

NPLM: Learning New Physics aware of systematic uncertainties

R.T. d'Agnolo¹, G. Grosso^{2,3}, M. Pierini³, A. Wulzer², M. Zanetti²

¹ Université Paris-Saclay and CEA, ² Università degli Studi di Padova and INFN, ³ CERN experimental department

In this talk

- How to include systematic uncertainties in NPLM
- NPLM application to HEP:
5D analysis of a di-body final state at the LHC

“Learning New Physics from an Imperfect Machine” [Eur. Phys. J. C](#)

More about NPLM:

- “Learning New Physics from a Machine” [Phys. Rev. D](#)
- “Learning Multivariate New Physics” [Eur. Phys. J. C](#)
- Previous IML talk [27/04/21](#)

NPLM algorithm

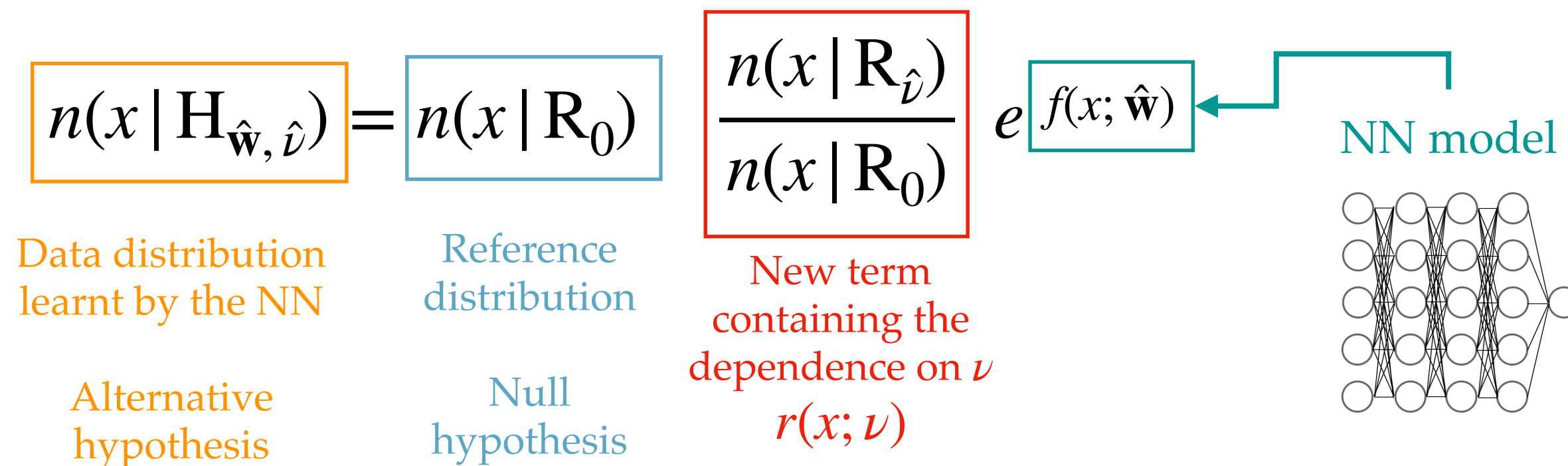
New Physics Learning Machine (NPLM)

Including systematic uncertainties

- Goal: performing a **log-likelihood-ratio hypothesis test**
(End-to-end strategy, from the data to a p -value for the discovery)

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

- Exploiting a Neural Network (NN) to **parametrize** the data distribution in terms of a Reference distribution (R_0)



R_{ν} : reference (null) hypothesis

$H_{\mathbf{w}, \nu}$: alternative hypothesis

\mathbf{w} : trainable parameters on the NN model

ν : set of nuisance parameters modelling the uncertainties effects

\mathcal{D} : data sample

\mathcal{A} : auxiliary sample (to constrain ν)

- **Signal-model-independent**: reduced assumptions on the signal hypothesis

“Learning New Physics from an Imperfect Machine” [Eur. Phys. J. C](#)

New Physics Learning Machine (NPLM)

Including systematic uncertainties

Maximum Likelihood from minimal loss:

Test statistic

$$t(\mathcal{D}, \mathcal{A}) = 2 \log \left[\frac{\max_{\mathbf{w}, \nu} \mathcal{L}(\mathbf{H}_{\mathbf{w}, \nu} | \mathcal{D}, \mathcal{A})}{\max_{\nu} \mathcal{L}(\mathbf{R}_{\nu} | \mathcal{D}, \mathcal{A})} \right] = 2 \log \left[\frac{\max_{\mathbf{w}, \nu} \mathcal{L}(\mathbf{H}_{\mathbf{w}, \nu} | \mathcal{D}) \mathcal{L}(\nu | \mathcal{A})}{\max_{\nu} \mathcal{L}(\mathbf{R}_{\nu} | \mathcal{D}) \mathcal{L}(\nu | \mathcal{A})} \right]$$

$$= \tau(\mathcal{D}, \mathcal{A}) - \Delta(\mathcal{D}, \mathcal{A})$$

\mathbf{w} : trainable parameters on the NN model
 ν : set of nuisance parameters modelling the uncertainties effects
 \mathcal{D} : data sample
 \mathcal{A} : auxiliary sample (used to constrain ν)

Tau term:

$$\tau(\mathcal{D}, \mathcal{A}) = 2 \max_{\mathbf{w}, \nu} \log \left[\frac{\mathcal{L}(\mathbf{H}_{\mathbf{w}, \nu} | \mathcal{D}) \mathcal{L}(\nu | \mathcal{A})}{\mathcal{L}(\mathbf{R}_0 | \mathcal{D}) \mathcal{L}(\mathbf{0} | \mathcal{A})} \right] = -2 \min_{\mathbf{w}, \nu} L \left[f(x, \mathbf{w}), \nu; \hat{\delta}(x) \right]$$

Contains the dependence on a NN model (universal approximator)

Built on the knowledge of the Reference model (purely SM term)

Delta term:

$$\Delta(\mathcal{D}, \mathcal{A}) = 2 \max_{\nu} \log \left[\frac{\mathcal{L}(\mathbf{R}_{\nu} | \mathcal{D}) \mathcal{L}(\nu | \mathcal{A})}{\mathcal{L}(\mathbf{R}_0 | \mathcal{D}) \mathcal{L}(\mathbf{0} | \mathcal{A})} \right] = -2 \min_{\nu} L \left[\nu; \hat{\delta}(x) \right]$$

$$r(x; \nu) = \frac{n(x | \mathbf{R}_{\nu})}{n(x | \mathbf{R}_0)}$$

Taylor's expansion learning:

$$\hat{r}(x; \nu) = \exp \left[\hat{\delta}_1(x) \nu + \hat{\delta}_2(x) \nu^2 + \dots \right]$$

NN 1 NN2 ...

"Learning New Physics from an Imperfect Machine" [Eur. Phys. J. C](#)

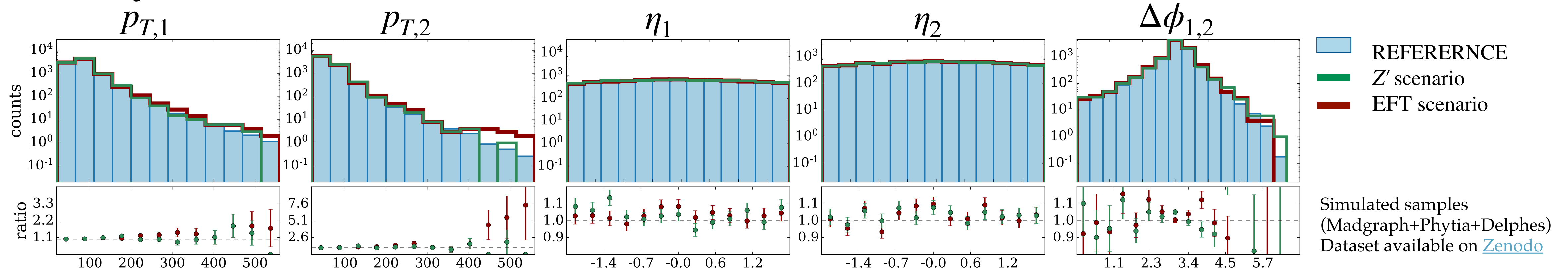
Application:

5D analysis of a di-body final state at the LHC

Di-body final state at the LHC

Dataset

5D analysis — Input variables:



Uncertainties on the reference sample (SM):

- Global normalization effect: $\sigma_N = 2.5\%$
- Momentum scale effect:

$$p_{T1,2}^{(b,e)} = \exp \left[\nu_s \sigma_s^{(b,e)} / \sigma_s^{(b)} \right] p_{T1,2}^{(b,e)} \quad \text{(b) barrel region } |\eta| < 1.2, \quad \text{(e) endcaps region } |\eta| \geq 1.2$$

- Muon-like regime: $\sigma_S^{(b)} = 0.05\%$, $\sigma_S^{(e)} = 0.15\%$
- Electron-like regime: $\sigma_S^{(b)} = 0.3\%$, $\sigma_S^{(e)} = 0.9\%$
- Tau-like regime: $\sigma_S^{(b)} = \sigma_S^{(e)} = 3\%$

Di-body final state at the LHC

Dataset

New Physics benchmarks:

Resonance in the two-body invariant mass

- **Z' scenario:** new vector boson with the same SM coupling as the Z boson and mass of 300 GeV.

