

# Proton structure at the LHC

*Nathan Hartland*



Doctor of Philosophy  
University of Edinburgh  
2014

*To my parents*

# Declaration

This thesis is the sole work of myself, and is a record of work performed as part of the NNPDF collaboration, and a separate collaboration with Steffen Schumann and my supervisor Luigi Del Debbio. Unless explicitly stated, the results presented in this thesis are either mine, or the product of collaboration in which I have made a significant contribution.

In chapters ?? and ?? the presented results are based upon work published in

- R. D. Ball *et al.*, Nucl. Phys. B **855** (2012) 608 [arXiv:1108.1758 [hep-ph]].
- L. Del Debbio *et al.*, Comput. Phys. Commun. **185** (2014) 2115 [arXiv:1312.4460].
- R. D. Ball *et al.*, JHEP **1304** (2013) 125 [arXiv:1211.5142 [hep-ph]].
- R. D. Ball *et al.*, Nucl. Phys. B **867** (2013) 244 [arXiv:1207.1303 [hep-ph]].

The discussion in Chapter ?? includes results published in

- N. P. Hartland *et al.*, PoS DIS **2013** (2013) 043 [arXiv:1307.2046 [hep-ph]].

My work presented in this thesis has not been used as the basis of any other professional qualification or degree.

---

Nathan Hartland

# Acknowledgements

After almost four years of studying in Edinburgh I find that there are many people I would like to thank for making my time here so memorable. First and foremost I would like to thank my supervisor Luigi Del Debbio for his advice and expertise, along with his encouragement and interest in my progress as physicist. I wouldn't have made it this far without his guidance and direction.

For providing such a great environment in which to learn, I'd like to thank the staff of the Edinburgh PPT group. Particular thanks go to Richard Ball, Arjun Berera, Einan Gardi, and Brian Pendleton for their consistently useful advice, for which I am greatly appreciative. Special thanks go to my fellow students, without whom the experience wouldn't have been anywhere near as interesting. I'd particularly like to thank the students of the 2010 contingent, Ashley Cooke, Tom Metcalf and Samuel Abreu for their companionship and humour. Outside of the Edinburgh physics group, I'd like to thank Xin He and James Kelly for many long and interesting discussions.

During my PhD I've had the pleasure of working within the NNPDF collaboration. From the collaboration I have received an incredible amount of support, learnt a great deal, and made many good friends. Special thanks go to Chris Deans who shared the experience with me as an NNPDF student in Edinburgh.

During my third year, Steffen Schumann provided me with the opportunity to work and study at the University of Göttingen, where I gained valuable experience outside the main topic of my PhD. I'm very grateful for the opportunity, and

I'd like to thank Steffen, along with Erik Gerwick, for making the visit such a welcoming and enjoyable experience.

For providing me with the chance to study at Edinburgh and Göttingen, thanks go to the STFC and the MCnet initial training network for their financial support.

Finally, thanks go to my girlfriend Sophie, and my family, particularly my parents who have supported me throughout my life and studies so far. This thesis is dedicated to them.

# Abstract

A determination of Parton Distribution Functions (PDFs) from a global fit to a dataset including measurements from the LHC has been performed for the first time. The determinations have been performed according to the NNPDF methodology, leading to a fit relatively free of parametrisation bias and with an accurate account of PDF uncertainty.

In this thesis the importance of QCD measurements at the LHC to PDF extraction are discussed, and we summarise some of the technical difficulties in their inclusion into PDF fits. A number of methods are presented that permit the efficient inclusion of these observables into PDF determinations.

Firstly a Bayesian reweighting procedure taking advantage of the Monte Carlo representation of PDF uncertainties in NNPDF sets is discussed. The utility of the Bayesian reweighting method is demonstrated by a study of the impact of early  $W$  production asymmetry measurements from ATLAS, CMS and LHCb upon an earlier PDF set.

A package for the fast computation of observables in an automated NLO framework is presented, providing an interface between Monte Carlo event generators and NLO interpolation tools.

Finally, a new method of combining PDF evolution with interpolating codes for hadronic observable computation is also described. This method largely overcomes the computational difficulties in performing fast perturbative QCD predictions for collider observables. The method has been applied to the

determination of PDFs from a global dataset including electroweak vector boson production data from LHCb, ATLAS and CMS along with inclusive jet data from ATLAS. The resulting set, NNPDF2.3 provides the most accurate determination of parton distributions via the NNPDF methodology to date.

Finally, the method of closure testing is introduced, and the method is applied to the study of the NNPDF methodology. A number of improvements are found in the minimisation and stopping procedures, which are adopted for the development of the next NNPDF release, NNPDF3.0. Alongside the sounder methodological basis, the NNPDF3.0 PDF set will provide a determination based upon an expanded dataset in order to produce a comprehensive upgrade to the NNPDF2.3 family of fits.

# Contents

# List of Figures

# List of Tables

# Introduction

The study of elementary particles and their behaviour relies on a great many sources of experimental information. In order to verify the predictions of the Standard Model (SM) of particle physics or indeed extensions beyond, precise and accurate measurements must be made of the fundamental properties of matter. Building upon decades of advances in the study of elementary particles, today the foremost source of cutting edge measurements in particle physics is the Large Hadron Collider (LHC) based at CERN in Switzerland. The LHC, through colossal scientific and human effort has opened up the study of the properties of nature to scales that were previously inaccessible.

The LHC probes the building blocks of nature by the collision of high energy protons. Maximising the physics potential of the LHC therefore requires a deep understanding of the composite nature of the proton. The short range dynamics of a proton's constituent particles can be described by perturbative Quantum Chromo-Dynamics (QCD), however an understanding of the low energy behaviour is impossible to obtain through perturbative methods, therefore making its determination by a calculation from first principles challenging. In practice the structure of the proton is understood through a comprehensive analysis of experimental data, and described in terms of Parton Distribution Functions (PDFs). These functions parametrise the unknown non-perturbative dynamics of the proton. As a universal property of protons, the PDFs may be determined from available experimental data and then applied in the calculation of predictions for

other experiments, therefore making the application of QCD in hadron collisions into a predictive theory which may be tested via comparison to data.

The accurate determination of parton densities in the proton is an ongoing effort, with several groups providing competing sets of parton distributions. The NNPDF collaboration provides a set of parton distributions determined through a rather different methodology than the standard procedures, resulting in a PDF set suffering from little of the parametrisation bias possible in competing approaches. Furthermore the NNPDF methodology has a unique treatment of the experimental uncertainty propagation, leading to a statistically sound estimation of the uncertainties in the resulting PDFs.

While a precise knowledge of the dynamics of the proton is vital for LHC studies of physics in the standard model and beyond, LHC data also has the potential to provide the most in depth information on parton densities to date. This thesis is based upon work conducted in the study of early LHC standard model measurements of particular sensitivity to parton distributions. The inclusion of such an experimental dataset into a fit in the NNPDF framework has necessitated the development of a number of tools for the efficient calculation of collider observables. These tools and their applications shall be discussed alongside the development of the NNPDF methodology to better handle the ever-enlarging LHC dataset.

This thesis is arranged as so. In Chapter One we shall provide a brief discussion of the theoretical structure of parton distributions, where they arise in the calculation of deep-inelastic scattering cross-sections and further theoretical background relevant to the reliable determination of PDFs from experimental data. Chapter Two is concerned with the practical extraction of PDFs and shall describe experimental observables of interest along with the different approaches used by major PDF collaborations to fit the data. In Chapter Three, the tools that have been developed to enable the inclusion of a large LHC dataset into

a computationally intensive fit such as the NNPDF procedure are introduced and described. A brief summary of experimental measurements at the LHC of interest to the determination of PDFs is provided in Chapter Four. In Chapter Five, we shall examine the impact of some of these measurements made by LHC collaborations upon PDF determinations, enabled by the tools developed in Chapter Three. The data impact will be assessed in the context of the two most recent public releases of the NNPDF collaboration; providing a summary of their datasets and the tools used in their extraction. Finally in Chapter Six we examine some of the methodological improvements that have been made in the NNPDF procedure in order to ensure the maximal efficiency in extracting new information on PDFs from future LHC measurements, and demonstrate their application in early prototypes of the NNPDF3.0 PDF set.

# Chapter 1

## Parton distribution functions

Parton distributions are one of the central pillars of perturbative QCD, factorising as they do the perturbatively incalculable long distance dynamics present in calculations involving hadronic initial states. Combined with the perturbative description of the short-distance cross-section what could seem at first a hopeless situation is alleviated, and QCD becomes a predictive and useful theory when applied to hadronic scattering.

In this chapter a brief overview of how parton distribution functions arise in QCD calculations will be presented. We shall explore the prototypical example of the deep inelastic scattering (DIS) of leptons off a hadronic target, first in the *naive parton model* arising before the advent of QCD and then with the QCD-improved parton model which allows for an excellent description of DIS measurements across a wide range of hard scales.

The treatment of heavy quarks in parton distributions is a particularly delicate issue and therefore will also be discussed in this introductory theory section. Finally there will be some exploration of the general properties of parton distributions in order to provide a summary of the available theoretical constraints upon PDFs.

## 1.1 Partons in deep inelastic scattering

We shall begin by introducing parton distribution functions as they arise in the early parton model. The model was originally introduced by Feynman and Bjorken [?, ?, ?, ?] in the late 1960's in an effort to understand the scattering behaviour of hadronic states and successfully describes many properties observed in early deep inelastic scattering experiments.

In this process, a charged lepton  $l$  probes a proton  $P$  by the exchange of a gauge boson. For simplicity we shall describe here the neutral current process where a photon is exchanged. In the inelastic regime where the momentum transfer to the target proton is large, the proton does not survive the scattering process and fragments into an arbitrary hadronic final state  $X$ . The process  $l(k) + P(p) \rightarrow l(k') + X$  is illustrated at tree level in Figure ??.

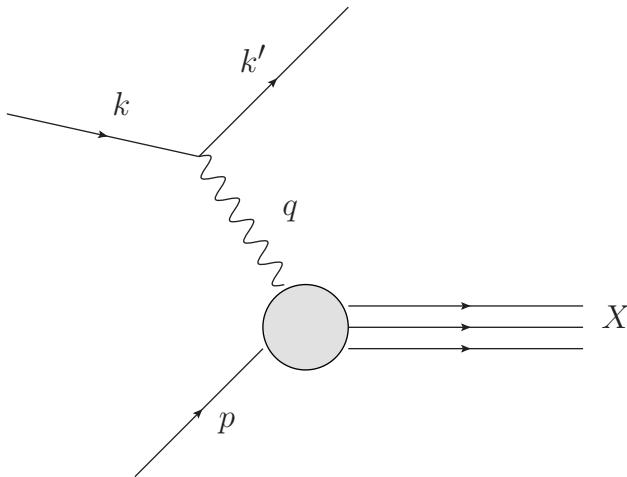


Figure 1.1: Deep inelastic scattering of a charged lepton with a proton target.

In this system we can define the standard DIS kinematic variables;  $Q^2$  denotes the momentum transfer from the electron to the target proton,  $\nu$  the energy transfer and  $y$  the measure of the reaction's inelasticity, or fractional energy

transfer. In the rest frame of the proton these are given by

$$Q^2 = -q^2 = -(k - k')^2, \quad (1.1)$$

$$\nu = M(E - E'), \quad (1.2)$$

$$y = (q \cdot p)/(k \cdot p), \quad (1.3)$$

where  $M$  refers to the mass of the proton, and the inelasticity ranges between 0 (elastic scattering) and 1.  $E$  and  $E'$  denote the energies associated with the four-momenta  $k$  and  $k'$  respectively. Additionally, we may introduce the Bjorken scaling parameter  $x$ , central to the parton model,

$$x = \frac{Q^2}{2\nu}. \quad (1.4)$$

Neglecting spin labels, the amplitude for this diagram in the Feynman gauge is given by

$$\mathcal{M} = ie^2 \bar{u}(k') \gamma^\mu u(k) \left( i \frac{g_{\mu\nu}}{Q^2} \right) \langle X | J_h^\nu | P \rangle, \quad (1.5)$$

where  $J_h^\nu$  represents the hadronic current. The fundamental difficulty in attempting to compute the cross section for this process is our ignorance of the wavefunction for the hadronic states  $|X\rangle$  and  $|P\rangle$ . To isolate the problem, we are able to factorise the spin averaged square of the amplitude in Equation ?? into a leptonic ( $L_{\mu\nu}$ ) and a hadronic ( $W^{\mu\nu}$ ) part

$$|\overline{\mathcal{M}}|^2 = \frac{1}{Q^2} L_{\mu\nu} W^{\mu\nu}, \quad (1.6)$$

where the leptonic tensor is straightforwardly calculable:

$$L_{\mu\nu} = e^2 \sum_{\text{spin}} \bar{u}(k') \gamma_\mu u(k) \bar{u}(k) \gamma_\nu u(k'), \quad (1.7)$$

$$= e^2 \text{tr} [\not{k'} \gamma_\mu \not{k} \gamma_\nu], \quad (1.8)$$

$$= 4e^2 [k_\mu k'_\nu + k_\nu k'_\mu - g_{\mu\nu} k \cdot k'], \quad (1.9)$$

where here we have neglected the fermion masses. The hadronic part of the calculation is considerably more difficult to evaluate, and indeed impossible to compute from first principles in perturbation theory as it is sensitive to the low-scale, and therefore strongly coupled dynamics of the proton target:

$$W^{\mu\nu} \sim \sum_X \langle P(p) | J_h^{\mu\dagger} | X \rangle \langle X | J_h^\nu | P(p) \rangle, \quad (1.10)$$

$$\sim \langle P(p) | J_h^{\mu\dagger} J_h^\nu | P(p) \rangle. \quad (1.11)$$

However, we can gain some insight into its structure by noting that the tensor must obey the conservation requirements of the hadronic current  $q_\mu W^{\mu\nu} = 0$  and  $q_\nu W^{\mu\nu} = 0$ . The tensor may therefore be parametrised without loss of generality by the following structure:

$$W_{\mu\nu} = - \left( g_{\mu\nu} - \frac{q_\mu q_\nu}{q^2} \right) F_1(x, Q^2) + \left( p_\mu - q_\mu \frac{p \cdot q}{q^2} \right) \left( p_\nu - q_\nu \frac{p \cdot q}{q^2} \right) \frac{1}{\nu} F_2(x, Q^2). \quad (1.12)$$

Here we have introduced the parameters in our tensor  $F_i$  which are known as the electromagnetic structure functions. For interactions involving parity-violating currents, there is a third contribution to the hadronic tensor arising through the  $F_3$  structure function. Here the only possible functional dependence for the structure functions is upon the quantities  $Q^2$  and  $x$ .

It is convenient now to define a projection vector  $n$  with the properties  $p \cdot n = 1$ ,  $n \cdot q = 0$ , and  $n^2 = p^2 = 0$ , where the assumption of negligible proton mass

has been made. Any vector can now be written as a combination of  $n$ ,  $p$  and a component transverse to the proton momentum as a *Sudakov decomposition*. Using this projection vector we may obtain the structure functions from the hadronic tensor as so:

$$F_2 = \nu n^\mu n^\nu W_{\mu\nu}, \quad (1.13)$$

$$F_L = F_2 - 2x F_1 = \frac{Q^4}{\nu^3} p^\mu p^\nu W_{\mu\nu}, \quad (1.14)$$

where the quantity in the second equation is known as the longitudinal structure function. So far, few assumptions have been made about the form of the EM hadronic tensor  $W_{\mu\nu}$ , we have simply parametrised it in terms of a Lorentz invariant tensor structure and structure functions. Feynman's parton model allows us to describe more of the hadronic tensor with perturbation theory by proposing a composite proton formed as a bound state of fundamental, spin-1/2 constituents: the *partons*.

The parton model approximation states that for a sufficiently hard interaction, the virtual photon only interacts with a single point-like parton inside the target proton and we can treat the partons as approximately free particles. The hadronic tensor then admits a probabilistic expansion in terms of Parton Distributions which encode the probability of the hard photon interacting with a constituent parton carrying a fraction  $\xi$  of the parent proton's momentum. The probability of interacting with a parton carrying between  $\xi$  and  $\xi + \delta\xi$  of the proton's momentum being given by  $f(\xi)\delta\xi$  where  $f(\xi)$  is the interaction probability for a parton with momentum  $\xi p$ . Diagrammatically we may therefore construct the photon-hadron interaction as a weighted sum of partonic diagrams:

$$\left| \begin{array}{c} \text{wavy line} \\ \text{--- circle} \\ \text{--- arrow} \\ \text{--- line labeled } p \end{array} \right|^2 = \sum_i^{N_{part}} f_i(\xi, Q^2) \otimes \left| \begin{array}{c} \text{wavy line} \\ \text{--- arrow} \\ \text{--- line labeled } \xi p \end{array} \right|^2 (\xi),$$

where we have introduced the multiplicative convolution

$$(f \otimes g)(x) = \int_0^1 \frac{d\xi}{\xi} f\left(\frac{\xi}{x}\right) g(\xi). \quad (1.15)$$

The hadronic tensor is then given in terms of a sum of individual hard scattering partonic tensors, denoted  $\widetilde{W}_{\mu\nu}^i(\xi)$  for a target parton of type  $i$ . Writing the hadronic tensor as the probabilistic sum over all constituent parton types we obtain

$$W_{\mu\nu} = \int_0^1 \frac{d\xi}{\xi} \sum_i f_i(\xi, Q^2) \widetilde{W}_{\mu\nu}^i(\xi, Q^2). \quad (1.16)$$

As the parton level tensors must obey the same conservation relations as the full hadronic tensor, we can once again form a general parameterization of  $\widetilde{W}_{\mu\nu}^i$ :

$$\widetilde{W}_{\mu\nu}^i = - \left( g_{\mu\nu} - \frac{q_\mu q_\nu}{q^2} \right) \widetilde{F}_1^i(\xi, Q^2) + \xi^2 \left( p_\mu - q_\mu \frac{p \cdot q}{q^2} \right) \left( p_\nu - q_\nu \frac{p \cdot q}{q^2} \right) \widetilde{F}_2^i(\xi, Q^2), \quad (1.17)$$

where the factors of  $\xi^2$  arise from taking  $p^\mu \rightarrow \xi p^\mu$ . Substituting this form for  $\widetilde{W}_{\mu\nu}^i(\xi)$  into Eqn ?? and comparing with the form in Eqn ??, we find two expressions for the proton EM structure functions,

$$F_1(x, Q^2) = \int_0^1 \frac{d\xi}{\xi} \sum_i f_i(\xi) \widetilde{F}_1^i(\xi, Q^2), \quad (1.18)$$

$$F_2(x, Q^2) = \int_0^1 \xi d\xi \sum_i f_i(\xi) \widetilde{F}_2^i(\xi, Q^2). \quad (1.19)$$

The naive parton level structure functions  $\widetilde{F}_1^i(\xi, Q^2)$  describe the hard scattering subprocess involving a parton of species  $i$  and may be computed by considering the parton level squared amplitude for the subprocess,  $\gamma^*(q) + q(\xi p) \rightarrow q(l)$  and projecting out the desired quantities with the operators defined previously. At

leading order, using the parton level version of the projector Eqn ??:

$$\mathcal{M}_\mu = -ie_{q^i}\bar{u}(l)\gamma^\mu u(\xi p), \quad (1.20)$$

$$\frac{n^\mu n^\nu}{\xi^2} \widetilde{W}_{\mu\nu}^i = \frac{n^\mu n^\nu}{\xi^2} \overline{\sum} |\mathcal{M}|_{\mu\nu}^2 = 4e_{q^i}^2, \quad (1.21)$$

where we have made the approximation that momenta transverse to the beam axis vanish. Including the phase space for the final state quark in the CM frame we obtain:

$$\widetilde{F}_2^i = 2e_{q^i}^2 \delta(l^2), \quad (1.22)$$

where the delta function can be rewritten in terms of  $\xi p$  and  $q$ :

$$\delta(l^2) = \delta((\xi p + q)^2) = \delta(2\xi\nu - Q^2) = \delta(2\nu(\xi - x)). \quad (1.23)$$

This is an interesting result of the analysis at leading order, the kinematical variable  $x$  actually describes the momentum fraction of the interacting parton. The parton level structure function  $\widetilde{F}_2^i$  is therefore given by:

$$\widetilde{F}_2^i = 2e_{q^i}^2 \delta(\xi - x). \quad (1.24)$$

The parton level longitudinal structure function is also straightforwardly projected out of the same amplitude,

$$\widetilde{F}_L^i = \frac{Q^4}{\xi\nu^3} p^\mu p^\nu \widetilde{W}_{\mu\nu}^i = \widetilde{F}_2^i - \frac{2x}{\xi^2} \widetilde{F}_1^i. \quad (1.25)$$

At leading order this projection, and therefore the longitudinal structure function, are exactly zero, consequently

$$\widetilde{F}_1^i = \frac{\xi^2}{2x} \widetilde{F}_2^i = e_{q^i}^2 \frac{\xi^2}{x} \delta(\xi - x). \quad (1.26)$$

We may therefore write the full EM proton structure functions in the naive parton model as

$$F_1(x, Q^2) = \int_0^1 d\xi \sum_i f_i(\xi) e_{q^i}^2 \frac{\xi}{x} \delta(\xi - x) = \sum_i f_i(x) e_{q^i}^2, \quad (1.27)$$

$$F_2(x, Q^2) = 2 \int_0^1 \xi d\xi \sum_i f_i(\xi) e_{q^i}^2 \delta(\xi - x) = 2x \sum_i f_i(x) e_{q^i}^2. \quad (1.28)$$

These results have a number of important features. Firstly in this model the structure functions have no dependence upon the resolution parameter  $Q^2$ , a phenomenon known as Bjorken scaling [?]. This scaling effect was an important achievement of the original parton model, as it was able to describe contemporary experimental results rather well. The lack of any scale dependence in the structure functions is a consequence of the model's assumptions treating interactions with the proton's constituent partons as point like, and consequently having no characteristic length scale.

Secondly we note that  $F_2(x) = 2xF_1(x)$ , which is known as the Callan-Gross relation [?]. It illustrates a fundamental property of spin-1/2 particles, that they are unable to absorb a longitudinally polarised photon [?].

## 1.2 QCD and the parton model

The naive parton model was able to provide a good phenomenological description of early DIS measurements. Its success also provided great support for QCD as the correct description of the strong interaction. The phenomenon of Bjorken scaling placed substantial constraints upon the theory governing the internal dynamics of the proton. The asymptotic freedom of QCD allows for a consistent description of Bjorken-scaling, where the constituents of the hadron can be viewed as independent, non-interacting point like particles at high values of the resolution parameter  $Q^2$ . The partons in Feynman's model were therefore quickly associated

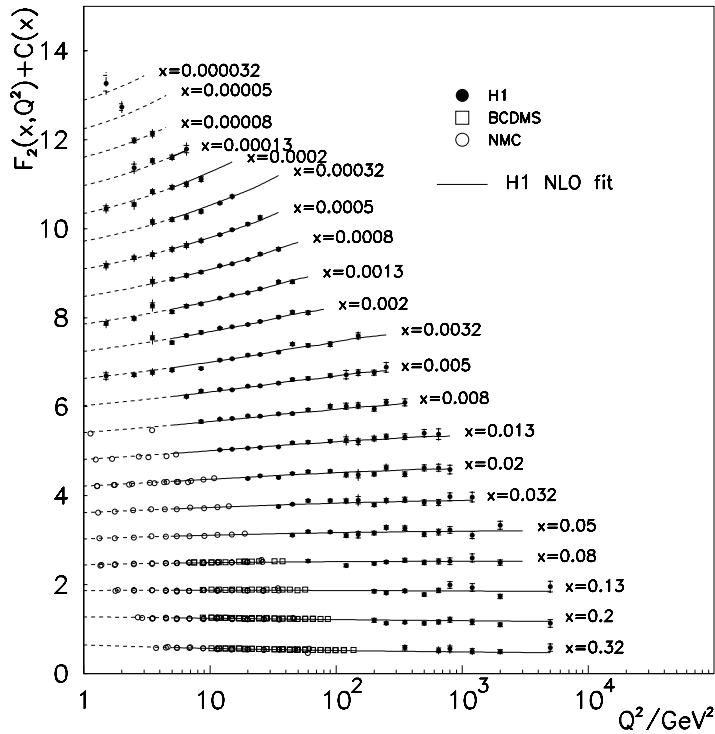


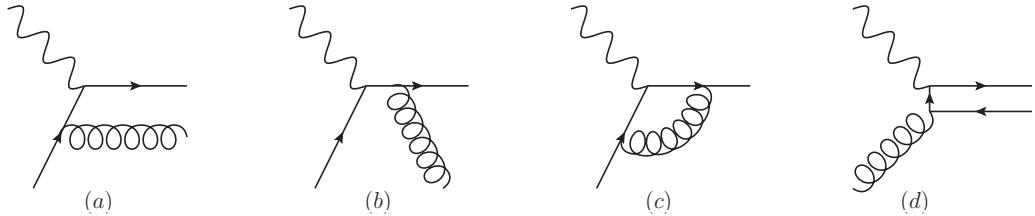
Figure 1.2: Scaling violations in the proton structure function  $F_2$ . Here each curve in  $x$  is scaled by a function  $C(x) = 0.6(i - 0.4)$  for presentation purposes, where  $i$  denotes the bin in  $x$ . Figure from [?].

with the quarks and gluons of QCD.

Despite the ‘snapshot’ picture of non-interacting partons at leading order in QCD, we cannot neglect the higher order corrections to the point vertex calculated in the previous section. These corrections introduce logarithms of  $Q^2$  which break the naive Bjorken scaling of the structure functions. Indeed, the measurement of such scaling violations provided one of the most powerful experimental verifications of QCD. Such violations are demonstrated in measurements of  $F_2$  in Figure ???. In this section we shall perform an overview of the extension of the parton model to  $\mathcal{O}(\alpha_s)$  in QCD.

At one loop order, there are three diagrams that contribute to the  $qq\gamma$  vertex studied in the previous section; the real emission of a gluon from the initial (a)

or final state (b) quarks, and the virtual correction diagram (c). Additionally at one loop order in QCD there arises a diagram initiated by a gluon splitting into a  $q\bar{q}$  pair (d).

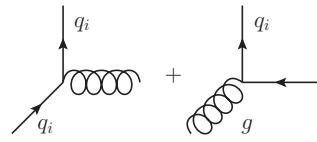


All four of these diagrams are separately divergent. When appropriately regularised however, the divergences in the final state real emission and virtual correction diagrams cancel explicitly as a consequence of the IR safety of QCD, yielding a finite contribution to the cross section. However the divergences present in the real emission diagrams from the initial state partons are not subject to the same cancellations, as they modify the momenta at the interaction vertex.

