

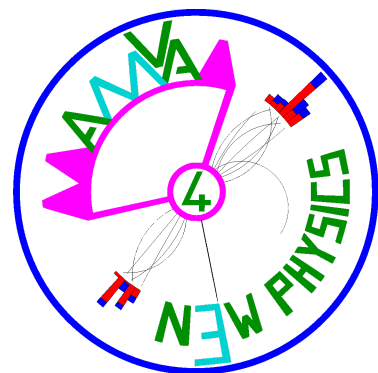
# DIRECT LEARNING OF SYSTEMATICS-AWARE SUMMARY STATISTICS

*Pablo de Castro (@pablodecm) and Tommaso Dorigo (@dorigo)*

6th August 2018 @ XIIIth QCHS Conference (Maynooth University - Ireland)

**Poster Award Lightning Talk**

*Poster & Parallel - Statistical Methods for Physics Analysis in the XXI Century*



AMVA4NewPhysics has received funding from European Union's Horizon 2020 Programme under Grant Agreement number 675440

# USE OF MACHINE LEARNING IN PHYSICS DATA ANALYSIS

Currently undergoing inflation

But beware ...

**Advances in Machine Learning & Applications to QCD**  
An Experimentalist's Perspective

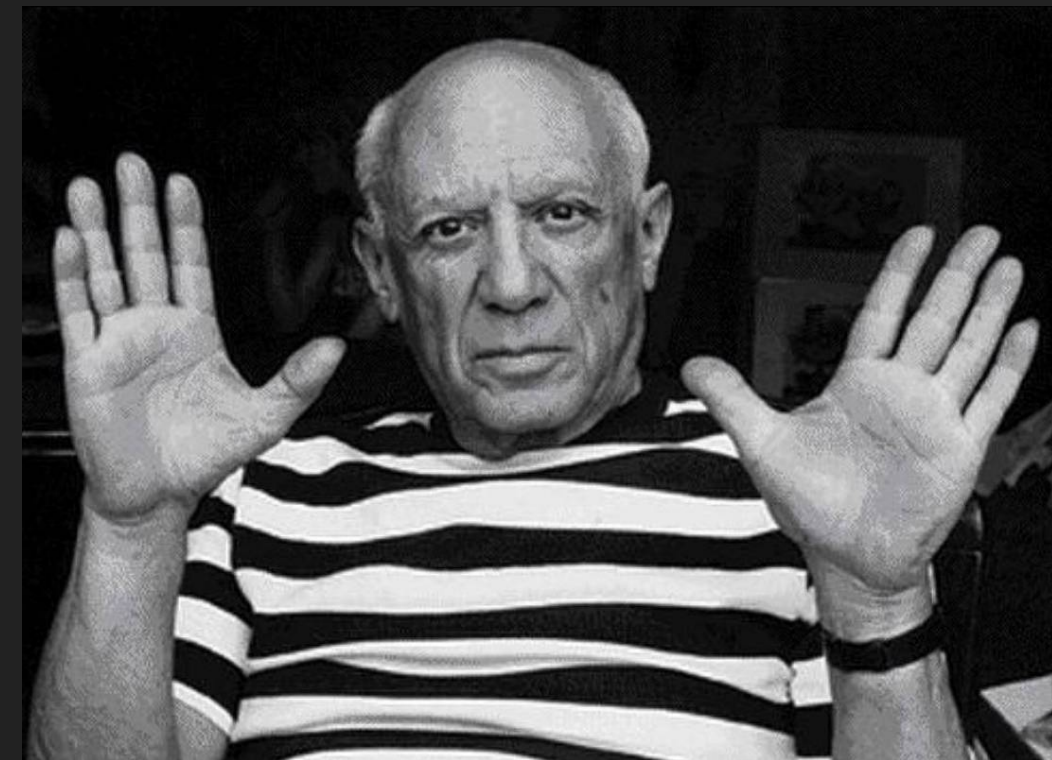
13<sup>th</sup> Quark Confinement and the Hadron Spectrum

Maynooth University

2<sup>nd</sup> August 2018

Lily Asquith

THE ROYAL SOCIETY US  
University of Sussex



Review by Lily Asquith on her [Advances in ML and applications to QCD](#) plenary talk

