





# Statistical issues in global fits: Lessons from PDF determinations

Juan Rojo Rudolf Peierls Center for Theoretical Physics University of Oxford

International Workshop on Global Fits to Neutrino Scattering Data and Generator Tuning (NuTune2016) Liverpool, 11/06/2016

Juan Rojo



# The inner life of protons : Parton Distribution Functions



#### Lepton vs Hadron Colliders

In high-energy **lepton colliders**, such as the **Large Electron-Positron Collider** (LEP) at CERN, the collisions involve **elementary particles** without substructure



**Cross-sections** in lepton colliders can be computed in perturbation theory using the **Feynman rules of the Standard Model Lagrangian** 

## Lepton vs Hadron Colliders

In high-energy **hadron colliders**, such as the LHC, the collisions involve **composite particles** (protons) with internal structure (quarks and gluons)



Calculations of cross-sections in hadron collisions require the combination of **perturbative**, **quark/gluon-initiated processes**, and **non-perturbative**, **parton distributions**, information

### Parton Distributions

The distribution of energy that **quarks and gluons carry inside the proton** is quantified by the **Parton Distribution Functions (PDFs)** 



**PDFs are** determined by **non-perturbative QCD dynamics:** cannot be computed from first principles, and need to be **extracted from experimental data** with a **global analysis** 

Energy conservation

$$\int_0^1 dx \left( g(x,Q) + \sum_q q(x,Q) \right) = 1$$

Dependence with quark/gluon collision energy Q determined in perturbation theory

$$\frac{\partial g(x,Q)}{\partial \ln Q} = P_g(\alpha_s) \otimes g(x,Q) + P_q(\alpha_s) \otimes q(x,Q)$$

Juan Rojo

#### The Factorization Theorem

The **QCD Factorization Theorem** guarantees **PDF universality: extract them from a subset of** process and use them to provide pure predictions for new processes

$$\sigma_{lp}\simeq\widetilde{\sigma}_{lq}\left(lpha_{s},lpha
ight)\otimes q(x,Q)$$

 $\sigma_{pp} \simeq \widetilde{\sigma}_{q\bar{q}} \left( \alpha_s, \alpha 
ight) \otimes q(x_1, Q) \otimes \bar{q}(x_2, Q)$ 



## The global PDF analysis

- Combine state-of-the-art theory calculations, the constraints from PDF-sensitive measurements from different processes and colliders, and a statistically robust fitting methodology
- Extract Parton Distributions at hadronic scales of a few GeV, where non-perturbative QCD sets in
- Use **perturbative evolution** to compute PDFs at high scales as **input to LHC predictions**



## The NNPDF approach

A **novel approach to PDF determination**, improving the limitations of the traditional PDF fitting methods with the use of **advanced statistical techniques** such as **machine learning** and **multivariate analysis** 

#### Non-perturbative PDF parametrization

- **Traditional approach**: based on **restrictive functional forms** leading to strong theoretical bias
- Solution: use Artificial Neural Networks as universal unbiased interpolants

PDF uncertainties and propagation to LHC calculations

- Traditional approach: Hessian method, limited to Gaussian/linear approximation
- NNPDF solution: based on the Monte Carlo replica method to create a probability distribution in the space of PDFs. Specially critical in extrapolation regions (i.e. high-*x*) for New Physics searches

Fitting technique

- **Fraditional approach**: deterministic minimization of  $\chi^2$ , flat directions problem
- NNPDF solution: Genetic Algorithms to explore efficiently the vast parameter space, with crossvalidation to avoid fitting stat fluctuations

## The Monte Carlo replica method

- Two main approaches to estimate PDF uncertainties: the Hessian method and the Monte Carlo method
- $\Im$  In the **Hessian method**, the  $\chi^2$  is expanded quadratically in the **fit parameters**  $\{a_n\}$  around the best fit

$$H_{lm}^{n} \equiv \frac{1}{2} \frac{\partial \chi_{n}^{2}}{\partial a_{l} \partial a_{m}} = \sum_{i=1}^{N_{\text{pts.}}} \frac{1}{(\sigma_{n,i}^{\text{uncorr.}})^{2} + \sum_{k} (\sigma_{n,k,i}^{\text{corr.}})^{2}} \frac{\partial T_{n,i}(\{a\})/\mathcal{N}_{n}}{\partial a_{l}} \frac{\partial T_{n,i}(\{a\})}{\partial a_{m}}$$