Like the real emission diagram of a gluon from an initial state quark, the initial state gluon diagram (d) suffers from an equivalent divergence mediated by a perturbatively calculable  $g \rightarrow q\bar{q}$  splitting function  $P_{gq}$ . Including all of the finite contributions from the other contributing diagrams as the coefficient  $W(x)$ , the parton level structure function at next to leading order in QCD is given by

$$\begin{aligned} \widetilde{F}_2^i(\xi, Q^2) &= 2e_i^2 [\delta(\xi - x) \\ &+ \frac{\alpha_S}{2\pi} \sum_j \left( P_{ij}(\xi) \log \frac{Q^2}{\kappa^2} + W_{ij}(\xi) \right) \\ &+ \mathcal{O}(\alpha_S^2)] . \end{aligned} \quad (1.29)$$

Here the  $i$  once again refers to the partonic species at the interaction vertex, and we have introduced an infrared cutoff  $\kappa$  to regulate the parton splitting. The sum over splitting functions arises from the multiple contributions from partonic species  $j$  splitting to  $i$ :



The splitting functions  $P_{ij}$  were known for some time at leading and next-to-leading accuracy [?, ?, ?, ?, ?, ?, ?, ?, ?, ?, ?], and more recently extended to next-next-to-leading order accuracy [?, ?]. After convoluting the parton level functions with the PDFs, we obtain the full structure function

$$\begin{aligned} F_2(x, Q^2) &= \sum_i x e_i^2 [ f_i(x) \\ &+ \frac{\alpha_S}{2\pi} \int_0^1 \frac{d\xi}{\xi} \sum_j \left( P_{ij} \left( \frac{x}{\xi} \right) \log \frac{Q^2}{\kappa^2} + W_{ij}(x) \right) f_j(\xi) \\ &+ \mathcal{O}(\alpha_S^2)] . \end{aligned} \quad (1.30)$$

Our expression for the parton level structure function still suffers from the IR divergence when we take the limit  $\kappa \rightarrow 0$ . This issue may be resolved by concluding that the singularity arises from a breakdown of the ability of perturbation theory to describe physics in the strongly-coupled infrared. We may therefore attempt to factorise out the long distance behaviour of the structure functions into some bare parameters of the theory; analogously to the treatment of ultraviolet divergences by renormalisation of the strong coupling. In this instance we shall absorb the divergences present in the parton level structure functions into our parton distribution functions by replacing the bare quantities  $f(x)$  with a physically accessible quantity measured at the *factorisation scale*  $\mu_f$ . We can express these in terms of an expansion in the bare PDFs as

$$f_i(x, \mu_F^2) = f_i(x) + \frac{\alpha_S}{2\pi} \int_0^1 \frac{d\xi}{\xi} \Delta_{ij}^{(1)} \left( \frac{x}{\xi}, \frac{\mu_F}{\kappa} \right) f_j(\xi) + \mathcal{O}(\alpha_S^2), \quad (1.31)$$

where the counter terms  $\Delta_{ij}^{(n)}$  are formed as a sum of a regular part  $\Delta_{r,ij}^{(n)}$  and

a singular part  $\Delta_{s,ij}^{(n)}$ , and the sum over the dummy index  $j$  is implicit. The singular part of these counterterms is uniquely specified by having to remove the divergence present in the structure functions due to the collinearly divergent parton splitting. Comparing to Eqn. ??, this divergence may be subtracted by setting

$$\Delta_{s,ij}^{(1)} = P_{ij} \left( \frac{x}{\xi} \right) \log \frac{\mu_F^2}{\kappa^2}. \quad (1.32)$$

Unlike the divergent part, the regular part of the counter-term is not uniquely defined by the factorisation procedure. The choice of a specific regular counter-term is known as a *factorisation scheme*; a choice consisting of shuffling terms between the regular part of the PDF definition and the coefficients present in the calculation. For example one may make a process-specific choice where all of the regular coefficients are absorbed into the PDF definition. In our example case of  $F_2$  this is known as the DIS scheme [?],  $\Delta_{r,ij}^{(1)} = W_{ij}(x)$ , in terms of which the form of the calculation becomes particularly simple:

$$F_2(x, Q^2) = 2 \int_0^1 \xi d\xi \sum_i f_i^{\text{DIS}}(\xi) e_i^2. \quad (1.33)$$

In practice this scheme choice is often rather unhelpful, as it does not permit a consistent definition of PDFs across multiple processes. With this in mind, the most common choice is the *Modified Minimal Subtraction* or  $\overline{\text{MS}}$  scheme where the only regular counterterms are a process independent  $\Delta_{r,ij}^{(1)} = \log 4\pi - \gamma_E$ . In the  $\overline{\text{MS}}$  scheme therefore our factorised PDFs are given by

$$f_i(x, \mu_F^2) = f_i(x) + \frac{\alpha_S}{2\pi} \sum_j \left[ \left( P_{ij}(x) \log \frac{\mu_F^2}{\kappa^2} + \log 4\pi - \gamma_E \right) \right] \otimes f_j(x) + \mathcal{O}(\alpha_S^2), \quad (1.34)$$

and the expression for  $F_2$  becomes

$$F_2(x, Q^2) = x \sum_i e_i^2 \left\{ f_i(x, \mu_F^2) + \frac{\alpha_S}{2\pi} \int_x^1 \frac{d\xi}{\xi} f_i(\xi, \mu_F^2) \widetilde{W}_i \left( \frac{x}{\xi}, \frac{Q^2}{\mu_F^2}, \alpha_S \right) \right\}, \quad (1.35)$$

where the  $\widetilde{W}_i$  are the finite contributions remaining after factorisation. While the relationship between the PDFs at the factorisation scale and the bare distributions is now divergent, the renormalised quantities may be measured at some scale and used in subsequent calculations, thus making the theory predictive. In general, under a universal factorisation scheme such as  $\overline{\text{MS}}$ , structure functions may be calculated as

$$F(x, Q^2) = \sum_i \int_x^1 \frac{d\xi}{\xi} C_i \left( \frac{x}{\xi}, \frac{Q^2}{\mu_F^2}, \alpha_S \right) f_i(\xi, \mu_F^2), \quad (1.36)$$

where the  $C_i$  are the finite Wilson coefficients determined perturbatively and the PDFs  $f_i$  encode the non-perturbative structure of the calculation. This differs from the naive parton model in that the Bjorken-scaling is now broken by logarithms of the hard scale  $Q^2$ , and the sum over parton species not only runs over spin-1/2 partons (the quarks of QCD), but also contains a contribution from an initial state gluon splitting into a quark-antiquark pair.

### 1.2.1 DGLAP and PDF evolution

As a measurable quantity, the structure function itself clearly must be independent of the unphysical factorisation scheme and scale choices. The requirement of scheme independence is of course met when the factorisation scheme is followed consistently for the definition of PDFs and Wilson coefficients in all subsequent calculations. The requirement of factorisation scale independence leads to a

renormalisation group equation (RGE) for the structure function

$$\mu_F \frac{d}{d\mu_F} F(x, Q^2) = 0, \quad (1.37)$$

and consequently RGEs for the parton distributions and Wilson coefficients, once again in terms of the Altarelli-Parisi splitting functions  $P_{ij}$

$$\mu_F \frac{d}{d\mu_F} f_i(y, \mu_F^2) = \sum_j \int_z^1 \frac{dz}{z} P_{ij} \left( \frac{y}{z}, \alpha_S \right) f_j(z, \mu_F^2), \quad (1.38)$$

$$\mu_F \frac{d}{d\mu_F} C_i \left( x, \frac{Q^2}{\mu_F^2}, \alpha_S \right) = - \sum_i \int_z^1 \frac{dy}{y} C_j \left( y, \frac{Q^2}{\mu_F^2}, \alpha_S \right) P_{ij} \left( \frac{x}{y}, \alpha_S \right). \quad (1.39)$$

These are known as the Altarelli-Parisi equations [?] or the Dokshitzer-Gribov-Lipatov-Altarelli-Parisi (DGLAP) equations [?, ?, ?], and they describe how PDFs change, or *evolve* with the factorisation scale. Identically as the RGE for the running of the strong coupling performs a resummation of contributions arising from self energy diagrams, the DGLAP equation resums scale logarithms arising from collinear parton splittings.

The equations may be greatly simplified by moving to a PDF basis that largely diagonalises the matrix of splitting functions  $P_{ij}$ . For example we may construct a basis of *non-singlet* PDFs, e.g the valence distributions

$$V_i = q_i - \bar{q}_i, \quad (1.40)$$

and differences between quark sea distributions  $q_s = q + \bar{q}$

$$T_3 = u_s - d_s, \quad (1.41)$$

$$T_8 = u_s + d_s - 2s_s, \quad (1.42)$$

$$T_{15} = u_s + d_s + s_s - 3c_s, \quad (1.43)$$

$$T_{24} = u_s + d_s + s_s + c_s - 4b_s, \quad (1.44)$$

$$T_{35} = u_s + d_s + s_s + c_s + b_s - 5t_s. \quad (1.45)$$

As QCD is flavour blind, the gluon contribution to the evolution of these PDFs cancels, therefore diagonalising the matrix of splitting functions in this basis. For the nonsinglet distributions the DGLAP equation reduces to

$$\mu_F \frac{d}{d\mu_F} f_i^{\text{NS}}(y, \mu_F^2) = \int_z^1 \frac{dz}{z} P_i^{\text{NS}} \left( \frac{y}{z}, \alpha_S \right) f_i^{\text{NS}}(z, \mu_F^2). \quad (1.46)$$

Completing this basis are the gluon and the flavour singlet  $\Sigma = \sum_i (q_i + \bar{q}_i)$  PDFs. These remain coupled leading to a  $2 \times 2$  matrix of integro-differential equations for their evolution:

$$\mu_F \frac{d}{d\mu_F} \begin{pmatrix} g(x, \mu_F) \\ \Sigma(x, \mu_F) \end{pmatrix} = \int_z^1 \frac{dz}{z} \begin{pmatrix} P_{gg} & P_{g\Sigma} \\ P_{\Sigma g} & P_{\Sigma\Sigma} \end{pmatrix} \begin{pmatrix} g(z, \mu_F) \\ \Sigma(z, \mu_F) \end{pmatrix}. \quad (1.47)$$

These equations may be solved for a PDF at some scale  $Q^2$  evolved from an initial scale  $Q_0^2$ . Solutions typically follow one of two procedures; arguably the most direct consists of solving the equations iteratively through numerical methods in  $x$ -space. This method is followed in codes such as HOPPET [?], QCDCNUM [?] and APFEL [?] which employ interpolation techniques to improve the speed of the solution. Alternatively the DGLAP equations may be solved by making use

of the Mellin convolution theorem

$$\mathcal{M}\{f \otimes g\} = \mathcal{M}\{f\} \cdot \mathcal{M}\{g\}, \quad (1.48)$$

whereby the multiplicative convolution present in equations ??, ?? is reduced to a product in Mellin space; the method employed by QCD-Pegasus [?]. In the Mellin space approach, the emphasis largely lies on a fast numerical implementation of the Mellin inversion integral.

Through either method, the solution of the DGLAP equations provides a perturbative description of the behaviour of parton distributions as they vary in scale. However we remain short of a full description of the distributions having not determined their dependence upon the momentum fraction  $x$ . Furthermore the precise behaviour of the PDF and structure function renormalisation may be complicated in the attempt to overcome some of the approximations we have made so far regarding the masses of quarks contributing to our parton model, which we shall address here before discussing how the  $x$  behaviour of the PDFs may be elucidated.

### 1.3 Treatment of heavy quarks

So far in our discussion of the QCD parton model we have made the assumption that all the quarks contributing in the theory are massless, an approximation that becomes increasingly untenable when investigating scattering processes with a hard scale approaching a quark's physical mass. A careful treatment of terms depending on quark masses is therefore vital for making theoretical predictions to a dataset that spans heavy quark mass thresholds.

Dealing with heavy quark mass effects is a delicate issue in that different treatments generally have different regions of applicability. The specific combination of approaches to quark masses used when confronting a dataset with a broad

reach in hard scale is known as a heavy quark *scheme*, although not necessarily in the spirit of factorisation or renormalisation schemes as the choice often lies in the particulars of the approximation rather than in some arbitrary shuffling of parameters. A heavy quark scheme choice can therefore potentially lead to differences with alternative calculations that do not in principle vanish in the limit of an all-orders calculation.

The space of heavy quark renormalisation schemes is bounded by two regimes where the treatment is fairly simple, the fixed flavour number scheme (FFNS) and the zero-mass variable flavour number scheme (ZM-VFNS). The remaining schemes, known as general-mass variable flavour number schemes (GM-VFNS) aim to interpolate between the FFNS and ZM-VFNS, reducing to the simpler calculations in certain kinematic limits. Motivated by observations suggesting that a more careful treatment of quark mass effects is phenomenologically relevant at the LHC [?], a number of such schemes have arisen in an attempt to better describe experimental data. These typically differ by sub-leading terms in the method of interpolation between the two limiting regimes. We shall now outline some of the available choices and their potential impact in the case of a deep-inelastic scattering analysis. For simplicity we shall discuss a theory with  $n_l$  light quarks, and attempt to introduce a single massive quark  $h$  with mass  $m_h$ .

### 1.3.1 The FFN and ZM-VFN schemes

We consider first the kinematical regime where the hard scale of our scattering problem is of similar order or smaller than our heavy quark mass;  $Q^2 \lesssim m_h^2$ . Making the assumption that the initial state proton has no intrinsic heavy quark component it is reasonable to treat the heavy quark as a purely final state particle, and the only partons in the theory are the  $n_l$  light quark flavours and the gluon. In this instance, setting the factorisation and renormalisation scales  $\mu_F^2 = \mu_R^2 = \mu^2$ ;

the calculation of a structure function in Eqn. ?? takes the form

$$F(n_l, Q^2, m_h^2) = \sum_i^{n_l} C_i \left( n_l, \frac{Q^2}{m_h^2}, \frac{\mu^2}{m_h^2}, \frac{Q^2}{\mu^2} \right) \otimes f_i(n_l, \mu^2), \quad (1.49)$$

where the sum is over light quark flavours only and the full mass dependence of the heavy quark is intact in the calculation. The structure function can be separated into a contribution where only light flavours are present,  $F^L$ , and a contribution including the heavy flavour  $F^H$  as,

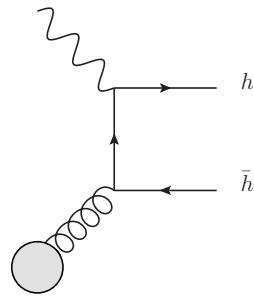
$$F(n_l, Q^2, m_h^2) = F^L(n_l, Q^2) + F^H(n_l, Q^2, m_h^2), \quad (1.50)$$

where

$$F^L(n_l, Q^2) = \sum_i^{n_l} L_i \left( n_l, \frac{Q^2}{\mu^2} \right) \otimes f_i(n_l, \mu^2), \quad (1.51)$$

$$F^H(n_l, Q^2, m_h^2) = \sum_o^{n_l} H_i \left( n_l, \frac{Q^2}{m_h^2}, \frac{\mu^2}{m_h^2}, \frac{Q^2}{\mu^2} \right) \otimes f_i(n_l, \mu^2). \quad (1.52)$$

Here  $L$  denotes the Wilson coefficients that do not contain heavy quark lines, and  $H$  includes only the diagrams that do. In this instance the heavy quark structure function first contributes at  $\mathcal{O}(\alpha_S)$  via the splitting of an initial state gluon into a  $h\bar{h}$  pair:



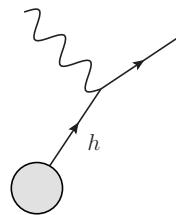
This approach is known as the *decoupling* or FFN scheme where the only quarks treated as partons are the  $n_l$  light quarks. The expression in Eqn. ?? is

unique up to terms of order  $m_l^2/Q^2$  in the light quark masses which are typically treated as part of the factorisation level corrections of  $\mathcal{O}(\Lambda_{\text{QCD}}^2/Q^2)$ . While accurate in the quark mass threshold region and below, this scheme suffers from unresummed logarithms of the ratio  $Q^2/m_h^2$  contained in the Wilson coefficients which can become large and damage the convergence of the perturbative series at scales much larger than the heavy quark mass.

These problems may be resolved in a scheme which treats the heavy quark as a massless parton above its mass threshold with the introduction of an associated heavy quark PDF. The subsequent renormalisation of the PDF resums the logarithmic contributions due to parton splitting via solution of the DGLAP equation, removing a significant disadvantage present in the FFN treatment. As this scheme is identical to the zero mass scheme discussed previously, but with an additional partonic flavour, this procedure is known as the Zero-Mass Variable Flavour Number (ZM-VFN) scheme. In the ZM-VFN a structure function calculation is simply

$$F(n_l + 1, x, Q^2) = \sum_i^{n_l+1} C_i \left( n_l + 1, \frac{Q^2}{\mu^2} \right) \otimes f_i(n_l + 1, \mu^2). \quad (1.53)$$

In this instance the heavy quark contribution to the structure function first arises now at leading order via diagrams of the type:



In the ZM-VFNS the heavy quark PDFs are set to zero below mass threshold and evolved as a massless parton according to the DGLAP equations for scales greater than the heavy quark mass. While this method alleviates the difficulties present in the FFN scheme at large scales, its treatment of heavy quarks only

in terms of massless partons completely ignores the massive contributions to the Wilson coefficients and is therefore no longer exact. The reliability of the ZM scheme is therefore particularly reduced in the region where powers of  $m_h^2/Q^2$  are significant.

### 1.3.2 General mass schemes

Analyses of QCD measurements are often performed by making a choice between using a suitable FFN scheme at scales in the region of heavy quark mass thresholds or a ZM scheme at high scales where the associated powers of  $m_h^2/Q^2$  can be safely neglected. In either case the treatment of heavy quarks is at least unambiguous, with the ZM approach yielding a simpler procedure as there is no requirement to calculate coefficient functions with the heavy quark masses intact.

For analyses of a large dataset, potentially spanning several heavy quark thresholds and extending to very high scales, the desire to improve the perturbative reliability of the calculations has led to the development of a number of hybrid or *general mass* schemes. In such schemes the treatments generally reduce to the FFN regime at low scales and the ZM treatment at high scales, with the intermediate regime handled via some interpolation between the two. Generally in a Variable Flavour Number (VFN) scheme one requires that

$$F^L(n_l, Q^2) + \lim_{Q^2 \gg m_h^2} [F^H(n_l, Q^2, m_h^2)] = F(n_l + 1, x, Q^2), \quad (1.54)$$

i.e. that the ZM-VFN and FFN calculations coincide at large scales, where the heavy quark mass dependance of the FFN Wilson coefficients can be neglected. The constraint in Eqn. ?? means that parton distributions in the two schemes may be related by a perturbatively calculable transformation.

$$f_i(n_l + 1, \mu^2) = \sum_j^{n_l} A_{ij} \left( n_l, \frac{\mu^2}{m_h^2} \right) \otimes f_j(n_l, \mu^2), \quad (1.55)$$

where the  $A$  are determined to NNLO in  $\alpha_S$  in Refs. [?, ?]. It should be noted that the  $A$  are not square matrices, with  $i$  running over the  $n_l + 1$  partons in the zero mass scheme, and the  $j$  running over the  $n_l$  partons in the FFN.

In general a GM-VFN operates as a tower of FFN-type schemes with increasing  $n_l$  as the scale increases over each quark mass threshold. In constructing a GM-VFN, the guiding principle is that physical observables should be continuous across these thresholds and therefore continuous across the  $n_l$  and  $n_l + 1$  regimes. Taking the heavy quark mass itself as the matching point between the two regimes, we demand that the GM-VFN structure function  $F^{\text{GM}}$  obeys

$$\begin{aligned} F^{\text{GM}}(m_h^2) &= \sum_j^{n_l} C_j^{\text{GM}}(n_l, m_h^2) \otimes f_j(n_l) \\ &= \sum_i^{n_l+1} C_i^{\text{GM}}(n_l + 1, m_h^2) \otimes f_i(n_l + 1). \end{aligned} \quad (1.56)$$

where the dependance upon the perturbative scales has been omitted for notational simplicity, and the GM superscripts refer to the coefficients in a general mass scheme. Using the relation in Eqn. ?? we can express the  $n_l + 1$  expression in the matching Eqn. ?? in terms of the  $n_l$  scheme PDFs, therefore obtaining the relation

$$\sum_j^{n_l} C_j^{\text{GM}}(n_l, m_h^2) \otimes f_j(n_l) \quad (1.57)$$

$$= \sum_i^{n_l+1} \sum_j^{n_l} C_i^{\text{GM}}(n_l + 1, m_h^2) \otimes A_{ij}(n_l, m_h^2) \otimes f_j(n_l). \quad (1.58)$$

Subsequently, we may make the identification

$$C_j^{\text{GM}}(n_l, m_h^2) = \sum_i^{n_l+1} C_i^{\text{GM}}(n_l + 1, m_h^2) \otimes A_{ij}(n_l, m_h^2), \quad (1.59)$$

which provides the minimal description for the construction of a GM-VFN

scheme [?]. Ensuring that Eqn. ?? is satisfied order by order in  $\alpha_S$ , we can construct the expression for the GM-VFN scheme coefficient functions above the heavy quark mass threshold. Taking the simplistic example case of Ref. [?] with a theory including only a gluon and a single heavy quark ( $h = \bar{h}$ ), the GM-VFN coefficients may be constructed to order  $\alpha_S$  as

$$C_g^{\text{LO}}(n_l + 1, m_h) = C_g^{\text{LO}}(n_l, m_h), \quad (1.60)$$

$$\begin{aligned} C_g^{\text{NLO}}(n_l + 1, m_h) &= C_g^{\text{NLO}}(n_l, m_h) \\ &\quad - C_h^{\text{LO}}(n_l + 1, m_h) \otimes A_{hg}^{\text{LO}}(n_l, m_h^2). \end{aligned} \quad (1.61)$$

where the GM superscript has been omitted, the new superscript specifying the order of the term in the perturbative expansions of the quantities  $C$  and  $A$ . Here the rightmost term in the  $\mathcal{O}(\alpha_S)$  expression Eqn. ?? is known as the *subtraction term* which ensures the cancellation of the IR-unsafe scale logs present in the FFN calculation. The ambiguity in the definition of a GM-VFNS arises upon noticing that terms proportional to powers of  $m_h/Q$  may be interchanged between the Wilson coefficients in Eqn. ?? without changing the final value of the structure function. In this respect changing the distribution of terms between the gluon and heavy quark initiated diagrams in Eqn. ?? provides the opportunity to perform a *scheme choice*, a freedom which has been exploited by several different GM-VFN scheme implementations.

The earliest complete description of a GM-VFNS was provided by the ACOT method [?] which ensures the continuity of physical quantities through Eqn. ??, but does not attempt to take advantage of the degeneracy in the GM-VFN procedure. An important result was achieved with the Simplified-ACOT or S-ACOT scheme [?, ?] which was able to exploit this ambiguity to considerably simplify the calculation of physical observables. In the S-ACOT scheme it was noted that shifts of the Wilson coefficients by their zero-mass limits may be

absorbed into a redefinition of the GM-VFNS. That is, terms such as

$$C_h(n_l + 1, m_h) - C_h(n_l + 1, 0), \quad (1.62)$$

vanish in the limit  $Q^2 \gg m_h^2$ , and therefore do not spoil the interpolation between the FFN and ZM schemes. This leads to the option of shifting to a simpler scheme where the massive heavy quark initiated coefficients may instead be evaluated with the heavy quark mass set to zero. Other options for the scheme definition were explored by Thorne and Roberts in the TR type schemes [?, ?], with the additional constraint that scale derivatives of heavy flavour structure functions should also be continuous at the matching scale.

### The FONLL approach

A more recent approach was developed by examining methods previously used to combine fixed order calculations with next-to-leading log resummation via the FONLL method [?]. The method was extended from the original application of studying the  $p_\perp$  spectrum in heavy flavour hadroproduction to the treatment of heavy quarks in DIS by Forte *et al.* [?]. The procedure begins by inverting the relationship in Eqn. ?? so as to express an  $n_l$  flavour structure function in terms of  $n_l + 1$  flavour PDFs,

$$F(n_l, Q^2) = \sum_i^{n_l} B_i \left( \frac{Q^2}{m_h^2} \right) \otimes f_i(n_l + 1, Q^2), \quad (1.63)$$

where it is important to note that the sum over flavours does not include the heavy flavour PDF, and the full heavy quark mass dependence is present in the coefficients  $B$ . To perform a matching with the massless scheme, the ZM result in Eqn. ?? can be expressed in terms of light flavour PDFs only, given the assumption that the heavy flavour PDF is generated perturbatively. In this case,

Eqn. ?? can be written

$$F(n_l + 1, Q^2) = \sum_i^{n_l} \tilde{C}_i \left( n_l + 1, \frac{Q^2}{m_h^2} \right) \otimes f_i(n_l + 1, \mu^2). \quad (1.64)$$

where once again, the sum runs over only light flavours, this time with the heavy flavour contribution being generated via DGLAP evolution included into the modified coefficient function  $\tilde{C}$ . To understand which terms are common in the two descriptions, the massive coefficient functions may be decomposed into terms logarithmically dependant upon the heavy quark mass, and terms suppressed by powers of  $m_h/Q$ :

$$B_i \left( \frac{Q^2}{m_h^2} \right) = \bar{B}_i \left( \frac{Q^2}{m_h^2} \right) + \mathcal{O} \left( \frac{m_h}{Q} \right). \quad (1.65)$$

As only the power suppressed terms vanish in the limit of  $Q^2 \gg m_h^2$ , the terms remaining must be common to both the ZM and massive scheme calculations. We can therefore express the massive structure function in a ‘massless’ limit, having dropped those terms in the coefficient functions that are suppressed by powers of  $m_h/Q$ :

$$\bar{F}(n_l, Q^2) = \sum_i^{n_l} \bar{B}_i \left( \frac{Q^2}{m_h^2} \right) \otimes f_i(n_l + 1, Q^2). \quad (1.66)$$

The FONLL result for the structure function is given by the sum of the massive calculation in Eqn. ??, and the massless calculation in Eqn. ?? with the asymptotic limit of the massive calculation in Eqn. ?? subtracted.

$$F^{\text{FONLL}}(Q^2) = [F(n_l, Q^2) + F(n_l + 1, Q^2)] - \bar{F}(n_l, Q^2). \quad (1.67)$$

With the subtraction ensuring the cancellation of terms which are double counted between the massive and massless calculations. Therefore in this expression the mass-suppressed terms present in the FFN calculation are fully accounted for in the GM scheme, with the duplicate terms subtracted. The simplicity of this approach helped to elucidate many of the differences between general mass

schemes.

It should be noted that while general mass schemes suffer from an ambiguity in their definition compared to the simpler fixed-flavour and zero mass schemes, the differences between them are always of higher order compared to the calculation at hand, as is the case in any true scheme choice. Indeed, a well-defined GM-VFNS will always reduce to the decoupled result at low scales and the zero-mass result at scales much higher than the quark mass, behaving effectively as a tower of fixed flavour schemes with increasing number of partonic quarks. The general-mass schemes therefore do not suffer from a significant loss of predictive power, and are able to provide considerable improvement over the simpler schemes when dealing with datasets spanning quark mass thresholds.

