# Applications of Likelihood Learning

**Benasque 2023**

**Alfredo Glioti & Jaco ter Hoeve**
IPhT & VU Amsterdam

# Recap: likelihood ratio from ML

‣ Starting from two balanced datasets $\mathscr{D}_{\mathrm{SM}}$ and $\mathscr{D}_{\mathrm{EFT}}$ drawn from $f(\boldsymbol{x}\,|\,\mathrm{SM})$ and $f(\boldsymbol{x}\,|\,\mathrm{EFT})$, we minimise e.g. the cross-entropy loss

$$L[g(\boldsymbol{x})] = -\frac{1}{N} \sum_{e \in \mathscr{D}_{\mathrm{EFT}}} w_e \log(1 - g(\boldsymbol{x}_e)) - \frac{1}{N} \sum_{\mathscr{D}_{\mathrm{SM}}} w_e \log g(\boldsymbol{x}_e)$$

Event weights

$\{m_{t\bar{t}}, \eta_l, \Delta\phi, \dots\}$

‣ The learned decision boundary $g(\boldsymbol{x})$ is one-to-one with the likelihood ratio (LR) as $N \to \infty$

$$\frac{\delta L}{\delta g} = 0 \implies \hat{g}(\boldsymbol{x}) = \left(1 + \frac{f(\boldsymbol{x}\,|\,\mathrm{EFT})}{f(\boldsymbol{x}\,|\,\mathrm{SM})}\right)^{-1} \equiv \frac{1}{1 + r(\boldsymbol{x})}$$

Parameterise with NNs

# The experimental pipeline

We are progressively moving through the simulation chain (latent space)

$$p(x|c) \sim \int dz_{\mathrm{det}} dz_{\mathrm{shower}} dz_{\mathrm{parton}} p(x|z_{\mathrm{det}}) p(z_{\mathrm{det}}|z_{\mathrm{shower}}) p(z_{\mathrm{shower}}|z_{\mathrm{parton}}) p(z_{\mathrm{parton}}|c)$$



**DELPHES** fast simulation

**MADGRAPH5** aMC@NLO

**SMEFT@NLO**

| **(Likelihood free) Inference** | Observed | $\longrightarrow$ | Hidden/latent variables | $\longrightarrow$ | POI |
|---|---|---|---|---|---|

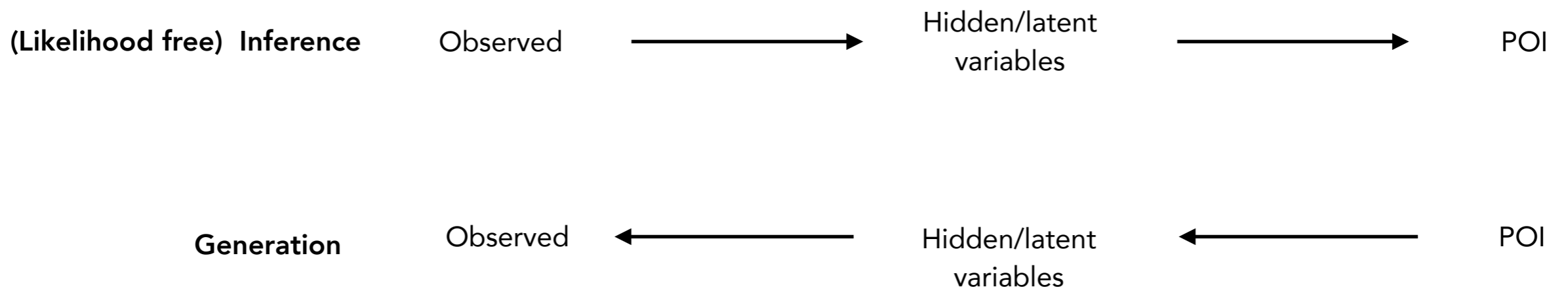| **Generation** | Observed | $\longleftarrow$ | Hidden/latent variables | $\longleftarrow$ | POI |
|---|---|---|---|---|---|

# The experimental pipeline

We are progressively moving through the simulation chain (latent space)

$$p(x|c) \sim \int dz_{\text{det}} dz_{\text{shower}} dz_{\text{parton}} p(x|z_{\text{det}}) p(z_{\text{det}}|z_{\text{shower}}) p(z_{\text{shower}}|z_{\text{parton}}) p(z_{\text{parton}}|c)$$
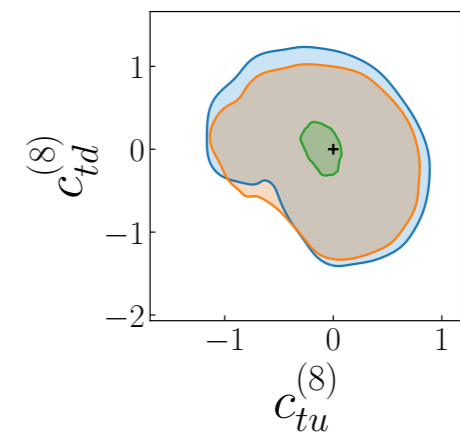


Unbinned unfolding, Omnifold [1911.09107]





**SMEFT@NLO**

# Why interesting for us?

‣ Global efforts **reinterpret** existing "SM measurements" in an EFT context

‣ But which measurements are the most **sensitive** to EFT parameters?

   ‣ Inclusive, single to multi-differential (which variables)

   ‣ Binned or unbinned, which binning?

Framework needed to integrate unbinned multivariate observables into **global** SMEFT fits

   ‣ **Optimal** bounds on the EFT parameters

   ‣ Useful **diagnosis tool** to assess information loss

# Applications of likelihood learning

Focusses on global EFT fits

Reweighting for more accurate learning

**Unbinned multivariate observables for global SMEFT analyses from machine learning**

Raquel Gomez Ambrosio,[1,2] Jaco ter Hoeve,[3,4] Maeve Madigan,[5] Juan Rojo,[3,4] and Veronica Sanz[6,7]

[1] Dipartimento di Fisica "G. Occhialini", Universita degli Studi di Milano-Bicocca,
and INFN, Sezione di Milano Bicocca, Piazza della Scienza 3, I – 20126 Milano, Italy
[2] Dipartimento di Fisica, Università di Torino, and INFN, Sezione di Torino, Via P. Giuria 1, 10125 Torino, Italy
[3] Department of Physics and Astronomy, VU Amsterdam, 1081HV Amsterdam, The Netherlands
[4] Nikhef Theory Group, Science Park 105, 1098 XG Amsterdam, The Netherlands
[5] DAMTP, University of Cambridge, Wilberforce Road, Cambridge CB3 0WA, UK
[6] Instituto de Física Corpuscular (IFIC), Universidad de Valencia-CSIC, E-46980 Valencia, Spain
[7] Department of Physics and Astronomy, University of Sussex, Brighton BN1 9QH, UK

