

Normalizing Flows for Physics Data Analyses

Jan Gavranovič (jan.gavranovic@ijs.si)^{1,2}

Borut Paul Kerševan (borut.kersevan@ijs.si)^{1,2}

¹Jožef Stefan Institute, Ljubljana, Slovenia

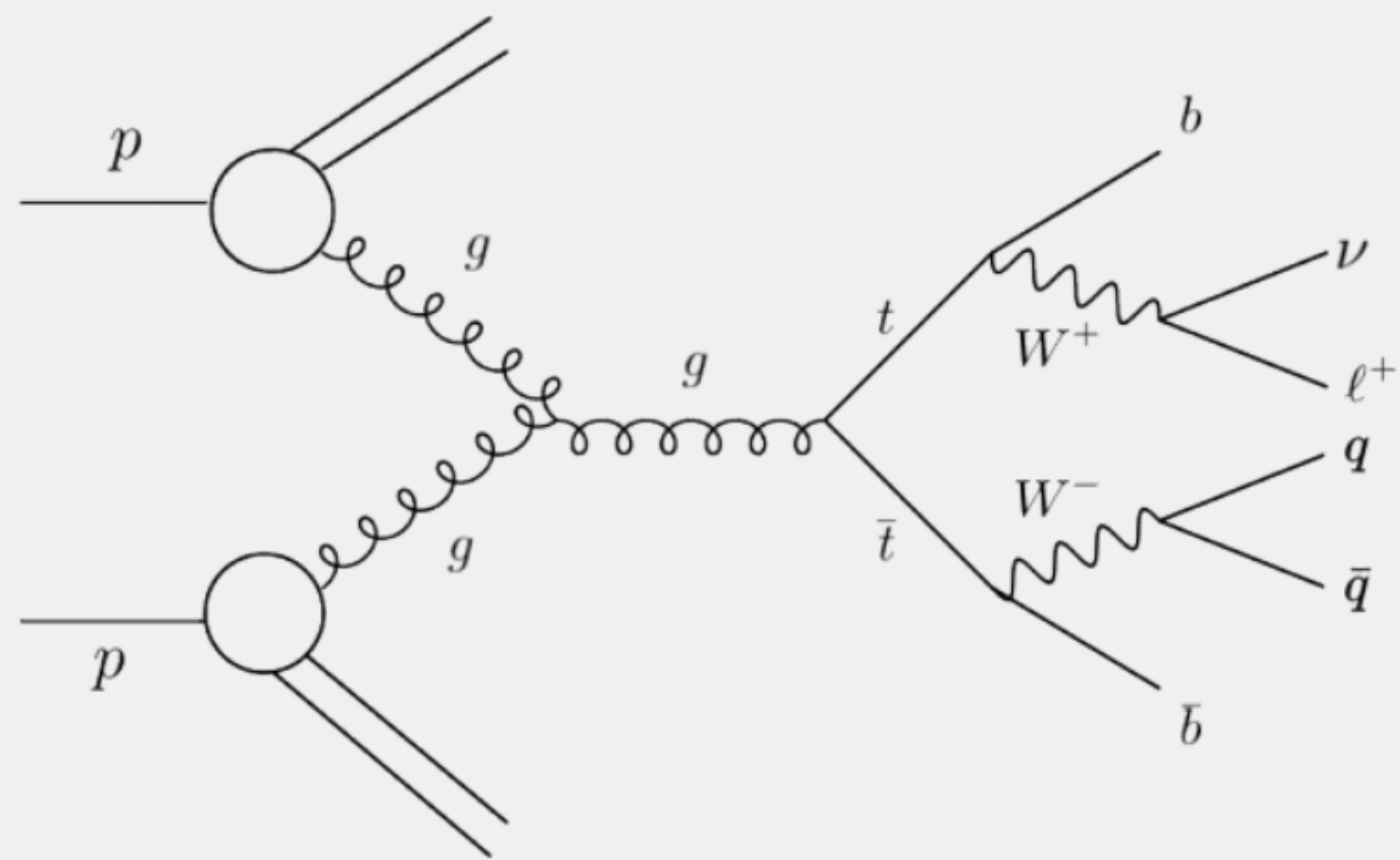
²Faculty of Mathematics and Physics, University of Ljubljana, Slovenia

