

# Determining parton distribution functions accurately and precisely with machine learning

Digital Twins for Nuclear and Particle physics - NPTwins 2024

Emanuele R. Nocera

Università degli Studi di Torino and INFN, Torino

18 December 2024



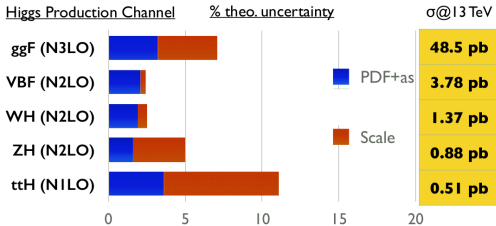
UNIVERSITÀ  
DI TORINO

# Parton Distribution Functions at the LHC

$$\sigma(Q^2, \tau, \mathbf{k}) = \sum_{ij} \int_{\tau}^1 \frac{dz}{z} \mathcal{L}_{ij}(z, Q^2) \hat{\sigma}_{ij} \left( \frac{\tau}{z}, \alpha_s(Q^2), \mathbf{k} \right) \quad \mathcal{L}_{ij}(z, Q^2) = (f_i^{h1} \otimes f_j^{h2})(z, Q^2)$$

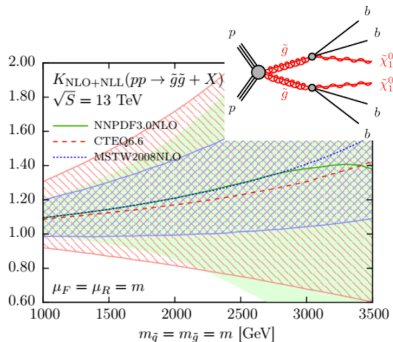
PDF uncertainty is often the dominant source of uncertainty in LHC cross sections

## Precision



Unc. [MeV]	Total	Stat.	Syst.	PDF	$A_i$	Backg.	EW	$e$	$\mu$	$u_T$	Lumi	$\Gamma_W$	PS
$p_T^c$	16.2	11.1	11.8	4.9	3.5	1.7	5.6	5.9	5.4	0.9	1.1	0.1	1.5
$m_T$	24.4	11.4	21.6	11.7	4.7	4.1	4.9	6.7	6.0	11.4	2.5	0.2	7.0
Combined	15.9	9.8	12.5	5.7	3.7	2.0	5.4	6.0	5.4	2.3	1.3	0.1	2.3

## Discovery



[CERN Yellow Report 2016; arXiv:2403.15085]

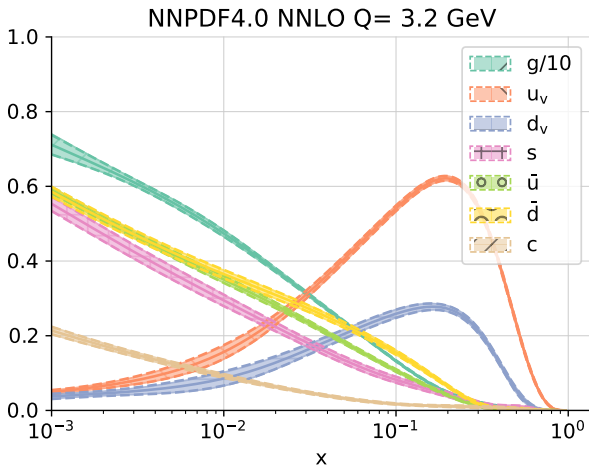
[EPJC 76 (2016) 53]

# Parton Distribution Functions

PDFs express the likelihood of a quark or gluon (partons) to enter a collision

That is,  $x \times$ PDFs are momentum fraction distributions for each parton

Dependence on  $x$  is non-perturbative (fit); dependence on  $Q^2$  is perturbative (DGLAP)



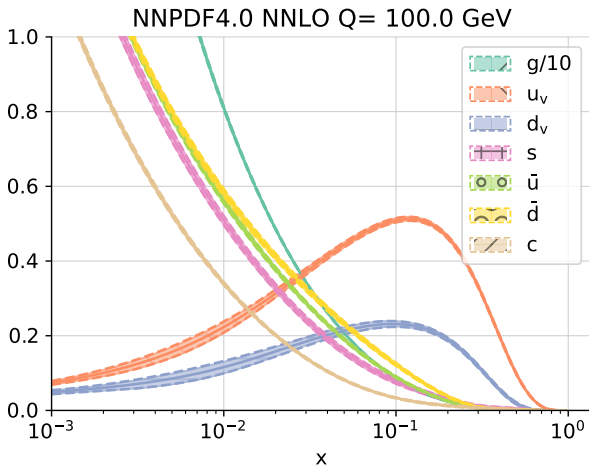
[Plot from the PDG Review of Particle Physics]

# Parton Distribution Functions

PDFs express the likelihood of a quark or gluon (partons) to enter a collision

That is,  $x \times$ PDFs are momentum fraction distributions for each parton

Dependence on  $x$  is non-perturbative (fit); dependence on  $Q^2$  is perturbative (DGLAP)



[Plot from the PDG Review of Particle Physics]

# PDF determination in statistical language

## Inverse problem

Given a set of data  $D$ , determine  $p(f|D)$  in the space of functions  $f : [0, 1] \rightarrow \mathbb{R}$ .