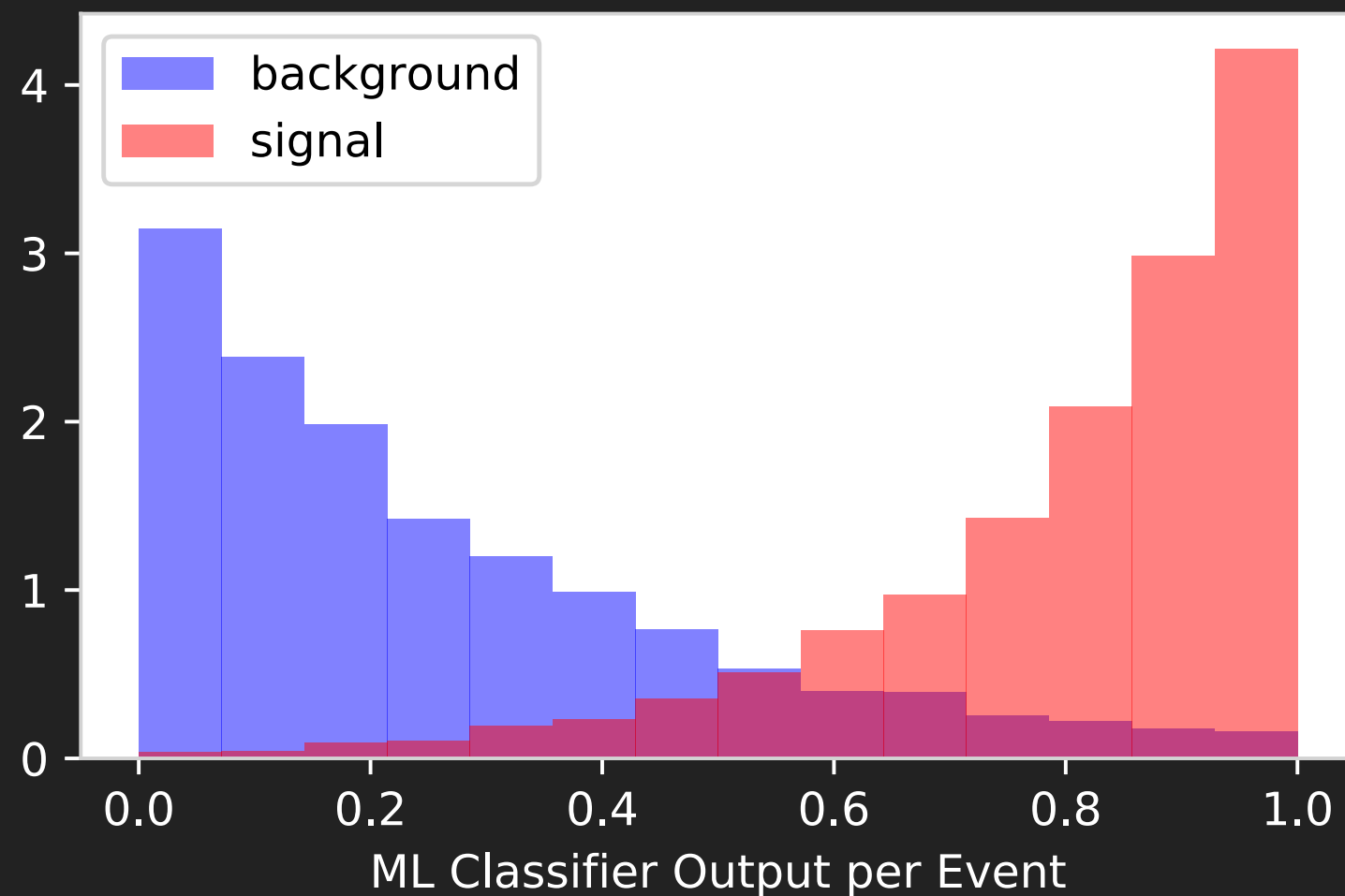
*Computers are useless. They can only give you answers.*

Pablo Picasso (1960s)

ARE ASKING THE RIGHT QUESTIONS AT THE LHC ANALYSES?