Introduction

- LHC produces **big data** \Rightarrow MC and analysis need to follow
- Can generative models be used to support physics modeling?**
- Problem:** do not know the true generating data distribution
- Objective:** approximate $p_{\text{data}}(\mathbf{x})$ to enable **infinite sampling**
- Learn true $p_{\text{data}}(\mathbf{x})$ from $\mathbf{x} \in \mathbb{R}^D$ using approximate $p_{\text{model},\theta}(\mathbf{x}) \approx p_{\text{data}}(\mathbf{x})$
- Focus on LHC *analysis-specific* distributions of *final* analysis variables

Higgs Benchmark Dataset

- Publicly available dataset with 11M events and 28 variables
- Binary classification problem: signal (BSM) vs. background ($t\bar{t}$)
- Use as test for LHC *final* event simulation with normalizing flows



- 21 low-level and 7 high-level variables
- Data preprocessing (*feature scaling*) is a crucial step in training
- Task:** train ML model to generate *new background events*

Normalizing Flows (Invertible Neural Networks)

- Two pieces:
 - base distribution $p_u(\mathbf{u})$, typically $\mathcal{N}(\mathbf{u}|\mathbf{0},\mathbf{I})$
 - differentiable transformation $\mathbf{x} = T(\mathbf{u})$ with an inverse $\mathbf{u} = T^{-1}(\mathbf{x})$
- Construct a **flow** by composing together many transformations T :

$$T = T_K \circ \dots \circ T_1 \quad \text{and} \quad T^{-1} = T_1^{-1} \circ \dots \circ T_K^{-1}$$
- Transformations T are (invertible) neural networks with parameters ϕ**
- Generative process:

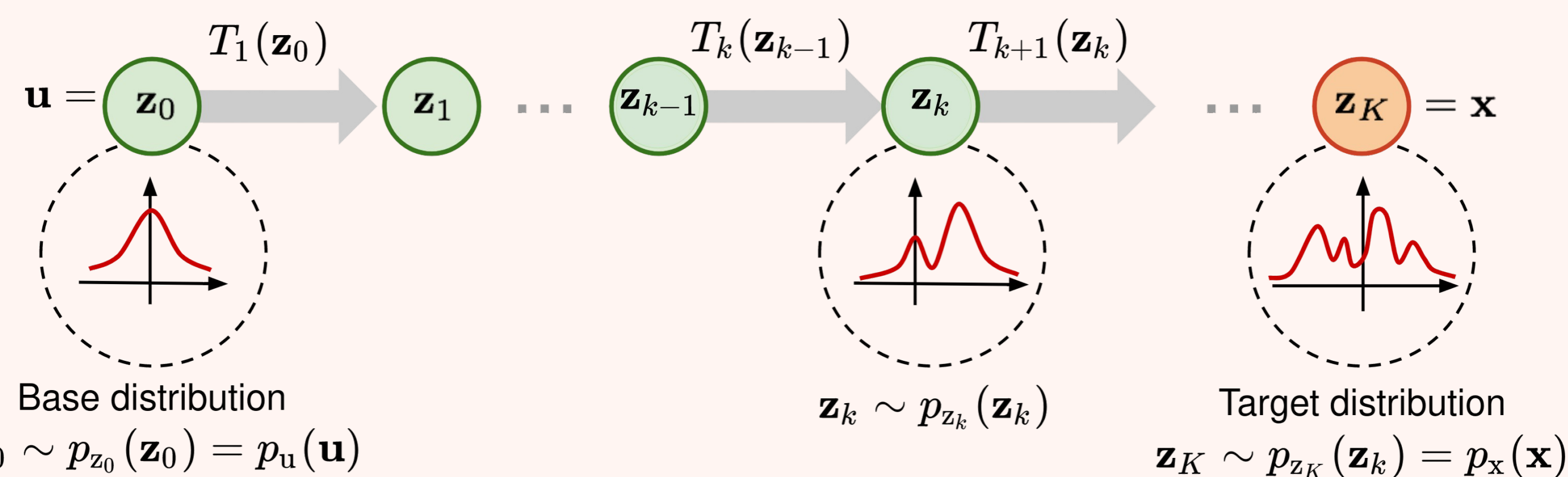
$$\mathbf{x} = T(\mathbf{u}) \approx p_x(\mathbf{x}) \quad \text{with sampling} \quad \mathbf{u} \sim p_u(\mathbf{u})$$

