



UNIVERSITÄT  
HEIDELBERG  
ZUKUNFT  
SEIT 1386



CLUSTER OF EXCELLENCE

QUANTUM UNIVERSE



Data Science in Hamburg  
HELMHOLTZ Graduate School  
for the Structure of Matter

# Realistic Calomplification

**Sebastian Bieringer**<sup>1</sup>, Anja Butter, Sascha Diefenbacher, Engin Enren,  
Frank Gaede, Daniel Hundshausen, Gregor Kasieczka,  
Benjamin Nachman, Tilman Plehn, Mathias Trabs

<sup>1</sup>Institut für Experimentalphysik, Universität Hamburg, Germany

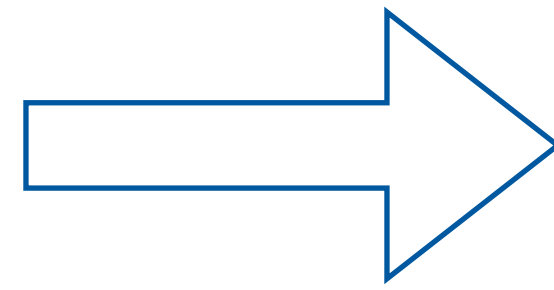
[sebastian.guido.bieringer@uni-hamburg.de](mailto:sebastian.guido.bieringer@uni-hamburg.de)

IDT-UM / ErUM-Data Meeting

# Introduction

Need to speed up MC

- Event generation
- Calorimeter simulation



Use generative machine learning models

- Generative Adversarial Networks (GANs)
- Variational Autoencoders (VAEs)
- Normalizing Flows

$$\text{simulation speed} = \frac{\# \text{ samples}}{\text{time}}$$

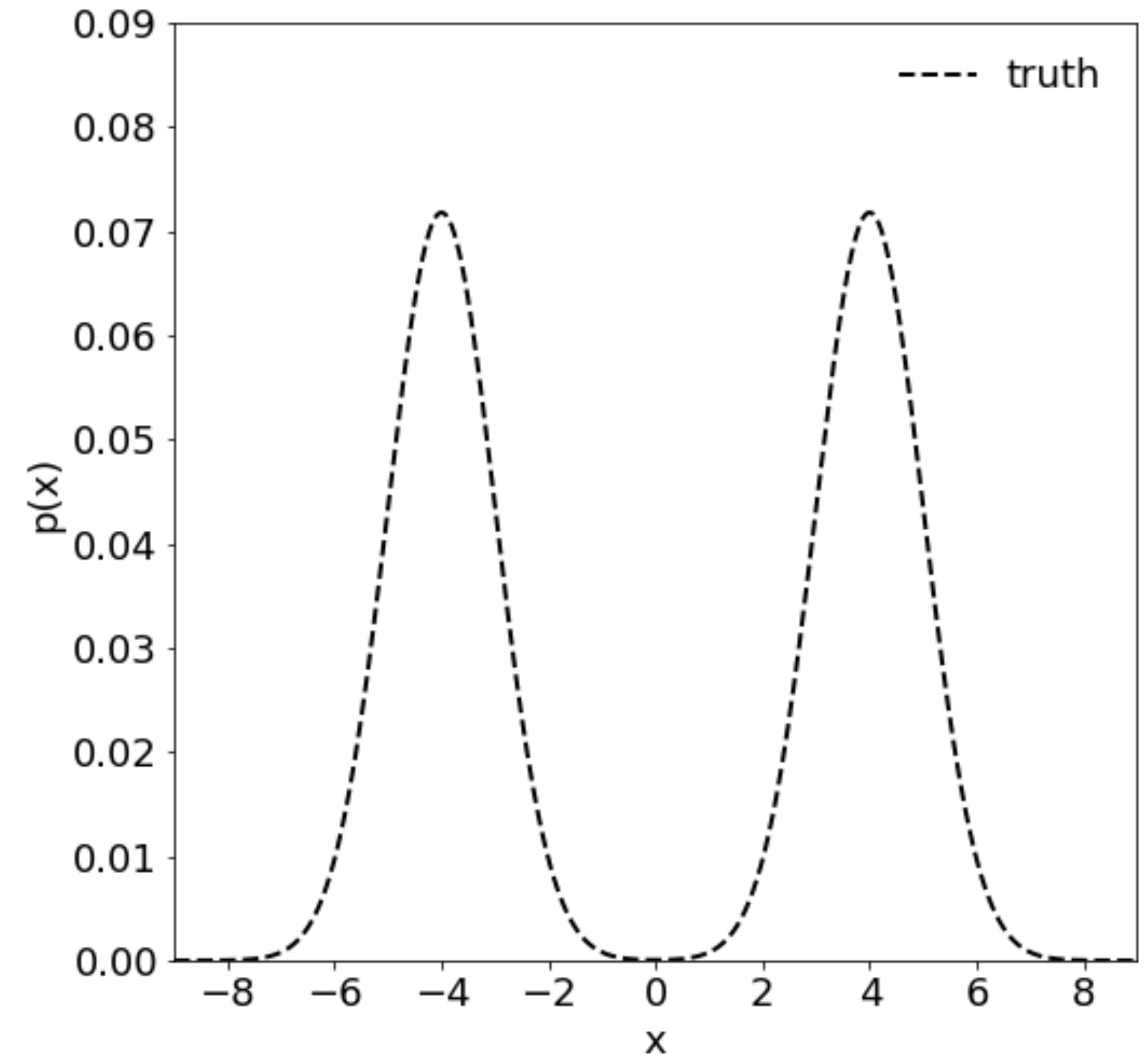
What about # samples? How many new points can a generative model generate?

Anja Butter et al. *GANplifying Event Samples*. 2021. arXiv: [2008.06545 \[hep-ph\]](https://arxiv.org/abs/2008.06545)

# Toy Model: Setup

- Camel back function:

$$P(x) = \frac{1}{2} \left( \mathcal{N}_{-4,1}(x) + \mathcal{N}_{4,1}(x) \right)$$

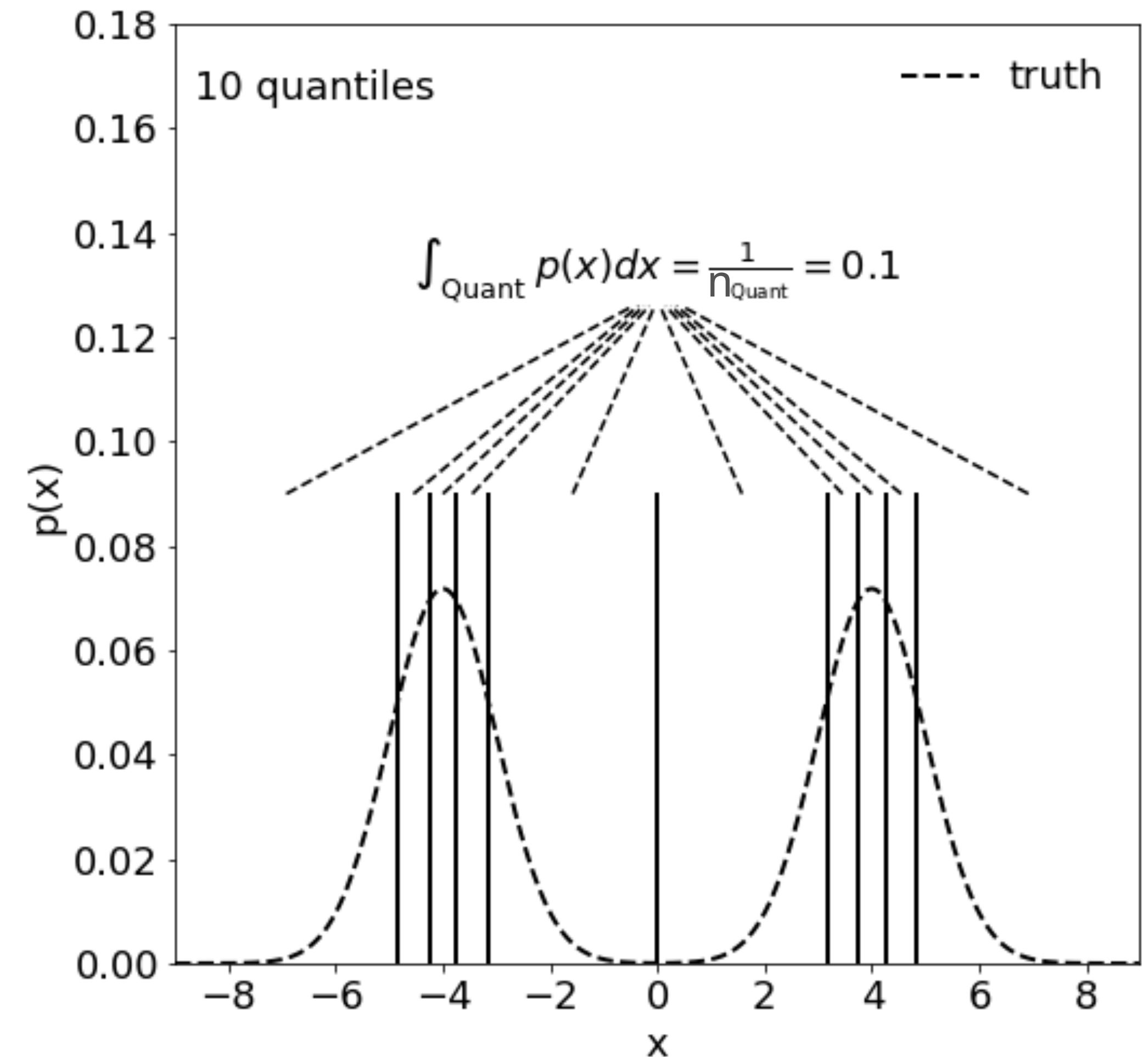


# Toy Model: Setup

- Camel back function:

$$P(x) = \frac{1}{2} \left( \mathcal{N}_{-4,1}(x) + \mathcal{N}_{4,1}(x) \right)$$

- "Pearson  $\chi^2$ -test":
  - Introduce equal probability quantiles



# Toy Model: Setup

- Camel back function:

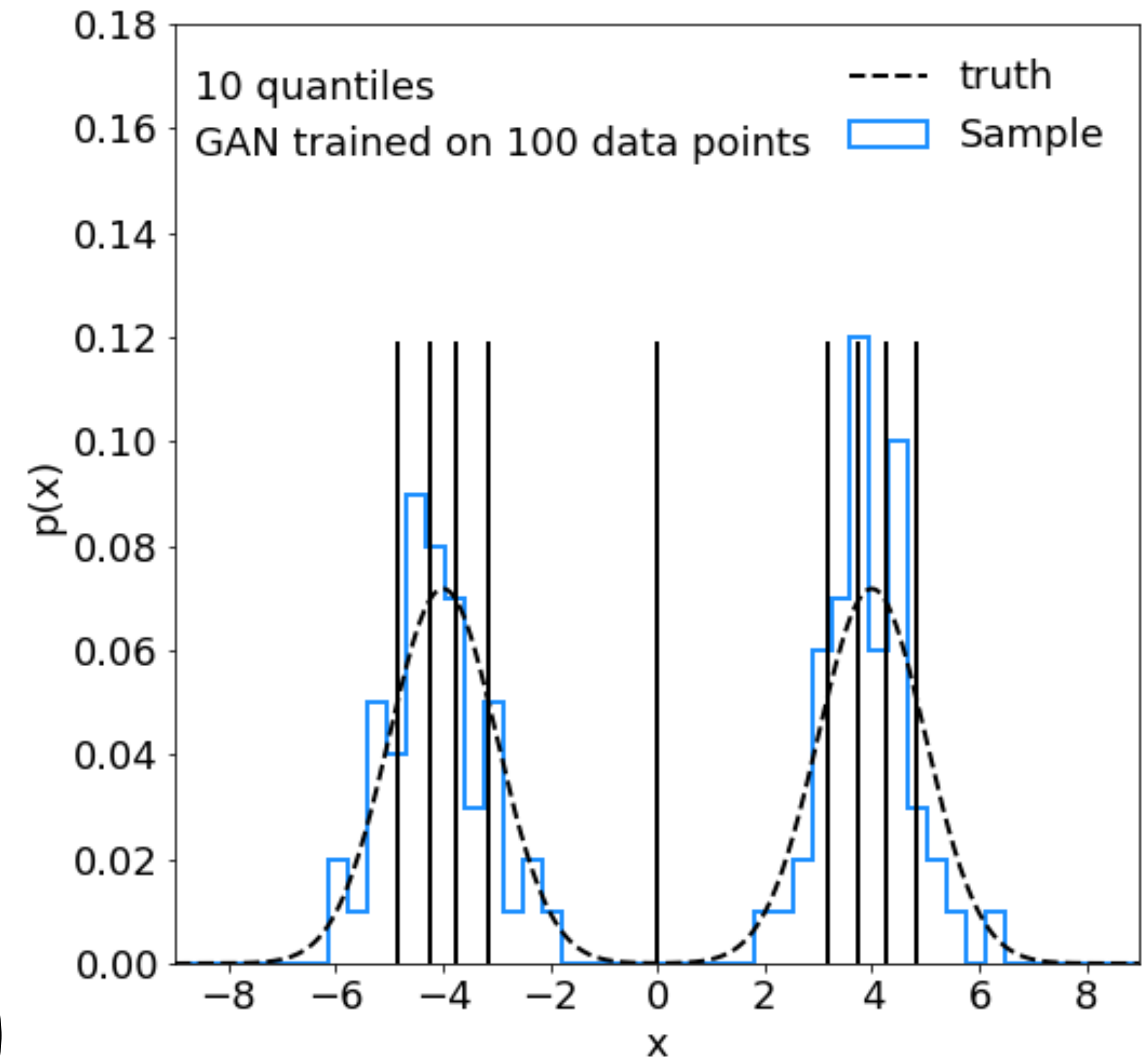
$$P(x) = \frac{1}{2} \left( \mathcal{N}_{-4,1}(x) + \mathcal{N}_{4,1}(x) \right)$$

- "Pearson  $\chi^2$ -test":

- Introduce equal probability quantiles
- Generate data and calculate

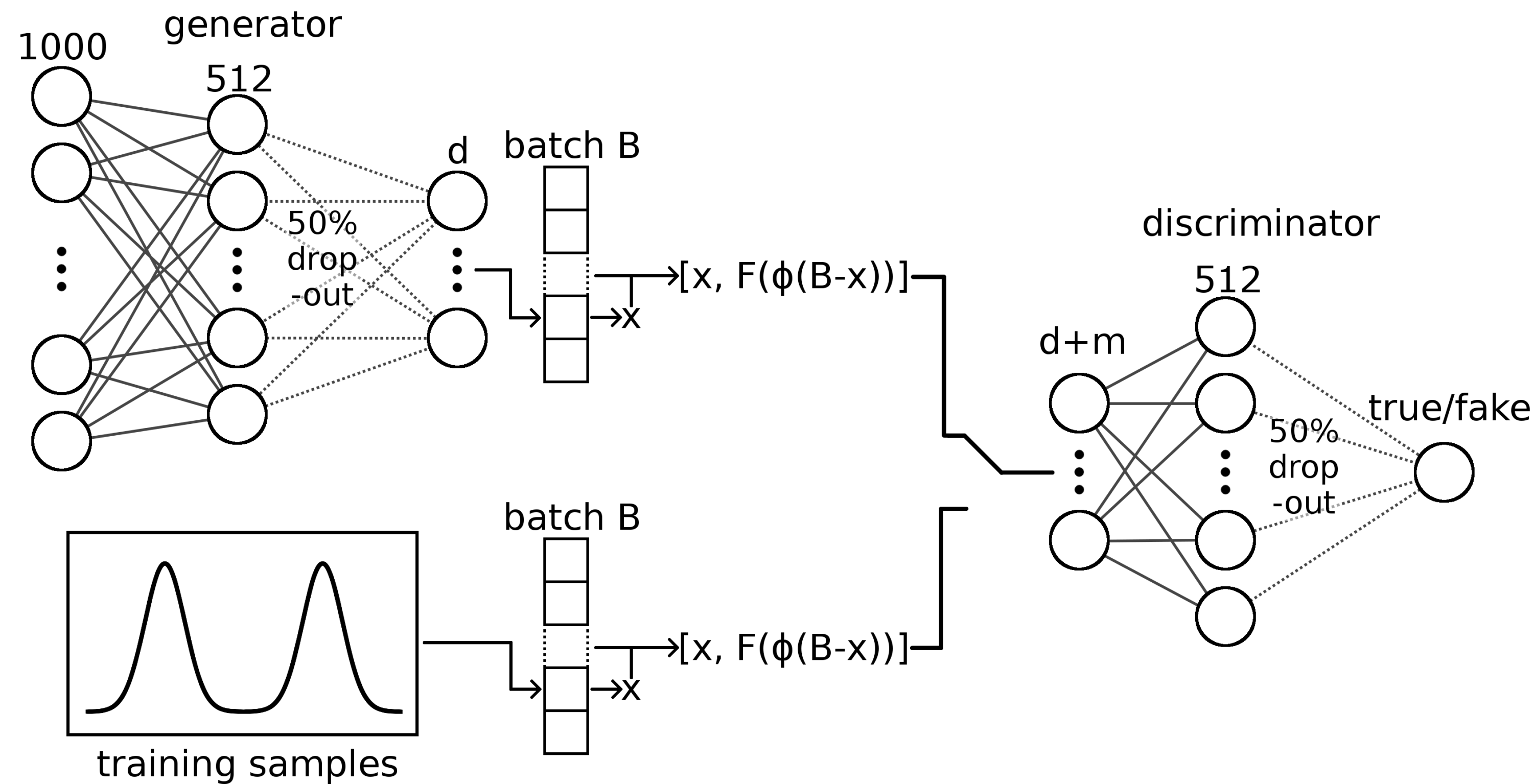
$$\hat{\chi}_{n_{\text{quant}}}^2 = n_{\text{quant}} \sum_{j=0}^{n_{\text{quant}}} \left( x_j - \frac{1}{n_{\text{quant}}} \right)^2$$

with  $\hat{\chi}_{n_{\text{quant}}}^2 \xrightarrow{n_{\text{quant}} \rightarrow \infty} \chi^2 \left( P(x), P_{\text{samples}} \right)$



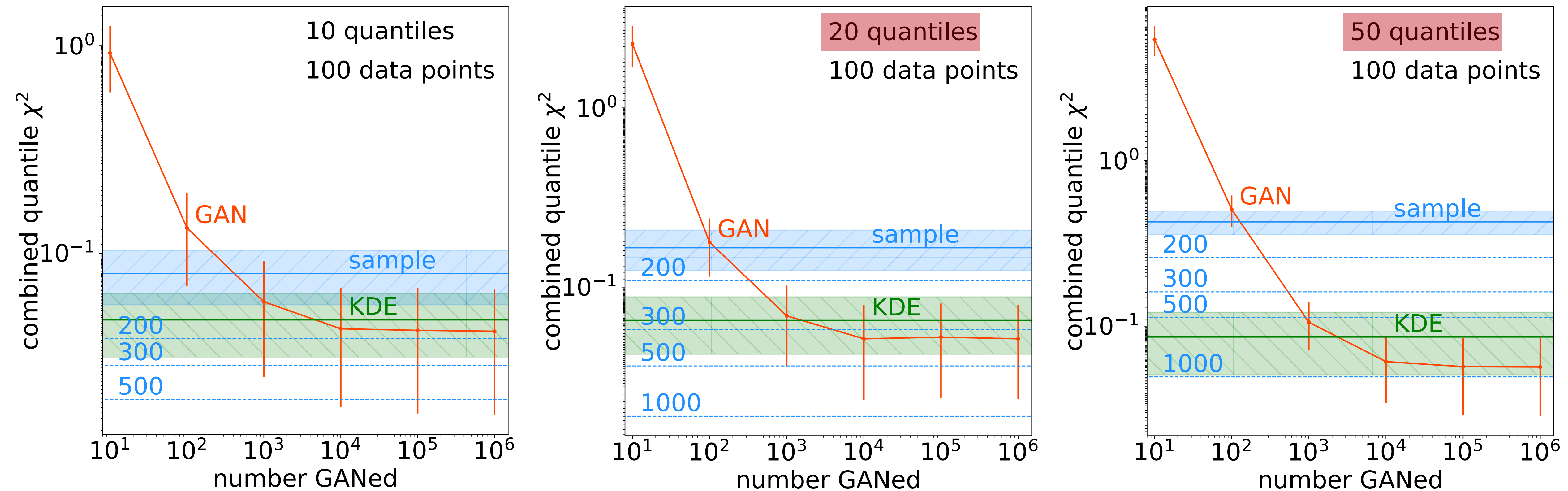
# Toy Model: Generative Network

- Train on  $n_{\text{data}} = 100$  data points generated from  $P(x)$
- Prone to mode-collapse and overfitting:
  - Dropout
  - Noise augmentation
  - Batch-statistics
- Generate high amounts of data from Network



