

Towards NNPDF4.0

Tommaso Giani, Michael Wilson and Shayan Iranipour
RADCOR-LoopFest

The road to NNPDF4.0

Progresses towards extending **data**, **theory** and **methodology**.

06/2017	NNPDF3.1	[EPJ C77 (2017) 663]
10/2017	NNPDF3.1 _{sx} : PDFs with small-x resummation	[EPJ C78 (2018) 321]
12/2017	NNPDF3.1 _{luxQED} : photon PDF	[SciPost Phys. 5 (2018) 008]
02/2018	NNPDF3.1+ATLAS _{photon} : inclusion of direct photon data	[EPJ C78 (2018) 470]
12/2018	NNPDF3.1 _{alphas} : α_s from correlated-replica method	[EPJ C78 (2018) 408]
05/2019	NNPDF3.1 _{th} : missing higher-order uncertainties in a fit	[EPJ C79 (2019) 838; <i>ibid.</i> 931]
07/2019	Gradient descent and hyperoptimisation in PDF fits	[EPJ C79 (2019) 676]
12/2019	NNPDF3.1 _{singletop} : inclusion of single top t-channel data	[JHEP 05 (2020) 067]
05/2020	NNPDF3.1 _{dijets} : comparative study of single- and di-jets	[EPJ C80 (2020) 797]
06/2020	Positivity of MS PDFs	[JHEP 11 (2020) 129]
08/2020	PineAPPL: fast evaluation of EW \times QCD corrections	[JHEP 12 (2020) 108]
08/2020	NNPDF3.1 _{strangeness} : assessment of strange-sensitive data	[EPJ C80 (2020) 1168]
11/2020	NNPDF3.1 _{deu} : deuteron uncertainties in a fit	[EPJ C81 (2021) 37]
03/2021	Future tests	[arXiv:2103.08606]
2021	NNPDF4.0	[to appear]

- ① Dataset and Methodology
- ② PDF validation
- ③ PDFs and Phenomenology

Dataset and Methodology

Experimental data in NNPDF4.0

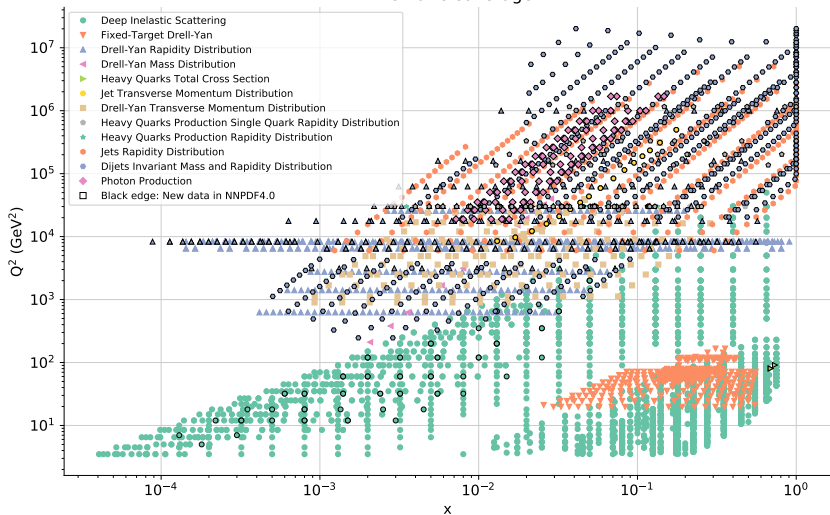
Data set	NNPDF4.0	NNPDF3.1	ABMP16	CT18	MSHT20
ATLAS W, Z 7 TeV (2010)	✓	✓	✓	✓	✓
ATLAS W, Z 7 TeV (2011)	✓	✓	✗	✓	✓
ATLAS low-mass DY 7 TeV	✓	✓	✗	✗	✗
ATLAS high-mass DY 7 TeV	✓	✓	✗	✗	✓
ATLAS W 8 TeV	✓	✗	✗	✗	✓
ATLAS DY 2D 8 TeV	✓	✗	✗	✗	✓
ATLAS high-mass DY 2D 8 TeV	✓	✗	✗	✗	✓
ATLAS $\sigma_{W,Z}$ 13 TeV	✓	✗	✓	✗	✗
ATLAS W^+ +jet 8 TeV	✓	✗	✗	✗	✓
ATLAS Z p_T 8 TeV	✓	✓	✗	✓	✓
ATLAS $\sigma_{H^0}^{tot}$ 7, 8 TeV	✓	✓	✓	✗	✗
ATLAS $\sigma_{H^0}^{tot}$ 13 TeV	✓	✓	✓	✗	✗
ATLAS $t\bar{t}$ lepton+jets 8 TeV	✓	✓	✗	✓	✓
ATLAS $t\bar{t}$ dilepton 8 TeV	✓	✗	✗	✗	✓
ATLAS single-inclusive jets 7 TeV, R=0.6	✗	✓	✗	✓	✓
ATLAS single-inclusive jets 8 TeV, R=0.6	✓	✗	✗	✗	✗
ATLAS dijets 7 TeV, R=0.6	✓	✗	✗	✗	✗
ATLAS direct photon production 13 TeV	✓	✗	✗	✗	✗
ATLAS single top R_t 7, 8, 13 TeV	✓	✗	✓	✗	✗
ATLAS single top diff. 7, 8 TeV	✓	✗	✗	✗	✗
ATLAS single top diff. 8 TeV	✓	✗	✗	✗	✗