Final The Hessian matrix is diagonalized, and PDF errors on cross sections F from linear error propagation

$$\Delta \chi_{\text{global}}^2 \equiv \chi_{\text{global}}^2 - \chi_{\text{min}}^2 = \sum_{i,j=1}^n H_{ij}(a_i - a_i^0)(a_j - a_j^0) \qquad \Delta F = \frac{1}{2} \sqrt{\sum_{k=1}^n \left[F(S_k^+) - F(S_k^-)\right]^2},$$

- In the **Monte Carlo replica method**, pseudo-data replicas with same fluctuations as real data are generated, and then a PDF fit is performed **in each individual replica**
- Leads to probability distribution in the space of PDFs, without linear/Gaussian approximations

$$D_{m,i} \rightarrow \left( D_{m,i} + R_{m,i}^{\text{uncorr.}} \sigma_{m,i}^{\text{uncorr.}} + \sum_{k=1}^{N_{\text{corr.}}} R_{m,k}^{\text{corr.}} \sigma_{m,k,i}^{\text{corr.}} \right) \cdot \left( 1 + R_m^{\mathcal{N}} \sigma_m^{\mathcal{N}} \right)$$

$$\Delta F = \sqrt{\frac{N_{\text{rep}}}{N_{\text{rep}} - 1}} \left( \langle F^2 \rangle - \langle F \rangle^2 \right).$$

Juan Rojo

## ANN for PDF parametrization

- ANNs are routinely exploited in **high-energy physics**, in most cases as **classifiers** to separate between interesting and more mundane events
- ANNs also provide universal unbiased interpolants to parametrize the non-perturbative dynamics that determines the size and shape of the PDFs from experimental data



$$g(x,Q_0) = A_g (1-x)^{a_g} x^{-b_g} \left(1 + c_g \sqrt{s} + d_g x + \dots\right)$$

$$g(x, Q_0) = A_g \text{ANN}_g(x)$$

$$ANN_{g}(x) = \xi^{(L)} = \mathcal{F}\left[\xi^{(1)}, \{\omega_{ij}^{(l)}\}, \{\theta_{i}^{(l)}\}\right]$$
$$\xi_{i}^{(l)} = g\left(\sum_{j=1}^{n_{l-1}} \omega_{ij}^{(l-1)} \xi_{j}^{(l-1)} - \theta_{i}^{(l)}\right)$$

- ANNs eliminate **theory bias** introduced in PDF fits from choice of *ad-hoc* functional forms
- NNPDF fits used O(400) free parameters, to be compared with O(10-20) in traditional PDFs. Results stable if O(4000) parameters used!
- Faithful extrapolation: PDF uncertainties blow up in regions with scarce experimental data

10

#### Artificial Neural Networks vs Polynomials

Geometric Compare a **benchmark PDF analysis** where **the same dataset** is fitted with **Artificial Neural Networks** and with **standard polynomials**, other settings identical)

ANNs avoid biasing the PDFs, faithful extrapolation at small-x (very few data, thus error blow up)





# **Combining Inconsistent Experiments in the PDF fit**

NNPDF3.0 NLO dataset



#### Experimental data

- Final Provide the second secon
  - **Type of high energy collision** (lepton-proton, proton-proton), **center-of-mass energy** of collision
  - Whether of not **experimental correlated systematics** are available, and if so, in **which format**
  - Mutually inconsistent datasets and datasets with few points but large constraining power vs datasets with many points but moderate constraining power