# MACHINE LEARNING WITHIN LHC ANALYSES

Most common approach → supervised learning classification trained on simulation



## Event-by-Event Signal vs Background

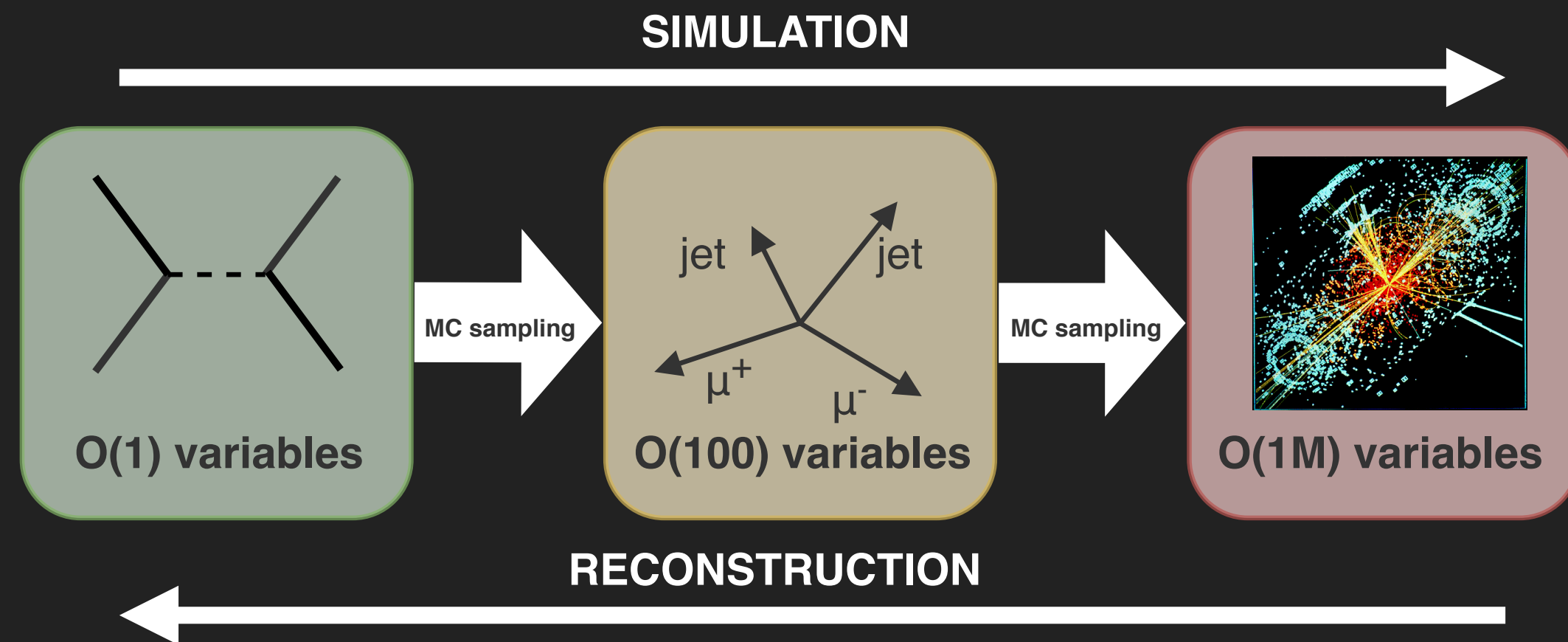
- [Higgs Kaggle challenge \(2014\)](#)
- Almost every other LHC analysis

**IS IT REALLY A CLASSIFICATION PROBLEM?**

**NO! IT IS A STATISTICAL INFERENCE PROBLEM**

# $p(\mathbf{x}|\text{model})$ IS NOT KNOWN AT LHC EXPERIMENTS

Samples under different hypotheses can be simulated via complex physics-based MC programs but  $p(\mathbf{x})$  cannot be directly evaluated  $\rightarrow$  **LIKELIHOOD-FREE INFERENCE**



good approximations of  $p(\mathbf{x})$  are unachievable due to curse of dimensionality

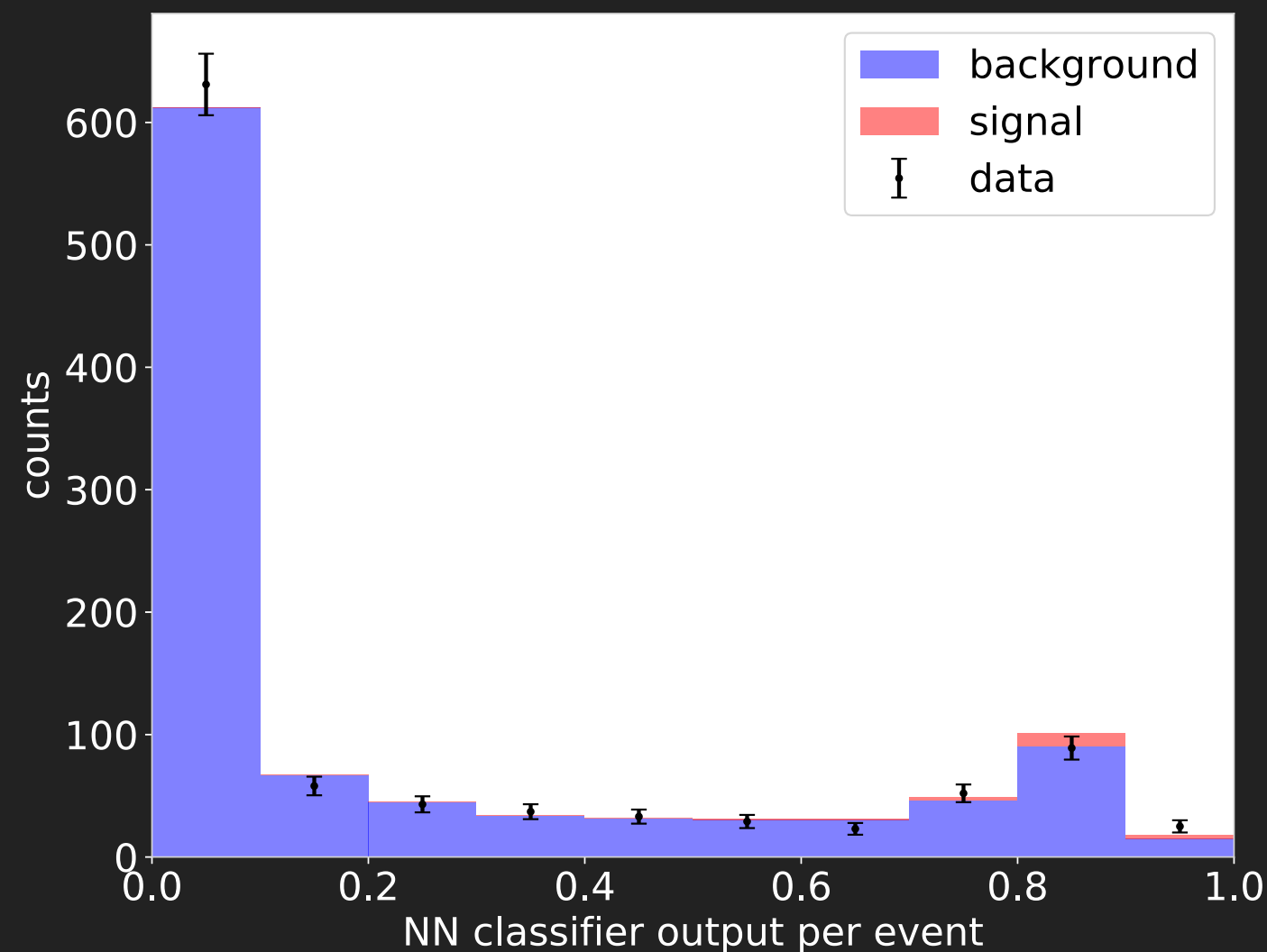
**DIMENSIONALITY REDUCTION  $\mathbf{R}^n \rightarrow \mathbf{R}^{O(1)}$  (SUMMARY STATISTIC)**