### Abstract

Theoretical interpretations of particle physics data, such as the determination of the Wilson coefficients of the Standard Model Effective Field Theory (SMEFT), often involve the inference of multiple parameters from a global dataset. Optimizing such interpretations requires the identification of observables that exhibit the highest possible sensitivity to the underlying theory parameters. In this work we develop a flexible open source framework, ML4EFT, enabling the integration of unbinned multivariate observables into global SMEFT fits. As compared to traditional measurements, such observables enhance the sensitivity to the theory parameters by preventing the information loss incurred when binning in a subset of final-state kinematic variables. Our strategy combines machine learning regression and classification techniques to parameterize high-dimensional likelihood ratios, using the Monte Carlo replica method to estimate and propagate methodological uncertainties. As a proof of concept we construct unbinned multivariate observables for top-quark pair and Higgs+$Z$ production at the LHC, demonstrate their impact on the SMEFT parameter space as compared to binned measurements, and study the improved constraints associated to multivariate inputs. Since the number of neural networks to be trained scales quadratically with the number of parameters and can be fully parallelized, the ML4EFT framework is well-suited to construct unbinned multivariate observables which depend on up to tens of EFT coefficients, as required in global fits.

1

**Boosting likelihood learning with event reweighting**

Siyu Chen[1], Alfredo Glioti[2], Giuliano Panico[3,4], and Andrea Wulzer[5,6]

[1] Institut de Théorie des Phénomenes Physiques, EPFL, Lausanne, Switzerland
[2] Université Paris-Saclay, CNRS, CEA, Institut de Physique Théorique, 91191, Gif-sur-Yvette, France
[3] Dipartimento di Fisica e Astronomia, Università di Firenze,
Via G. Sansone 1, 50019 Sesto Fiorentino, Italy
[4] INFN, Sezione di Firenze, Via G. Sansone 1, 50019 Sesto Fiorentino, Italy
[5] Institut de Física d'Altes Energies (IFAE), The Barcelona Institute of Science and Technology (BIST),
Campus UAB, 08193 Bellaterra, Barcelona, Spain
[6] ICREA, Institució Catalana de Recerca i Estudis Avançats, Passeig de Lluís Companys 23, 08010
Barcelona, Spain

### Abstract

Extracting maximal information from experimental data requires access to the likelihood function, which however is never directly available for complex experiments like those performed at high energy colliders. Theoretical predictions are obtained in this context by Monte Carlo events, which do furnish an accurate but abstract and implicit representation of the likelihood. Strategies based on statistical learning are currently being developed to infer the likelihood function explicitly by training a continuous-output classifier on Monte Carlo events. In this paper, we investigate the usage of Monte Carlo events that incorporate the dependence on the parameters of interest by reweighting. This enables more accurate likelihood learning with less training data and a more robust learning scheme that is more suited for automation and extensive deployment. We illustrate these advantages in the context of LHC precision probes of new Effective Field Theory interactions.

Alfredo

# Applications of likelihood learning

Focusses on global EFT fits

Reweighting for more accurate learning

## Unbinned multivariate observables for global SMEFT analyses from machine learning

Raquel Gomez Ambrosio,[1,2] Jaco ter Hoeve,[3,4] Maeve Madigan,[5] Juan Rojo,[3,4] and Veronica Sanz[6,7]

[1] Dipartimento di Fisica "G. Occhialini", Universita degli Studi di Milano-Bicocca,
and INFN, Sezione di Milano Bicocca, Piazza della Scienza 3, I – 20126 Milano, Italy
[2] Dipartimento di Fisica, Università di Torino, and INFN, Sezione di Torino, Via P. Giuria 1, 10125 Torino, Italy
[3] Department of Physics and Astronomy, VU Amsterdam, 1081HV Amsterdam, The Netherlands
[4] Nikhef Theory Group, Science Park 105, 1098 XG Amsterdam, The Netherlands
[5] DAMTP, University of Cambridge, Wilberforce Road, Cambridge CB3 0WA, UK
[6] Instituto de Física Corpuscular (IFIC), Universidad de Valencia-CSIC, E-46980 Valencia, Spain
[7] Department of Physics and Astronomy, University of Sussex, Brighton BN1 9QH, UK

### Abstract

Theoretical interpretations of particle physics data, such as the determination of the Wilson coefficients of the Standard Model Effective Field Theory (SMEFT), often involve the inference of multiple parameters from a global dataset. Optimizing such interpretations requires the identification of observables that exhibit the highest possible sensitivity to the underlying theory parameters. In this work we develop a flexible open source framework, ML4EFT, enabling the integration of unbinned multivariate observables into global SMEFT fits. As compared to traditional measurements, such observables enhance the sensitivity to the theory parameters by preventing the information loss incurred when binning in a subset of final-state kinematic variables. Our strategy combines machine learning regression and classification techniques to parameterize high-dimensional likelihood ratios, using the Monte Carlo replica method to estimate and propagate methodological uncertainties. As a proof of concept we construct unbinned multivariate observables for top-quark pair and Higgs+$Z$ production at the LHC, demonstrate their impact on the SMEFT parameter space as compared to binned measurements, and study the improved constraints associated to multivariate inputs. Since the number of neural networks to be trained scales quadratically with the number of parameters and can be fully parallelized, the ML4EFT framework is well-suited to construct unbinned multivariate observables which depend on up to tens of EFT coefficients, as required in global fits.

arXiv:2211.02058v2 [hep-ph] 23 May 2023

## Boosting likelihood learning with event reweighting

Siyu Chen[1], Alfredo Glioti[2], Giuliano Panico[3,4], and Andrea Wulzer[5,6]

[1] Institut de Théorie des Phénomènes Physiques, EPFL, Lausanne, Switzerland
[2] Université Paris-Saclay, CNRS, CEA, Institut de Physique Théorique, 91191, Gif-sur-Yvette, France
[3] Dipartimento di Fisica e Astronomia, Università di Firenze,
Via G. Sansone 1, 50019 Sesto Fiorentino, Italy
[4] INFN, Sezione di Firenze, Via G. Sansone 1, 50019 Sesto Fiorentino, Italy
[5] Institut de Física d'Altes Energies (IFAE), The Barcelona Institute of Science and Technology (BIST),
Campus UAB, 08193 Bellaterra, Barcelona, Spain
[6] ICREA, Institució Catalana de Recerca i Estudis Avançats, Passeig de Lluís Companys 23, 08010
Barcelona, Spain

### Abstract

Extracting maximal information from experimental data requires access to the likelihood function, which however is never directly available for complex experiments like those performed at high energy colliders. Theoretical predictions are obtained in this context by Monte Carlo events, which do furnish an accurate but abstract and implicit representation of the likelihood. Strategies based on statistical learning are currently being developed to infer the likelihood function explicitly by training a continuous-output classifier on Monte Carlo events. In this paper, we investigate the usage of Monte Carlo events that incorporate the dependence on the parameters of interest by reweighting. This enables more accurate likelihood learning with less training data and a more robust learning scheme that is more suited for automation and extensive deployment. We illustrate these advantages in the context of LHC precision probes of new Effective Field Theory interactions.

arXiv:2308.05704v1 [hep-ph] 10 Aug 2023

Alfredo

# The ML4EFT framework

Open-source NN-based python framework for the integration of unbinned multivariate observables into global SMEFT fits

‣ **Goal**: provide optimal constraints on the SMEFT

‣ **Diagnostic tool:** what is the information loss incurred by a particular choice of bins?