## Solution: parametric regression

Approximate  $p(f|D)$  with its projection in the space of parameters  $p(\boldsymbol{\theta}|D)$

$$x f_i(x, Q_0^2) = A_{f_i} x^{a_{f_i}} (1-x)^{b_{f_i}} \mathcal{F}(x, \{c_{f_i}\})$$

Determine  $p(\boldsymbol{\theta}|D) \propto p(D|\boldsymbol{\theta})p(\boldsymbol{\theta})$  as MAP  $\boldsymbol{\theta}^* = \arg \max_{\boldsymbol{\theta}} p(\boldsymbol{\theta}|D)$

$$\chi^2 = \sum_{i,j}^{N_{\text{dat}}} [T_i[\boldsymbol{\theta}] - D_i] (\text{cov}^{-1})_{ij} [T_j[\boldsymbol{\theta}] - D_j]$$

Use a prescription to compute expectation values and uncertainties of observables

$$E[\mathcal{O}] = \int \mathcal{D}f \mathcal{P}(f|D) \mathcal{O}(f) \quad V[\mathcal{O}] = \int \mathcal{D}f \mathcal{P}(f|D) [\mathcal{O}(f) - E[\mathcal{O}]]^2$$

Monte Carlo:  $\mathcal{P}(f|D) \rightarrow \{f_k\}$

Maximum likelihood:  $\mathcal{P}(f|D) \rightarrow f_0$

$$E[\mathcal{O}] \approx \frac{1}{N} \sum_k \mathcal{O}(f_k)$$

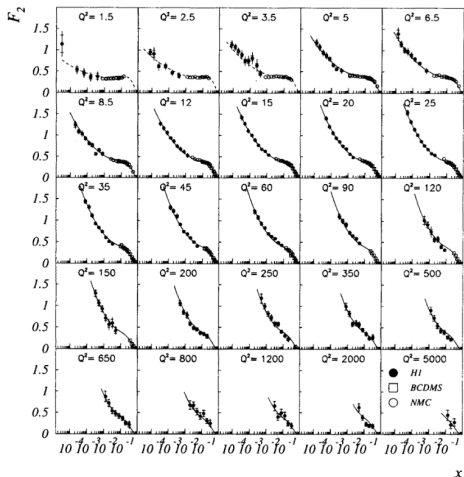
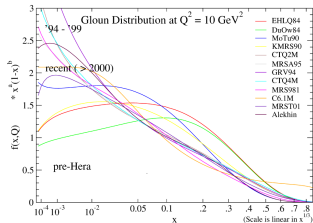
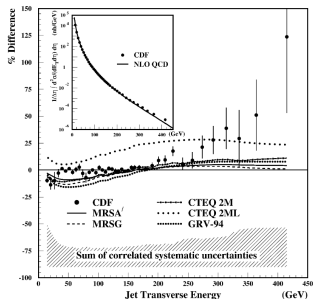
$$E[\mathcal{O}] \approx \mathcal{O}(f_0)$$

$$V[\mathcal{O}] \approx \frac{1}{N} \sum_k [\mathcal{O}(f_k) - E[\mathcal{O}]]^2$$

$$V[\mathcal{O}] \approx \text{Hessian}, \Delta\chi^2 \text{ envelope}, \dots$$

Interplay between DATA, THEORY, and METHODOLOGY

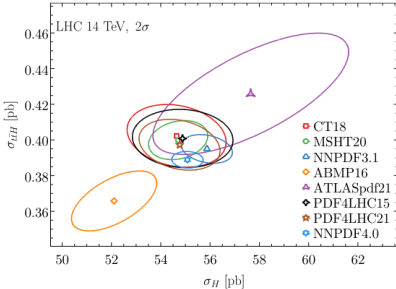
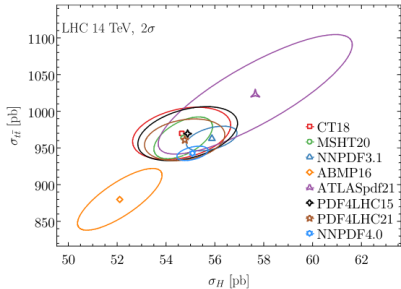
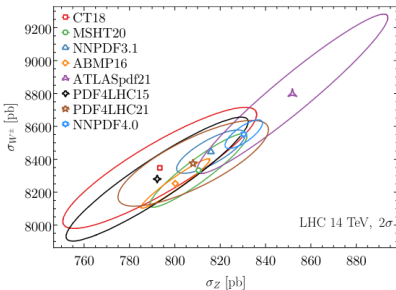
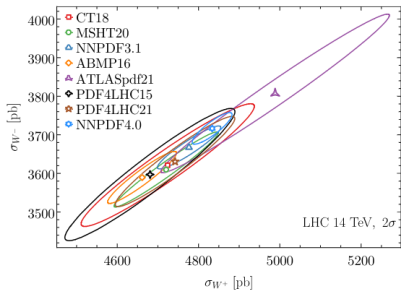
# Why is the methodology important?



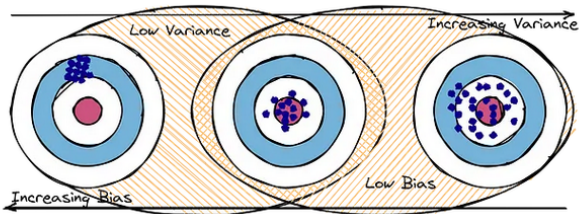
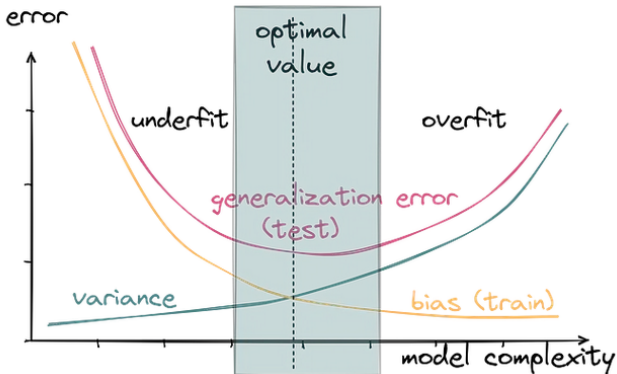
circa 1995: small- $x$  rise of HERA  $F_2^D$  and CDF jet discrepancy