## 1.4 General features of parton distributions

While we have now described how the parton distributions functions at an experimental scale  $Q^2$  may be found by evolving parton distributions from an initial scale, and discussed briefly how the renormalisation of heavy quark distributions may be accomplished, the issue of determining the functional dependence of the parton distributions upon the momentum fraction  $x$  at some initial scale  $f_i(x, Q_0^2)$  remains.

The number of independent PDFs to be determined is dependent upon the choice of initial scale, as quark distributions that can be considered *heavy* with respect to  $Q_0^2$  may be generated perturbatively through the DGLAP procedure outlined previously. The typical choice is to determine the parton distributions at some scale  $m_s^2 < Q_0^2 \leq m_c^2$  such that the flavours  $c, b, t$  are produced by evolution. These scale choices minimise the number of distributions to be determined while remaining perturbatively reliable.

As the remaining seven distributions<sup>1</sup> are fundamentally a parametrisation of

---

<sup>1</sup>The gluon, the  $u, d, s$  quarks and their antiquarks.

the nonperturbative dynamics of the proton, they are by definition out of reach of a perturbative analysis. There are however some general statements that may be made of their  $x$ -dependence that are independent of the hard scale. The most important of which are the parton distribution *sum rules* which constrain the relative normalisation of PDFs.

Firstly, the *momentum sum rule* (MSR) ensures that the parton distributions' fractional momenta sum to the momentum of the parent proton

$$\int_0^1 dx [x\Sigma(x, Q^2) + xg(x, Q^2)] = 1, \quad (1.68)$$

where  $\Sigma$  is the singlet distribution defined previously. Following this are the quark valence sum rules. These fix the quark distributions such that the resulting proton has the appropriate quantum numbers,

$$\text{up-valence: } \int_0^1 dx (f_u(x, Q^2) - f_{\bar{u}}(x, Q^2)) = 2, \quad (1.69a)$$

$$\text{down-valence: } \int_0^1 dx (f_d(x, Q^2) - f_{\bar{d}}(x, Q^2)) = 1, \quad (1.69b)$$

$$\text{strange-valence: } \int_0^1 dx (f_s(x, Q^2) - f_{\bar{s}}(x, Q^2)) = 0. \quad (1.69c)$$

From these rules we may infer additional constraints upon individual PDFs. The MSR suggests a form for the large- $x$  behaviour of the distributions, in that they should parametrically tend to zero as  $x \rightarrow 1$ . The number sum rules in Eqns. ?? require the valence-type distributions to be integrable over the whole  $x$ -range. While there is no requirement for the singlet and gluon distributions to be integrable, their first moments must be, as required by the MSR. Combining these three constraints we may parametrise the large and small- $x$  behaviour of

both valence-like and gluon or singlet-like distributions as:

$$\begin{aligned} f_V(x, Q_0^2) &= N_V x^{\alpha_V} (1-x)^{\beta_V} r_V(x), \\ f_\Sigma(x, Q_0^2) &= N_\Sigma x^{\alpha_\Sigma} (1-x)^{\beta_\Sigma} r_\Sigma(x). \end{aligned} \quad (1.70)$$

In these expressions, the parameters  $\alpha$  and  $\beta$  control the small and large- $x$  PDF behaviour respectively. The  $\beta$  should be such that the PDFs tend to zero smoothly at large- $x$ , and the  $\alpha$  such that the valence distributions are integrable, and the first moment of the gluon and singlet are integrable. The overall PDF normalisations  $N$  being constrained via the appropriate sum rules.

Finally, what remains in the determination of the distributions are the remainder terms  $r(x)$  which describe the PDFs between the two  $x$ -limits. Their determination is considerably more complex and is a ongoing source of research. Much of this thesis will be dedicated to discussing the determination of these remainder functions.

# Chapter 2

## Review of PDF determination

Understanding the functional structure of parton distributions is a complex task that has been subject to a number of approaches over the years. As nonperturbative quantities describing the behaviour of QCD bound states, in principle they may be subject to analysis using Lattice QCD methods. While a great deal of effort and progress has been made in understanding PDFs through nonperturbative methods [?, ?, ?, ?], results remain short of providing distributions for practical application at hadron colliders.

The majority of PDF analyses are therefore performed analogously to the determination of many other QCD parameters; via a fit to appropriate experimental data. The fundamental difficulty in PDF fits being that they are determinations of *functions* rather than single parameters and therefore one must attempt to find some optimum solution in an (in principle) infinite-dimensional functional parameter space. This is of course complicated by having only a finite set of experimental data points upon which to perform a fit. Moreover as the applications involving PDFs have become more precise, a detailed understanding of the uncertainties in the determination of PDFs has become vital. The problem of PDF fitting is therefore one of finding a reliable estimator for a probability distribution in a space of functions.

The complexity of the task, along with the inherent ambiguities in the QCD treatment of data, led to the emergence of several competing methodologies and determinations. Today there are a diverse array of fitting groups producing sets of parton distribution functions, the most important of which being the ABM [?, ?] (formerly ABKM [?]), CTEQ/CJ [?, ?, ?, ?], JR/GJR [?, ?], HERAPDF [?, ?], MSTW [?, ?] (formerly MRST [?, ?, ?, ?]) and NNPDF [?, ?, ?, ?, ?, ?] groups. Typically PDF sets are provided for a variety of theory input parameters such as perturbative order, and value of the strong coupling. All modern PDF sets now include a quantitative assessment of their associated uncertainties. In this chapter we shall review the ingredients and methods utilised in a modern PDF determination, primarily focusing on the methodology of the three global PDF fits recommended for LHC phenomenology by the PDF4LHC working group [?], namely the procedures of the CTEQ, MSTW and NNPDF collaborations.

These three groups produce PDF sets determined from a fit to a wide range of experimental data, including DIS, Drell-Yan and inclusive jet cross sections. The CTEQ and MSTW determinations follow a similar fitting procedure and method of uncertainty estimation, with the NNPDF group taking a rather different approach to both. We will now describe the basic fitting procedure of these groups, with an eye to detailing areas where the groups have different solutions.

## 2.1 Experimental data on parton distributions before the LHC

The most important ingredient in the determination of parton distributions is naturally the selection of the dataset from which to extract PDF constraints. The first step in performing a PDF fit is therefore to identify which datasets are most sensitive to input parton distributions, and offer precise and reliable data. As PDF determinations to date have relied only upon fixed-order perturbation

theory, the dataset chosen should probe sufficiently inclusive observables which are therefore relatively insensitive to resummation effects. In general PDF fitting collaborations also require data to be taken at a sufficiently high scale that leading-twist factorisation remains reliable, although there are some exceptions which we shall discuss later in the section. Here we shall briefly discuss some of the most important processes in terms of PDF sensitivity, and review some of the most relevant experimental measurements. For this section we shall restrict ourselves to data available before the start of LHC operation in order to provide a background for the methodological developments made in the light of LHC data.

### Fixed-Target and collider DIS

Deep inelastic scattering data provides the backbone for much of a PDF analysis, and data is available from a wide array of sources. Precise electron-proton scattering data from HERA provides the cleanest probe of proton structure function data, while high-luminosity fixed-target experiments can provide important constraints, at the expense of potentially having to deal with additional data corrections due to nuclear and higher-twist effects. As DIS is one of the best understood scattering processes in QCD, precise theoretical predictions are available up to 3-loop order in the zero-mass scheme [?, ?] and 2-loop order with full heavy quark masses intact [?, ?, ?, ?, ?, ?, ?].

At leading order, neutral current DIS measurements from a proton target directly probe the quark sea distributions  $q_i + \bar{q}_i$ , with the relative power of each flavour contribution mediated via its coupling to  $\gamma, Z$ . Charged current, and  $Z$ -mediated neutral current data can provide some constraint upon PDF flavour separation via the  $F^3$  structure function.

In addition to proton structure function measurements, data obtained from scattering off deuterium targets can be important in constraining light quark flavour separation under the assumption of isospin symmetry. Data may be

presented as direct measurements of  $F_d$  structure functions or as the ratio  $F_d/F_p$ . A simultaneous fit to deuterium and proton data may therefore provide important constraints upon the  $u - d$  and  $u/d$  PDF combinations. Data determined via deuteron scattering are subject to nuclear corrections e.g. shadowing effects [?] which may be estimated as part of the theoretical treatment or neglected; the corrections to be considered part of the theory uncertainty.

Alongside the direct information on quark distributions, scaling violations present in structure function data provide constraints upon the gluon. While rather indirect, the wealth of DIS measurements available at a wide range of scales provides a great deal of information on the structure of the gluon distribution.

DIS data may be presented either as experimental cross sections, or separated into structure functions. Fixed target structure function data on  $F_2$  from muon scattering is available for both proton and deuteron targets from the BCDMS [?, ?], NMC [?, ?] and Fermilab E665 [?] experiments. Electron scattering  $F_2$  data is also available from SLAC data on both proton and deuteron targets [?]. The longitudinal structure function  $F_L$  is measured in fixed target experiments also by SLAC [?] , BCDMS [?] and NMC [?].

In addition to the large datasets available from fixed target experiments, HERA data provides a clean probe of DIS properties, although with HERA data the separation of cross-sections into structure functions is typically not performed. Neutral current cross-section data is provided by ZEUS [?, ?, ?, ?] and H1 [?, ?, ?]. Charged-current DIS data is also provided by the HERA collaborations [?, ?] along with information on the longitudinal structure function  $F_L$  [?, ?]. Information on charm hadroproduction in DIS is available via  $F_2^{\text{charm}}$  measurements at HERA also [?, ?, ?, ?, ?, ?, ?]. This data provides particular constraint upon the gluon PDF, and has been an important testing ground for heavy quark flavour schemes. The clean  $ep$  environment means that data is unaffected by nuclear or deuteron corrections, although low energy datapoints may still suffer from substantial

higher-twist corrections. These corrections are typically kept under control by kinematic cuts on the affected points, however some groups (notably the ABM/CJ groups) include the affected data and attempt to model the corrections.

HERA measurements from the two collaborations have been examined as a combined analysis and dataset, so far resulting in two studies of direct interest to PDF determination; a combination of HERA-1 inclusive DIS data [?], and of charm production cross-sections [?].

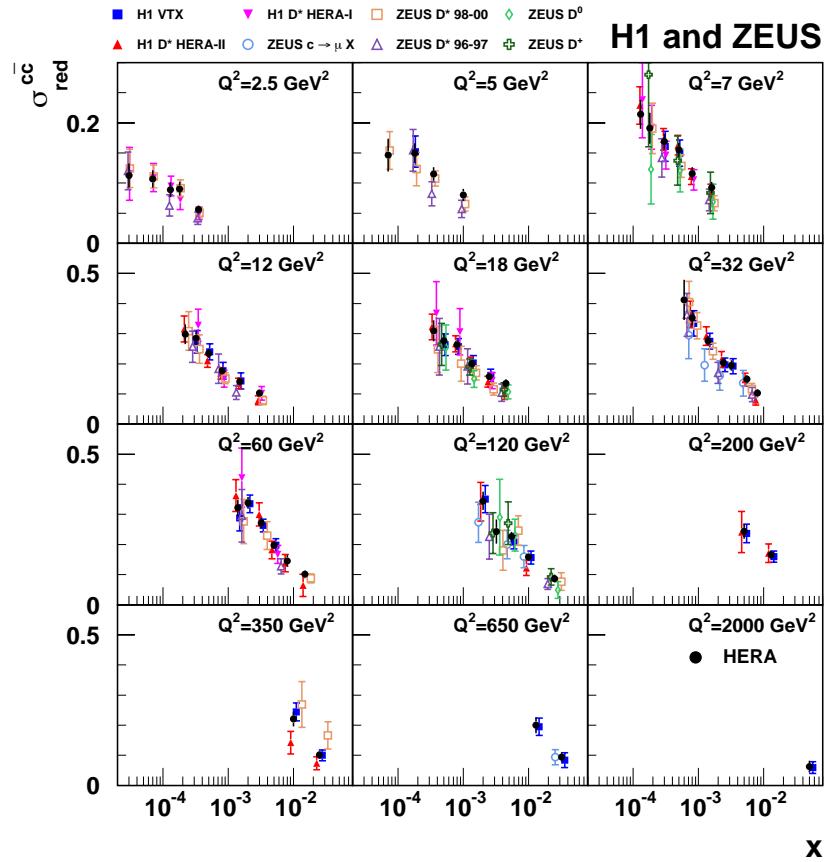


Figure 2.1: Reduced charm cross-section data from the HERA combined measurement. Data from the measurements contained in the combination analysis is shown for comparison. Figure from [?].

The very large quantity of deep-inelastic scattering measurements performed at a variety of experimental facilities means that generally DIS data forms the backbone for PDF fits, providing a substantial proportion of the experimental

data points used in a fit.

### Neutrino DIS

There are a number of measurements available for the scattering of neutrino beams from heavy nuclear targets. For example the NuTeV [?] and CHORUS [?] data on neutrino  $F_2$  and  $F_3$ . Assuming an approximately isoscalar target, and neglecting CKM factors, the PDF dependence of the neutrino structure function data at leading order is given by [?]

$$F_2^\nu(x) = x(u^+(x) + d^+(x) + 2s(x) + 2\bar{c}(x)), \quad (2.1)$$

$$F_2^{\bar{\nu}}(x) = x(u^+(x) + d^+(x) + 2\bar{s}(x) + 2c(x)), \quad (2.2)$$

and for the  $F_3$  structure function,

$$F_3^\nu(x) = x(u^-(x) + d^-(x) + 2s - 2\bar{c}), \quad (2.3)$$

$$F_3^{\bar{\nu}}(x) = x(u^-(x) + d^-(x) - 2\bar{s}(x) + 2c(x)). \quad (2.4)$$

A simultaneous fit of these data points therefore provides a good handle upon the valence quark distributions  $q - \bar{q}$ . These datasets are relatively precise; however they are subject to potentially large nuclear corrections which introduce an uncertainty that is poorly understood.

Neutrino DIS becomes particularly valuable for PDF determination when considering the semi-inclusive DIS dimuon production process  $\nu N \rightarrow \mu\mu X$  illustrated in Figure ???. In this process the contribution from initial state strangeness is Cabibbo favoured, therefore providing a direct handle on the strange distribution whose contribution is ordinarily difficult to discern from total structure function measurements. Measurements of this process are therefore commonly used as a strangeness probe, and data has been provided by the NuTeV/CCFR collaborations [?].

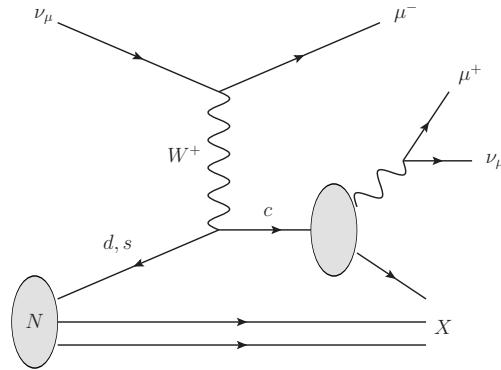


Figure 2.2: Leading order diagram for dimuon production in neutrino DIS.

### Fixed-target and collider Drell-Yan

After DIS measurements, the production of electroweak vector bosons in hadronic collisions provides the next most important contribution to the constraint of parton densities, with precise predictions available at NNLO in QCD [?, ?, ?]. At leading order the neutral current Drell-Yan process is moderated by the PDF combination

$$q(x_1)\bar{q}(x_2) + \bar{q}(x_1)q(x_2), \quad (2.5)$$

and provides a direct probe of various partonic combinations depending upon the experimental configuration. In the Drell-Yan process the relevant kinematic variables are the invariant mass of the lepton pair

$$M_{ll}^2 = (E_1 + E_2)^2 - (\mathbf{p}_1 + \mathbf{p}_2)^2, \quad (2.6)$$

and the intermediate boson's rapidity, given in the detector frame by

$$y = \frac{1}{2} \log \frac{E + p_L}{E - p_L}, \quad (2.7)$$

where  $E$  is the detector frame energy of the intermediate boson, and  $p_L$  its longitudinal momentum. in terms of which the parton- $x$  is given by;

$$x_{\pm} = M_{ll} e^{\pm y} / \sqrt{s}, \quad (2.8)$$

where  $s$  is the centre-of-mass energy squared of the reaction and the  $\pm$  denotes the parton direction with respect to the beam frame. High rapidity measurements therefore constrain PDFs at both high and low- $x$ .

Additionally the charged-current process  $qq' \rightarrow l^{\pm}\nu_l$  provides information on quark flavour separation in the initial state hadrons. While the rapidity of the lepton pair resulting from  $Z/\gamma$  decay in neutral current Drell-Yan is experimentally straightforward to distinguish, the presence of a neutrino in the final state of  $W$  production processes complicates the direct resolution of the  $W$  rapidity. Therefore data is often presented in the pseudorapidity of the detected lepton,

$$\eta = -\log \tan \theta, \quad (2.9)$$

defined in terms of the angle  $\theta$  between the final state lepton and the beam axis. It can therefore be measured without knowledge of the particle mass and momentum. The pseudorapidity coincides with the standard rapidity in the case of massless particles where  $E = |\bar{\mathbf{p}}|$ .

Lepton asymmetries are another common form for experimental results in Drell-Yan, defined in terms of  $W^{\pm} \rightarrow l^{\pm}\nu_l$  differential cross-sections  $d\sigma_{l^{\pm}}/d\eta_l$  as

$$A_W^l = \frac{d\sigma_{l^+}/d\eta_l - d\sigma_{l^-}/d\eta_l}{d\sigma_{l^+}/d\eta_l + d\sigma_{l^-}/d\eta_l}, \quad (2.10)$$

such measurements also benefit from the cancellation of shared systematic uncertainties. Measurements of lepton pair production from proton beams incident upon heavy nuclear targets, such as the E605 [?] experiment determining dimuon production from a copper target are useful for the constraint of the light

quark sea  $q + \bar{q}$ . These measurements are typically very precise but suffer from poorly determined nuclear corrections. Several approaches have been performed to study the extent of these corrections [?, ?, ?, ?], although the effects are typically small and may sometimes be discounted in comparison to experimental uncertainties [?]. Contributions from initial state heavy quarks and strangeness are typically suppressed in these measurements due to the relatively low scales.

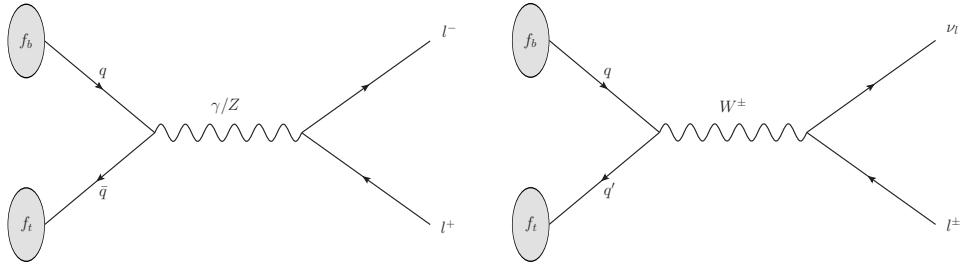


Figure 2.3: Drell-Yan process at leading order, initiated by beam protons with PDF  $f_b$  and target protons with PDF  $f_t$ . The neutral current process is shown on the left, and the charged current process on the right.

Fixed target experiments upon hydrogen or deuterium targets provide a relatively clean probe and the ratio of Drell-Yan cross sections in proton to deuteron targets can provide crucial information on the  $u/d$  PDF combination. While relatively free of nuclear effects, deuteron data still suffers from poorly understood corrections, which have been the subject of extensive study [?, ?, ?, ?]. Experimental measurements from the Fermilab NuSea/E866 collaboration are commonly used, providing data from  $pp$  [?] and  $pd/pp$  [?] experiments.

The theoretically cleanest environment to examine the Drell-Yan process is at high scales at a collider. Several measurements are available from the Tevatron collaborations which provide information free of nuclear or deuteron corrections. As a  $p\bar{p}$  collider, neutral-current Drell-Yan at the Tevatron targets the quark valence contribution and asymmetry data provides information on the  $u/d$  ratio. A measurement of the  $Z$  rapidity distribution is available from D0 [?], and several measurements are available for  $W$  lepton asymmetries from both Tevatron

collaborations [?, ?, ?, ?].

In order to obtain a handle on the contribution of initial state strange quarks to the Drell-Yan process it is once again necessary to examine less inclusive processes. Of particular interest are measurements of  $W$  production in association with a charm jet, analogous to the usefulness of dimuon measurements in neutrino DIS where the strange contribution is favoured in terms of CKM elements. Measurements of this process were initially made at the Tevatron by both CDF [?] and D0 [?]. More precise determinations can be obtained by normalisation with respect to the total  $W +$  jets rate [?].

### Jet production data

While DIS data provides constraints upon the gluon distribution via scaling violations and contribution to heavy quark and longitudinal structure functions, DIS and Drell-Yan data do not provide a substantial direct constraint upon gluon densities. The most constraining datasets for the gluon, particularly in the uncertain large- $x$  region, are those of jet production measurements. The large strong coupling of the gluon combined with a high gluon luminosity in the proton at high scales results in  $gg$  initiated diagrams being the dominant sub channels for the production of inclusive jet and dijet events.

Cross-section calculations for inclusive jet and dijet data in hadron-hadron collisions are available at NLO in QCD [?, ?, ?, ?], however a great deal of progress has been made in the determination of the NNLO corrections [?, ?, ?], with the exact gluon-gluon sub channel calculation recently determined [?]. For the full calculation however, only approximate NNLO results are available via threshold resummation techniques [?, ?, ?]. Jet data may therefore only be included into an NNLO PDF fit through an approximate treatment if at all.

Jet data must be included via some clustering algorithm which takes a QCD final state and identifies suitable jet-like structures. Earlier measurements were

performed with so-called cone algorithms, although these are potentially very sensitive to infrared and collinear effects. More recent experiments typically utilise sequential-combination algorithms such as the Cambridge-Aachen [?, ?],  $k_T$  [?] or anti- $k_T$  [?] algorithms, often used as implemented in the efficient FastJet [?] package.

The CDF collaboration has published precise measurements of inclusive jet [?, ?] and dijet [?] cross sections. Data is also available from the D0 experiment, once again for inclusive [?] and dijet [?] quantities. Figure ?? shows the results of an inclusive jet measurement at CDF using the  $k_T$  clustering algorithm.

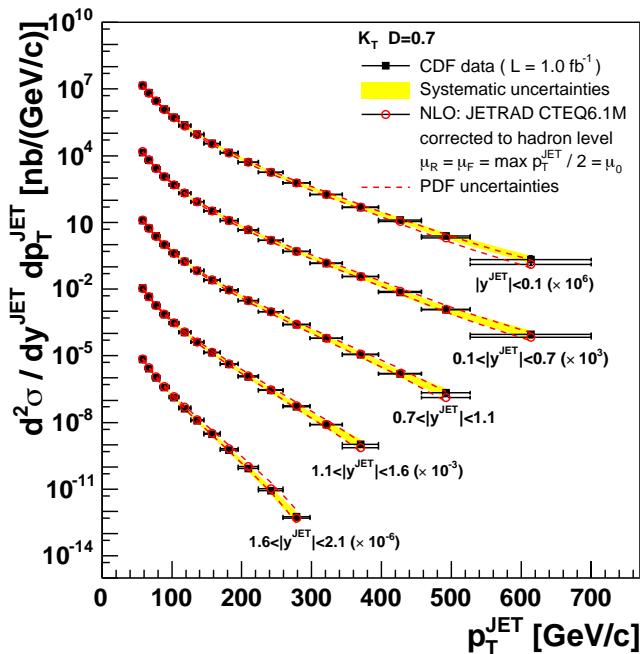


Figure 2.4: Inclusive jet data from CDF using the  $k_T$  jet clustering algorithm, compared to predictions from the CTEQ6.1M PDF set. Figure from [?].

### Prompt photon measurements

Complementary to the data on jet production, measurements of prompt photon processes  $pp/p\bar{p} \rightarrow \gamma X$  can also provide an important handle on the gluon. The term *prompt* photon refers to the production of a photon in the hard scatter

rather than in subsequent emissions. Prompt photons in the final state can originate either from Compton scattering processes  $gq \rightarrow \gamma q$  or annihilation events  $q\bar{q} \rightarrow \gamma g$ , processes denoted *direct* photon production. Alternatively prompt photons may be produced via the fragmentation of final state hadrons into photons via so-called fragmentation functions [?, ?]. In  $p\bar{p}$  collisions the Compton scatter is typically the dominant process, particularly at higher scales where the fragmentation contribution is suppressed. For  $p\bar{p}$  events the annihilation contribution becomes more important due to the enhanced  $q\bar{q}$  PDF luminosity. Figure ?? demonstrates the relative fraction of these contributions to the cross-section for a range of photon transverse energy  $E_T$ .

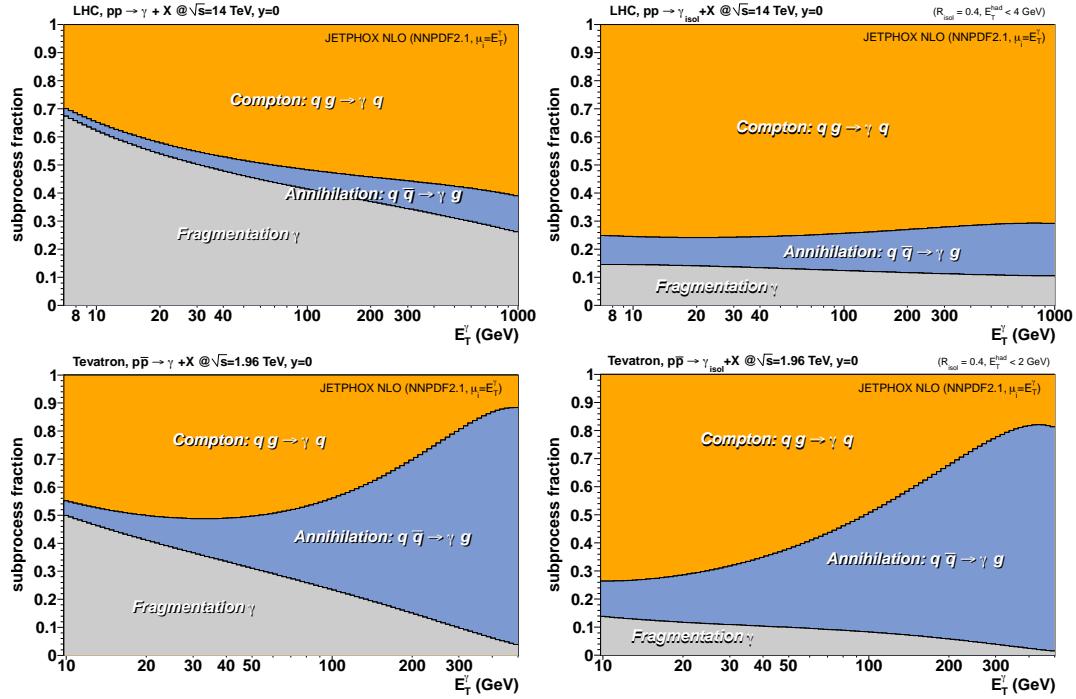


Figure 2.5: Relative contribution of partonic subprocesses to  $pp/p\bar{p} \rightarrow \gamma X$ . Figures on the left refer to the inclusive case, and on the right to the observable after isolation cuts on the final state photon. Figure from [?].