- Muon-like, electron-like regimes:
 $M_{12} > 100 \text{ GeV}, L = 0.35 \text{ fb}^{-1}, N(S) = 120$
- Tau-like regime:
 $M_{12} > 120 \text{ GeV}, L = 1.1 \text{ fb}^{-1}, N(S) = 210$

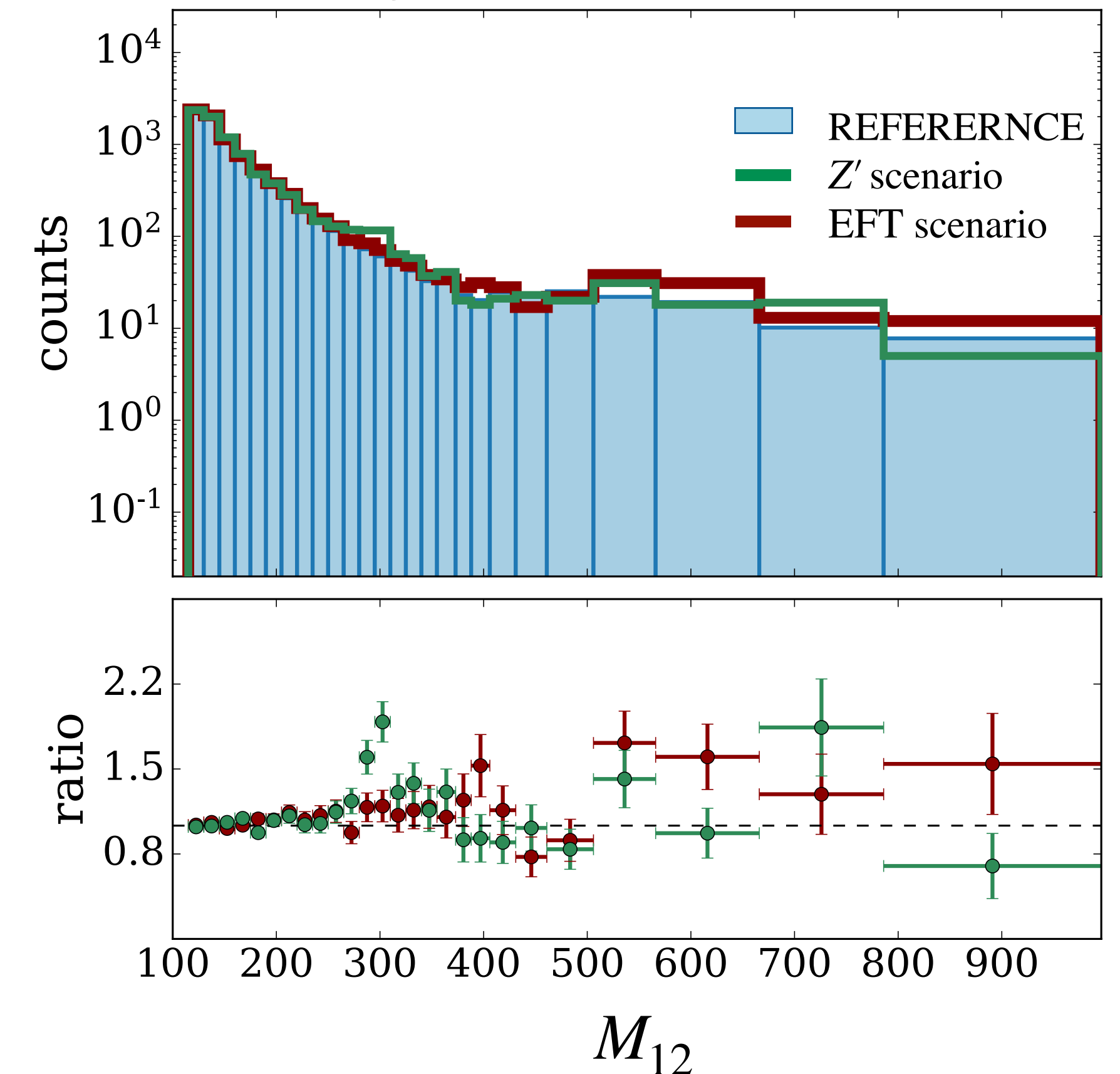
Non resonant excess in the tail of the two-body invariant mass

- **EFT scenario:** dimension-6 4-fermions contact operator:

$$\frac{c_W}{\Lambda} J_{L\mu}^a J_{La}^\mu$$

- Muon-like, electron-like regimes:
 $M_{12} > 100 \text{ GeV}, L = 0.35 \text{ fb}^{-1}, c_W = 1.0 \text{ TeV}^{-2}$
- Tau-like regime:
 $M_{12} > 120 \text{ GeV}, L = 1.1 \text{ fb}^{-1}, c_W = 0.25 \text{ TeV}^{-2}$

Example:
Tau-like regime



NOTE:

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

Di-body final state at the LHC

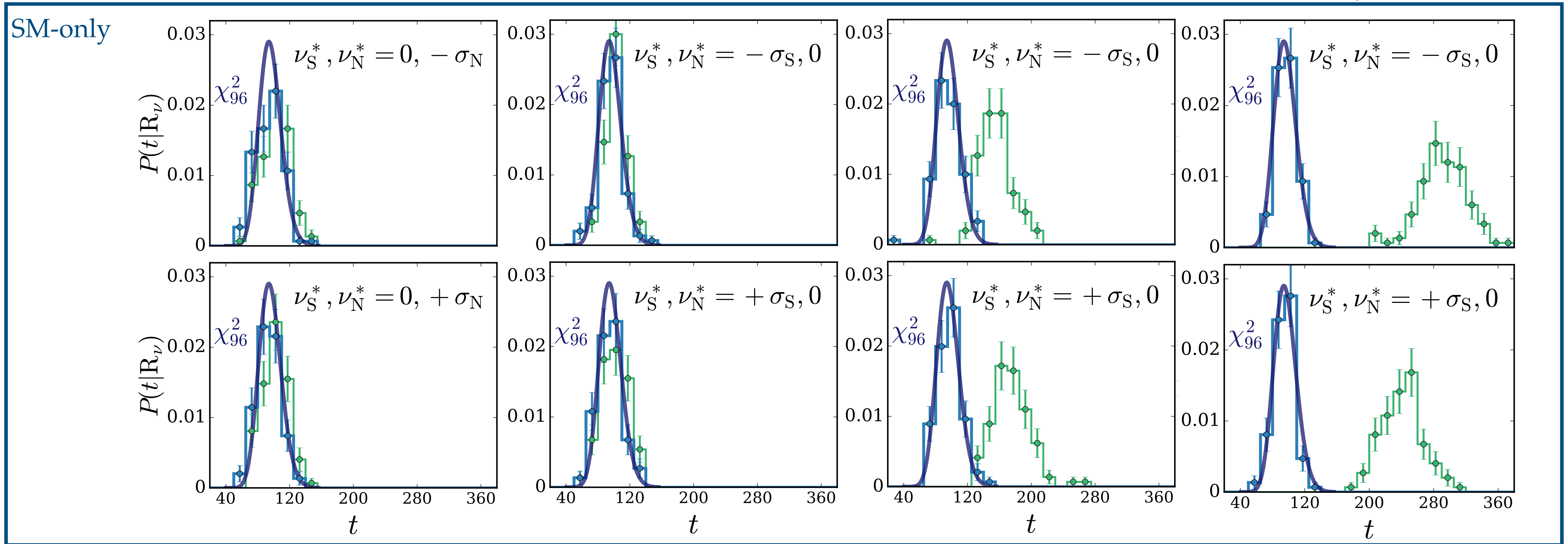
$\tau - \Delta$ validation

Negligible
systematic uncertainties

Muon-like regime

Electron-like regime

Tau-like regime



$\bullet t(D, A) = \tau(D, A) - \Delta(D, A)$

$\bullet \tau(D, A)$

DNN [5-5-5-5-1], #trainable parameters = 96, weight clipping = 2.16

Di-body final state at the LHC

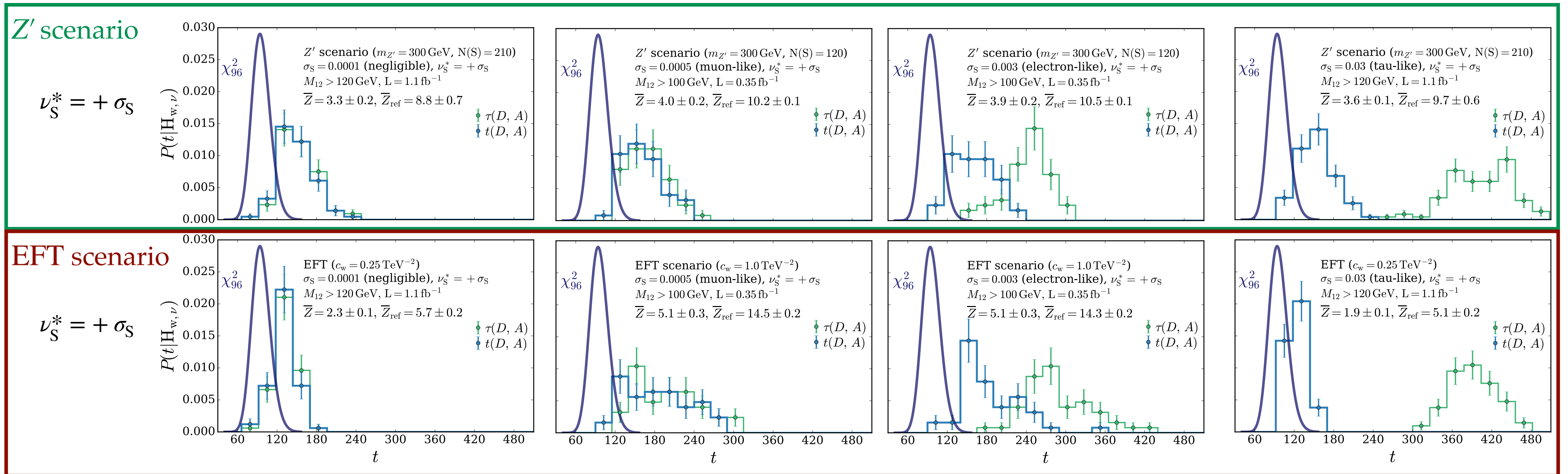
Sensitivity to New Physics scenarios

Negligible
systematic uncertainties

Muon-like regime

Electron-like regime

Tau-like regime



Z-score: $Z = \Phi^{-1} [1 - p]$

DNN [5-5-5-5-1], #trainable parameters = 96, weight clipping = 2.16

$t(D, A) = \tau(D, A) - \Delta(D, A)$

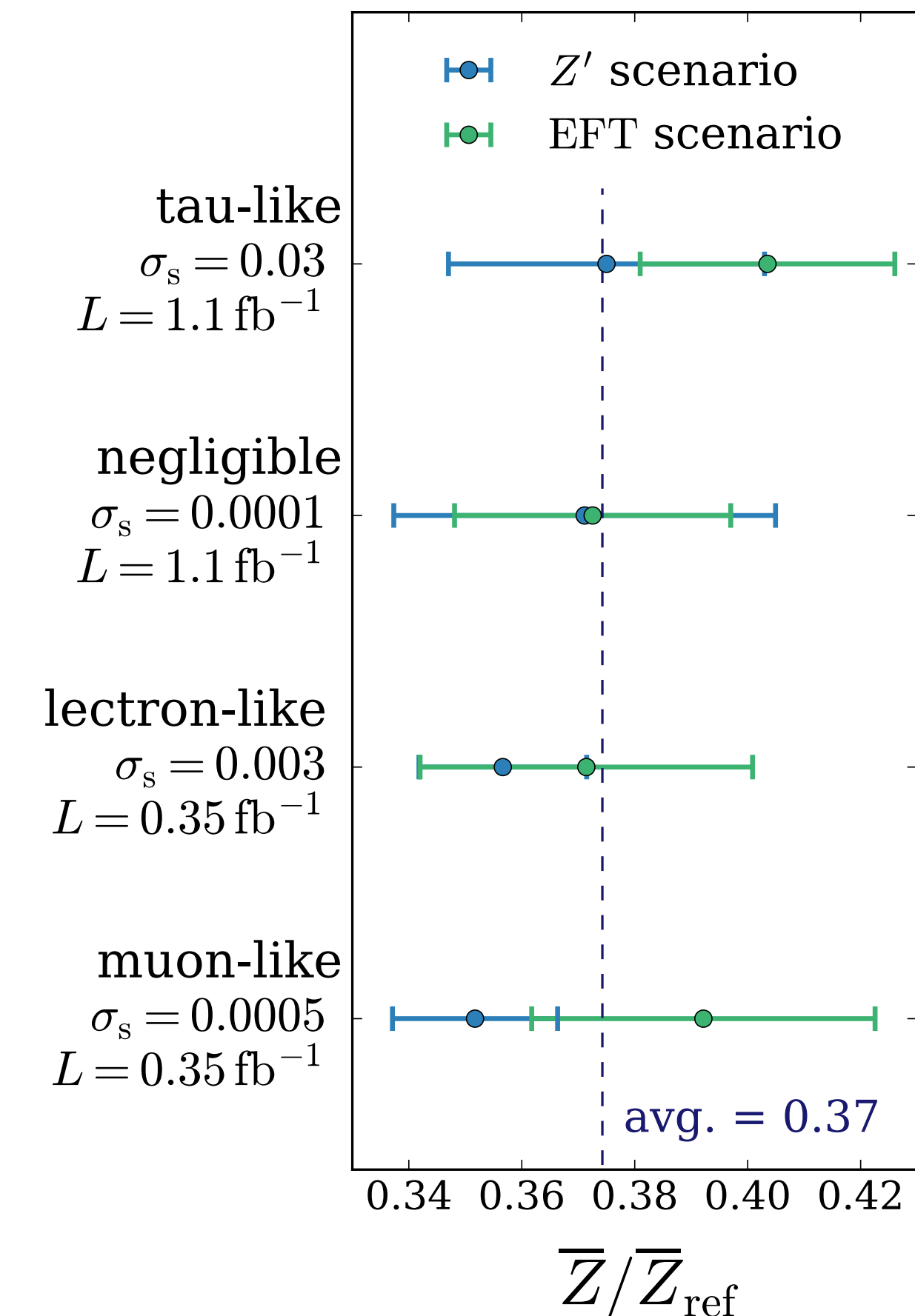
$\tau(D, A)$

Di-body final state at the LHC

Sensitivity to New Physics scenarios

Summary of the results:

- Comparable performances in the resonant and non-resonant scenarios:
 - NPLM is **simultaneously sensitive to multiple sources of New Physics**;
- Comparable performances at different systematic uncertainties regimes:
 - NPLM is robust against the presence of systematic uncertainties;
 - the presence of systematic uncertainties affects NPLM in the same measure as any other hypothesis test;
- **No information** about the New Physics **signal** has been provided to the algorithm at any step of its implementation:
 - The performances of NPLM are lower than any model-dependent strategy by construction ($\bar{Z}/\bar{Z}_{\text{ref}} = 0.37$);



