







European Research Council

Established by the European Commission

MARIA UBIALI UNIVERSITY OF CAMBRIDGE

STATISTICAL ASPECTS IN THE DETERMINATION OF THE PROTON'S STRUCTURE (PART II)

PHYSTAT SEMINAR SERIES

28TH MAY 2025

<u>OUTLINE</u>

- Physics background: The Standard Model, the LHC, factorisation theorem & parton distributions functions
- Statistical formulation of the inverse problem of determining the proton structure from the LHC data and **the CT solution**
- Part I Statistical formulation of the inverse problem of determining the proton structure from the LHC data and **the NNPDF solution**

PAVEI

MF

 Part II Statistical closure, experimental inconsistencies and simultaneous parameter determination

➡ Part III (Some) open questions and discussion

Part I Statistical formulation of the inverse problem of determining the proton structure from the LHC data and **the NNPDF solution**

A STATISTICAL FORMULATION OF THE PROBLEM

- Data D vector of dimension N_{dat} ~ O(5000), is obtained from many correlated collider measurements:
 (ep) with one proton or nucleus in the initial state and (pp) with two protons in the initial state.
- Experimentalists tell us that results are approximately distributed as a multivariate normal distribution, and provide us with a **mean value** and an experimental **covariance matrix** Σ_{exp} for the data.

$$D \sim \mathcal{N}(D_0, \Sigma_{\exp})$$

Model predictions for the data T depend on a <u>vector of physics parameters</u> c (coupling constants, heavy quark and boson masses...) plus some <u>continuous functions</u> f_i that parametrise the proton nuclear substructure in terms of PDFs. Note that PDFs themselves depend on physics parameters.

$$T^{(\mathrm{ep})}(\{f\}, \vec{c}) = \sum_{i=1}^{n_{\mathrm{partons}}} \int_{x}^{1} dy f_{i}(y; \vec{c}) \sigma_{i}\left(\frac{x}{y}; \vec{c}\right) \equiv \sum_{i=1}^{n_{\mathrm{partons}}} \left[f_{i} \otimes \sigma_{i}\right](x; \vec{c})$$
$$T^{(\mathrm{pp})}(\{f\}, \vec{c}) = \sum_{i,j=1}^{n_{\mathrm{partons}}} \left[f_{i} \otimes f_{j} \otimes \sigma_{ij}\right](x; \vec{c})$$

<u>A STATISTICAL FORMULATION OF THE PROBLEM</u>

• For now let's leave the parameters **c** = **c*** fixed to some optimal value (e. g. PDG value).

$$T(\{f\}, \vec{c} = \vec{c}^*) \equiv T(\{f\})$$

• What theory constraints/indications on the functions **{f}** do we have from the theory?

*Observable predictions are necessarily positive, so $T(\{f\}) \geq 0$

* Parton distributions are **smooth** functions $f \in C^{\infty}[0,1]$.

- * Parton distributions obey the '**Regge' limits**: as $x \to 0$ and as $x \to 1$, the **scaling behaviour** of the PDFs is a **power law**, $f(x) \sim x^{\alpha}$ and $f(x) \sim (1 x)^{\beta}$ respectively, where α, β are unknown.
- * When we consider all flavours of PDFs, there exist **sum rules** which constrain integrals of linear combinations of the PDFs (e.g. $f_u f_{\bar{u}}$ must integrate to 2).

<u>A STATISTICAL FORMULATION OF THE PROBLEM</u>

- Given these constraints and given O(5000) experimental data points D want to determine the set of functions {f} (with i = 1, ..., n_{partons}) and estimate their uncertainty
- Want to find a infinite-dimensional object from a finite number of information.
- Hence, PDF inference is an example of a **non-linear infinite-dimensional** inverse problem.
- Mapping D into f is mathematically ill defined, nobody knows the true f = f*
- Best we can do it to find best **f** given the data **D**.



A STATISTICAL FORMULATION OF THE PROBLEM

- Given these constraints and given O(5000) experimental data points D want to determine the set of functions {f} (with i = 1, ..., n_{partons}) and estimate their uncertainty
- Want to find a infinite-dimensional object from a finite number of information.
- Hence, PDF inference is an example of a **non-linear infinite-dimensional** inverse problem.
- Mapping D into f is mathematically ill defined, nobody knows the true f = f*
- Best we can do it to find best **f** given the data **D**.