For the purposes of PDF determination direct photon measurements which are free of the additional uncertainties introduced when performing calculations with photon fragmentation functions are the ideal measurement. While performing

selection cuts to measure only the direct photon contribution is experimentally challenging, the relative contribution of fragmentation photons may be suppressed by making isolation cuts upon the final state photon. These cuts admit only photons with no hadronic material in close proximity. Smooth-cone cuts such as the Frixione isolation criterion [?] in principle can remove entirely the fragmentation contribution. However these cuts remain challenging to implement experimentally, with experimental data usually obtained with simpler isolation cuts which aim to suppress rather than eliminate fragmentation photons.

Theoretical predictions are available at NLO for the Compton process [?, ?] and commonly used as implemented in the JETPHOX program [?, ?, ?]. While inclusive data is challenging to include in a PDF determination due to contamination by fragmentation photons, results are available from a wide range of isolated photon measurements. Isolated data is available from UA1/UA2 at the SppS [?, ?, ?], PHENIX at RHIC [?], CDF [?, ?, ?, ?, ?] and D0 [?, ?, ?].

### **Top quark pair production data**

The production of top-antitop pairs is potentially a process of great interest in the determination of PDFs, with calculations available up to NNLO for the total cross-section [?, ?, ?, ?]. The impact of the total top pair production cross-section upon PDFs is quite sensitive to the kinematics of the collider, with Tevatron data probing directly the quark content of the proton, while data from colliders with higher centre of mass energies being dominated by the gluon-gluon channel. Precise data from the Tevatron is available in the form of a combined D0-CDF analysis [?].

### **Experimental cuts**

A simple cut is typically performed on the hard scale  $Q^2$  and for DIS the final state invariant mass  $W^2$  to ensure the reliability of perturbative predictions. The

MSTW2008 parton fit uses an initial scale for evolution of  $Q_0^2 = 1 \text{ GeV}^2$ , CT10 uses  $Q_0^2 = 1.3 \text{ GeV}^2$  and NNPDF2.3  $Q_0^2 = 2 \text{ GeV}^2$ . Most of the data included in global parton fits has a minimum of  $Q^2 \sim 2$  to  $5 \text{ GeV}^2$  [?].

## 2.2 Methodological elements

### 2.2.1 Parametrisation

Given an experimental dataset, one must choose a convenient and effective parametrisation of the parton distribution functions such that their predictions may be compared to data. Nominally there are a total of 13 PDFs, six quarks, six antiquarks and a gluon. However as mentioned in the previous section, the heavy quarks  $c, b, t$  are determined perturbatively. There are therefore typically seven free PDFs remaining to be fitted. The parton parametrisation basis is chosen for ease of fitting and perturbative evolution; a basis close to the DGLAP basis in Eqn. ?? is desirable for efficiency. However often a different basis is chosen to avoid fitting quantities that are poorly defined by the experimental dataset.

For example, MSTW2008 [?] uses the following basis for their determination:

$$\begin{aligned} g, \\ q_v &\equiv q - \bar{q}, \\ \Delta &\equiv \bar{d} - \bar{u}, \\ S &\equiv 2(\bar{u} + \bar{d}) + s + \bar{s}, \\ s^\pm &\equiv s \pm \bar{s}, \end{aligned} \tag{2.11}$$

where  $g$  is the gluon PDF and the  $q_v$  correspond to the  $u, d$  quark valence PDFs. These fully parameterise the degrees of freedom to be determined. A functional form in  $x$  is then chosen for each of the distributions (the value of  $Q^2$  is kept fixed at the input scale for fitting). While all groups include the limiting- $x$  description

of Eqn. ??, the choice of parametrisation for the remainder function  $r$  varies substantially between fitting groups. As an example, the valence quark PDF  $q_v$  parametrisation in MSTW2008 is provided by the expression

$$xq_v(x, Q_0^2) = ax^b(1-x)^c(1+d\sqrt{x}+ex), \quad (2.12)$$

and the equivalent parametrisation in CT10 [?] is

$$xq_v(x, Q_0^2) = ax^b(1-x)^b \exp(cx + dx^2 + e\sqrt{x}), \quad (2.13)$$

where the  $(a, \dots, e)$  are the parameters to be determined in the fit. In total the MSTW08 basis has 30 free parameters (taking into account sum rule constraints), the CT10 parametrisation is a little less flexible, having 26 free parameters. The problem is now reduced to finding the optimum parameters for the 7 PDFs that minimise some measure of fit quality, the differing versions of which we shall discuss later in the chapter.

The NNPDF procedure is markedly different from that of the other PDF fitting groups and the first major difference lies in the choice of parametrisation. Unlike in the general procedure outlined above, neural networks are used to provide the functional  $x$  dependence of the PDFs. Neural networks are a typical computational tool in machine learning environments, often used in regression applications where flexibility and a lack of bias with respect to a conventional fixed parametrisation are desired. A typical neural network in a fitting context will usually have considerably more functional freedom (and therefore parameters) than a normal parametric model, with the neural network compensating for its relative generality with respect to the problem by having much greater flexibility.

The use of neural networks as applied to the determination of the proton structure function  $F_2^p$  was first suggested in Ref. [?] and subsequently developed in [?]. The approach was later extended to the determination of quark

distributions [?] before becoming a global analysis of PDFs as of NNPDF2.0 [?] as part of the wider NNPDF methodology.

In the NNPDF approach the specific networks used in the parametrisation are multi-layer feed forward neural networks configured with 2-5-3-1 architecture. This architecture applied over seven PDFs results in a fit with a total of 259 free parameters, considerably more than in competing approaches. The architecture chosen in fact has considerable redundancy to minimise potential bias due to inflexibility or choice of architecture. The flexibility of the approach was demonstrated in Ref. [?] where the architecture was modified considerably, with no significant change in the fit results.

Due to the redundant parametrisation provided by the neural networks, there is a great deal of freedom in the choice of the input parton distribution basis. In the more recent NNPDF analyses: sets NNPDF 2.1 and NNPDF 2.3, the basis is chosen for simplicity of evolution as:

$$\begin{aligned}
 \text{gluon} & \quad g, \\
 \text{singlet} & \quad \Sigma \equiv \sum_{i=1}^{n_f} (q_i + \bar{q}_i), \\
 \text{valence} & \quad V \equiv \sum_{i=1}^{n_f} (q_i - \bar{q}_i), \\
 \text{triplet} & \quad T_3 \equiv (u + \bar{u}) - (d + \bar{d}), \\
 \text{sea asymmetry} & \quad \Delta \equiv \bar{d} - \bar{u}, \\
 \text{strange sea/valence} & \quad s\pm \equiv s \pm \bar{s}.
 \end{aligned} \tag{2.14}$$

The equivalent functional forms for the fitting in terms of the Neural Networks

are;

$$\begin{aligned}
\Sigma(x, Q_0^2) &= x^{-\alpha_\Sigma} (1-x)^{\beta_\Sigma} \text{NN}_\Sigma(x), \\
V(x, Q_0^2) &= A_V x^{-\alpha_V} (1-x)^{\beta_V} \text{NN}_V(x), \\
T3(x, Q_0^2) &= x^{-\alpha_{T3}} (1-x)^{\beta_{T3}} \text{NN}_{T3}(x), \\
\Delta(x, Q_0^2) &= A_\Delta x^{-\alpha_\Delta} (1-x)^{\beta_\Delta} \text{NN}_\Delta(x), \\
g(x, Q_0^2) &= A_g x^{-\alpha_g} (1-x)^{\beta_g} \text{NN}_g(x), \\
s^+(x, Q_0^2) &= x^{-\alpha_{s^+}} (1-x)^{\beta_{s^+}} \text{NN}_{s^+}(x), \\
s^-(x, Q_0^2) &= x^{-\alpha_{s^-}} (1-x)^{\beta_{s^-}} \text{NN}_{s^-}(x) - s_{\text{aux}}(x, Q_0^2),
\end{aligned} \tag{2.15}$$

where the NN denote the 2-5-3-1 neural network parametrisations and the  $A$  are set by enforcing the appropriate sum rules. In the NNPDF approach the treatment of the limiting exponents  $\alpha$ ,  $\beta$  is rather different. These factors are introduced in order to speed up the convergence of the neural network fitting, with the intention of providing a rough preprocessing function as a backbone for the neural networks to deviate from, and ensuring that the functions have the correct behaviour under integration. These exponents are therefore randomised within an optimised range at the start of the fit and are not modified by the fitting procedure. The final results should therefore be reasonably independent of the preprocessing factor and of the coefficients involved.

While determinations with fixed parametrisations typically design the strange valence functional form such that the strange valence sum rule is automatically satisfied, this cannot be done with a neural net parametrisation. In the determinations up to NNPDF2.3 the strange auxiliary term  $s_{\text{aux}}(x, Q_0^2)$  in Eqn. ?? is therefore introduced to ensure the strange valence sum rule is followed, and has the form [?]:

$$s_{\text{aux}}(x, Q_0^2) = A_{s^-} (x^r (1-x))^s. \tag{2.16}$$

### 2.2.2 Fit quality and minimisation

With an experimental dataset selected and a choice made for the parametrisation of the PDFs, the optimal fit should be determined by varying fit parameters and attempting to minimise some measure of fit quality. Different groups make quite different choices not only in the minimisation method but also in the measure used to determine fit quality. The most general statement that can be made is that the global fit quality (generally denoted  $\chi^2$ ) is built from the quality of fit to individual datasets as

$$\chi^2 = \sum_k^n \chi_k^2, \quad (2.17)$$

for a fit with  $n$  data sets, each with a consistent normalisation. In the NNPDF approach the full covariance matrix of the data is used in determining the quality of fit, including all appropriate correlations within and between datasets. The  $\chi^2$  measure for a set of data with common correlations is then given by

$$\chi_k^2 = \sum_{i,j=1}^{N_{\text{dat}}} \frac{(D_{k,i} - T_{k,i})(D_{k,j} - T_{k,j})}{\text{Cov}[i,j]}. \quad (2.18)$$

Here the  $T$  are the theoretical predictions for the experimental data points  $D$  calculated from the neural network parametrisation, and  $\text{Cov}[i,j]$  is the covariance between data points  $i$  and  $j$ . In practice there is a ensemble of neural networks each associated with a single Monte Carlo sample of the experimental data, for the purposes of error propagation. This point will be discussed in more detail later in the chapter. In NNPDF determinations the full experimental correlations should be available for a dataset to be included into the determination.

Other groups take a different strategy, often with the suggestion that correlation effects are small to negligible with the exception of overall normalisations.

Adopting the same practice as earlier MRST fits, the MSTW2008 fit uses an uncorrelated  $\chi^2$  measure over much of its dataset [?], with the normalisation of the theory predictions set by a fitted parameter  $\mathcal{N}$

$$\chi_k^2 = \sum_{i=1}^{N_{\text{dat}}} \frac{(D_{k,i} - T_{k,i}/\mathcal{N}_k)^2}{\text{Var}[i]} + \left( \frac{1 - \mathcal{N}_k}{\sigma_k^{\mathcal{N}}} \right)^4, \quad (2.19)$$

where the final quartic penalty is intended to prevent the normalisation deviating too far from the experimental normalisation uncertainty  $\sigma_{\mathcal{N}}$ , and the variance  $\text{Var}[i]$  is constructed by the sum in quadrature of the statistical and uncorrelated systematic errors. The CT series of fits utilise a  $\chi^2$  measure that includes systematic uncertainties in terms of explicit shifts [?, ?]. In this arrangement, the fit quality measure is given by

$$\chi_k^2 = \sum_{i=1}^{N_{\text{dat}}} \frac{1}{\text{Var}[i]} \left( D_{k,i} - T_{k,i} - \sum_{n=1}^{N_{\text{corr}}} r_n \sigma_{k,n,i}^{\text{corr}} \right)^2 + \sum_{n=1}^{N_{\text{corr}}} r_n^2, \quad (2.20)$$

where here the  $\sigma^{\text{corr}}$  are the  $N_{\text{corr}}$  correlated systematic uncertainties. In this procedure the theory predictions  $T$  are shifted parametrically by the variables  $r$ . The optimal shift values are found by minimising the  $\chi^2$  with respect to the  $r$  analytically at each stage of the fit. This procedure was introduced to accommodate for overall shifts in the CT10 distributions. A similar method which was adopted in MSTW2008 for a limited number of datasets where correlations were deemed to be important, with the normalisations also determined in the fit as per the uncorrelated case.

### Normalisation uncertainty

A key point that must be addressed when constructing a measure of fit quality is the treatment of normalisation uncertainties, or multiplicative uncertainties in general. Even using the same definition of the fit quality measure, substantial

deviations may be produced by defining the covariance matrix and therefore the breakdown into systematic errors, differently.

The full experimental uncertainty information is characterised by the sum of all uncorrelated errors for a datapoint  $\sigma^{\text{unc}}$ ; the set of  $N_{\text{add}}$  correlated additive systematics  $\sigma^{\text{add}}$ ; and the set of  $N_{\text{mul}}$  correlated multiplicative systematics  $\sigma^{\text{mul}}$ . Given this information one may naively define an ‘*experimental*’ prescription [?] for constructing a covariance matrix as

$$\text{Cov}[i, j] = \delta_{ij} \sigma_i^{\text{unc}} \sigma_j^{\text{unc}} + \sum_{k=1}^{N_{\text{add}}} \sigma_{i,k}^{\text{add}} \sigma_{j,k}^{\text{add}} + \left( \sum_{k=1}^{N_{\text{mul}}} \sigma_{i,k}^{\text{mul}} \sigma_{j,k}^{\text{mul}} \right) D_i D_j, \quad (2.21)$$

where once again the  $D$  represent the experimental data points. This method of constructing the covariance matrix is therefore unambiguously defined by the experimental results. While a perfectly valid definition for analysing the description of data after a PDF determination, it is unreliable for use directly within a fitting procedure. The use of the experimental definition has for some time been understood to result in a *d’Agostini bias* [?]. That is, the theoretical values determined via a minimisation of a  $\chi^2$  function with the experimental covariance matrix are systematically shifted lower than the true value, an effect which only worsens as the number of data points subject to a common multiplicative error increases. The bias is generated by downward statistical fluctuations of data, if these low data points are used to generate the normalisation uncertainty, the result is a smaller uncertainty for the lower points, causing the fit to systematically undershoot the data.

The typical method employed to avoid the d’Agostini bias proceeds by including the normalisation as a fitted parameter and penalising large deviations as shown in Eqn. ???. This procedure largely corrects for the problem, although when applied to a dataset with several different normalisation uncertainties it still suffers from a bias. This effect was demonstrated by the NNPDF collaboration

in Ref. [?]. The bias can be avoided by using the so-called  $t_0$  prescription [?] for defining the covariance matrix. In this method the covariance matrix is constructed using the predictions from a previous fit rather than the experimental data values, to multiply with the multiplicative uncertainties.

$$\text{Cov}^{t_0}[i, j] = \delta_{ij} \sigma_i^{\text{unc}} \sigma_j^{\text{unc}} + \sum_{k=1}^{N_{\text{add}}} \sigma_{i,k}^{\text{add}} \sigma_{j,k}^{\text{add}} + \left( \sum_{k=1}^{N_{\text{mul}}} \sigma_{i,k}^{\text{mul}} \sigma_{j,k}^{\text{mul}} \right) T_i T_j, \quad (2.22)$$

where here the  $T$  are theory predictions for the associated datapoint, generated by some prior (fixed) PDF set. The prior, or  $t_0$  set should be determined self-consistently via an iterative procedure in which the  $t_0$  set is obtained from the previous result for the full fit. As the theory predictions are not subject to the same fluctuations as the data, the fit is not subject to the aforementioned bias. This effect can be seen explicitly in a fit to artificial pseudodata, performed with the experimental and  $t_0$  covariance matrix definitions in Figure ??.

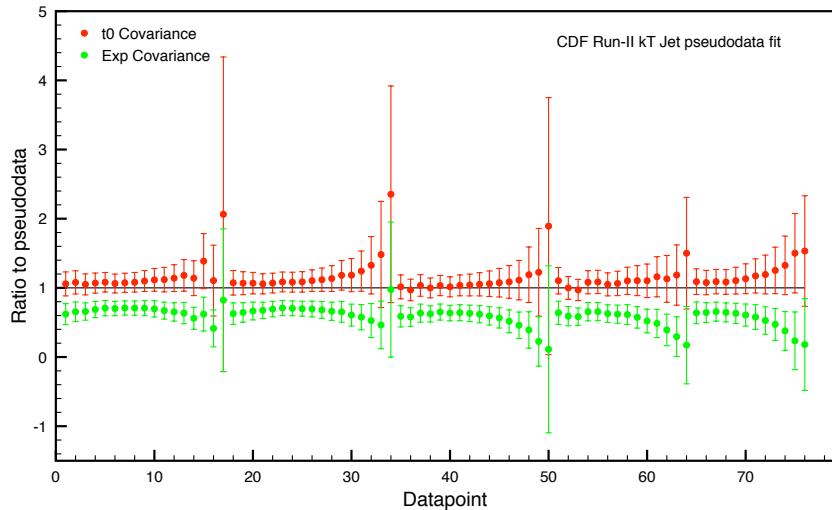


Figure 2.6: Demonstration of d’Agostini bias in a fit to pseudodata generated according to the kinematics of CDF inclusive jet data. Fit results are shown as a ratio to the ‘true’ value used to generate the pseudodata. The fit performed with the experimental definition of the covariance matrix results in predictions shifted systematically downwards with respect to the underlying law. The predictions from the fit using a  $t_0$  covariance matrix do not suffer from such a bias.

## Minimisation

With a figure of merit constructed, the PDF determination now becomes a problem of varying the free parameters in the PDF basis to minimise said measure. Even for those groups utilising a fixed parametrisation, performing a minimisation of the global  $\chi^2$  for a large,  $n \sim \mathcal{O}(1000)$  dataset with a fairly large number of free parameters (approximately 50 in the MSTW analysis once normalisation uncertainties are added as free parameters) is a challenging numerical task. For performing the minimisation, the MINUIT [?] package is a common choice, although other function minimisation methods are applied such as the Levenberg-Marquardt [?, ?] method as used in the MSTW fits.

In the NNPDF case the minimisation is complicated by the very large number of parameters and highly nonlocal behaviour in the error function, making conventional methods of minimisation difficult. These difficulties are overcome in the NNPDF methodology by the use of *genetic algorithms*, which are particularly efficient at exploring large parameter spaces. The implementation of the genetic algorithm is discussed in detail in Refs. [?, ?].

In addition to the basic difficulty of minimisation in a large parameter space, there is a further issue that arises when considering the fitting of a function with a great deal of redundant flexibility. Because of the flexibility of the parametrisation, it is possible that training the neural networks so that each reaches the global minimum in the error function actually results in the networks fitting to statistical noise. This effect is known as *overlearning* and is a problem often encountered in the training of large neural networks [?, ?]. In previous NNPDF determinations, the widely used *cross-validation* technique [?, ?] was employed in order to identify when overlearning occurred.

In this method the experimental data set is split into two separate sets. The first, a fitting set which is used for the minimisation of the error function, and a second validation set which is not used directly in the fitting procedure. For each

iteration in the genetic algorithm minimisation the error function is computed between the neural network predictions and both data sets. In the early stages of the training both error functions should decrease. However in the latter stages of the training where statistical noise begins to become an important contribution, the goodness-of-fit calculated to the fitting data set may continue to decrease while the same value calculated to the validation set has stopped decreasing or even begun to increase. This is a clear signal of overlearning, where fitting to statistical noise in the fitting set means that the fit to the validation set is no longer improving. At this stage the training of the neural networks is stopped. A typical signal of overlearning in a cross-validated fit can be seen in Figure ?? which compares the fit quality for both the training and validated sets over a number of fit iterations.

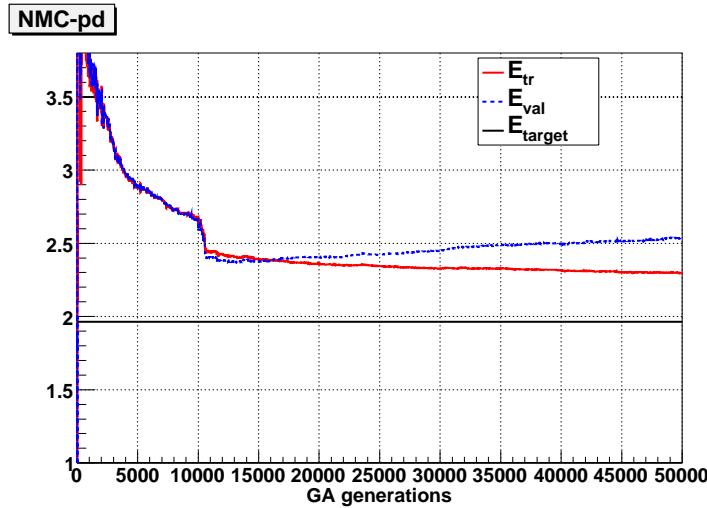


Figure 2.7: A typical signal of overlearning in a neural network fit.  $E_{\text{tr}}$  and  $E_{\text{val}}$  represent the training and validation figures of merit respectively. As the number of genetic algorithm generations proceeds, eventually the network begins to fit statistical noise in the training set and the validation fit quality begins to decrease. Figure from [?].

### 2.2.3 Error propagation

In order to undertake precision QCD studies, some estimate of the uncertainty on PDFs is required for a meaningful interpretation of the measured observables. The need for PDF sets with quantified uncertainties has been long recognised, and all modern determinations provide sets with at least experimental uncertainty estimation. While performing a comprehensive quantification of the theoretical uncertainty in a PDF fit is challenging, many methods have been developed in order to propagate the uncertainty from the dataset to the fitted PDFs. Ideally, one would like to determine a representation of the probability distribution in the whole functional space. That is given a dataset  $d$ , we would like to find the probability of a certain PDF candidate  $f$  such that our fitted PDF central value is given by

$$\langle f \rangle (x) = \int \mathcal{D}f f(x) \mathcal{P}(f|d), \quad (2.23)$$

and the uncertainty by

$$\text{Var}[f](x) = \int \mathcal{D}f [f(x) - \langle f \rangle(x)]^2 \mathcal{P}(f|d). \quad (2.24)$$

The probability distribution for an observable  $\mathcal{O}$  is then simply  $\mathcal{O}[f] \mathcal{P}(f|d)$ , in terms of which an observable's central value and PDF uncertainty can be calculated by

$$\langle \mathcal{O} \rangle = \int \mathcal{D}f \mathcal{O}[f] \mathcal{P}(f|d), \quad (2.25)$$

$$\text{Var}[\mathcal{O}] = \int \mathcal{D}f (\mathcal{O}[f] - \langle \mathcal{O} \rangle)^2 \mathcal{P}(f|d). \quad (2.26)$$

The probability distribution  $\mathcal{P}(f|d)$  is however a difficult quantity to determine. In this section we shall examine a number of the methods used in the literature to provide an estimate of PDF uncertainties.

### The Hessian method

The Hessian method is the most widely used method of uncertainty determination in PDFs. In essence, the method involves examining how the fit quality  $\chi^2$  varies when the  $n$  fit parameters  $a$  are perturbed about the values which minimise the  $\chi^2$ , here denoted by  $a^{\min}$ . A tolerance in the  $\chi^2$  variation is then chosen, and the error on an observable is determined geometrically from observables calculated with parameters perturbed by the selected tolerance. To examine this quantitatively, we first define the difference in  $\chi^2$  from the minimum value

$$\Delta\chi^2(a) \equiv \chi^2(a) - \chi^2(a^{\min}) = \sum_{i,j=1}^n H_{ij}(a_i - a_i^{\min})(a_j - a_j^{\min}), \quad (2.27)$$

where the  $a_i$  represent the  $i$ th component of the parameter set  $a$  (and likewise, for the minimised set  $a^{\min}$ ). Here we assume that the variation around the  $\chi^2$  minimum is approximately quadratic. The Hessian matrix  $H$  has values determined by

$$H_{ij} = \frac{1}{2} \left. \frac{\partial^2 \chi^2(a)}{\partial a_i \partial a_j} \right|_{\min}, \quad (2.28)$$

where the  $\min$  subscript refers to the parameters obtained at the  $\chi^2$  minimum. Early Hessian uncertainty estimates [?, ?] were based upon the standard formula for linear error propagation

$$(\Delta F)^2 = T^2 \sum_{i,j=1}^n \frac{\partial F}{\partial a_i} C_{ij} \frac{\partial F}{\partial a_j}, \quad (2.29)$$

where  $T^2 = \Delta\chi^2$  is the tolerance in  $\chi^2$  variation and  $C = H^{-1}$  is the inverse Hessian matrix. This procedure is however a little inconvenient due to the requirement of the partial derivatives of the observable with respect to the fit parameters. There are also numerical issues relating to this method which give rise to peculiar uncertainty estimates [?]. In order to overcome these issues the geometrical method outlined above was developed by the CTEQ

collaboration [?, ?].

For this method it is convenient to work in a rescaled orthogonal eigenbasis for the covariance matrix. The orthonormal eigenbasis is defined in the usual way

$$Hv_i = \lambda_i v_i, \quad (2.30)$$

and the rescaled eigenbasis is defined as  $e_i = 1/\sqrt{\lambda_i}v_i$ . The difference between a parameter set  $a$  and  $a^{\min}$  can now be expanded as

$$a_i - a_i^{\min} = \sum_{k=1}^n e_{ik} z_k, \quad (2.31)$$

where  $e_{ik}$  is the  $i$ th component of the  $k$ th rescaled eigenvector, and the  $z_k$  are the coefficients for the expansion of the parameter difference onto the rescaled eigenbasis. Therefore the expression for  $\Delta\chi^2$  reduces to

$$\Delta\chi^2(a) = \sum_{k=1}^n z_k^2 \quad \text{or,} \quad \chi^2(a) = \chi^2(a^{\min}) + \sum_{k=1}^n z_k^2. \quad (2.32)$$

This defines a hypersphere in the parameter space of radius  $\Delta\chi^2$  centred around  $a^{\min}$ , which corresponds to the variation in the parameters that is consistent with the tolerance  $T = \sqrt{\Delta\chi^2}$  in the quadratic approximation. It is now possible to construct an ensemble of  $2n$  PDF sets corresponding to the fits on the boundaries of the volume. A PDF set  $S_k^\pm$  therefore has the parameter set

$$a_i(S_k^\pm) = a_i^{\min} \pm t e^{ik}, \quad (2.33)$$

i.e. each parameter is perturbed by  $t$  in the direction of the  $e_k$  eigenvector. In the quadratic approximation  $t = T$ , when the approximation breaks down  $t$  can be determined by an iterative procedure to obtain the desired  $\Delta\chi^2$ . The error on

an observable  $F$  is then given simply by Pythagoras' theorem on the hypersphere

$$(\Delta F)^2 = \frac{1}{2} \sum_i^n (F(S_i^+) - F(S_i^-))^2. \quad (2.34)$$

In this procedure there is something of an ambiguity in the determination of the tolerance (and hence, the volume of the sphere in parameter space). Ideally the difference in  $\chi^2$  values should be exactly one for a confidence level of one-sigma<sup>1</sup>. In the case of PDF fits, this tolerance often leads to uncertainties far lower than expected. In practice, the CTEQ group uses a value of  $\Delta\chi^2 \sim 100$  and MSTW uses a value  $\sim 50$ . The more recent MSTW PDF sets have uncertainties calculated with a dynamically determined tolerance. More specialised fits such as ABM11 or the HERAPDF series, based upon relatively restrictive datasets may use the standard tolerance of  $\Delta\chi^2 = 1$ . Their use of a more restrictive dataset perhaps leading to fewer conflicts between experimental datasets that could require a more flexible tolerance.