#### Lepton-Hadron collisions

| Experiment             | Dataset                   | Ref.     | Sys. Unc. |        |     | $N_{dat}$ no cuts | Kinematics                                      |  |  |  |
|------------------------|---------------------------|----------|-----------|--------|-----|-------------------|---|--|--|--|
|                        |                           |          |           |        |     | (NLO/NNLO cuts)   |   |  |  |  |
| NMC                    |                           |          |           |        |     |                   | 1   |  |  |  |
|                        | NMC $d/p$                 | [19]     | A         | full   | 1   | 289(132/132)      | I   |  |  |  |
|                        | NMC ( <sup>NC,p</sup>     | [20]     | Δ         | full   |     | 211 (224/224)     |   |  |  |  |
| SLAC                   | THE U                     | [20]     | 11        | Tun    |     | 211 (224/224)     |   |  |  |  |
| BLAU                   | STAC -                    | [99]     | I A       | none   |     | 101 (27/27)       | I   |  |  |  |
|                        | SLAC p                    | [23]     |           | none   | a   | 191 (37/37)       |   |  |  |  |
| DODMO                  | SLAC d                    | [23]     | A         | none   | a   | 191 (37/37)       |   |  |  |  |
| BCDMS                  | DODMO                     | 10.01    |           |        |     | and (222 (222)    | L   |  |  |  |
|                        | BCDMS $p$                 | [21]     | A         | full   | b   | 351 (333/333)     |   |  |  |  |
|                        | BCDMS d                   | [22]     | A         | full   | b   | 254(248/248)      |   |  |  |  |
| CHORUS                 |                           |          |           |        |     |                   |   |  |  |  |
|                        | CHORUS $\nu$              | [35]     | A         | full   | с   | 572(431/431)      |   |  |  |  |
|                        | CHORUS $\bar{\nu}$        | [35]     | A         | full   | с   | 572(431/431)      |   |  |  |  |
| NuTeV                  |                           |          | -         | -      | -   |                   |   |  |  |  |
|                        | NuTeV $\nu$               | [36, 37] | A         | none   |     | 45 (41/41)        |   |  |  |  |
|                        | NuTeV $\bar{\nu}$         | [36, 37] | A         | none   |     | 44 (38/38)        |   |  |  |  |
| HERA-I                 |                           |          |           |        |     |                   |   |  |  |  |
|                        | HERA-I NC $e^+$           | [24]     | M         | full   | d   | 434 (379/379)     |   |  |  |  |
|                        | HERA-I NC $e^-$           | [24]     | M         | full   | d   | 145 (145/145)     |   |  |  |  |
|                        | HERA-I CC $e^+$           | [24]     | M         | full   | d   | 34 (34/34)        |   |  |  |  |
|                        | HERA-I CC $e^-$           | [24]     | M         | full   | d   | 34(34/34)         |   |  |  |  |
| ZEUS HERA-II           |                           | []       |           |        | -   | (//               |   |  |  |  |
| ZEOS HEIGT H           | ZEUS-II NC e              | [33]     | M         | full   |     | 90 (90/90)        | I   |  |  |  |
|                        | ZEUS-II CC e <sup>-</sup> | [34]     | M         | full   |     | 37 (37/37)        |   |  |  |  |