Two ways of finding the prior so far:
Explicit parametrization (parametrical modelling)
Function is projected on a vector space of parameters θ
Non-parametrical inference (NNs, functional space sampling,...): use data to also infer probable 'smoothness' of f(x) and thus infer probability function P(f) throughout the space of functions

NNPDF APPROACH

 NNPDF produces regular fits of PDFs, since the NNPDF4.0 global fit in 2021 public code with up-to-date documentation that includes all theory and experimental inputs to reproduce all results and produce new ones <u>https://github.com/NNPDF/nnpdf</u> & <u>https://docs.nnpdf.science/</u>



(1) NN PARAMETRIZATION

- Instead of using fixed parametrisation, choose parametrisation so large that in principle can fit any conceivable function f.
- One option is to use NNs [Forte, Latorre 2002]
- First PDF fit using NN to parametrise full set of PDFs [Ball, Del Debbio, Forte, Guffanti, Piccione, Rojo, MU, 2008]



number of iterations

(1) DEEP NN PARAMETRIZATION

- NNPDF4.0: Single deep Neural Network (763 parameters) [Ball et al, 2021 arXiv:2109.02653]
- Hyper-parameter optimisation via K-fold procedure: scan parameter space and optimise K-fold loss.



(1) DEEP NN PARAMETRIZATION

- NNPDF4.0: Single deep Neural Network (763 parameters) [Ball et al, 2021 arXiv:2109.02653]
- Hyper-parameter optimisation via K-fold procedure: scan parameter space and optimise K-fold loss.
- Each fold reproduces features of the full dataset, loss = average of the non-fitted folds => good generalisation

	Fold 1	
CHORUS σ_{CC}^{ν}	HERA I+II inc NC e^+p 920 GeV	BCDMS p
LHCb Z 940 pb	ATLAS W, Z 7 TeV 2010	CMS Z p_T 8 TeV (p_T^{ll}, y_{ll})
DY E605 σ_{DY}^{p}	CMS Drell-Yan 2D 7 TeV 2011	CMS 3D dijets 8 TeV
ATLAS single- $\bar{t} y$ (normalised)	ATLAS single top R_t 7 TeV	CMS $t\bar{t}$ rapidity $y_{t\bar{t}}$
CMS single top R_t 8 TeV		
	Fold 2	
HERA I+II inc CC e^-p	HERA I+II inc NC e^+p 460 GeV	HERA comb. $\sigma_{b\bar{b}}^{red}$
NMC p	NuTeV σ_c^p	LHCb $Z \rightarrow ee \ 2 \ fb$
CMS W asymmetry 840 pb	ATLAS Z p_T 8 TeV (p_T^{ll}, M_{ll})	D0 $W \rightarrow \mu \nu$ asymmetry
DY E886 σ_{DY}^{p}	ATLAS direct photon 13 TeV	ATLAS dijets 7 TeV, R=0.6
ATLAS single antitop y (normalised)	CMS σ_{tt}^{tot}	CMS single top $\sigma_t + \sigma_{\tilde{t}}$ 7 TeV
	Fold 3	
HERA I+II inc CC e^+p	HERA I+II inc NC e^+p 575 GeV	NMC d/p
NuTeV σ_c^{ν}	LHCb $W, Z \rightarrow \mu$ 7 TeV	LHCb $Z \rightarrow ee$
ATLAS W, Z 7 TeV 2011 Central selection	ATLAS W^+ +jet 8 TeV	ATLAS HM DY 7 TeV
CMS W asymmetry 4.7 fb	DYE 866 $\sigma_{DY}^d / \sigma_{DY}^p$	CDF Z rapidity (new)
ATLAS σ_{tt}^{tot}	ATLAS single top y_t (normalised)	CMS σ_{tt}^{tot} 5 TeV
CMS $t\bar{t}$ double diff. $(m_{t\bar{t}},y_t)$		
	Fold 4	
CHORUS $\sigma_{CC}^{\bar{\nu}}$	HERA I+II inc NC e^+p 820 GeV	LHCb $W, Z \rightarrow \mu 8 \text{ TeV}$
LHCb $Z \rightarrow \mu \mu$	ATLAS W, Z 7 TeV 2011 Fwd	ATLAS W^- +jet 8 TeV
ATLAS low-mass DY 2011	ATLAS Z p_T 8 TeV (p_T^{ll}, y_{ll})	CMS W rapidity 8 TeV
D0 Z rapidity	CMS dijets 7 TeV	ATLAS single top y_t (normalised)
ATLAS single top R_t 13 TeV	CMS single top R_t 13 TeV	