The uncertainties produced via the Hessian procedure are difficult to analyse in a statistical sense due to the (occasional) inflation of the  $\Delta\chi^2$  and the approximations made in the procedure. It is therefore difficult to find a representation in the Hessian approach of the full probability distribution  $\mathcal{P}(f|d)$ . Furthermore the uncertainty in the choice of functional form, or estimation of parametrisation bias, is not typically take account of. The HERAPDF family of fits however do attempt to estimate this uncertainty by performing a series of fits with slightly modified parametrisations.

### Lagrange multiplier method

Another method of error propagation that has been explored is the Lagrange multiplier method. The method has the advantage of not assuming that the  $\chi^2$

---

<sup>1</sup>It should be noted that this is only the case when, either the data errors are uncorrelated, or when the correlations are included in the definition of the global goodness-of-fit  $\chi^2$  [?]

function is quadratic around the global minimum. We shall briefly discuss the method applied to the PDF error determination as suggested by Pumplin [?] and Stump [?]. A description of the process can also be found in [?, ?].

Firstly, a general global fit is performed to the data as described above. This yields a set of parameters  $a^{\min}$  which minimise the  $\chi^2$  measure. Using these parameters we calculate the best fit prediction for the observable in question  $F(a^{\min})$ . A new PDF fit can now be performed, where instead of minimising the  $\chi^2$  the following function is minimised

$$\Psi = \chi^2(a) + \lambda(F(a) - F(a^{\min})) \quad (2.35)$$

i.e. we introduce the observable  $F$  as a parameter in the fitting procedure and constrain the fit so that the minimal  $\Psi$  occurs when  $F(a) = F(a^{\min})$ . The value  $\lambda$  in this function is the Lagrange multiplier. The fit above is performed for many values of  $\lambda$ , each time leading to a parameter set that depends on that particular value of  $\lambda$ , this parameter set will be denoted  $a_\lambda$ . Using these parameters, we now calculate values for  $\chi^2(a_\lambda)$  and  $\mathcal{O}(a_\lambda)$ .

At this stage we now have a set of values for  $\chi^2(a_\lambda)$  and  $\mathcal{O}(a_\lambda)$  over a large range of  $\lambda$  values. This allows a determination of the relationship between the goodness-of-fit and the prediction for  $F$  via the parameter  $\lambda$ . We obtain an approximate function  $\chi^2(F)$  over a range of observable values, with a minimum at  $F = F(a^{\min})$  i.e  $\lambda = 0$  and  $a_\lambda = a^{\min}$ . Also we have a set of the  $a$  parameters for every point on the curve which are optimised for the best fit to the observable  $F$ . This means that we have a set of fully optimised parameters for any arbitrary confidence level determined by the  $\Delta\chi^2$  that we select as a tolerance. Uncertainties for the PDFs can therefore be given in a way that utilises the whole of the  $a$  parameter space, rather than just perturbing around the global minimum as in the Hessian approach.

Of course, the disadvantage of this method is that the PDF uncertainties

must be calculated for each observable in a rather computationally intensive process. The errors are naturally optimised for the particular observable, but the process is inconvenient for a PDF end-user, and so it is not widely-used in error determination. In this sense the Lagrange multiplier approach suggests a method of estimating  $\mathcal{P}(\mathcal{O}|a)$ , or the probability density of an observable in the space of parameters. The Lagrange multiplier method also relies on the same somewhat arbitrary choice of tolerance in  $\chi^2$  as the Hessian method. The method has however been applied as a cross-check to the Hessian results [?, ?].

### Monte Carlo method

Another quite distinct method of PDF uncertainty determination is the Monte Carlo method, first suggested by Giele and Keller [?, ?] where a Monte Carlo procedure in the space of fit parameters was outlined. The NNPDF collaboration uses a similar method in all of its fits, although with the Monte Carlo performed in the space of experimental data. The method is designed to faithfully represent the uncertainties present in the initial data, and to propagate the errors in a way that does not assume anything of the nature of the error propagation. The Monte Carlo approach was also analysed and compared to the results of a Hessian fit by the MSTW group in [?].

In the Monte Carlo procedure an ensemble of  $N_{\text{rep}}$  artificial data replicas is produced for every data point in the fit, generated according to the probability distribution of the initial data. Typically this distribution is multi-Gaussian with central values and variances provided by experimental results, but any probability distribution may be used if required. If we use  $F_p^{(\text{art})(k)}$  to represent a single element  $k$  of the pseudo-data sample (the *art* superscript designates the data as an artificial sample) of the observable  $F$  at the kinematical point  $\{x_p, Q_p^2\}$ . Then

we can generate such a pseudo-data element as in [?] by

$$F_p^{(\text{art})(k)} = S_{p,N}^{(k)} F_p^{(\text{exp})} \left( 1 + \sum_{l=1}^{N_c} r_{p,l}^{(k)} \sigma_{p,l} + r_p^{(k)} \sigma_{p,s} \right), \quad (2.36)$$

where the  $r$  are independent Gaussian random numbers centred upon the experimental central value. The  $\sigma_{p,s}$  term contains the uncorrelated systematic uncertainties and the statistical uncertainty added in quadrature. The  $\sigma_{p,l}$  are the correlated errors for the data provided by experiment. The normalisation of the probability distribution is fixed by the term  $S_{p,N}$ . Provided a large enough quantity of these artificial replicas ( $N_{\text{rep}}$ ) is generated, this form of the generating distribution for the Monte Carlo ensemble reproduces all of the statistical qualities of the original experimental data. In Ref. [?] it is demonstrated that  $N_{\text{rep}} = 1000$  is sufficient to reproduce the experimental central values and variances to an accuracy of better than one percent.

Now that a good Monte Carlo sample of the experimental data is available, instead of performing just the one fit to the data,  $N_{\text{rep}}$  independent fits are performed, one for each of the data replicas. At the end of the fitting procedure we obtain an ensemble of  $N_{\text{rep}}$  equally probable PDFs which reliably describe the probability distribution of the PDFs based upon the original experimental uncertainties. The central values and uncertainties of an observable can be simply obtained by computing the average and the variance over the ensemble of PDFs.

$$\langle F \rangle = \frac{1}{N_{\text{rep}}} \sum_i^{N_{\text{rep}}} F^{(k)}, \quad (2.37)$$

$$\sigma^2[F] = \frac{1}{N_{\text{rep}} - 1} \sum_{i=1}^{N_{\text{rep}}} (F^{(k)} - \langle F \rangle)^2, \quad (2.38)$$

where  $F^{(k)}$  denotes the observable  $F$  computed using PDF replica  $k$ .

The Monte Carlo method therefore propagates the errors from the experi-

mental data through to the PDFs in a natural way, without the need for a linear propagation of errors assumption, or the need for an inflated tolerance in the  $\chi^2$  distribution. Figure ?? demonstrates a Monte Carlo ensemble of PDF replicas for the gluon distribution.

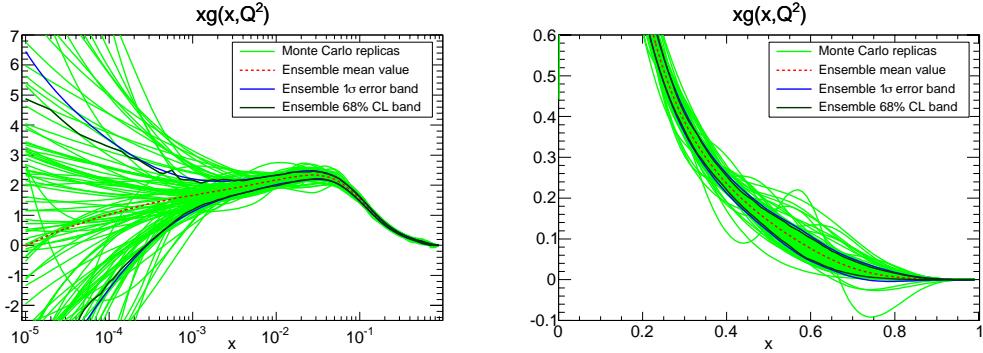


Figure 2.8: A Monte Carlo representation of the gluon PDF probability distribution. Individual PDF replicas are shown as green lines, and the ensemble average, standard deviation and 68% confidence level are shown.

## 2.3 Status of PDF determination before the LHC

In preparation for the application of parton distributions at the LHC, extensive studies were performed in order to benchmark and understand areas of agreement and discrepancy across fitting collaborations [?, ?]. While agreement had generally improved as the level of sophistication applied in parton fits increased, there were still notable regions where PDF fits from the widest datasets remained in disagreement at levels greater than their quoted uncertainties. Figure ?? illustrates the situation for two important PDF luminosities before the LHC. These discrepancies extended not only to so far unmeasured quantities such as Higgs production cross sections, but also to PDF standard candle observables such as  $W$  boson production (c.f. Figure ??).

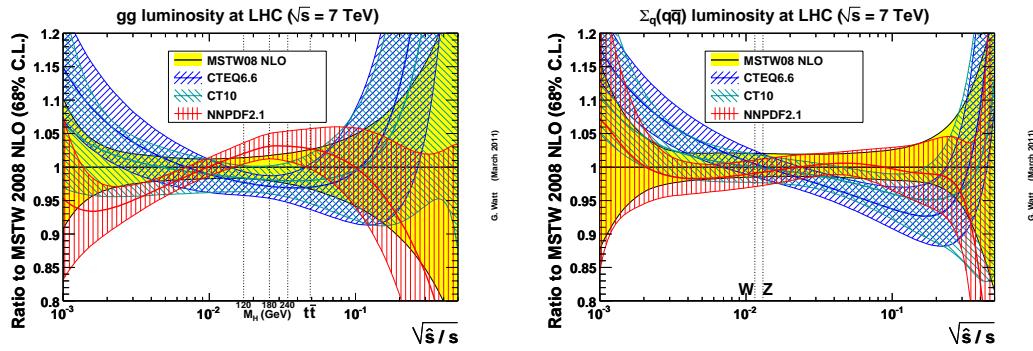


Figure 2.9: Luminosities for  $gg$  (left) and  $q\bar{q}$  (right) PDF combinations at the 7 TeV LHC. Figure from [?].

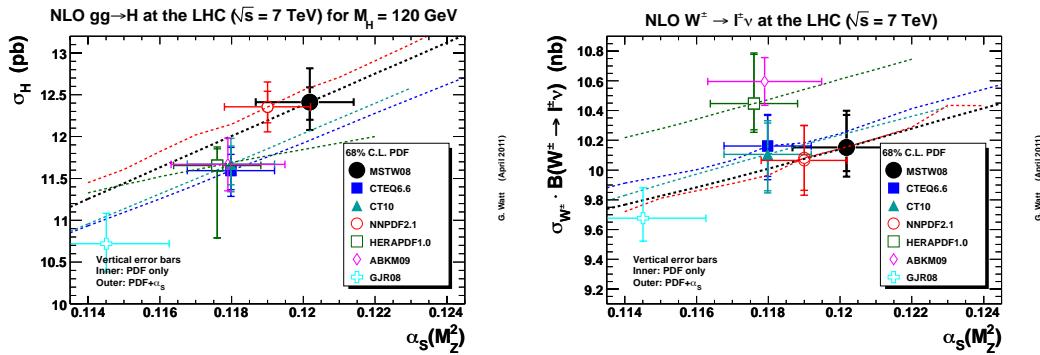


Figure 2.10: Predictions for LHC processes based upon a number of PDF determinations. Left figure: cross section for Higgs production in gluon fusion. Right figure: cross section for the production of  $W$  bosons. Figure from [?].

The Les Houches benchmark exercise [?] helped to elucidate the methodological source of many of these differences by testing fits from various methodologies to a standard dataset.

Many of the observed discrepancies arise due to differences in the theoretical description of data, with the choice of flavour number scheme providing the largest differences. Dataset choice and methodological choices introducing significant differences also. These differences led to the conservative PDF4LHC recommendation for observables to be calculated as the central contour of the CTEQ-MSTW-NNPDF uncertainty envelope. Despite the differences, for the LHC Run-I the range of available sets allowed for experimental collaborations to

effectively explore the differences in the resulting predictions.

While providing accurate determinations for use at the LHC has been the primary concern in the years leading up to the LHC’s first operation, there was substantial interest in the potential of the LHC to provide constraints upon PDFs and potentially provide discriminating power between sets. Data from the LHC provides the best opportunity for distinguishing the most effective approaches both theoretically and methodologically. Additionally LHC data provides particularly valuable input in the field of collider-only determinations, which aim to provide a cleaner description of data by avoiding the inclusion of nuclear-corrected and low energy data. The inclusion of a large LHC dataset into PDF fits is however a challenging problem, and one which has inspired a great deal of progress in the efficient calculation of collider observables. The remainder of this work will therefore deal with the both the technical inclusion of LHC data into parton distribution fits and the subsequent phenomenological results.

# Chapter 3

## Tools for the LHC

Including data from a wide range of LHC or collider sources into a global PDF determination provides several challenges, particularly in the context of the computationally intensive NNPDF methodology. In this chapter we shall discuss some of the methods that have been developed in order to study the impact of collider data, and include their constraints into PDF fits.

Firstly we shall describe the method of Bayesian reweighting of Monte Carlo error sets, along with the associated set of tools made available by a Bayesian study of PDF sets and their uncertainties. Secondly the FastKernel method developed by the NNPDF collaboration for the fast evolution of PDFs will be introduced, along with its extension to the fast computation of experimental observables in the FK method. Finally we shall discuss the application of interpolation methods such as FastKernel to the automated calculation of cross sections at next-to-leading order accuracy in QCD. To this end we shall perform a brief overview of such calculations in the context of general purpose event generators.

### 3.1 Bayesian reweighting

When examining the statistical properties of PDF fits it is important to note that in the Monte Carlo approach, not only the uncertainties on PDFs are provided, but a full representation of the probability distribution. As described in Section ??, the integral over the PDF probability distribution is approximated by a sum over replicas,

$$\begin{aligned} \langle f \rangle(x, Q^2) &= \int f(x, Q^2) \mathcal{P}(f|d) \mathcal{D}f \\ &\approx \frac{1}{N_{\text{rep}}} \sum_i^{N_{\text{rep}}} f_i(x, Q^2), \end{aligned} \quad (3.1)$$

where the subscript  $i$  here refers to the PDF replica in the Monte Carlo ensemble. This correspondence leaves PDFs in the Monte Carlo representation open to standard statistical analysis methods. One of the most important of which is the *Bayesian reweighting* technique, first proposed by Giele and Keller alongside the original Monte Carlo procedure [?] and then subsequently developed by the NNPDF collaboration [?, ?]. The problem that reweighting seeks to address is the rapid addition of experimental data into an existing parton determination. The method is particularly useful in cases where there are no fast implementations of a calculation, and allows for the fast assessment of experimental impact upon PDFs and their uncertainties.

Given a probability distribution for PDFs, Bayes' theorem suggests that we can update the experimental information in an existing PDF fit, here denoted  $\mathcal{P}(f)$  by determining the conditional probability of the PDF given the new dataset  $y$ ,

$$\mathcal{P}(f|y) \mathcal{D}f = \frac{\mathcal{P}(y|f)}{\mathcal{P}(y)} \mathcal{P}(f) \mathcal{D}f. \quad (3.2)$$

However it was noted in Ref. [?] that the probability of a PDF given the new data is not strictly what a fitting procedure would obtain. Rather the fitting procedure

aims to find the probability distribution of the PDFs given some measure of fit quality to the new data, e.g  $\chi^2$ . Therefore to obtain a distribution statistically equivalent to a refit, one should attempt to determine

$$\mathcal{P}(f|\chi)\mathcal{D}f = \frac{\mathcal{P}(\chi|f)}{\mathcal{P}(\chi)} \mathcal{P}(f)\mathcal{D}f, \quad (3.3)$$

where  $\mathcal{P}(\chi)$  may be marginalised over to obtain the correct normalisation for  $\mathcal{P}(f|\chi)$ . Armed with such a distribution, we may then compute our predictions for a general observable given the information contained in the new dataset,

$$\begin{aligned} \langle \mathcal{O} \rangle_{\text{new}} &= \int \mathcal{O}[f] \mathcal{P}(f|\chi) Df \\ &= \int \mathcal{O}[f] \frac{\mathcal{P}(\chi|f)}{\mathcal{P}(\chi)} \mathcal{P}(f) Df, \end{aligned}$$

where  $\langle \mathcal{O} \rangle_{\text{new}}$  is the central value prediction for the observable  $\mathcal{O}$  provided by a PDF distribution updated with the new experimental data. Given this probability distribution we can form a Monte Carlo representation in terms of PDF replicas once again,

$$\begin{aligned} \langle \mathcal{O} \rangle_{\text{new}} &= \frac{1}{N_{\text{rep}}} \sum_{i=1}^{N_{\text{rep}}} \frac{\mathcal{P}(\chi|f_i)}{\mathcal{P}(\chi)} \mathcal{O}[f_i], \\ &= \frac{1}{N_{\text{rep}}} \sum_{i=1}^{N_{\text{rep}}} w_i \mathcal{O}[f_i]. \end{aligned} \quad (3.4)$$

The weights  $w_i$  for the individual replicas encoding the information from the new dataset, may be obtained from the  $\chi^2$  goodness-of-fit measure to the new data

$$w_i = \frac{\mathcal{P}(\chi|f_i)}{\mathcal{P}(\chi)} \propto \chi_i^{n-1} e^{-\frac{1}{2}\chi_i^2}. \quad (3.5)$$

Where  $n$  denotes the number of new datapoints. The new data may therefore be included into an existing MC parton set by the simple calculation of a  $\chi^2$  for each

replica in the set. In comparison to a fitting procedure where many thousands of  $\chi^2$  computations are required, this procedure is extremely fast. Furthermore, as a purely statistical exercise this PDF reweighting does not suffer from any of the inherent vagaries of a fitting procedure.

The reweighting technique does however come at a cost in that it may reduce the overall efficiency of the Monte Carlo ensemble's representation of the underlying probability distribution. As can be seen from Eqn. ??, replicas in the prior distribution which do not provide a good description of the new experimental data and therefore have a large  $\chi^2$  value are penalised by small weights. For a sufficiently large or constraining dataset this can mean that many of the replicas are effectively switched out of the distribution, leaving a smaller number of *effective* replicas. The efficiency of the representation can be quantified by the Shannon entropy, which provides the number of effective replicas as

$$N_{\text{eff}} \equiv \exp \left( \frac{1}{N_{\text{rep}}} \sum_{i=1}^{N_{\text{rep}}} w_i \ln(N_{\text{rep}}/w_i) \right). \quad (3.6)$$

As the constraining power of the new dataset increases, so the Shannon entropy  $N_{\text{eff}}$  decreases. Consequently a larger number of replicas sampling the prior distribution are required to maintain a fixed level of ensemble accuracy. Despite this limitation, reweighting can provide an extremely useful method for analysis of a typical experimental dataset.

### 3.1.1 Error rescaling parameter

A Bayesian analysis of the Monte Carlo probability representation opens up other avenues of investigation. Of particular interest is the examination of an *error rescaling parameter*. When examining the impact of an experimental measurement, we can study how the constraints are modified under a global

rescaling of the experimental error, i.e the  $\chi^2$  values

$$\chi_k^2 \rightarrow \chi_{k,\alpha}^2 = \chi_k^2/\alpha^2, \quad (3.7)$$

where  $\alpha$  is the rescaling parameter. In our reweighting exercise the weights are subsequently given by

$$w_k(\alpha) \propto (\chi_{k,\alpha}^2)^{(n-1)/2} e^{-\chi_{k,\alpha}^2/2}. \quad (3.8)$$

Our Bayesian expression for the updated probability density is now also a function of the rescaling parameter  $\alpha$ . A further application of Bayes' theorem inverts this relationship, and allows us to form a probability density for the rescaling parameter itself.

$$\mathcal{P}(\alpha|\chi^2) \propto \frac{1}{\alpha} \sum_{k=1}^N w_k(\alpha). \quad (3.9)$$

This probability distribution provides an estimate as to whether the experimental errors in the new dataset may have been under or overestimated, based upon agreement with the prior distribution. An experimental result where the uncertainties have accurately estimated leads to a  $\mathcal{P}(\alpha)$  distribution peaked at  $\alpha = 1$ , whereby an over(under)-estimated set of uncertainties leads to a lower(higher) peak in the distribution. This is a particularly useful tool for analysing experimental uncertainties, and can provide some differentiation between inconsistent and constraining data in cases where  $N_{\text{eff}}$  is small.

### 3.1.2 PDF unweighting

While the PDF reweighting approach is a powerful method for the addition of new data to an existing set, a reweighted PDF set is unsuitable for general distribution. For use in typical calculational codes, a standard interface is required through packages such as LHAPDF. Therefore the provision of a PDF ensemble with

an associated set of weights would require the retooling of codes in which a reweighted calculation is desired. To alleviate this a method was developed in order to present a reweighted distribution as a standard MC replica ensemble [?].

This is done by representing the reweighted set upon a cumulative line of weights as in Figure ???. Each line segment corresponds to the weight of an individual replica. The total cumulant line therefore being normalised to  $N_{\text{rep}}$ , the number of replicas in the reweighted distribution. Replicas in an ‘unweighted’ set are then chosen by distributing evenly  $N'_{\text{rep}}$  replicas across this cumulant line. When one of these replicas falls into the weight segment of a corresponding reweighted replica, that PDF is selected for inclusion in the unweighted set. Importantly, the same reweighted replica may be selected more than once to appear in the unweighted set.

As an example, consider the case where there are four replicas in an initial distribution, with weights  $w_i = \{1, 2, 3, 4\}$ . The cumulant line formed by these weighted replicas is shown on the left side of Figure ???. This line is subdivided into  $N'_{\text{rep}} + 1$  intervals. With  $N'_{\text{rep}} = 20$  as shown in the Figure, two unweighted replicas fall in the first weighted segment, three in the second, six in the third and nine in the fourth. Therefore the unweighted ensemble is formed by duplicating the original weighted replicas with a frequency dictated by how many unweighted replicas fall in their respective line segment.

The weights of the original set are therefore approximately represented as replica multiplicities in the unweighted set, with low-weight replicas selected few times (if at all), and large weight replicas selected multiple times. In this way a conventional MC ensemble can be formed with the usual LHAPDF interface, this time including duplicate replicas for those with high weights and excluding replicas with weights that fall under the unweighted set’s resolution. Therefore the unweighting procedure can provide an exact representation of the reweighted ensemble in the limit  $N'_{\text{rep}} \rightarrow \infty$ .

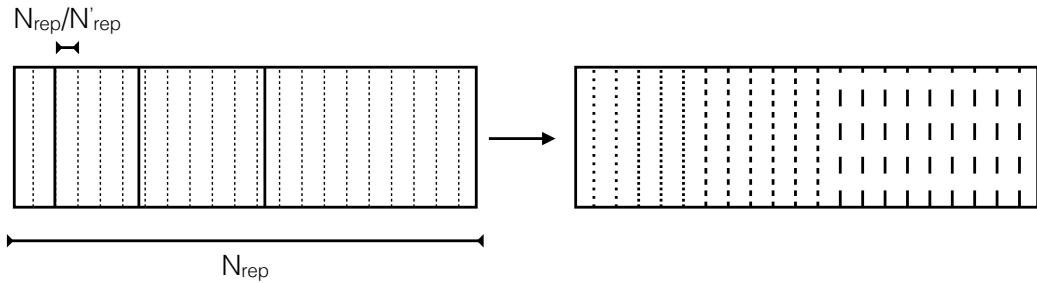


Figure 3.1: The unweighting of a Bayesian reweighted Monte Carlo PDF set. The left hand figure shows the weight cumulant segments for the original weighted set, with four replicas of weight  $w_i = \{1, 2, 3, 4\}$ . The line is subdivided by  $N'_{\text{rep}} = 20$  lines. The right hand figure illustrates the *unweighted* set in this case. Here each replica in the unweighted set has equal weight, with different line strokes denoting different replicas from the weighted distribution.

However in practice a number of unweighted replicas of the order of the number of effective replicas  $N_{\text{eff}}$  is typically sufficient for a good level of accuracy in the reproduction.

### 3.1.3 Reweighting validation

The Bayesian reweighting procedure has been extensively validated by the NNPDF collaboration in a number of highly non-trivial tests of the methodology. As the method has been designed to update a prior distribution with new information analogously to the approach used in an ideal fit, the first test is to ensure that a PDF set reweighted with a new dataset is statistically equivalent to a new set refitted from scratch utilising the new data. This was first performed in [?] by reweighting an NNPDF 2.0 fit which included only DIS and Drell-Yan data with information from Tevatron inclusive jet measurements. The reweighted set was compared to the full NNPDF 2.0 fit including the data. As Figure ?? demonstrates, the reweighted set is able to reproduce the refitted set up to the level of statistical fluctuation.

The development of the unweighting method as outlined in the previous section, allowed for further tests of the reweighting method. A series of tests were

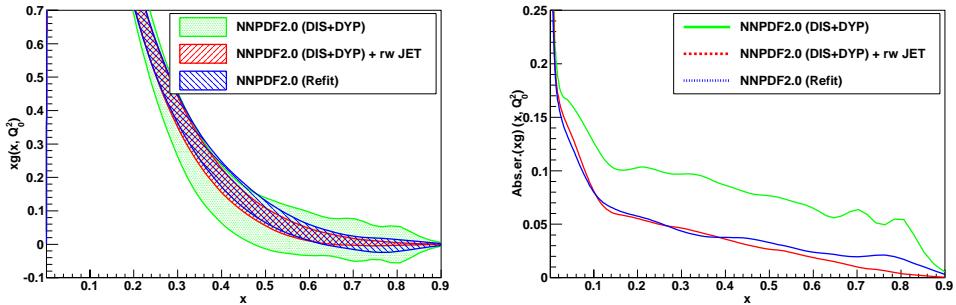


Figure 3.2: The validation of Bayesian reweighting by the inclusion of Tevatron jet data. The left figure demonstrates the prior distribution along with the reweighted and refitted distributions upon the addition of Tevatron jet data. The right plot shows the absolute error upon the PDFs for the three sets. Figures are from [?].

carried out in order to assess the behaviour of PDFs under successive reweighting operations.