‣ **Projections:** how will SMEFT constraints improve if unbinned data are made available?



Modular structure, easy to maintain, well documented

# Anticipating global fits

‣ Global EFT fits typically feature ~50 WCs and thus efficient scaling with the number of WCs becomes essential

‣ **ML4EFT 1.0:** learn the coefficient functions separately and combine afterwards

Assumes no sign flips in interferences
Fix is part of **ML4EFT2.0**

$$r(\boldsymbol{x}, \boldsymbol{c}) = 1 + \sum_{j=1}^{n_{\text{eft}}} r^{(j)}(\boldsymbol{x}) c_j + \sum_{j=1}^{n_{\text{eft}}} \sum_{k \geq j}^{n_{\text{eft}}} r^{(j,k)}(\boldsymbol{x}) c_j c_k$$
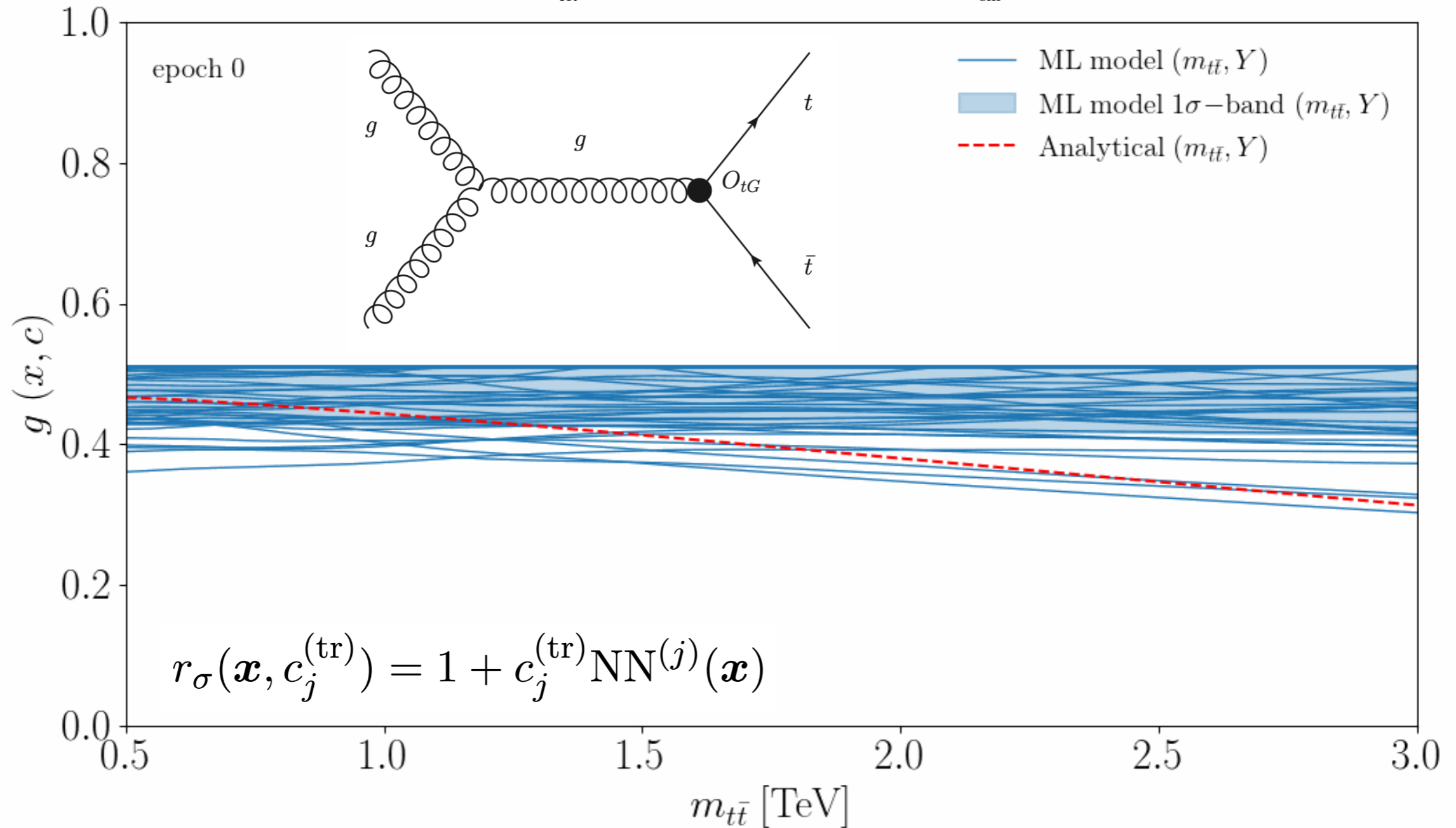
**Example**: to learn a single $r^{(j)}$, generate $\mathscr{D}_{\text{sm}}$ and $\mathscr{D}_{\text{eft}}$ at $c_j$ up to $\mathscr{O}(\Lambda^{-2})$. Then $r(\boldsymbol{x}, \boldsymbol{c}) = 1 + r^{(j)}(\boldsymbol{x}) c_j^{(\text{tr})}$ and training means