The methodology is crucial if we aim at percent-level accurate PDFs

# Making predictions with PDFs



# Accuracy vs precision or bias vs variance





# A lot of progress in the last (two) years

ANL-186400

## A Markov Chain Monte Carlo determination of Proton PDF uncertainties at NNLO

Learning PDFs through Interpretable Latent Representations in Mellin Space

Brandon Kristen and T. J. Hobbs

High Energy Physics Division, Argonne National Laboratory, Lemont, IL 60439

(Dated: June 21, 2024)

Peter Risse<sup>1,2,\*</sup>, Nasim Derakhshanian<sup>1,2,†</sup>, Tomas Ježo<sup>1,2,‡</sup>, Karol Kovarik<sup>1,2,§</sup> and Aleksander Kusina<sup>1,2,¶</sup>

<sup>1</sup>Institut für Theoretische Physik, Universität Münster, Wilhelm-Klemm-Straße 9, D-48149 Münster, Germany  
<sup>2</sup>Institute of Nuclear Physics Polish Academy of Sciences, PL-31342 Kraków, Poland  
E-mail: [risse.ph.uni-muenster.de](mailto:risse.ph.uni-muenster.de)

Err. Phys. J. C 2024: 84:716  
<https://doi.org/10.1140/epjc/s10052-024-13100-1>

Regular Article - Theoretical Physics

THE EUROPEAN  
PHYSICAL JOURNAL C



MSUHEP-24-002

## Bayesian inference with Gaussian processes for the determination of parton distribution functions

Alessandro Caidolo<sup>1</sup>, Luigi Del Debbio<sup>2</sup>, Tommaso Giani<sup>1,3,4,5</sup>, Giacomo Petrillo<sup>5</sup>

<sup>1</sup>Theoretical Physics Department, CERN, 1211 Geneva 23, Switzerland  
<sup>2</sup>Higgs Centre for Theoretical Physics, School of Physics and Astronomy, Peter Guthrie Tait Road, Edinburgh EH9 3FD, UK  
<sup>3</sup>Department of Physics and Astronomy, Vrije Universiteit, 1081 HV Amsterdam, The Netherlands  
<sup>4</sup>Nikhef Theory Group, Science Park 105, 1098 XG Amsterdam, The Netherlands  
<sup>5</sup>Dipartimento di Statistica, Informatica, Applicazioni "Giuseppe Peano" (DISIA), Università di Firenze, Viale Morgagni 59, 50134 Firenze, Italy

Received: 7 May 2024 / Accepted: 1 July 2024 / Published online: 22 July 2024  
© The Author(s) 2024



PUBLISHED FOR SISSA BY SPRINGER

RECEIVED: April 06, 2024  
REVISED: September 18, 2024  
ACCEPTED: November 5, 2024  
PUBLISHED: December 10, 2024



PUBLISHED FOR SISSA BY SPRINGER

RECEIVED: August 5, 2024  
ACCEPTED: October 11, 2024  
PUBLISHED: November 5, 2024

## A critical study of the Monte Carlo replica method

Mark N. Costantini<sup>1,2,\*</sup>, Maeva Madigan<sup>1,2,†</sup>, Luca Mantani<sup>1,2,‡</sup> and James M. Moore<sup>1,2,§</sup>

<sup>1</sup>DAMTP, University of Cambridge, Wilberforce Road, Cambridge, CB3 0WA, U.K.  
<sup>2</sup>Institut für Theoretische Physik, Universität Heidelberg, Philosophenweg 16, D-69120, Heidelberg, Germany  
E-mail: [mnc330@cam.ac.uk](mailto:mnc330@cam.ac.uk), [madigan@thphys.uni-heidelberg.de](mailto:madigan@thphys.uni-heidelberg.de), [luca.mantani@maths.can.ac.uk](mailto:luca.mantani@maths.can.ac.uk), [jm2320@cam.ac.uk](mailto:jm2320@cam.ac.uk)

## Explainable AI classification for parton density theory

Brandon Kristen<sup>1,\*</sup>, Jonathan Gomprecht<sup>1,2,†</sup> and T.J. Hobbs<sup>1,2,‡</sup>

<sup>1</sup>High Energy Physics Division, Argonne National Laboratory, Lemont, IL 60439, U.S.A.  
<sup>2</sup>Department of Physics, University of Arizona, Tucson, AZ 85721, U.S.A.  
E-mail: [bkristen@anl.gov](mailto:bkristen@anl.gov), [jgomprecht@arizona.edu](mailto:jgomprecht@arizona.edu), [tin@anl.gov](mailto:tin@anl.gov)

[EPJ C84 (2024) 716; JHEP 11 (2024) 007; JHEP 12 (2024) 064; arXiv:2312.02278; arXiv:2406.01664; arXiv:2407.12377]

# NNPDF4.0: a machine learning methodology

uncertainty representation

Monte Carlo sampling of experimental uncertainties

**what is the statistical meaning of uncertainties?**

parametrisation

neural network(s)

**is there a bias due to the parametrisation?**

optimisation

(adaptive) gradient descent

**is the parameter space explored efficiently?**

delivery

GAN enhancement and compression

**can the number of replicas be reduced?**

uncertainty characterisation and validation

closure tests (what happens if I know in advance the underlying law that I am fitting?)