When attempting to include multiple datasets into a PDF fit via reweighting, there are three possibilities. One can reweight with the combined  $\chi^2$  values for the two experiments, or one can reweight first with one experiment, unweight the PDF ensemble, then reweight with the second. The resulting PDFs should be reasonably independent of the method chosen, and of the order in which the successive reweighting is performed. This requirement is a stringent test of the Monte Carlo PDF representation, as it determines whether or not the ensemble truly behaves as a probability distribution. More pragmatically, the test verifies whether the loss of ensemble efficiency in one reweighting operation is not so great as to prevent a further reweighing. This investigation was carried out in Ref. [?] with a DIS only prior. The E605 Drell-Yan experiment and CDF/D0 inclusive jet measurements were included into this set by reweighting. As the E605 experiment provides global fits with rather stringent constraints compared to the moderate effect of the jet data, this is a rather asymmetrical and therefore effective test.

In Figure ?? these reweighting procedures are compared for the case of the

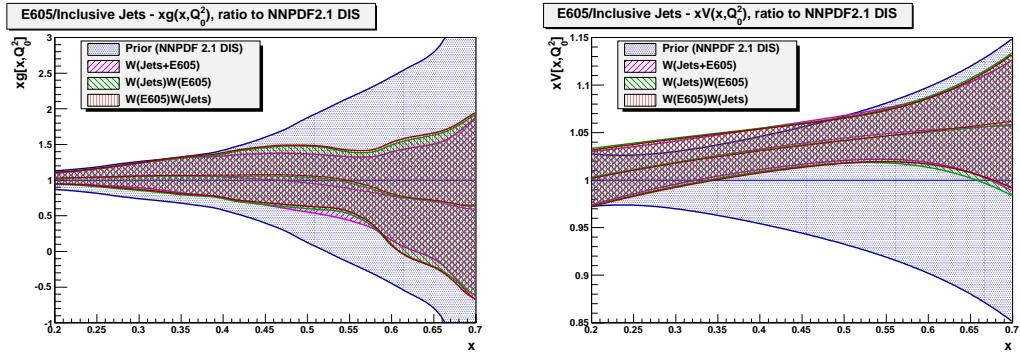


Figure 3.3: Test of Bayesian reweighting under a successive reweighting operation. Inclusive Jet and Drell-Yan data are added as a combined dataset, and as individual reweightings separated by an unweighting operation. The resulting distributions, the gluon PDF on the left and the valence distribution on the right, show excellent agreement between the different procedures. All curves are normalised to the prior, NNPDF 2.1 DIS result. The figures are taken from [?].

gluon and valence distributions of the NNPDF2.0 DIS only fit. It is clear that while the impact of the data upon the prior is substantial, the three reweighting methods hardly differ in their results. There is therefore a strong confirmation of the statistical properties of both the Monte Carlo representation of PDFs, and of the reweighting method.

## 3.2 The FastKernel method

The method of Bayesian reweighting provides an extremely fast and efficient method of including new data into a determination. However as described previously, the method is ill-suited to the addition of a large or very constraining dataset as the required size of the prior distribution in replicas rapidly becomes unmanageable. Therefore the standard fitting methodology remains the most important procedure in the determination of parton distributions.

The primary issue in the standard methodology upon the addition of a large LHC dataset is the computational time required to perform the theoretical predictions for experimental data. Not only must the standard double convolution

over the two parton densities be performed, but also each PDF must be evolved from some initial fitting scale to the scale of the experimental data by yet another set of convolutions. We shall first describe the methods used for fast PDF evolution, before going on to discuss the extension to the calculation of physical observables at colliders.

### 3.2.1 Fast PDF evolution

While there are many methods for performing the evolution of parton distributions, the technique used in NNPDF fits must be particularly efficient due to the computational complexity of the NNPDF procedure. The evolution of a flavour basis PDF of flavour  $i$  from an initial scale  $Q_0^2$  to a target scale  $Q_\tau^2$  can be expressed as

$$f_i(x_\alpha, Q_\tau^2) = \sum_j^{N_f} \int_{x_\alpha}^1 d\xi \Gamma_{ij} \left( \frac{x_\alpha}{\xi}, \frac{Q_\tau^2}{Q_0^2} \right) f_j(\xi, Q_0^2), \quad (3.10)$$

where the  $\Gamma$  are found by solution of the DGLAP equation as shown in Eqn. ???. In order to take advantage of the sparse nature of the DGLAP evolution kernels, we work in the evolution basis defined in Section ??,

$$N_i(x_\alpha, Q_\tau^2) = \sum_j^{N_f} \int_{x_\alpha}^1 d\xi \tilde{\Gamma}_{ij} \left( \frac{x_\alpha}{\xi}, \frac{Q_\tau^2}{Q_0^2} \right) N_j(\xi, Q_0^2), \quad (3.11)$$

where here we have introduced the notation  $N$  for the evolution basis PDFs; related to the flavour basis by a simple rotation

$$f_i(x, Q_\tau^2) = \sum_j^{N_f} R_{ij} N_j(x, Q_\tau^2). \quad (3.12)$$

Having to perform many instances of the convolution integral in Eqn. ?? would be prohibitively expensive in most fitting applications, and so an alternative approach must be used. In the NNPDF framework this is based upon the

**FastKernel** interpolation method introduced in Ref. [?], and shares the general approach with other interpolating methods, while maintaining a hybrid  $x$  and Mellin space solution. Here we shall outline the general method used in all interpolating tools.

The first step is to expand the initial-state PDFs upon some set of interpolating basis functions  $\mathcal{I}$ ,

$$f_i(x, Q_0^2) \approx \sum_{\beta}^{N_{\text{fn}}} c_i^{(\beta)} \mathcal{I}^{(\beta)}(x), \quad (3.13)$$

with the coefficients of this expansion calculable through the usual overlap integral

$$c_i^{(\beta)} = \int_0^1 dx f_i(x, Q_0^2) \mathcal{I}^{(\beta)}(x). \quad (3.14)$$

Substituting the interpolated version of the initial state PDF into the evolution equation and applying the inverse transformation of Eqn. ?? to work in the evolution basis we obtain

$$N_i(x_\alpha, Q_\tau^2) = \sum_{j,k}^{N_f} \sum_{\beta}^{N_{\text{fn}}} \int_{x_\alpha}^1 d\xi \tilde{\Gamma}_{ij} \left( \frac{x_\alpha}{\xi}, \frac{Q_\tau^2}{Q_0^2} \right) R_{jk}^{-1} c_k^{(\beta)} \mathcal{I}^{(\beta)}(x). \quad (3.15)$$

In this expression, we can actually factorise the PDF-dependent expansion coefficients  $c$  from the integral, and perform the convolution over the interpolating functions

$$N_i(x_\alpha, Q_\tau^2) = \sum_k^{N_f} \sum_{\beta}^{N_{\text{fn}}} E_{ik\alpha\beta}^\tau c_k^{(\beta)}, \quad (3.16)$$

where the *evolution tables*  $E$  are given by

$$E_{ik\alpha\beta}^\tau = \sum_j^{N_f} \int_{x_\alpha}^1 d\xi \tilde{\Gamma}_{ij} \left( \frac{x_\alpha}{\xi}, \frac{Q_\tau^2}{Q_0^2} \right) R_{jk}^{-1} \mathcal{I}^{(\beta)}(x). \quad (3.17)$$

While it may seem as if we have simply moved the problem from the convolution

with the DGLAP kernel to the overlap integral required to compute the coefficients  $c$ , this can be avoided via a careful choice in the interpolating functions. A suitable choice of interpolating function yields the following identification for the coefficients

$$c_k^{(\beta)} = f_k(x_\beta, Q_0^2), \quad (3.18)$$

that is, the interpolants effectively pick out the value of the PDF at some point  $\beta$  in an  $x$ -grid. Providing the grid in  $\beta$  is dense enough the interpolation accuracy can still be very high. With such a choice of functional basis, the full evolution product becomes particularly simple

$$f_i(x_\alpha, Q_\tau^2) = \sum_{j,k}^{n_f} \sum_{\beta}^{N_x} R_{ij} E_{\alpha\beta jk}^\tau f_k(x_\beta, Q_0^2), \quad (3.19)$$

$$= \sum_k^{n_f} \sum_{\beta}^{N_x} A_{\alpha\beta ik}^\tau f_k(x_\beta, Q_0^2). \quad (3.20)$$

The convolution required by the initial solution to the DGLAP equation has now been reduced via interpolation methods to a simple product over a rotated evolution table  $A$ .

### 3.2.2 Fast calculation of collider observables

Similar methods to what we have discussed for fast PDF evolution have also been applied to the calculation of collider observables. For a typical observable with two partons in the initial state, a full calculation is given by a double convolution over two parton densities,

$$\sigma_{pp \rightarrow X} = \left( \frac{\alpha_s(Q^2)}{2\pi} \right)^p \int dx_1 dx_2 f_i(x_1, Q^2) d\hat{\sigma}_{ij \rightarrow X} f_j(x_2, Q^2). \quad (3.21)$$

The double convolution can once again be avoided by inserting interpolated versions of the PDFs, and performing the convolution over the interpolating functions.

$$\sigma_{pp \rightarrow X} = \left( \frac{\alpha_s(Q^2)}{2\pi} \right)^p \sum_{\alpha,\beta}^{N_X} f_i(x_\alpha, Q^2) W_{\alpha\beta,ij} f_j(x_\beta, Q^2). \quad (3.22)$$

The weight grid  $W$  is calculated analogously to the evolution tables in Eqn. ??,

$$W_{\alpha\beta,ij} = \int dx_1 dx_2 \mathcal{I}^{(\alpha)}(x_1) d\hat{\sigma}_{ij \rightarrow X} \mathcal{I}^{(\beta)}(x_2). \quad (3.23)$$

Identical methods can be used to interpolate over the hard scale  $Q^2$  in multi-scale processes. These techniques are used in publicly available tools such as **APPLgrid** [?] and **FastNLO** [?]. In the **APPLgrid** framework, the full product used to calculate a hadronic observable is

$$\sigma = \sum_p \sum_s \sum_{\alpha,\beta}^{N_{\text{sub}}} \sum_{\tau}^{N_x} \sum_{\tau}^{N_Q} W_{\alpha\beta\tau}^{(p)(s)} \left( \frac{\alpha_s(Q_\tau^2)}{2\pi} \right)^p F^{(s)}(x_\alpha, x_\beta, Q_\tau^2), \quad (3.24)$$

where the interpolation over a grid of points in hard scale runs over the index  $\tau$ , and the perturbative order of the contributions is separated by the index  $p$ . The initial state parton combinations have been grouped into the appropriate QCD subprocesses  $s$ , according to a table of coefficients  $C$ ,

$$F^{(s)}(x_\alpha, x_\beta, Q_\tau^2) = \sum_{i,j}^{13} C_{ij}^{(s)} (f_i(x_\alpha, Q_\tau^2) f_j(x_\beta, Q_\tau^2)). \quad (3.25)$$

The resulting product in Eqn ?? allows for the simple variation of PDFs, strong coupling and perturbative scales in a fast calculation; the product taking typically of order milliseconds rather than the hours to days required to obtain reliable statistics in an NLO code.

Despite the dramatic speed improvement, the **APPLgrid**/ **FastNLO** products

represent a considerable computational expense when introducing a large dataset. The NNPDF methodology in particular is extremely sensitive to the convolution speed due to the nature of the genetic algorithm minimisation, orders of magnitude more convolutions are required than in competing approaches. Therefore in order to practically include a large collider dataset into an NNPDF fit more work must be done on improving the convolution algorithm.

### 3.2.3 Combined evolution and observable calculation

The `APPLgrid/` `FastNLO` approach maintains a great deal of flexibility, in that scale,  $\alpha_S$  and PDF variations are all possible within the same framework. In a PDF fit the only requirement is an efficient variation of input parton distributions. We can therefore try to improve the efficiency of the calculation at the cost of some of the flexibility available in the fast convolution tools. The FK procedure and toolchain has therefore been developed, implementing a combined PDF evolution and collider observable calculation.

Recalling the fast PDF evolution method in Eqn. ?? with the suitable grids precomputed, PDF evolution can be performed simply as

$$f_i(x_\alpha, Q_\tau^2) = \sum_k^{n_f} \sum_\beta^{N_x} A_{\alpha\beta ik}^\tau f_k(x_\beta, Q_0^2). \quad (3.26)$$

The evolution of the `APPLgrid` subprocess in Eqn. ?? from an initial state distribution is therefore

$$F^{(s)}(x_\alpha, x_\beta, Q_\tau^2) = \sum_{i,j}^{13} \sum_{k,l}^{n_f} \sum_{\delta,\gamma}^{N_x} C_{ij}^{(s)} [A_{\alpha\delta ik}^\tau f_k(x_\delta, Q_0^2) A_{\beta\gamma jl}^\tau f_l(x_\gamma, Q_0^2)] \quad (3.27)$$

$$= \sum_{k,l}^{n_f} \sum_{\delta,\gamma}^{N_x} \tilde{C}_{kl,\alpha\beta\gamma\delta}^{(s),\tau} f_k(x_\delta, Q_0^2) f_l(x_\gamma, Q_0^2), \quad (3.28)$$

where the evolved subprocess coefficients are

$$\tilde{C}_{kl,\alpha\beta\gamma\delta}^{(s),\tau} = \sum_{i,j}^{13} C_{ij}^{(s)} A_{\alpha\delta ik}^\tau A_{\beta\gamma jl}^\tau. \quad (3.29)$$

Substituting the expression for the subprocess in terms of initial state PDFs, Eqn. ??, into the `APPLgrid` expression for the full convolution shown in Eqn. ?? we obtain

$$\sigma = \sum_p \sum_s \sum_{\alpha,\beta}^{N_{\text{sub}}} \sum_{\tau}^{N'_x} \sum_{\tau}^{N_Q} W_{\alpha\beta\tau}^{(p)(s)} \left( \frac{\alpha_s(Q_\tau^2)}{2\pi} \right)^p \sum_{k,l}^{n_f} \sum_{\delta,\gamma}^{N_x} \tilde{C}_{kl,\alpha\beta\gamma\delta}^{(s),\tau} f_k(x_\delta, Q_0^2) f_l(x_\gamma, Q_0^2). \quad (3.30)$$

where the number of points in the `APPLgrid`  $x$ -grid is denoted  $N'_x$  to indicate that the grid is different to the input parton  $x$ -grid which runs over  $\gamma, \delta$  up to  $N_x$  points. Now that the PDF evolution has been factorised into the coefficients  $\tilde{C}$ , much more of this sum may now be precomputed. Specifically we are now able to sum over the indices for subprocess  $s$ , perturbative order  $p$ , hard scale  $\tau$ , and the `APPLgrid`  $x$ -grids  $\alpha$  and  $\beta$ . The resulting expression for the combined evolution and observable calculation is therefore

$$\sigma = \sum_{k,l}^{n_f} \sum_{\delta,\gamma}^{N_x} \tilde{W}_{kl\delta\gamma} f_k(x_\delta, Q_0^2) f_l(x_\gamma, Q_0^2), \quad (3.31)$$

with the combined grid, which may be precomputed and stored, given by

$$\tilde{W}_{kl\delta\gamma} = \sum_p \sum_s \sum_{\alpha,\beta}^{N_{\text{sub}}} \sum_{\tau}^{N'_x} \sum_{\tau}^{N_Q} W_{\alpha\beta\tau}^{(p)(s)} \left( \frac{\alpha_s(Q_\tau^2)}{2\pi} \right)^p \tilde{C}_{kl,\alpha\beta\gamma\delta}^{(s),\tau}. \quad (3.32)$$

The quantity  $\tilde{W}_{kl\delta\gamma}$  is the FK table for the observable  $\sigma$  and encodes all of the theoretical treatment of the observable. The product in Eqn. ?? is therefore completely agnostic with regards to all theory parameters such as process, scales, perturbative order and strong coupling value. This makes the

FK table particularly simple to implement in a fitting procedure, and allows a clean separation of theory concerns from the calculation.

The FK convolution also benefits from requiring considerably fewer floating point operations than a typical `APPLgrid` convolution. This is particularly evident when studying multi-scale processes, where the sum over the scale grid is precomputed. The product over PDF flavours is now also limited to the  $n_f$ , typically seven, light partons rather than the general 13 parton basis. Of course in the FK procedure the ability to vary scales and the strong coupling with a single grid is lost, and new FK tables  $\widetilde{W}$  must be generated for different theoretical treatments.

The procedure outlined above for generating FK tables from `APPLgrid` or `FastNLO` files has been implemented in a C++ framework, alongside a comprehensive toolchain for performing FK table I/O and optimisation. The convolution in Eqn. ?? has been implemented for a general PDF input (for example Neural Network or LHAPDF) and extensively optimised. The optimisation ensures only the relevant parton sub channels and  $x$ -grid entries enter the product, which is performed as a memory-aligned scalar product with the use of SSE intrinsics [?]. Table ?? compares the relative speed improvement compared to the `APPLgrid` calculation of the basic FK convolution and the optimised version, using PDFs obtained through the LHAPDF library.

Observable	APPLgrid	FK	optimised FK
Total $W^+$ xsec	1.03 ms	0.41 ms (2.5X)	0.32 ms (3.2X)
Jet distribution	2.45 ms	20.1 $\mu$ s (120X)	6.57 $\mu$ s (370X)

Table 3.1: Typical timings per observable for several convolution methods. Two observables are presented, the total cross-section for  $W^+$  production and the inclusive jet  $p_\perp$  distribution. Values are given per datapoint. In brackets the relative speed-up compared to the native `APPLgrid` convolution is shown. For this test, the timings were calculated with a 2.9 GHz Intel Core i7 processor.

In Table ?? a good speed improvement is evident for the example single-scaled

process of  $W^+$  production, and a very significant improvement on the multi-scale jet production observable. It is important to note that these figures were obtained via convolutions with LHAPDF parton densities, and in a PDF fit a considerably greater speed advantage is gained via the FK procedure as no additional operation is required to evolve the PDFs.

While for most applications, the original `APPLgrid` convolution speed is more than sufficient, these speed improvements make the inclusion of a large LHC dataset possible, rather than prohibitively expensive in the NNPDF methodology. For example, in a typical NNPDF fit of 20,000 genetic algorithm generations, including a 100 datapoint jet dataset via the `APPLgrid` interface would add several days of additional computer time to each individual replica fit. With the FK procedure this additional cost is reduced to minutes.

The speed improvement is achieved without any loss of accuracy, as the interpolation procedure used to perform the PDF evolution is required in both the `APPLgrid` and FK convolutions. The two methods were benchmarked in Ref. [?], with the results shown in Table ???. The relative discrepancy  $\epsilon$  noted in the table is largely due to the additional interpolation in hard scale  $Q^2$  from LHAPDF required in the `APPLgrid` convolution that is not present in the FK method, as evolution is performed directly to the required scale.

### 3.3 Interpolating tools for automated NLO

Tools such as `FastNLO/ APPLgrid` and their extension for fast PDF fitting in the FK method, are invaluable in the analysis of collider data. Their usefulness is not limited to applications such as fitting, but can also be used to perform thorough QCD analysis with rigorous theory uncertainty estimation in situations where obtaining sufficient statistics with an NLO code or event generator would be extremely expensive computationally. Despite this, at the outset of LHC data

$ \eta_l $	W <sup>+</sup> distribution [pb]			W <sup>-</sup> distribution [pb]		
	FK	APPLgrid	$\epsilon_{\text{rel}}$	FK	APPLgrid	$\epsilon_{\text{rel}}$
0.00–0.21	617.287	617.345	0.01%	456.540	456.819	0.06%
0.21–0.42	616.988	617.062	0.01%	453.045	453.315	0.06%
0.42–0.63	620.237	620.290	0.01%	448.902	449.172	0.06%
0.63–0.84	624.192	624.235	0.01%	441.789	442.045	0.06%
0.84–1.05	630.235	630.286	0.01%	432.206	432.435	0.05%
1.05–1.37	636.835	636.886	0.01%	419.027	419.222	0.05%
1.37–1.52	642.800	642.861	0.01%	403.908	404.084	0.04%
1.52–1.74	642.499	642.569	0.01%	390.564	390.724	0.04%
1.74–1.95	642.351	642.437	0.01%	377.328	377.473	0.04%
1.95–2.18	628.592	628.693	0.02%	359.373	359.498	0.03%
2.18–2.50	590.961	591.079	0.02%	337.255	337.366	0.03%

$ y $	Z distribution [pb]		
	FK	APPLgrid	$\epsilon_{\text{rel}}$
0.0–0.4	124.634	124.633	0.001%
0.4–0.8	123.478	123.488	0.01%
0.8–1.2	121.079	121.108	0.02%
1.2–1.6	118.057	118.108	0.04%
1.6–2.0	113.512	113.549	0.03%
2.0–2.4	106.552	106.562	0.01%
2.4–2.8	93.7637	937.838	0.02%
2.8–3.6	55.8421	558.538	0.02%

$p_T$ (GeV)	ATLAS 2010 jets [pb]		
	FK	APPLgrid	$\epsilon_{\text{rel}}$
20–30	$6.1078 \times 10^6$	$6.1090 \times 10^6$	0.02%
30–45	986285	98654	0.03%
45–60	190487	190556	0.04%
60–80	48008.7	48029.7	0.04%
80–110	10706.6	10710.4	0.03%
110–160	1822.62	1822.87	0.01%
160–210	303.34	303.443	0.03%
210–260	76.1127	76.1338	0.03%

Table 3.2: Benchmark of the FK result for datasets with different underlying processes, all generated according to ATLAS experimental kinematics and acceptances. The APPLgrid and FK results are presented along with the relative discrepancy between the two. Table from [?].

taking the amount of codes interfaced to such interpolating tools was extremely limited. Additionally the need for separate interfaces to existing codes meant a great deal of duplication in terms of analysis tools and software. The APPLgrid group provided a direct interface to the NLO codes MCFM [?] and `nlojet++` [?, ?]. FastNLO provided a set of precomputed scenarios generated through a private

interface to NLO codes. More recently, a public toolkit was released to allow for the interfacing of `FastNLO` to external calculations.

A conspicuous absence was an interface to tools providing automated NLO calculations via computer algebra suitable one-loop methods [?, ?, ?, ?, ?, ?, ?] and their implementations in parton level Monte Carlo codes such as `MadGraph` [?], `HELAC` [?] and `SHERPA` [?, ?]. In this section we shall discuss the implementation of a fast interface to such codes, the `MCgrid` [?] package; developed with the aid of funding from the MCnet initial training network.

### 3.3.1 Reweighting Monte Carlo calculations

Recalling Eqn. ??, a hadronic observable calculation proceeds via

$$\sigma_{pp \rightarrow X} = \sum_p \left( \frac{\alpha_s(Q^2)}{2\pi} \right)^p \int dx_1 dx_2 F_l(x_1, x_2, Q^2) d\hat{\sigma}_{l \rightarrow X}^{(p)}, \quad (3.33)$$

where the initial state PDFs have been grouped according to Eqn. ?? and the sum over subprocesses is implicit. In an event generator this integral is performed via Monte Carlo integration. At leading order this is a relatively straightforward procedure,

$$\sigma_{pp \rightarrow X}^{\text{LO}} = \sum_{e=1}^{N_{\text{evt}}} \tilde{w}_e(k_e) = \sum_{e=1}^{N_{\text{evt}}} \left( \frac{\alpha_s(k_e)}{2\pi} \right)^{p_{\text{LO}}} w_e(k_e) F_{l_e}(k_e), \quad (3.34)$$

where  $\tilde{w}$  is the full event weight and the  $w$  are the matrix element weights generated via importance sampling of the integrand of Eqn. ???. The  $F_{l_e}$  refer to the parton density of the event's subprocess  $l_e$ . Each event is generated according to a set of kinematics

$$k_e = \left\{ p_1, \dots, p_n, x_1, x_2, \frac{\mu_F^2}{Q^2}, \frac{\mu_R^2}{Q^2} \right\}. \quad (3.35)$$

As a full re-run of the event generator for every parameter variation is

extremely expensive, the variation is typically performed via an event-by-event reweighting procedure. The full set of events is stored in a common format such as `HepMC` [?] and re-processed by dividing out the appropriate factors of the old PDFs and  $\alpha_s$  and multiplying in the desired new values.

$$\tilde{w}_e(k_e) \rightarrow \left( \frac{\alpha'_s(k_e)}{\alpha_s(k_e)} \right)^{p_{\text{LO}}} \frac{F'_{l_e}(k_e)}{F_{l_e}(k_e)} \tilde{w}_e(k_e), \quad (3.36)$$

where the primed quantities denote the new, reweighted strong coupling and PDF choices. Having the full generated event sample stored also has the advantage of being able to rerun analysis software with varying parameters/selections without the need to rerun the potentially expensive event generation.

The reweighting situation in an NLO calculation is considerably more complicated. In order to be able to solve the integral numerically, a divergence-subtraction scheme e.g Catani-Seymour [?] or Frixione-Kunst-Signer (FKS) [?,?] must be employed. These subtraction algorithms separate the calculation into distinct sections which are to be numerically evaluated individually. Here we shall discuss the implementation in terms of a Catani-Seymour dipole scheme. The four contributions to the total NLO cross section are

$$\sigma_{pp \rightarrow X}^{\text{NLO}} = \int d\hat{\sigma}^B + \int d\hat{\sigma}^V + \int d\hat{\sigma}^I + \int d\hat{\sigma}^{RS}. \quad (3.37)$$

The terms  $B$ ,  $V$ ,  $I$  and  $RS$  refer to the Born (B), Virtual (V), Integrated subtraction (I) and Real Subtracted (RS) cross section elements respectively. Neglecting terms used in the variation of perturbative scales, the  $B$ ,  $V$  and  $RS$  terms may be integrated via a Monte Carlo procedure equivalently to Eqn. ???. The integrated dipole term however has a rather more complicated dependance upon the initial state PDFs, originating as it does through the splitting of an initial state parton. The Monte Carlo solution to the  $I$  integral is given by

$$\int d\hat{\sigma}^I = \sum_{e=1}^{N_{\text{evt}}} \left( \frac{\alpha_s(\mu_R^2)}{2\pi} \right)^{p_{\text{NLO}}} \left\{ F_{l_e}(x_1, x_2, \mu_F^2) w_e^{(0)} + \sum_s^{N_{\text{sub}}} F_s(x_1/x'_1, x_2, \mu_F^2) \tilde{w}_{e,s}^{(1)} + \sum_s^{N_{\text{sub}}} F_s(x_1, x_2/x'_2, \mu_F^2) \tilde{w}_{e,s}^{(2)} \right\}, \quad (3.38)$$

in which the weight  $w^{(0)}$  arises through the usual Born-like PDF dependance, and the weights  $w^{(1/2)}$  arise from integration over parton- $x$  from the first or second parton in the initial state splitting. To reweight such events these weights must therefore be properly distinguished in the event record. While this is not the case in the standard `HepMC` layout, a format based upon `ROOT NTuples` was designed by the `BlackHat-Sherpa` group for the reweighting of NLO event weights [?]. In the `BlackHat NTuple` format the weights that must be distinguished for the accurate treatment of scale variations are also stored.