$$g(\boldsymbol{x}, c_j^{(\text{tr})}) = \left( 1 + \left[ 1 + c_j^{(\text{tr})} \cdot \text{NN}^{(j)}(\boldsymbol{x}) \right] \right)^{-1} \qquad \text{NN}^{(j)}(\boldsymbol{x}) \rightarrow r^{(j)}(\boldsymbol{x})$$
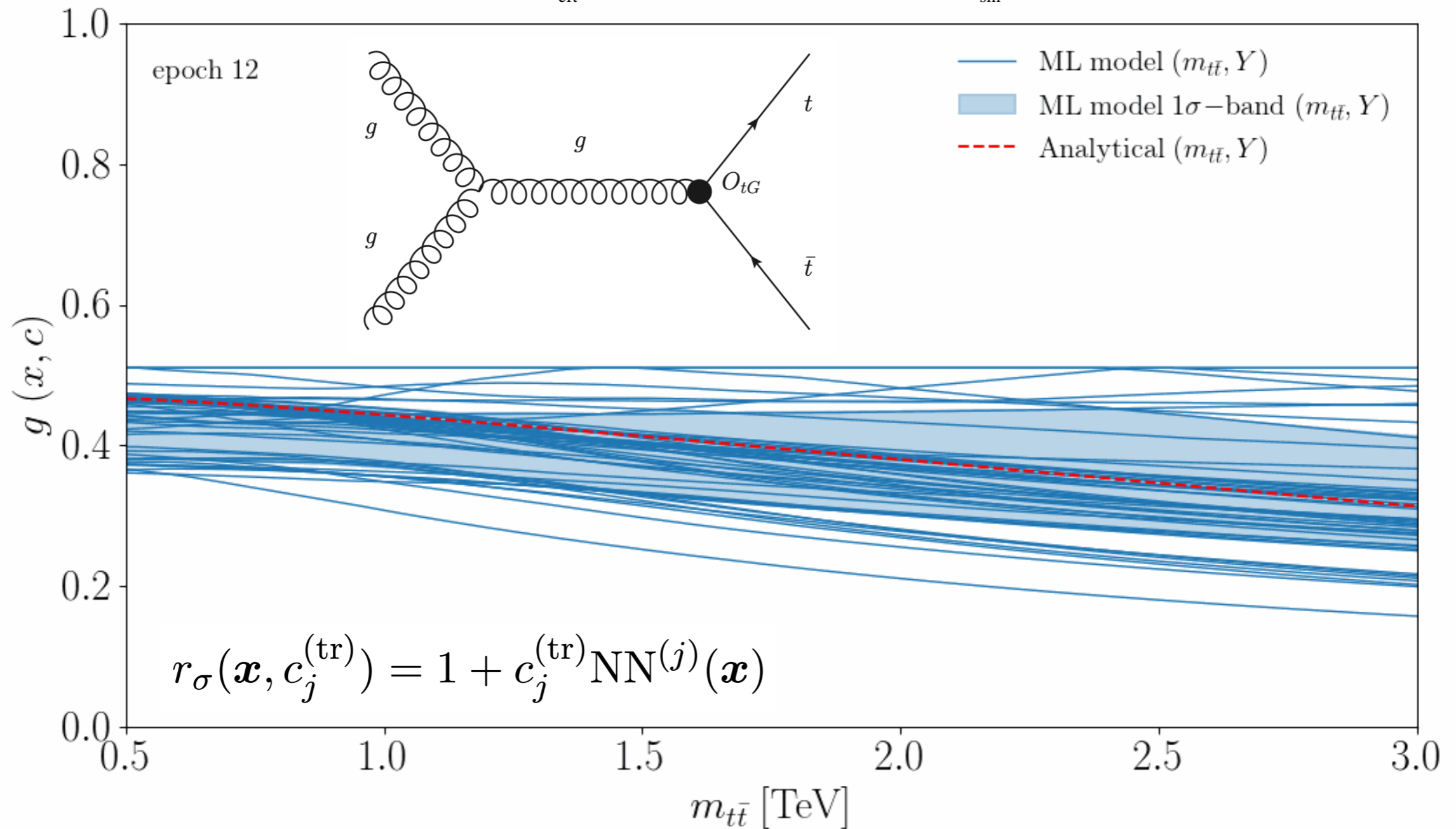
# Toy example

$$L[g(\boldsymbol{x}, \boldsymbol{c})] = -\frac{1}{N} \sum_{e \in \mathscr{D}_{\text{eft}}} w_e \log(1 - g(\boldsymbol{x}_e, \boldsymbol{c})) - \frac{1}{N} \sum_{e \in \mathscr{D}_{\text{sm}}} w_e \log g(\boldsymbol{x}_e, \boldsymbol{c})$$



$$r_\sigma(\boldsymbol{x}, c_j^{(\text{tr})}) = 1 + c_j^{(\text{tr})} \text{NN}^{(j)}(\boldsymbol{x})$$

# Toy example

$$L[g(\boldsymbol{x}, \boldsymbol{c})] = -\frac{1}{N} \sum_{e \in \mathscr{D}_{\mathrm{eft}}} w_e \log(1 - g(\boldsymbol{x}_e, \boldsymbol{c})) - \frac{1}{N} \sum_{e \in \mathscr{D}_{\mathrm{sm}}} w_e \log g(\boldsymbol{x}_e, \boldsymbol{c})$$

# Toy example



$$L[g(\boldsymbol{x}, \boldsymbol{c})] = -\frac{1}{N} \sum_{e \in \mathscr{D}_{\text{eft}}} w_e \log(1 - g(\boldsymbol{x}_e, \boldsymbol{c})) - \frac{1}{N} \sum_{e \in \mathscr{D}_{\text{sm}}} w_e \log g(\boldsymbol{x}_e, \boldsymbol{c})$$

epoch 29

ML model $(m_{t\bar{t}}, Y)$

ML model $1\sigma-$band $(m_{t\bar{t}}, Y)$

Analytical $(m_{t\bar{t}}, Y)$

$$r_\sigma(\boldsymbol{x}, c_j^{(\text{tr})}) = 1 + c_j^{(\text{tr})} \text{NN}^{(j)}(\boldsymbol{x})$$

# Toy example

$$L[g(\boldsymbol{x}, \boldsymbol{c})] = -\frac{1}{N}\sum_{e \in \mathscr{D}_{\text{eft}}} w_e \log(1 - g(\boldsymbol{x}_e, \boldsymbol{c})) - \frac{1}{N}\sum_{e \in \mathscr{D}_{\text{sm}}} w_e \log g(\boldsymbol{x}_e, \boldsymbol{c})$$



$$r_\sigma(\boldsymbol{x}, c_j^{(\text{tr})}) = 1 + c_j^{(\text{tr})} \text{NN}^{(j)}(\boldsymbol{x})$$

# Toy example

$$L[g(\boldsymbol{x}, \boldsymbol{c})] = -\frac{1}{N} \sum_{e \in \mathscr{D}_{\text{eft}}} w_e \log(1 - g(\boldsymbol{x}_e, \boldsymbol{c})) - \frac{1}{N} \sum_{e \in \mathscr{D}_{\text{sm}}} w_e \log g(\boldsymbol{x}_e, \boldsymbol{c})$$



$$r_\sigma(\boldsymbol{x}, c_j^{(\text{tr})}) = 1 + c_j^{(\text{tr})} \text{NN}^{(j)}(\boldsymbol{x})$$

