

# Goodness of fit by Neyman-Pearson testing

G. Grosso<sup>1</sup>, M. Letizia, A. Wulzer, M. Pierini

<sup>1</sup>University and INFN of Padova

Mainly based on:

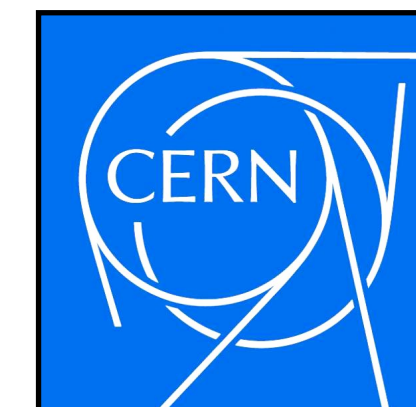
**!! NEW !!** [2305.14137](#) (Grosso, Letizia, Pierini, Wulzer)

**!! NEW !!** [arXiv:2303.05413](#) (Lai, Grosso, Letizia et al.)

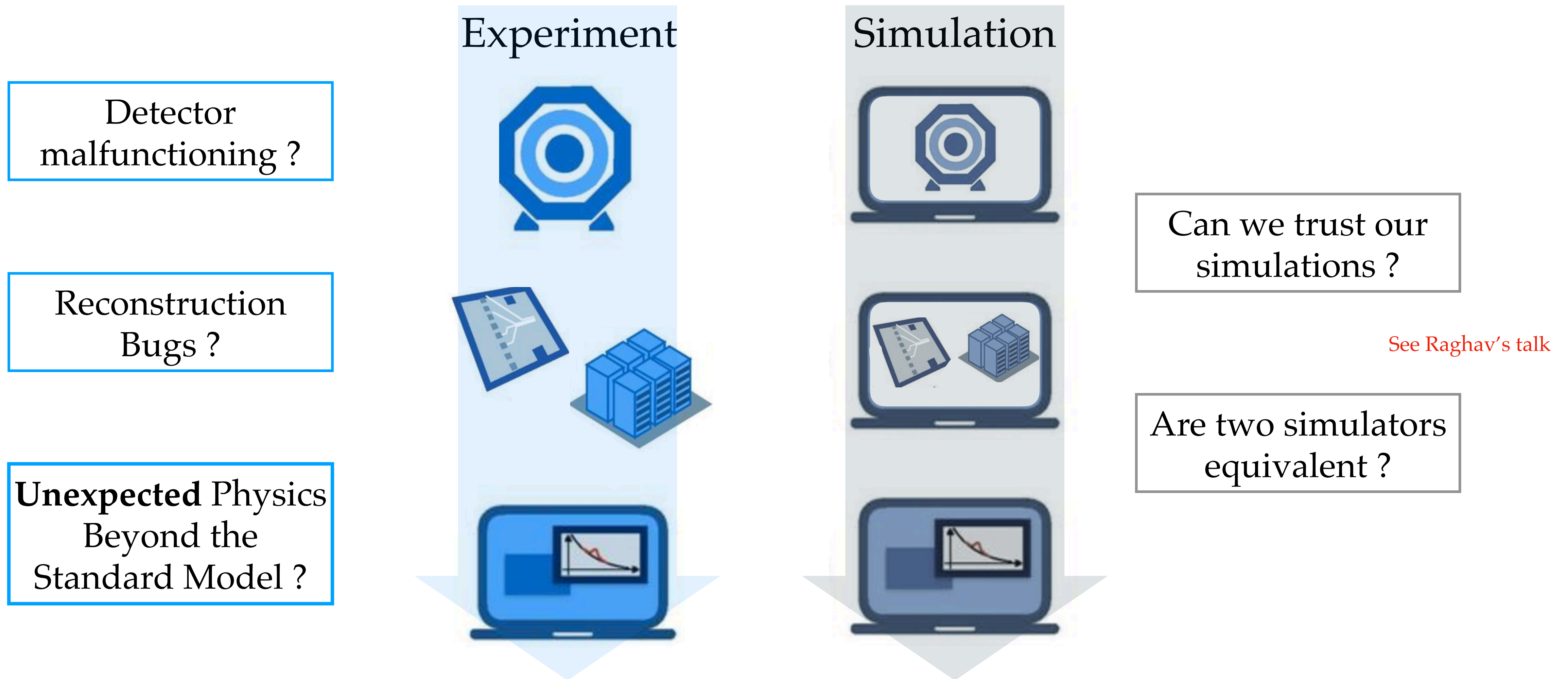
[Eur. Phys. J. C 82, 879 \(2022\)](#) (Letizia, Grosso, Pierini Wulzer et al.)

[Eur. Phys. J. C 81, 89 \(2021\)](#) (d'Agnolo, Grosso, Pierini, Wulzer, Zanetti)

[Phys. Rev. D 99, 015014](#) (d'Agnolo, Wulzer)

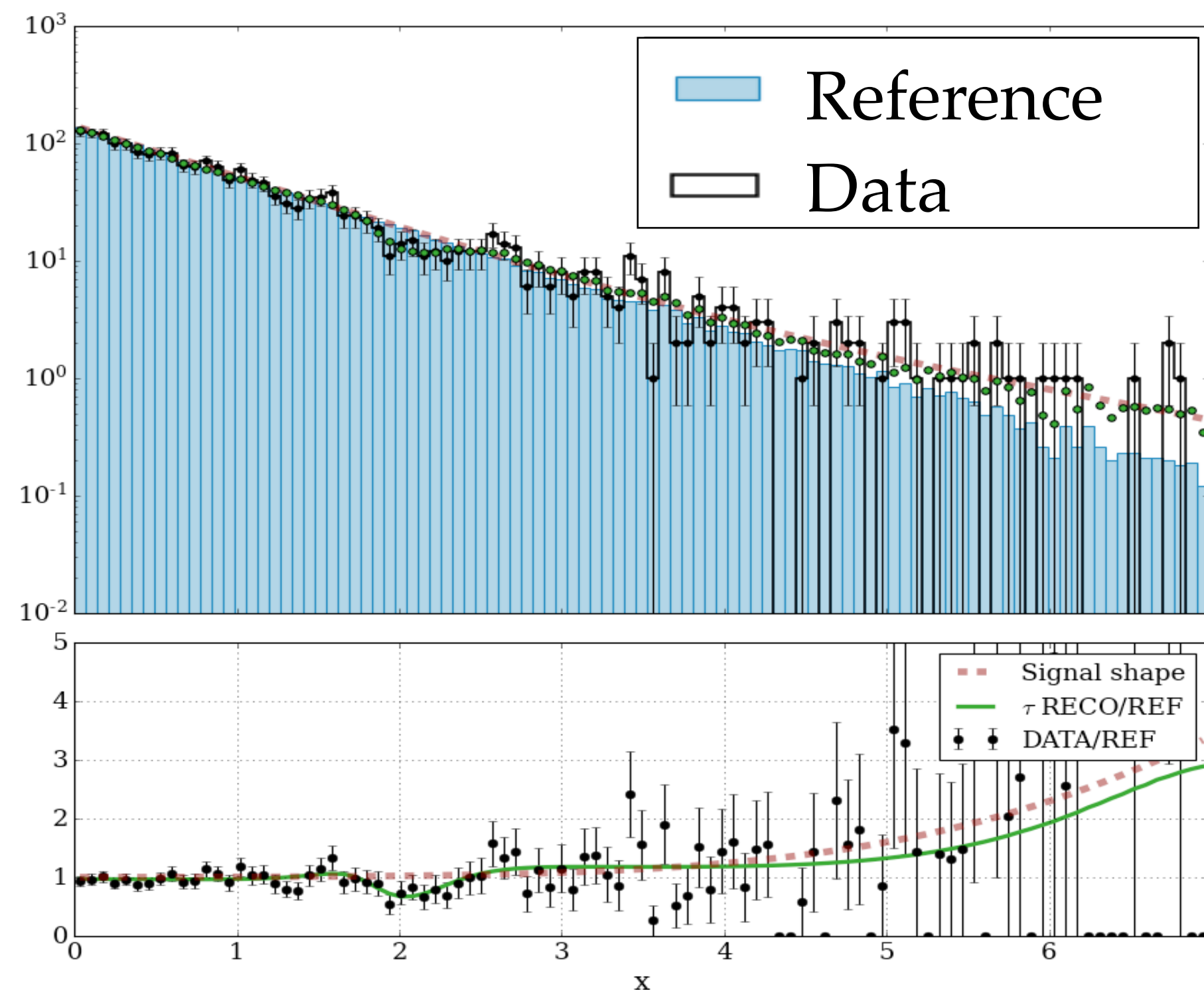


# How well do we understand the LHC data?



# How well do we understand the LHC data?

These are problems of **Goodness of Fit**:



- **Data:** experimental measurements of the natural process  $\{x_i\}_{i=1}^{N_D}$

- **Reference model:** expected nominal behaviour of the data

(Standard Model, normal operating condition of a detector...).

Most of the times not known in close form:  
Reference sample  $\{y_i\}_{i=1}^{N_R} \rightarrow$  **2sample test**

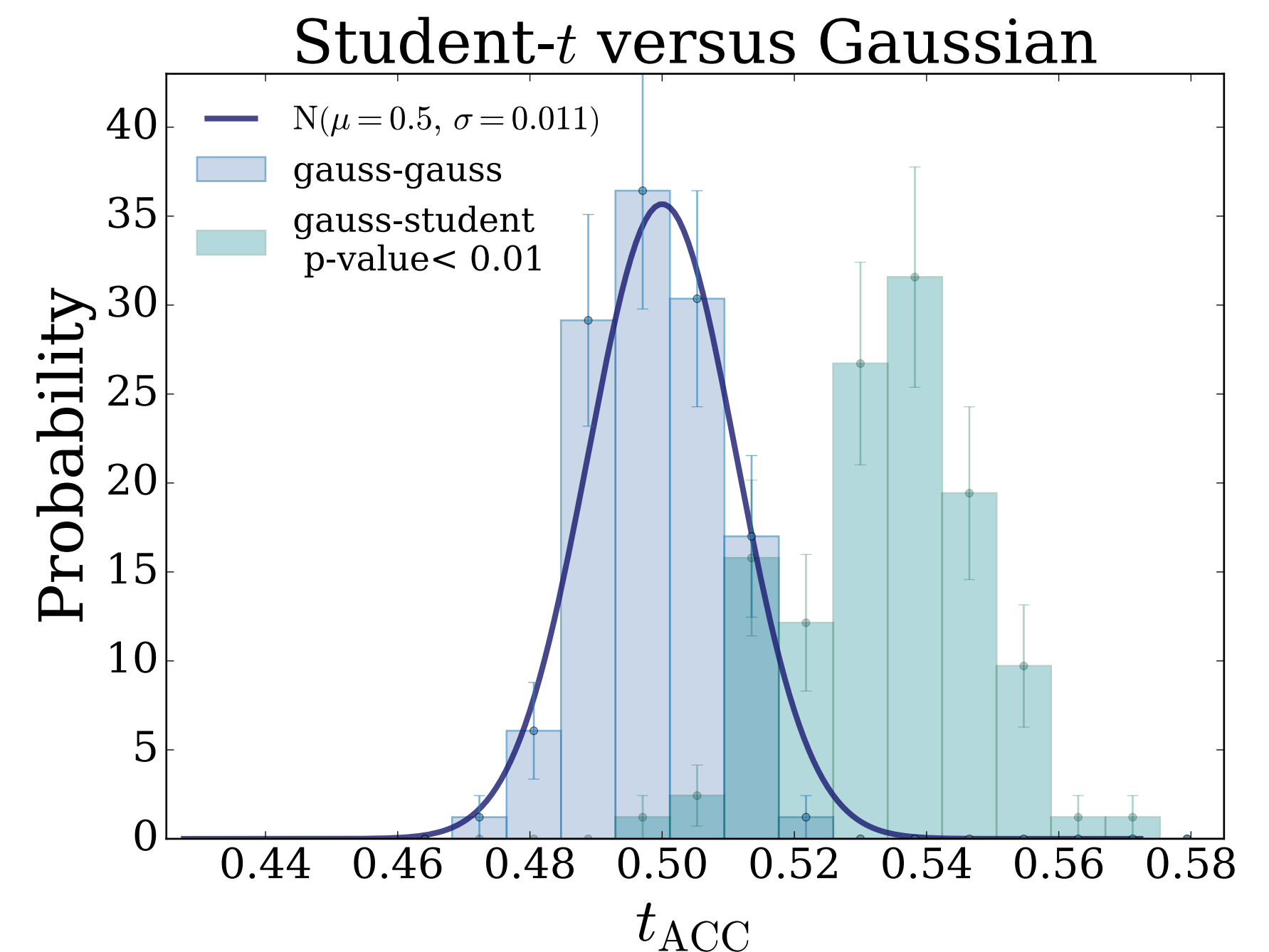
At LHC: multi-dimensional, large statistics samples  $\rightarrow$  **Machine Learning approaches**

# Classifier-based GOF approaches

[Friedman \(2003\)](#): Training a classifier to tell apart two samples (data and reference samples).

Main features:

- Balanced problem ( $N_D = N_R$ )
- Train-test split (out of sample evaluation)
- Test:
  - Classification metrics: accuracy, AUC ([Charkavarti et al. \(2021\)](#), [Lopez et al. \(2017\)](#) )
  - Standard 1D GOF on the classifier output: classifier for dimensionality reduction ([Friedman \(2003\)](#) )
- Calibration and  $p$ -value:
  - Toy experiments (training a new model each time)
  - The single value of test statistic is not sufficient without a fair comparison to the test statistic distribution under the null