While the event reweighting approach is considerably faster than an entire rerun of the Monte Carlo, a reweight of a full event sample can still take a considerable amount of computer time. The key issue being that the statistical accuracy of the calculation is limited by the number of events in the sample, and therefore for a more accurate calculation, more computational expense is incurred. This dependence on the event loop is not removed by the event reweighting procedure. The dependance can however be removed by applying the interpolation methods of the previous section, by providing an interface for event generators to interpolating packages such as `APPLgrid`.

### 3.3.2 An interpolation interface for automated NLO

The **MCgrid** project began as a direct interface for the **SHERPA** event generator framework to the **APPLgrid** interpolation package. The development of a new interface between an event generator and the **APPLgrid** framework in principle requires the implementation of an analysis suite to provide the categorisation of event final states into appropriate observable bins. However the **MCgrid** interface is built upon standard analysis tools and formats to provide a more general interface between standards-compliant Catani-Seymour event generators to the **APPLgrid** package.

**MCgrid** is written as a set of additional tools for the **Rivet** MC analysis system. The **Rivet** system implements a wide range of experimental analysis tools and provides the flexibility for the user to define their own selection criteria and processing tools to operate on an event final state. Writing the **APPLgrid** interface as a **Rivet** extension therefore removes the need to implement a separate toolchain, and allows a degree of generator agnosticism. As **Rivet** operated upon events in the standard **HepMC** format, any generator equipped to output events in this format may potentially be interfaced to **APPLgrid** through **MCgrid**.

The interface requires additional information over the standard data available in **HepMC**, as the information on the weight breakdown as per Eqn. ?? must be available. However this can be straightforwardly appended in the **HepMC** user defined weights fields. The interface then provides the correct handling of initial state parton mappings from the PDF basis used in the Catani-Seymour process to the **APPLgrid** flavour basis.

With the appropriate mapping to initial state parton flavours performed, the weights must be converted to the appropriate subprocess basis. The minimal initial state PDF basis can be automatically determined by a set of packaged scripts. General purpose Monte Carlo codes such as **SHERPA** will typically generate events with the full initial parton flavours explicit rather than generating events

based upon QCD subprocesses. Weights originating from different flavour basis channels are generated via importance sampling of the distribution in order to ensure an efficient description of the most important channels. Accordingly there is a selection weight present in each event weight, given by

$$w_e(i, j, k_e) = \mathcal{N}_{ij} d\hat{\sigma}_{l_e \rightarrow X}(k_e) \Pi_{\text{ps}}(k_e) \Theta(k_e - k_{\text{cuts}}), \quad (3.39)$$

where the factor  $\mathcal{N}_{ij}$  is approximately given by

$$\mathcal{N}_{ij} \sim \frac{N_{\text{tot}}}{N_{ij}}, \quad (3.40)$$

where  $N_{\text{tot}}$  denotes the total number of events in the sample, and  $N_{ij}$  is the number of events initiated by partons of flavour  $i$  and  $j$ . In Eqn. ?? we use  $\Pi$  to represent the phase space weight associated with the kinematics  $k_e$ , and the  $\Theta$  as a step function implementing the desired kinematic cuts in the analysis. The selection weights  $\mathcal{N}$  must be converted into the appropriate subprocess selection weight to prevent the statistical uncertainty in poorly-sampled distributions from overwhelming the subprocess. **MCgrid** monitors the relative population of the channels and subprocesses in order to provide a statistically sound subprocess combination. The selection weight in Eqn. ?? must be converted into the appropriate subprocess selection weight as

$$\mathcal{N}_{ij} \rightarrow \mathcal{N}_l = \frac{N_{\text{tot}}}{N_l}, \quad (3.41)$$

where  $N_l$  is the number of events falling into the initial state subprocess  $l$ . Converting the selection weights to the appropriate subprocess selection weight is therefore a matter of multiplying each event weight by a factor  $N_{ij}/N_l$ .

In this way the fully exclusive predictions given in a typical Monte Carlo event generator may be effectively converted into the relevant subprocess basis. With

this accomplished, and the correct weight conversion performed according to the exact PDF dependance of the Catani-Seymour counterterms, the weights may be filled directly into an **APPLgrid** type weight grid. Figure ?? demonstrates the application of the **MCgrid** package when used in conjunction with **Sherpa** and **BlackHat** as a one-loop generator. The tools are applied to the test cases of Drell-Yan and inclusive jet production, with the resulting **APPLgrid** applied to the estimation of scale and  $\alpha_S$  uncertainties alongside standard PDF error, requiring a very large number of replicas.

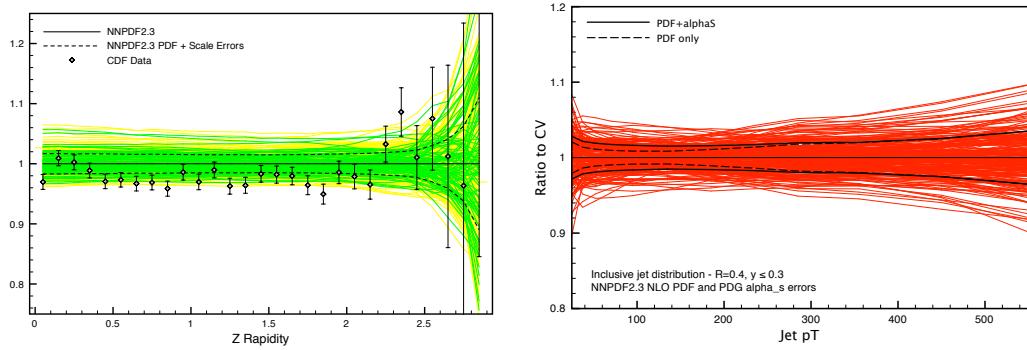


Figure 3.4: An example of the output of the **MCgrid** package. A  $Z$  boson rapidity distribution plot is shown on the left, with scale error estimation. The plot on the right demonstrates the grids applied to inclusive jet data, with  $\alpha_S$  error estimation. Both plots are normalised to their central values, to demonstrate the level of uncertainty.

The **MCgrid** project is publicly available<sup>1</sup>, and allows for the first time calculations from automated NLO event generators to be interfaced to interpolation tools, for potential application in PDF fits, or indeed fast parameter variation studies in phenomenological applications.

---

<sup>1</sup>The software is available at <http://mcgrid.hepforge.org>.

# Chapter 4

## LHC Data for Parton Determinations

The Large Hadron Collider has the ability to provide a comprehensive examination of QCD and electroweak physics at a wide range of scales. The requirement of precise and reliable determinations of proton structure is clear in order to fully exploit the LHC’s potential. LHC data also has the potential to provide deep new insights into parton distributions, examining hitherto poorly determined flavours and kinematic regimes. A great deal of effort has therefore been expended in providing and validating tools for the inclusion of LHC data in an efficient manner into NNPDF fits.

In this section the Standard Model measurements of relevance to PDF determination so far performed by the LHC shall be briefly summarised. While the general processes have been described previously, here we shall look directly at the experimental data along with a brief examination of the areas of agreement or discrepancy with regard to PDF sets made available before the first data runs of the LHC.

## 4.1 Jet measurements

At the LHC, data on the production of collimated jets of particles originating from partonic final states provides valuable information on proton structure and additional constraints for  $\alpha_S$  determinations. The LHC’s centre-of-mass energies mean that jets with transverse momenta in the TeV range are observable for the first time. Forward jets probing the very large- $x$  gluon that has suffered from poor constraints prior to the LHC. As the prototypical QCD measurement, data is available from both of the general purpose LHC experiments, and preliminary data on jets in the forward region is available from LHCb [?]. LHC measurements are based upon modern infrared and collinear safe jet-finding algorithms such as anti- $k_T$  [?]. In PDF fits the jet quantity of interest is typically the inclusive measurement rather than dijet data. In principle dijet measurements offer more discriminating power over the parton distributions, however they typically suffer from larger scale uncertainties and often must be corrected for higher order effects, typically modelled through parton showers.

Here we shall summarise the relevant jet measurements at the LHC with a focus on the data most relevant to PDF determination.

The first ATLAS inclusive jet and dijet measurements were based upon a partial analysis of  $17 \text{ nb}^{-1}$  of data available from the 2010 data run at a centre of mass energy of 7 TeV [?]. This result was then updated to the full 2010 dataset of  $37 \text{ pb}^{-1}$  [?]. The full 2010 measurement presents the inclusive jet cross section differentially in both the jet  $p_T$  and rapidity. Data is available for the  $20 \leq p_T < 1500 \text{ GeV}$  range for jets with rapidity  $|y| < 4.4$ , and is available for two choices of the anti- $k_T$  cone size,  $R = 0.4$  and  $R = 0.6$ .

Figure ?? from the ATLAS  $37 \text{ pb}^{-1}$  result demonstrates the level of agreement of the fixed-order NLO inclusive jet computation present in `NLOJet++` with the experimental data given four choices of PDFs: CT10, MSTW2008, NNPDF2.1 and HERAPDF1.5. Predictions from the four sets largely agree within their PDF

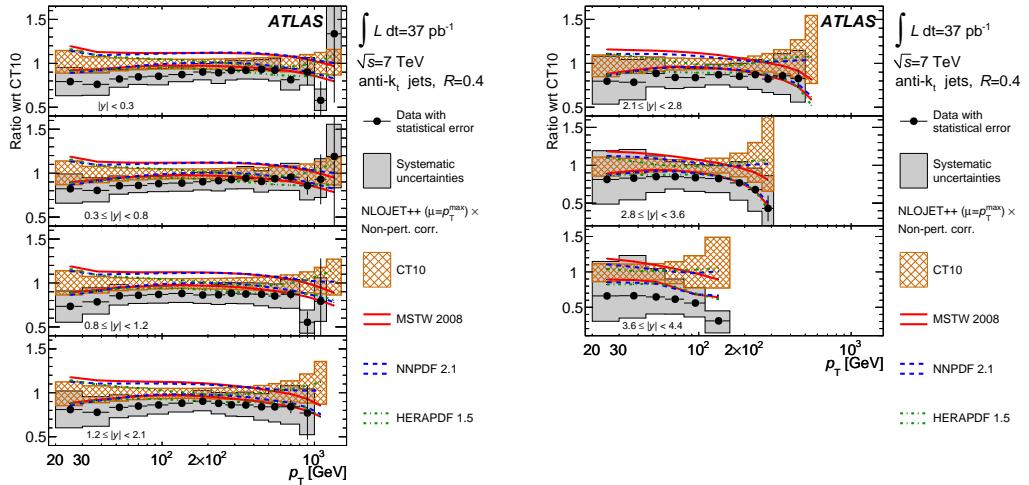


Figure 4.1: ATLAS inclusive jet data with anti- $k_T$  algorithm  $R = 0.4$  from the 2010 dataset. Figures from [?]. Predictions are shown based upon MSTW2008, NNPDF2.1 and HERAPDF1.5 PDFs, with all data and theory normalised to the CT10 central value.

uncertainties, and the experimental data also shows good agreement for most of the data range. Some evidence of a systematic discrepancy is visible at large  $p_T$ , an effect that becomes more noticeable in the larger rapidity bins (and therefore more extreme values of parton- $x$ ).

ATLAS has also published data on the inclusive jet cross-sections at  $\sqrt{s} = 2.76$  GeV measured during the 2011 run [?]. The data provides an important link between jet measurements at lower centre-of-mass energies at the Tevatron and the higher scale measurements previously published. In addition, the ratio of the  $\sqrt{s} = 2.76$  GeV data to the 2010  $\sqrt{s} = 7$  GeV measurement is presented. The ratio offers additional important constraints in that the dominant uncertainties upon the jet measurements are systematic across both datasets, and therefore largely cancel in the ratio. Figure ?? demonstrates the reduced uncertainty in the measurement, and therefore the additional constraint that the data may provide parton fits.

CMS has published three measurements of inclusive and dijet observables to

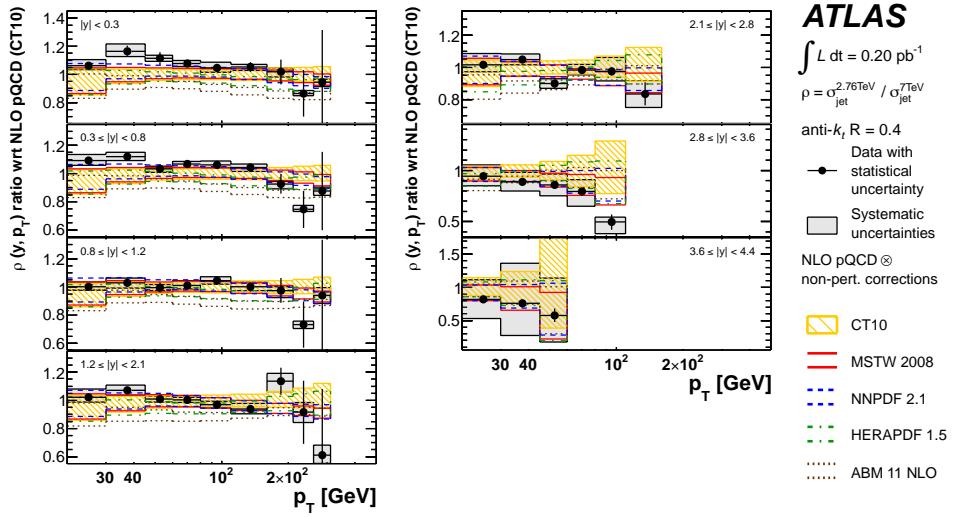


Figure 4.2: ATLAS inclusive jet ratio between  $\sqrt{s} = 2.76$  GeV and 7 GeV data, anti- $k_T$   $R = 0.4$ . Figures from [?].

date. The first provided data in the  $18 < p_T < 1100$  GeV interval for jets with  $|y| < 3$  based upon  $34\text{pb}^{-1}$  of 2010 data [?]. This was followed up by a study of jets in the forward region [?], examining inclusive jets with pseudorapidities  $3.2 < |\eta| < 4.7$ , and dijets with one forward jet and one central  $|\eta| < 2.8$  jet. A study of 2011 data totalling  $5.0\text{fb}^{-1}$  was also performed of jets in the central  $|y| < 2.5$  region up to very high jet transverse momenta  $p_T < 2$  TeV [?]. CMS also utilises the anti- $k_T$  clustering algorithm, with cone sizes  $R = 0.5$  and  $R = 0.7$ .

Figure ?? shows the inclusive data from the CMS central region jet measurement normalised to the NNPDF2.1 central value. Results are once again largely consistent with PDFs determined with pre-LHC data.

## 4.2 $W/Z$ boson production

The measurement of electroweak vector boson production and Drell-Yan cross sections are standard candle measurements for the LHC, and have been widely studied by ATLAS, CMS and LHCb in the first run.

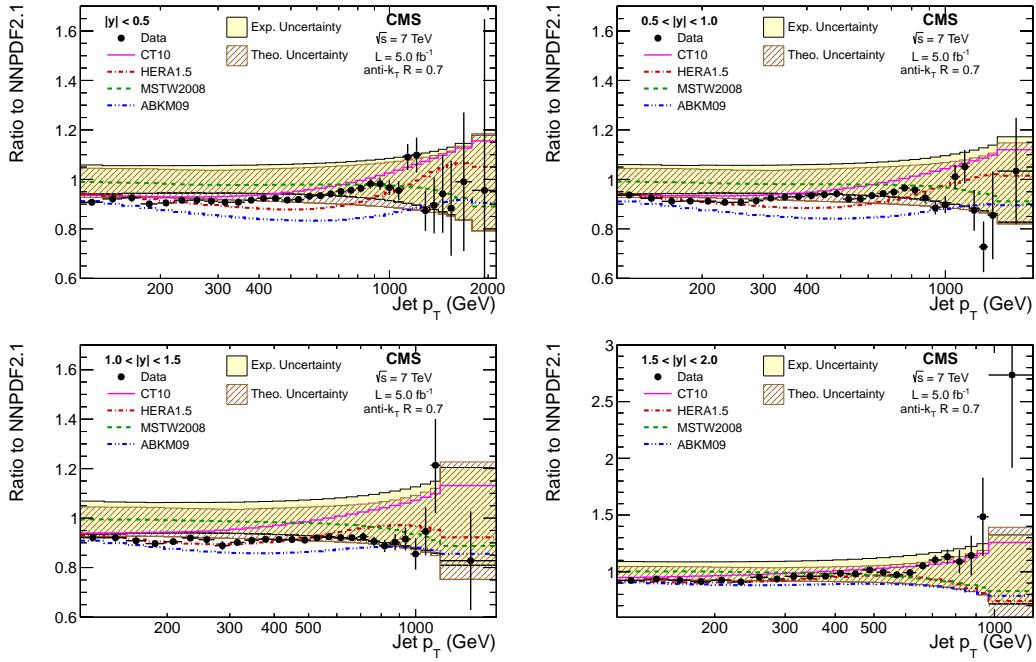


Figure 4.3: CMS inclusive jet data with anti- $k_T$  algorithm  $R = 0.7$  from the 2011 dataset. Figures from [?]. Predictions are shown based upon MSTW2008, NPDF2.1 and HERAPDF1.5 PDFs, with all data and theory normalised to the NPDF2.1 central value.

CMS has presented measurements of the  $Z$  boson  $p_T$  and rapidity distributions, initially upon  $36\text{pb}^{-1}$  of 7 TeV 2010 data [?], and more recently a preliminary study of 8 TeV data on  $Z$  decay to dimuons [?, ?]. The first differential measurements of  $W$  boson production at CMS were lepton charge asymmetry measurements based upon 2010 data [?], which were superseded by the muon asymmetry measurement based upon  $840\text{pb}^{-1}$ , and then  $4.6\text{pb}^{-1}$  of 2011 data [?, ?]. In Figure ?? the 2010 data  $W$  asymmetry measurement of CMS is shown, demonstrating the constraining power of the earlier CMS result, where agreement is generally good with the pre-LHC parton distributions with the exception of the MSTW 2008 description.

ATLAS initially published a study of the  $W$  muon asymmetry distribution with  $31 \text{ pb}^{-1}$  of 7 TeV data [?]. This was followed by studies of the  $Z$  [?] and  $W$  [?]  $p_T$  distributions. The most recent data is provided by a combined study

of the  $W$  and  $Z$   $p_T$  distributions based upon the full 2010 dataset [?].

The LHCb detector has a window upon electroweak vector boson production in the very forward region, a kinematic regime that cannot be explored by the general-purpose detectors.  $W$  and  $Z$  to muon production data based upon an integrated luminosity sample  $37\text{pb}^{-1}$  was published in Ref. [?], where data was taken in the pseudorapidity range  $2.0 < |\eta| < 4.5$  and presented differentially in the (pseudo)rapidity of the detected lepton (pair). Figure ?? shows the main result of the LHCb  $W/Z$  study and demonstrates the good agreement of the theoretical predictions, within the limited statistical precision available in the forward data sample.

### 4.3 Prompt photon data

Constraints upon the gluon distribution are possible through measurements made of direct photon production at the LHC. Both CMS and ATLAS have published prompt photon data. ATLAS provides inclusive data in photon pseudorapidity intervals of  $|\eta| < 1.37$  and  $1.52 \leq |\eta| < 2.37$ , for transverse energies  $45 \leq E_T < 400$  GeV [?], the data showing excellent agreement with predictions from CTEQ6.6 and JETPHOX. Additionally data is available for isolated prompt photon data in association with a jet, based upon the same dataset [?], where once again NLO predictions provide a good description of the data, albeit with a small discrepancy arising for photons with  $E_T < 45$  GeV.

CMS has performed an isolated photon measurement based upon the same 2010 data run, in the pseudorapidity range  $|\eta| < 2.3$  for photons with  $25 < E_T < 400$  GeV [?]. The CMS result is plotted in Figure ??, which shows the agreement between the NLO calculation and the experimental data. The figure demonstrates clearly the precision available of the experimental measurement, however the theoretical predictions clearly suffer from relatively large scale uncertainties. The

inclusion of such data into PDF determinations is therefore likely to be challenging without further theoretical progress.

## 4.4 Top pair production data

LHC collaborations have made extensive measurements of the top pair production cross-section, building upon the combined Tevatron analysis of [?]. Unlike at the Tevatron where the  $qq$  initiated channel is favoured,  $t\bar{t}$  data at the LHC is primarily a probe of the gluon content of the proton through the  $gg \rightarrow t\bar{t}$  subprocess. The ATLAS collaboration has published measurements of the  $t\bar{t}$  cross section in a number of channels, with combination results available at both 7 TeV [?] and 8 TeV [?] centre of mass energies. Likewise CMS have published combined  $t\bar{t}$  analyses at 7 [?] and 8 [?] TeV. These results are compared to the theoretical prediction obtained from NNPDF2.3 at NNLO+NNLL with `top++ v2.0` [?] in Figure ??.

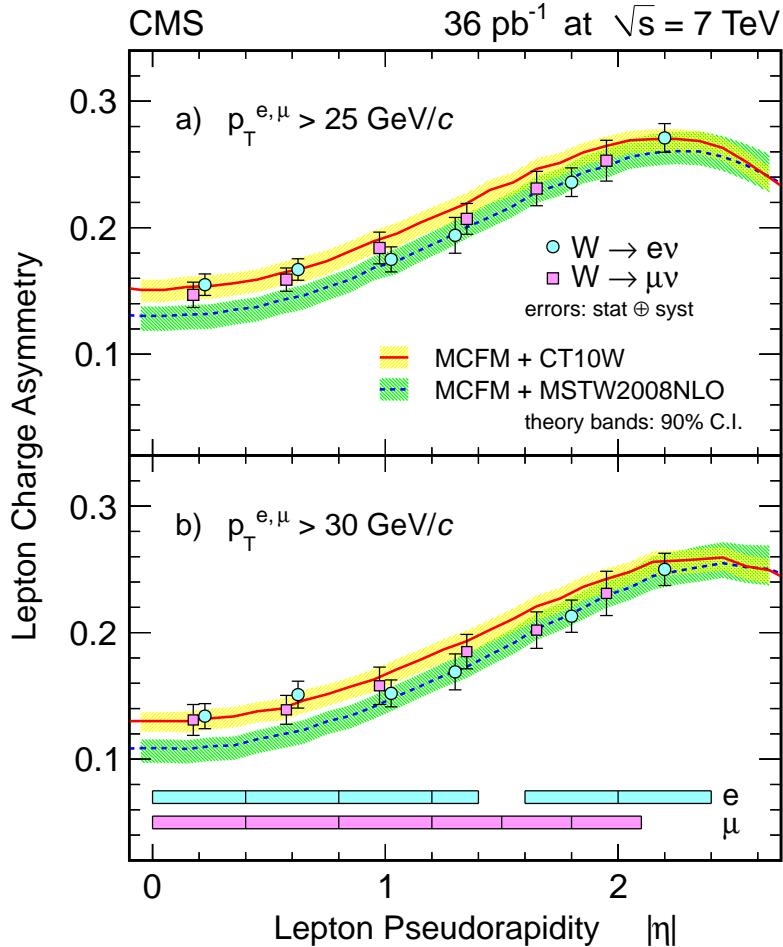


Figure 4.4: CMS 2010  $W$  asymmetry data from the electron and muon decay channels. Figure from [?]. The figure demonstrates the good agreement of the CT10 PDF set with the experimental data, and the somewhat poorer agreement of the MSTW2008 set, a typical feature of the MSTW fit with LHC electroweak data.

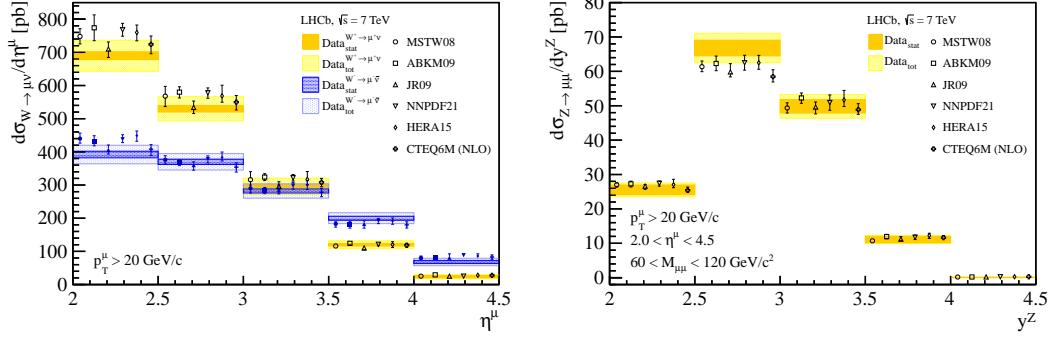


Figure 4.5: LHCb  $W/Z$  boson (pseudo)rapidity data. The left panel shows the  $W^\pm$  distributions in the pseudorapidity of the detected lepton compared to leading PDF sets. The right panel shows the  $Z$  in the rapidity of the resulting lepton pair. Figures from [?].

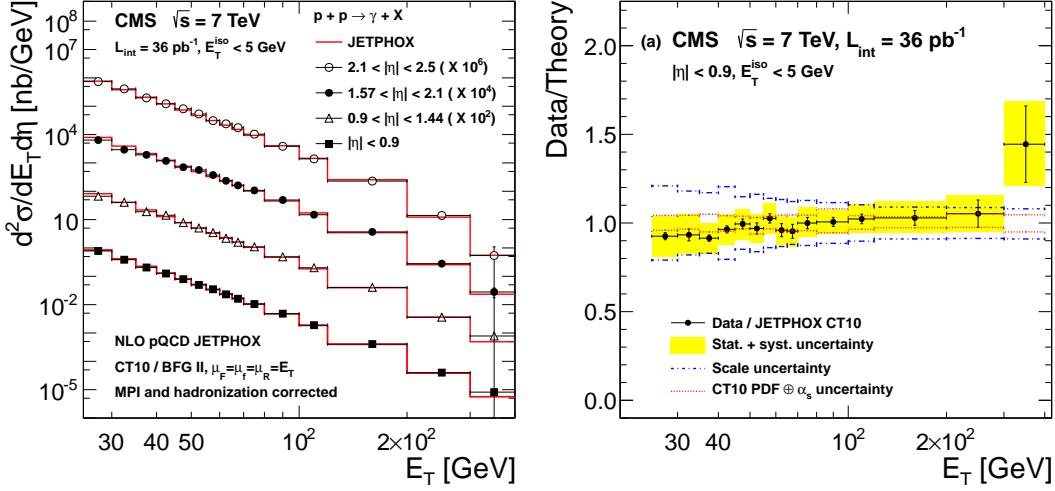


Figure 4.6: Figures from the CMS isolated prompt photon measurement [?]. The left panel shows the full dataset compared to the CT10/JETPHOX predictions. The right panel shows the result for the lowest pseudorapidity bin as a ratio to the CT10 central value.