# Toy example

$$L[g(\boldsymbol{x}, \boldsymbol{c})] = -\frac{1}{N} \sum_{e \in \mathscr{D}_{\mathrm{eft}}} w_e \log(1 - g(\boldsymbol{x}_e, \boldsymbol{c})) - \frac{1}{N} \sum_{e \in \mathscr{D}_{\mathrm{sm}}} w_e \log g(\boldsymbol{x}_e, \boldsymbol{c})$$



$$r_\sigma(\boldsymbol{x}, c_j^{(\mathrm{tr})}) = 1 + c_j^{(\mathrm{tr})} \mathrm{NN}^{(j)}(\boldsymbol{x})$$
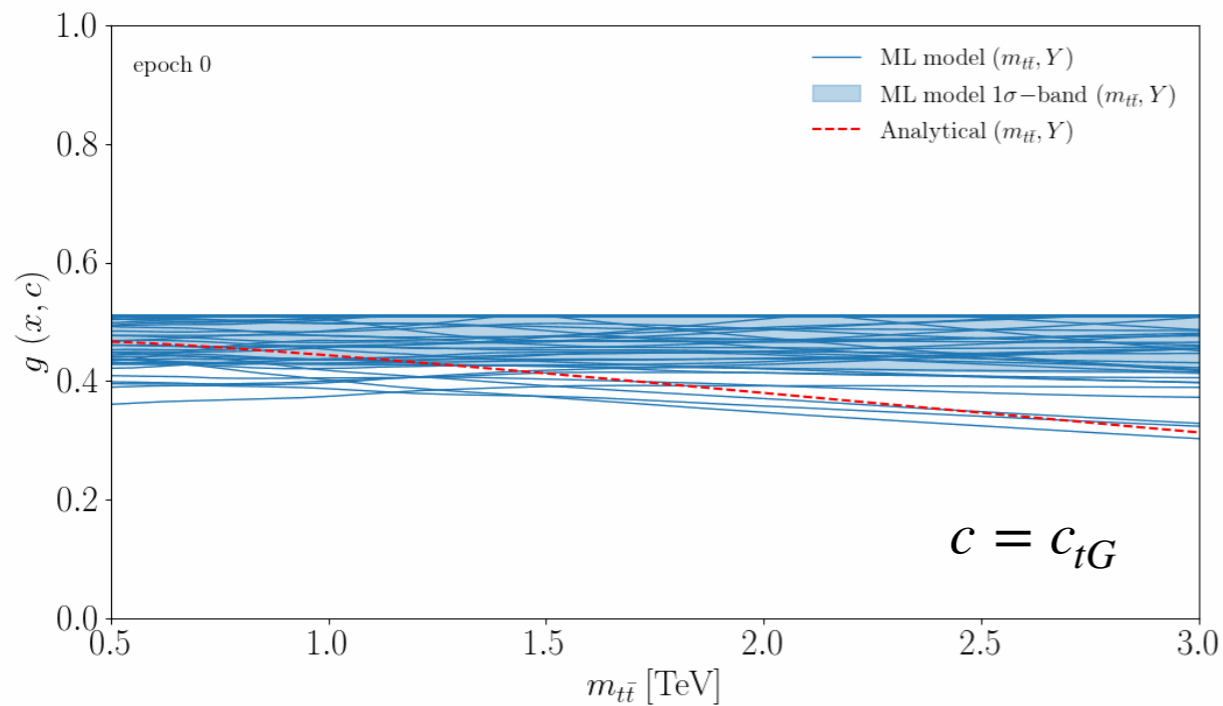
# Toy example

2 kinematic features

# Uncertainty treatment

▸ We only have finite training data and NNs are subject to methodological uncertainties

▸ Propagate uncertainties as well as finite training set effects to the space of models by training multiple replicas

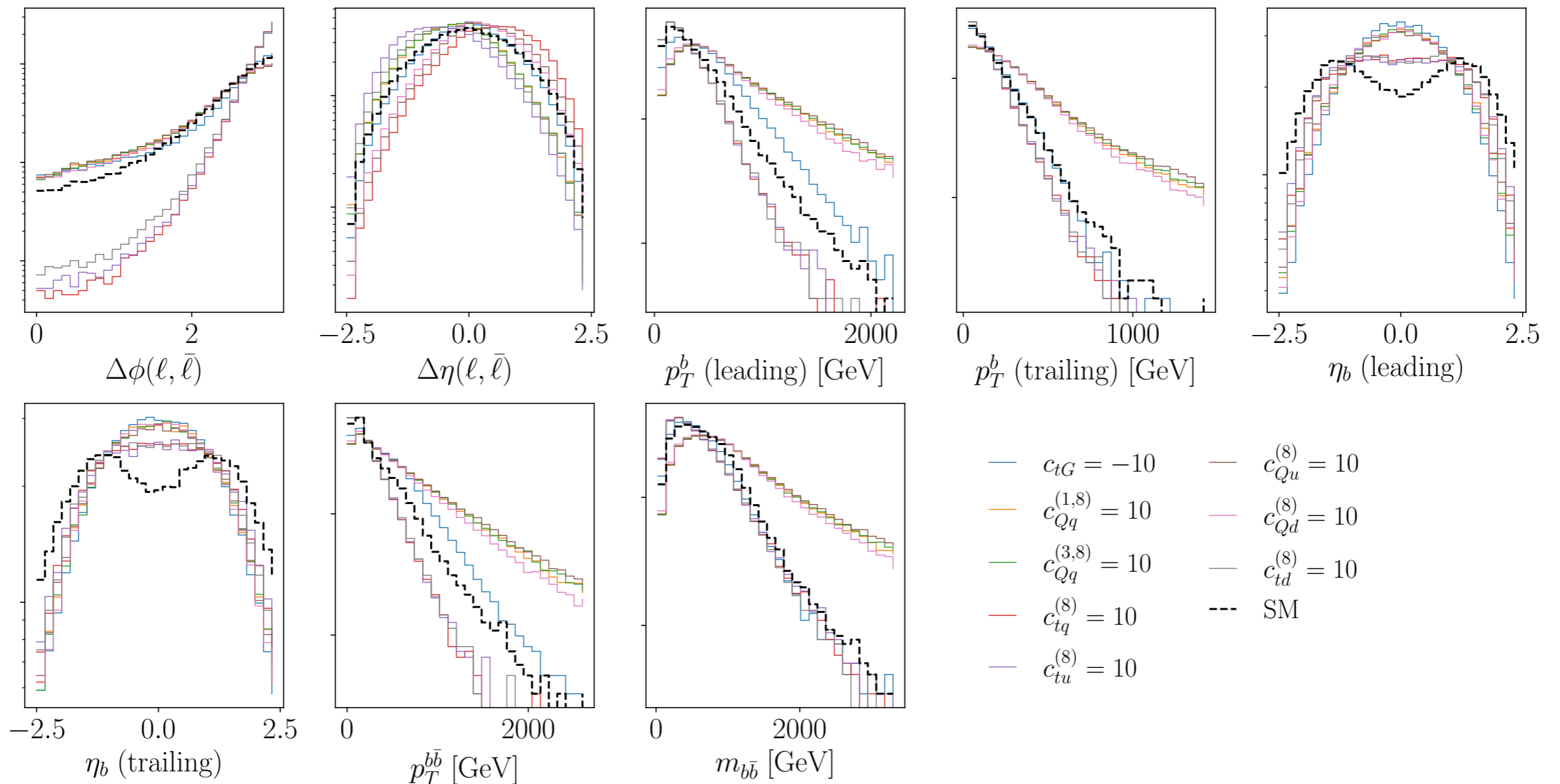$$\hat{r}^{(i)}(\boldsymbol{x}, \boldsymbol{c}) \equiv 1 + \sum_{j=1}^{n_{\text{eft}}} \text{NN}_i^{(j)}(\boldsymbol{x}) c_j + \sum_{j=1}^{n_{\text{eft}}} \sum_{k \geq j}^{n_{\text{eft}}} \text{NN}_i^{(j,k)}(\boldsymbol{x}) c_j c_k, \qquad i = 1, \ldots, N_{\text{rep}}$$