FROM DNN TO THEORY PREDICTIONS



 α , β are pre-processing exponents randomised within a range for each parton combination and iterated until range is stable



The convolution of the PDFs with the theory functions σ is done by computing PDFs on a grid of $x^{(k)}$ points and multiplying the by a precomputed FK table

(2) THE MONTE CARLO SAMPLING FOR ERROR PROPAGATION

- In such a large parameter space Hessian approach not applicable.
- Use Monte Carlo or bootstrap error propagation by importance sampling

[Giele, Kosover 1993] [Forte et al, 2006]

$$\begin{array}{c} D_0 = t + \eta \\ \uparrow & \uparrow & \eta \sim \mathcal{N}(0, \Sigma_{\mathrm{exp}}) \end{array}$$

Vector of Vector of "true", Observational noise experimental values values

For now let's only consider experimental covariance matrix. In the NNPDF approach a theory covariance matrix is also included and added to the experimental covariance matrix

$$\boldsymbol{\epsilon}^{(k)} \sim \mathcal{N}(0, \Sigma_{\mathrm{exp}})$$

$$\mu^{(k)} = D_0 + \epsilon^{(k)} = t + \eta + \epsilon^{(k)}$$

Pseudo-data replicas, k = 1, ..., N_{rep}

(2) THE MONTE CARLO SAMPLING FOR ERROR PROPAGATION

- Visually, fitting pseudo data = throwing random pseudo data points about the experimental data D₀, according to a multivariate normal distributions centred on the experimental data D₀ with experimental covariance matrix Σ_{exp}.
- For each pseudo data compute the optimal point on the theory surface based on training-validation splitting and minimisation stopping procedure and obtain associated parameter values.
- Repeating gives an approximation to the parameter distribution by importance sampling.



Costantini et al, arXiv:2404.10056

(2) THE MONTE CARLO SAMPLING FOR ERROR PROPAGATION



11/30

(2) THE REPLICA DISTRIBUTION

- Using results of the optimisation for each replica (PDF replicas) we can make predictions for observables that depend on PDFs. In prediction space: Gaussian distributions.
- Data uncertainty => PDF replica fluctuations
- Interpolation, extrapolation and functional uncertainties
 => best fit degeneracy
- No correlation between fit quality and position in the Z-H plane => uniform fit quality



NNPDF4.0



Ball et al. 2211.12961

(2) THE REPLICA DISTRIBUTION

- Using results of the optimisation for each replica (PDF replicas) we can make predictions for observables that depend on PDFs. In prediction space: Gaussian distributions.
- Data uncertainty => PDF replica fluctuations
- Interpolation, extrapolation and functional uncertainties
 => best fit degeneracy
- No correlation between fit quality and position in the Z-H plane => uniform fit quality





Fit quality of each replica to the central data D_0 statistically distributed.