Z-score: $Z = \Phi^{-1} [1 - p]$

- \bar{Z} : Z-score from NPLM

- \bar{Z}_{ref} : Z-score from a model-dependent (optimized) test statistics

Conclusions

Outlook on future perspectives

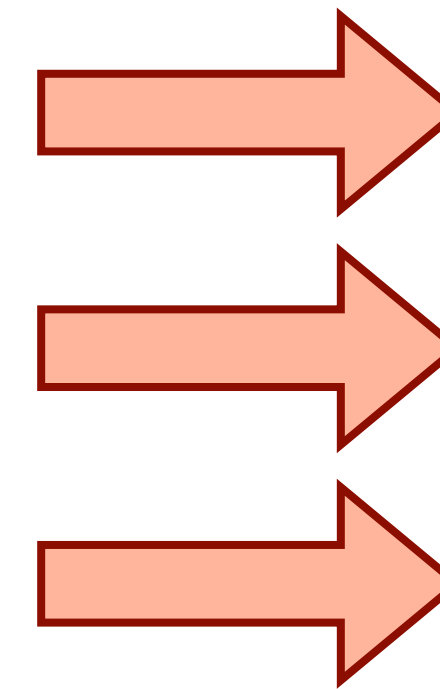
Current limitations and future developments:

- Accuracy and size of the Reference sample
- Accuracy in the (multivariate) modelling of the nuisance effects
- Training time

A solution from Kernel Methods (“Learning new physics efficiently with nonparametric methods” [2204.02317](#))

- **Optimisation** of NPLM sensitivity performances:

how do we choose the NN architecture? Is the regularization heuristic optimal? (ongoing work)



Set a **limit** on the actual **luminosity** that we are allowed to inspect, but do not obstacle the applicability of NPLM.

NPLM is ready to be performed on a real analysis at the LHC!

- ✓ Heuristic method to setup **multivariate** analysis
- ✓ Strategy to account for **systematic uncertainties**

Outlook on future perspectives

Getting started with NPLM

- [NPLM package](#): python-based package to run the NPLM analysis strategy
- [Tutorial](#) on 1D toy model for getting started

NPLM 0.0.6

Latest version

Released: Feb 1, 2022

```
pip install NPLM
```

package to run the New Physics Learning Machine (NPLM) algorithm.

Navigation

- Project description
- Release history
- Download files

Project links

- Homepage

Statistics

GitHub statistics:

- Stars: 1
- Forks: 1
- Open issues/PRs: 0

View statistics for this project via [Libraries.io](#), or by using [our public dataset on Google BigQuery](#)

Project description

NPLM_package

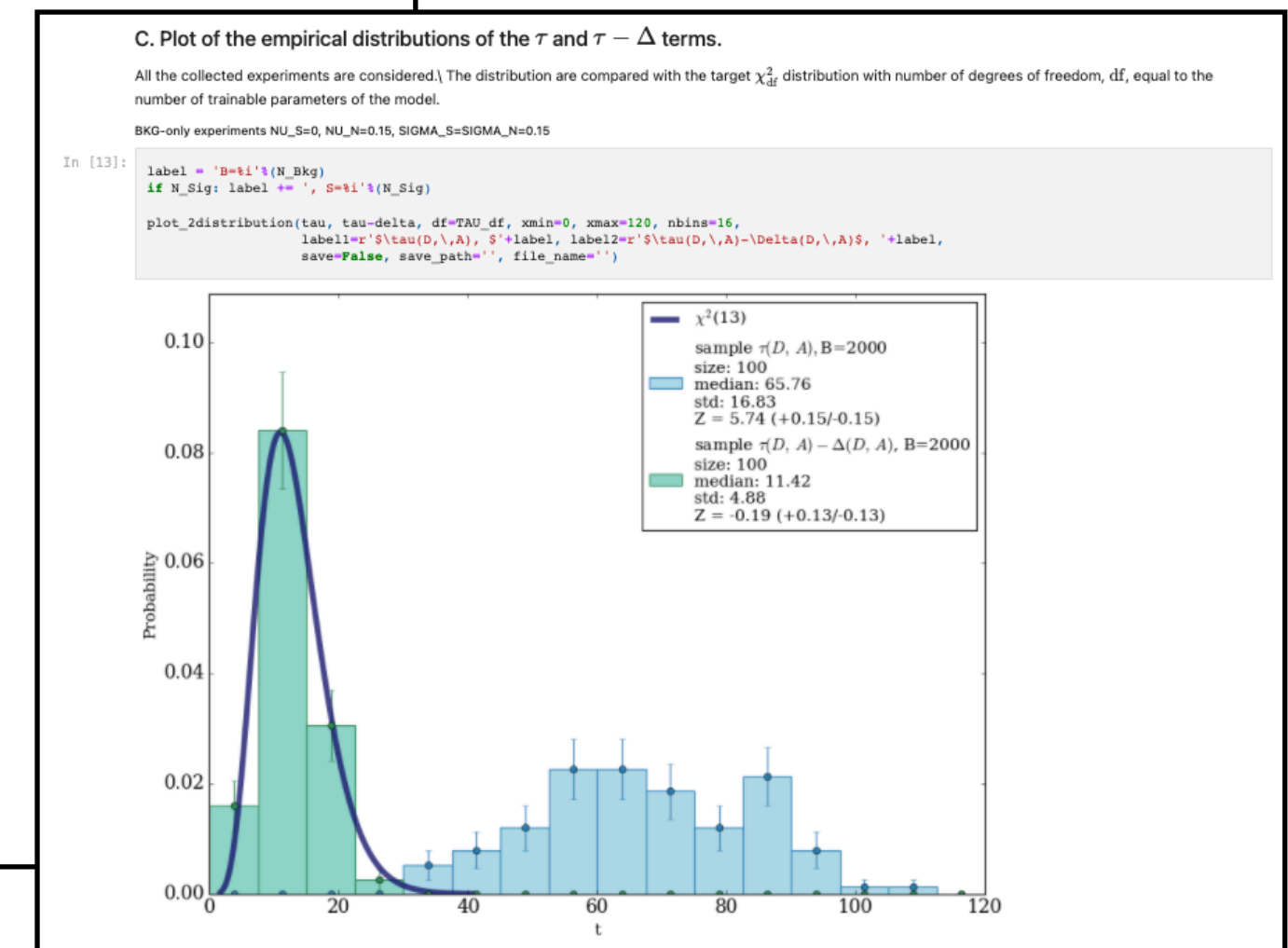
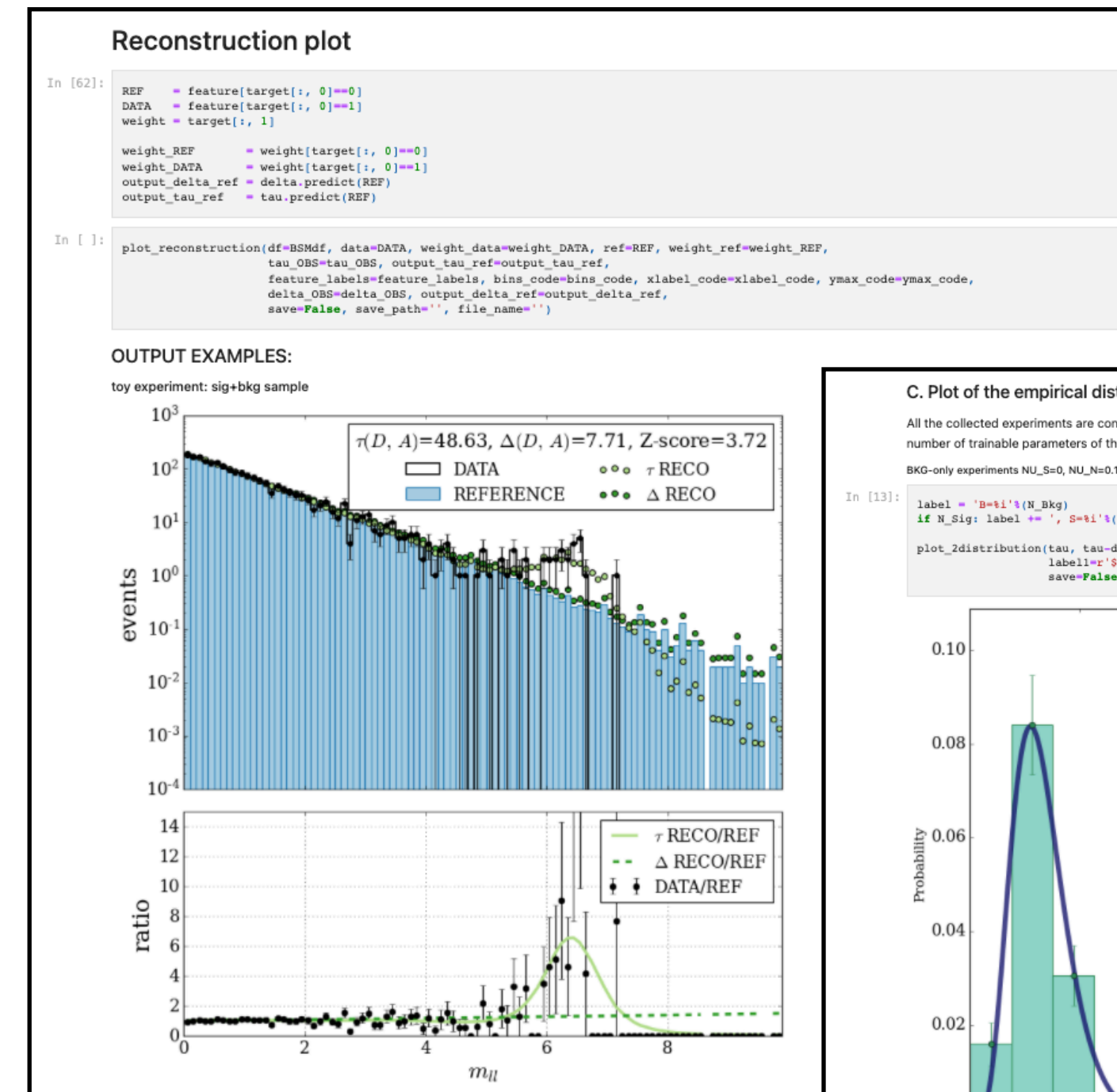
a package to implement the New Physics Learning Machine (NPLM) algorithm