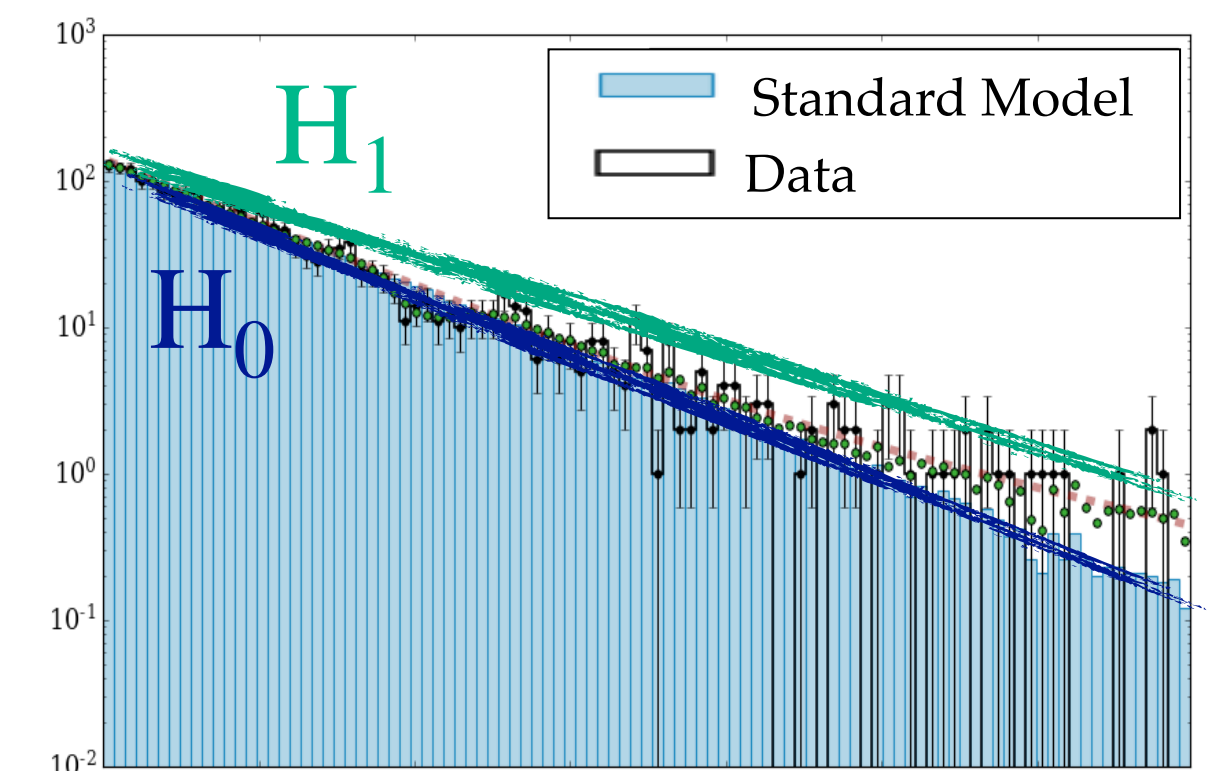
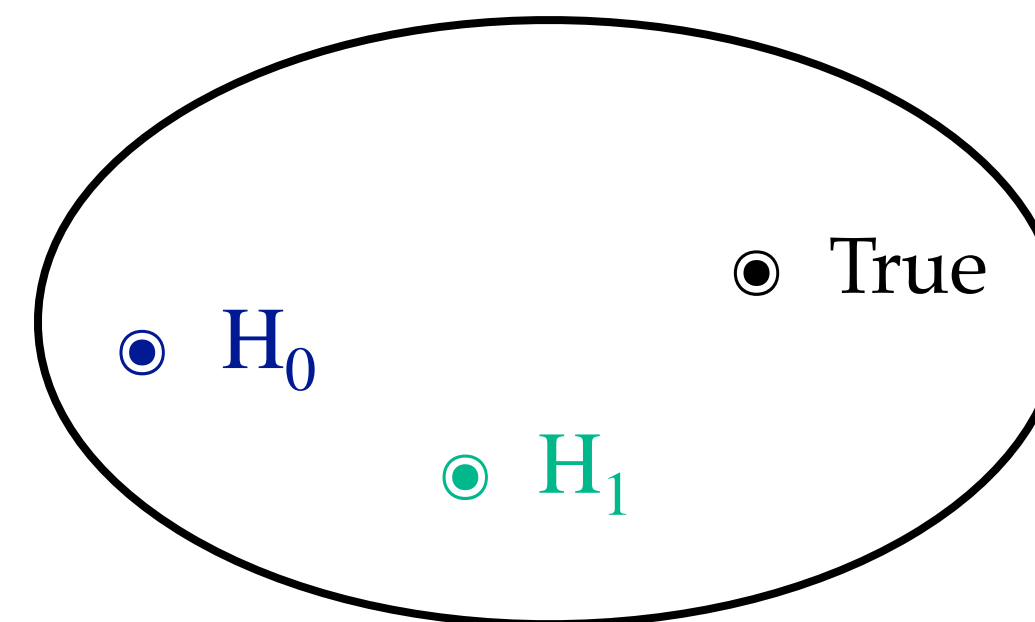


# Neyman-Pearson testing

Comparing the null hypothesis (Reference) with an alternative.  
Traditionally employed in BSM searches at the LHC.

$$t(\mathcal{D}) = 2 \log \frac{\mathcal{L}(\mathcal{D} | H_1)}{\mathcal{L}(\mathcal{D} | H_0)}$$

- Choosing an alternative is mandatory
- The alternative defines the landscape signals that the test is sensitive to.
- Sensitivity (and **optimality**) are guaranteed (according to Neyman and Pearson) only if the data do follow the chosen alternative



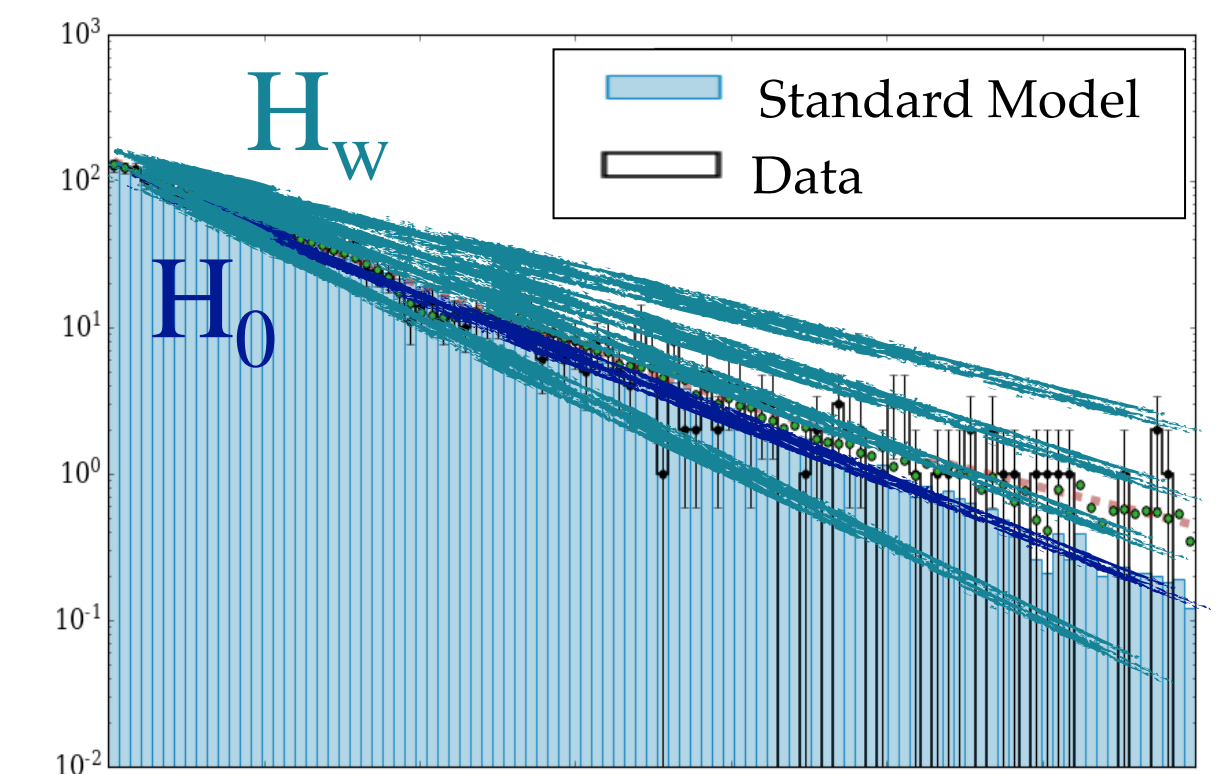
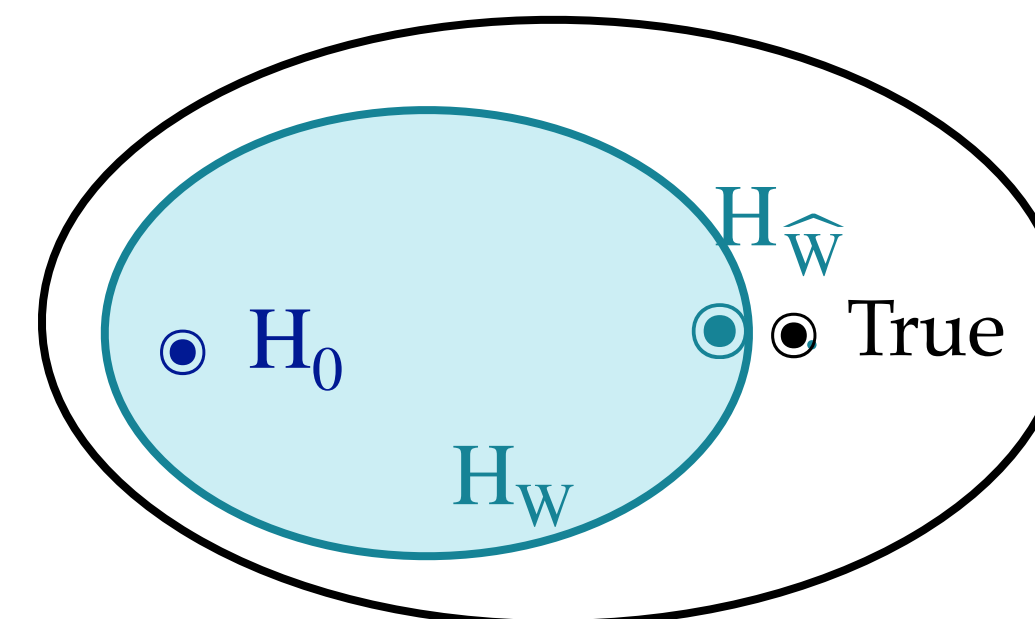
# Neyman-Pearson testing

Comparing the null hypothesis (Reference) with an alternative.

**Release** the assumptions on the alternative

$$t(\mathcal{D}) = \max_{\mathbf{w}} \left[ 2 \log \frac{\mathcal{L}(\mathcal{D} | H_{\mathbf{w}})}{\mathcal{L}(\mathcal{D} | H_0)} \right]$$

- **Expand** the family of alternatives to increase the chance of containing the True data distribution



# Neyman-Pearson testing

## The problem of sensitivity

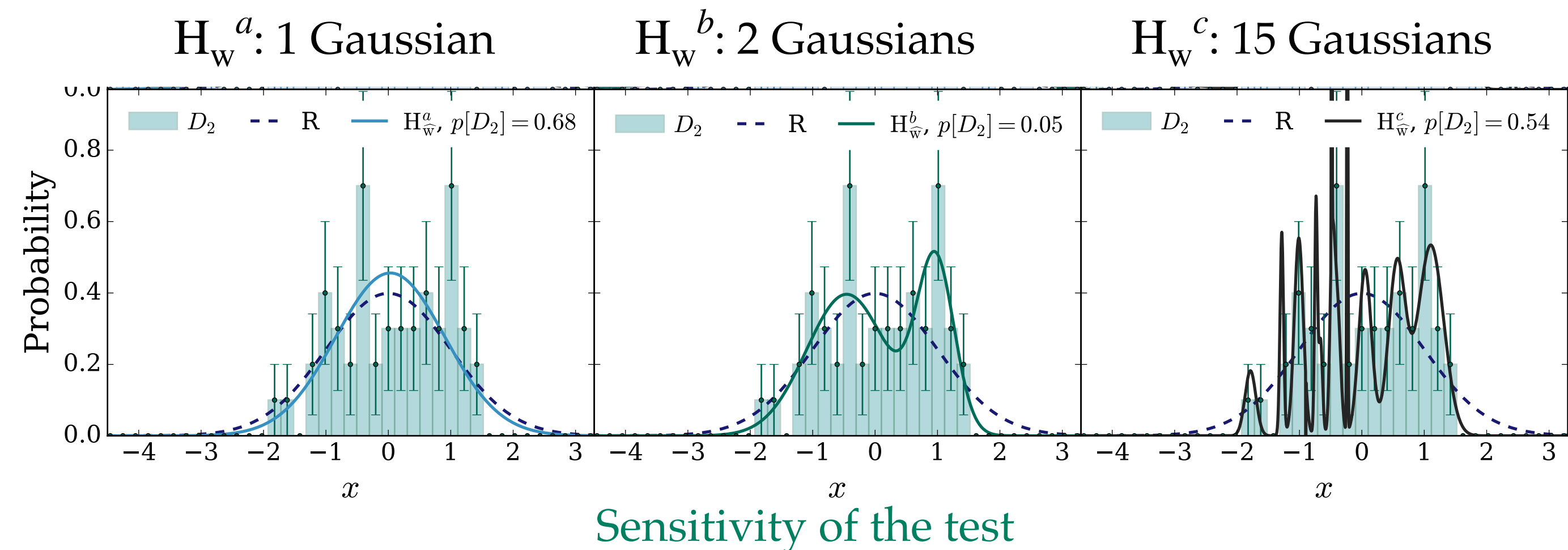
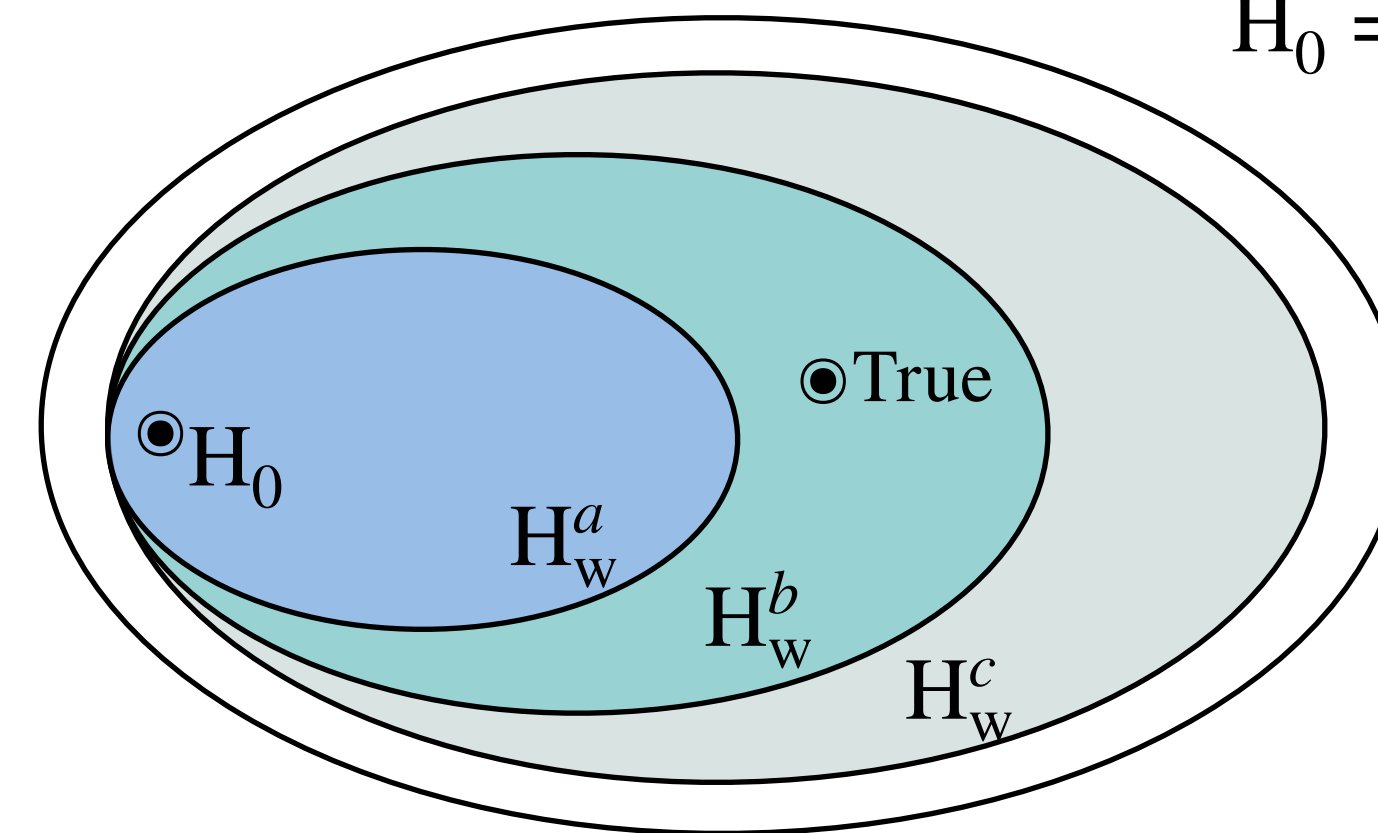
The **specificity** of the alternative determines the sensitivity of the test.