**KEEPING AS MUCH USEFUL INFORMATION FOR INFERENCE AS POSSIBLE**

# CLASSIFIER-BASED INFERENCE

A ML classifier trained on simulation  $d(\mathbf{x})$  is an approximation of  $p_s(\mathbf{x})/p_b(\mathbf{x})$

How can it be used for statistical inference from observed data  $\mathcal{D}$ ?



1-D  $\rightarrow$  cut or histogram to build a Poisson counts non-parametric likelihood

$$\mathcal{L}(\mu) = \prod_{i \in \text{bins}} \text{Pois}(n_i | \mu \cdot s_i + b_i)$$

which can be used for further inference, such as measuring  $\mu$  given observed  $\mathcal{D}$

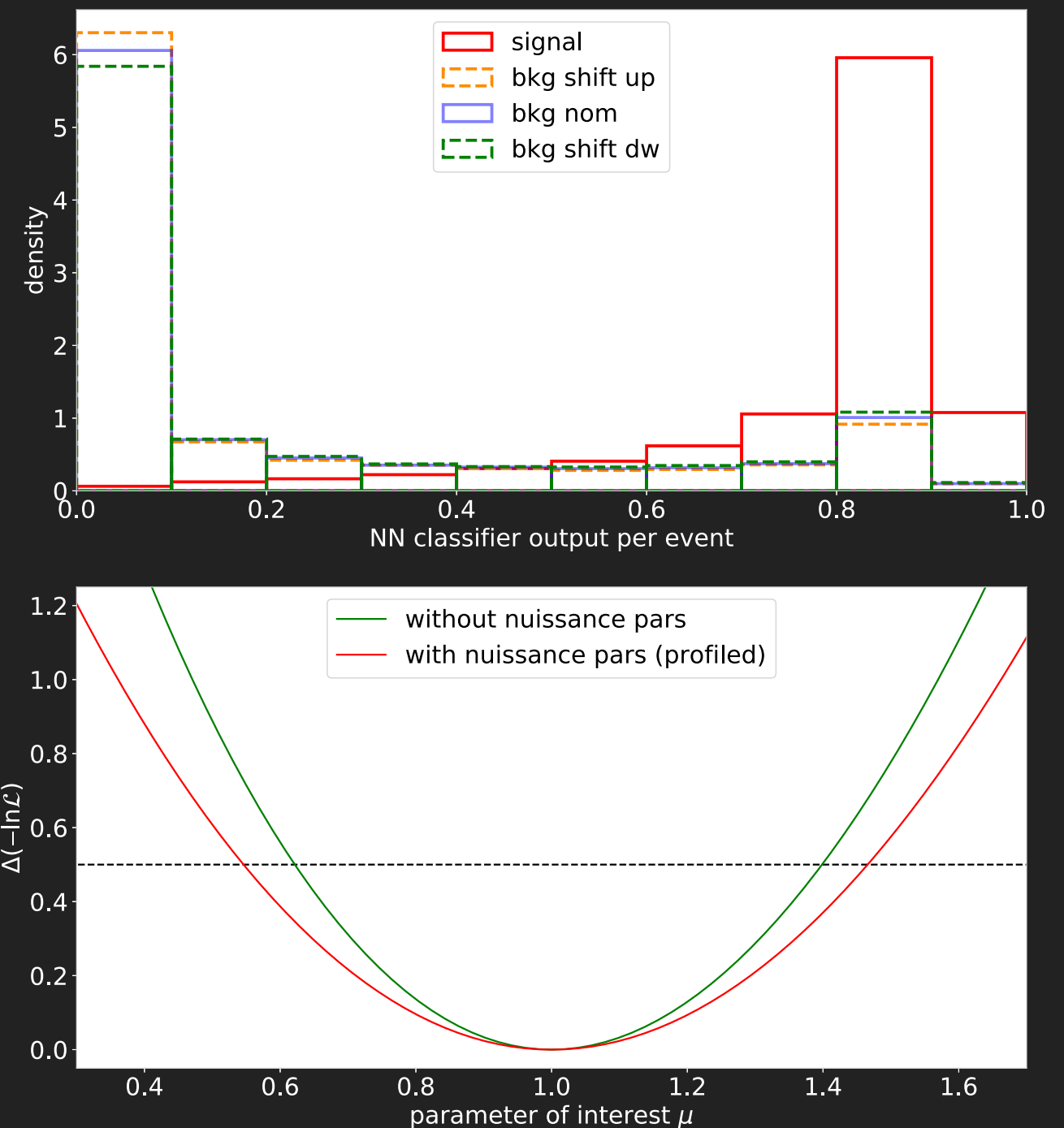
A more direct exploitation of the likelihood ratio approximation for inference is carried out in "[Approximating Likelihood Ratios with Calibrated Discriminative Classifiers](#)" by K. Cranmer et al. That work was further extended in [arXiv:1805.12244](#) (and cited articles therein) by J. Brehmer et al to use also the joint score.