- Density evaluation using change of variables formula:

$$p_x(\mathbf{x}) = p_u(T^{-1}(\mathbf{x})) \left| \det \frac{\partial T^{-1}(\mathbf{x})}{\partial \mathbf{x}} \right|$$

Forward and Inverse Directions

- Forward direction:** $\mathbf{z}_k = T_k(\mathbf{z}_{k-1})$ for $k = 1, \dots, K$ with $\mathbf{z}_0 = \mathbf{u}$ (*infer*)
- Inverse direction:** $\mathbf{z}_{k-1} = T_k^{-1}(\mathbf{z}_k)$ for $k = K, \dots, 1$ with $\mathbf{z}_K = \mathbf{x}$ (*train*)



- Similar to autoencoder: forward mode \Leftrightarrow decoder, backward mode \Leftrightarrow encoder
- Loss function** has two terms (**log-likelihood + log-determinant**):

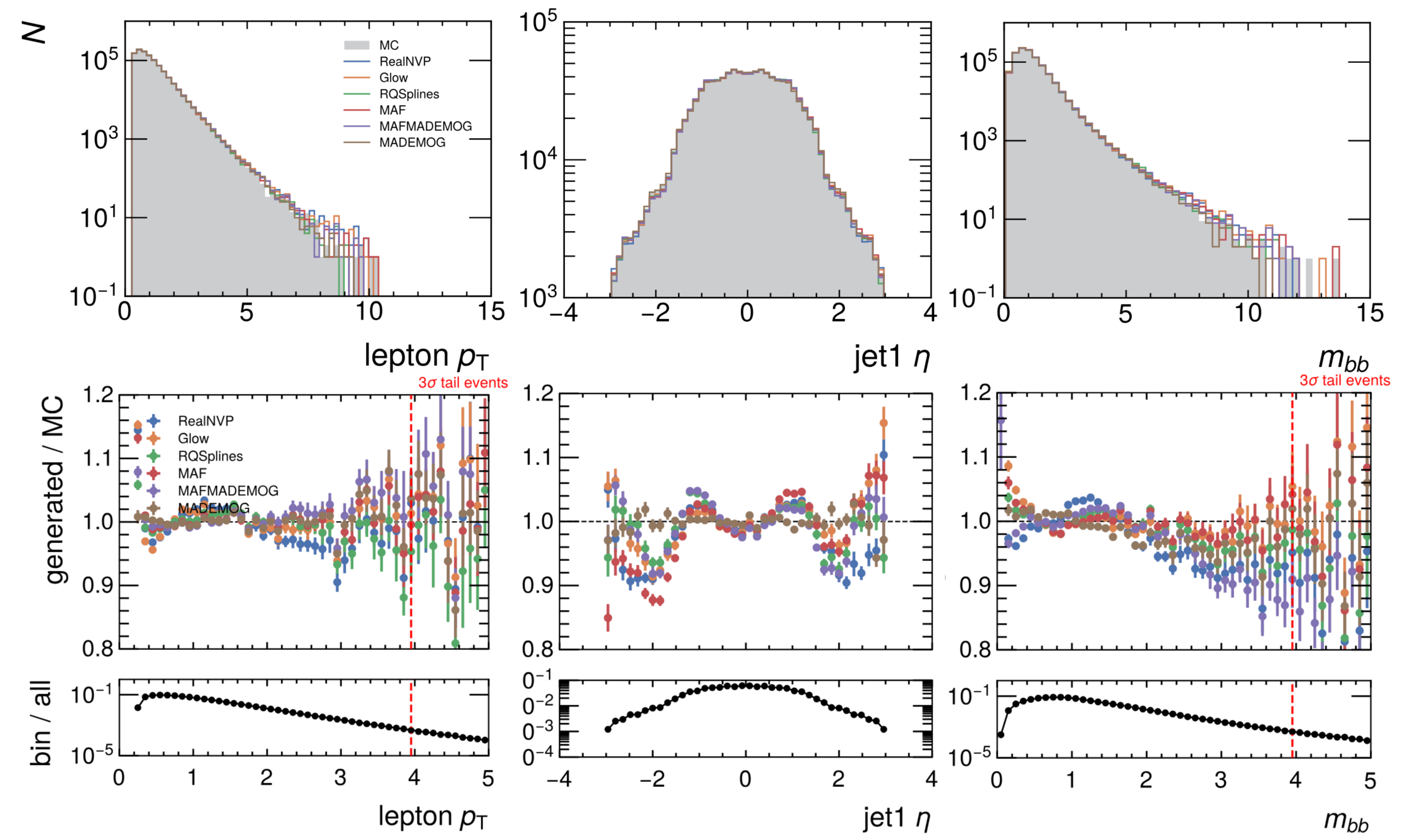
$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{n=1}^N [\log p_u(T^{-1}(\mathbf{x}_n; \phi); \psi) + \log |\det J_{T^{-1}}(\mathbf{x}_n; \phi)|]$$