# Toy Model: 1D

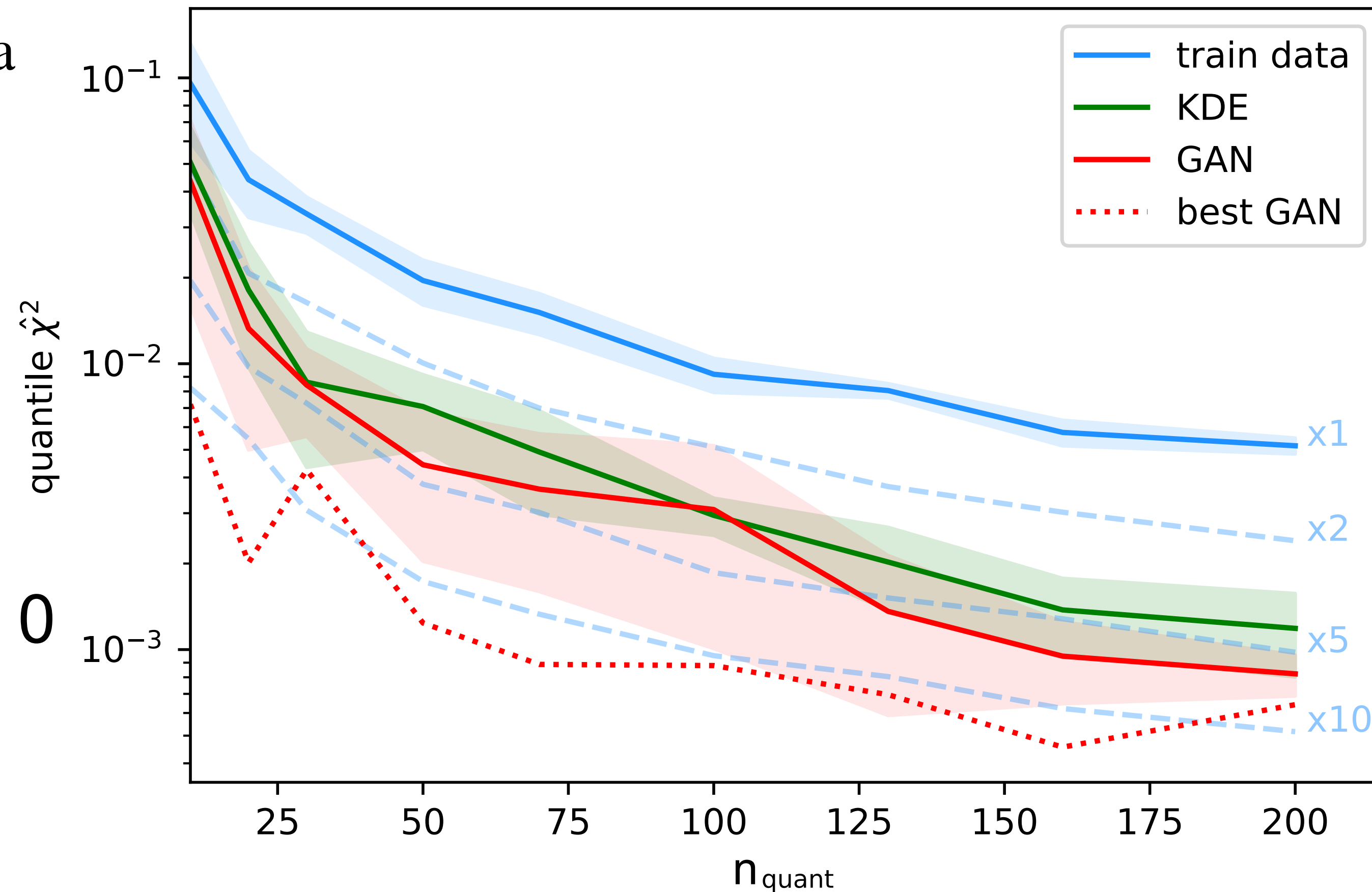
- Interpolation more noticeable for sparse data



- However,** quantile measure breaks down for sparse data

# Toy Model: 1D

- Examine high  $n_{\text{quant}}$  and high  $n_{\text{data}}$ 
  - Train on  $n_{\text{data}} = n_{\text{quant}}^2$
  - Generate  $100 \cdot n_{\text{data}}$
- Examine which data converges to 0 (fastest)
- GAN amplifies data by a factor  $\sim 5$



# Calorimeter Simulations: Data

- 269k photon showers at 50 GeV in ILD detector [1]

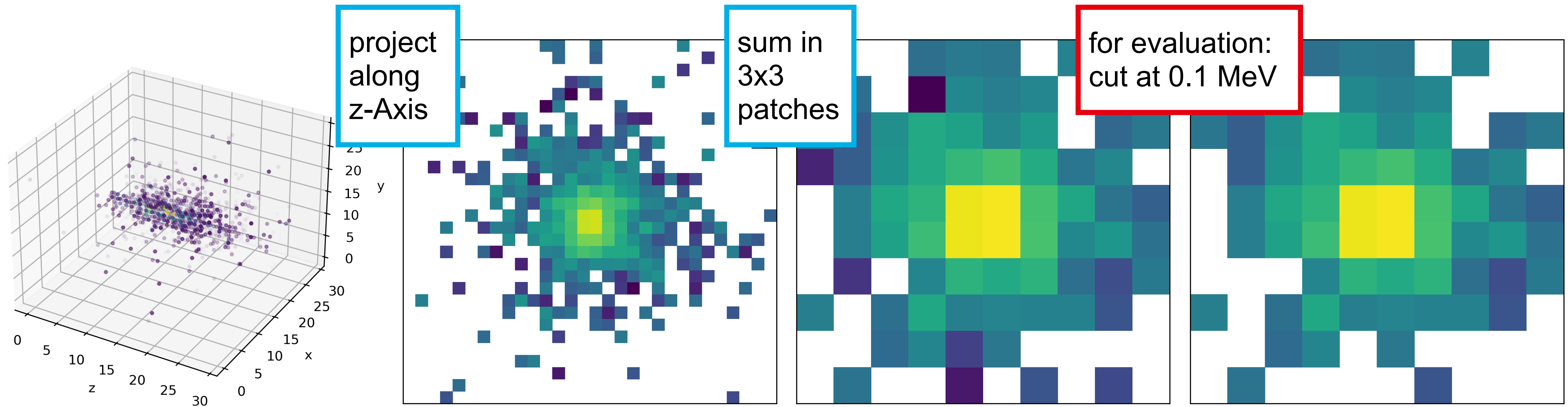


Image-shaped data

Unknown true distribution, limited data

Harder learning task → training on multiple training set sizes unfeasible

# Calorimeter Simulations: Architecture

- Change to *location aware* VAE-GAN architecture [2,3,4]

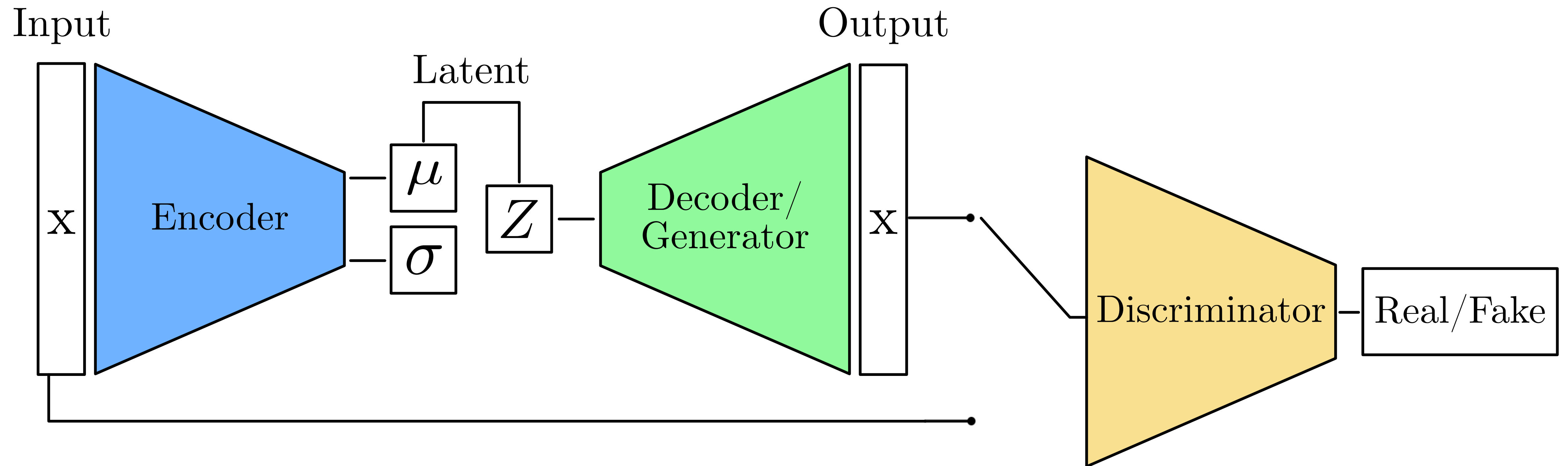


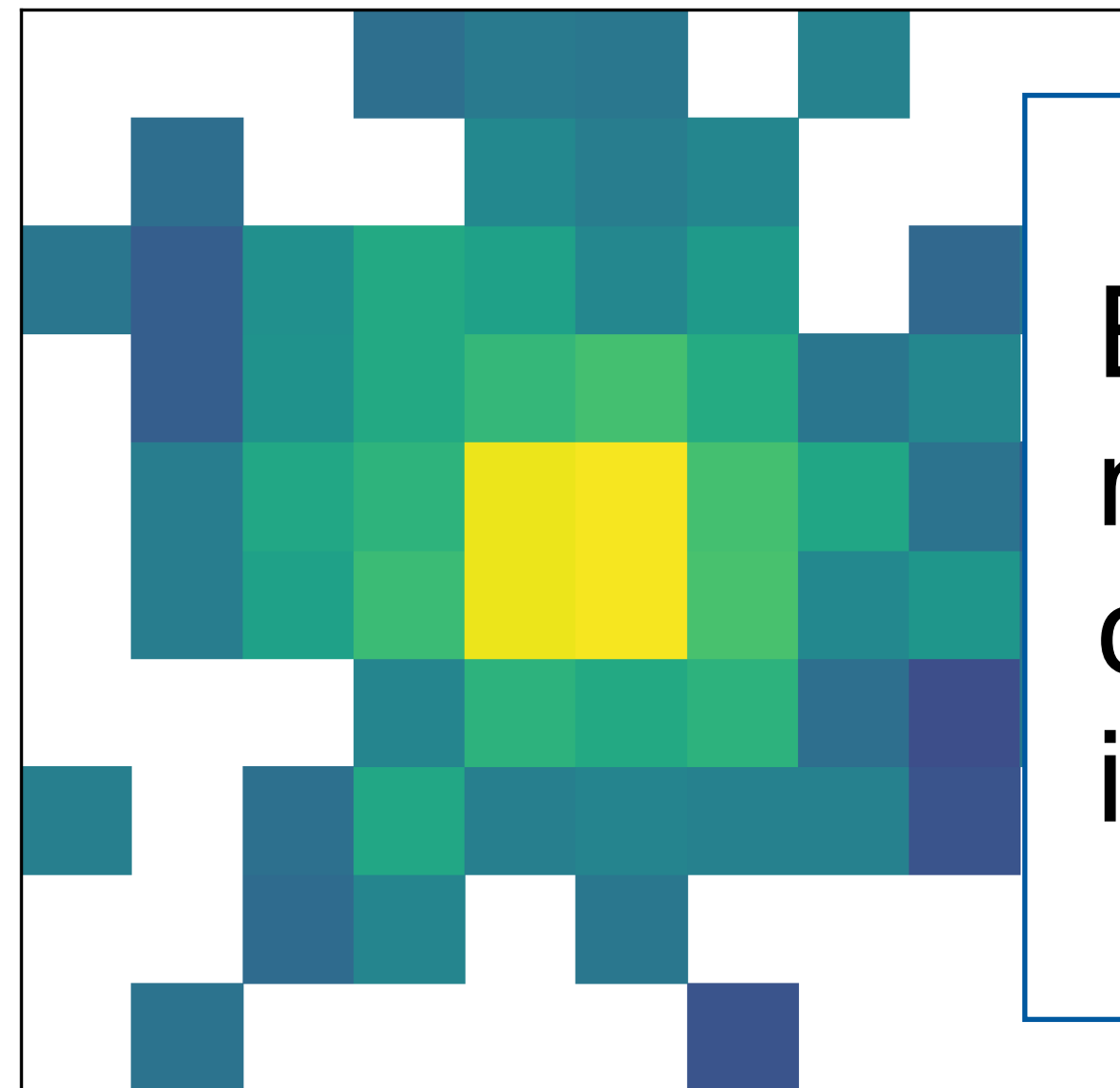
Image-shaped data

Unknown true distribution, limited data

Harder learning task  $\rightarrow$  training on multiple training set sizes unfeasible