Average replica (best fit PDF) not necessarily the one with lowest chi2

Ball et al, 2211.12961

THE ROLE OF DATA AND METHODOLOGY IN A PDF FIT



Ball et al, 2021 arXiv:2109.02653

Shift in parton luminosities mostly due to inclusion of O(500) more data points

• Parton luminosities based on same dataset are consistent with each other but 4.0 methodology displays smaller uncertainty than 3.1 methodology

➡ Part II Statistical closure test, experimental inconsistencies and simultaneous parameter fits

(3) CLOSURE TESTS

- What checks are done on the methodology?
- Assume a given underlying law of Nature: e.g. NNLO predictions for partonic cross section and a given PDF set (for example a random NNPDF replica)
- Generate data with central values given by the "true" law of Nature, and distributed according to experimental covariance matrix.
- Run a fit with NNPDF methodology
- Do statistics on "runs of the universe": generate D₀ N_{fit} times with different random noise drawn from experimental distribution



L2 pseudo-data (to propagate uncertainties
in Monte Carlo fits)
$$\mu^{(k)} = D_0 + \epsilon^{(k)} = t + \eta + \epsilon^{(k)}$$
$$k = 1, ..., N_{\rm rep}$$

(3) STATISTICAL ESTIMATORS

$$u_{*,k} = \operatorname*{argmin}_{u_k} \left[\chi_{\mathrm{val}}^{2(k)} | \operatorname*{argmin}_{u_k} \chi_{\mathrm{tr}}^{2(k)}
ight], \quad k = 1, ..., N_{\mathrm{rep}}$$

Best fit for each data replica in Monte Carlo approach

15/30

 $(\Delta_{\text{PDF}})_{i} = \mathbb{E}_{\epsilon} \left[\mathcal{G}(u_{*,k})_{i} - \mathbb{E}_{\epsilon} \mathcal{G}(u_{*,k})_{i} \right] \longrightarrow \text{PDF uncertainty of prediction i}$ $(C_{\text{PDF}})_{ij} = \frac{N_{\text{reps}}}{N_{\text{reps}} - 1} \mathbb{E}_{\epsilon} \left[\left(\mathcal{G}(u_{*,k})_{i} - \mathbb{E}_{\epsilon} \mathcal{G}(u_{*,k})_{i} \right) \left(\mathcal{G}(u_{*,k})_{j} - \mathbb{E}_{\epsilon} \mathcal{G}(u_{*,k})_{j} \right) \right]$ $\overset{\text{PDF covariance matrix of predictions i and j}}{(\text{PDF induced correlation among observables})}$

 $(\bar{\mathcal{G}})_i = \mathbb{E}_{\epsilon} \left[\mathcal{G}(u_{*,k})_i \right] \longrightarrow \text{Central value of prediction i}$

• Key estimate is the normalised <u>bias</u>: measure mean square deviation of predictions from the "truth" in units of predicted standard deviation

$$B^{(l)}(C_{\rm PDF}) = \frac{1}{N_{\rm data}} \sum_{i,j=1}^{N_{\rm data}} \left(\mathbb{E}_{\epsilon} \mathcal{G}(u_{*,k}^{(l)})_i - f_i \right) (\overline{C}_{\rm PDF})_{ij}^{-1} \left(\mathbb{E}_{\epsilon} \mathcal{G}(u_{*,k}^{(l)})_j - f_j \right)$$

 $R_b = \sqrt{\mathbb{E}_{\eta} B^{(l)}(\overline{C}_{\rm PDF})}$

Bias averaged over N_{fit} L1 data, i.e. N_{fit} runs of the universe, I = 1, ..., N_{fit}



 $R_b \sim 1$ faithfully estimated PDF uncertainties $R_b < 1$ overestimated PDF uncertainties $R_b > 1$ underestimated PDF uncertainties

(3) CONSISTENT CLOSURE TEST

NNPDF4.0 global fit: distributions of Δi



Barontini et al, 2503.17447

• Normalised <u>bias</u> in the basis of the eigenvectors of the PDF covariance matrix

$$B^{(l)}(\overline{C}_{\rm PDF}) = \frac{1}{N_{\rm data}} \sum_{i=1}^{N_{\rm data}} \frac{(\Delta_i^{(l)})^2}{(\sigma_{\rm PDF}_i)^2}$$