- **Too complex** models are over-sensitive to statistical fluctuations, independently on the nature of the data
- **Too simple** models may not contain an element that approximate the True hypothesis well enough

A successful GoF strategy based on Neyman-Pearson testing should balance between **flexibility** and **regularisation**.

$$D \sim \frac{1}{2}\text{Gauss}(\mu = -1, \sigma = 1) + \frac{1}{2}\text{Gauss}(\mu = 1, \sigma = 1)$$

$$H_0 = \text{Gauss}(\mu = 0, \sigma = 1)$$



Sensitivity of the test

# Neyman-Pearson testing

Comparing the null hypothesis (Reference) with an alternative.

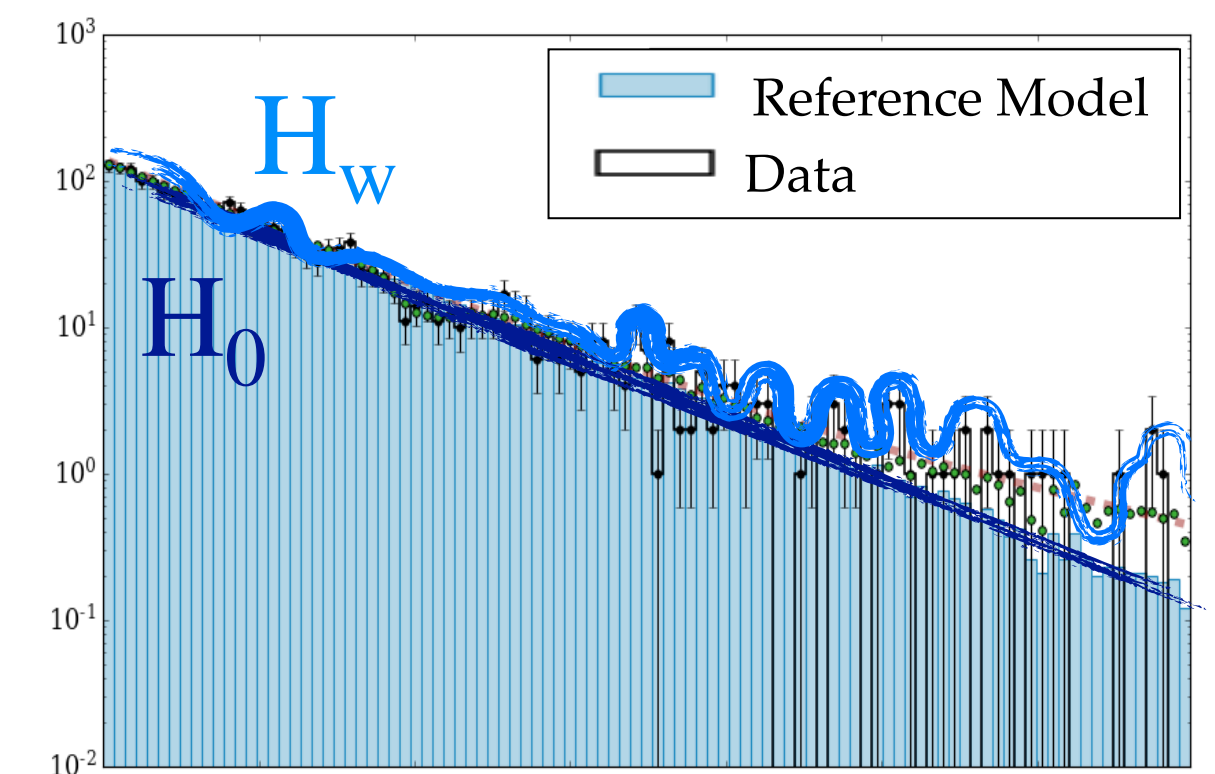
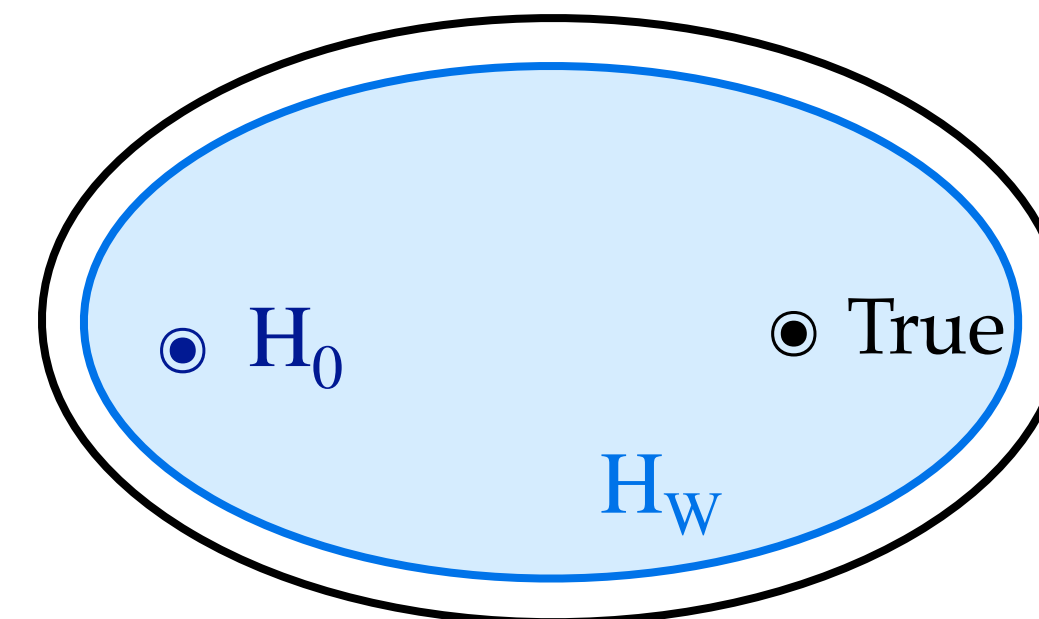
**Release** the assumptions on the alternative

$$t(\mathcal{D}) = \max_{\mathbf{w}} \left[ 2 \log \frac{\mathcal{L}(\mathcal{D} | H_{\mathbf{w}})}{\mathcal{L}(\mathcal{D} | H_0)} \right]$$

NPLM algorithm

Universal approximator  
(NN, kernel methods, ...)

$$n(x | H_{\mathbf{w}}) = e^{f(x; \mathbf{w})} n(x | R_0)$$



**Alternative  $\equiv$  Data**

(The *saturated model* proposed by Baker and Cousins — [Nucl.Instrum.Meth. 221 \(1984\)](#) — is a binned version of this approach)



# New Physics Learning Machine (NPLM)

## Maximum Likelihood from Minimal Loss

Test statistic

$$\bar{t}(\mathcal{D}) = 2 \max_{\mathbf{w}} \log \left[ \frac{\mathcal{L}(\mathcal{D} | H_{\mathbf{w}})}{\mathcal{L}(\mathcal{D} | R_0)} \right] = 2 \max_{\mathbf{w}} \left\{ \log \left[ \frac{e^{-N(\mathbf{w})}}{e^{-N(\mathcal{R})}} \prod_{i=1}^{N_D} \frac{n(x_i | \mathbf{w})}{n(x_i | \mathcal{R})} \right] \right\}$$

$$= -2 \min_{\mathbf{w}} \{ \bar{L}[f(\cdot; \mathbf{w})] \}$$

$\mathbf{w}$ : trainable parameters on the NN model  
 $\mathcal{D}$ : data sample  
 $\mathcal{R}$ : reference sample (built according to the  $R_0$  hypothesis); could be weighted ( $w$ )

Assumptions:

- $N_{\mathcal{R}} \gg N_{\mathcal{D}}$  the statistical fluctuations of the reference sample are negligible.
- the weights of the reference sample ( $w$ ) are such that the reference sample is normalised to match the data sample luminosity

Loss function

$$\bar{L}[f(x; \mathbf{w})] = - \sum_{x \in \mathcal{D}} f_{\mathbf{w}}(x) + \sum_{x \in \mathcal{R}} \frac{N(\mathcal{R})}{N_{\mathcal{R}}} (e^{f_{\mathbf{w}}(x)} - 1)$$

# New Physics Learning Machine (NPLM)

Likelihood ratio from [Binary Cross Entropy](#)

Test statistic

$$\bar{t}(\mathcal{D}) = 2 \max_{\mathbf{w}} \log \left[ \frac{\mathcal{L}(\mathcal{D} | H_{\mathbf{w}})}{\mathcal{L}(\mathcal{D} | R_0)} \right] = 2 \max_{\mathbf{w}} \left\{ \log \left[ \frac{e^{-N(\mathbf{w})}}{e^{-N(\mathbf{R})}} \prod_{i=1}^{N_D} \frac{n(x_i | \mathbf{w})}{n(x_i | \mathbf{R})} \right] \right\}$$

$$= 2 \sum_{x \in \mathcal{D}} f_{\mathbf{w}}(x) - 2 \sum_{x \in \mathcal{R}} \frac{N(\mathbf{R})}{N_{\mathcal{R}}} \left[ e^{f(x; \mathbf{w})} - 1 \right]$$

$$f(x; \hat{\mathbf{w}}) = \log \left[ \frac{n(x | H_{\hat{\mathbf{w}}})}{n(x | \mathbf{R})} \right]$$

$\mathbf{w}$ : trainable parameters on the NN model  
 $\mathcal{D}$ : data sample  
 $\mathbf{R}$ : reference sample (built according to the  $R_0$  hypothesis); could be weighted ( $w$ )

Assumptions:

- $N_{\mathcal{R}} \gg N_{\mathcal{D}}$  the statistical fluctuations of the reference sample are negligible.
- the weights of the reference sample ( $w$ ) are such that the reference sample is normalised to match the data sample luminosity

Loss function

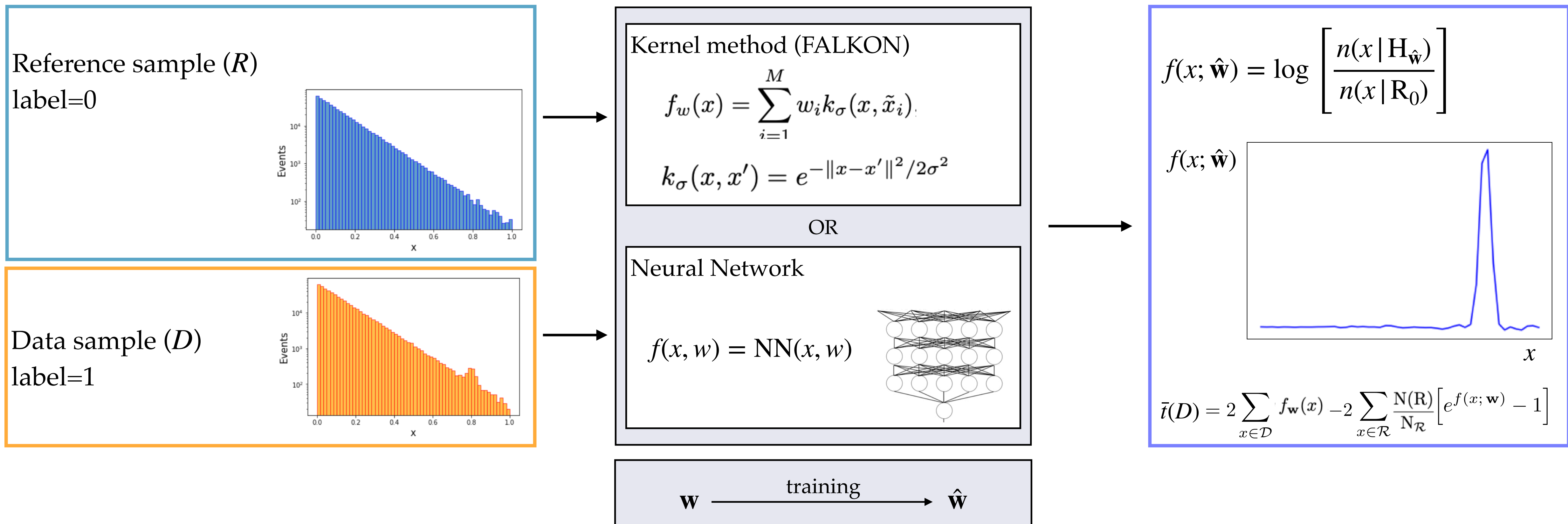
$$\bar{L} [f(x; \mathbf{w})] = - \sum_{x \in \mathcal{D}} \log \left[ 1 + e^{-f_{\mathbf{w}}(x)} \right] + \sum_{x \in \mathcal{R}} \frac{N(\mathbf{R})}{N_{\mathcal{R}}} \log \left[ 1 + e^{f_{\mathbf{w}}(x)} \right]$$

# New Physics Learning Machine (NPLM)

INPUT

MODEL  $f(\cdot, \mathbf{w})$

OUTPUT

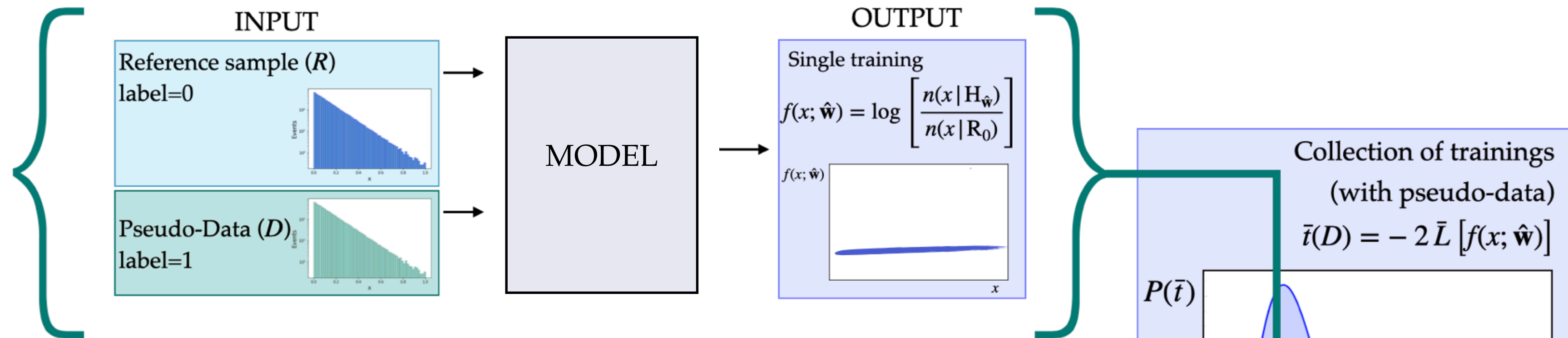


“Learning New Physics from a Machine” [Phys. Rev. D](#)