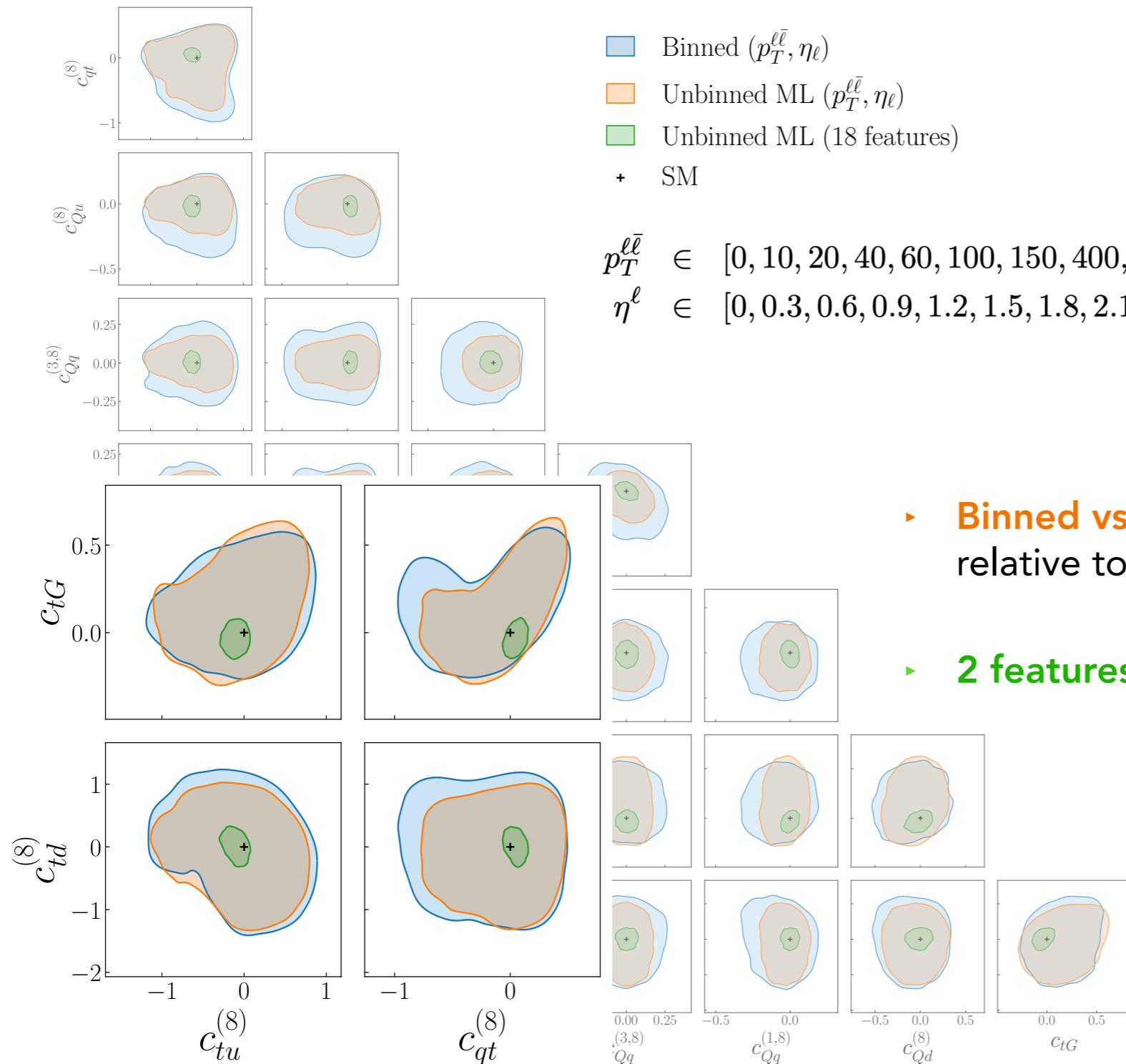| Process | $N_{\text{rep}}$ | $\widetilde{N}_{\text{ev}}$ (per replica) | $N_{\text{nn}}$ | #trainings |
|---|---|---|---|---|
| $pp \rightarrow t\bar{t}$ | 50 | $10^5$ | 4 | 200 |
| $pp \rightarrow t\bar{t} \rightarrow b\bar{b}\ell^+\ell^-\nu_\ell\bar{\nu}_\ell$ | 25 | $10^5$ | 40 | 1000 |

# Let's go multivariate

- $pp \to t\bar{t} \to b\bar{b}\ell^+\ell^-\nu_\ell\bar{\nu}_\ell$ : 18 features, 8 EFT coefficients

- $pp \to hZ \to b\bar{b}\ell^+\ell^-$ : 7 features, 7 EFT coefficients

Marginalised 95 % C.L. intervals, $\mathcal{O}\left(\Lambda^{-4}\right)$ at $\mathcal{L} = 300\ \text{fb}^{-1}$

Binned $(p_T^{\ell\bar{\ell}}, \eta_\ell)$

Unbinned ML $(p_T^{\ell\bar{\ell}}, \eta_\ell)$

Unbinned ML (18 features)