# Calorimeter Simulations: Setup

DASHH



Evaluate 1D metrics, calculated on images

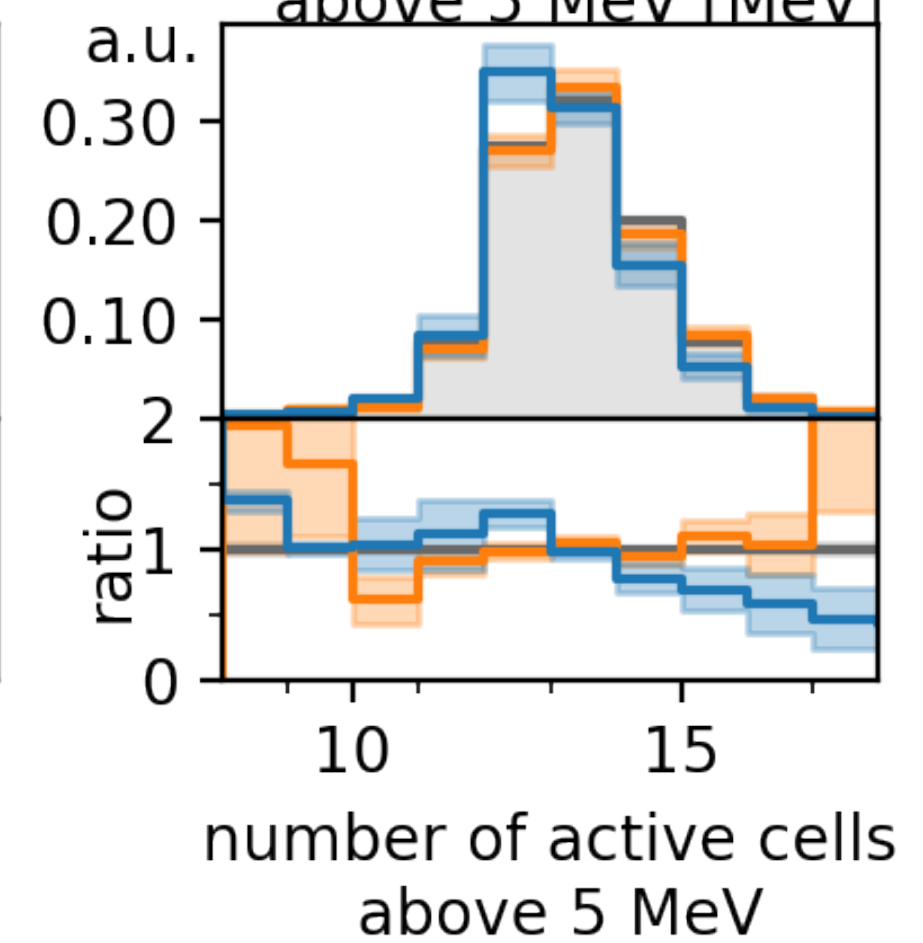
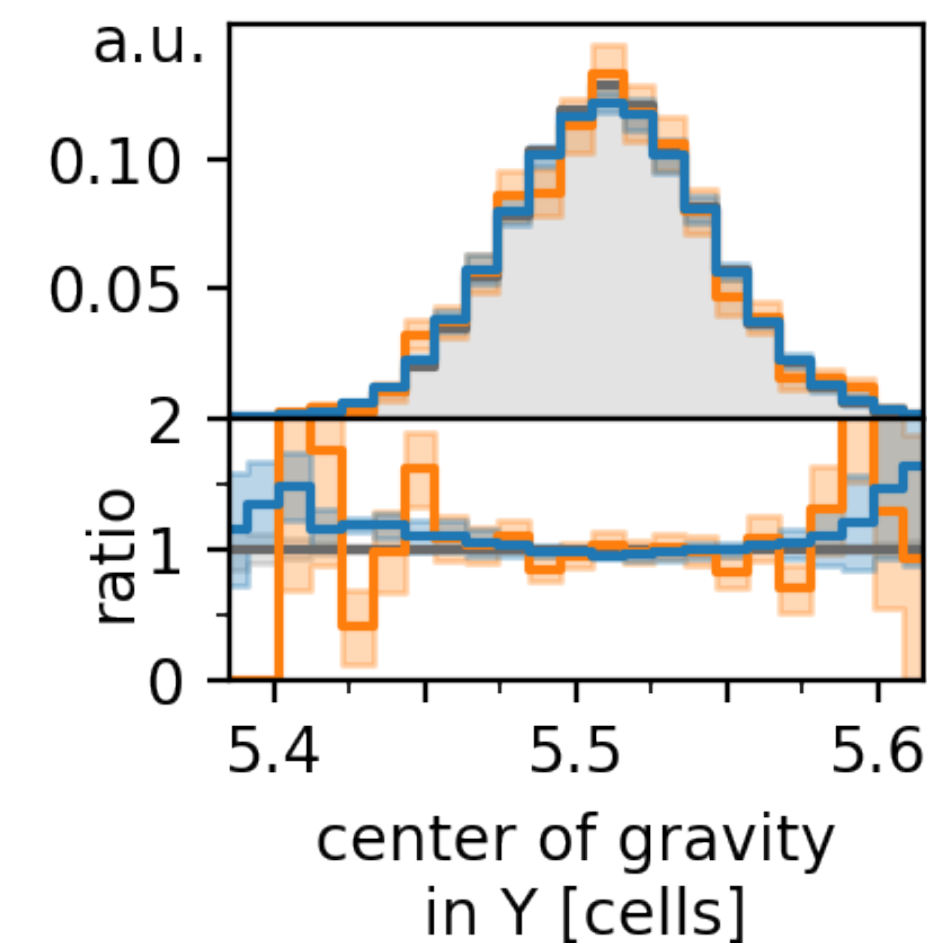
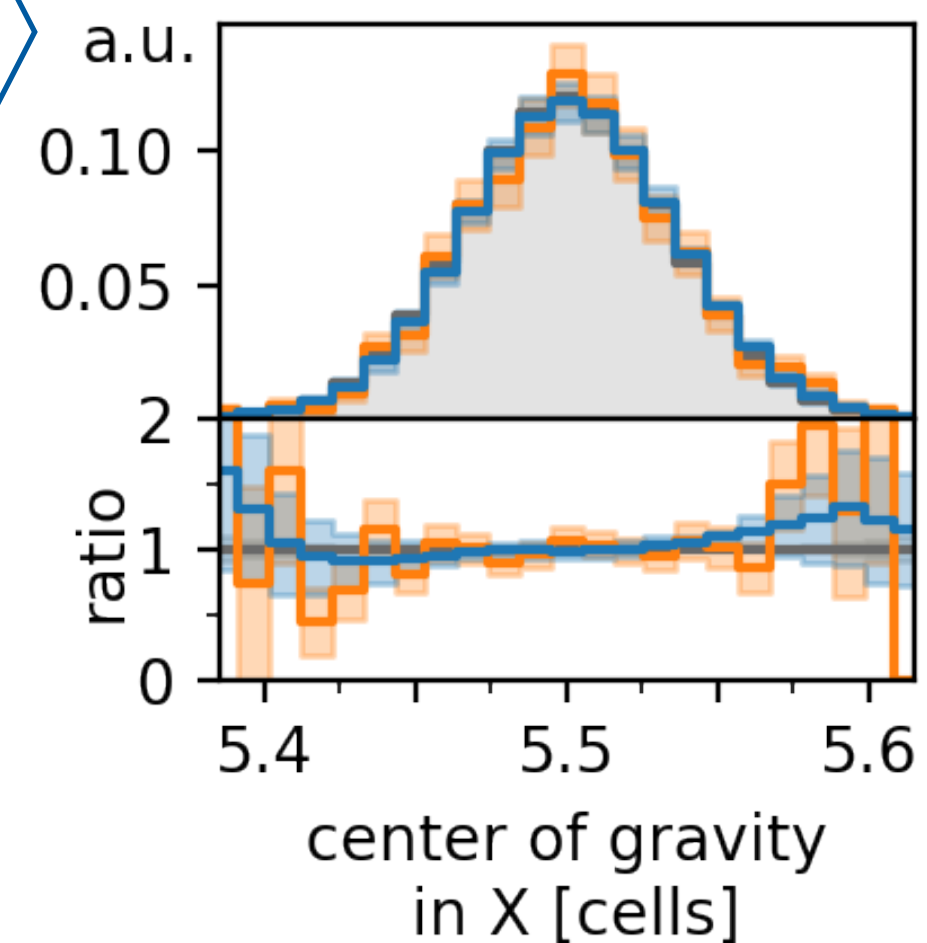
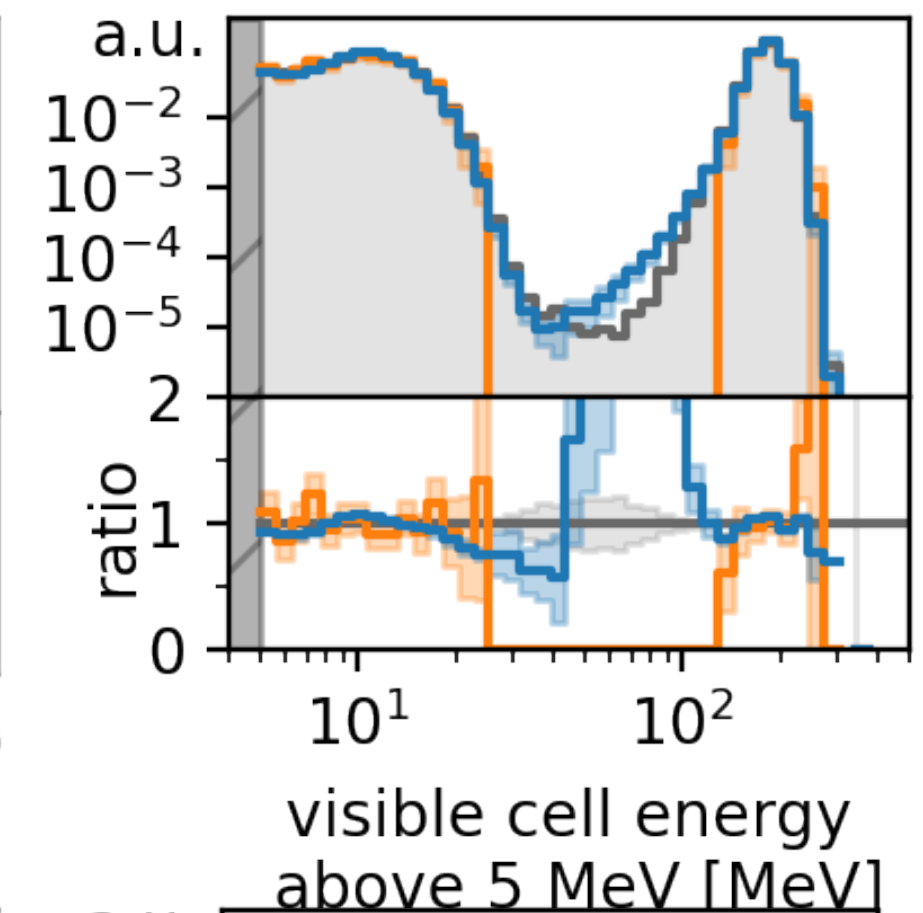
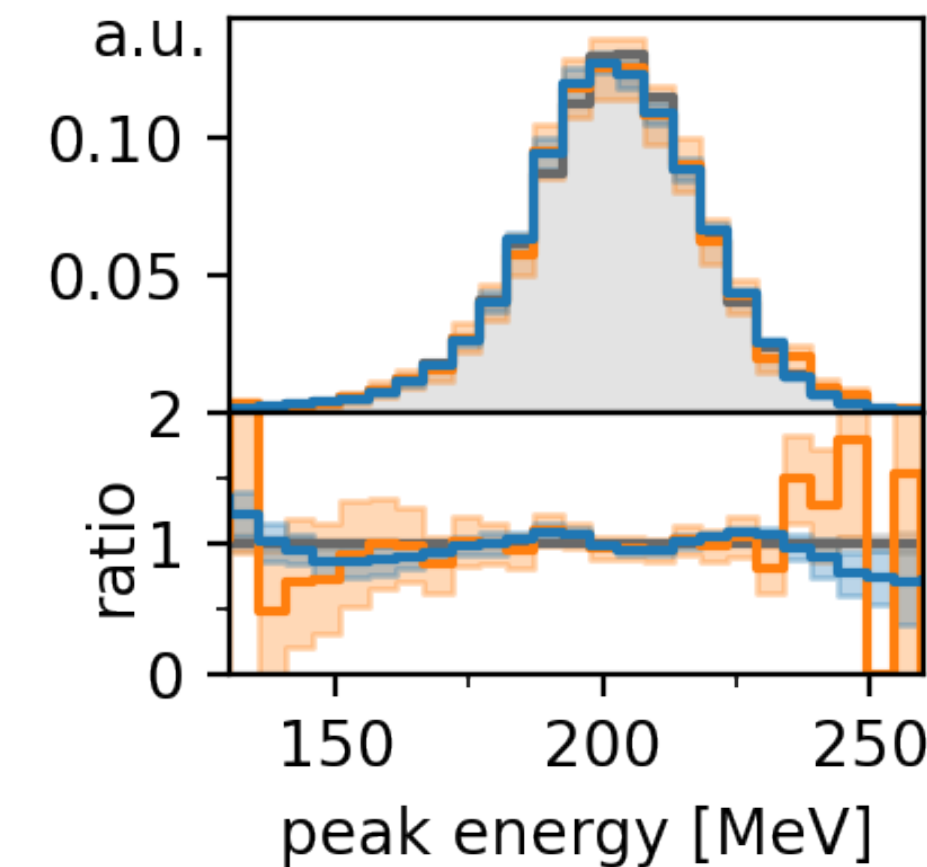
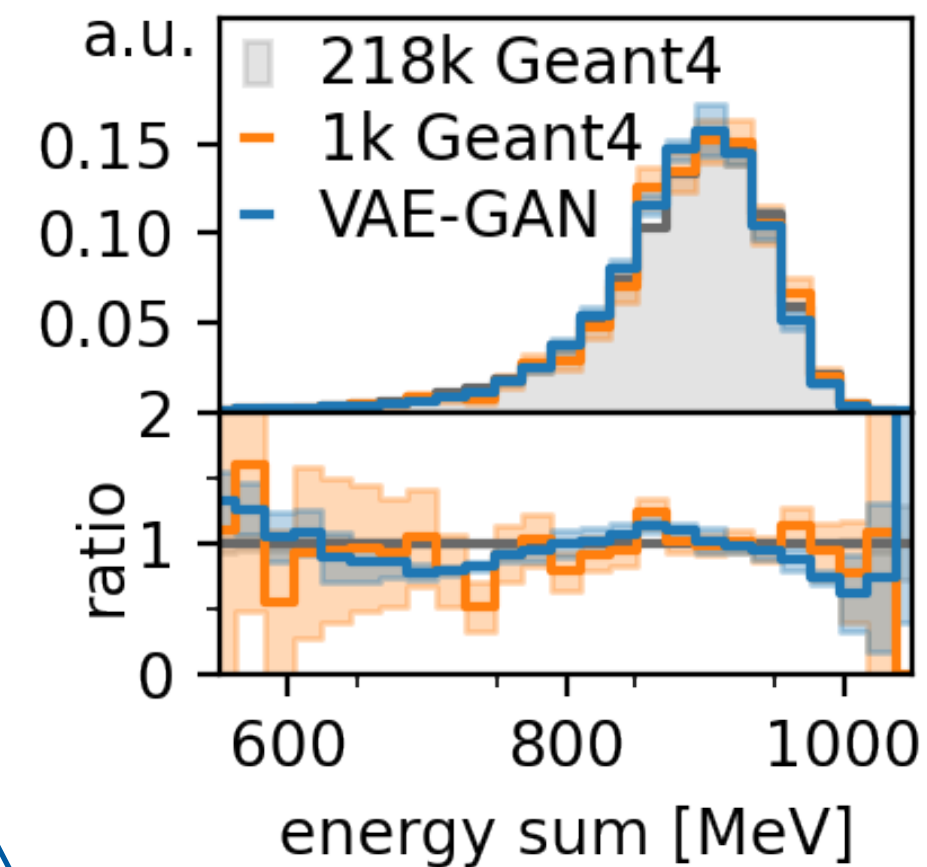