Data set	NNPDF4.0	NNPDF3.1	ABMP16	CT18	MSHT20
CMS W electron asymmetry 7 TeV	✓	✓	✗	✓	✓
CMS W muon asymmetry 7 TeV	✓	✓	✓	✓	✗
CMS Drell-Yan 2D 7 TeV	✓	✓	✗	✗	✓
CMS W rapidity 8 TeV	✓	✓	✓	✓	✓
CMS Z p_T 8 TeV	✓	✓	✗	✓	✗
CMS $W + c$ 7 TeV	✓	✓	✗	✗	✓
CMS $W + c$ 13 TeV	✓	✗	✗	✗	✗
CMS single-inclusive jets 2.76 TeV	✗	✓	✗	✗	✓
CMS single-inclusive jets 7 TeV	✗	✓	✗	✓	✓
CMS dijets 7 TeV	✓	✗	✗	✗	✗
CMS single-inclusive jets 8 TeV	✗	✗	✗	✓	✓
CMS 3D dijets 8 TeV	✓	✗	✗	✗	✗
CMS $\sigma_{H^0}^{tot}$ 5 TeV	✓	✗	✓	✗	✗
CMS $\sigma_{H^0}^{tot}$ 7, 8 TeV	✓	✓	✓	✗	✓
CMS $\sigma_{H^0}^{tot}$ 13 TeV	✓	✓	✓	✗	✗
CMS $t\bar{t}$ lepton+jets 8 TeV	✓	✓	✗	✗	✓
CMS $t\bar{t}$ 2D dilepton 8 TeV	✓	✗	✗	✓	✓
CMS $t\bar{t}$ lepton+jet 13 TeV	✓	✗	✗	✗	✗
CMS $t\bar{t}$ dilepton 13 TeV	✓	✗	✗	✗	✗
CMS single top $\sigma_t + \sigma_{\bar{t}}$ 7 TeV	✓	✗	✓	✗	✗
CMS single top R_t 8, 13 TeV	✓	✗	✓	✗	✗

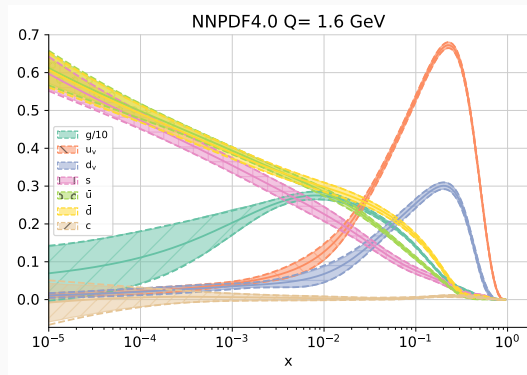
Data set	NNPDF4.0	NNPDF3.1	ABMP16	CT18	MSHT20
LHCb Z 940 pb	✓	✓	✗	✗	✓
LHCb $Z \rightarrow ee$ 2 fb	✓	✓	✓	✓	✓
LHCb $W, Z \rightarrow \mu$ 7 TeV	✓	✓	✓	✓	✓
LHCb $W, Z \rightarrow \mu$ 8 TeV	✓	✓	✓	✓	✓
LHCb $Z \rightarrow \mu\mu, ee$ 13 TeV	✓	✗	✗	✗	✗

- $\mathcal{O}(50)$ datasets investigated
- $\mathcal{O}(400)$ new datapoints wrt NNPDF3.1
- New processes included: **single top, W +jet, isolated photon, di-jets.**

Kinematic coverage



Data set	N_{dat}	χ^2/N_{dat}
Fixed target DIS	1881	1.10
Fixed target DY	189	1.00
HERA	1208	1.21
CDF	28	1.31
D0	37	1.00
ATLAS	621	1.18
Drell-Yan, 7, 8, 13 TeV	153	1.32
W+jet, 8 TeV	32	1.15
single top, 7, 8, 13 TeV	14	0.36
di-jets, 7 TeV	90	1.93
jets, 8 TeV	171	0.61
top pair, 7, 8, 13 TeV	16	2.30
ZpT, 8 TeV	92	0.86
direct photon, 13 TeV	53	0.72
CMS	411	1.40
Drell-Yan, 7, 8 TeV	154	1.34
single top, 7, 8, 13 TeV	3	0.43
di-jets, 7 TeV	54	1.67
di-jets, 8 TeV	122	1.50
top pair, 5, 7, 8 TeV	29	0.84
top pair, 13 TeV	21	0.67
ZpT, 8 TeV	28	1.42
LHCb	116	1.53
Total	4991	1.17



- overall good fit quality
- slight tension between ATLAS top-pair and di-jets data

Key differences with respect to the 3.1 methodology

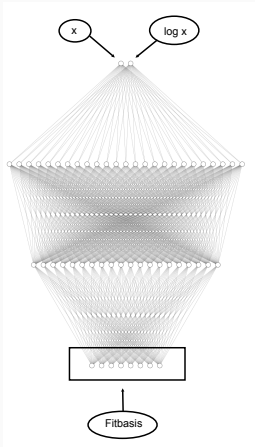
NNPDF 3.1 code

- C++ monolithic codebase
- In-house Machine Learning optimization framework
- One network per flavour
- Genetic Algorithm optimizer
- Fitting times of up to various days
- Fit parameters manually chosen (manual optimization of hyperparameters)

NNPDF 4.0 code

- Python object oriented codebase
- Freedom to use external libraries (default: TensorFlow)
- One network for all flavours
- Gradient Descent optimization
- Results available in less than an hour
- Fit parameters chosen automatically (hyperparameter scan)

NNPDF4.0 methodology



Evolution basis

$$q_k(x, Q_0) \propto x^{-\alpha_k} (1-x)^{\beta_k} NN(x),$$
$$q_k = \{V, V_3, V_8, T_3, T_8, T_{15}, \Sigma, g\}.$$

Flavour basis

$$q_k(x, Q_0) \propto (1-x)^{\beta_k} NN(x),$$
$$q_k = \{u, \bar{u}, d, \bar{d}, s, \bar{s}, c, g\}.$$

Physical constraints:

- ✓ PDFs positivity [\[JHEP 11 \(2020\) 129\]](#)
- ✓ Integrability of nonsinglet distributions

We will use the evolution basis, but as an additional check of our methodology we have explicitly checked fitbasis independence.