Short description:

NPLM is a strategy to detect data departures from a given reference model, with no prior bias on the nature of the new physics model responsible for the discrepancy. The method employs neural networks, leveraging their virtues as flexible function approximants, but builds its foundations directly on the canonical likelihood-ratio approach to hypothesis testing. The algorithm compares observations with an auxiliary set of reference-distributed events, possibly obtained with a Monte Carlo event generator. It returns a p-value, which measures the compatibility of the reference model with the data. It also identifies the most discrepant phase-space region of the dataset, to be selected for further investigation. Imperfections due to mis-modelling in the reference dataset can be taken into account straightforwardly as nuisance parameters.

Related works:

- "Learning New Physics from a Machine" ([Phys. Rev. D](#))
- "Learning Multivariate New Physics" ([Eur. Phys. J. C](#))
- "Learning New Physics from an Imperfect Machine" ([arXiv](#))



Backup

New Physics Learning Machine (NPLM)

Main Concepts (negligible uncertainties)

- Goal: performing a **log-likelihood-ratio hypothesis test**
(End-to-end strategy, from the data to a p -value for the discovery)

$$t(\mathcal{D}) = \max_{\mathbf{w}} \left[2 \log \frac{\mathcal{L}(H_{\mathbf{w}} | \mathcal{D})}{\mathcal{L}(R_0 | \mathcal{D})} \right]$$

R_0 : reference (null) hypothesis
 $H_{\mathbf{w}}$: alternative hypothesis

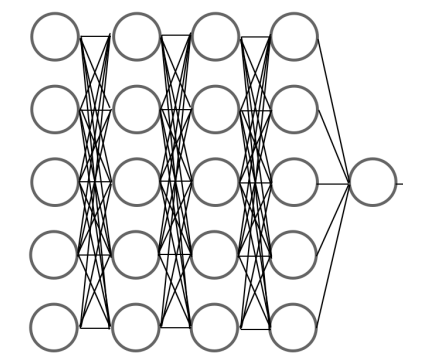
- Exploiting a Neural Network (NN) to **parametrize** the data distribution in terms of a Reference distribution (R_0)

$$n(x | T) \approx n(x | H_{\hat{\mathbf{w}}}) = n(x | R_0) e^{f(x, \hat{\mathbf{w}})}$$

True (T) data distribution

Data distribution learnt by the NN

Reference distribution



NN model

Unknown

Alternative hypothesis

Null hypothesis (SM)

- **Signal-model-independent**: reduced assumptions on the signal hypothesis

New Physics Learning Machine (NPLM)

Main Concepts (negligible uncertainties)

Maximum Likelihood from minimal loss:

Test statistic

$$\bar{t}(\mathcal{D}) = 2 \max_{\mathbf{w}} \log \left[\frac{\mathcal{L}(\mathbf{H}_{\mathbf{w}}|\mathcal{D})}{\mathcal{L}(\mathbf{R}_0|\mathcal{D})} \right] = -2 \min_{\mathbf{w}} \{ \bar{L}[f(\cdot; \mathbf{w})] \}$$

Loss function

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

\mathbf{w} : trainable parameters on the NN model

\mathcal{D} : data sample

\mathcal{R} : reference sample (built according to the \mathbf{R}_0 hypothesis); could be weighted (w_x)

Assumptions:

- $N_R \gg N_D$ 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 $\sum_{x \in \mathcal{R}} w_x = N(\mathbf{R}_0)$

New Physics Learning Machine (NPLM)

Main Concepts (negligible uncertainties)

Maximum Likelihood from minimal loss:

Test statistic

$$\bar{t}(\mathcal{D}) = 2 \max_{\mathbf{w}} \log \left[\frac{\mathcal{L}(\mathbf{H}_{\mathbf{w}} | \mathcal{D})}{\mathcal{L}(\mathbf{R}_0 | \mathcal{D})} \right] = -2 \min_{\mathbf{w}} \{ \bar{L}[f(\cdot; \mathbf{w})] \}$$

Loss function

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

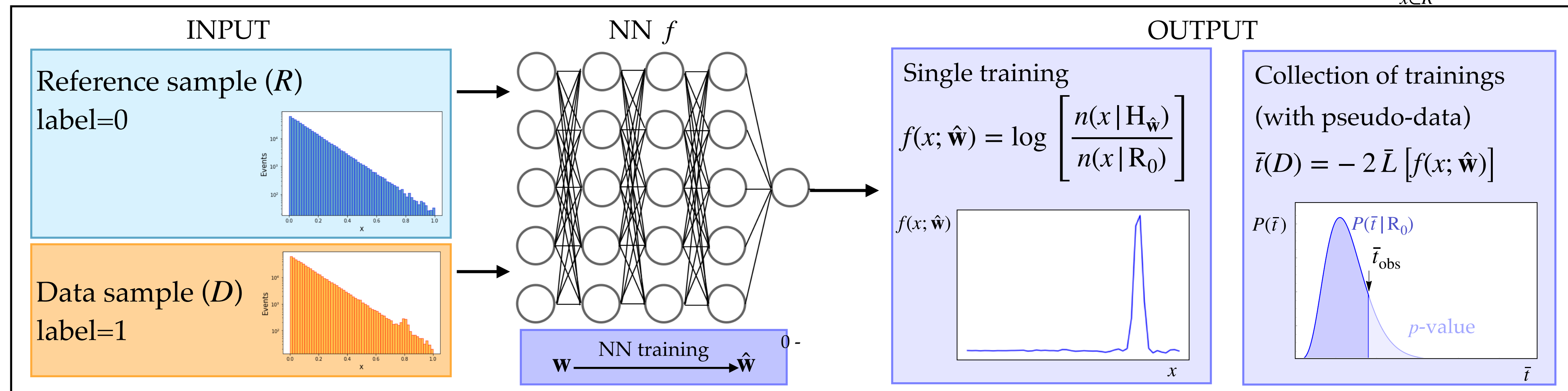
\mathbf{w} : trainable parameters on the NN model

\mathcal{D} : data sample

\mathcal{R} : reference sample (built according to the \mathbf{R}_0 hypothesis); could be weighted (w_x)

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 $\sum_{x \in \mathcal{R}} w_x = N(\mathbf{R}_0)$