Image-shaped data

Unknown true distribution, limited data

Harder learning task → training on multiple training set sizes unfeasible

# Calorimeter Simulations: Setup

- Split into **218k validation data points** and **50k evaluation data points**
- Generate quantiles by dividing the validation set into equal-sized parts

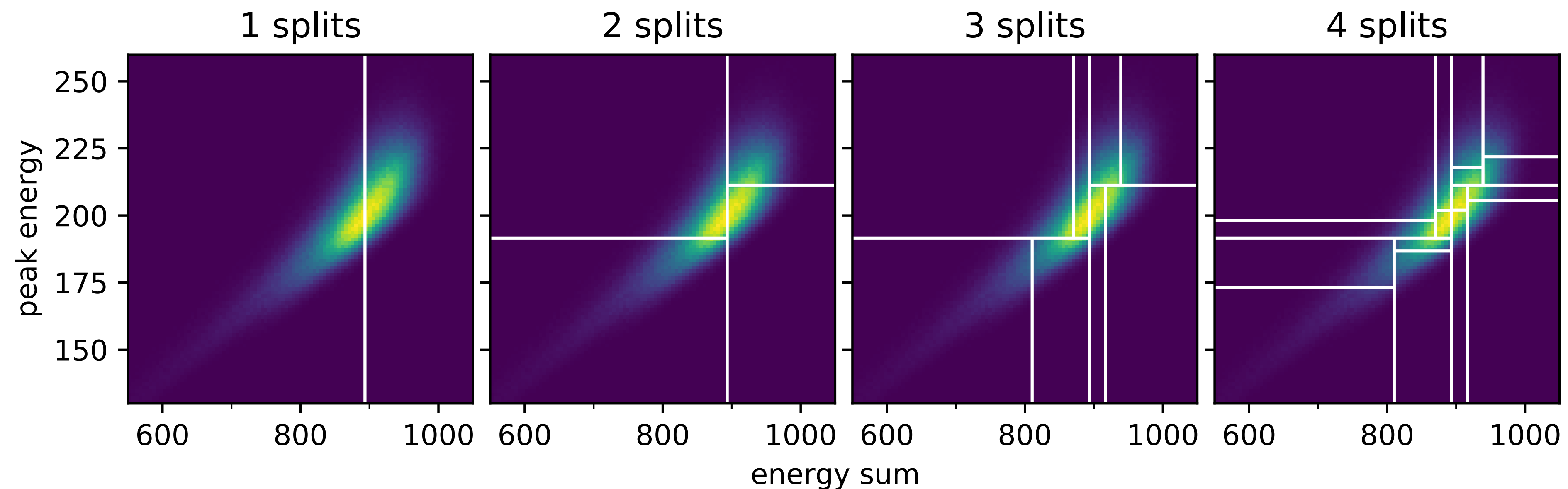


Image-shaped data

Unknown true distribution, limited data

Harder learning task → training on multiple training set sizes unfeasible

# Calorimeter Simulations: Setup

$$\text{Calculate } \overline{\text{JSD}}(g || p) = \frac{1}{2} \sum_{Q_i \in Q} \left( g_i \log \frac{g_i}{\frac{1}{2}(g_i + p_i)} + p_i \log \frac{p_i}{\frac{1}{2}(g_i + p_i)} \right).$$

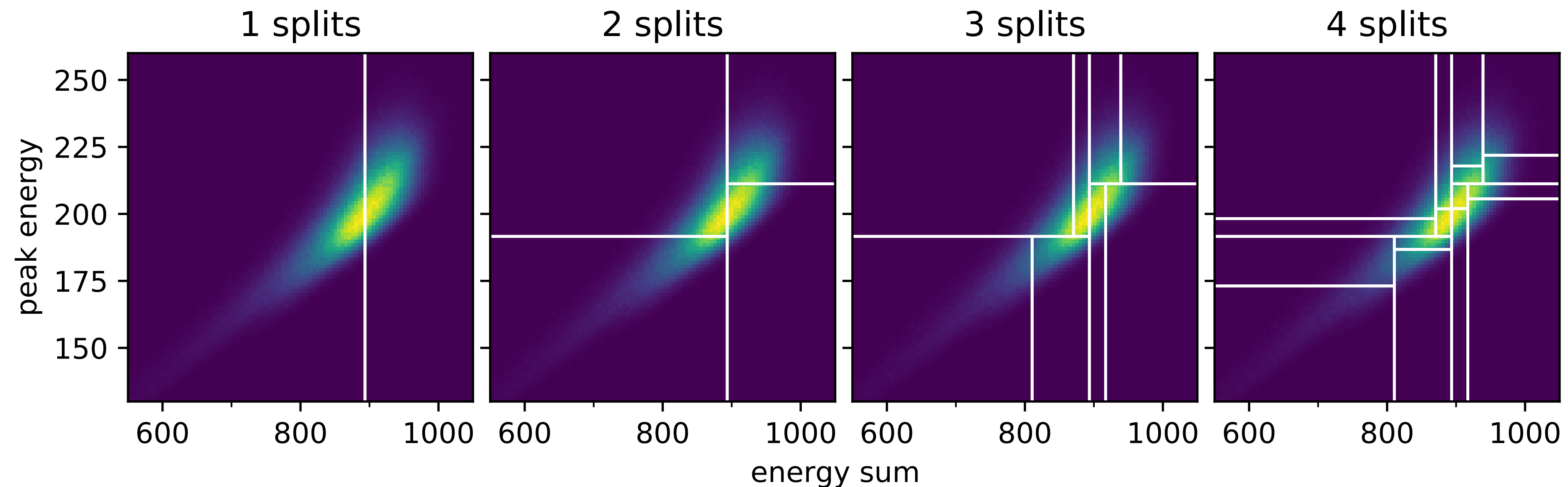


Image-shaped data

Unknown true distribution, limited data

Harder learning task → training on multiple training set sizes unfeasible

# Calorimeter Simulations: Results



- Evaluate for fixed training set size
- Add multiples of the training set size
- Use less than  $n_{\text{data}}/10$  bins