$+$  SM

$$p_T^{\ell\bar{\ell}} \in [0, 10, 20, 40, 60, 100, 150, 400, \infty)\ \text{GeV},$$
$$\eta^\ell \in [0, 0.3, 0.6, 0.9, 1.2, 1.5, 1.8, 2.1, 2.5],$$

▸ **Binned vs unbinned** in $(p_T^{\ell\bar{\ell}}, \eta_\ell)$ small improvement relative to binned setup

▸ **2 features vs 18 features:** big increase in sensitivity

# Unbinned observables in Higgs + Z associated production

Marginalised 95 % C.L. intervals, $\mathcal{O}\left(\Lambda^{-4}\right)$ at $\mathcal{L} = 300 \text{ fb}^{-1}$



Legend:
- $p_T^Z \in [75, 150, 250, 400, \infty)$ [GeV]
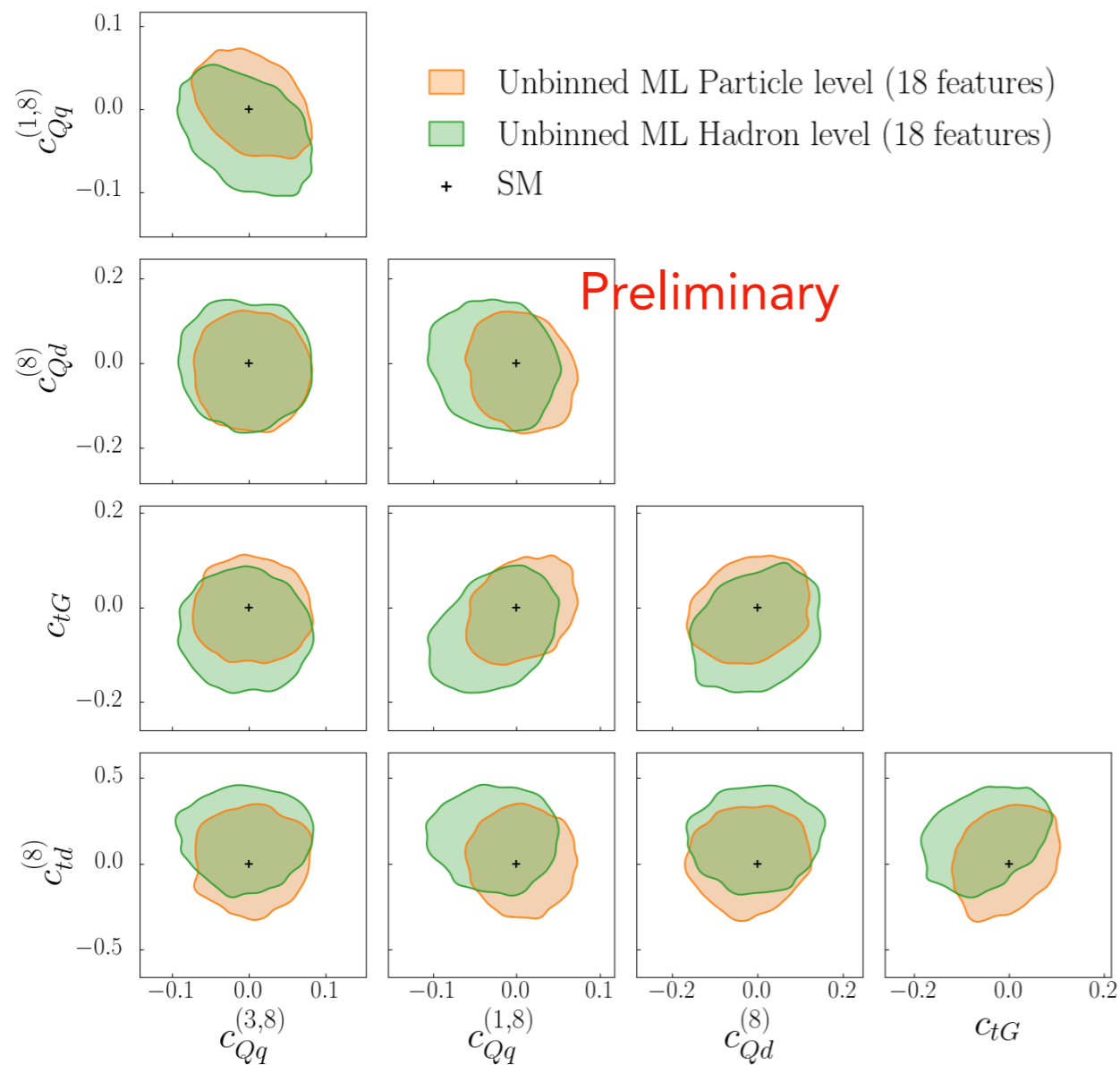- Unbinned ML ($p_T^Z$)
- Unbinned ML (7 features)
- SM

$$pp \to hZ \to b\bar{b}\ell^+\ell^-$$

‣ Unbinned multivariate data is advantageous to constrain the EFT parameter space

‣ Degeneracies get lifted
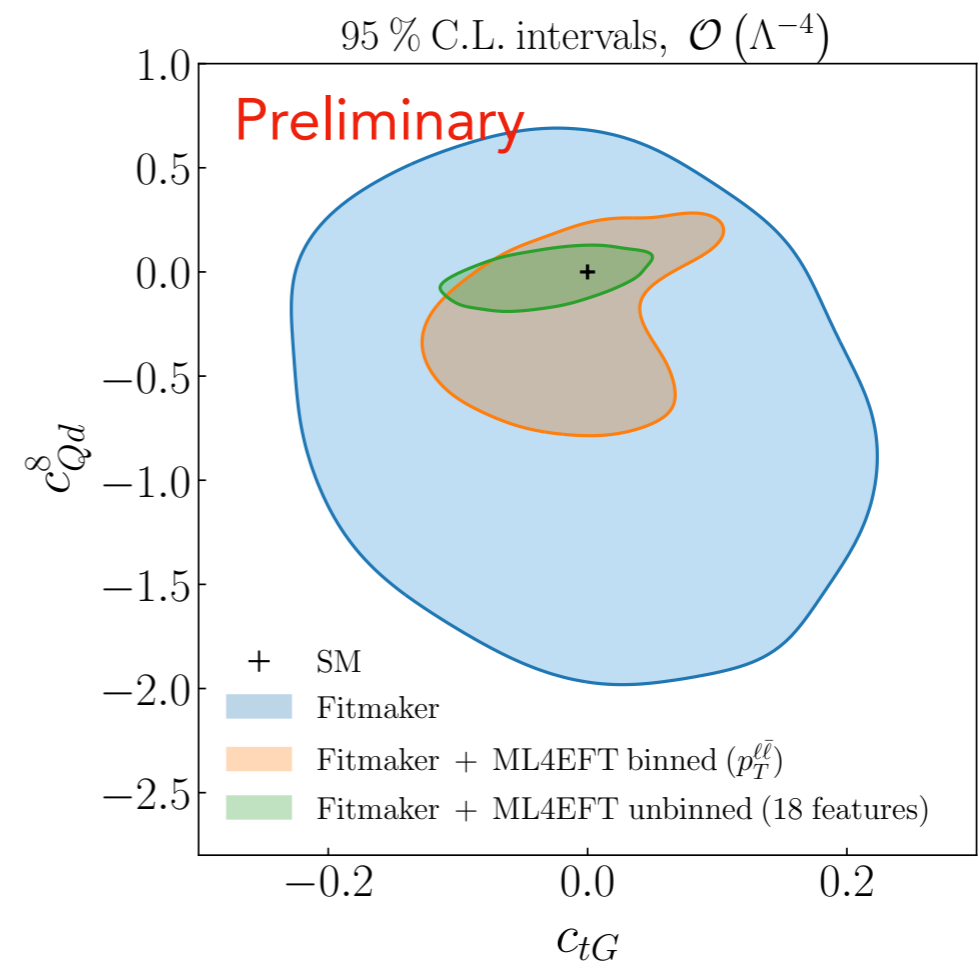
# Ongoing efforts