• With Δi being the projection of the Bias along each normalised eigenvector of C_{PDF}

$$\Delta_i^{(l)} = \sum_{j=1}^{N_{ ext{data}}} \left(\mathbb{E}_\epsilon \mathcal{G}(u_{*,k}^{(l)})_j - f_j
ight) v_j^{(i)}$$

$$B^{(l)}(C_{\rm PDF}) = \frac{1}{N_{\rm data}} \sum_{i,j=1}^{N_{\rm data}} \left(\mathbb{E}_{\epsilon} \mathcal{G}(u_{*,k}^{(l)})_i - f_i \right) (\overline{C}_{\rm PDF})_{ij}^{-1} \left(\mathbb{E}_{\epsilon} \mathcal{G}(u_{*,k}^{(l)})_j - f_j \right)$$

(4) MODELLING EXPERIMENTAL INCONSISTENCIES

- How to model an inconsistency of experimental origin = some underestimated experimental systematics?
- Generate L₁ data with "true" experimental covariance matrix

 $(C_{\exp})_{ij} = \delta_{ij}\sigma_i^{(\text{uncorr})}\sigma_j^{(\text{uncorr})} + \sum_{k=1}^{N_{\text{corr}}}\sigma_{i,k}^{(\text{corr})}\sigma_{j,k}^{(\text{corr})}$

• Fit the data using the rescaled experimental covariance matrix (both in pseudo data generation and in the loss function)

$$(C_{\exp}^{\lambda})_{ij} = \delta_{ij}\sigma_i^{(\text{uncorr})}\sigma_j^{(\text{uncorr})} + \sum_{k=1}^{N_{\text{corr}}} \lambda_{i,k}\sigma_{i,k}^{(\text{corr})}\lambda_{j,k}\sigma_{j,k}^{(\text{corr})}$$



(4) BULK INCONSISTENCY

- DIS only fit, in-sample HERA NC data are inconsistent.
- 860 out of 2576 inconsistent datapoints, with all systematic uncertainties underestimated by the same factor λ.





- Highly non linear behaviour: for λ ≥ 0.4 despite inconsistency, PDF uncertainties remain faithful, but then sharply rises.
- In-sample and out-of-sample datasets behave in a similar way => NN model effective at generalising.

(4) BULK INCONSISTENCY

- DIS only fit, in-sample HERA NC data are inconsistent.
- 860 out of 2576 inconsistent datapoints, with all systematic uncertainties underestimated by the same factor λ.





 $\lambda \ge 0.4$: Model corrects for underestimated uncertainty in the inconsistent data, PDF uncertainties do not decrease despite the reduced data uncertainty. λ < 0.4: PDF uncertainty shrinks
=> underestimated PDF
uncertainties, largest shifts in the
gluon and quark singlet
combinations.

- Highly non linear behaviour: for λ ≥ 0.4 despite the inconsistency, PDF uncertainties remain faithful, but then sharply rises.
- In-sample and out-of-sample datasets behave in a similar way => NN model effective at generalising.



- So far learn that NN model cures experimental inconsistencies until they are not too big
- If strong inconsistencies model fails, as either PDF uncertainties shrink or PDFs shift far from underlying law.
- In real life no access to the "truth", normalised bias cannot be computed and only a single run of the Universe.
- How to spot inconsistencies in a actual PDF fit?

Large $R_b \rightarrow$ Large χ^2 of the inconsistent dataset and of consistent datasets correlated to it by the PDFs



Distribution of n_{σ} values across N_{fits} in DIS case with maximal inconsistency ($\lambda = 0$)

Large $R_b \rightarrow$ Large χ^2 of the inconsistent dataset and of consistent datasets correlated to it by the PDFs



Distribution of n_{σ} values across N_{fits} in DIS case with maximal inconsistency ($\lambda = 0$)

A single criterion S_1 given by n_σ threshold Z not enough to minimise false positive and false negative

$$n_{\sigma}^{(i)} = \frac{\chi_i^2 - N_{\text{data}}^{(i)}}{\sqrt{2/N_{\text{data}}^{(i)}}}$$

$$\chi_{\text{weighted}}^{2(i)} = \frac{1}{N_{\text{data}} - N_{\text{data}}^{(j)}} \sum_{j=1, j \neq i}^{N_{\text{exp}}} N_{\text{data}}^{(j)} \chi_j^2 + w^{(i)} \chi_i^2$$

$$w^{(i)} = N_{\text{data}} / N_{\text{data}}^{(i)}$$

Extra criterion: what happens if a large weight is given to inconsistent/consistent dataset? Expect that the inconsistent dataset, if given extra weight in the fit, will either fail to improve or will improve but spoil χ^2 of other consistent datasets, while the consistent dataset will not.