# MODELLING UNCERTAINTIES DEGRADE INFERENCE

Simulations are imperfect, mainly due to the limited information of the system being modelled

Lack of knowledge for inference accounted by additional uncertain parameters (nuisance parameters  $\nu$ )

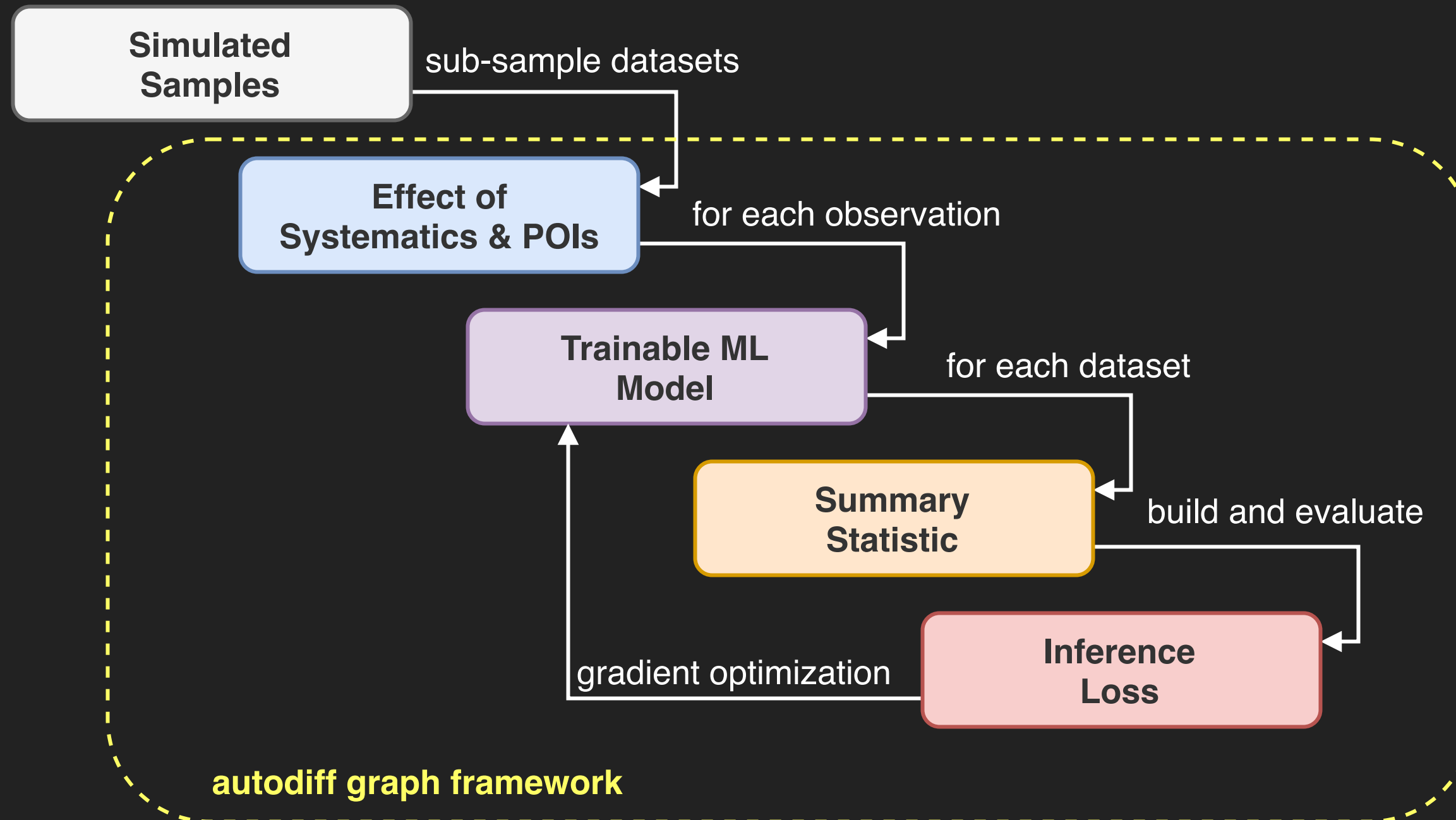
Causes a degradation of classifier-based inference, leading to larger measurement error  $\rightarrow$  **systematic uncertainties**