**are interpolation and extrapolation uncertainties statistically faithful?**

[EPJC 82 (2022) 428]

The NNPDF code is public, see <https://github.com/NNPDF> [EPJ C81 (2021) 958]

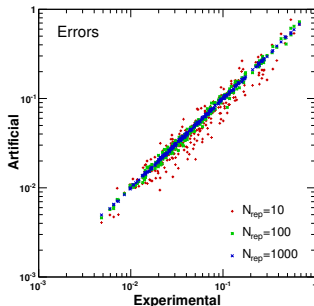
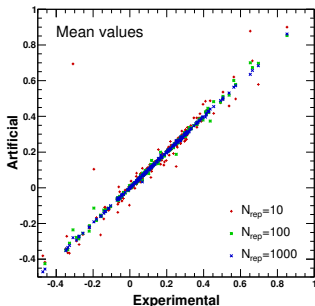
# Uncertainty representation

Generate Monte Carlo replicas and perform a fit to each replica

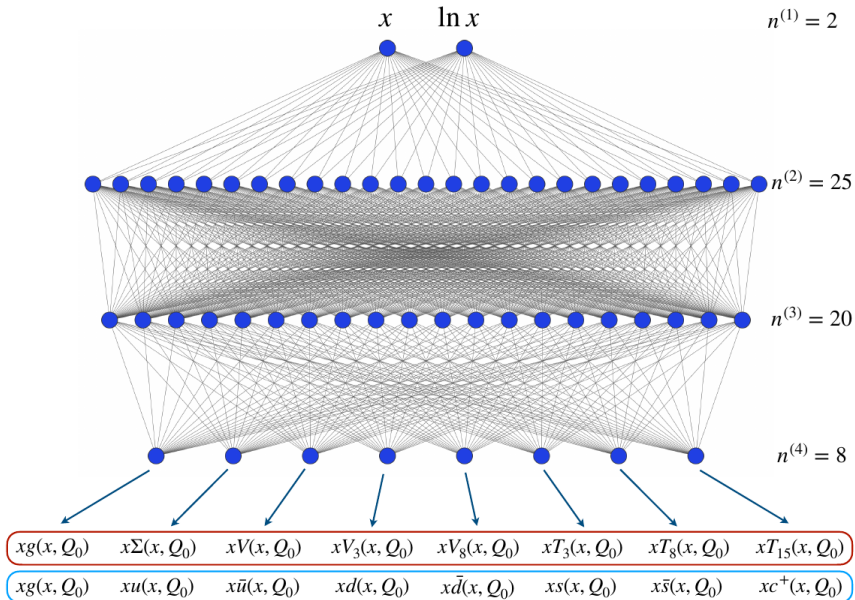
$$\mathcal{O}_i^{(art)(k)} = \mathcal{O}_i^{(exp)} + L_{ij} r_i^{(k)} \quad i, j = 1 \dots N_{\text{dat}} \quad k = 1 \dots, N_{\text{rep}} \quad \text{cov} = L \cdot L^T$$

$$\langle \mathcal{O}[f(x, Q^2)] \rangle = \frac{1}{N_{\text{rep}}} \sum_{k=1}^{N_{\text{rep}}} \mathcal{O}[f^{(k)}(x, Q^2)]$$

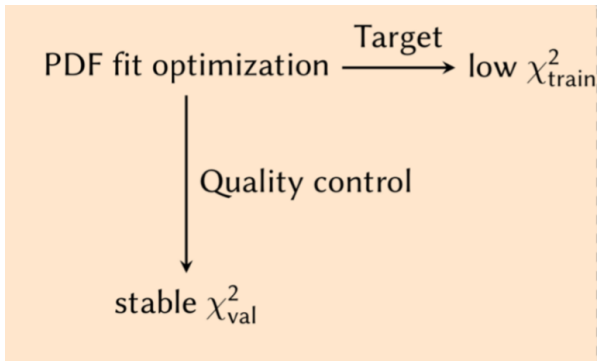
$$\sigma_{\mathcal{O}[f(x, Q^2)]} = \left[ \frac{1}{N_{\text{rep}} - 1} \sum_{k=1}^{N_{\text{rep}}} \left( \mathcal{O}[f^{(k)}(x, Q^2)] - \langle \mathcal{O}[f(x, Q^2)] \rangle \right)^2 \right]^{1/2}$$



# Parametrisation



# Optimisation

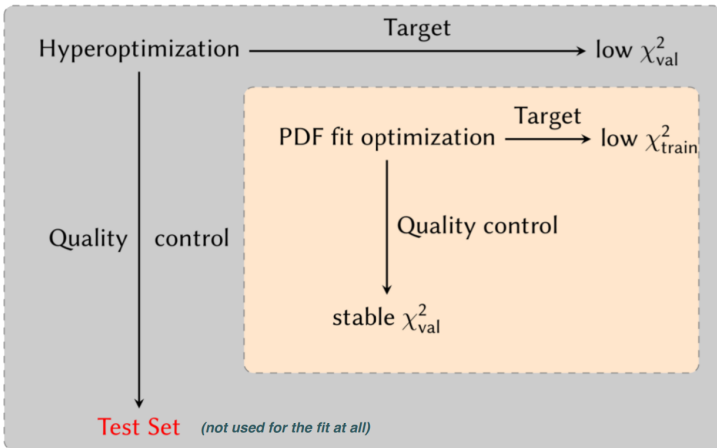


$$\chi^2 = \sum_{i,j}^{N_{\text{dat}}} (T_i[\{\boldsymbol{\theta}\}] - D_i) (\text{cov}^{-1})_{ij} (T_j[\{\boldsymbol{\theta}\}] - D_j)$$

$$(\text{cov})_{ij} = \delta_{ij} s_i^2 + \left( \sum_{\alpha}^{N_c} \sigma_{i,\alpha}^{(c)} \sigma_{j,\alpha}^{(c)} + \sum_{\alpha}^{N_{\mathcal{L}}} \sigma_{i,\alpha}^{(\mathcal{L})} \sigma_{j,\alpha}^{(\mathcal{L})} \right) D_i D_j$$

stochastic gradient descent with backpropagation

# Hyperoptimisation: fitting the methodology



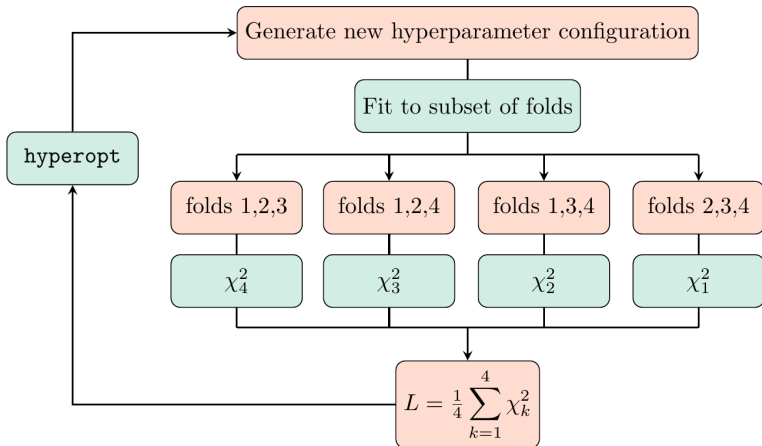
Compare to a Test Set (new set of data previously not used at all)