$$\mathbf{S_1} \quad \mu_i > Z \qquad \mathbf{S_2} \quad n_{\sigma}^{\text{weighted},(i)} > Z \qquad \mathbf{S_3} \quad n_{\sigma}^{\text{weighted},(j)} - n_{\sigma}^{(j)} > Z \qquad \forall j \neq i$$

(4) OPTIMAL INCONSISTENCY DETECTION

 $C_1 {:} \ condition \ S_1 \ satisfied$

C₂: condition S₁ satisfied, and in a weighted fit either S₂ or S₃ are satisfied (NNPDF4.0 criterion) C₃: in weighted fit either S₂ or S₃ are satisfied



(4) OPTIMAL INCONSISTENCY DETECTION

 $C_1: \text{condition } S_1 \text{ satisfied}$

C₂: condition S₁ satisfied, and in a weighted fit either S₂ or S₃ are satisfied (NNPDF4.0 criterion) C₃: in weighted fit either S₂ or S₃ are satisfied



In this case for $Z \approx 0.5$, 90% probability of NOT flagging a consistent dataset as inconsistent and 95% probability of flagging an inconsistent dataset as inconsistent.

(5) EXTRACTING PARAMETERS FROM DATA

$$\chi^{2} = \sum_{a,b=1}^{N_{\text{dat}}} (T_{a}(\lbrace f \rbrace, \vec{c}) - D_{a}) \Sigma_{ab}^{-1} (T_{b}(\lbrace f \rbrace, \vec{c}) - D_{b})$$

$$T^{(\text{ep})}(\lbrace f \rbrace, \vec{c}) = \sum_{i=1}^{n_{\text{partons}}} \int_{x}^{1} dy f_{i}(y; \vec{c}) \sigma_{i} \left(\frac{x}{y}; \vec{c}\right) \equiv \sum_{i=1}^{n_{\text{partons}}} [f_{i} \otimes \sigma_{i}] (x; \vec{c})$$

$$SM \text{ parameters } (\alpha S(Mz), Mw, \theta w \dots), \text{ or BSM parameters } (SMEFTWCs \dots)$$

$$T(\lbrace f \rbrace, c) = T^{\text{SM}}(\lbrace f \rbrace, c = 0)(1 + k_{\text{lin}} \cdot c + k_{\text{quad}} \cdot c^{2})$$

✓ In a PDF fit typically

$$T(\{f\}) = T^{SM}(\{f\}, c = 0)$$

✓ In a fit of SMEFT Wilson Coefficients

$$T(c) = T^{SM}(\{f\} = \{f\}^*)(1 + k_{lin} \cdot c + k_{quad} \cdot c^2)$$

(5) SIMULTANEOUS EXTRACTION OF PARAMETERS AND PDF

$$\chi^{2} = \sum_{a,b=1}^{N_{\text{dat}}} (T_{a}(\lbrace f \rbrace, \vec{c}) - D_{a}) \Sigma_{ab}^{-1} (T_{b}(\lbrace f \rbrace, \vec{c}) - D_{b})$$

$$T^{(\text{ep})}(\lbrace f \rbrace, \vec{c}) = \sum_{i=1}^{n_{\text{partons}}} \int_{x}^{1} dy f_{i}(y; \vec{c}) \sigma_{i} \left(\frac{x}{y}; \vec{c}\right) \equiv \sum_{i=1}^{n_{\text{partons}}} [f_{i} \otimes \sigma_{i}] (x; \vec{c})$$

$$\text{SM parameters (} \alpha_{\text{S}}(\text{Mz}), \text{Mw, } \theta_{\text{W}} \dots), \text{ or BSM parameters (SMEFT WCs } \dots)$$

$$T(\lbrace f \rbrace, c) = T^{\text{SM}}(\lbrace f \rbrace, c = 0)(1 + k_{\text{lin}} \cdot c + k_{\text{quad}} \cdot c^{2})$$