Beyond the PDF fit: fitting the methodology

The NNPDF methodology is defined by a set of hyperparameters: depth of the network, number of nodes, optimizer, training length, . . .

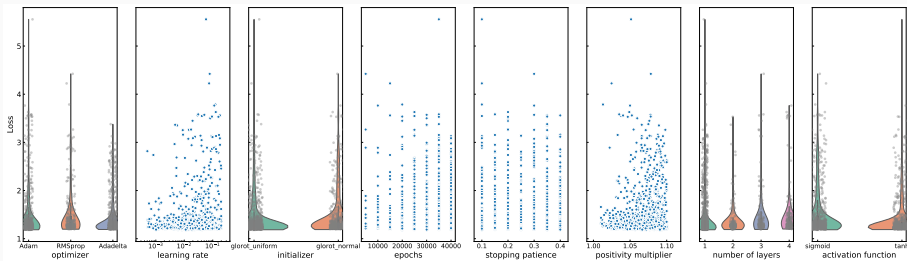
- ✗ Selecting manually the best set of parameters is a slow process and systematic success is not guaranteed

- ✓ Hyperparameter scan: let the computer decide automatically
 - Define a methodology (a specific hyperparameter combination)
 - Define a reward function to grade the methodology
 - Scan over thousands of hyperparameter combinations and select the best one

In NNPDF4.0 hyperparameter scan is implemented for the first time to select the best methodology.

Hyperparameter scan

Each blue dot corresponds to a fit of a different set of hyperparameters:



Thousands of fits for the hyperoptimization algorithm to choose:

- ✓ Optimizer
- ✓ Initializer
- ✓ Stopping Patience
- ✓ Number of Layers
- ✓ Learning Rate
- ✓ Epochs
- ✓ Positivity Multiplier
- ✓ Activation Function

K-folding algorithm: used to define the reward function and select the best hyperparameters combination (backup slides)

PDF validation

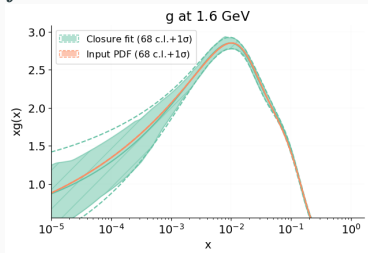
Closure Tests

Fit replicas to pseudodata in usual way

$$(1) \quad \begin{aligned} \mathbf{y} &= \mathbf{f} + \boldsymbol{\eta} + \boldsymbol{\epsilon} \\ &= \mathbf{z} + \boldsymbol{\epsilon}, \end{aligned}$$

where $\boldsymbol{\eta} \sim \mathcal{N}(0, C)$ and $\boldsymbol{\epsilon} \sim \mathcal{N}(0, C)$ are sampled independently.

Use predictions from an input PDF as proxy for \mathbf{f} .



Example closure fit and input PDF.

First presented in NNPDF3.0 [arxiv:1410.8849](https://arxiv.org/abs/1410.8849)

Allows testing of methodology, if the input assumptions hold.

For example:

Bias: (squared) difference between central prediction and true observable

Variance: variance of model predictions

Bias is a stochastic variable. If PDF uncertainty is faithful then

$$\mathbf{E}_{\boldsymbol{\eta}}[\text{bias}] = \text{variance} \quad (2)$$

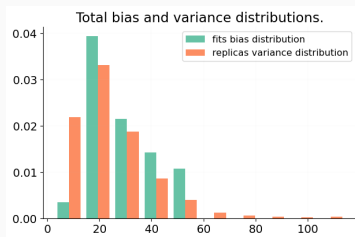
High demand on resources (many fits!) - made feasible with next generation fitting code.

Closure test results

Compare first moments:

$\sqrt{\mathbf{E}_\eta[\text{bias}]/\mathbf{E}_\eta[\text{variance}]}$	
Total	1.03 ± 0.5

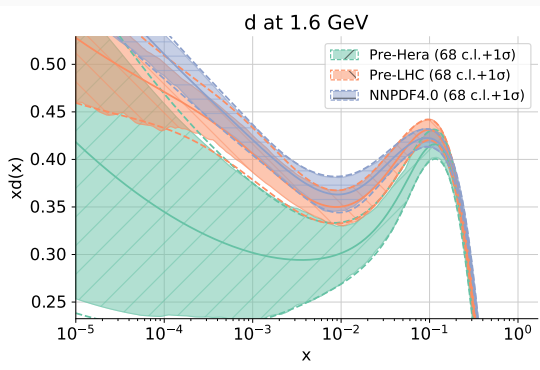
Alternatively look at the respective distributions of differences (averaged across data).



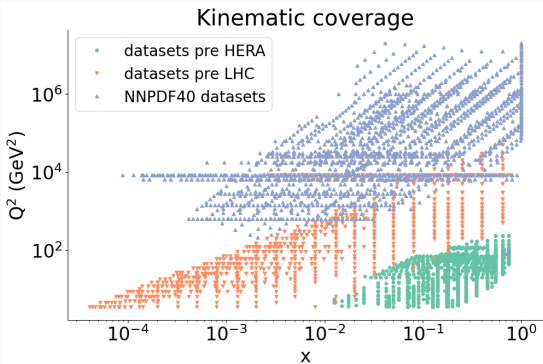
Bias distribution sampled with 25 fits, 40 replicas each.

Trusting uncertainties in extrapolation region.