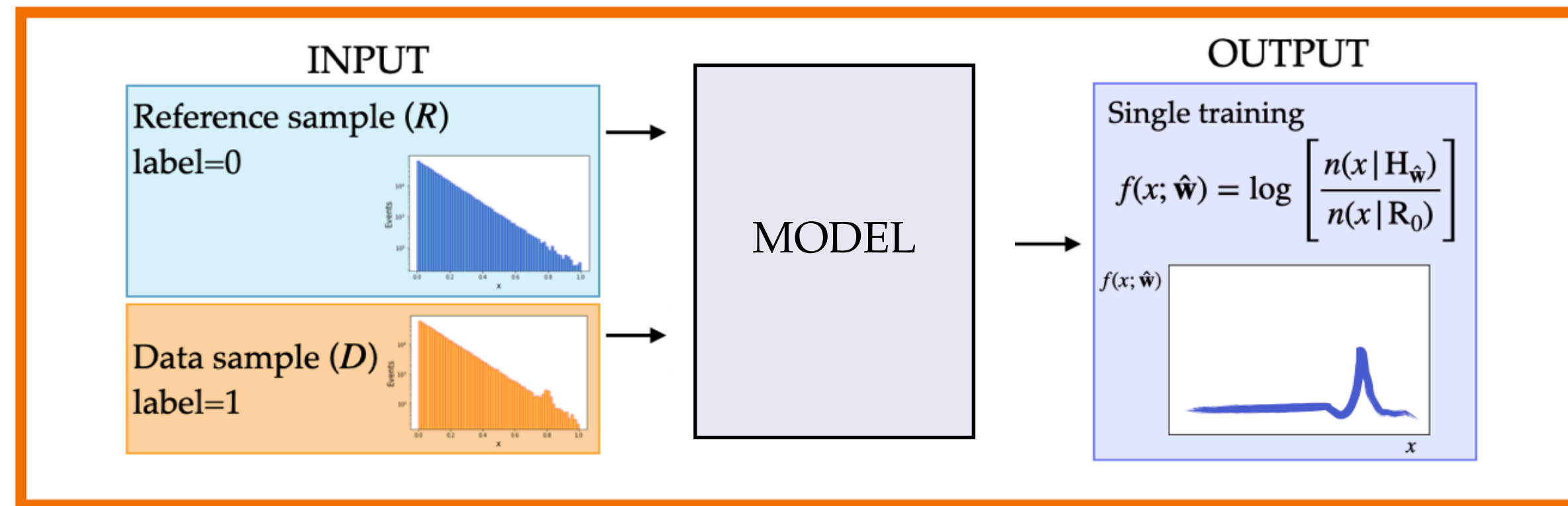
# New Physics Learning Machine (NPLM)

Frequentist  $p$ -value (aka calibration):

1. Run NPLM on toy experiments to simulated the response under the null hypothesis



2. Run NPLM the data of interest and check where the test outcome falls to compute an exclusion  $p$ -value

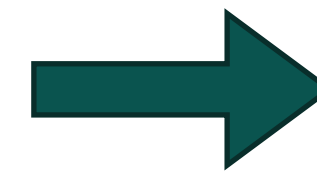


# New Physics Learning Machine (NPLM)

Controlling type I errors: NN model regularisation

## Weight clipping parameter:

Upper boundary to the magnitude that each trainable parameter can assume during the training.

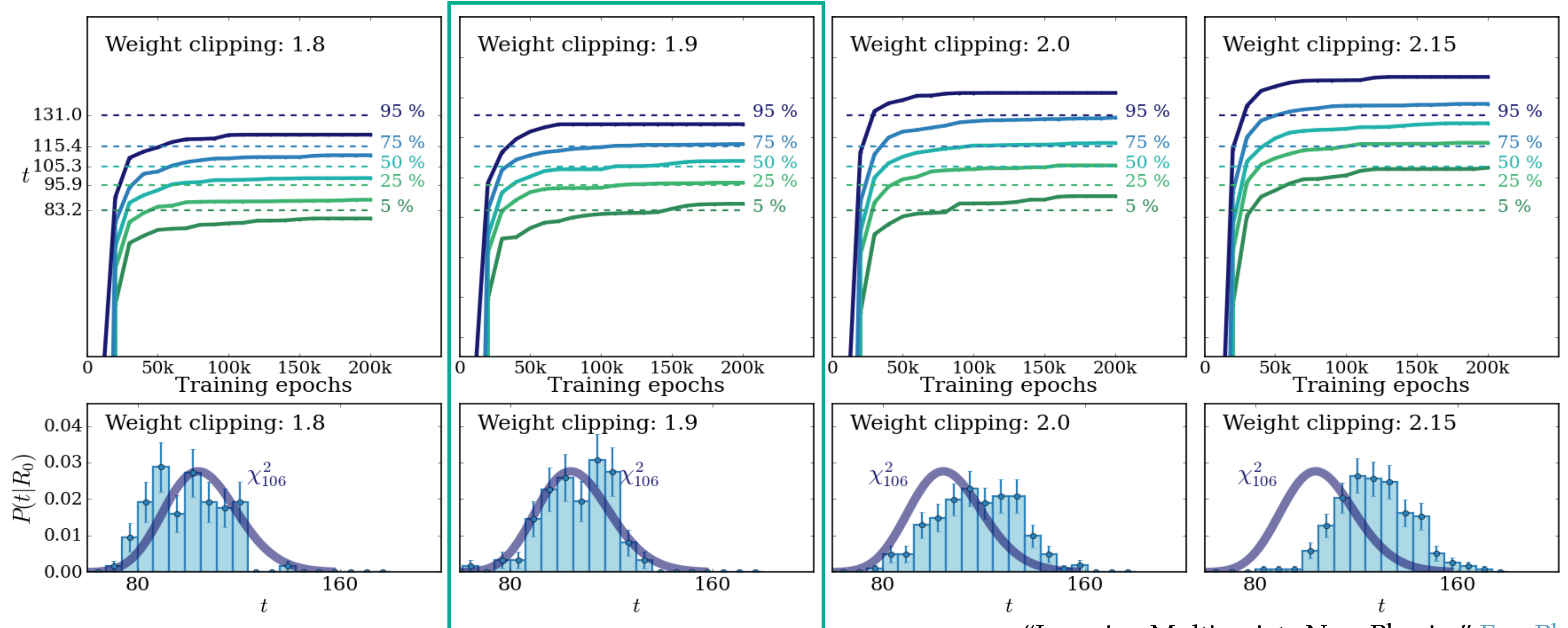


For a chosen NN architecture, tuning the weight clipping allows to recover a good agreement of the empirical distribution of  $t$  under  $R_0$  with a target  $\chi^2_{|w|}$  distribution.

Example:  
 NN model: 5-7-7-1,  
 Number of parameters: 106

### Legend:

- Percentiles of the empirical  $\bar{t}$  distribution under  $R_0$
- Percentiles of the target  $\chi^2_{|w|}$
- Empirical  $\bar{t}$  distribution under  $R_0$
- Target  $\chi^2_{|w|}$

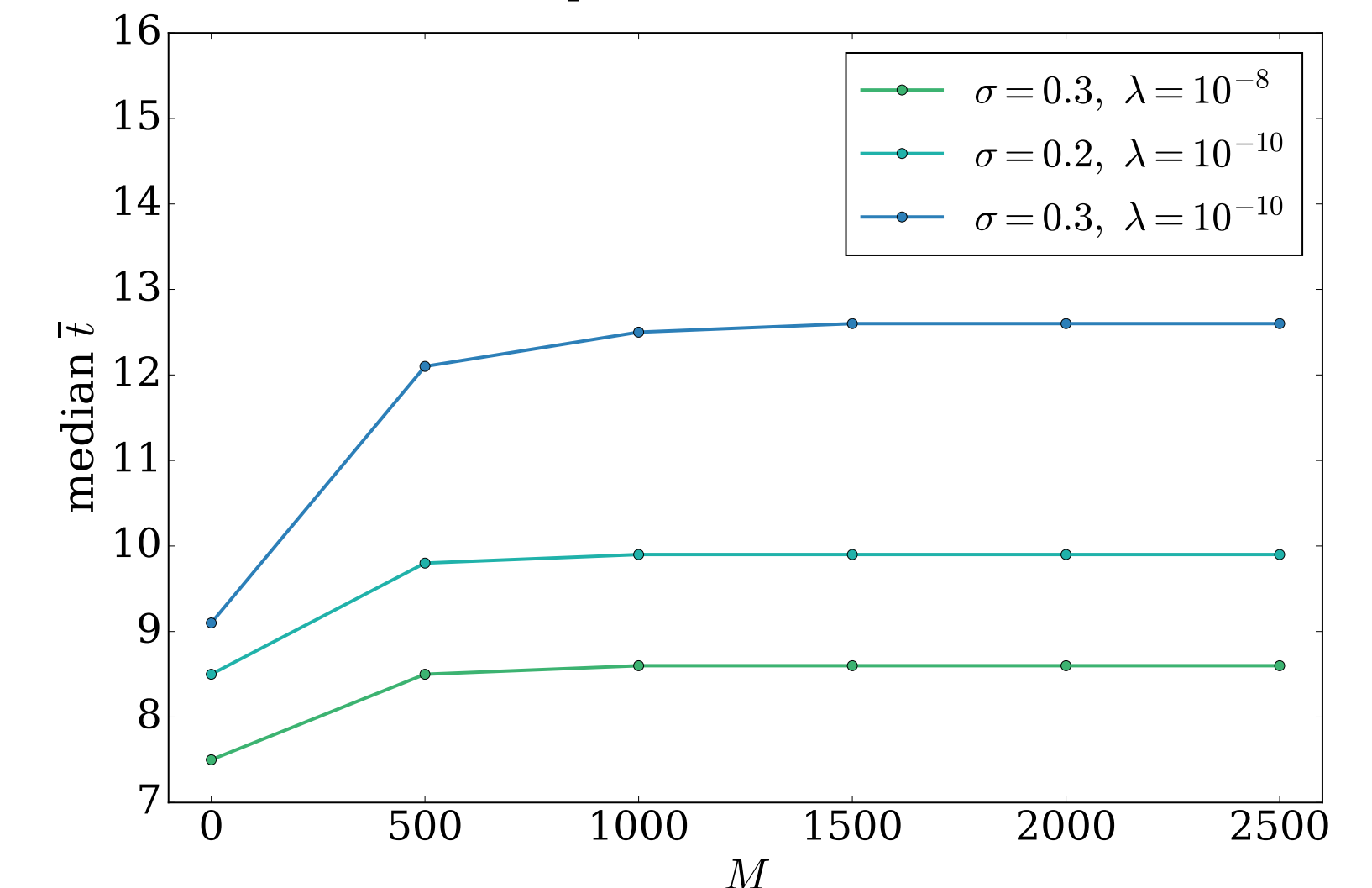
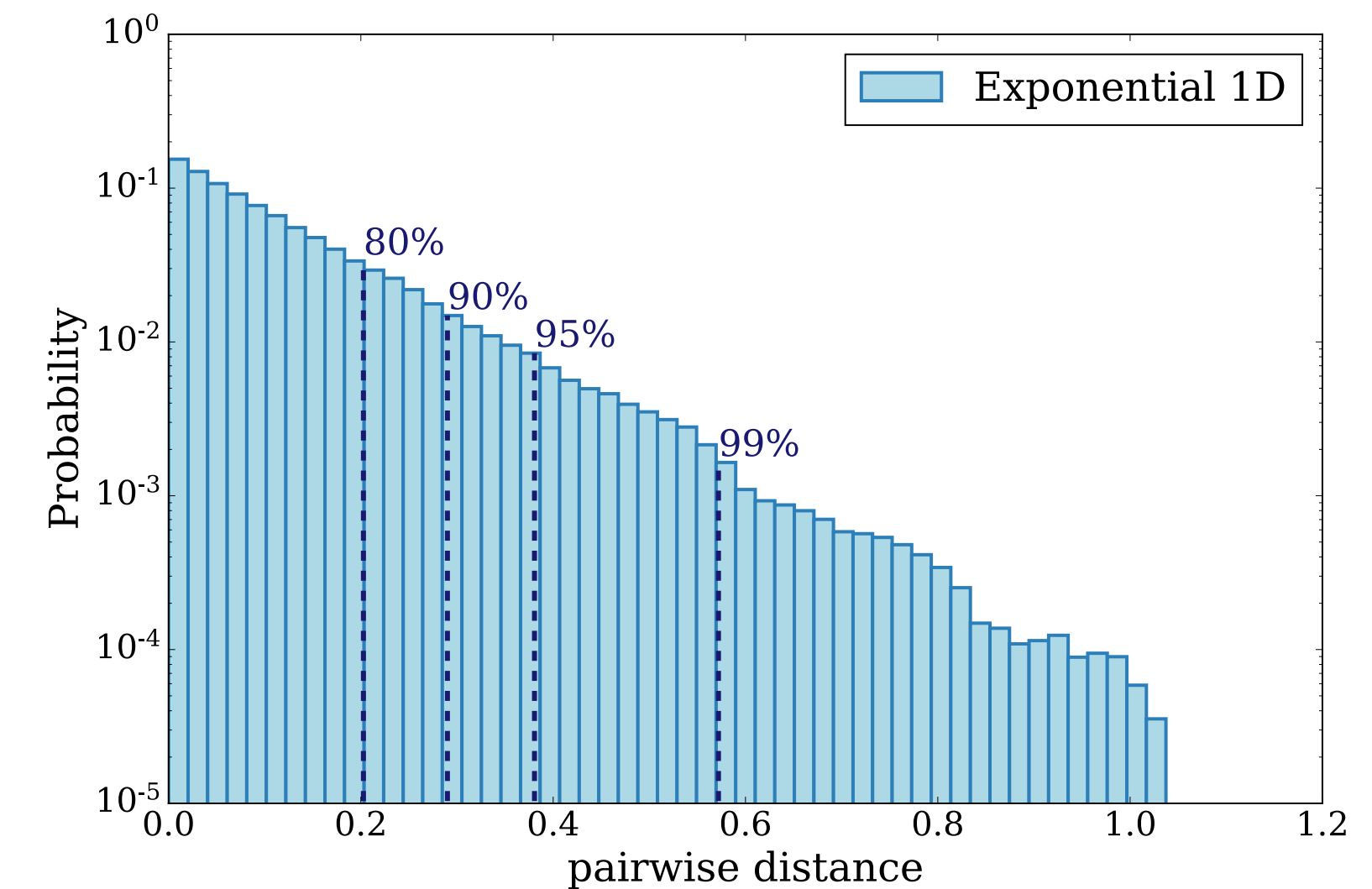


“Learning Multivariate New Physics” [Eur. Phys. J. C](#)

# New Physics Learning Machine (NPLM)

## Controlling type I errors: Kernel Method regularisation

- $\chi^2$  is found for each choice of  $(M, \sigma, \lambda)$ , provided that  $N_R \gg N_D$ .
- **Number of centres  $M$** : controls the expressive power of the model and therefore it should be as large as possible (at least as large as  $\sqrt{N}$  to achieve statistically optimal bounds of the training convergence).
- **Gaussian width  $\sigma$** : is selected as the 90th percentile of the pairwise distance between reference-distributed data points (after standardisation).
- **Regularisation parameter  $\lambda$** : is kept as small as possible while keeping training stable, i.e. avoiding large training times or non-numerical outputs.