|                        | ZEUS II NC e <sup>+</sup> | [54]     | M         | full   | f   | 00 (00/00)        | $5 10^{-3} \le m \le 0.40$                      |  |  |  |
|                        | ZEOD-II NO E              | [00]     | 1111      | Iun    | 1   | 30 (30/30)        | $200 < O^2 < 210^4 C_0 V^2$                     |  |  |  |
|                        | ZEUS IL CC at             | [E 4]    | M         | 6.11   | F   | 25 (25/25)        | $200 \le Q \le 310$ GeV                         |  |  |  |
|                        | ZE03-11 CC e              | [04]     | 11/1      | Iun    | 1   | 30 (30/30)        | $1.810 \leq x \leq 0.42$                        |  |  |  |
| III HEDA H             |                           |          |           |        |     |                   | $280 \leq Q \leq 310$ GeV                       |  |  |  |
| HI HERA-II             |                           | (mail    |           |        | i . | 100 (100 (100)    | 0.10-3 < < 0.05                                 |  |  |  |
|                        | HI-II NC e                | [51]     | M         | [ Tull | g   | 139 (139/139)     | $210^{-1} \le x \le 0.65$                       |  |  |  |
|                        |                           |          |           |        |     |                   | $120 \le Q^2 \le 4 \ 10^2 \ \text{GeV}^2$       |  |  |  |
|                        | H1-II NC $e^+$            | [51]     |           | full   | g   | 138(138/138)      | $210^{-3} \le x \le 0.65$                       |  |  |  |
|                        |                           |          |           |        |     |                   | $120 \le Q_{0}^{2} \le 410^{4} \text{ GeV}^{2}$ |  |  |  |
|                        | H1-II CC $e^-$            | [51]     | M         | full   | g   | 29(29/29)         | $810^{-3} \le x \le 0.40$                       |  |  |  |
|                        |                           |          |           |        |     |                   | $300 \le Q^2 \le 3  10^4  \text{GeV}^2$         |  |  |  |
|                        | H1-II CC $e^+$            | [51]     | M         | full   | g   | 29(29/29)         | $8  10^{-3} \le x \le 0.40$                     |  |  |  |
|                        |                           |          |           |        |     |                   | $300 \le Q^2 \le 3  10^4   {\rm GeV^2}$         |  |  |  |
|                        | H1-II low $Q^2$           | [52]     | M         | full   |     | 136(124/124)      | $2.810^{-5} \le x \le 0.015$                    |  |  |  |
|                        |                           |          |           |        |     |                   | $1.5 \le Q^2 \le 90 \text{ GeV}^2$              |  |  |  |
|                        | H1-II high y              | [52]     | Μ         | full   |     | 55 (52/52)        | $2.910^{-5} \le x \le 510^{-3}$                 |  |  |  |
|                        |                           |          |           |        |     |                   | $2.5 \le Q^2 \le 90 \text{ GeV}^2$              |  |  |  |
| HERA $\sigma_{NC}^{c}$ |                           | [55]     | Μ         | full   |     | 52 (47/47)        | $310^{-5} \le x \le 0.05$                       |  |  |  |
|                        |                           |          |           |        |     | × ′ ′             | $2.5 < Q^2 < 2\overline{10}^3 { m GeV}^2$       |  |  |  |
| L                      |                           |          |           |        |     |                   | _ · · _   |  |  |  |

