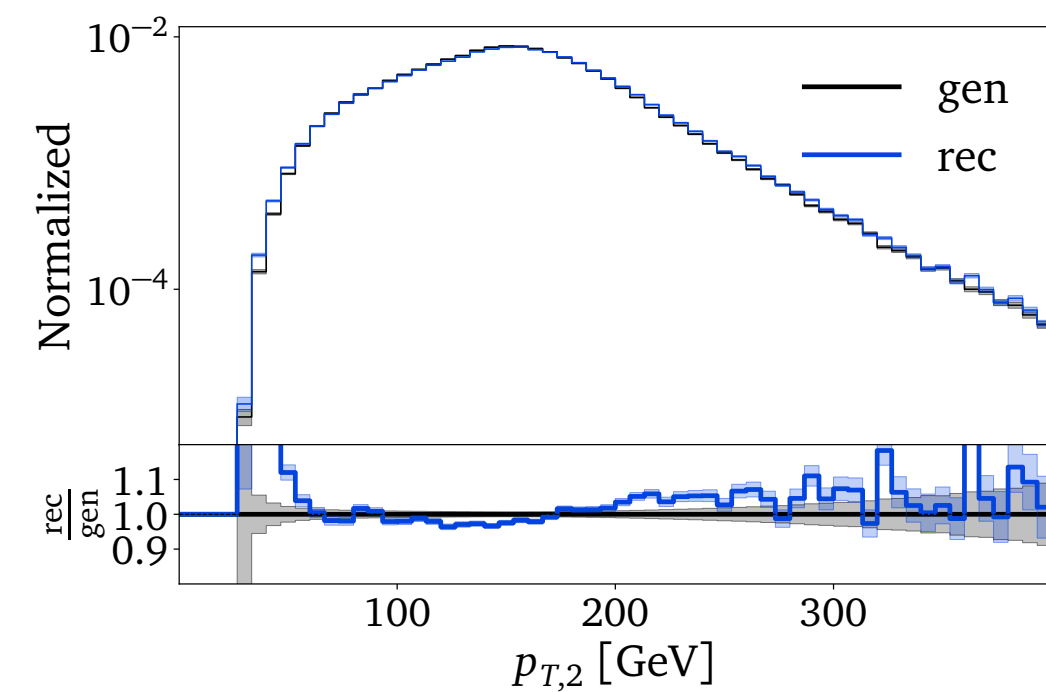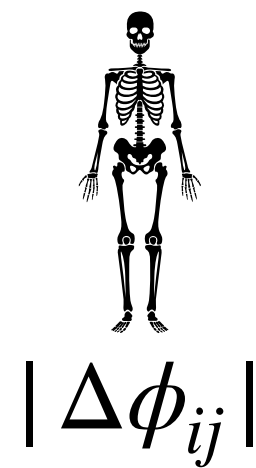
$$
\left.
\begin{aligned}
p_1 &= (p_{T,1}, M_{12}, \eta_1, m_1) \\
p_2 &= (p_{T,2}, M_{23}, \eta_2, m_2) \\
p_3 &= (p_{T,3}, M_{13}, \eta_3, m_3)
\end{aligned}
\right\}
\qquad
M_{jjj}^2 = \sum_{ij,i>j} M_{ij}^2 - \sum_i m_i^2
$$



21

**3. The least naive**

$$p_1 = (p_{T,1}, M_{12}, \eta_1, m_1)$$
$$p_2 = (p_{T,2}, M_{23}, \eta_2, m_2)$$
$$p_3 = (p_{T,3}, M_{13}, \eta_3, m_3)$$