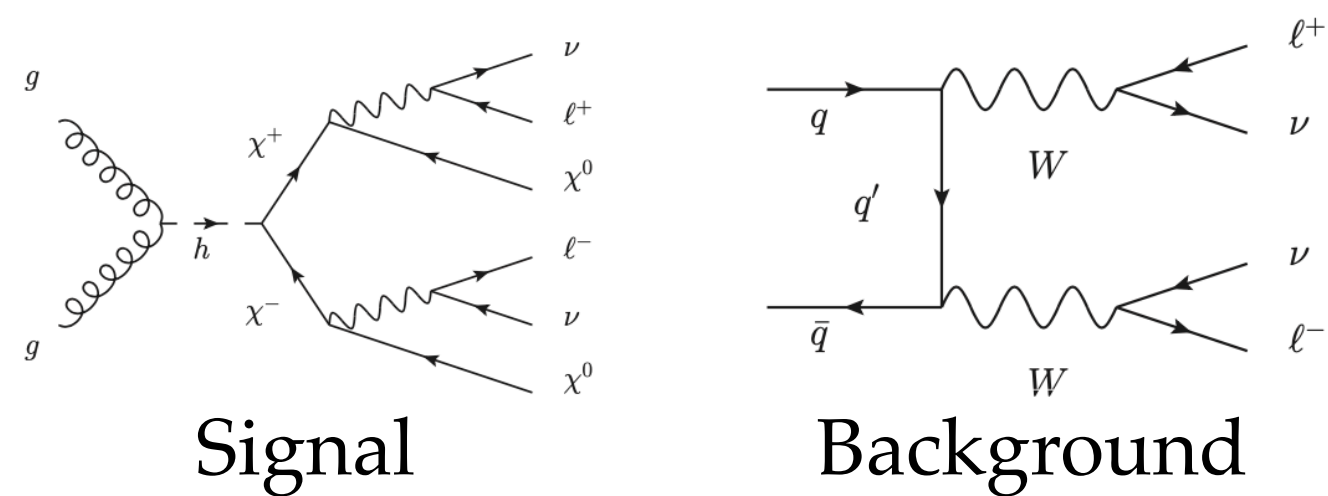


“Learning New Physics Efficiently with non-parametric models” [Eur. Phys. J. C](#)

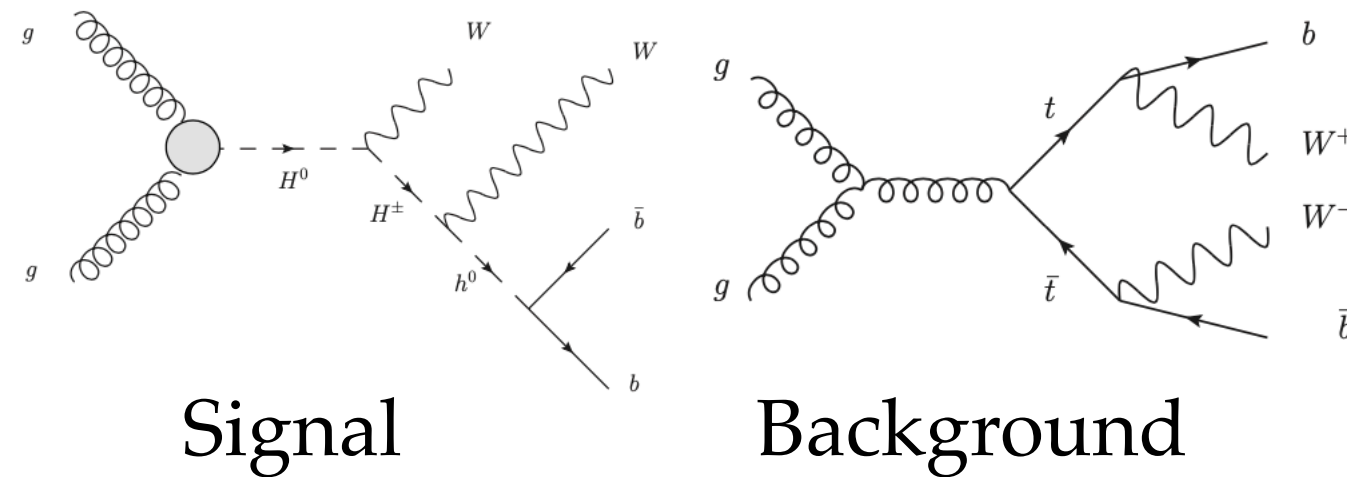
# NPLM: NN vs. Kernel methods

## Benchmarks:

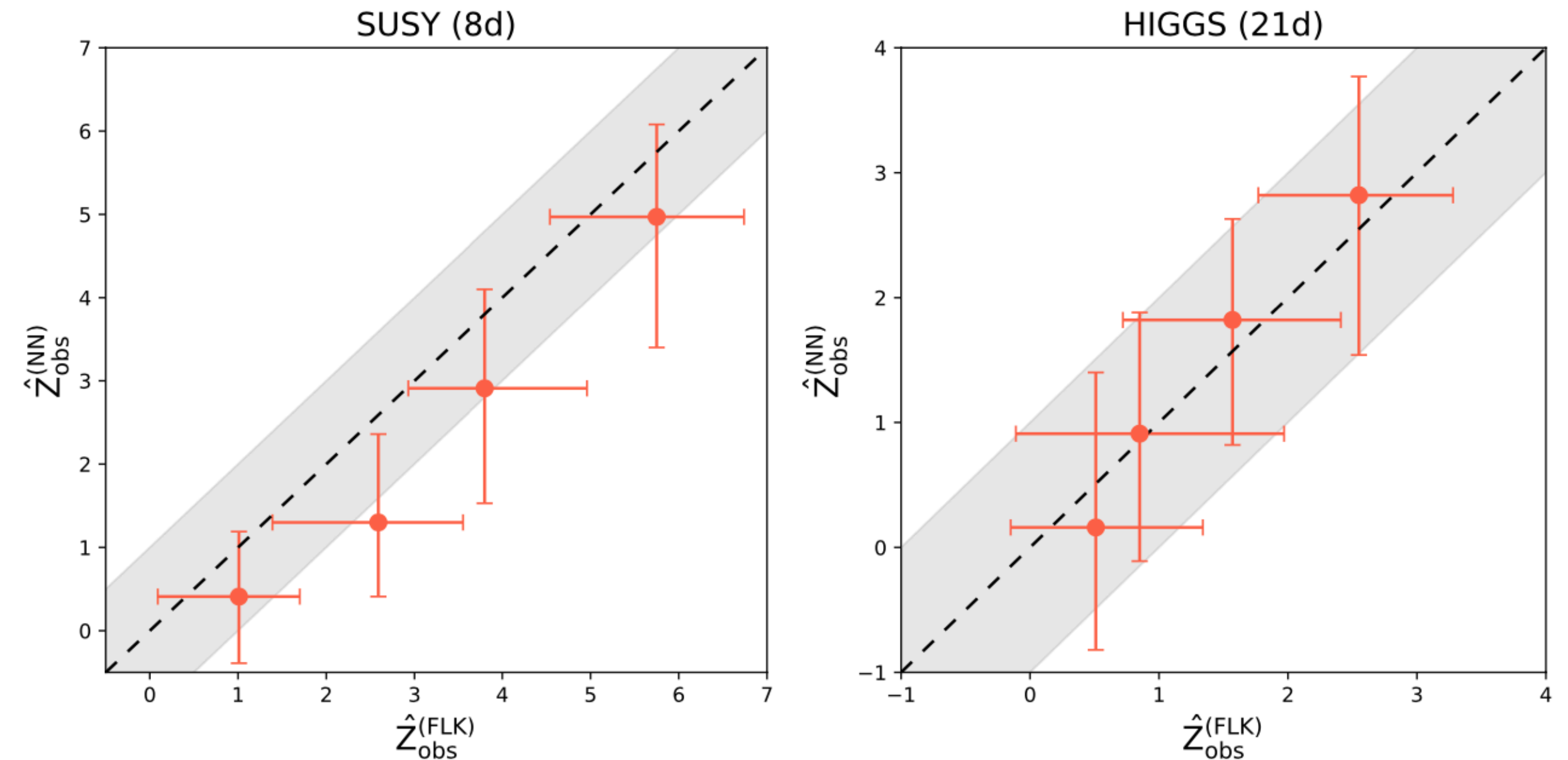
- 8D problem: (SUSY dataset\*)



- 21D problem: (HIGGS dataset\*\*)



## Falkon vs. NN: sensitivity performance



## Falkon vs. NN: execution time

Model	SUSY	HIGGS
<b>FLK</b>	<b>(18.2 ± 1.2) s</b>	<b>(22.7 ± 0.4) s</b>
<b>NN</b>	<b>(73.1 ± 10) h</b>	<b>(112 ± 9) h</b>

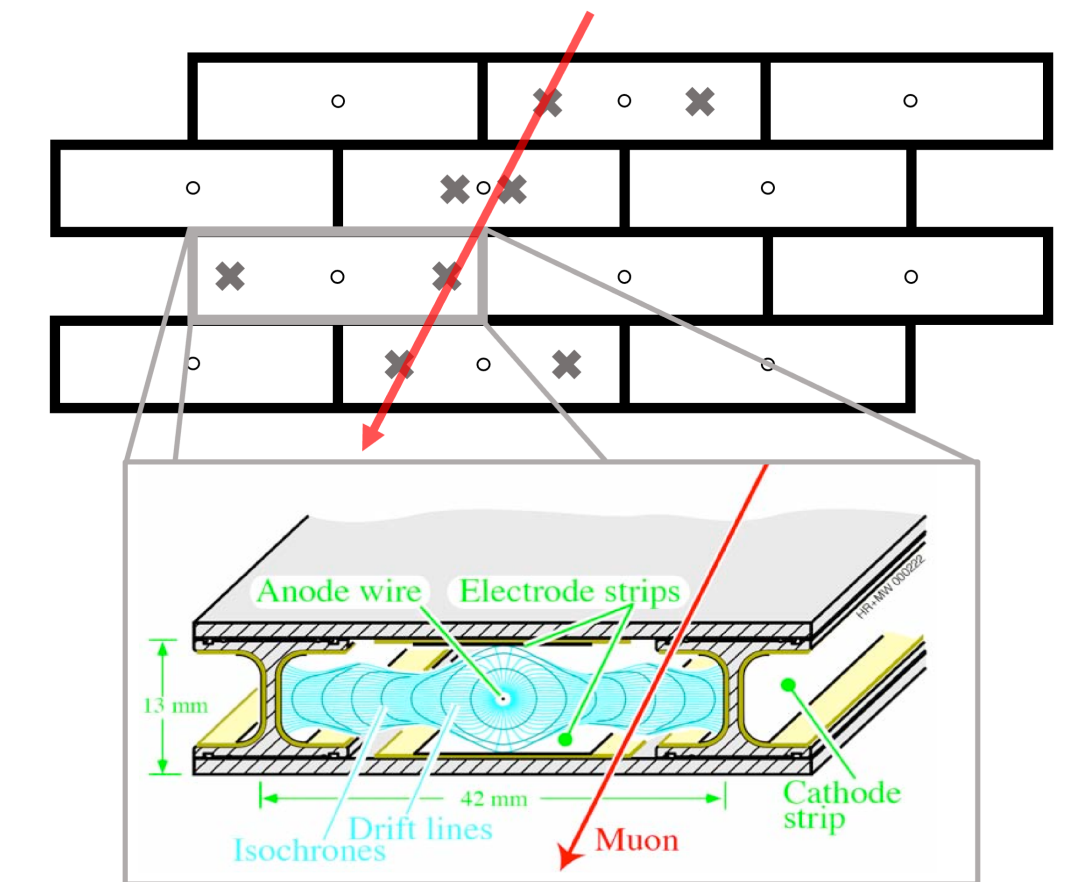
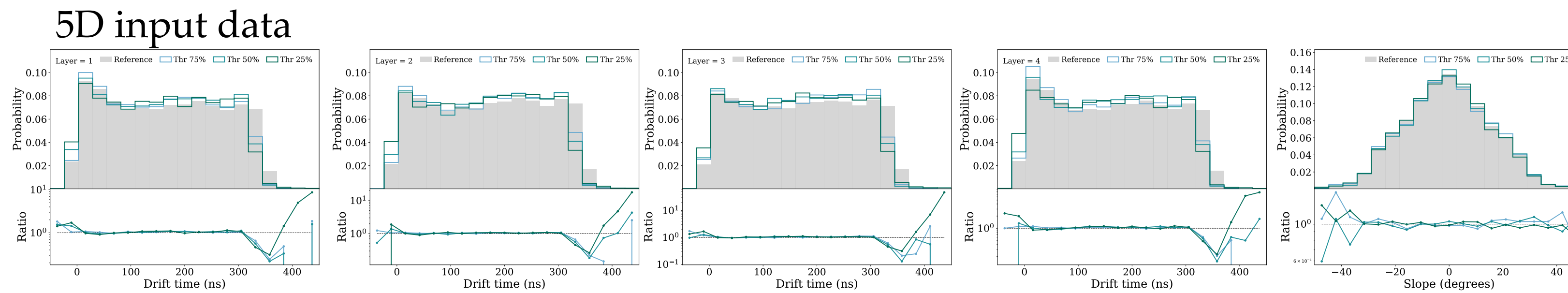
\* SUSY dataset: <https://archive.ics.uci.edu/ml/datasets/SUSY>

\*\* HIGGS dataset: <https://archive.ics.uci.edu/ml/datasets/HIGGS>

# NPLM with Kernel methods

## Efficient GOF computation on GPUs with Falkon

Example: Online Data Quality monitoring of a Drift Tube Chamber



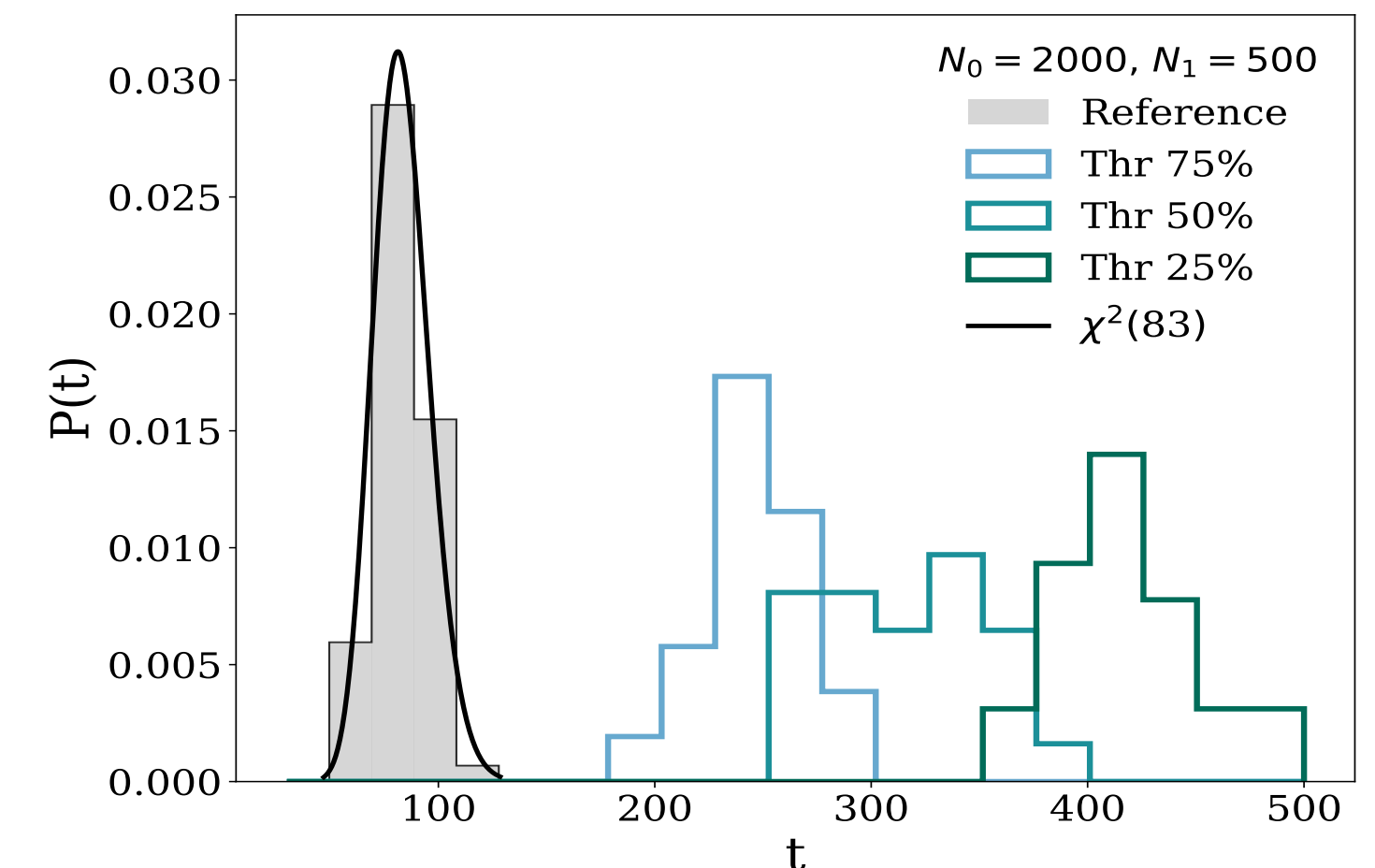
- **Reference sample:** long run in optimal conditions
- **Anomalous samples:** short runs acquired in presence of a controlled anomaly in the value of the **threshold tension** of the DT chamber

$$M = 2000, \sigma = 4.5, \lambda = 10^{-7}$$

$$N(D) = 500, N_{\text{ref}} = 2000$$

Execution time:  $\sim 0.5$  s

More about this in our recent preprint [arXiv:2303.05413](https://arxiv.org/abs/2303.05413)





# Neyman-Pearson testing with NPLM for GOF

Comparing the null hypothesis (Reference) with an alternative.