- Smaller PDF uncertainties driven by methodological improvements and increased data.
- uncertainties tested in "data region"
- what about kin. regions not covered by data?
- Ideally, test "extrapolation region" with new experimental data - **not an option.**
- Retroactively order subsets of data chronologically, named "future tests" - for more information see [arxiv:2103.08606](https://arxiv.org/abs/2103.08606)

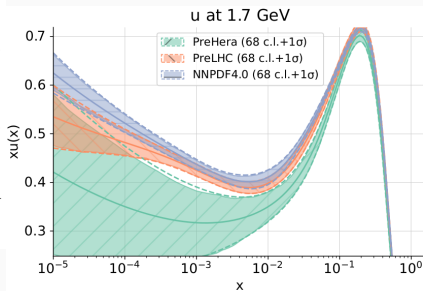


Future tests

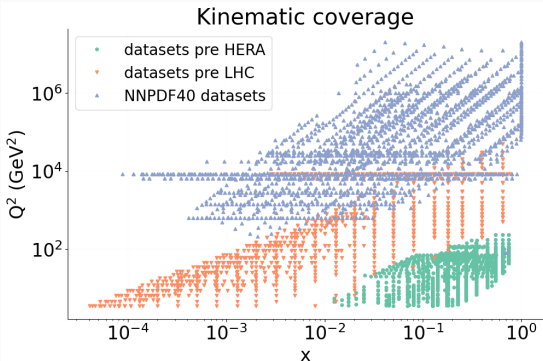


χ^2/N (only exp. covmat)

(dataset)	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

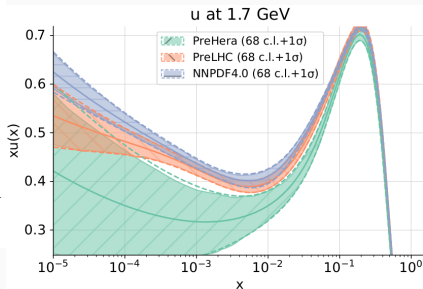


Future tests



χ^2/N (exp. and PDF covmat)

(dataset)	NNPDF4.0	pre-LHC	pre-Hera
pre-HERA			0.86
pre-LHC		1.17	1.22
NNPDF4.0	1.12	1.30	1.38

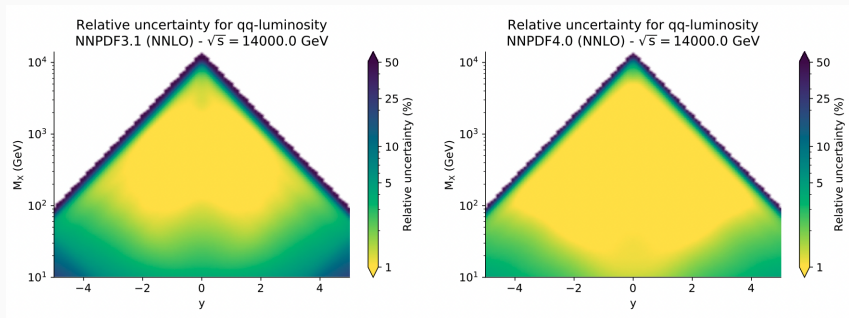


Total uncertainty increases up to a factor ~ 20 , and accommodates for difference between predictions and new data.

PDFs and Phenomenology

The road to 1%: the qq channel

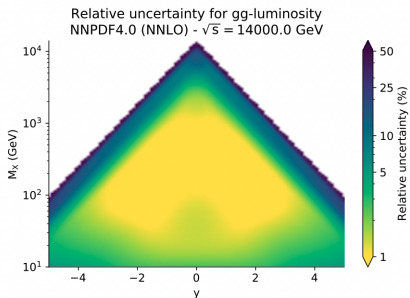
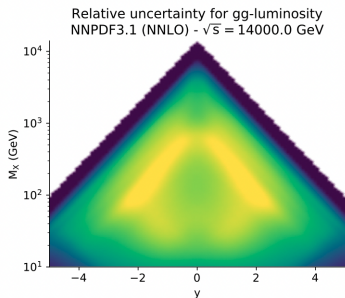
$$\mathcal{L}_{ij}(M_X, y, \sqrt{s}) = \frac{1}{s} \sum_{i,j} f_i \left(\frac{M_X e^y}{\sqrt{s}}, M_X \right) f_j \left(\frac{M_X e^{-y}}{\sqrt{s}}, M_X \right)$$



Plan to include MHO (arXiv:1905.0431) in a future study.

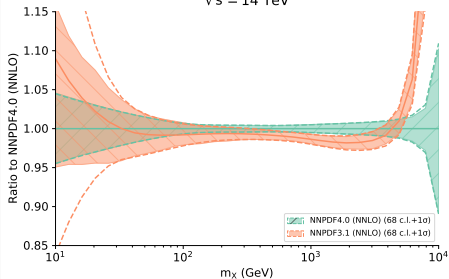
The road to 1%: the gg channel

$$\mathcal{L}_{ij}(M_X, y, \sqrt{s}) = \frac{1}{s} \sum_{i,j} f_i \left(\frac{M_X e^y}{\sqrt{s}}, M_X \right) f_j \left(\frac{M_X e^{-y}}{\sqrt{s}}, M_X \right)$$