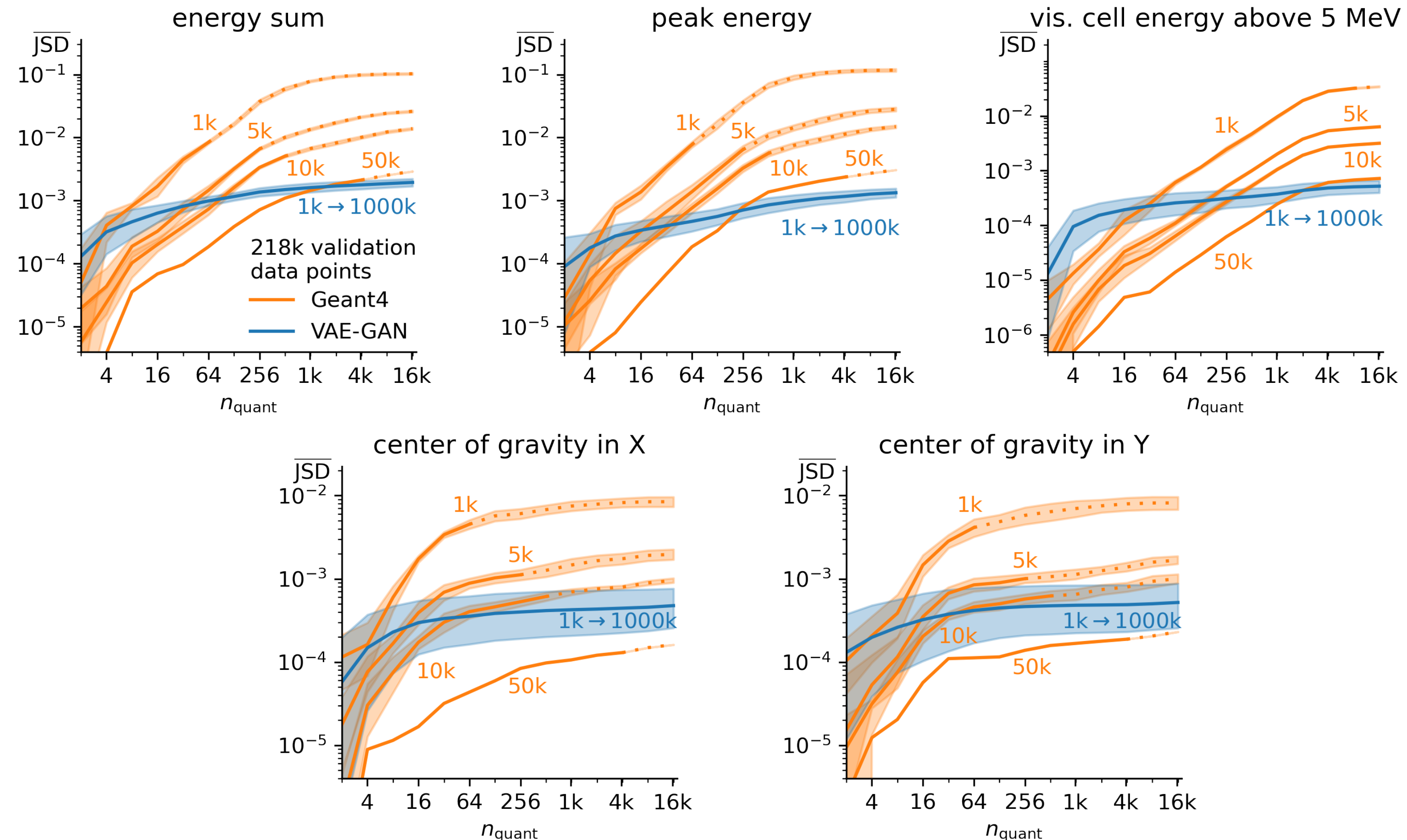


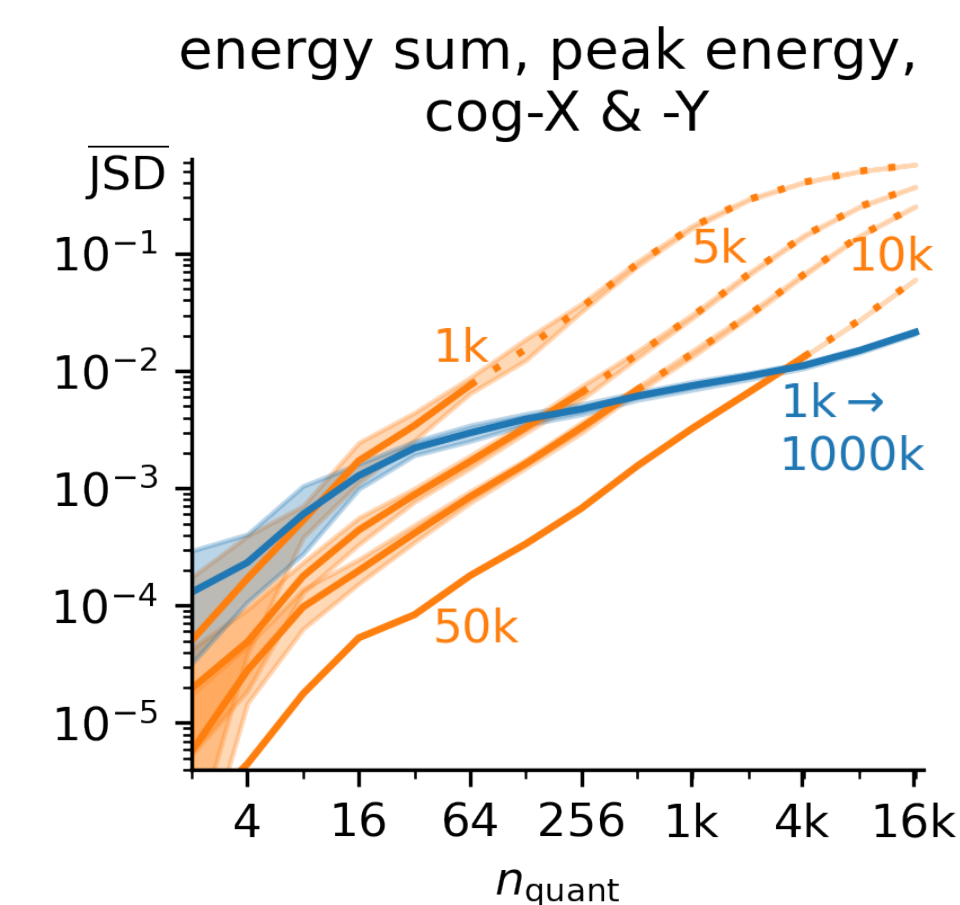
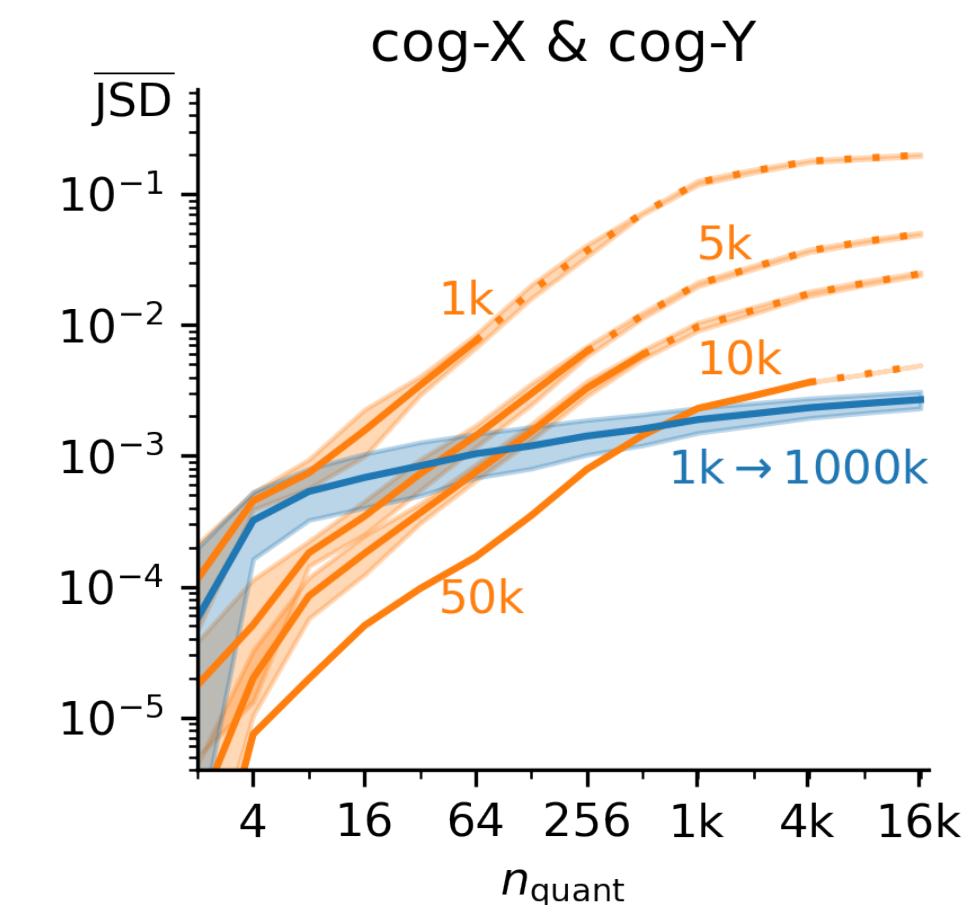
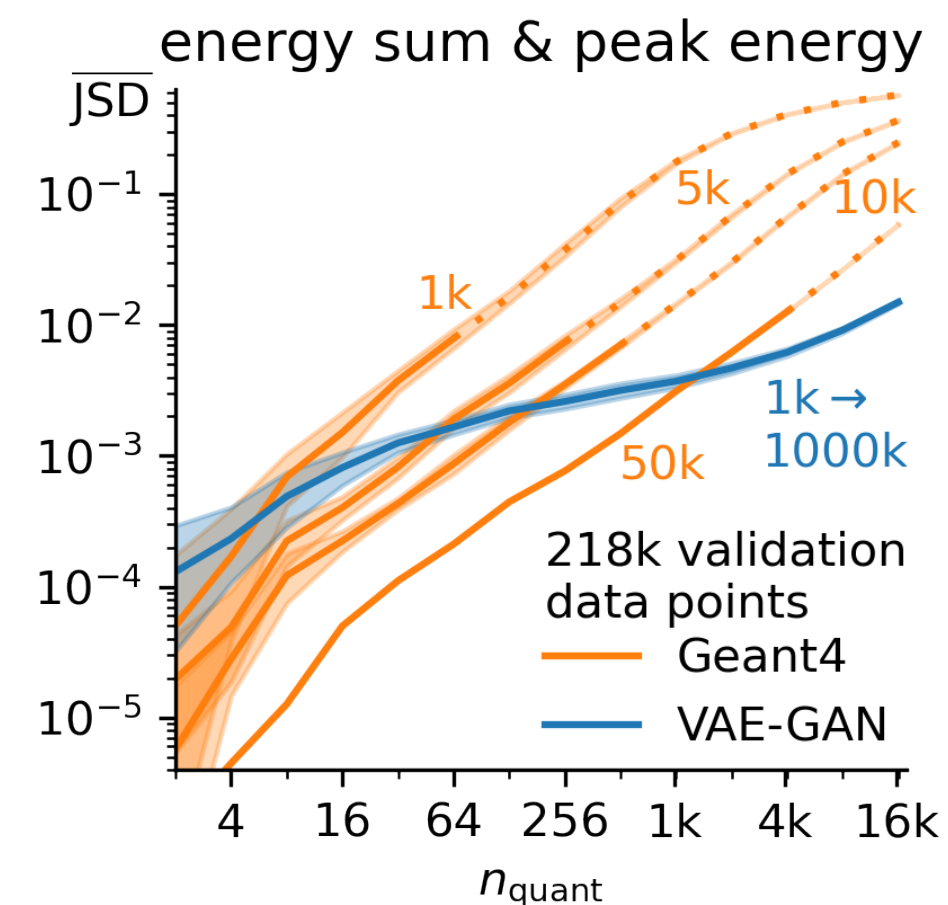
Image-shaped data

Unknown true distribution, limited data

Harder learning task  $\rightarrow$  training on multiple training set sizes unfeasible

# Calorimeter Simulations: Results

- Evaluate for fixed training set size
- Add multiples of the training set size
- Use less than  $n_{\text{data}}/10$  bins



- High scale features: limited by amount of training data
- Low scale features: GAN estimation can not be matched by adding more data

Image-shaped data

Unknown true distribution, limited data

Harder learning task → training on multiple training set sizes unfeasible

# Calorimeter Simulations: Results

How good is the density estimation actually?