NNPDF3.0 dataset

Hadron-Hadron collisions

| Experiment         | Dataset                  | Ref.      | Sys. Unc. |        |     | $N_{dat}$ no cuts | Kinematics                                      |  |
|--------------------|--------------------------|-----------|-----------|--------|-----|-------------------|---|--|
|                    |                          |           |           |        |     | (NLO/NNLO cuts)   |   |  |
| DY E866            |                          |           |           |        |     |                   |   |  |
|                    | DY E866 $d/p$            | [41]      | Μ         | M none |     | 15(15/15)         |   |  |
|                    | DY E866 $p$              | [39, 40]  | Μ         | none   |     | 184(184/184)      |   |  |
| DY E605            |                          | [38]      | Μ         | none   |     | 119(119/119)      |   |  |
| CDF                |                          |           |           |        |     |                   |   |  |
|                    | CDF Z rap                | [43]      | Μ         | full   | h   | 29(29/29)         |   |  |
|                    | CDF Run-II $k_t$ jets    | [83]      | Μ         | full   | h   | 76 (76/52)        |   |  |
| D0                 |                          |           |           |        |     |                   |   |  |
|                    | D0 $Z$ rap               | [44]      | Μ         | full   |     | 28(28/28)         |   |  |
| ATLAS              |                          |           |           |        |     |                   |   |  |
|                    | ATLAS $W, Z$ 2010        | [47]      | Μ         | full   | i   | 30(30/30)         |   |  |
|                    | ATLAS 7 TeV jets 2010    | [50]      | Μ         | full   | i,j | 90(90/9)          |   |  |
|                    | ATLAS 2.76 TeV jets      | [63]      | Μ         | full   | j   | 59(59/3)          | $20 \le p_T^{\text{jet}} \le 200 \text{ GeV}$   |  |
|                    |                          |           |           |        |     |                   | $0 \leq  \eta^{ m jet}  \leq 4.4$               |  |
|                    | ATLAS high-mass DY       | [56]      | Μ         | full   |     | 11(5/5)           | $116 \le M_{ll} \le 1500 \text{ GeV}$           |  |
|                    | ATLAS $W p_T$            | [57]      | Μ         | full   |     | 11 (9/-)          | $0 \le p_T^W \le 300 \text{ GeV}$               |  |
| CMS                |                          |           |           |        |     | •                 |   |  |
|                    | CMS W electron asy       | [48]      | Μ         | COV    |     | 11 (11/11)        |   |  |
|                    | CMS W muon asy           | [58]      | Μ         | COV    |     | 11 (11/11)        | $0 \le  \eta_l  \le 2.4$                        |  |
|                    | CMS jets 2011            | [62]      | Μ         | full   |     | 133(133/83)       | $114 \le p_T^{\text{jet}} \le 2116 \text{ GeV}$ |  |
|                    |                          |           |           |        |     |                   | $0 \le  \eta^{\text{jet}}  \le 2.5$             |  |
|                    | CMS $W + c$ total        | [60]      | Μ         | COV    |     | 5(5/5)            | $0 \le  \eta_l  \le 2.1$                        |  |
|                    | CMS $W + c$ ratio        | [60]      | Μ         | COV    |     | 5(5/5)            | $0 \le  \eta_l  \le 2.1$                        |  |
|                    | CMS 2D DY 2011           | [59]      | Μ         | COV    |     | 124 (88/110)      | $20 \le M_{ll} \le 1200 \text{ GeV}$            |  |
|                    |                          |           |           |        |     |                   | $0 \le  \eta_{ll}  \le 2.4$                     |  |
| LHCb               |                          |           |           |        |     |                   |   |  |
|                    | LHCb W rapidity          | [49]      | Μ         | COV    |     | 10(10/10)         |   |  |
|                    | LHCb $Z$ rapidity        | [61]      | Μ         | cov    |     | 9(9/9)            | $2.0 \le \eta_l \le 4.5$                        |  |
| $\sigma(t\bar{t})$ |                          |           | -         | -      |     |                   |   |  |
|                    | ATLAS $\sigma(t\bar{t})$ | [65-67]   | Μ         | none   |     | 3(3/3)            | -   |  |
|                    | CMS $\sigma(t\bar{t})$   | [68 - 70] | Μ         | none   |     | 3 (3/3)           | -   |  |
| Total              |                          |           |           |        |     | 5179 (4276/4078)  |   |  |

#### Juan Rojo

13

### Experimental data

- Final provide the second secon
  - **Type of high energy collision** (lepton-proton, proton-proton), **center-of-mass energy** of collision
  - ▶ Whether of not **experimental correlated systematics** are available, and if so, in **which format**
  - Mutually inconsistent datasets and datasets with few points but large constraining power vs datasets with many points but moderate constraining power



#### NNPDF3.0 NLO dataset

The kinematical coverage of the experiments included in NNPDF3.0 **span several orders** of magnitude both in x and Q<sup>2</sup>

14

#### Inconsistent data

- What it is usually meant by *inconsistent data*?
- $\stackrel{\circ}{=}$  Not a unique definition. Typically when one experiment that when added into a global fit leads to  $\chi^2 >> N_{dat}$
- Many possible reasons for this:
  - **Underestimated systematic uncertainties**?
  - **☑** Incomplete / partial **theory calculations**?
  - **Methodological limitations**, ie, too restrictive PDF fitting functional forms?

Genuine statistical pull between different experiments in the global fit? (This is not inconsistency!)

- Dealing with potentially inconsistent experiments in the global PDF fit is very delicate. On the one hand, it is not advisable to bias *a priori* the fit with a subjective selection of which experiments are more *reliable*. On the other hand, once wants to achieve statistically sound results, and in particular PDF uncertainties that truly quantify our genuine lack of information. So there are two complementary avenues:
  - Attempt to **understand how the inconsistencies arise**, and when possible fix them (for example using a better theory)
  - Devise a fitting methodology that can deal with inconsistent experiments, regardless of the origin of the inconsistency
- Note that some of the older fixed-target DIS experiments do not provide the full breakdown of systematics, but this is now a small weight in the global fit