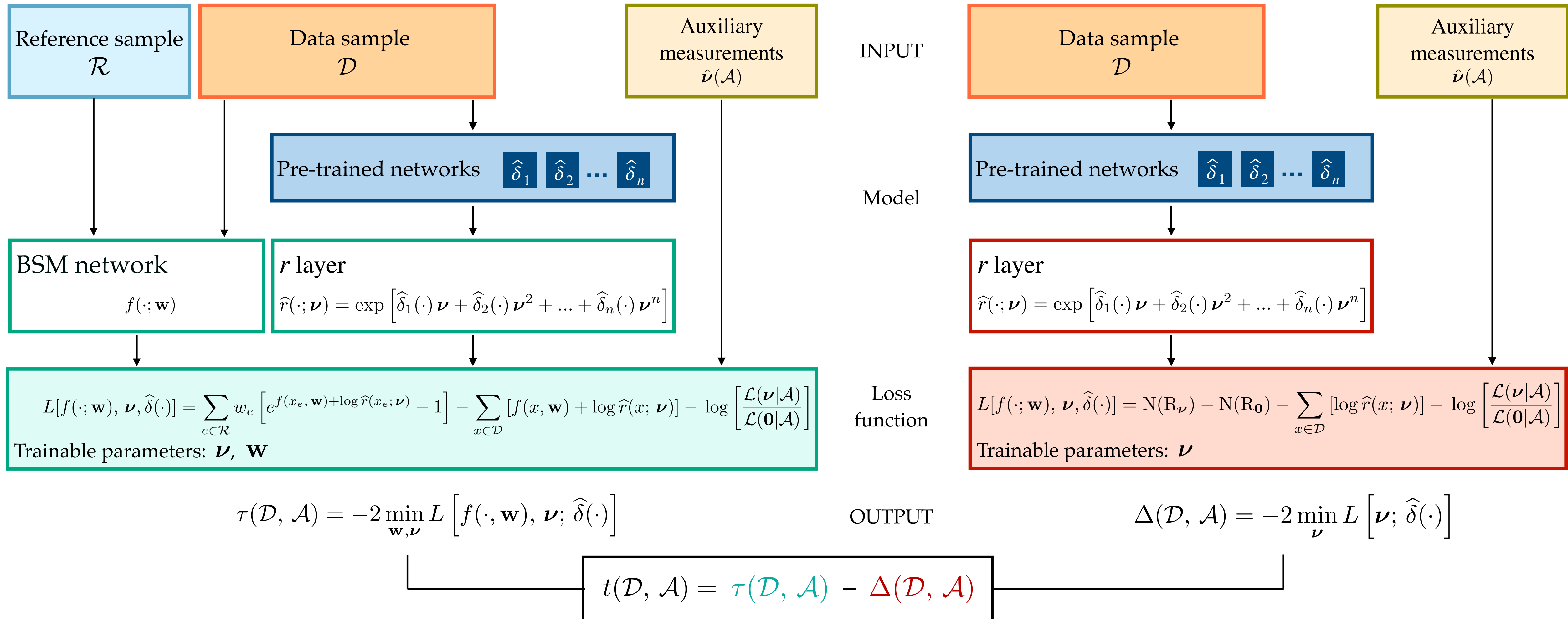
"Learning New Physics from a Machine" [Phys. Rev. D](#)

New Physics Learning Machine (NPLM)

Including systematic uncertainties

τ term

Δ term



New Physics Learning Machine (NPLM)

Including systematic uncertainties

Validation of the $(\tau - \Delta)$ procedure

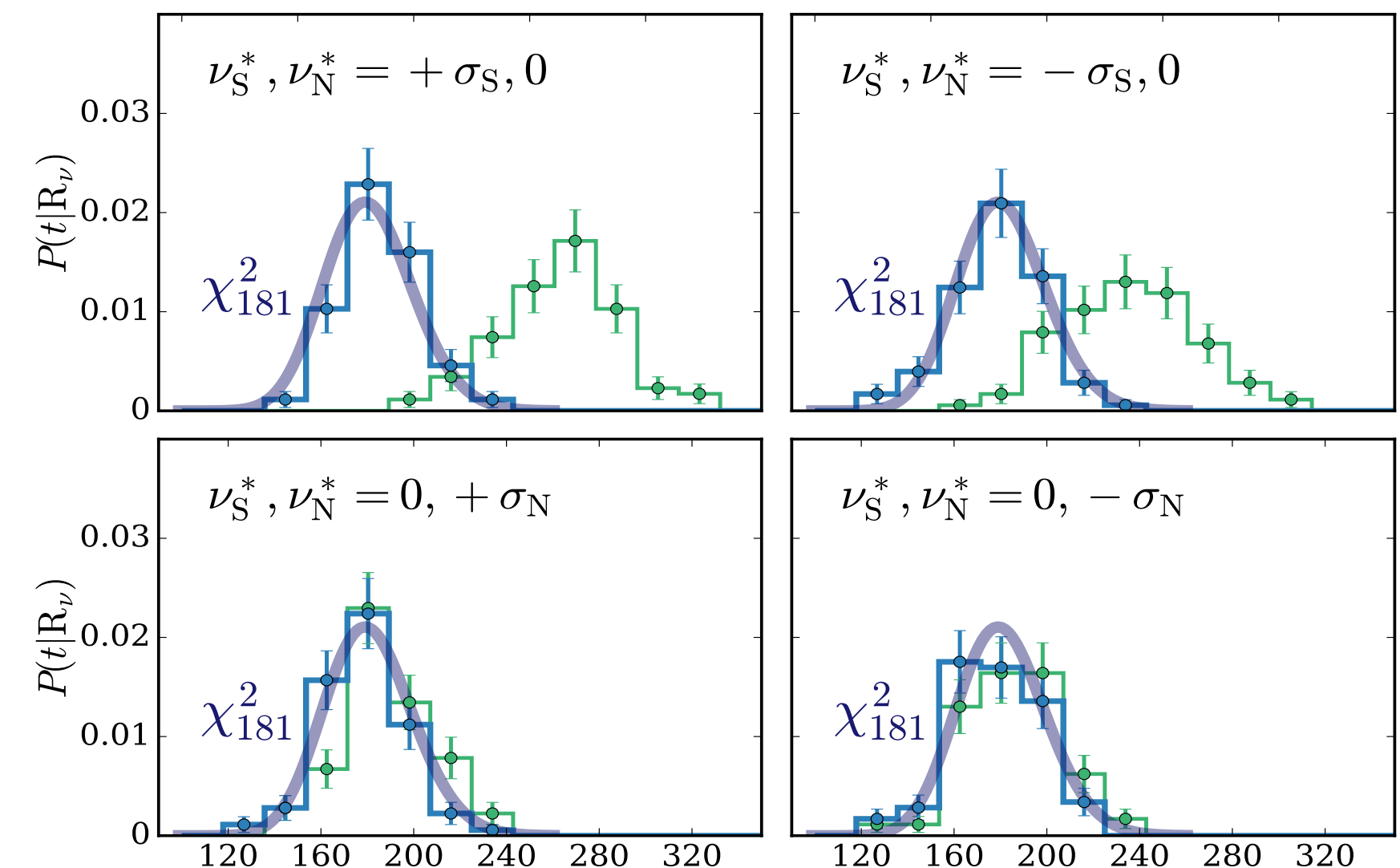
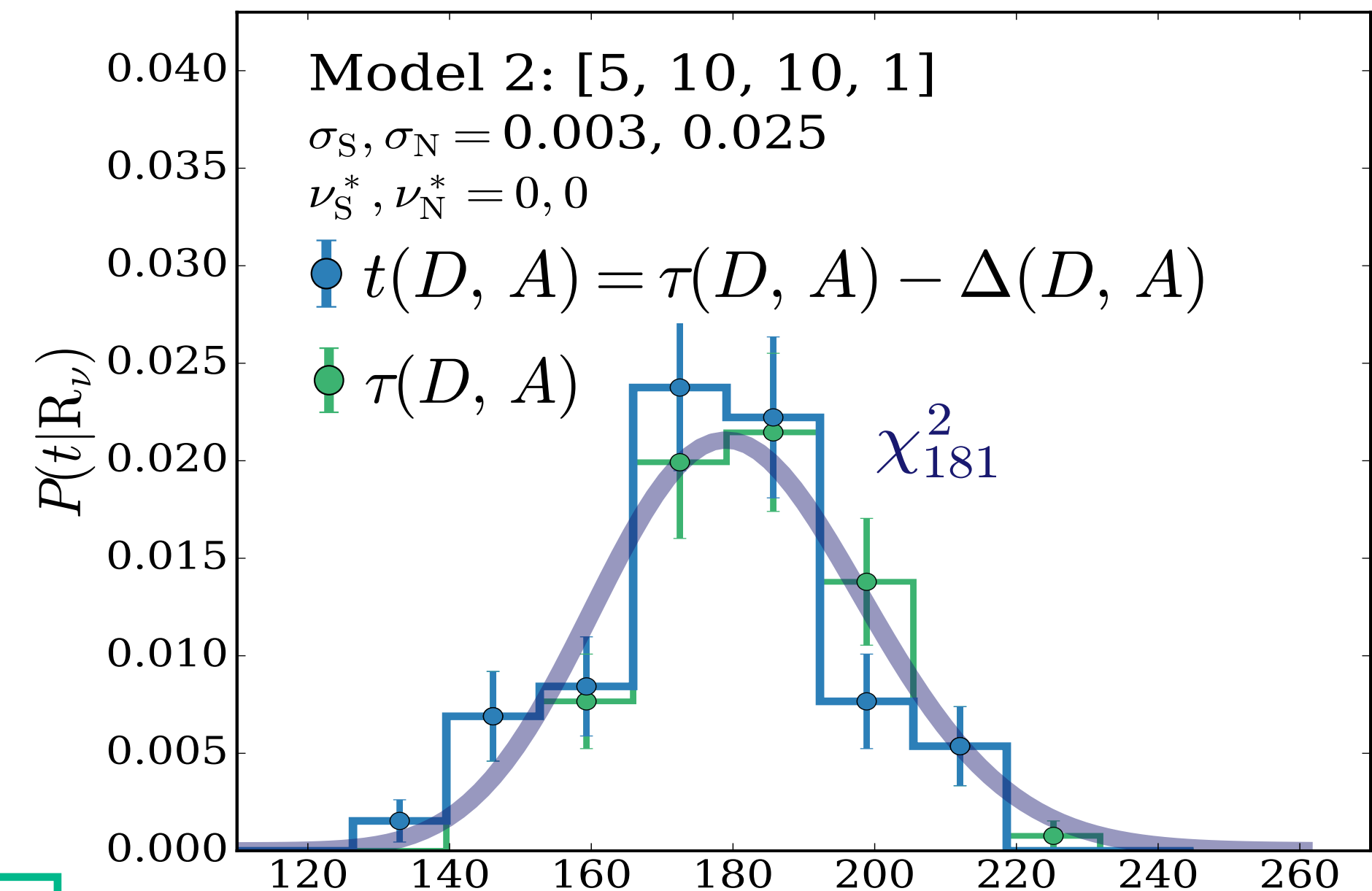
“Toy Data”: test the procedure on simulated toys following the Reference (SM) hypothesis with generation value for the nuisance parameters $\nu^* = \pm\sigma_\nu$:

$$\mathcal{D} \sim R_{\nu^*}, \quad \nu^* = \pm\sigma_\nu$$

The \bar{t} distribution under the reference hypothesis R_{ν^*} is **compatible with the target** $\chi^2_{|w|}$ for values of the true nuisance parameters within the uncertainty ($\nu^* = \pm\sigma_\nu$).

\bar{t} is **independent** of the true value of the nuisance parameters!

We can build a *frequentistic* test statistic relying on the asymptotic $\chi^2_{|w|}$.



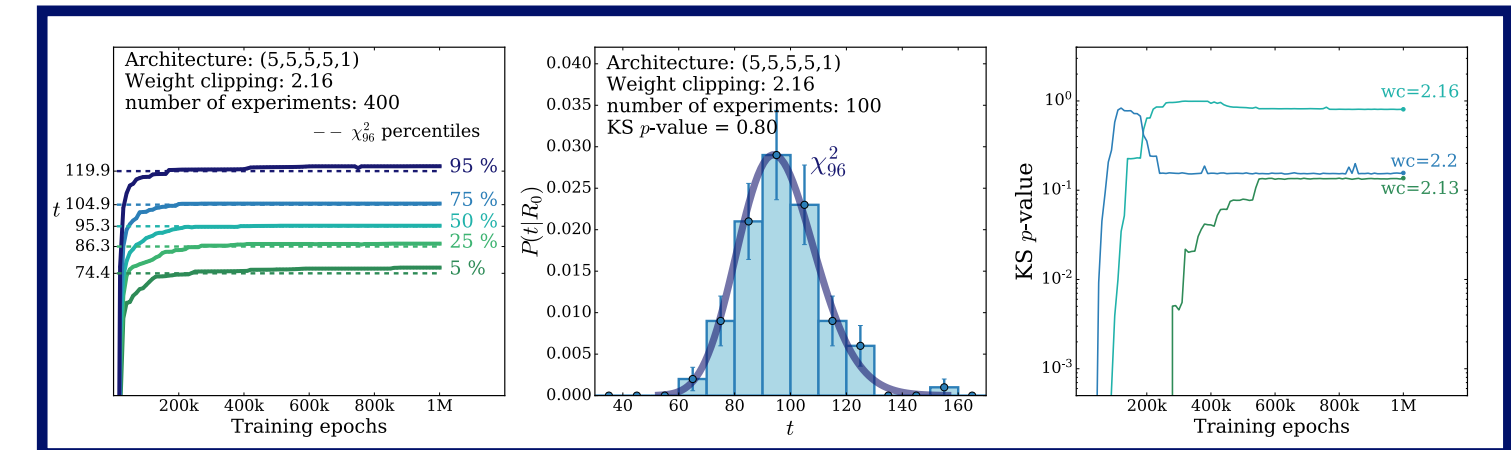
New Physics Learning Machine (NPLM)

Including systematic uncertainties

Final procedure in steps:

1. NN (f) REGULARIZATION:

weight clipping tuning \rightarrow target $\chi^2_{|w|}$;



2. NUISANCE TAYLOR'S EXPANSION LEARNING:

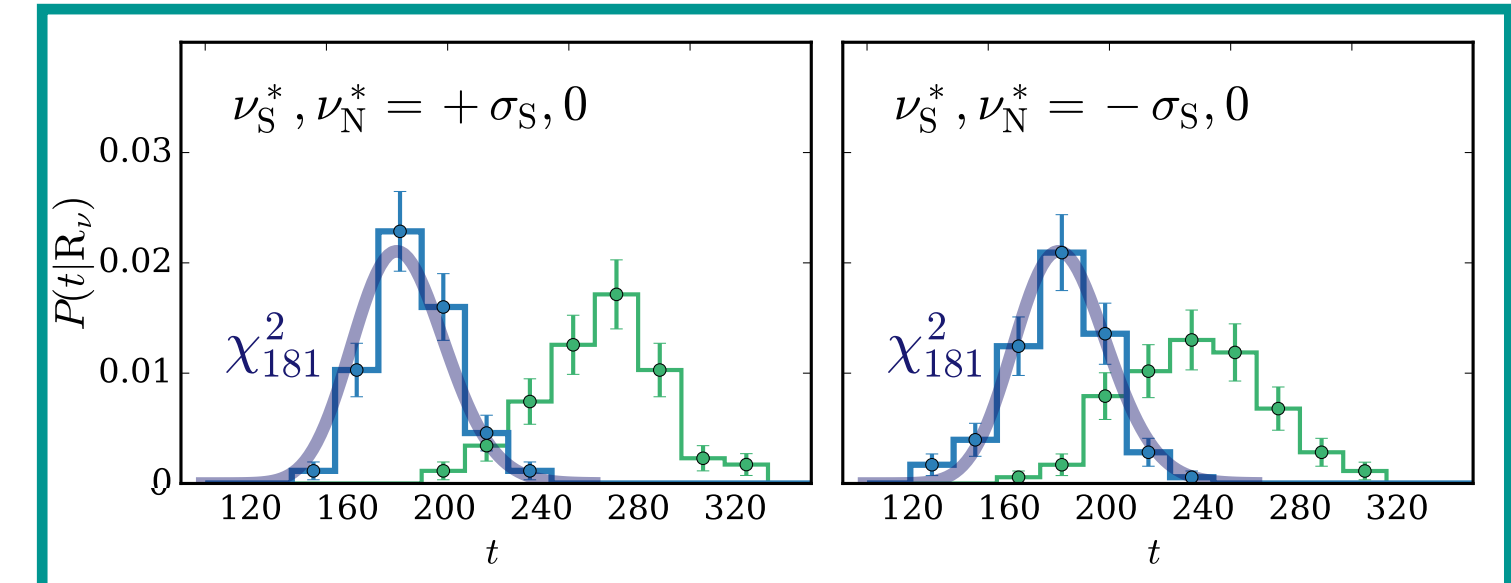
modelling $\hat{r}(x; \nu) = \exp \left[\hat{\delta}_1(x) \nu + \hat{\delta}_2(x) \nu^2 + \dots \right]$;

$$\hat{r}(x; \nu) = \exp \left[\underbrace{\hat{\delta}_1(x)}_{\text{NN 1}} \nu + \underbrace{\hat{\delta}_2(x)}_{\text{NN 2}} \nu^2 + \dots \right]$$