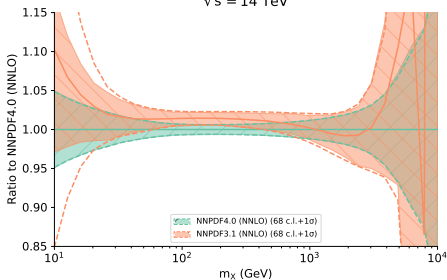


Precision: NNPDF4.0 vs NNPDF3.1

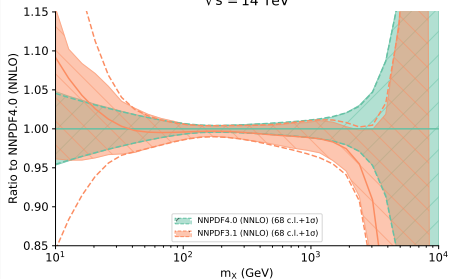
qq luminosity
 $\sqrt{s} = 14$ TeV



gg luminosity
 $\sqrt{s} = 14$ TeV



q\bar{q} luminosity
 $\sqrt{s} = 14$ TeV

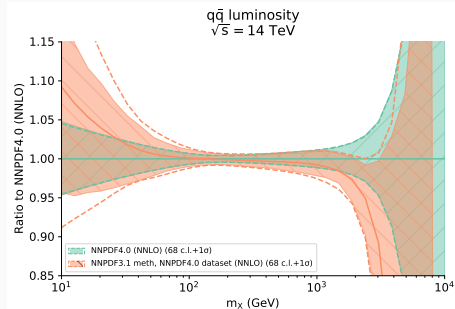
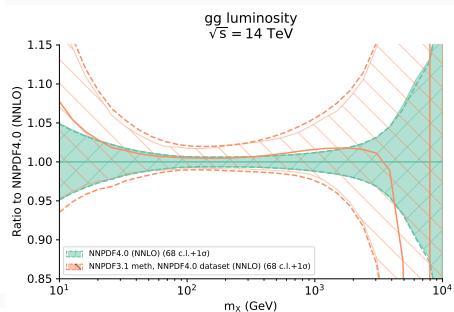
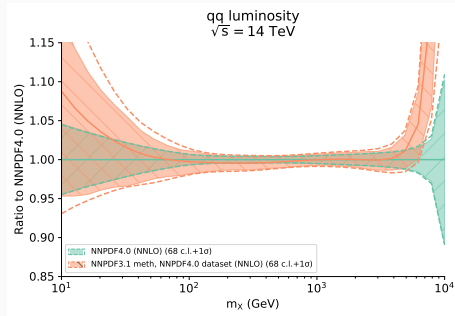


Methodology

Dataset	NNPDF3.1	NNPDF4.0
NNPDF3.1 (4093)	1.19	1.12
NNPDF4.0 (4491)	1.25	1.17

Shift in PDF luminosity due to addition of ~ 400 new data points. Reduced uncertainties highlighted in smaller error bands: NNPDF4.0 is more precise.

Accuracy: NNPDF4.0 data - NNPDF3.1 vs NNPDF4.0 methodology

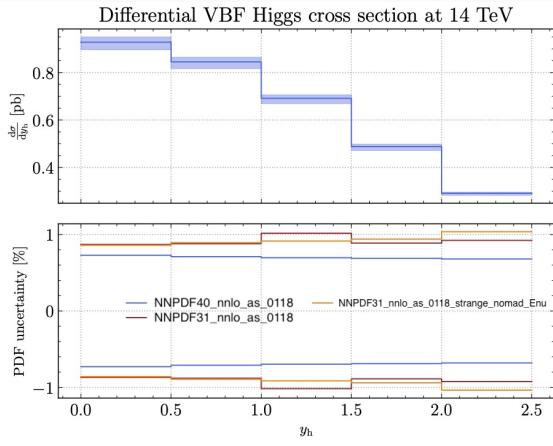


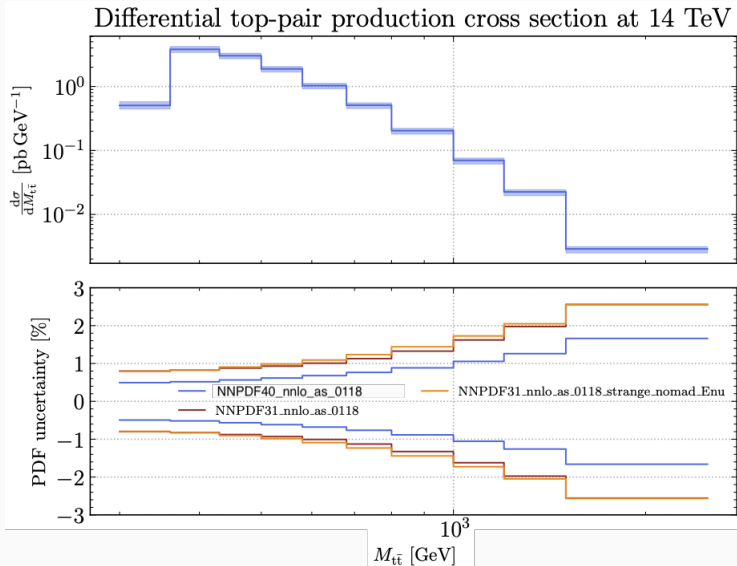
	Methodology	
Dataset	NNPDF3.1	NNPDF4.0
NNPDF3.1 (4093)	1.19	1.12
NNPDF4.0 (4491)	1.25	1.17

Luminosities fitted to the same dataset are compatible. Achieve a better fit to data:
 NNPDF4.0 is more accurate.

Accuracy and Precision

Reduced PDF luminosity uncertainties lead to $< 1\%$ uncertainty in various binnings at the observable level.





Plan to publicly release the code alongside code paper and in depth documentation at `docs.nnpdf.science`:

- Entire `n3fit/nnfit` fitting code: transparent methodology for extensions to other PDF fitting applications
- The `validphys` analysis suite: allowing for fully reproducible plots from various NNPDF papers.
- `Buildmaster` rawdata parsing code: ability for users to their own data to PDF fits.

Science using runcards and reproducibility

User exposed to human readable YAML format for analysis and fitting.

Runcards yield *reproducibility*.

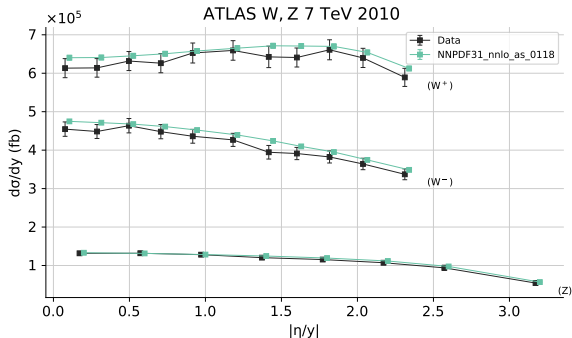
```
# runcard.yaml
# Any valid LHAPDF
pdf: NNPDF31_nnlo_as_0118
```

```
# Specify dataset to use
dataset: ATLASWZRAP36PB
```

```
# Theory settings
theoryid: nnlo
```

```
# Data cuts to use
cuts: internal
```

```
# The action we wish the code to perform
actions:
- data_theory_comparison
```