- Use *gradient descent* to get the best parameters:

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \mathcal{L}(\theta), \quad \theta \equiv \{\phi, \psi\}$$

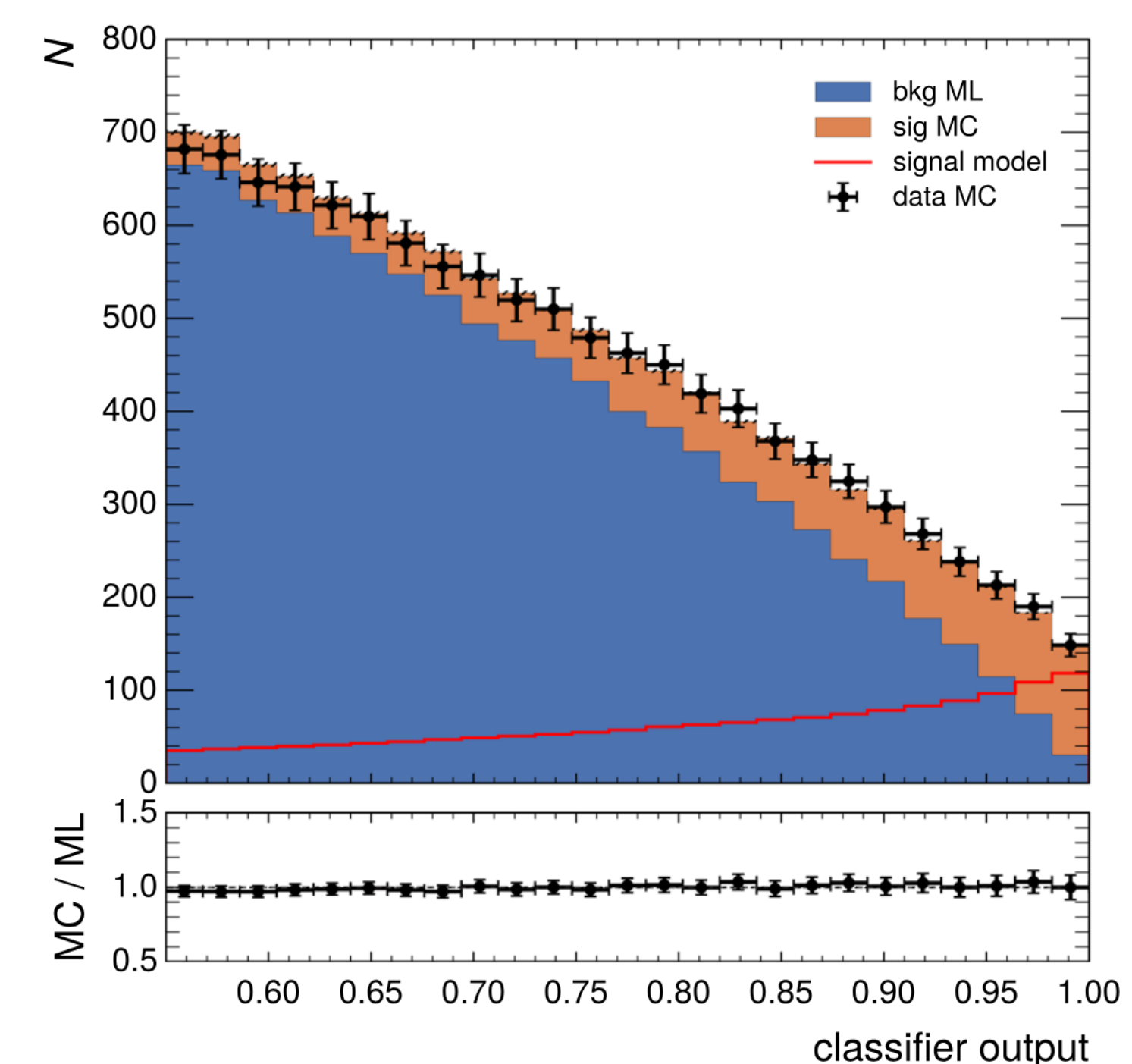
Performance Evaluation

- Comparison of ML *generated* and MC *simulated* distributions
- Best model was selected for the final analysis
- Performance was measured using **statistical distances** and **classifier two sample testing**

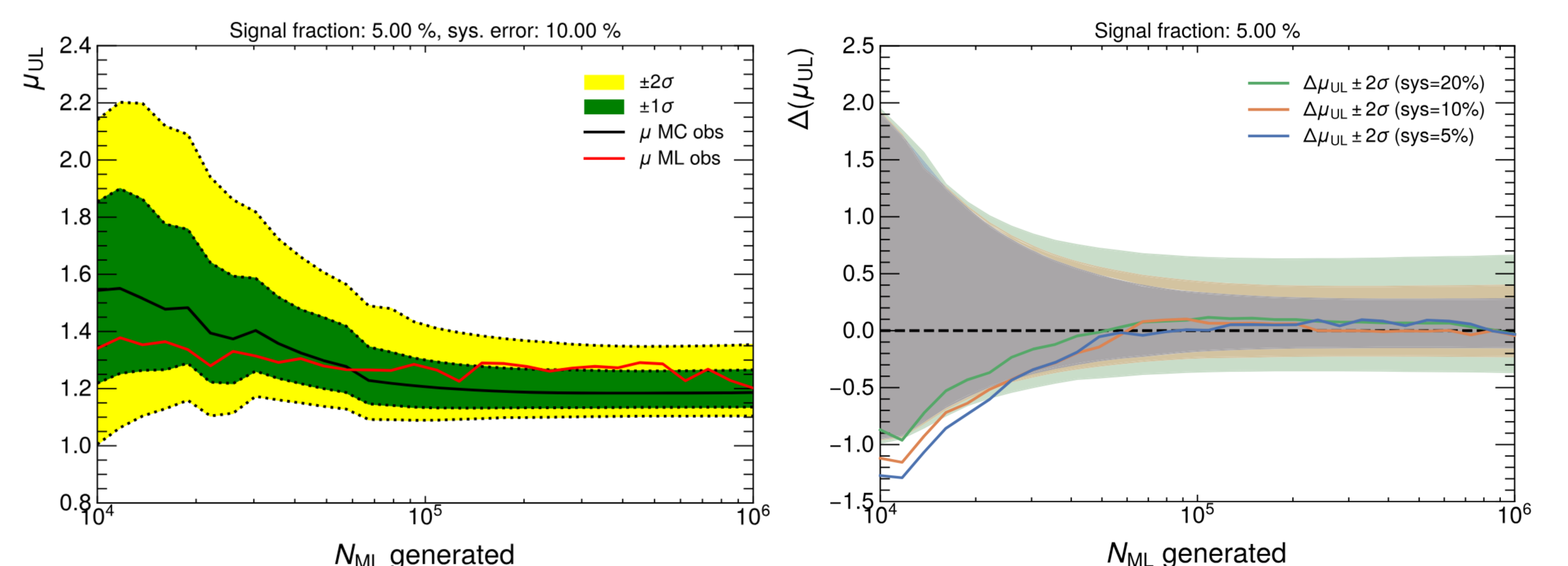


Physics Analysis

- A simplified analysis was performed, involving preselection (baseline cuts) and a **NN-based classifier final selection**



- Upper limits on the signal strength μ** for the likelihood fit to the classifier score distribution as a function of ML-generated events were calculated



Summary

- Generative modeling is a **promising new tool** for physics data analyses
- Needs careful performance evaluation and validation for physics use cases

References

- [1] J. Gavranovič and B. P. Kerševan. Systematic evaluation of generative machine learning capability to simulate distributions of observables at the large hadron collider. *Eur. Phys. J. C*, 84(9):911, 2024.