Release the assumptions on the alternative (**Alternative  $\equiv$  Data**)

$$t(\mathcal{D}) = \max_{\mathbf{w}} \left[ 2 \log \frac{\mathcal{L}(\mathcal{D} | H_{\mathbf{w}})}{\mathcal{L}(\mathcal{D} | H_0)} \right]$$

Peculiar aspects due to the nature of the test:

- Unbalanced problem (ideally  $N_D \ll N_R$ )
- Regularisation scheme to control type I errors (no need for train-test splitting)
- In-sample evaluation of the test

How do we compare to other GOF methods?

# NPLM vs. C2ST

## Classifier 2 Sample Test (C2ST, [Lopez et al. \(2017\)](#)):

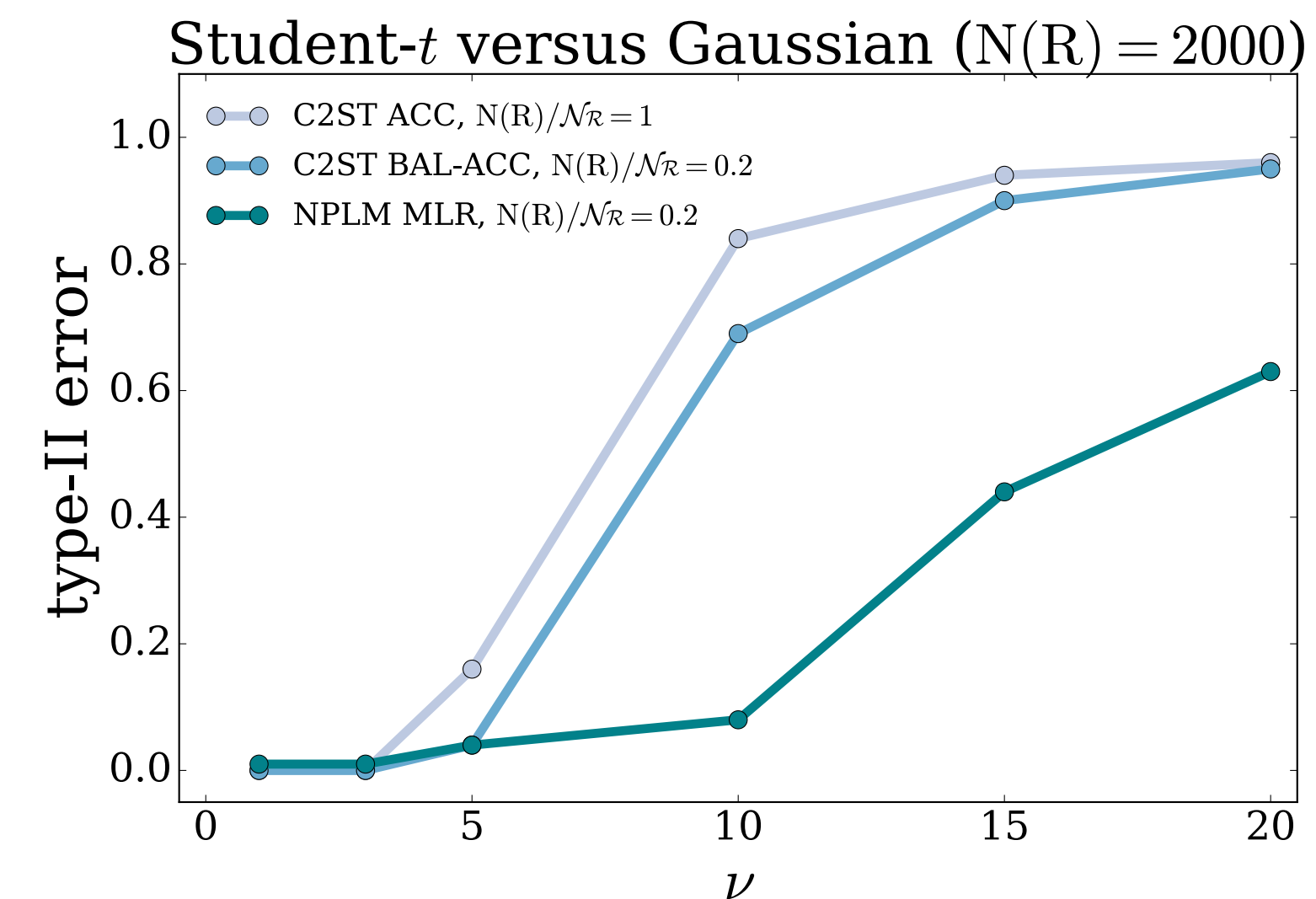
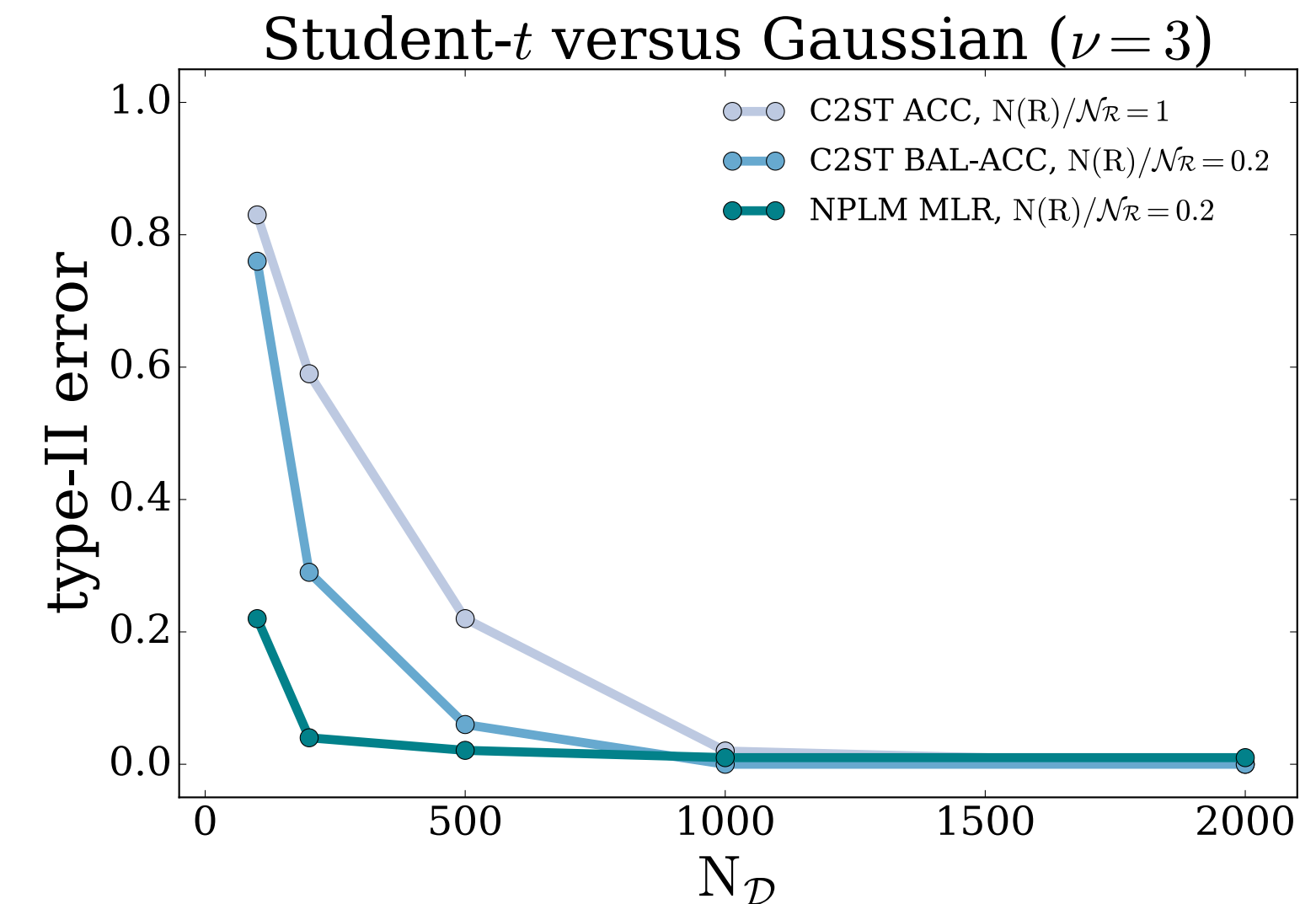
- Train-test split:  $R = R_{tr} \cup R_{te}$  and  $D = D_{tr} \cup D_{te}$
- Metric:
  - The original C2ST uses the accuracy (ACC)
  - We replace it with a modified version of the balanced accuracy (BAL-ACC) which is sensitive to normalization effects and accounts for unbalanced samples

$$t'_{\text{BACC}} = \frac{2}{N(R) + N_D} \left[ \frac{N(R)}{N_{\mathcal{R}}} \sum_{x \in \mathcal{R}_{te}} \mathbb{I}[c_{\hat{\mathbf{w}}}(x) < 1/2] + \sum_{x \in \mathcal{D}_{te}} \mathbb{I}[c_{\hat{\mathbf{w}}}(x) > 1/2] \right]$$

- Out of sample test statistic
- $p$ -value from empirical test statistic distribution under the null

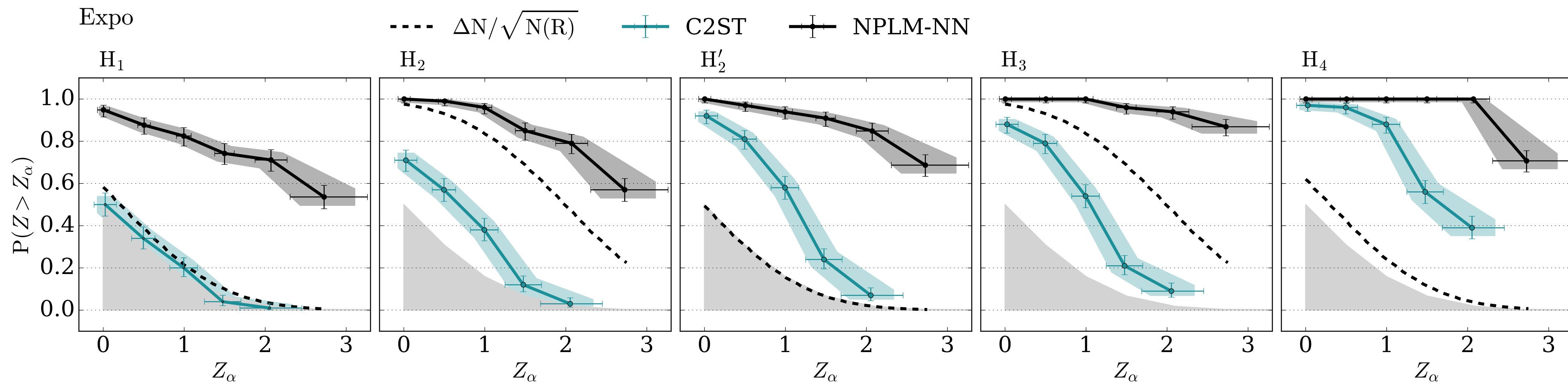
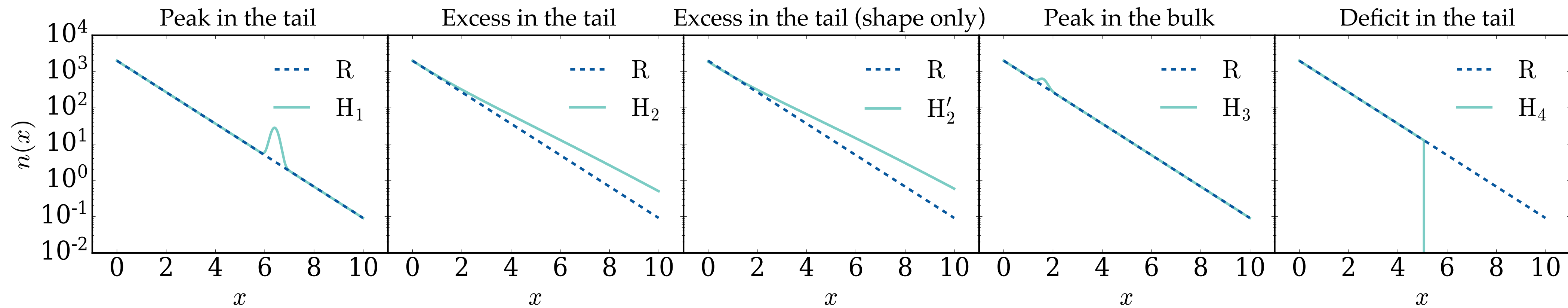
### Comments:

- Unbalanced problem ( $N_R > N_D$ ): performance improvement
- additional improvement NPLM wrt C2ST, why?
  - train-test splitting vs. in-sample ?
  - ACC vs. LRT ?