## UPPER LIMIT OF ML USEFULNESS IN LHC ANALYSES

Classifiers can be made pivotal as described in "[Learning to Pivot](#)" by G. Louppe et al. A review/benchmarks on how to deal with systematics when using machine learning can be found in [Adversarial learning to eliminate systematic errors: a case study in High Energy Physics](#) by Victor Estrade et al NIPS2017.

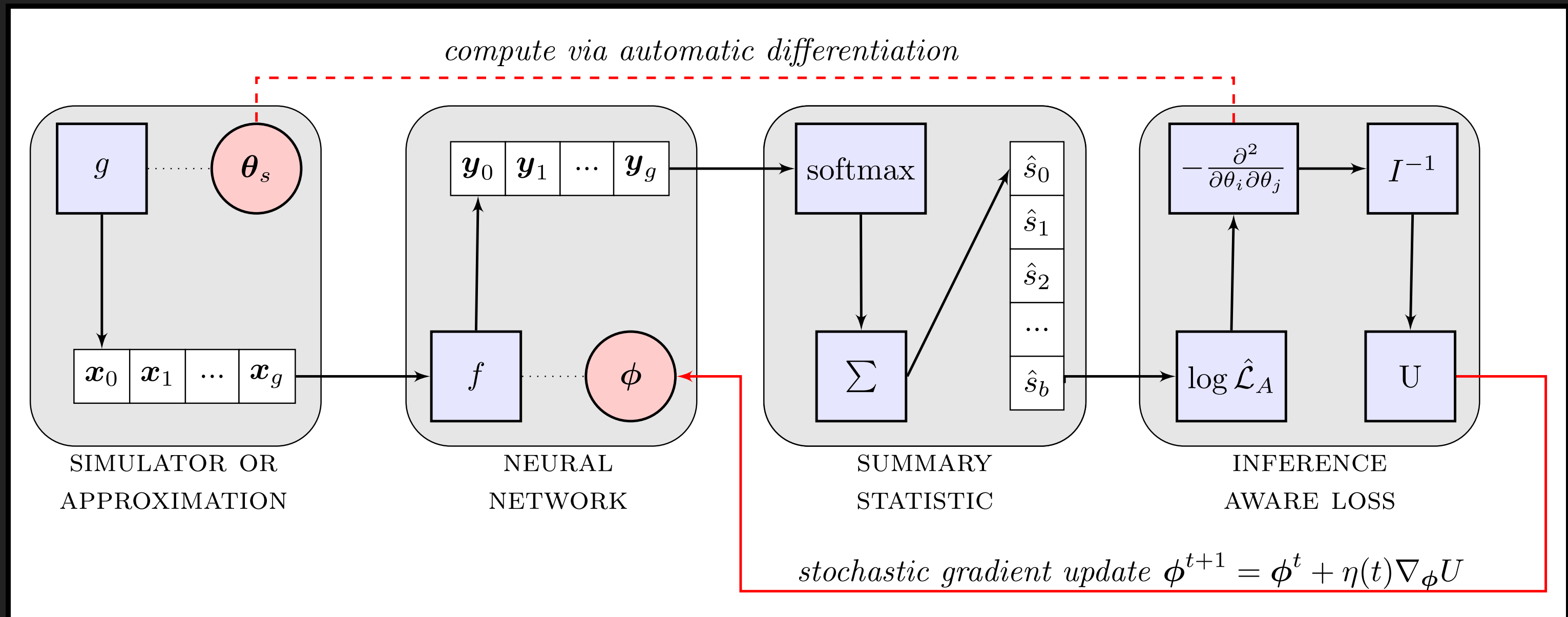
# END-TO-END DIFFERENTIABILITY FOR LHC ANALYSES



Within this general framework, several approaches are possible, focus here is

**DIRECT LEARNING OF SYSTEMATICS-AWARE SUMMARY STATISTICS**

# INFERNO: INFERENCE-AWARE NEURAL OPTIMISATION



*differentiable approx. covariance matrix of statistical model  $\rightarrow$  SGD optimisation*  
 check [arxiv.org/abs/1806.04743](https://arxiv.org/abs/1806.04743) for a detailed mathematical description