$$M_{jjj}^2 = \sum_{ij,i>j} M_{ij}^2 - \sum_i m_i^2$$

**3. The least naive**

$$p_1 = (p_{T,1}, M_{12}, \eta_1, m_1)$$
$$p_2 = (p_{T,2}, M_{23}, \eta_2, m_2)$$
$$p_3 = (p_{T,3}, M_{13}, \eta_3, m_3)$$
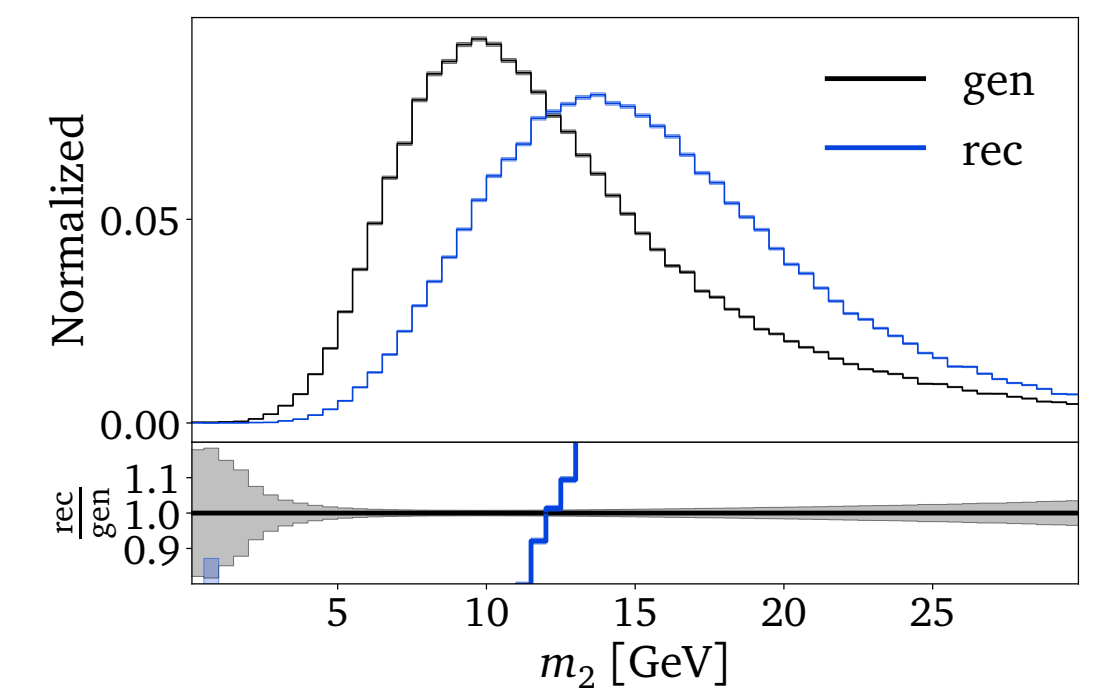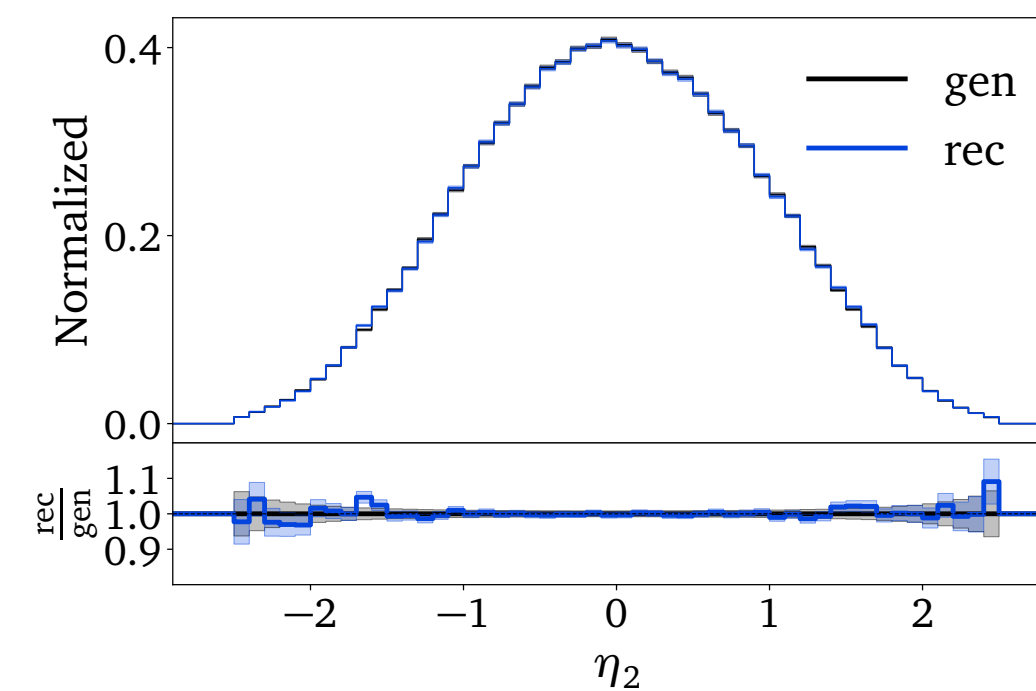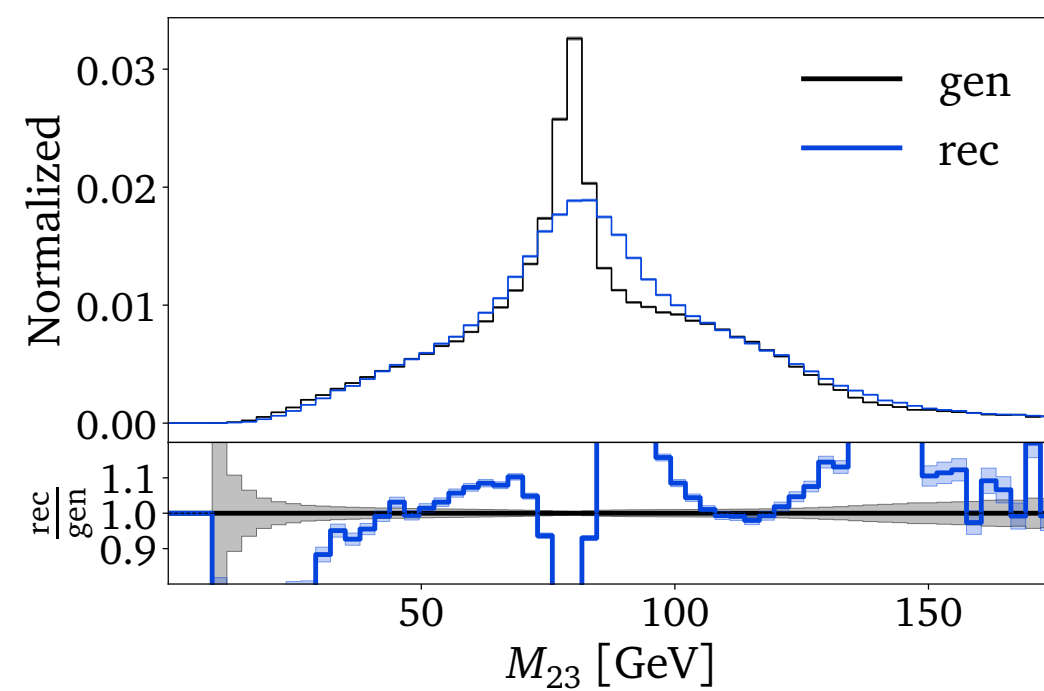
$$M_{jjj}^2 = \sum_{ij, i>j} M_{ij}^2 - \sum_i m_i^2$$

For mass measurement, we only use
6 dimensional subset of phase space
to increase network performance

Correct migration learned?

Train with full CMS simulation with
$$m_t = 172.5 \text{ GeV}$$

Unfolded distribution of triple jet mass within
$\mathcal{O}(1\%)$ of truth gen level

BUT: Test data also simulation with
$$m_t = 172.5 \text{ GeV}$$

Train with full CMS simulation with
$m_t = 172.5$ GeV

Unfolded distribution of triple jet mass within
$\mathcal{O}(1\%)$ of truth gen level

BUT: Test data also simulation with
$m_t = 172.5$ GeV

For pseudo-data with different top masses :
Algorithm falls back to prior ($m_t = 172.5$ GeV)



26

Train with full CMS simulation with
$m_t = 172.5$ GeV

Unfolded distribution of triple jet mass within
$\mathcal{O}(1\%)$ of truth gen level

BUT: Test data also simulation with
$m_t = 172.5$ GeV

For pseudo-data with different top masses :
Algorithm falls back to prior ($m_t = 172.5$
GeV)



27

$$p_{sim}(x_{gen} \mid m_s)$$

$$p_{unfold}(x_{gen} \mid m_s, m_d)$$

$$p(x_{reco} \mid x_{gen})$$

$$p_{model}(x_{gen} \mid x_{reco}, m_s)$$

correspondance

$$p_{sim}(x_{reco} \mid m_s)$$

$$p_{data}(x_{reco} \mid m_d)$$

$$p_{sim}(x_{gen}|m_s) \qquad\qquad\qquad p_{unfold}(x_{gen}|m_s, \cancel{m_d})$$

$$p(x_{reco}|x_{gen}) \qquad\qquad\qquad p_{model}(x_{gen}|x_{reco}, m_s)$$

correspondance

$$p_{sim}(x_{reco}|m_s) \longleftrightarrow p_{data}(x_{reco}|m_d)$$

$\rightarrow$ **Solution: Strengthen $m_d$ dependence, but how?**

$$p_{sim}(x_{gen} | m_s)$$

$$p_{unfold}(x_{gen} | m_s, \cancel{m_d})$$

$$p(x_{reco} | x_{gen})$$

$$p_{model}(x_{gen} | x_{reco}, m_s)$$

correspondance

$$p_{sim}(x_{reco} | m_s)$$

$$p_{data}(x_{reco} | m_d)$$

$\rightarrow$ **Solution: Strengthen $m_d$ dependence, but how?**

1. Augment training data with simulation from different top masses
2. Estimate batch-wise $m_d \approx$ weighted-median($M_{jjj}^{batch}$) on reco level

Train with full CMS simulation with
$m_t = [172.5 \text{ GeV}, 169.5 \text{ GeV}, 175.5 \text{ GeV}]$

Test by unfolding simulation with
$m_t = 171.5 \text{ GeV}$ & $173.5 \text{ GeV}$

Unfolded distribution of triple jet mass within $\mathcal{O}(1\%)$ of truth gen level **without** bias

Train with full CMS simulation with
$m_t = [172.5 \text{ GeV}, 169.5 \text{ GeV}, 175.5 \text{ GeV}]$