Who picks the Test Set? Automatic generalisation based on K foldings

Divide the data into  $n$  representative sets, fit  $n - 1$  sets and use the  $n$ -th set as test set

Hyperoptimise on mean and standard deviation of  $\chi_{test,i}^2$ ,  $i = 1 \dots n$

# Hyperoptimisation: $K$ -folding



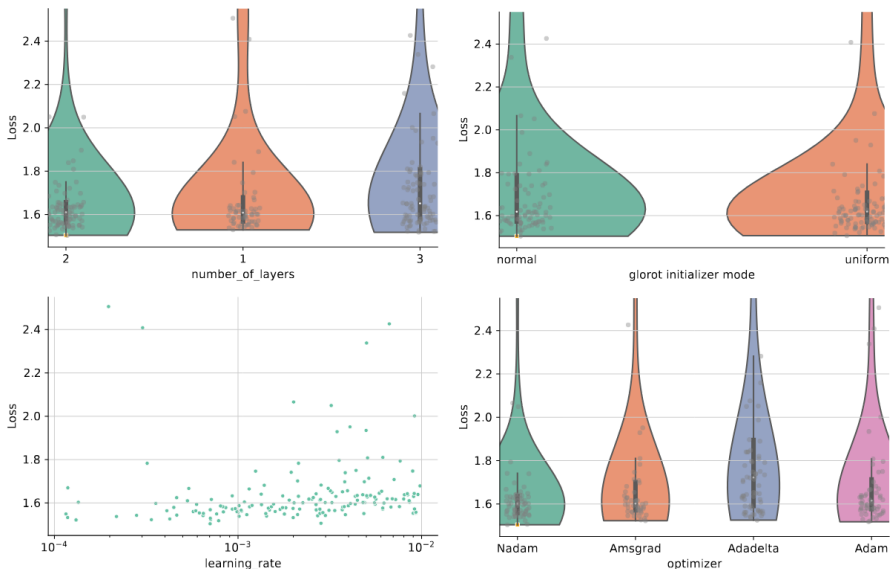
Compare to a Test Set (new set of data previously not used at all)

Who picks the Test Set? Automatic generalisation based on  $K$  foldings

Divide the data into  $n$  representative sets, fit  $n - 1$  sets and use the  $n$ -th set as test set

Hyperoptimise on mean and standard deviation of  $\chi_{\text{test},i}^2$ ,  $i = 1 \dots n$

# Hyperparameters



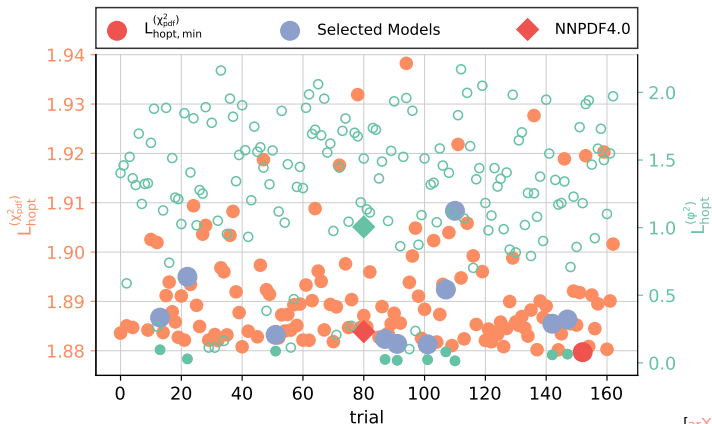


# Hyperoptimisation: metrics

$$\varphi_{\chi^2}^2 \equiv \langle \chi^2 [T^{(k)}, D^{(0)}] \rangle_{\text{rep}} - \chi^2 [T_{\text{rep}}, D^{(0)}] = \frac{1}{n_{\text{dat}}} \sum_{i,j}^{\text{dat}} (\text{cov})^{-1} T_{ij} \quad L_{\text{hopt}}^{(\varphi^2)} = \left( \frac{1}{n_{\text{fold}}} \sum_p^{\text{fold}} \varphi_{\chi^2}^{2(p)} \right)^{-1}$$

Select hyperparameters leading to best  $\chi^2$  and largest PDF errors in non-fitted data

Sample over the acceptable hyperparameters displaying comparable performance



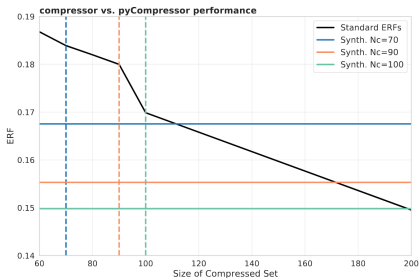
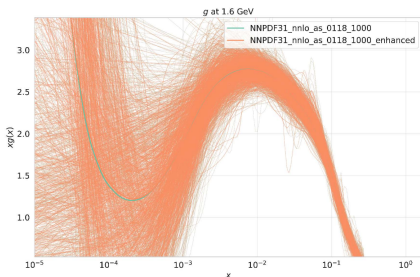
[arXiv:2410.16248]