- Compare to KDE and histogram estimators (maximizing log-likelihood of cross-validation sets)

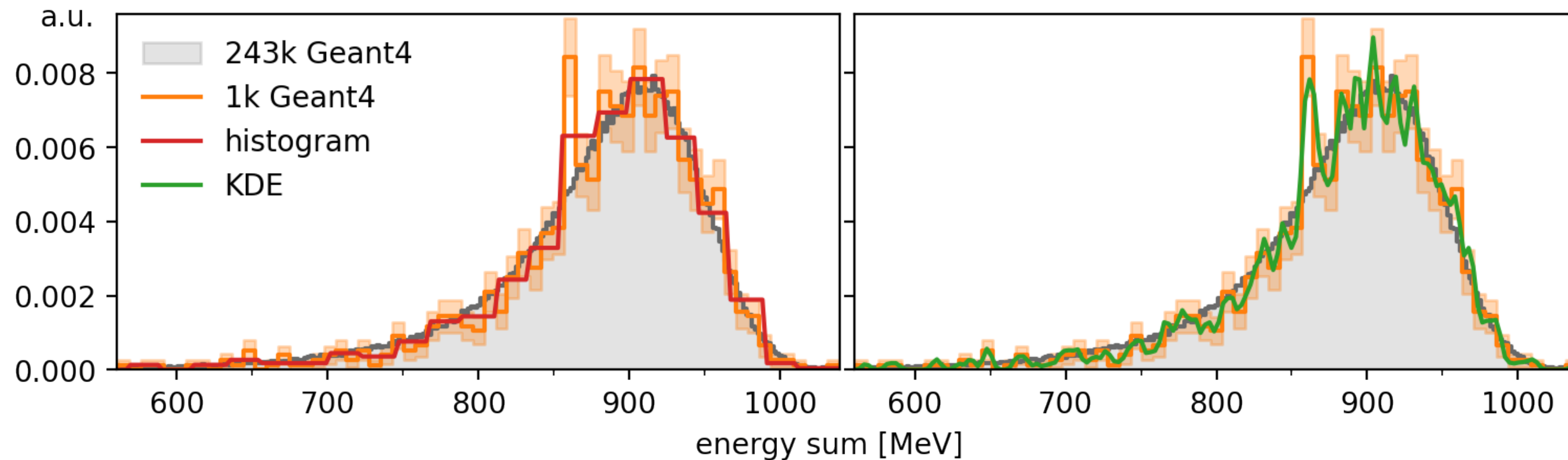


Image-shaped data

Unknown true distribution, limited data

Harder learning task → training on multiple training set sizes unfeasible

# Calorimeter Simulations: Results



- Generate  $10^6$  samples from every density estimator

- GAN outperforms standard density estimators

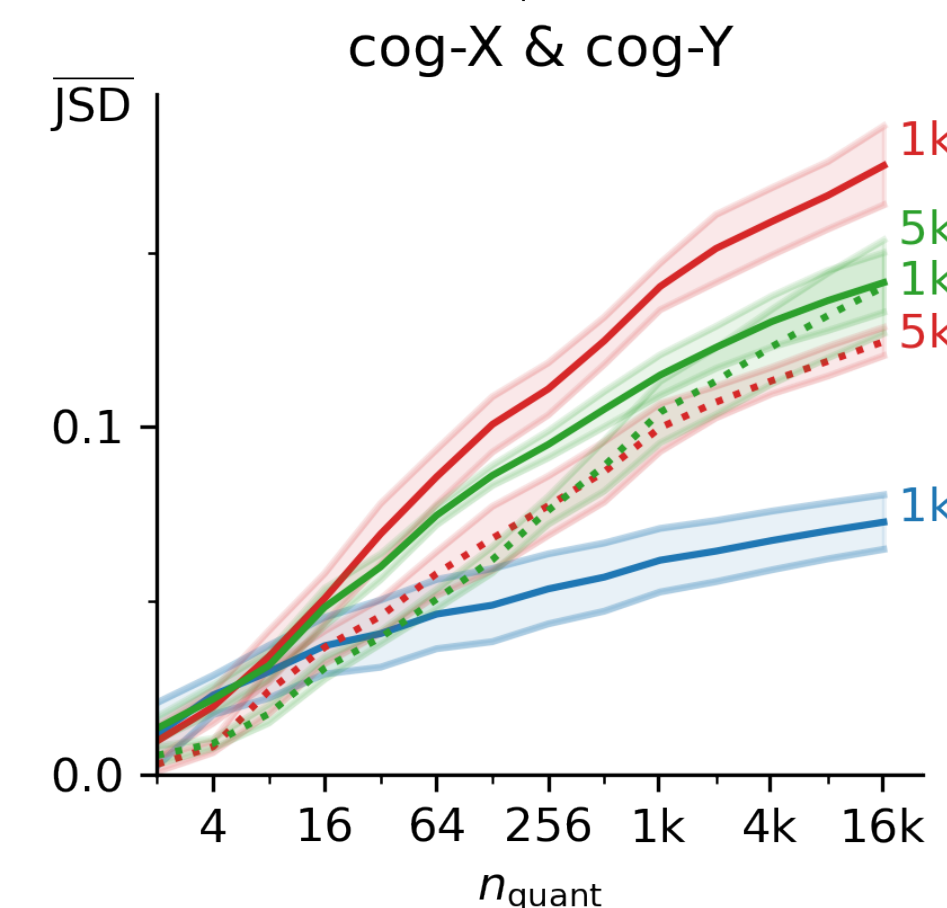
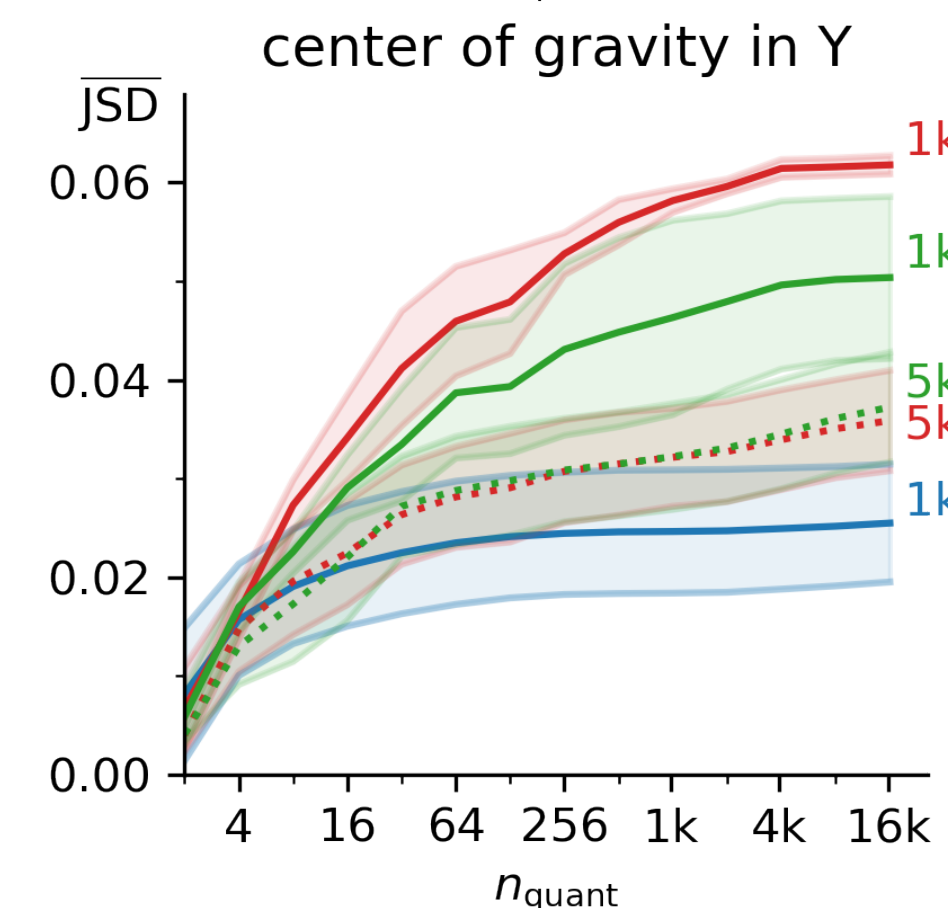
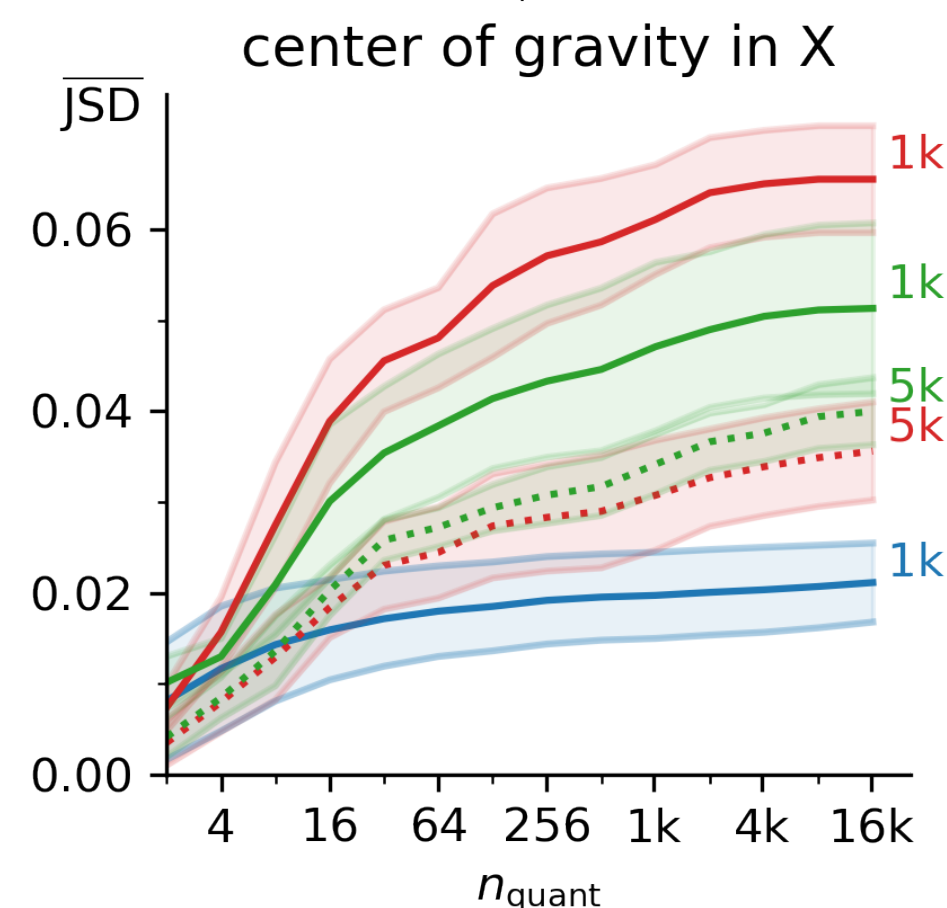
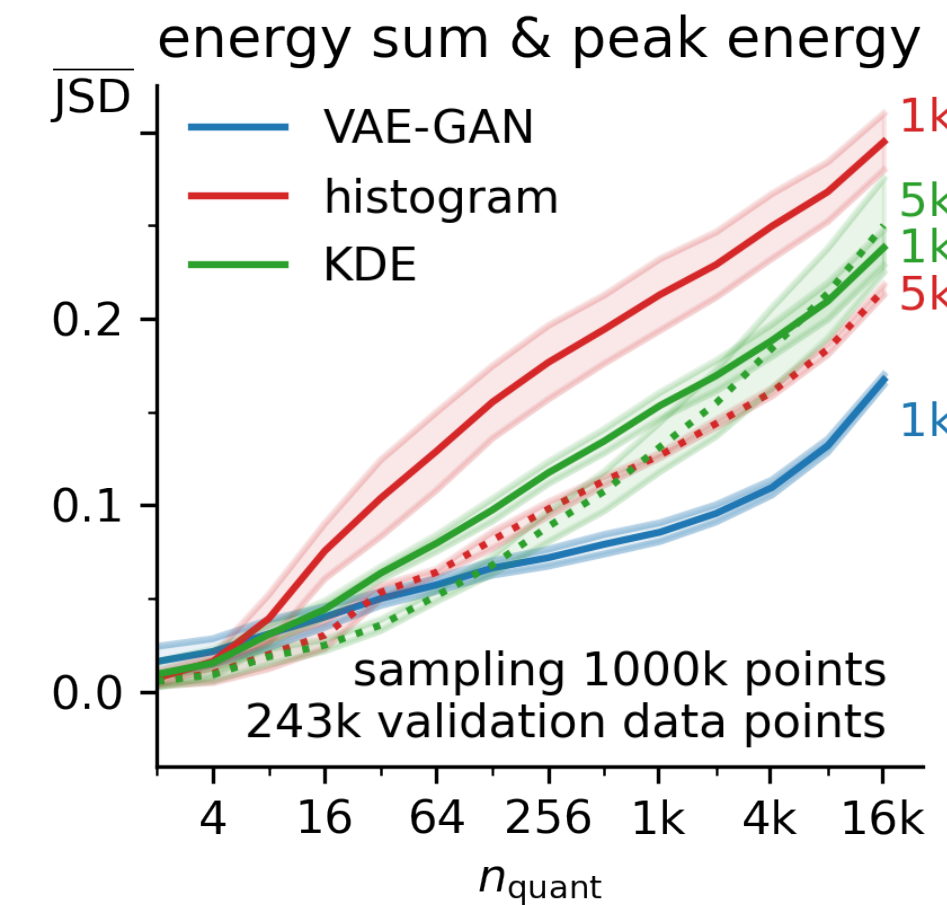
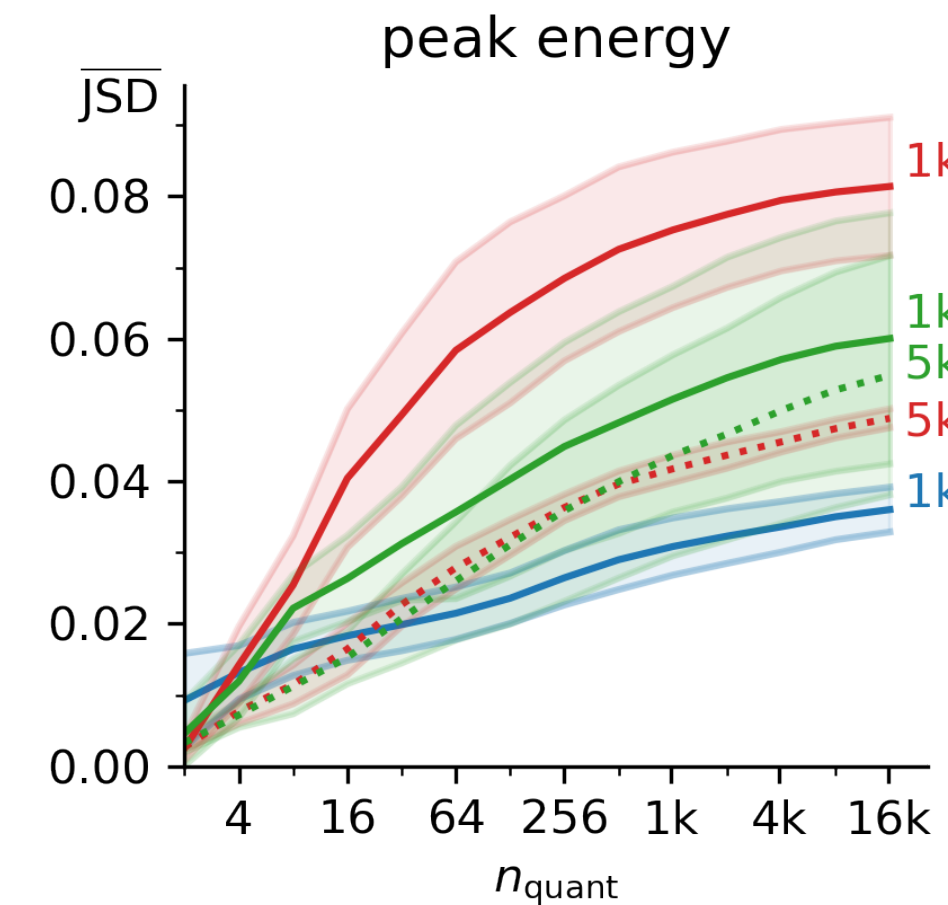
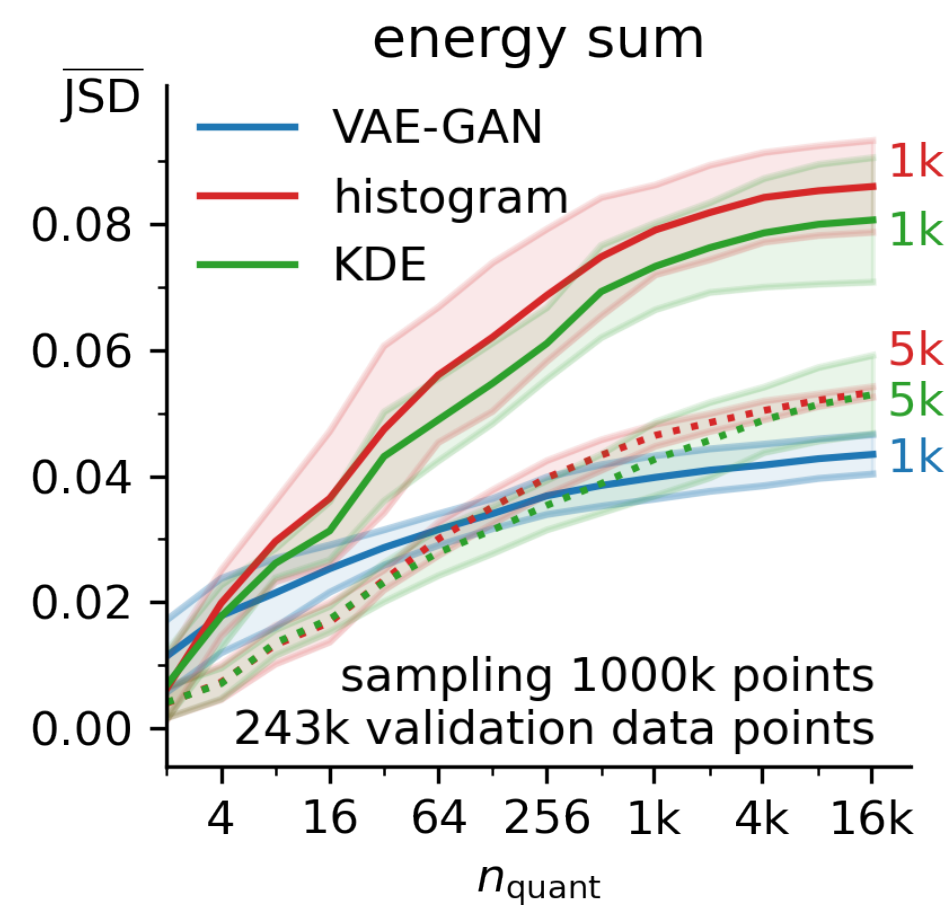


Image-shaped data

Unknown true distribution, limited data

Harder learning task  $\rightarrow$  training on multiple training set sizes unfeasible

# Conclusion

- What about # samples? How many new points should we generate from a generative model?
  - Depends on GAN setup and problem

- For high-scale observables (e.g. *mean, standard deviation, low moments*) generative network limited to the amount of training data
- For a smooth interpolation (e.g. *segments of the distribution, integrated quantities*) a generative networks outperform even higher numbers of data