Test by unfolding simulation with
$m_t = 171.5 \text{ GeV} \& 173.5 \text{ GeV}$

Unfolded distribution of triple jet mass within
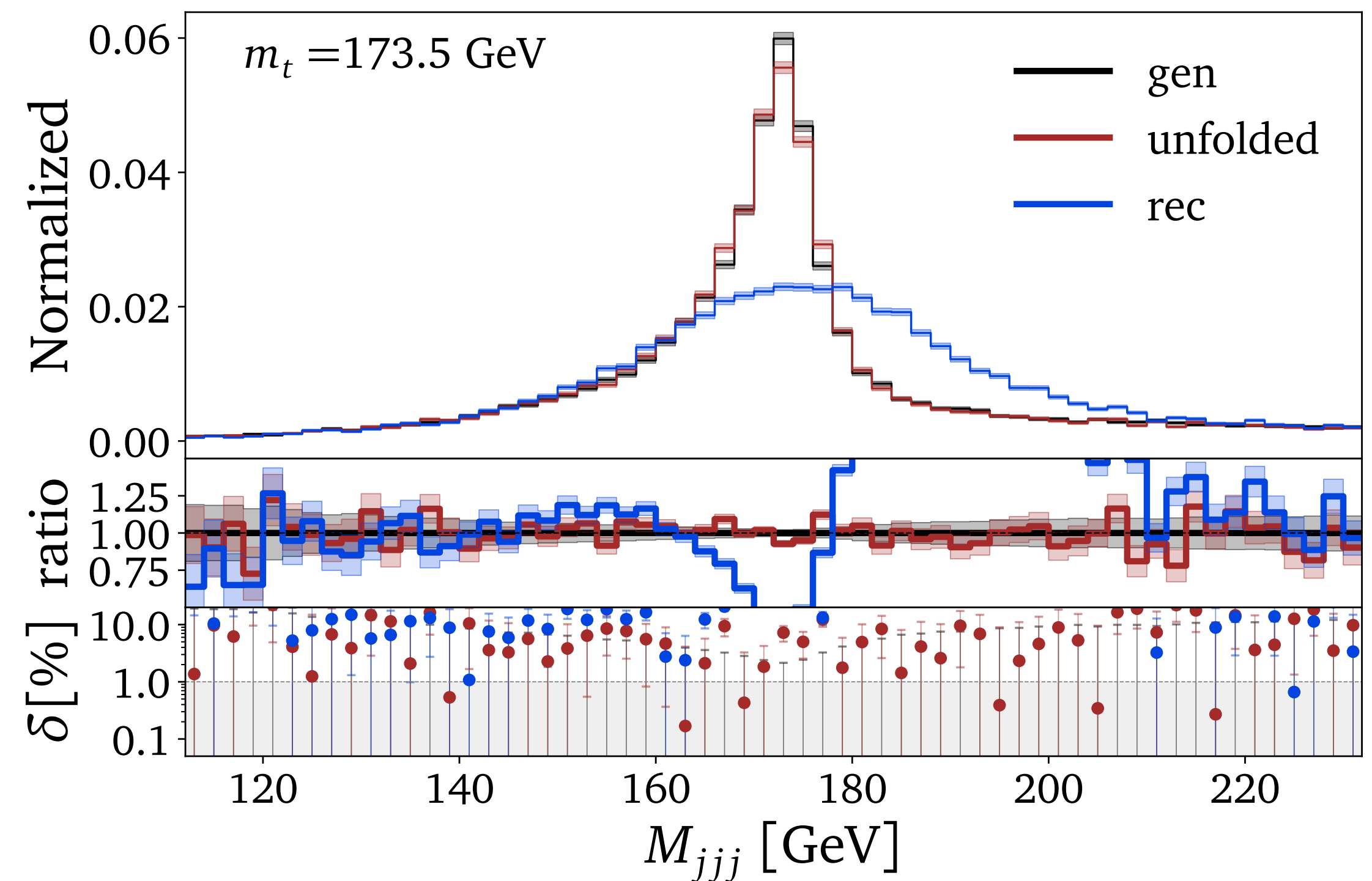$\mathcal{O}(1\%)$ of truth gen level **without** bias



⚠️ ML task becomes much harder

For a fixed top mass:

Choose subset of test data of 41000 reco level events

Unfolded 1000 bootstrapped replicas

Estimate covariance matrix and mean by 1000 different unfolded distributions



- CFM, 4d, 5 bins, $m_t = 172.53^{+0.26}_{-0.26}$ GeV
- CFM (stat. only), 4d, 5 bins, $m_t = 172.58^{+0.25}_{-0.25}$ GeV
- CFM, 6d, 5 bins, $m_t = 172.49^{+0.27}_{-0.27}$ GeV
- CFM (stat. only), 6d, 5 bins, $m_t = 172.60^{+0.25}_{-0.25}$ GeV
- TUnfold, $m_t = 172.50^{+0.32}_{-0.32}$ GeV
- TUnfold (stat. only), $m_t = 172.51^{+0.21}_{-0.21}$ GeV

For a fixed top mass:

Choose subset of test data of 41000 reco level events

Unfolded 1000 bootstrapped replicas

Estimate covariance matrix and mean by 1000 different unfolded distributions



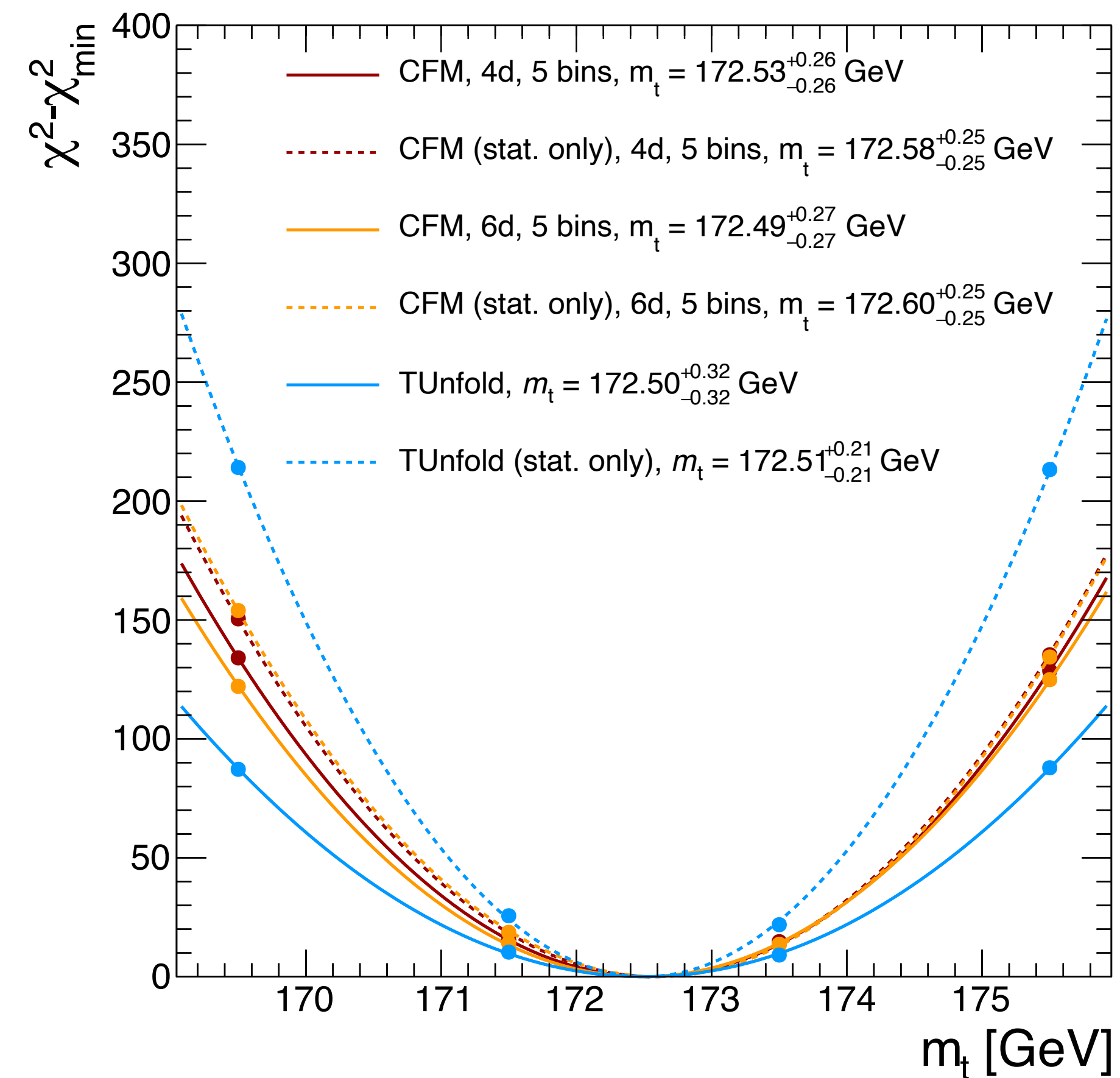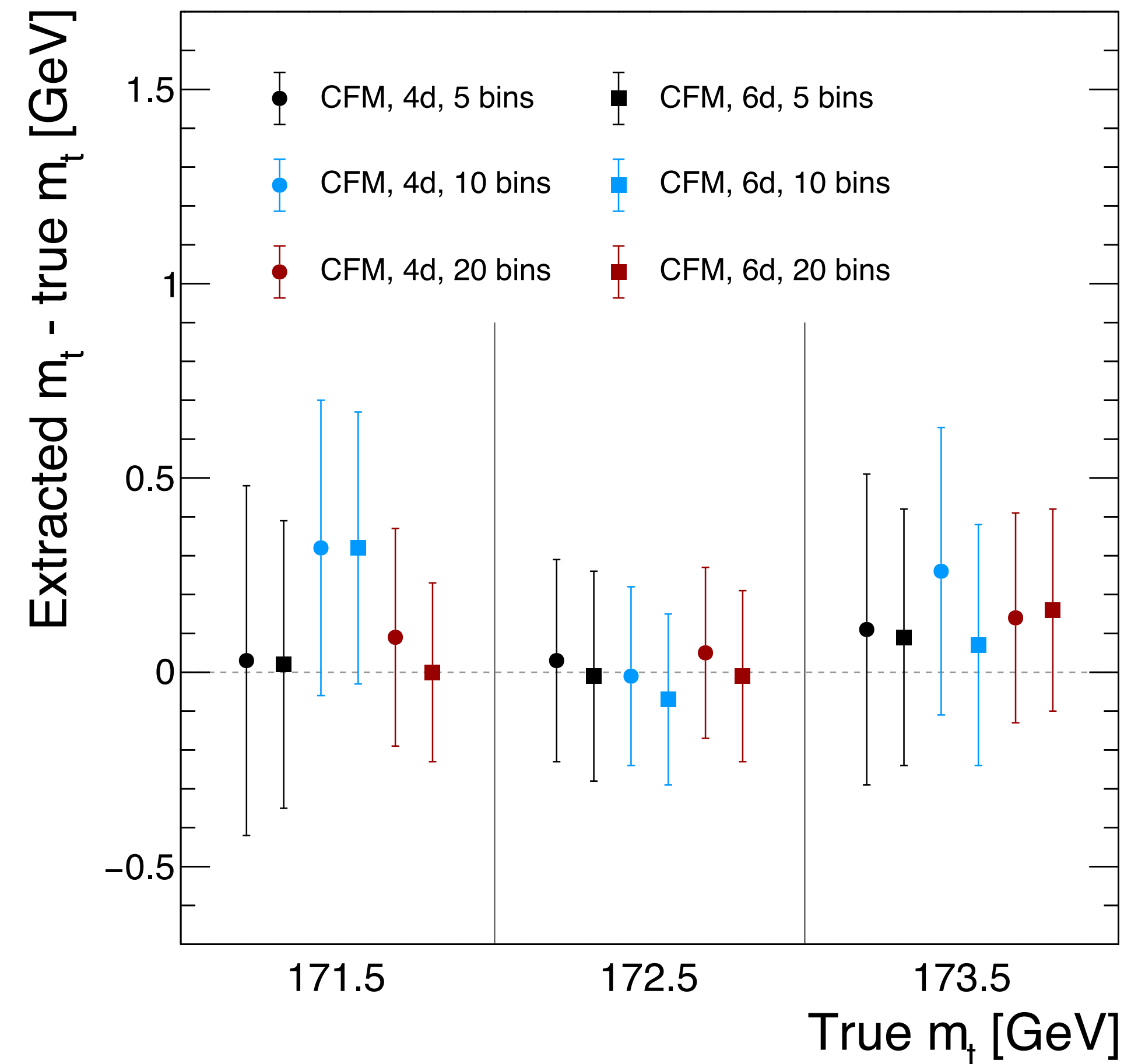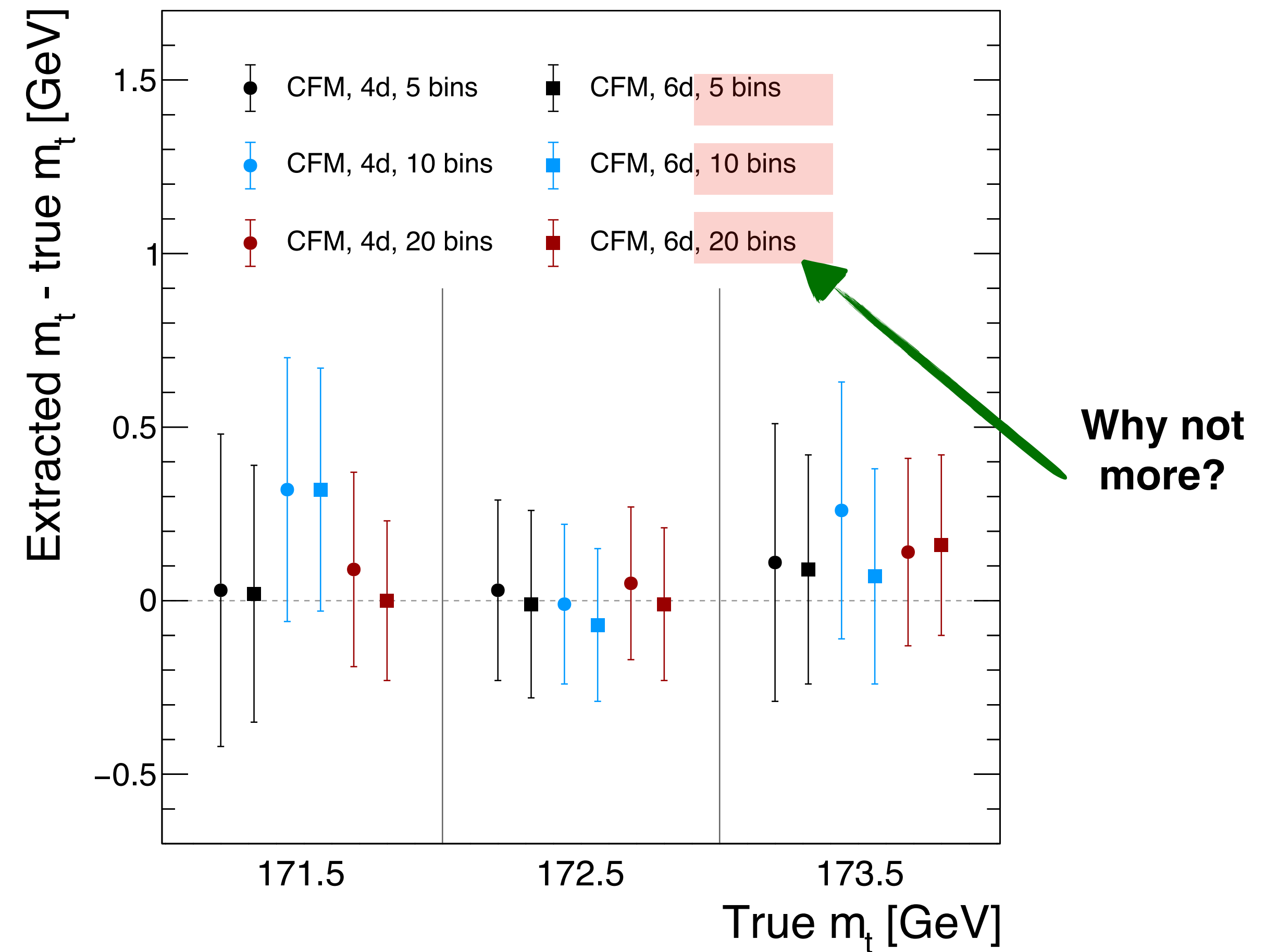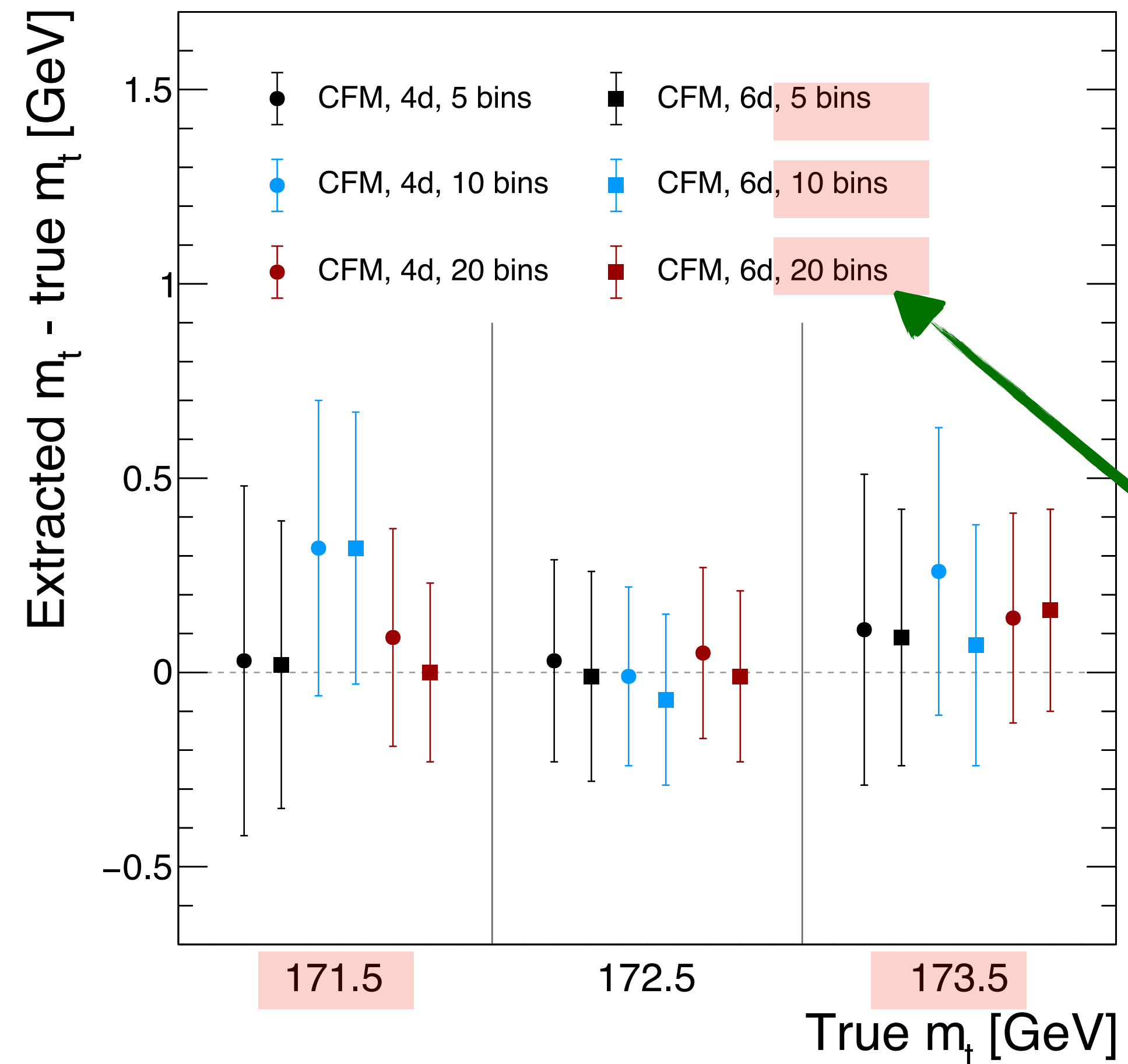→ Reliably unfold triple jet mass without bias