## 1. Hadronised level

Marginalised 95 % C.L. intervals, $\mathcal{O}\left(\Lambda^{-4}\right)$ at $\mathcal{L} = 300\,\mathrm{fb}^{-1}$



Unbinned ML Particle level (18 features)
Unbinned ML Hadron level (18 features)
+ SM

Preliminary

MSc project by Pim Herbschleb

## 2. Integration into global fits

$$\log\mathcal{L}(c) = \sum_{k=1}^{N_D^{(\mathrm{unbinned})}} \log\mathcal{L}_k^{\mathrm{unbinned}}(c) + \sum_{k=1}^{N_D^{(\mathrm{binned})}} \log\mathcal{L}_k^{\mathrm{binned}}(c)$$

95 % C.L. intervals, $\mathcal{O}\left(\Lambda^{-4}\right)$

Preliminary



+ SM
Fitmaker
Fitmaker + ML4EFT binned $(p_T^{\ell\bar{\ell}})$
Fitmaker + ML4EFT unbinned (18 features)

# Applications of likelihood learning

## Focusses on global EFT fits



Nikhef-2022-015

**Unbinned multivariate observables for global SMEFT analyses from machine learning**

Raquel Gomez Ambrosio,[1,2] Jaco ter Hoeve,[3,4] Maeve Madigan,[5] Juan Rojo,[3,4] and Veronica Sanz[6,7]

[1] Dipartimento di Fisica "G. Occhialini", Universita degli Studi di Milano-Bicocca, and INFN, Sezione di Milano Bicocca, Piazza della Scienza 3, I – 20126 Milano, Italy
[2] Dipartimento di Fisica, Università di Torino, and INFN, Sezione di Torino, Via P. Giuria 1, 10125 Torino, Italy
[3] Department of Physics and Astronomy, VU Amsterdam, 1081HV Amsterdam, The Netherlands
[4] Nikhef Theory Group, Science Park 105, 1098 XG Amsterdam, The Netherlands
[5] DAMTP, University of Cambridge, Wilberforce Road, Cambridge CB3 0WA, UK
[6] Instituto de Física Corpuscular (IFIC), Universidad de Valencia-CSIC, E-46980 Valencia, Spain
[7] Department of Physics and Astronomy, University of Sussex, Brighton BN1 9QH, UK

**Abstract**

Theoretical interpretations of particle physics data, such as the determination of the Wilson coefficients of the Standard Model Effective Field Theory (SMEFT), often involve the inference of multiple parameters from a global dataset. Optimizing such interpretations requires the identification of observables that exhibit the highest possible sensitivity to the underlying theory parameters. In this work we develop a flexible open source framework, ML4EFT, enabling the integration of unbinned multivariate observables into global SMEFT fits. As compared to traditional measurements, such observables enhance the sensitivity to the theory parameters by preventing the information loss incurred when binning in a subset of final-state kinematic variables. Our strategy combines machine learning regression and classification techniques to parameterize high-dimensional likelihood ratios, using the Monte Carlo replica method to estimate and propagate methodological uncertainties. As a proof of concept we construct unbinned multivariate observables for top-quark pair and Higgs+$Z$ production at the LHC, demonstrate their impact on the SMEFT parameter space as compared to binned measurements, and study the improved constraints associated to multivariate inputs. Since the number of neural networks to be trained scales quadratically with the number of parameters and can be fully parallelized, the ML4EFT framework is well-suited to construct unbinned multivariate observables which depend on up to tens of EFT coefficients, as required in global fits.

arXiv:2211.02058v2 [hep-ph] 23 May 2023

1

## Reweighting for more accurate learning



**Boosting likelihood learning with event reweighting**

Siyu Chen[1], Alfredo Glioti[2], Giuliano Panico[3,4], and Andrea Wulzer[5,6]

[1] Institut de Théorie des Phénomenes Physiques, EPFL, Lausanne, Switzerland
[2] Université Paris-Saclay, CNRS, CEA, Institut de Physique Théorique, 91191, Gif-sur-Yvette, France
[3] Dipartimento di Fisica e Astronomia, Università di Firenze, Via G. Sansone 1, 50019 Sesto Fiorentino, Italy
[4] INFN, Sezione di Firenze, Via G. Sansone 1, 50019 Sesto Fiorentino, Italy
[5] Institut de Física d'Altes Energies (IFAE), The Barcelona Institute of Science and Technology (BIST), Campus UAB, 08193 Bellaterra, Barcelona, Spain
[6] ICREA, Institució Catalana de Recerca i Estudis Avançats, Passeig de Lluís Companys 23, 08010 Barcelona, Spain

**Abstract**

Extracting maximal information from experimental data requires access to the likelihood function, which however is never directly available for complex experiments like those performed at high energy colliders. Theoretical predictions are obtained in this context by Monte Carlo events, which do furnish an accurate but abstract and implicit representation of the likelihood. Strategies based on statistical learning are currently being developed to infer the likelihood function explicitly by training a continuous-output classifier on Monte Carlo events. In this paper, we investigate the usage of Monte Carlo events that incorporate the dependence on the parameters of interest by reweighting. This enables more accurate likelihood learning with less training data and a more robust learning scheme that is more suited for automation and extensive deployment. We illustrate these advantages in the context of LHC precision probes of new Effective Field Theory interactions.

arXiv:2308.05704v1 [hep-ph] 10 Aug 2023