#### Dealing with inconsistent data

- In the global PDF fit, different experiments will prefer different solutions, not always compatible
- Also, the number of datapoints (statistical weight) of each experiment can be quite different, and one wants to describe **both** exps with many data points and those with few
- Using textbook statistics, 68% CL uncertainties in the PDF fit parameters should be determined from the  $\Delta \chi^2 = 1$  criteria
- However it has been shown that this criterion is not adequate in the global fits with many experiments
- So global Hessian fits use effectively an increased tolerance  $\Delta \chi^2 >> 1$ , to ensure that all fitted experiments are reasonably described



#### Dealing with inconsistent data

- In the MSTW approach, a dynamical tolerance criterion is used where different individual experiments determine the allowed upper and lower variations of each eigenvectors
- Final Section 2 This also indicates which datasets are more **sensitive to each eigenvector**



#### MSTW 2008 NLO PDF fit

#### Dealing with inconsistent data

Fin the NNPDF approach, **no need to modify the fitting methodology** in the presence of inconsistencies

When new (compatible) experiments are added, then **PDF errors decrease**. If inconsistent experiments included, fit essentially unaffected and PDF errors not modified (since no new information added)

Fitting methodology also unchanged even for large variations of the fitted dataset



LHC 13 TeV, ac=0.118, MadGraph5\_aMC@NLO fNLO



## **Bayesian Reweighting**

19



*NNPDF, arXiv:1012.0836 NNPDF, arXiv:1108.1758* 

#### Reweighting as an alternative to fitting

When a **new dataset becomes available**, instead of updating the global fit, it is possible include this new information on a prior PDF set using **Bayes' Theorem** 

Free weight (likelihood) in the presence of the new experiment corresponding to each Monte Carlo replica *k* is given in terms of the  $\chi^{2}_{k}$  between data and the theory computed with this replica

$$w_k = \frac{(\chi_k^2)^{\frac{1}{2}(n-1)} e^{-\frac{1}{2}\chi_k^2}}{\frac{1}{N} \sum_{k=1}^N (\chi_k^2)^{\frac{1}{2}(n-1)} e^{-\frac{1}{2}\chi_k^2}}$$

Free Bayesian reweighting technique also allows to quantify the overall consistency of the new experiment with those already included of the global fit by defining

$$\mathcal{P}(\alpha) \propto \frac{1}{\alpha} \sum_{k=1}^{N} w_k(\alpha).$$

where  $\omega_k(\alpha)$  are the **weights**  $\omega_k$  now with  $\chi^2$  **rescaled as**  $\chi^2/\alpha$ , that is, they correspond to the case where the new experimental dataset has **uncertainties rescaled by a factor**  $\alpha^{1/2}$ 

Any inconsistent experiment can be brought in agreement with the global fit by a **suitable rescaling of its uncertainties** (though this is not necessary in the NNPDF framework)

#### Reweighting as an alternative to fitting

Adding new experiments, in this case the **Tevatron inclusive jet data**, by reweighting leads to results that are **statistically consistent with refitting** 

Main benefit of RW is that it can be performed using only public tools (PDF sets and codes for cross-section calculation) **without any input from the PDF fitters** 



21

#### Reweighting as an alternative to fitting

<sup> $\square$ </sup> The distribution of the  $\chi^2$  rescaling parameter  $\alpha$  allows to **quantify the level of (in)consistency of a new experiment** with those already included in the global fit



Juan Rojo

22

#### **Conservative Partons**

- To study the robustness of the global fit results, it is possible to define parton distributions based on a maximally consistent dataset: the conservative partons
- Include in the conservative fit only those experiments which in the global fit have their P(α) distribution peaked at α < α<sub>max</sub>
- modifying this threshold allows to tune the PDF fit to be more or less conservative
- Quantify impact of *known* dataset inconsistencies on the global fit PDFs
- This is not merely a conceptual detail: assessing robustness of PDF errors in LHC cross-section is central *ie* for the *characterisation of the Higgs boson*