# Mass Measurement

For a fixed top mass:

Choose subset of test data of 41000 reco level events

Unfolded 1000 bootstrapped replicas

Estimate covariance matrix and mean by 1000 different unfolded distributions



**Why not more?**

→ Reliably unfold triple jet mass without bias

For a fixed top mass:

Choose subset of test data of 41000 reco level events

Unfolded 1000 bootstrapped replicas

Estimate covariance matrix and mean by 1000 different unfolded distributions
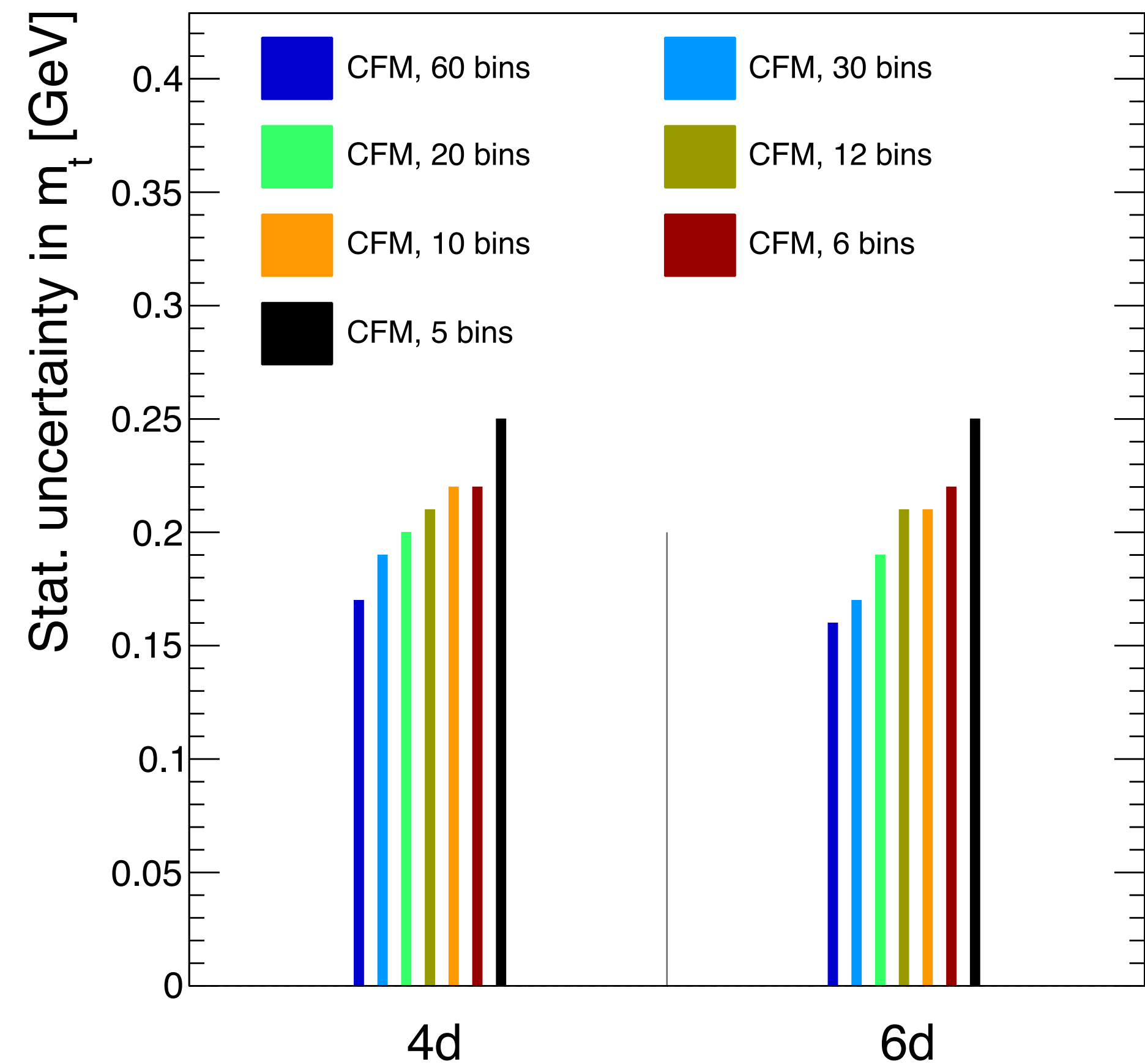


**Why not more?**

Limited by discrete grid of available $m_t$ simulations

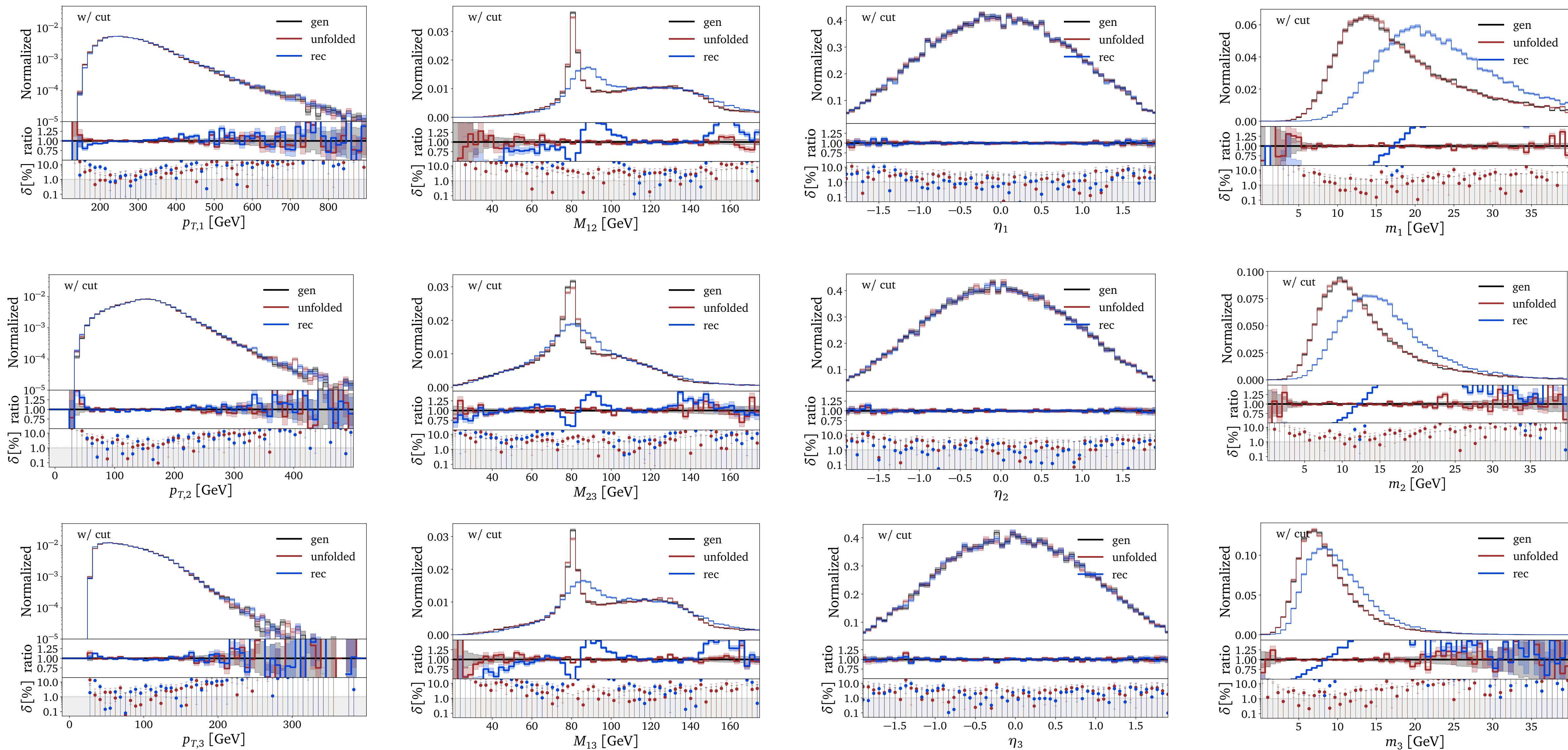→ Reliably unfold triple jet mass without bias