3. $\tau - \Delta$ VALIDATION:

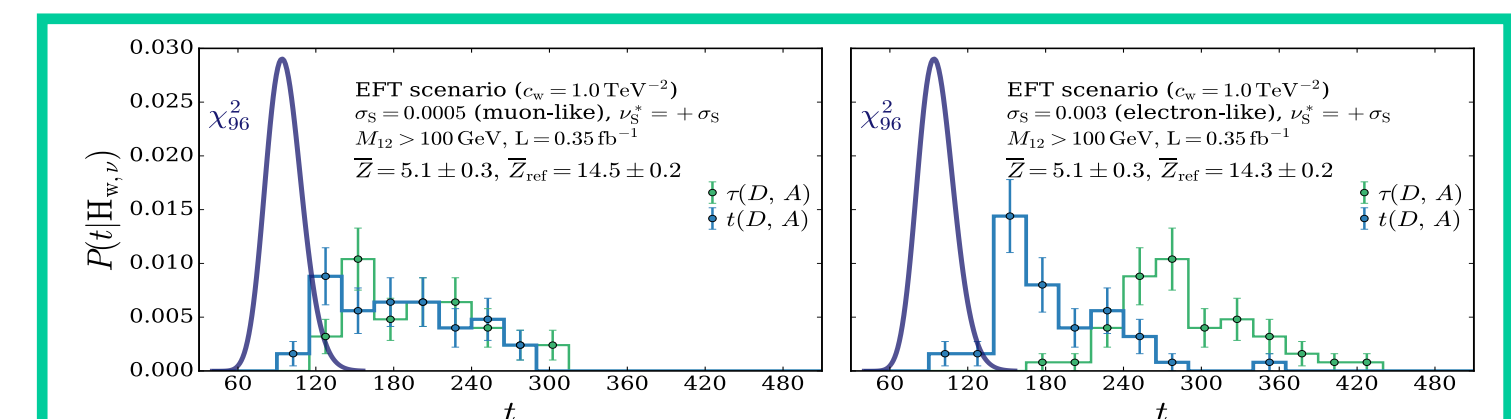
$$\mathcal{D} \sim R_{\nu^*}, \quad \nu^* = \pm \sigma_\nu$$

Verifying that the target $\chi^2_{|w|}$ is always recovered;



4. TESTING THE DATA:

running the procedure on real data.



Di-body final state at the LHC

Sensitivity to New Physics scenarios

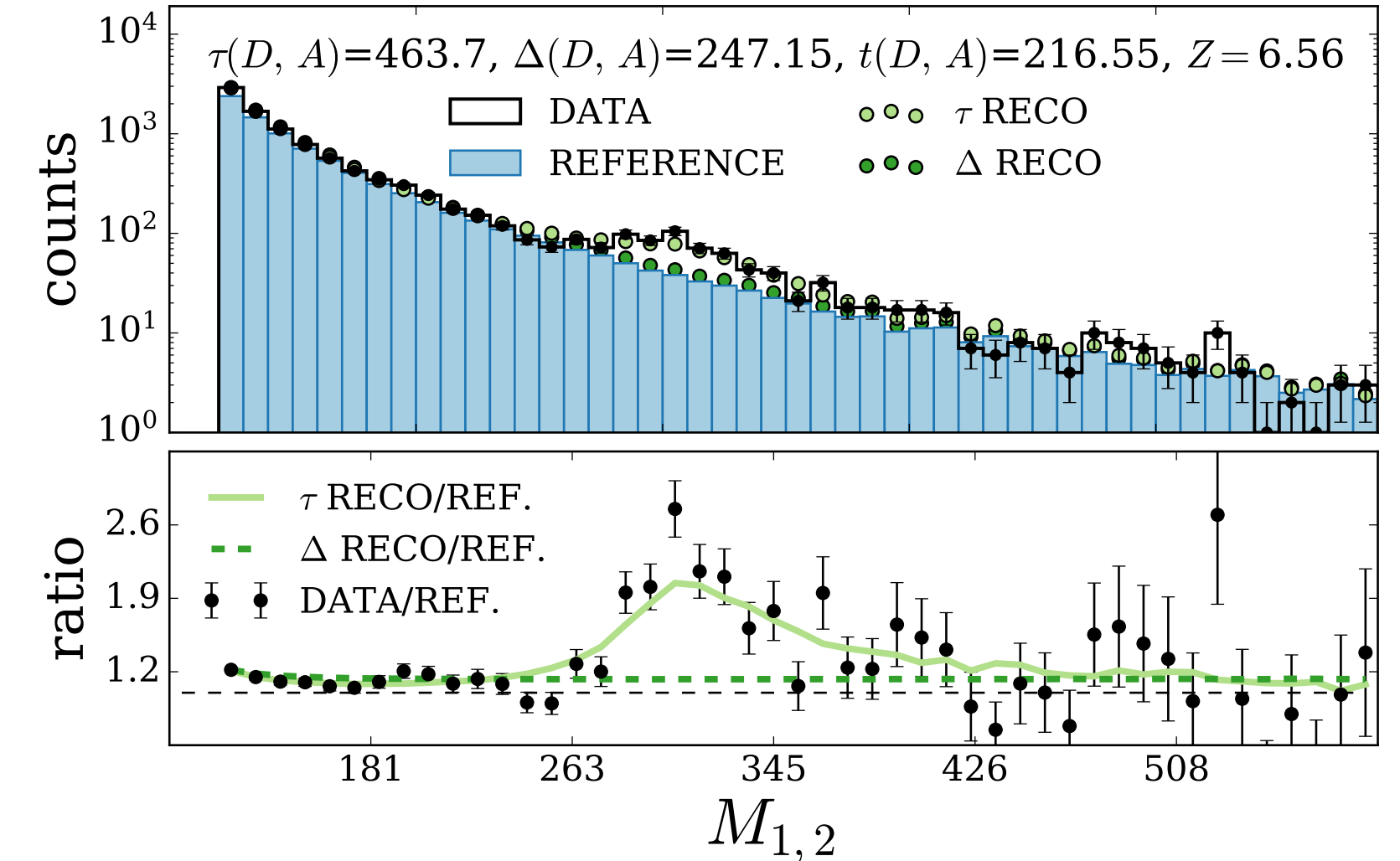
Summary of the results:

- Comparable performances in the resonant and non-resonant scenarios:
 - NPLM is **simultaneously sensitive to multiple sources of New Physics**;
- Comparable performances at different systematic uncertainties regimes:
 - NPLM is robust against the presence of systematic uncertainties;
 - the presence of systematic uncertainties affects NPLM in the same measure as any other hypothesis test;
- **No information** about the New Physics **signal** has been provided to the algorithm at any step of its implementation:
 - The performances of NPLM are lower than any model-dependent strategy by construction ($\bar{Z}/\bar{Z}_{\text{ref}} = 0.37$);
- NPLM is able to *learn* non trivial combinations of the input variables and point to the source of the significant discrepancy.

$$\tau \text{ reconstruction: } n(x | H_{\hat{w}, \hat{v}}) = n(x | R_0) \frac{n(x | R_{\hat{v}})}{n(x | R_0)} e^{f(x; \hat{w})}$$

$$\Delta \text{ reconstruction: } n(x | R_{\hat{v}})$$

Z' scenario (tau-like regime), $m = 300 \text{ GeV}$, $N(S) = 210$



EFT scenario (tau-like regime), $c_W = 0.25 \text{ TeV}^{-2}$