|                          | $\alpha_{\max}$   | = 1.1              | $\alpha_{\max}$    | = 1.2              | $\alpha_{\max}$    | = 1.3              | Glol              | bal fit            |
|--------------------------|-------------------|--------------------|--------------------|--------------------|--------------------|--------------------|-------------------|--------------------|
|                          | $\chi^2_{ m nlo}$ | $\chi^2_{ m nnlo}$ | $\chi^2_{\rm nlo}$ | $\chi^2_{ m nnlo}$ | $\chi^2_{\rm nlo}$ | $\chi^2_{ m nnlo}$ | $\chi^2_{ m nlo}$ | $\chi^2_{ m nnlo}$ |
| Total                    | 0.96              | 1.01               | 1.06               | 1.10               | 1.12               | 1.16               | 1.23              | 1.29               |
| NMC $d/p$                | 0.91              | 0.91               | 0.89               | 0.89               | 0.88               | 0.89               | 0.92              | 0.93               |
| NMC $\sigma^{NC,p}$      | -                 | -                  | -                  | -                  | -                  | -                  | 1.63              | 1.52               |
| SLAC                     | -                 | -                  | -                  | -                  | 1.77               | 1.19               | 1.59              | 1.13               |
| BCDMS                    | -                 | -                  | 1.11               | 1.15               | 1.12               | 1.16               | 1.22              | 1.29               |
| CHORUS                   | -                 | -                  | 1.06               | 1.02               | 1.09               | 1.07               | 1.11              | 1.09               |
| NuTeV                    | 0.35              | 0.34               | 0.62               | 0.64               | 0.70               | 0.70               | 0.70              | 0.86               |
| HERA-I                   | 0.97              | 0.98               | 1.02               | 1.00               | 1.02               | 0.99               | 1.05              | 1.04               |
| ZEUS HERA-II             | -                 | -                  | -                  | -                  | 1.41               | 1.48               | 1.40              | 1.48               |
| H1 HERA-II               | -                 | -                  | -                  | -                  | -                  | -                  | 1.65              | 1.79               |
| HERA $\sigma_{\rm NC}^c$ | -                 | -                  | 1.21               | 1.32               | 1.20               | 1.31               | 1.27              | 1.28               |
| E886 $d/p$               | 0.30              | 0.30               | 0.43               | 0.40               | 0.44               | 0.46               | 0.53              | 0.48               |
| E886 $p$                 | -                 | -                  | 1.18               | 1.40               | 1.27               | 1.53               | 1.19              | 1.55               |
| E605                     | 1.04              | 1.10               | 0.74               | 0.83               | 0.75               | 0.88               | 0.78              | 0.90               |
| CDF Z rapidity           | -                 | -                  | -                  | -                  | -                  | -                  | 1.33              | 1.53               |
| CDF Run-II $k_t$ jets    | -                 | -                  | 1.01               | 2.01               | 1.04               | 1.84               | 0.96              | 1.80               |
| D0 $Z$ rapidity          | 0.56              | 0.61               | 0.62               | 0.71               | 0.60               | 0.69               | 0.57              | 0.61               |
| ATLAS $W, Z$ 2010        | -                 | -                  | 1.19               | 1.13               | 1.19               | 1.17               | 1.19              | 1.23               |
| ATLAS 7 TeV jets 2010    | 0.96              | 1.65               | 1.08               | 1.58               | 1.10               | 1.54               | 1.07              | 1.36               |
| ATLAS 2.76 TeV jets      | 1.03              | 0.38               | 1.38               | 0.36               | 1.35               | 0.35               | 1.29              | 0.33               |
| ATLAS high-mass DY       | -                 | -                  | -                  | -                  | -                  | -                  | 2.06              | 1.45               |
| ATLAS $W p_T$            | -                 | -                  | -                  | -                  | -                  | -                  | 1.13              | -                  |
| CMS W electron asy       | 0.98              | 0.84               | 0.82               | 0.72               | 0.85               | 0.73               | 0.87              | 0.73               |
| CMS W muon asy           | -                 | -                  | -                  | -                  | -                  | -                  | 1.81              | 1.72               |
| CMS jets 2011            | 0.90              | 2.09               | 0.96               | 2.09               | 0.99               | 2.10               | 0.96              | 1.90               |
| CMS W + c total          | -                 | -                  | -                  | -                  | -                  | -                  | 0.96              | 0.84               |
| CMS $W + c$ ratio        | -                 | -                  | -                  | -                  | -                  | -                  | 2.02              | 1.77               |
| CMS 2D DY 2011           | -                 | -                  | -                  | -                  | 1.20               | 1.30               | 1.23              | 1.36               |
| LHCb $W$ rapidity        | -                 | -                  | 0.69               | 0.65               | 0.74               | 0.69               | 0.71              | 0.72               |
| LHCb $Z$ rapidity        | -                 | -                  | 1.23               | 1.78               | 1.11               | 1.58               | 1.10              | 1.59               |
| $\sigma(t\bar{t})$       | _                 | -                  | _                  | -                  | -                  | -                  | 1.43              | 0.66               |