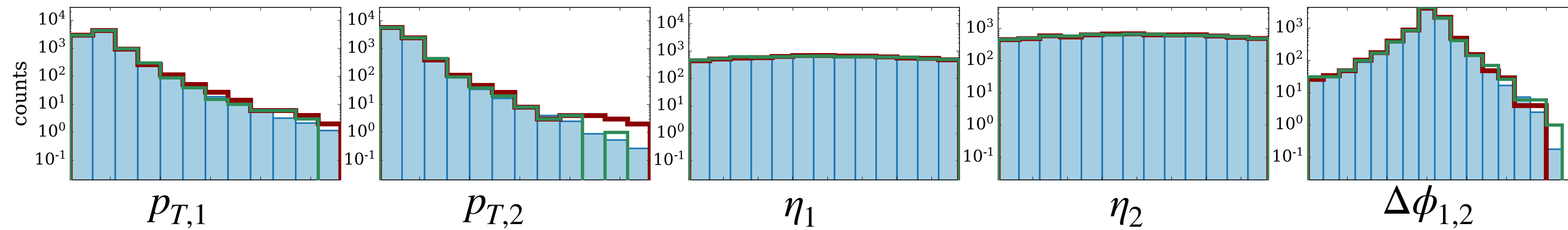
# NPLM vs. C2ST

1D toy model: Reference model  $n(x|R) = N(R)e^{-x}$



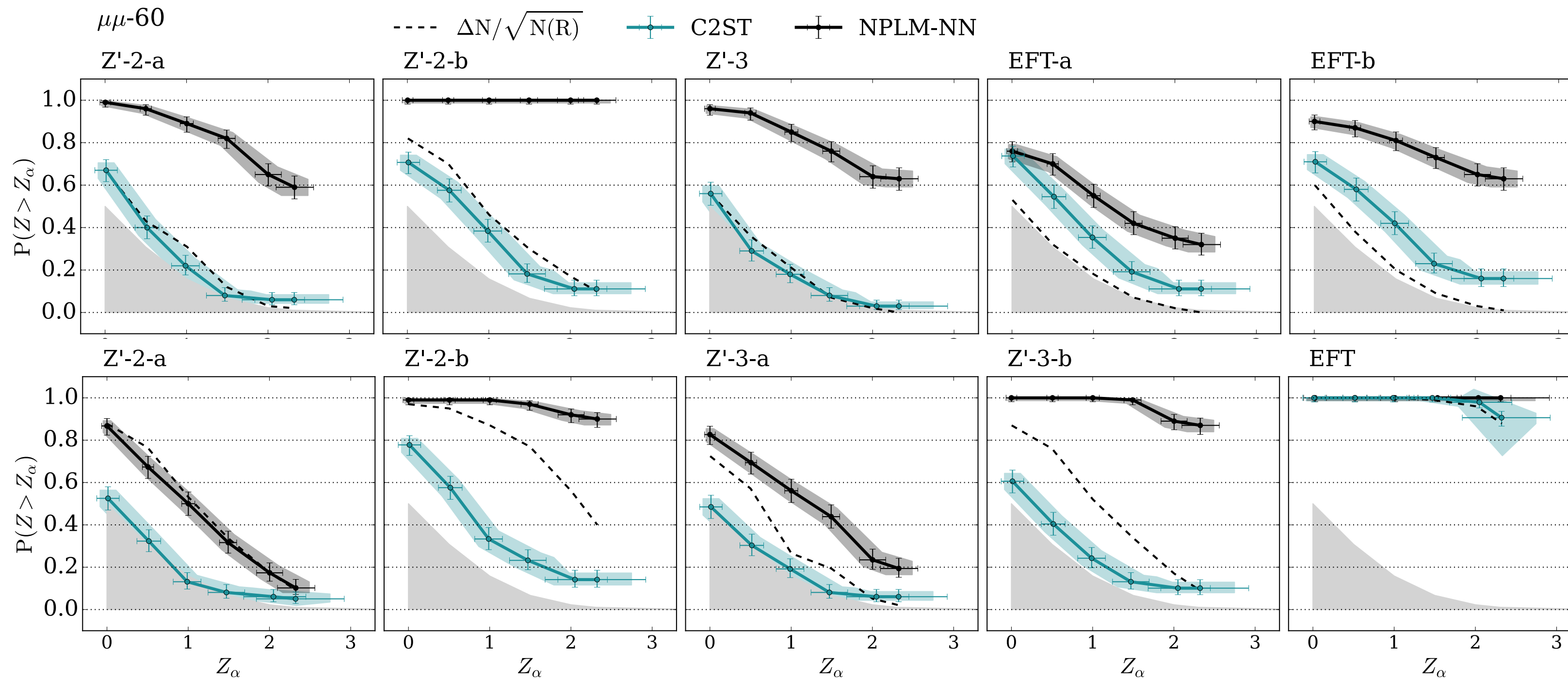
# NPLM vs. C2ST

5D  $\mu\mu$  final state:



█ REFERENCE  
█ Z' scenario  
█ EFT scenario

NOTE:  
 $M_{12}$  is **not** given as an input to the algorithm!



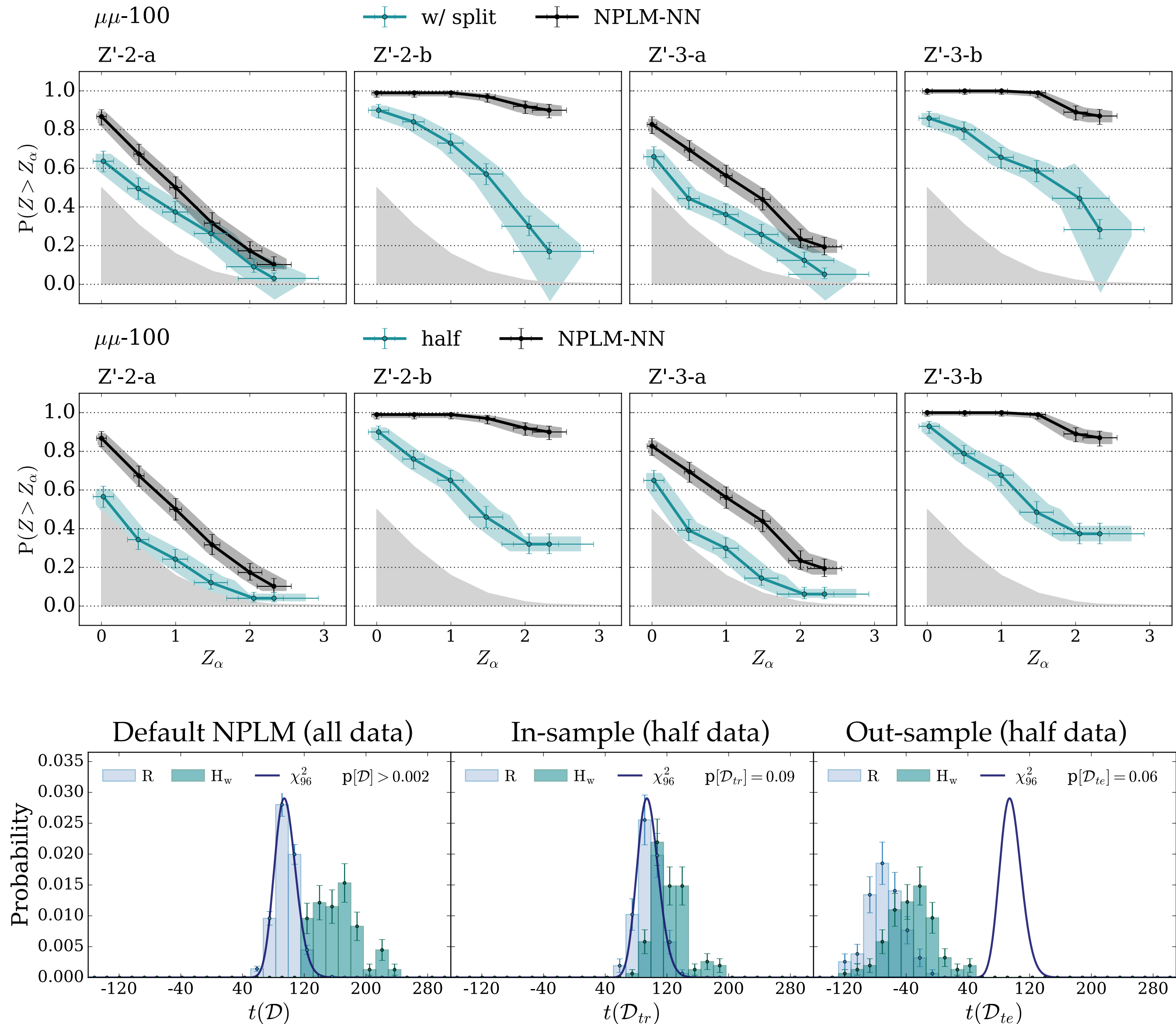
$M_{12} > 60 \text{ GeV } (\mu\mu\text{-60}).$

$M_{12} > 100 \text{ GeV } (\mu\mu\text{-100}).$

# Train-test splitting

- With the NPLM method, out of sample evaluation is almost equivalent to using half of the data sample.
- NPLM is regularised to “overfit” the specific fluctuations of the training data. By construction, the learnt model does not generalise well to a new data sample.
- $t(D_{te})$  is negative: the alternative learnt from  $D_{tr}$  is worst than the Reference

Train-test splitting is disadvantageous for NPLM



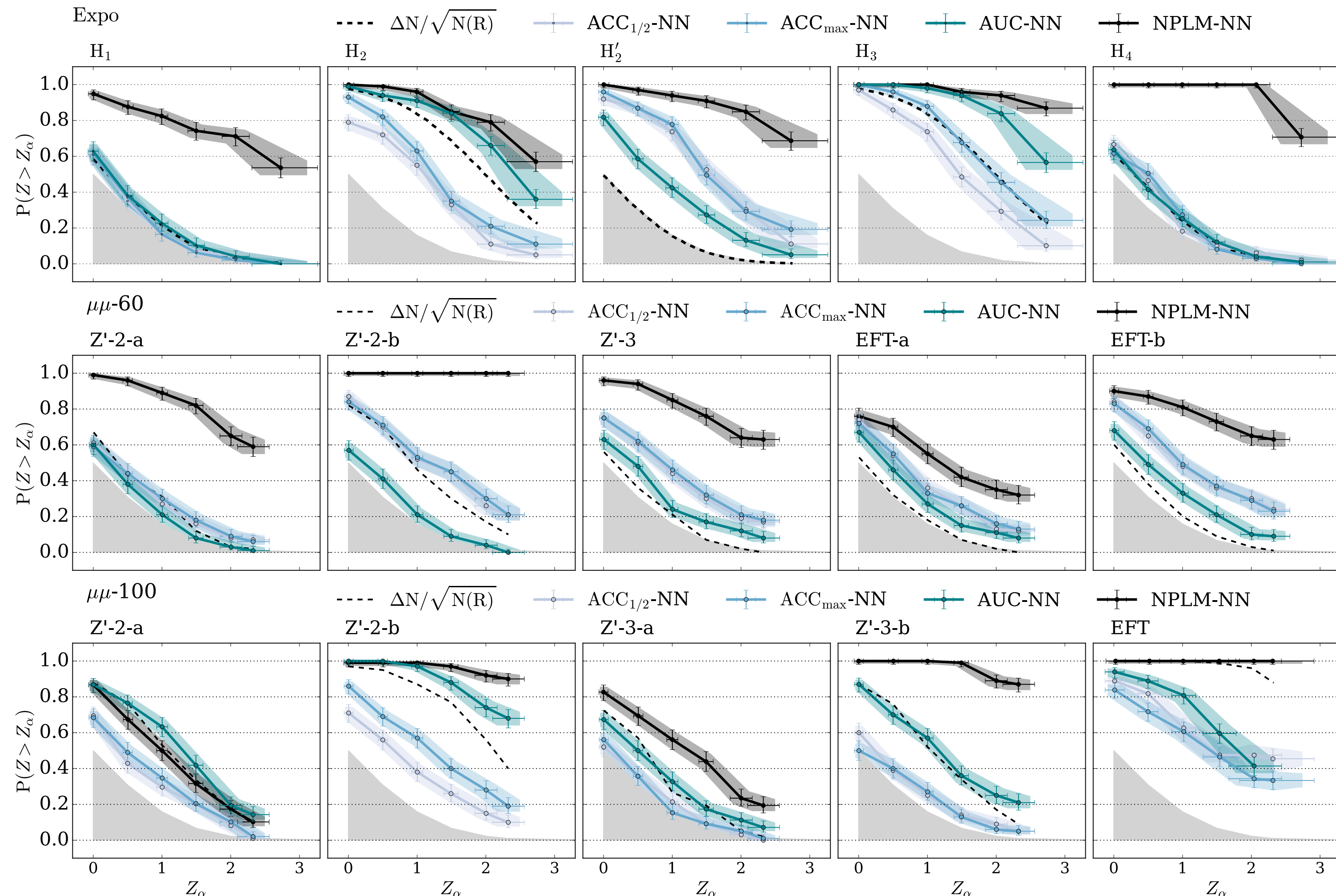
# Classifier-inspired variants of the test statistics in NPLM

Replace LRT with other 1D GOF tests

- Classifier-based tests

- Modified\* Balanced Accuracy (ACC)
- Modified\* Area Under the ROC Curve (AUC)