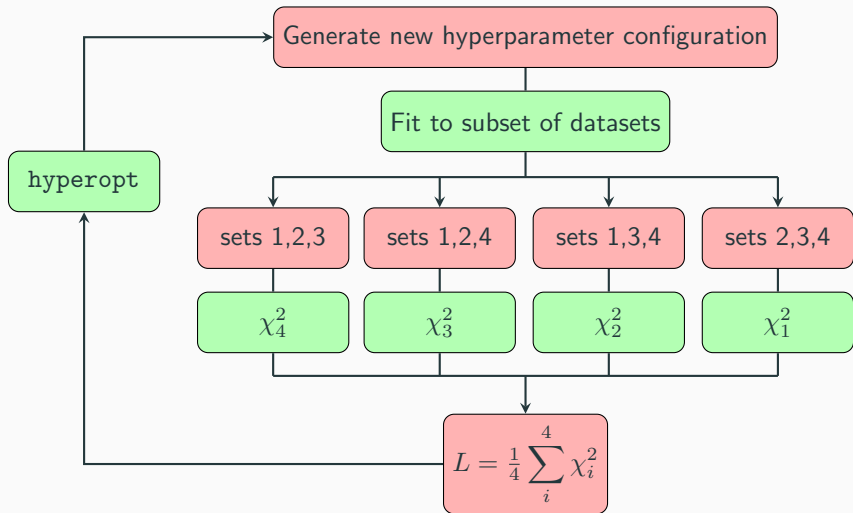


Conclusions

- Added host of new datasets and processes ~ 400 new datapoints
- Improved methodology using TensorFlow and hyperoptimization
- Faithfulness validated through closure and future tests
- Achieving 1% accuracy over impressively broad kinematic range
- Plan to release code publicly
- NNPDF4.0 is our recommended PDF set for LHC phenomenology: superseding previous iterations

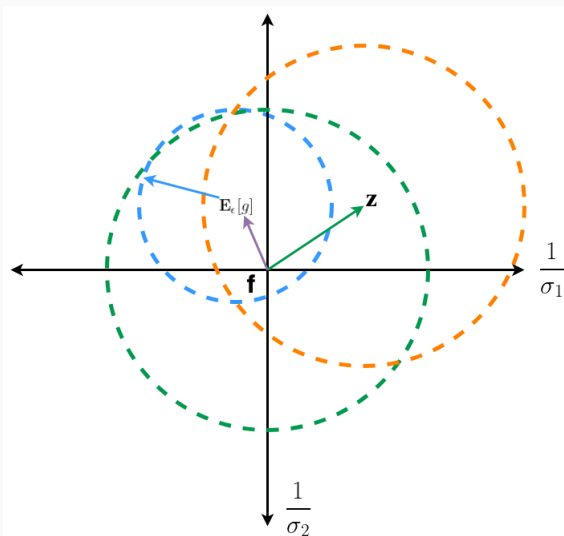
Additional slides

K-folding



Geometric Interpretation

Consider 2 data points on an axis in the basis which diagonalises C normalised by the square root of the eigenvalues:



Statistical estimators - more detail

Decompose the expectation value of the likelihood function, χ^2 , by completing the square.
Exposing some statistical indicators

$$\begin{aligned}\mathbf{E}_\epsilon[\chi^2(g; y)] &= \frac{1}{N_{\text{data}}} \mathbf{E}_\epsilon[(\mathbf{g} - \mathbf{y})^T C^{-1} (\mathbf{g} - \mathbf{y})] \\ &= \text{bias} + \text{variance} + \text{noise} - \text{crossterm}\end{aligned}\tag{3}$$

focus on the first two terms:

$$\text{bias} = \frac{1}{N_{\text{data}}} (\mathbf{f} - \mathbf{E}_\epsilon[\mathbf{g}])^T C^{-1} (\mathbf{f} - \mathbf{E}_\epsilon[\mathbf{g}])\tag{4}$$

$$\text{variance} = \mathbf{E}_\epsilon \left[\frac{1}{N_{\text{data}}} (\mathbf{g} - \mathbf{E}_\epsilon[\mathbf{g}])^T C^{-1} (\mathbf{g} - \mathbf{E}_\epsilon[\mathbf{g}]) \right]\tag{5}$$

where $\mathbf{E}_\epsilon[\cdot]$ is the expectation across replicas.

- Faithful uncertainties if $(\mathbf{g} - \mathbf{E}_\epsilon[\mathbf{g}])$ and $(\mathbf{f} - \mathbf{E}_\epsilon[\mathbf{g}])$ have same distribution.
- Sample $(\mathbf{g} - \mathbf{E}_\epsilon[\mathbf{g}])$ through sampling ϵ - usual MC replica procedure
- Sample $(\mathbf{f} - \mathbf{E}_\epsilon[\mathbf{g}])$ distribution through sampling η - only possible in closure test!

Closure test results

Breakdown of $\mathbf{E}_\eta[\text{bias}]/\mathbf{E}_\eta[\text{variance}]$ for out of sample data by experiment
- fitted on NNPDF3.1 dataset and validated on additional datasets to be included in NNPDF4.0

	$\sqrt{\mathbf{E}_\eta[\text{bias}]/\mathbf{E}_\eta[\text{variance}]}$
ATLAS	1.04 ± 0.4
CMS	1.04 ± 0.6
LHCb	0.82 ± 0.6
Total	1.03 ± 0.5

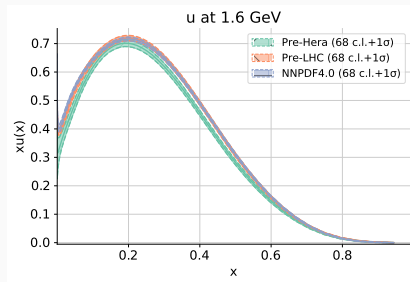
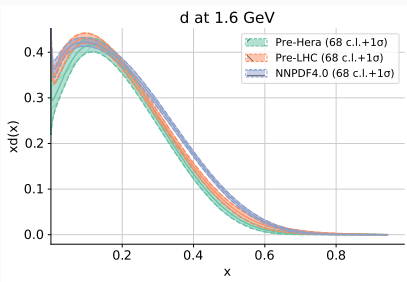
Full Future test subset specification