# Delivery: compression

Find a subset of replicas that describe the underlying probability distribution as accurately as the original ensemble of replicas

$$\text{ERF} = \frac{1}{N_{\text{EST}}} \sum_k \frac{1}{N_k} \sum_i \left( \frac{C^{(k)}(x_i) - P^{(k)}(x_i)}{P^{(k)}(x_i)} \right)^2$$

Use GANs to enhance the number of replicas before minimising ERF



GAN enhancement allows one to achieve the same ERF with fewer replicas

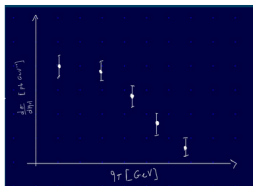
[EPJ C81 (2021) 530]

# Validation: closure tests

Fit PDFs to pseudodata generated assuming a known underlying law

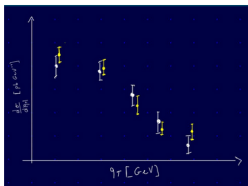
Level 0

no fluctuations



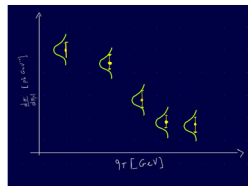
Level 1

Gaussian fluctuation

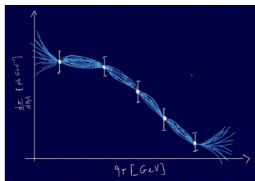


Level 2

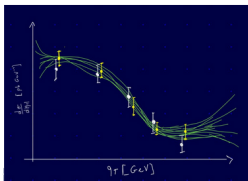
Monte Carlo replicas



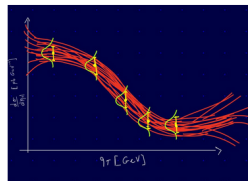
interpolation uncertainty



functional uncertainty



data uncertainty



# Closure tests at work

Data region: closure tests

Fit PDFs to pseudodata generated assuming a known underlying law

Define bias and variance

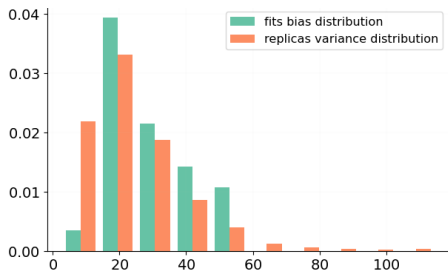
**bias** difference of central prediction and truth

**variance** uncertainty of replica predictions

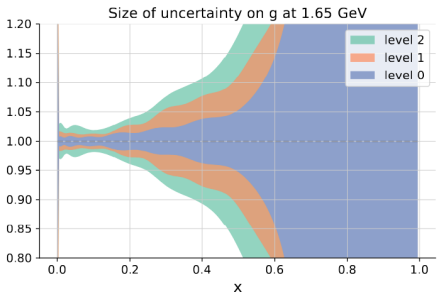
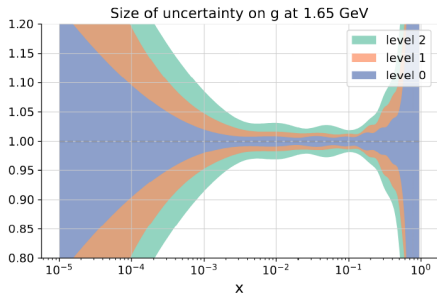
If PDF uncertainty faithful, then

$$E[\text{bias}] = \text{variance}$$

25 fits, 40 replicas each



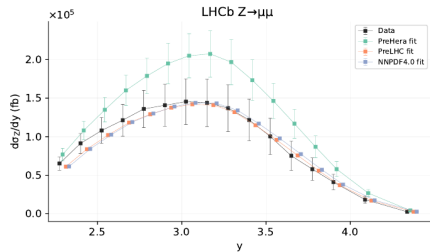
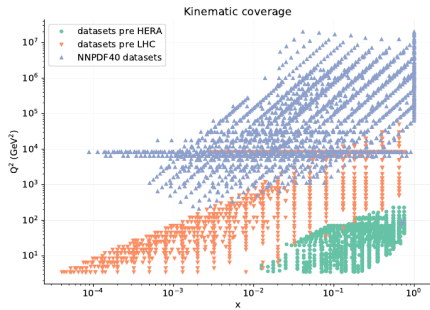
[EPJ C77 (2017) 663; EPJ C82 (2022) 330]



# Future tests

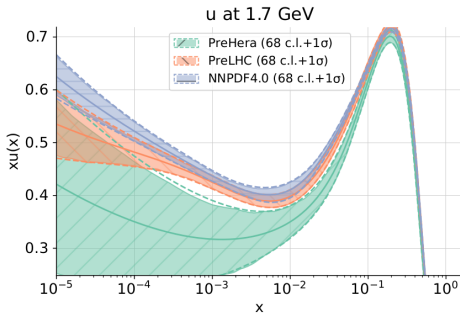
Extrapolation regions: future test

Test PDF uncertainties on data sets not included in a given PDF fit that cover unseen kinematic regions



Data set	NNPDF4.0	pre-LHC	pre-HERA
pre-HERA	1.09	1.01	0.90
pre-LHC	1.21	1.20	23.1
NNPDF4.0	1.29	3.30	23.1

Only exp. cov. matrix

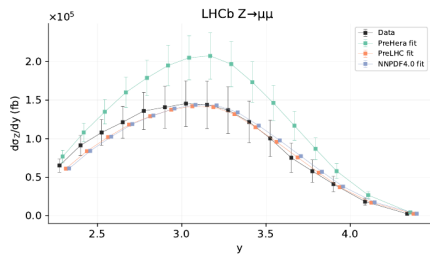
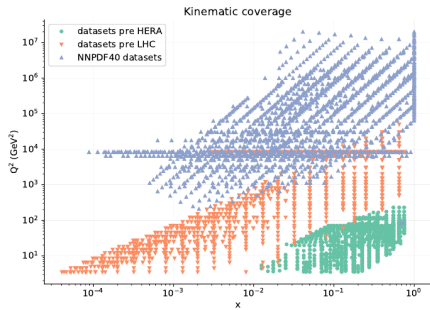