# References



- [0]: P. Calafiura, J. Catmore, D. Costanzo, and A. Di Girolamo, “ATLAS HLLHC Computing Conceptual Design Report,” CERN, Geneva, Tech. Rep., Sep 2020. [Online]. Available: <https://cds.cern.ch/record/2729668>
- [1]: ILD Concept Group, H. Abramowicz et al., *International Large Detector: Interim Design Report*, 3, 2020.
- [2]: L. de Oliveira, M. Paganini, and B. Nachman, “Learning particle physics by example: Location-aware generative adversarial networks for physics synthesis,” *Computing and Software for Big Science*, vol. 1, no. 1, Sep 2017. [Online]. Available: <http://dx.doi.org/10.1007/s4178101700046>
- [3]: M. Paganini, L. de Oliveira, and B. Nachman, “Calogan: Simulating 3d high energy particle showers in multilayer electromagnetic calorimeters with generative adversarial networks,” *Physical Review D*, vol. 97, no. 1, Jan 2018. [Online]. Available: <http://dx.doi.org/10.1103/PhysRevD.97.014021>
- [4]: A. B. L. Larsen, S. K. Sønderby, H. Larochelle, and O. Winther, “Autoencoding beyond pixels using a learned similarity metric,” in *Proceedings of the 33rd International Conference on International Conference on Machine Learning Volume 48*. JMLR.org, 2016, p. 1558–1566.