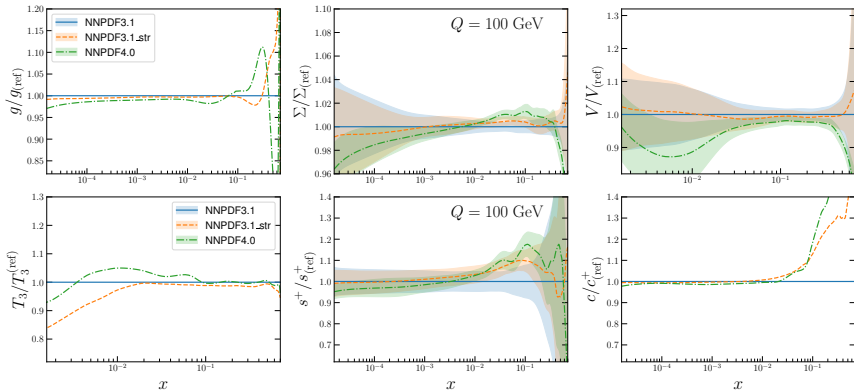
pre-HERA	Ref.	N_{pts}	NNPDF4.0 (nominal)	Ref.	N_{pts}
SMC d/μ	[14]	123	ATLAS $W^+ \rightarrow \mu e + \text{TV}$	[30]	36
SMC p	[14]	204	ATLAS $Z \text{ jet } 8 \text{ TV}$ ($p_{T,1}^Z, M_{12}$)	[30]	42
HLAC p	[17]	22	ATLAS $Z \text{ jet } 8 \text{ TV}$ ($p_{T,1}^Z, m_{12}$)	[30]	40
HLAC Z	[17]	22	ATLAS $Z \text{ jet } 7, 8, 13 \text{ TV}$	[37, 39]	33
ICEDRIP p	[18]	333	ATLAS d/μ nonresonant 1 TV	[30]	8
ICEDRIP d	[18]	333	ATLAS d/μ nonresonant 8 TV	[30]	4
CMORIS e_{had}^+	[19]	436	ATLAS of nonresonant (e, μ) dileptons 8 TV	[30]	3
CMORIS e_{had}^-	[19]	436	ATLAS jets 8 TV , $R=0.6$	[44]	171
TeVW e_{had}^+	[21]	29	ATLAS jets 7 TV , $R=0.6$	[42]	80
TeVW e_{had}^-	[21]	27	ATLAS direct photon production 13 TV	[43]	33
EV 906 d_{had}^+ / e_{had}^+	[22]	33	ATLAS single top $R=7 \text{ TV}$	[45]	1
EV 906 d_{had}^- / e_{had}^-	[22]	33	ATLAS single top $R=13 \text{ TV}$	[45]	1
EV 906 e_{had}^+	[22]	36	ATLAS single top jet cross. $7, 8 \text{ TV}$	[44, 46]	6
Total pre-HERA		2193	ATLAS single topology cross. $7, 8 \text{ TV}$	[44, 46]	6
pre-LHC	Ref.	N_{pts}	CMX W electron asymmetry 7 TV	[47]	11
HERA 1-II inclusive NC e^+p	[23]	139	CMX W muon asymmetry 7 TV	[48]	11
HERA 1-II inclusive NC e^+p 800 GeV	[23]	204	CMX Drell-Yan 13 TV	[49]	110
HERA 1-II inclusive NC e^+p 175 GeV	[23]	204	CMX $Z \text{ jet}$ ($p_{T,1}^Z, m_{12}$) 8 TV	[30]	22
HERA 1-II inclusive NC e^+p 820 GeV	[23]	79	CMX jets 7 TV	[51]	54
HERA 1-II inclusive NC e^+p 820 GeV	[23]	377	CMX $3D$ jets 8 TV	[51]	123
HERA 1-II inclusive CC e^+p	[23]	42	CMX $d^*/s, 8, 13 \text{ TV}$	[34, 50]	3
HERA 1-II inclusive CC e^+p	[23]	29	CMX of jet $p_{T,1}^Z$ 13 TV	[30]	9
HERA results d^*/s	[24]	27	CMX $d^*/s, 8 \text{ TV}$	[51]	1
HERA results d^*/s	[24]	26	CMX of $3D$ $(m_{12}, p_T) 8 \text{ TV}$	[51]	36
EW Z asymmetry	[25]	26	CMX of absolute m_{12} jet 13 TV	[51]	10
DZ Z asymmetry	[26]	26	CMX of absolute (m_{12}, p_T) 13 TV	[51]	10
DZ $W^+ \rightarrow \mu e$ asymmetry	[26]	9	CMX single top $m_{12} \rightarrow e_{\text{had}}^+ 7 \text{ TV}$	[45]	1
Total pre-LHC		1273	CMX single top $R=8, 9 \text{ TV}$	[45]	1
NNPDF4.0	Ref.	N_{pts}	CMX single top $R=13 \text{ TV}$	[45]	1
ATLAS $W, Z 7 \text{ TV}$	[30]	36	CMX Z 900 jets	[34]	39
ATLAS HJ $EV 7 \text{ TV}$	[30]	3	HLHC $Z \rightarrow \mu e \rightarrow 3 \text{ jet}$	[30]	17
ATLAS beamline $EV 7 \text{ TV}$	[30]	6	HLHC $W, Z \rightarrow \mu e \rightarrow 3 \text{ TV}$	[30]	20
ATLAS $W, Z 7 \text{ TV}$ central selection	[30]	36	HLHC $W, Z \rightarrow \mu e \rightarrow 8 \text{ TV}$	[37]	30
ATLAS $W, Z 7 \text{ TV}$ forward selection	[30]	33	HLHC $Z \rightarrow \mu e$ 13 TV	[30]	10
ATLAS $EV 20 8 \text{ TV}$	[30]	48	HLHC $Z \rightarrow ee$ 13 TV	[30]	15
ATLAS W, Z inclusive 13 TV	[30]	3	Total NNPDF4.0	1140	
ATLAS $W^+ \rightarrow \mu e 8 \text{ TV}$	[30]	36	Closed total	4891	