[Acta Phys. Polon. B52 (2021) 243]

# Future tests

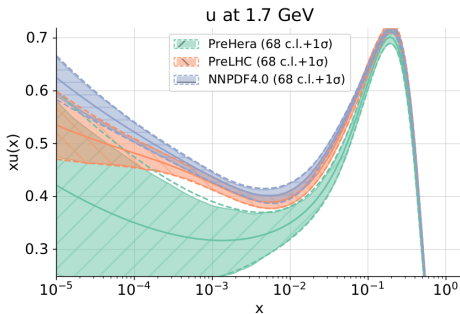
Extrapolation regions: future test

Test PDF uncertainties on data sets not included in a given PDF fit that cover unseen kinematic regions



Data set	NNPDF4.0	pre-LHC	pre-HERA
pre-HERA			0.86
pre-LHC		1.17	<b>1.22</b>
NNPDF4.0	1.12	<b>1.30</b>	<b>1.38</b>

Exp+PDF cov. matrix

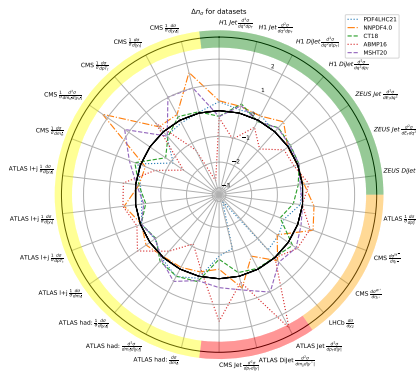
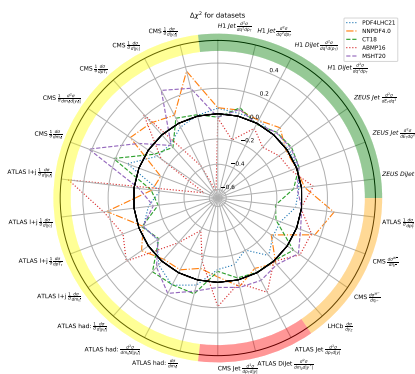


[Acta Phys. Polon. B52 (2021) 243]

# To conclude: are all PDF sets equally accurate?

$$\Delta\chi^2^{(i)} = \frac{\chi_{\text{exp+th}}^{2(i)} - \langle \chi_{\text{exp+th}}^2 \rangle_{\text{pdfs}}}{\langle \chi_{\text{exp+th}}^2 \rangle_{\text{pdfs}}}$$

$$\Delta n_{\sigma}^{(i)} = \frac{\chi_{\text{exp+th}}^{2(i)} - \langle \chi_{\text{exp+th}}^2 \rangle_{\text{pdfs}}}{\sqrt{2/n_{\text{data}}}}$$

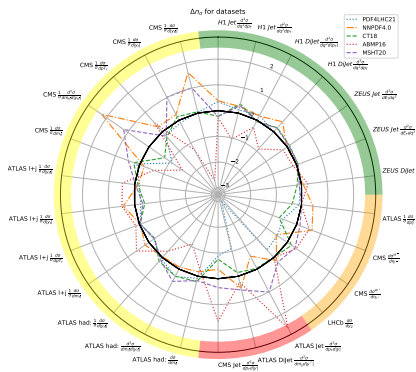
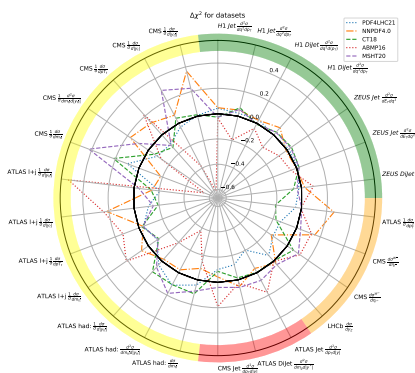


The PDF set obtained with ML is more precise and more accurate than all the others

# To conclude: are all PDF sets equally accurate?

$$\Delta\chi^2^{(i)} = \frac{\chi_{\text{exp+th}}^{2(i)} - \langle \chi_{\text{exp+th}}^2 \rangle_{\text{pdfs}}}{\langle \chi_{\text{exp+th}}^2 \rangle_{\text{pdfs}}}$$

$$\Delta n_{\sigma}^{(i)} = \frac{\chi_{\text{exp+th}}^{2(i)} - \langle \chi_{\text{exp+th}}^2 \rangle_{\text{pdfs}}}{\sqrt{2/n_{\text{data}}}}$$



The PDF set obtained with ML is more precise and more accurate than all the others