# Mass Measurement

For $m_t = 172.5$ GeV, we have a close grid of available simulations ($\pm 1$ GeV)

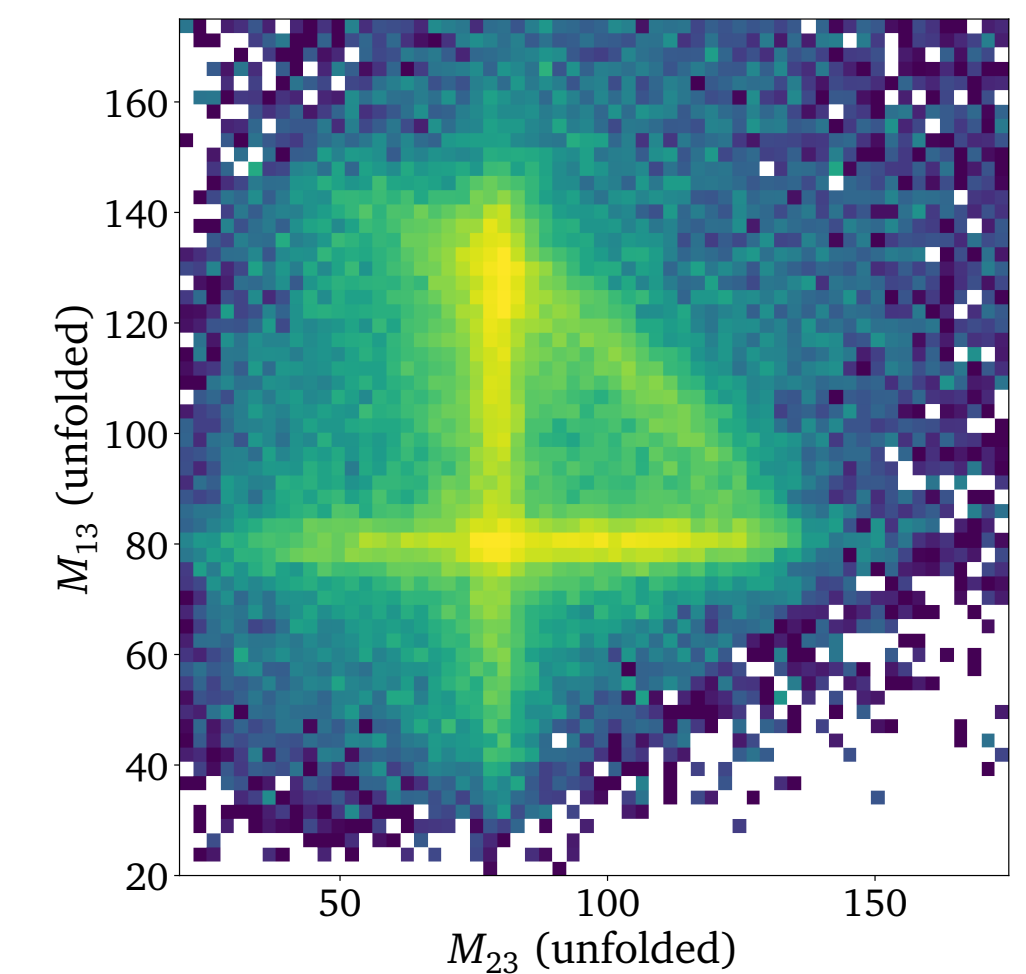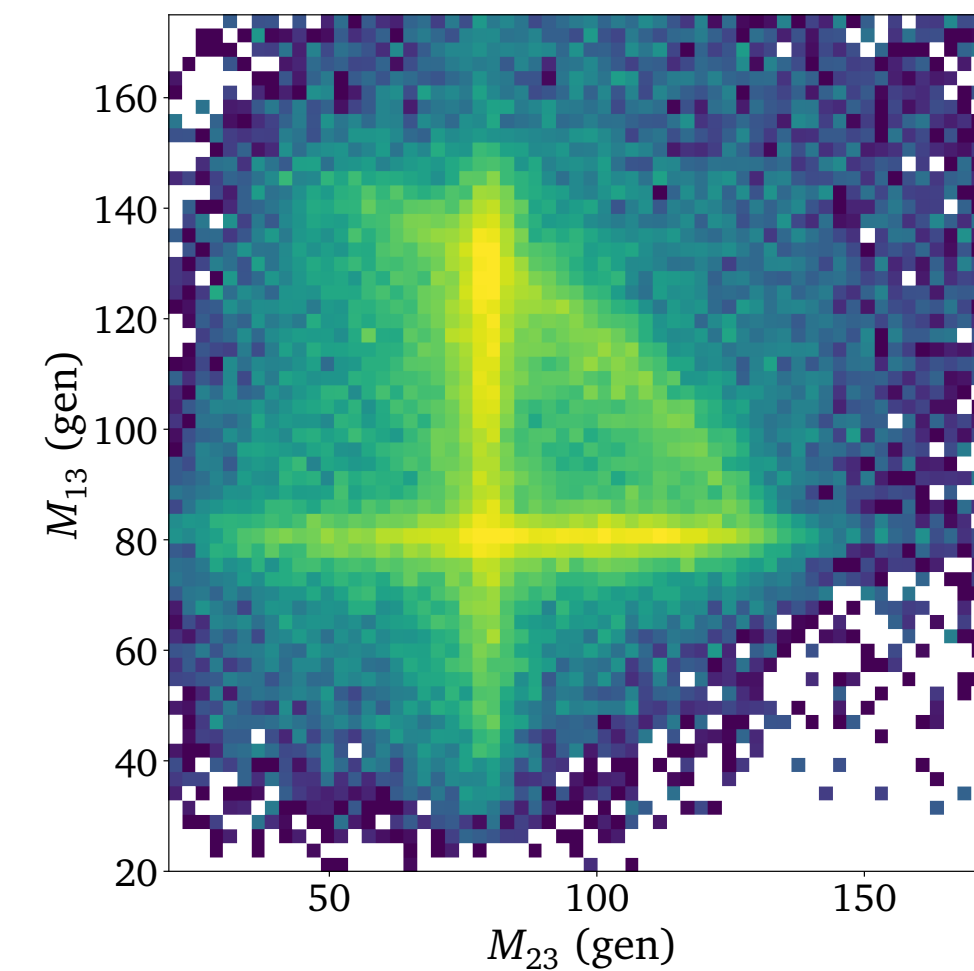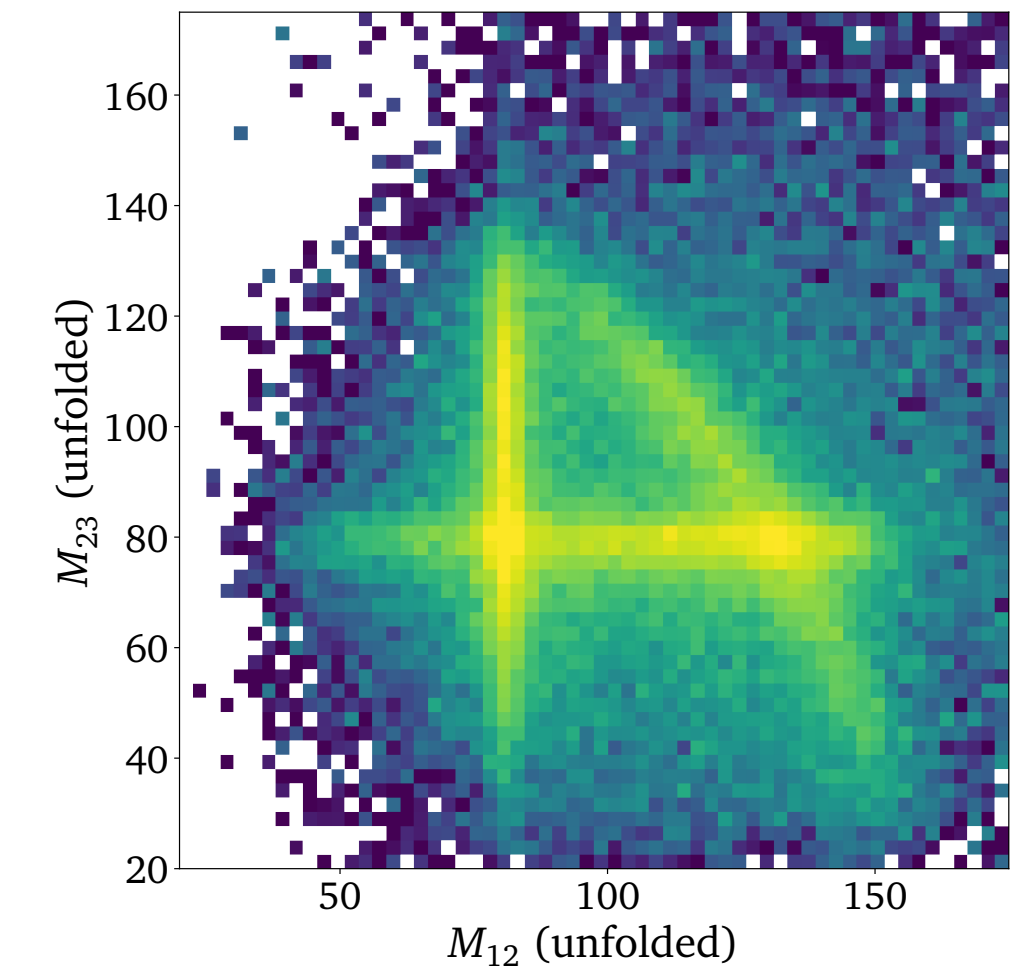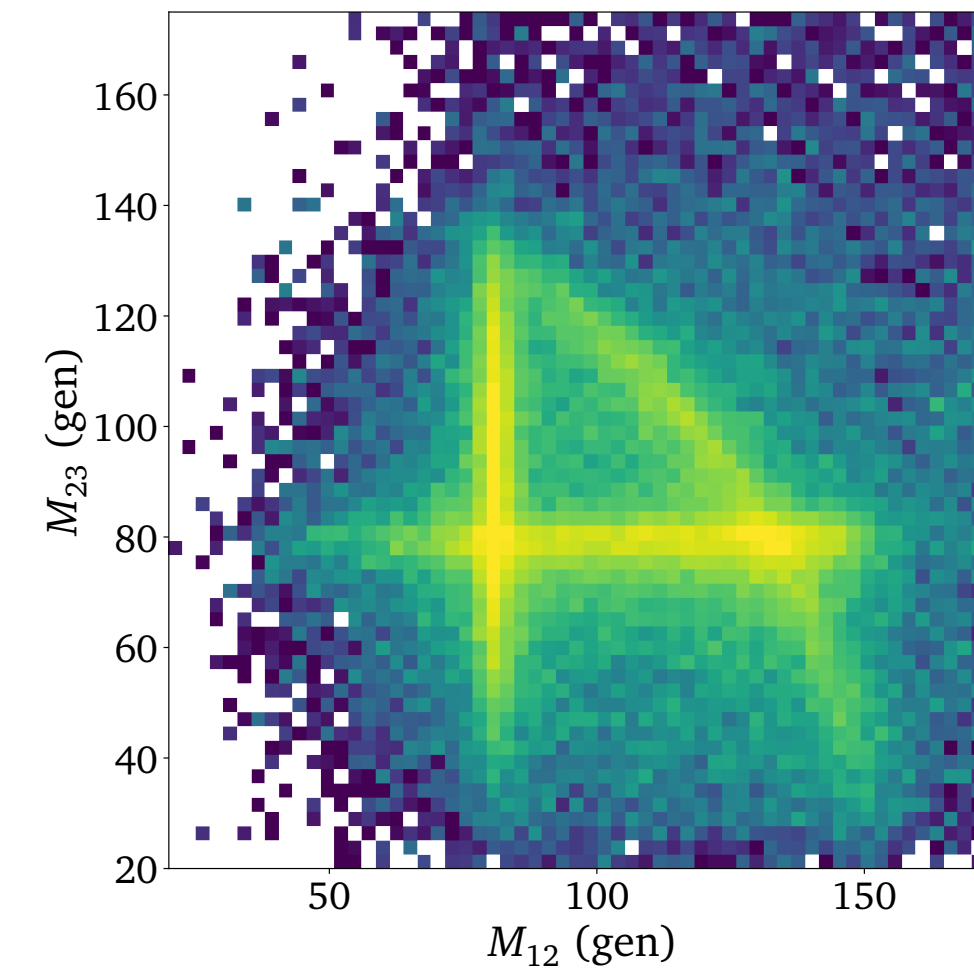Statistical uncertainty for 60 bins decreases by 36%

# And now what?

Generative machine learning allows for unbinned, high dimensional unfolding

Unbiased networks can enhance precision in e.g. top mass measurement

Crucial step to build generative unfolding into existing LHC analysis

Proposal of analysis pipeline:

1. Event Selection
2. Unfold subset
3. Jet calibration
4. Measure top mass
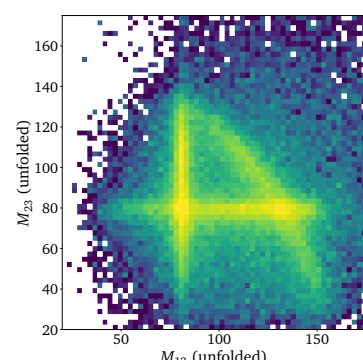5. Resimulate
6. Unfold full phasespace

# And now what?

Generative machine learning allows for unbinned, high dimensional unfolding

Unbiased networks can enhance precision in e.g. top mass measurement

Crucial step to build generative unfolding into existing LHC analysis

Proposal of analysis pipeline:

1. Event Selection
2. Unfold subset
3. Jet calibration
4. Measure top mass
5. Resimulate
6. Unfold full phasespace

re there any questions?