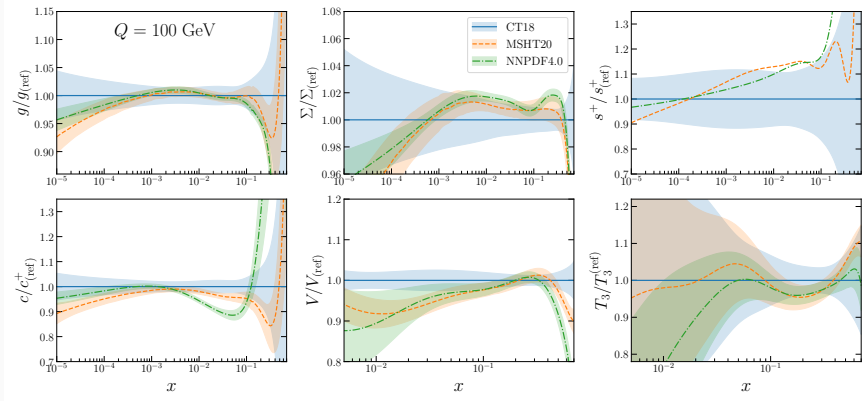
Table 1: The pre-HERA, pre-LHC and NNPDF4.0 datasets: for each experiment the reference to the original publication and the number of datapoints are given. Pre-HERA PDFs are fitted to pre-HERA data, pre-LHC PDFs are fitted to the pre-HERA and pre-LHC data, and NNPDF4.0 PDFs are fitted to the union of the three datasets.

For more information see
[arxiv:2103.08606](https://arxiv.org/abs/2103.08606)

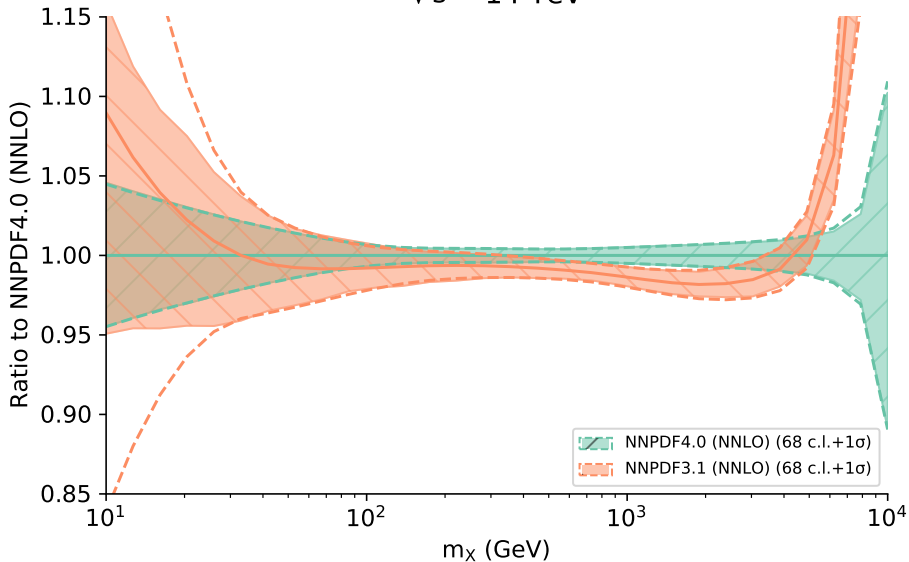
Future Test linear scale quarks



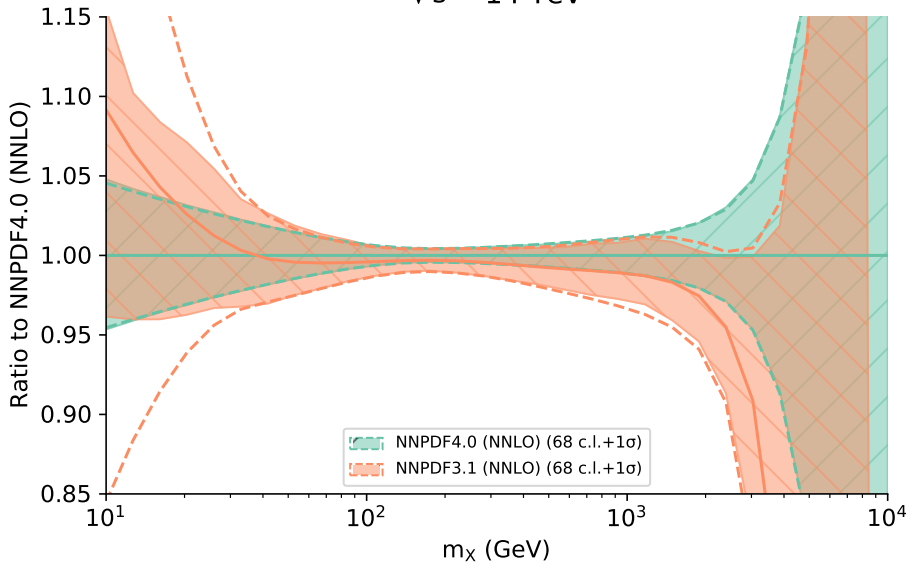




qq luminosity
 $\sqrt{s} = 14$ TeV



$q\bar{q}$ luminosity
 $\sqrt{s} = 14$ TeV



gg luminosity
 $\sqrt{s} = 14$ TeV