#### **Conservative Partons**

- At the level of LHC cross-sections, **conservative PDFs consistent with global fit PDFs** within uncertainties
- Conservative PDFs affected by larger uncertainties due to reduced dataset
- Non-trivial validation of the robustness of the global fit results

![](_page_23_Figure_4.jpeg)

![](_page_24_Figure_0.jpeg)

# Statistical Methodology Validation with Closure Testing

![](_page_24_Figure_2.jpeg)

Juan Rojo

**PDF uncertainties** have been often criticised by a **potential lack of statistical interpretation** 

Within NNPDF, we performed a systematic **closure tests analysis** based on pseudo-data, and verified that **PDF uncertainties** exhibit **a statistically robust behaviour** 

![](_page_25_Figure_3.jpeg)

For instance, if the **pseudo-data is generated without statistical fluctuations** (that is, identical to the input theory) then the agreement with theory by construction should become **arbitrarily good** 

And indeed it does: as the minimization advances, the  $\chi^2$  decreases monotonically, and the PDF uncertainties as well are reduced, as the fitted theory collapses to the underlying law

![](_page_26_Figure_3.jpeg)

 $\varphi_{\chi^2} \equiv \sqrt{\langle \chi^2[\mathcal{T}[f_{\text{fit}}], \mathcal{D}_0] \rangle - \chi^2[\langle \mathcal{T}[f_{\text{fit}}] \rangle, \mathcal{D}_0]}$ 

Measure of **PDF uncertainties** in units of **data uncertainties** 

Another important advantage of closure testing the global PDF analysis is that it allows to **disentangle the various components of the total PDF uncertainty** 

![](_page_27_Figure_2.jpeg)

In the closure tests, it is possible to validate new techniques, such as the Bayesian reweighting, in a **clean environment where everything is under control** (free in particular of potential data inconsistencies)

![](_page_28_Figure_2.jpeg)

Closure testing the global fit allows disentangling **methodological issues of principle** (in an ideal world with perfectly consistent datasets, does my fitting methodology give the result it should?) with those of **practice** (how to deal with inconsistent experiments or with incomplete theory?)

#### Adding artificial inconsistencies

To test the fitting methodology in a realistic situation, it is possible to **generate pseudo-data adding artificial inconsistencies** and study how the resulting PDFs are modified.

In the MMHT approach, adding artificial inconsistencies in a closure test leads to **modified PDFs in most** cases in agreement with the global fit PDF uncertainties

This is only the case if their **dynamical tolerance criterion**  $\Delta \chi^2 >> 1$  is used, as opposed to  $\Delta \chi^2 = 1$ 

![](_page_29_Figure_4.jpeg)

#### Summary and outlook

- Parton Distributions are an essential ingredient for LHC phenomenology
- At the LHC, precision PDFs are required for many analysis from the characterisation of the Higgs sector to BSM searches and Monte Carlo event generators
- The global QCD analysis aims to extract parton distributions from a diverse experimental dataset using state-of-the-art theory and methodology
- This involves having to deal with several non-trivial statistical issues, in particular with **potential inconsistencies between fitted datasets**, that can arise from various sources: partial theory, limited fitted methodology or underestimated systematic uncertainties
- To deal with these problems, a number of techniques have been developed, which allow to validate the robustness of our PDF uncertainty estimates for high-precision LHC phenomenology