# INFERENCE-AWARE LOSS FUNCTION

## Inference-Aware loss:

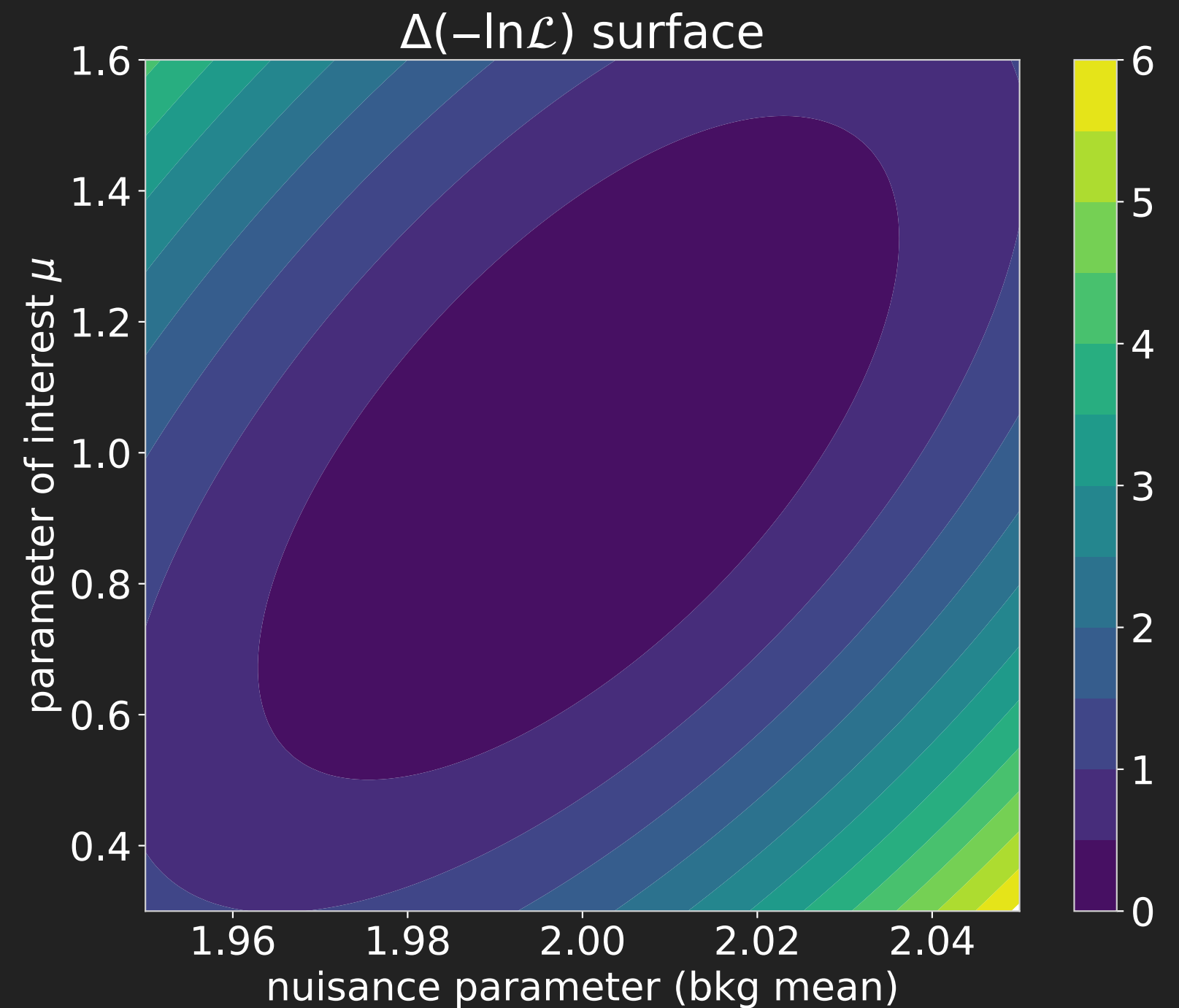
$$\text{loss} \approx \text{Var}(\mu) \quad (\text{expected})$$

- final analysis figure-of-merit
- accounts for systematics unc.
- extended for other parameters

## Classification-based loss:

$$\text{loss} = - \sum_i k_i \log y_i \quad (\text{c. entropy})$$

- approximates  $p_s(\mathbf{x})/p_b(\mathbf{x})$
- does not account for systematic unc.



# SYNTHETIC EXAMPLE RESULTS → IT WORKS!

Applied on 2D Gaussian two-component mixture toy problem, with an unknown background mean in one of the coordinates → **one nuisance parameter  $\nu$**