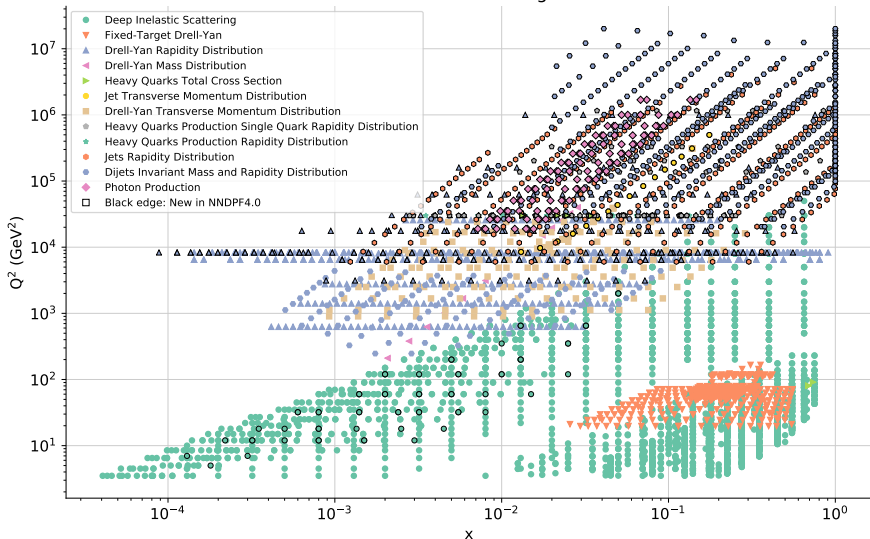
## Thank you



# Extra material

# Overview of experimental data

Kinematic coverage

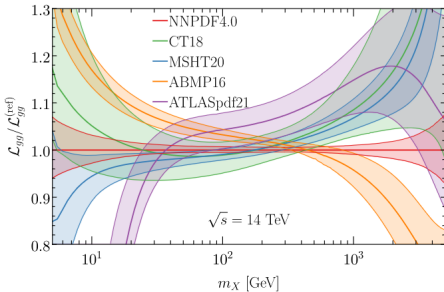
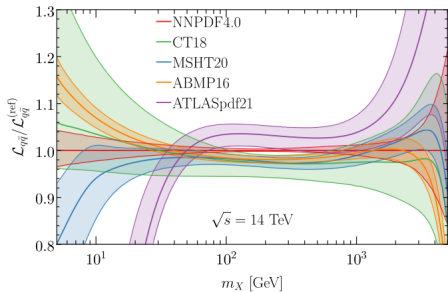
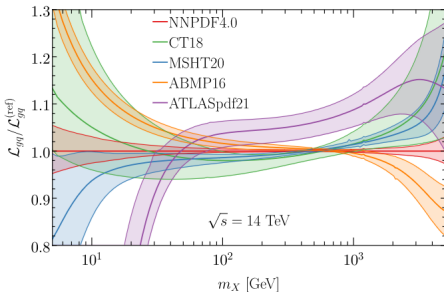
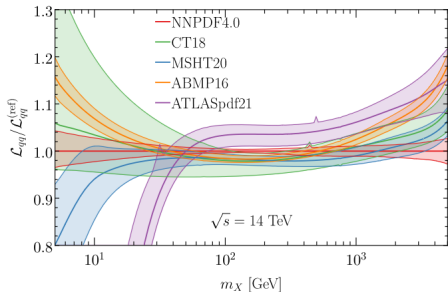


$N_{\text{dat}} = 4618$

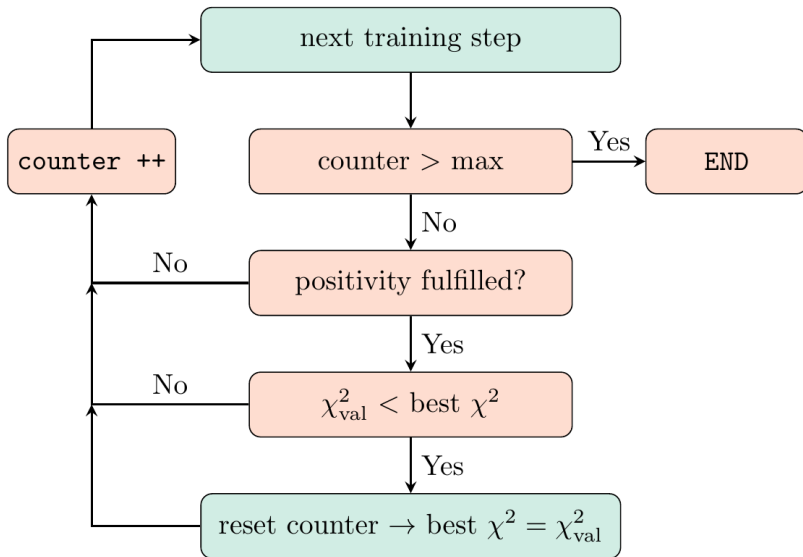
# Overview of current PDF determinations

	NNPDF4.0	MSHT20	CT18	HERAPDF2.0	CJ22	ABMP16
Fixed-target DIS	✓	✓	✓	✗	✓	✓
JLAB	✗	✗	✗	✗	✓	✗
HERA I+II	✓	✓	✓	✓	✓	✓
HERA jets	✓	✗	✗	✓	✗	✗
Fixed target DY	✓	✓	✓	✗	✓	✓
Tevatron $W, Z$	✓	✓	✓	✗	✓	✓
LHC vector boson	✓	✓	✓	✗	✓	✓
LHC $W + c Z + c$	✓	✗	✗	✗	✗	✗
Tevatron jets	✓	✓	✓	✗	✓	✗
LHC jets	✓	✓	✓	✗	✗	✗
LHC top	✓	✓	✗	✗	✗	✓
LHC single $t$	✓	✗	✗	✗	✗	✗
LHC prompt $\gamma$	✓	✗	✗	✗	✗	✗
statistical treatment	Monte Carlo	Hessian $\Delta\chi^2$ dynamical	Hessian $\Delta\chi^2$ dynamical	Hessian $\Delta\chi^2 = 1$	Hessian $\Delta\chi^2 = 1.645$	Hessian $\Delta\chi^2 = 1$
parametrisation	Neural Network	Chebyshev pol.	Bernstein pol.	polynomial	polynomial	polynomial
HQ scheme	FONLL	TR'	ACOT- $\chi$	TR'	ACOT- $\chi$	FFN
accuracy	aN <sup>3</sup> LO	aN <sup>3</sup> LO	NNLO	NNLO	NLO	NNLO
latest update	EPJ C82 (2022) 428	EPJ C81 (2021) 341	PRD 103 (2021) 014013	EPJ C82 (2022) 243	PRD 107 (2023) 113005	PRD 96 (2017) 014011

# Comparing PDF sets



# Optimisation



stochastic gradient descent with backpropagation