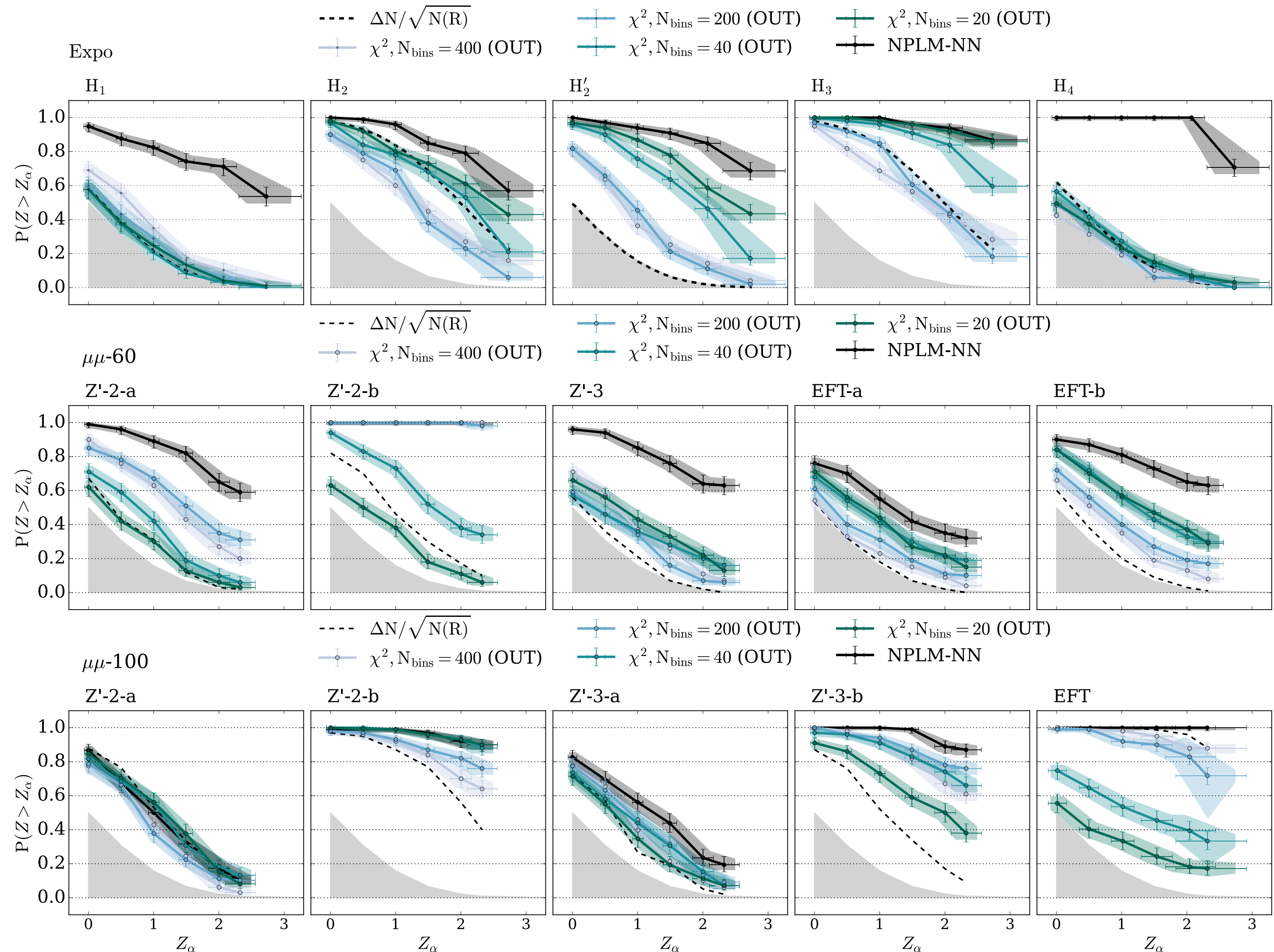
\* sensitive to n normalisation effects



# Classifier-inspired variants of the test statistics in NPLM

Replace LRT with other 1D GOF tests

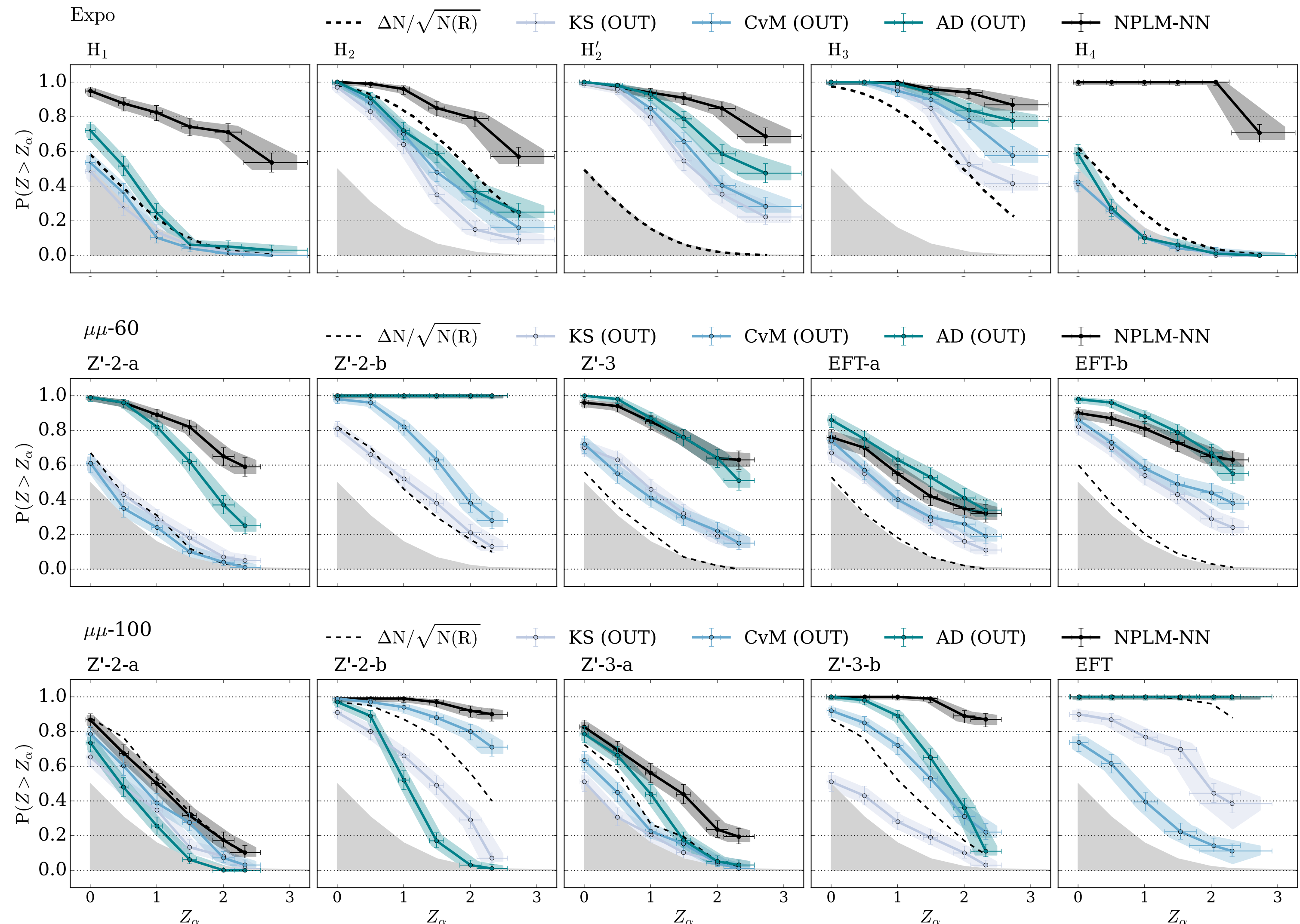
- Classifier-based tests:
  - Modified Balanced Accuracy (ACC)
  - Modified Area Under the ROC Curve (AUC)
  
- Standard 1D GoF tests:
  - $\chi^2$  tests
  - EDF tests
    - Kolmogorov-Smirnov (KS)
    - Cramer-von Mises (CvM)
    - Anderson-Darling (AD)
  - Spacing statistics:
    - Moran (M)
    - Recursive Product Spacing (RPS)



# Classifier-inspired variants of the test statistics in NPLM

Replace LRT with other 1D GOF tests

- Classifier-based tests:
  - Modified Balanced Accuracy (ACC)
  - Modified Area Under the ROC Curve (AUC)
  
- Standard 1D GoF tests:
  - $\chi^2$  tests
  - **EDF tests**
    - Kolmogorov-Smirnov (KS)
    - Cramer-von Mises (CvM)
    - Anderson-Darling (AD)
  - Spacing statistics:
    - Moran (M)
    - Recursive Product Spacing (RPS)

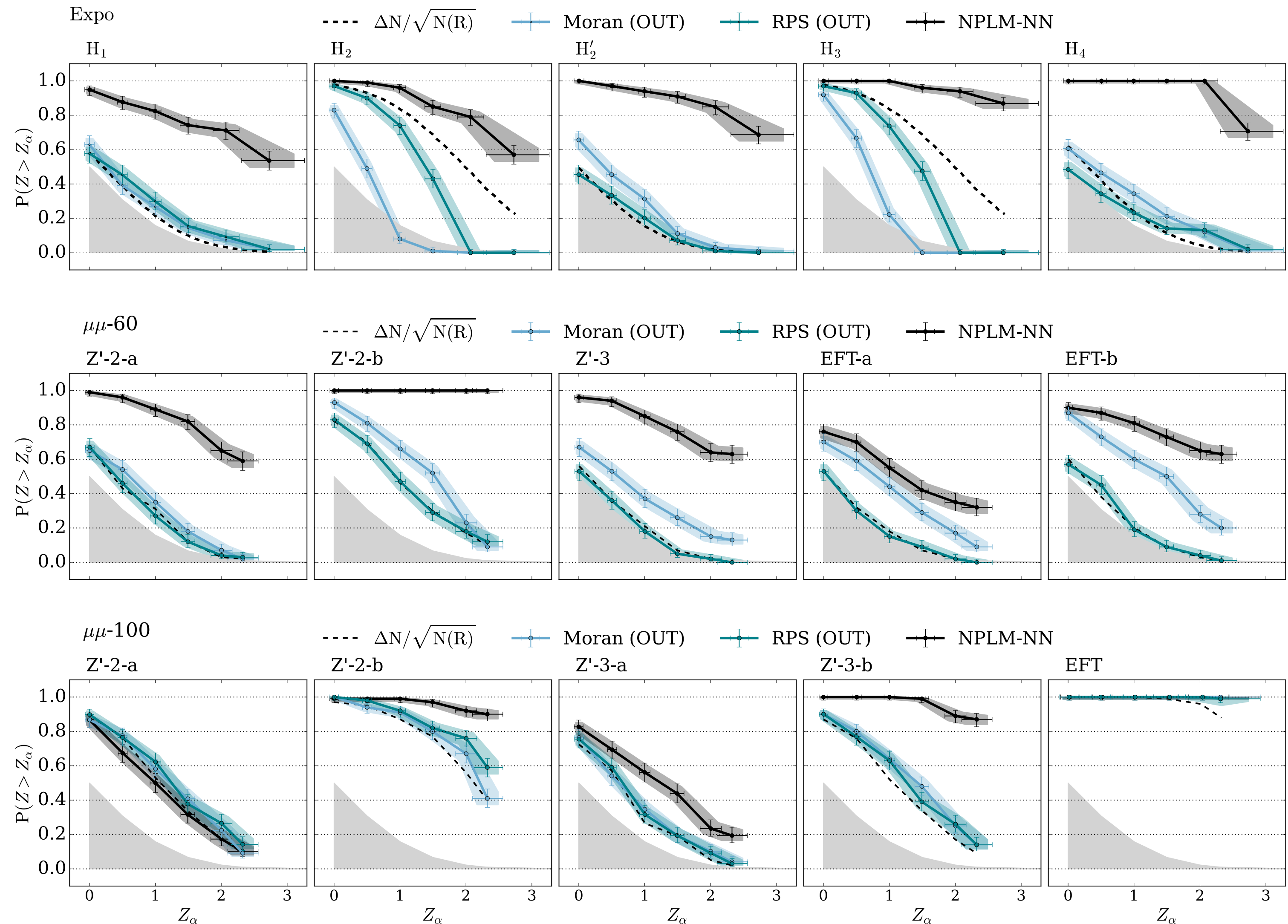




# Classifier-inspired variants of the test statistics in NPLM

Replace LRT with other 1D GOF tests

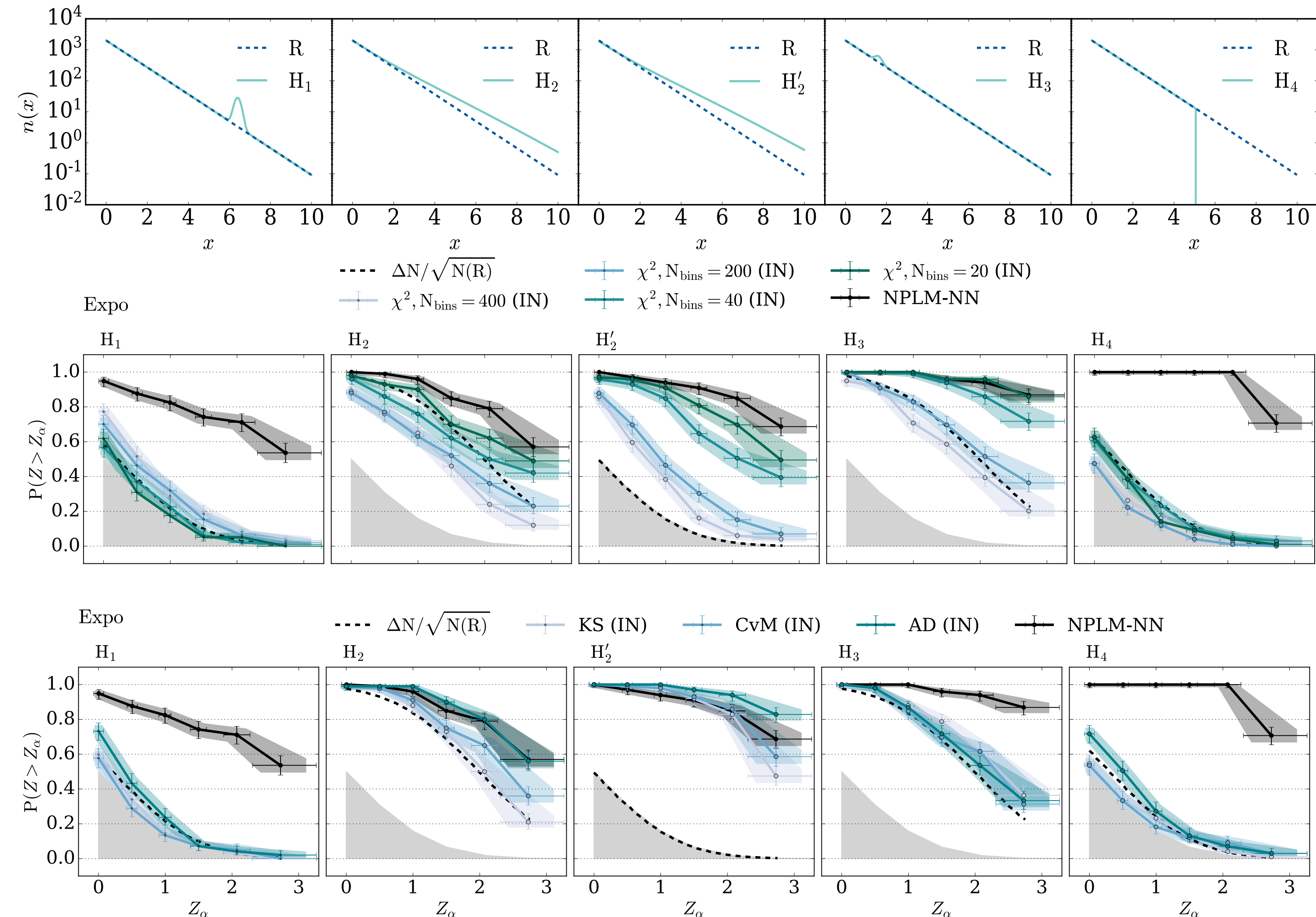
- Classifier-based tests
  - Modified Balanced Accuracy (ACC)
  - Modified Area Under the ROC Curve (AUC)
- Standard 1D GoF tests:
  - $\chi^2$  tests
  - EDF tests
    - Kolmogorov-Smirnov (KS)
    - Cramer-von Mises (CvM)
    - Anderson-Darling (AD)
  - Spacing statistics
    - Moran (M)
    - Recursive Product Spacing (RPS)



# NPLM vs. standard 1D GOF tests

Comparison with standard 1D GoF tests

- $\chi^2$  tests
- EDF tests
  - Kolmogorov-Smirnov (KS)
  - Cramer-von Mises (CvM)
  - Anderson-Darling (AD)



The likelihood-ratio test has, on average, the highest power against the list of considered anomalies

# Summary

NPLM as a general tool for GOF, beyond New Physics searches.

- **Signal-agnostic** test
- **Global p-value**: one test to detect them all
- **Multivariate**: allow to look at the data in a more inclusive way, no need to compress the information in one variable and hence to make assumptions on which observables are relevant
- **Systematic uncertainties** for “imperfect” Reference models [Eur. Phys. J. C 82, 275 \(2022\)](#) (crucial for LHC analysis)
- **Fast execution**: kernel methods on GPUs

# Summary

NPLM as a general tool for GOF, beyond New Physics searches.

- **Comparisons with a classifier-based approach (C2ST)** show on average better performances of NPLM.
  - Strong advantages in NPLM come from:
    - in-sample evaluation of the test
    - LRT as the test statistic
  - Comparisons to other ML-based approaches to be done
- **Comparisons with standard 1D GOF** show on average better performances of NPLM.
- Choice of **benchmarks**:
  - defines the figure of merit for the comparison.
  - never going to be exhaustive, but we should aim at the best possible representation of the landscape of possible anomalies
  - Our selected benchmarks are biased by the initial problem of BSM searches and should be extended (DQM use case is a first step in that direction)