✓ But PDFs and parameters are correlated and PDFs can absorb shifts in the parameters (absorb new physics) [E. Hammou, et al 2307.10370]: can we do simultaneous fits of PDFs and parameters?

$$\chi^{2} = \sum_{a,b=1}^{N_{dat}} (T_{a}(\{f\},\vec{c}) - D_{a}) \Sigma_{ab}^{-1} (T_{b}(\{f\},\vec{c}) - D_{b})$$

(5) THE SIMUNET SOLUTION

Hidden

PDF

Convolution

SM

- The idea: take a PDF fit based on NNPDF4.0 methodology and make dependence of observables on physics parameters {c_i} explicit before computing the loss function (e.g. adding SMEFT corrections, or expanding observables in terms of SM precision parameters)
- Perform minimisation of loss function over $\hat{\theta} = \theta \bigcup \{c_i\}$ by adding new layer to the deep neural network used in NNPDF4.0



SMEFT

S. Iranipour, MU - arXiv: 2201.07240

Input

Hidden

(5) EXAMPLE: DRELL-YAN DATA @HL-LHC

S. Iranipour, MU - arXiv: 2201.07240

	SM PDFs	SMEFT PDFs
$W\times 10^5~(68\%~{\rm CL})$	$\left[-1.1, 0.5 ight]$	[-2.4, 1.5]
$W\times 10^5~(95\%~{\rm CL})$	$\left[-2.0, 1.4 ight]$	$\left[-4.3,3.4 ight]$
$Y \times 10^5 (68\% \text{ CL})$	[-0.4, 5.2]	[0.6, 8.0]
$Y\times 10^5~(95\%~{\rm CL})$	$\left[-3.2,8.1\right]$	[-3.1, 11.7]

x 2.3 broadening of bounds for W x 1.3 broadening of bounds for Y

✓ Simultaneous analysis of PDFs and W&Y SMEFT coefficient using simuNET method shows that at HL-LHC the effect of interplay becomes important as WCs bounds broaden and PDF uncertainties change significantly once SMEFT effects allowed in theory predictions entering PDF fit





(5) EXAMPLE: DRELL-YAN DATA @HL-LHC

S. Iranipour, MU - arXiv: 2201.07240



✓ Simultaneous analysis of PDFs and W&Y SMEFT coefficient using simuNET method shows that at HL-LHC the effect of interplay becomes important as WCs bounds broaden and PDF uncertainties change significantly once SMEFT effects allowed in theory predictions entering PDF fit

✓ In this study the dependence on c was **linear**

29/30

(5) ISSUE OF MONTE CARLO METHOD WITH QUADRATIC DEPENDENCE



- In the quadratic SMEFT fit observed disagreement between MC method and Bayesian sampling method. Very different posterior (hence different CLs)
- Mathematical analysis of the problem in [Costantini, Madigan, Mantani, Moore arXiv 2404.10056]
- Towards a general Bayesian methodology for simultaneous fits

[Colibri: Costantini, Mantani, Moore, Sanchez, MU - in progress]



CONCLUSIONS AND OUTLOOK

- In an era of precision, need careful assessment of PDF uncertainties and rigorous claims about faithfulness of PDF uncertainties. Key is to strive for deeper understanding.
- Some random pressing questions from me:
 - In Bayesian terms, do we really have control of the prior probability that we assume in our PDF fits and on how the outcome depends on it?
 - Can we have a truly Bayesian determination of PDFs?
 - And a Bayesian determination of PDFs and SM/BSM parameters?
 - What can be used as a complete basis of the PDF space? Dimensionality is a problem for sampling?
 - Hyper-parameter determination using K-loss: should one use a single configuration or several?
 - Closure tests can be used as a controlled test set where to test faithfulness of PDF uncertainties: is the analysis of the normalised bias enough? Is the inconsistency modelled in a sound way? How to determine a set of maximally consistent datasets?
 - At LHC we have huge amount of data and, while lots of progress has been done in determining individual quantities from these (PDFs, alphaS, ...) we should strive for simultaneous fits, thus including the correlations between parameters and PDFs. How can this be done in a clever and robust way?