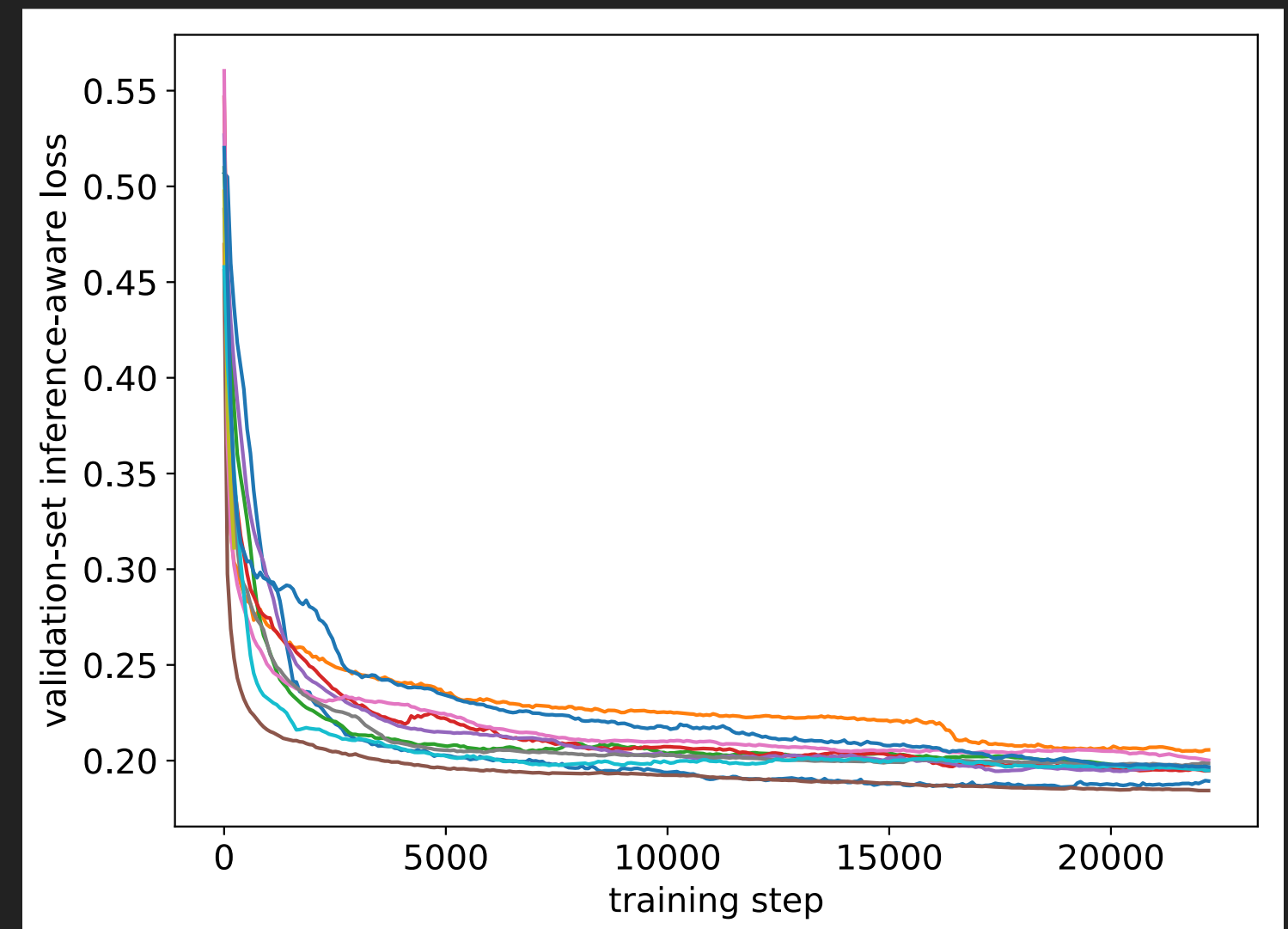
Approach converges consistently to low  $\text{Var}(\mu)$  solutions ind. of initialisation

Expected signal strength uncertainties computed using the validation set (average of 10 random initialisations):

- cross-entropy:  $0.444 \pm 0.003$
- **inference-aware:  $0.437 \pm 0.008$**

Now working on higher-dimensional problems with more nuisance parameters

**BETTER/EQUAL THAN CLASSIFICATION**



# THANK YOU FOR YOUR ATTENTION!

The screenshot shows the arXiv.org interface for the preprint 'INFERNO: Inference-Aware Neural Optimisation'. At the top left is the Cornell University Library logo. The top right features a search bar with the text 'Search or Article ID' and a dropdown menu set to 'All fields'. Below the search bar is a navigation breadcrumb: 'arXiv.org > stat > arXiv:1806.04743'. The main content area is titled 'Statistics > Machine Learning' and contains the preprint title 'INFERNO: Inference-Aware Neural Optimisation' by Pablo de Castro and Tommaso Dorigo, submitted on 12 Jun 2018. The abstract describes complex computer simulations for data modeling. The subjects listed are Machine Learning (stat.ML), Machine Learning (cs.LG), High Energy Physics - Experiment (hep-ex), Data Analysis, Statistics and Probability (physics.data-an), and Methodology (stat.ME). The submission history shows it was sent from Pablo De Castro Manzano on Tue, 12 Jun 2018. On the right side, there are sections for 'Download:' (PDF and other formats), 'Current browse context:' (stat.ML with navigation links), 'Change to browse by:' (listing categories like cs, hep-ex, physics, etc.), 'References & Citations' (listing INSPIRE HEP and NASA ADS), and 'Bookmark' with social media icons.

if interested on technique, more details on preprint [arxiv.org/abs/1806.04743](https://arxiv.org/abs/1806.04743)  
feedback is greatly welcomed (DM [@pablodecm](https://twitter.com/pablodecm) or [pablo.decastro@cern.ch](mailto:pablo.decastro@cern.ch) )