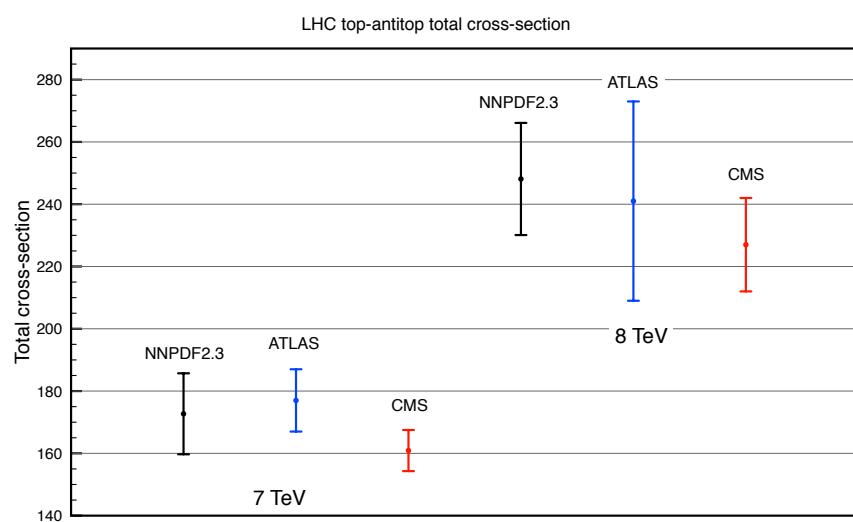


Figure 4.7: LHC 7 and 8 TeV  $t\bar{t}$  total cross-section data and predictions from NNPDF2.3 at NNLO+NNLL precision. Values in the figure are taken from [?].

# Chapter 5

## The impact of LHC data on PDFs

The series of measurements made to date during the first runs of the LHC have been studied to assess the impact upon PDFs and their uncertainties, and where appropriate, have been included into PDF fits through the NNPDF methodology. Early LHC measurements serve not only as useful constraints in their own right, but also as a testing ground for tools developed to include such data into PDF determinations, ready for future datasets with even higher precision.

In this section we shall provide an overview of the work performed in the inclusion of LHC data and some of the results obtained. The methods introduced in Chapter ?? are applied to some of the datasets in Chapter ?? and the resulting PDFs discussed and compared to results obtained before the LHC era.

### 5.1 The NNPDF2.2 parton set

The NNPDF2.2 parton set [?] was the first practical demonstration of the Bayesian reweighing and unweighting methods. These methods were applied to the inclusion of a series of  $W$  boson charge asymmetry measurements made

by the ATLAS, D0 and CMS experiments. In this way LHC data was included for the first time in a public parton set.

### 5.1.1 NNPDF2.2 dataset

The dataset studied included the 1.96 TeV  $p\bar{p}$  data from D0 on both the  $W$  electron [?] and muon asymmetries [?]. The LHC dataset consisted of the 2010 run  $W$  lepton asymmetry measurement of CMS [?] and ATLAS [?]. LHCb asymmetry data with a full covariance matrix was not available at the time and so was not included in the dataset. Agreement for the LHC data points is generally reasonable for PDF sets obtained without LHC data, as shown in Figure ???. While NNPDF2.1 and CT10 obtain good overall agreement, the MSTW2008 prediction tends to be systematically lower than the data.

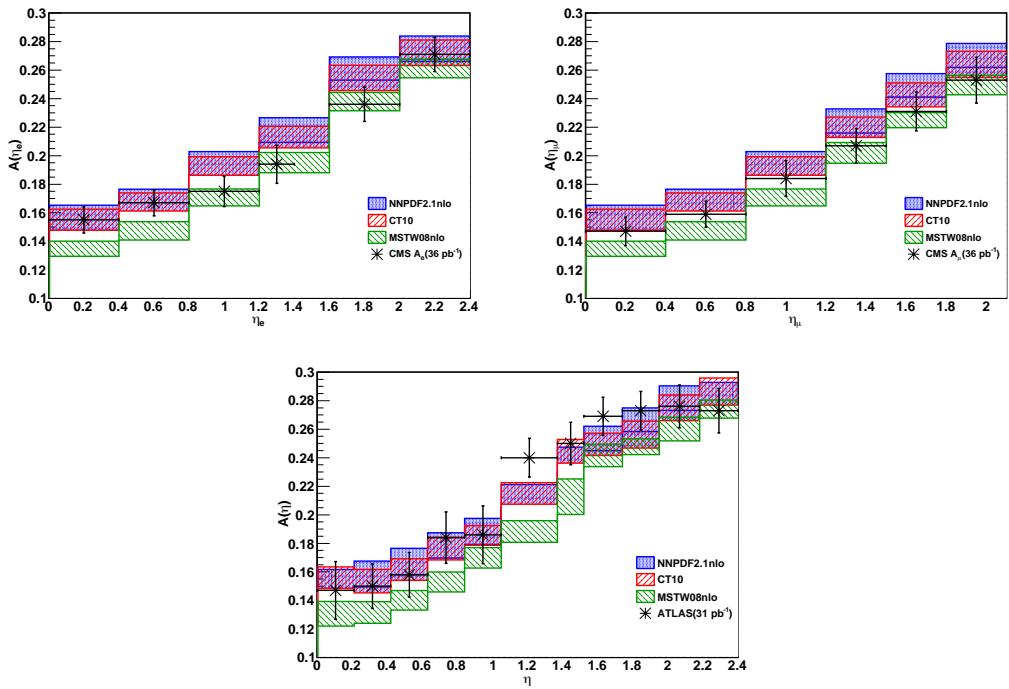


Figure 5.1: Plot of LHC data to be included in the NNPDF2.2 determination. Data from the CMS electron (top-left) and muon (top-right) data is given alongside the ATLAS muon asymmetry data (bottom). Figure from [?].

The level of agreement taking into account systematic uncertainties is of course most clearly quantified with a  $\chi^2$  measure. In Table ?? we compare the fit quality of NNPDF2.1, CT10 and MSTW2008 NLO sets to the new data. The result show that while generally the consistency with the new datasets is good in NNPDF2.1, there is certainly room for improvement.

	$N_{\text{dat}}$	NNPDF2.1	CT10	MSTW08
ATLAS(31pb $^{-1}$ )	11	0.77	0.77	3.32
CMS(36pb $^{-1}$ ) electron $p_T > 25$ GeV	6	1.83	1.19	1.70
CMS(36pb $^{-1}$ ) muon $p_T > 25$ GeV	6	1.24	0.73	0.77
D0(0.3fb $^{-1}$ ) muon $p_T > 20$ GeV	10	1.48	-	-
D0(0.75fb $^{-1}$ ) electron $E_T > 25$ GeV	12	4.39	-	-

Table 5.1: Table of  $\chi^2$  values for new data included in NNPDF2.2.

The combined goodness of fit value for the LHC and Tevatron datasets for NNPDF2.1 is  $\chi^2/N_{\text{data}} = 2.22$  which suggests a less than ideal description of the data, largely due to the precise D0 electron asymmetry measurement.

For this dataset the reweighting technique presented an ideal method for the data inclusion. With a total of 45 points the dataset is relatively small, and together with the fair agreement of the prior PDF set (NNPDF 2.1) the reweighting can be accomplished with a reasonable number of prior replicas. Also the lack of a fast method of determining these asymmetries within the NNPDF framework at the time meant that the data could not be included via a conventional fit, necessitating the reweighting approach.

### 5.1.2 NNPDF2.2 results

The LHC and Tevatron  $W$  boson asymmetry datasets were included into the NNPDF 2.1 determination by a reweighting both individually and upon the combined  $\chi^2$  figure for the whole dataset. To ensure maximal final ensemble efficiency, an NNPDF 2.1 prior with  $N_{\text{rep}} = 1000$  replicas was used for the

reweighting. Theoretical predictions for the various datasets were computed at NLO using DYNNLO [?]. After computing the  $\chi^2$  values using the  $t_0$  method to ensure consistency with the fit procedure, the number of effective replicas remaining in the ensemble is given by the Shannon entropy (Eqn. ??). For the different reweighting combinations attempted, the number of effective replicas is given in Table ??.

	ATLAS	CMS	LHC	LHC + TeV
$N_{\text{eff}}$	928	531	619	181

Table 5.2: Number of effective replicas for each dataset reweighting in NNPDF 2.2. Figures are given for the ATLAS and CMS experiments, along with their combination (LHC) and their combination with the Tevatron data (LHC+TeV).

Both ATLAS and CMS show good consistency with the prior in the reweighting, with the CMS data providing the greater constraint and resulting in a lower number of replicas surviving the reweighting process. The reweighting with ATLAS data only leading to 928 effective replicas and the CMS reweighting resulting in 531. The D0 data goes further to provide a great deal of extra constraint. In the final combined reweighting, roughly one fifth of the prior replicas remain active, a figure which demonstrates that the  $W$  asymmetry data available at the time was able to provide a great deal of additional information on parton distributions.

The PDF set resulting from the reweighting with the combined dataset was then unweighted to 100 replicas via the mechanism described in Chapter ?. The unweighted set forms the NNPDF 2.2 determination, available as part of the LHAPDF platform. In Table ?? the full  $\chi^2$  breakdown for every experiment in the NNPDF2.2 dataset is shown. It is clear from the table that a great deal of improvement in the new  $W$  asymmetry data is achieved by the addition of the new data, and there is no associated cost to the  $\chi^2$  values for the rest of the dataset, suggesting that the new data maintains a good consistency with the

measurements already utilised in NNPDF 2.1. The global fit quality therefore has a modest improvement from  $\chi^2/N_{\text{data}} = 1.165$  to 1.157.

Experiment	$N_{\text{dat}}$	NNPDF2.1	NNPDF2.2
NMC-pd	132	0.97	0.97
NMC	221	1.73	1.72
SLAC	74	1.33	1.28
BCDMS	581	1.24	1.23
HERAI-AV	592	1.07	1.07
CHORUS	862	1.15	1.15
FLH108	8	1.37	1.37
NTVDMN	79	0.79	0.70
ZEUS-H2	127	1.29	1.28
ZEUSF2C	50	0.78	0.78
H1F2C	38	1.51	1.51
DYE605	119	0.84	0.86
DYE886	199	1.25	1.27
CDFWASY	13	1.85	1.81
CDFZRAP	29	1.66	1.70
D0ZRAP	28	0.60	0.58
CDFR2KT	76	0.98	0.96
D0R2CON	110	0.84	0.83
ATLASmuASY	11	[0.77]	1.07
CMSeASY	6	[1.83]	1.08
CMSmuASY	6	[1.24]	0.56
D0eASY	12	[4.39]	1.38
D0muASY	10	[1.48]	0.35
Total		1.165	1.157

Table 5.3: The global  $\chi^2/N_{\text{dat}}$  values to all experiments included in the NNPDF 2.2 fit. Values presented within square brackets were not included in the associated fit, and do not contribute to the total at the end of the table. Values from [?].

Examining the NNPDF 2.2 PDFs directly, the largest differences with respect to the prior arise as expected in the light quark PDFs, the most relevant initial states for the  $W$  asymmetry. Figure ?? demonstrates the effect that the new data has upon the PDFs. For all of the light quark distributions a substantial reduction

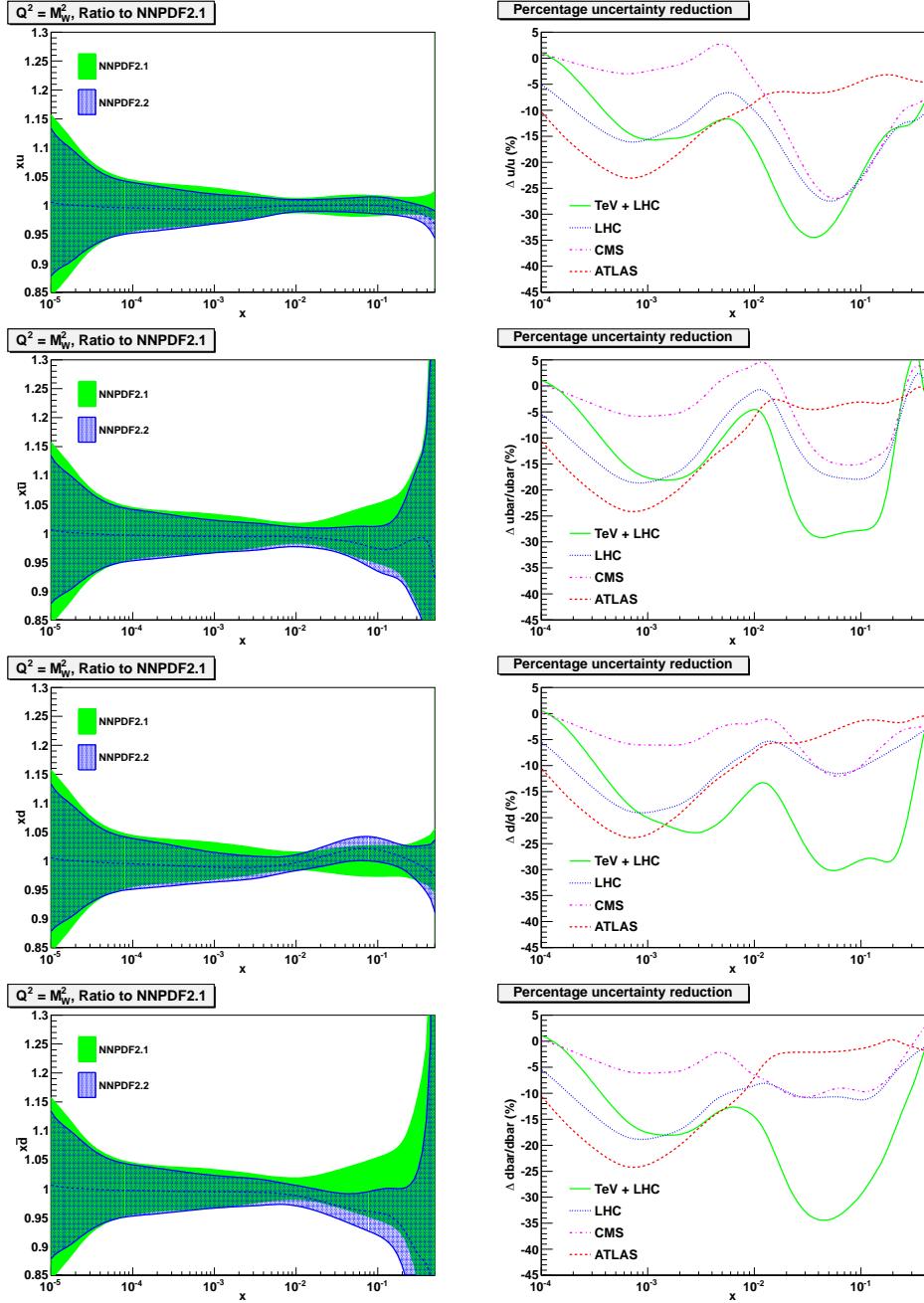


Figure 5.2: Impact of the LHC and Tevatron  $W$  asymmetry data upon PDFs. On the left, the NNPDF 2.1 (prior) and NNPDF 2.2 (reweighted) distributions are shown for the light quarks  $u, \bar{u}, d, \bar{d}$ . On the right are the relative uncertainty changes in the equivalent PDFs under a reweighting with the various dataset options, with the green lines indicating the final NNPDF2.2 result. The plots on the right therefore demonstrate the impact of the new data upon light quark PDF uncertainties. Figures from [?].

of the uncertainties can be observed, with a typical reduction of around 25%. The PDF central values also undergo a slight shift in the large- $x$  region, typically demonstrating a preference for softer light quarks. Phenomenologically these improvements will manifest in reduced uncertainties for observables sensitive to light/valence quarks over a large kinematic range, and a slightly tweaked distribution for those observables probing high- $x$  physics, such as the high rapidity observable region.

The NNPDF2.2 parton set was used in the CMS  $840\text{pb}^{-1}$   $W$  electron asymmetry measurement [?], where excellent agreement was demonstrated alongside the high precision available for electroweak observables with the 2.2 set. Figure ?? taken from the CMS paper illustrates the level of agreement in comparison to the CT10, MSTW 2008 and HERAPDF 1.5 predictions.

## 5.2 The NNPDF2.3 parton set

The NNPDF2.2 fit demonstrated the constraining power of early LHC measurements, and provided a showcase for the reweighting technique as a method of analysing the impact of new data and indeed producing a new PDF set including the data’s constraints. Nevertheless, the rapid pace of new experimental measurements meant that the data included in the set was soon superseded with higher integrated luminosity samples, and datasets sampling other processes of interest were being explored at the LHC. As the reweighting exercise in NNPDF2.2 had demonstrated, the inclusion of much more data into the fit would require priors with a rather unwieldy number of replicas, needing in excess of a thousand to include even a modest additional dataset. Therefore to include a large set of up to date measurements from the LHC into a parton fit, the conventional fitting methodology must still be applied.

The development of the FK method and associated toolchain enabled these

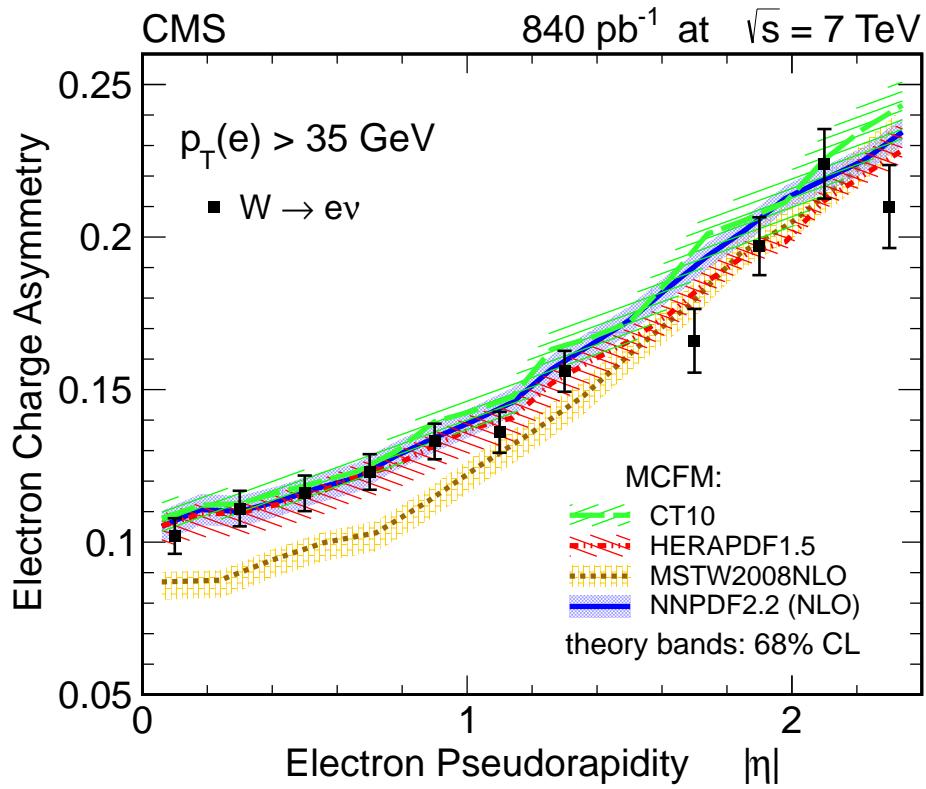


Figure 5.3: Comparison of NNPDF2.2 predictions with updated CMS  $W$  asymmetry measurement at  $840 \text{ pb}^{-1}$ . The comparison also includes the theory predictions from CT10, MSTW 2008 and HERAPDF 1.5. Agreement is generally very good for the PDF sets, although the MSTW2008 set demonstrates a significant discrepancy. Figure from [?].

fits to be performed without the requirements of extremely long fitting times, potentially requiring weeks of computer time per replica for a standard fit on a typical 2.4GHz Intel Xeon processor with the earlier technology. In this section we shall discuss the NNPDF2.3 fit, the successor to the NNPDF2.2 fit in that an updated and enlarged LHC dataset is included in a full NNPDF fit. We shall outline the datasets included in the determination, along with a discussion of methodological improvements made, as several optimisations were enabled by the faster fitting framework.

### 5.2.1 NNPDF2.3 dataset and methodology

For the NNPDF2.3 determination, the electroweak data included in NNPDF2.2 has been upgraded. From CMS the  $840\text{pb}^{-1}$   $W$  electron asymmetry data [?] replaces the previous measurement. The full  $W/Z$  (pseudo)rapidity distributions replace the asymmetry measurements for ATLAS, based upon  $35\text{ pb}^{-1}$  of 2010 data [?]. From LHCb, the  $W^\pm$  distributions in the forward region were included [?]. Beyond the electroweak sector, the ATLAS 2010 inclusive jet data was also included to obtain an additional handle upon the gluon. At the time of publication, the NNPDF2.3 dataset included all relevant published LHC data with publicly available covariance matrices. Theoretical predictions for these observables were implemented as FK tables obtained via APPLgrid files from MCFM for the electroweak processes, and `nlojet++` for the jet data.

Figure ?? demonstrates the additional reach of the NNPDF2.3 dataset upon the addition of the LHC data. The electroweak measurements extend those performed at the Tevatron to considerably lower values of parton- $x$ . The inclusive jet data spans a large range in kinematics, providing points at large and small- $x$  across a wide range of scales. Examining the description provided by earlier PDF sets, Table ?? demonstrates the agreement at NLO and NNLO of the previous 2.1 PDF set to the new experimental data. While fair agreement is reached for most sets the description is often sub-optimal therefore the data can provide useful additional constraints. This is particularly evident for the ATLAS electroweak data at NNLO ( $\chi^2/N_{\text{dat}} = 2.21$ ) and the CMS  $W$  electron asymmetry data at NLO ( $\chi^2/N_{\text{dat}} = 2.02$ ).

In the 2.3 fit, the theoretical prediction mechanism for all previously included observables was converted to the FK procedure, leading to a substantial decrease in fitting times. These speed improvements were exploited in order to perform a more aggressive fitting procedure. The NNPDF minimisation procedure involves a genetic algorithm where the best fit network per iteration undergoes a set of

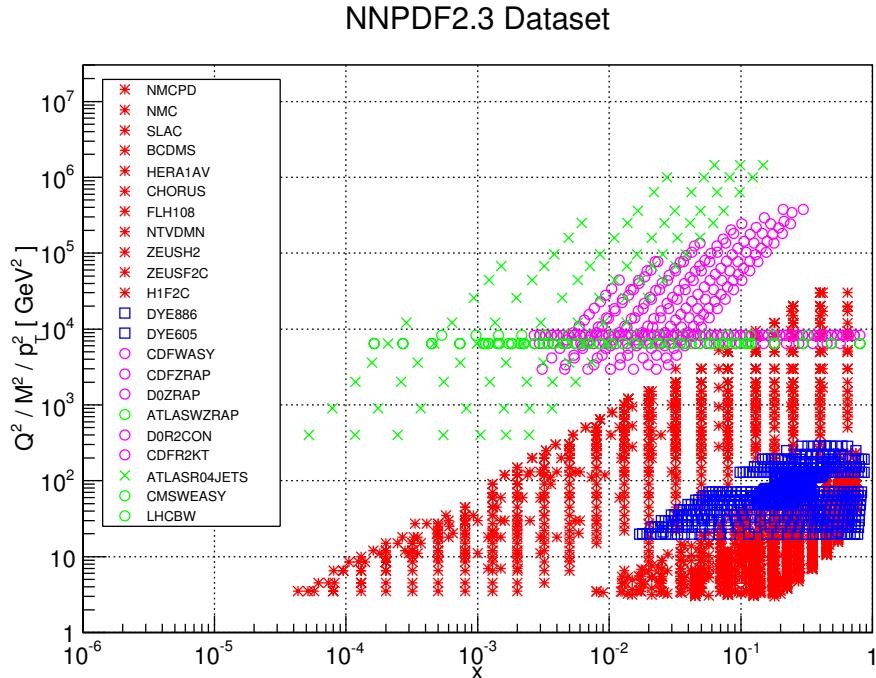


Figure 5.4: Kinematic distribution of points in the NNPDF2.3 analysis. The green points show the LHC data which was added to the analysis over the NNPDF2.1 dataset, and demonstrates the additional kinematic reach of the dataset. DIS datapoints are represented in red, the fixed-target Drell-Yan data in blue, and Tevatron data in pink. In the case of data with two hadrons in the initial state, the smaller parton- $x$  value is plotted.

random adjustments or ‘mutations’ , the best of which is selected for the next iteration. In the NNPDF2.1 NLO fits two genetic algorithm epochs are used. The first, or ‘ $a$ ’ phase with  $N_{\text{mut}}^a = 80$  mutants and the second ‘ $b$ ’ phase with  $N_{\text{mut}}^b = 10$  mutants per generation. This was upgraded to the more explorative settings of  $N_{\text{mut}}^b = 30$  mutants in the second epoch. The maximum number of training generations was extended to  $N_{\text{gen}}^{\text{max}} = 50,000$  generations from the 30,000 used in the NNPDF2.1 series. For mutation rates, the number of mutations  $N_{\text{mut}}$  were increased for a number of PDF combinations in order to better explore the fit quality minima, and the mutation sizes  $\eta$  optimised on a PDF by PDF basis. In Table ?? we summarise the modifications made in the genetic algorithm minimisation in terms of the parameters that have been modified.

	NNPDF2.1	
	NLO	NNLO
ATLAS $W/Z$	1.57	2.21
LHCb $W$	0.89	1.13
CMS $We$ Asy	2.02	1.27
ATLAS Jets	1.06	0.95

Table 5.4: Description of the NNPDF2.3 LHC dataset provided by the NNPDF2.1 PDF set, provided as  $\chi^2$  per degree of freedom,  $\chi^2/N_{\text{dat}}$ . The new data shows good consistency with the previous data available in NNPDF2.1 however there is room for improvement upon the inclusion of the data.

Additionally the parameters controlling the dynamical stopping criterion were tightened, requiring a clearer overlearning signal from the cross-validation.

Other small methodological changes included the addition of a maximum  $\chi^2$  criterion, whereby replicas with a fit quality outside a  $4\sigma$  band in  $\chi^2$  are vetoed from the ensemble as outliers. The training/validation split used in the cross-validation was also modified for experiments with smaller than 30 data points, where as of NNPDF2.3 all of these points enter in the training set to prevent them from *underlearning* or being ignored in the fit in favour of larger datasets.

With these methodological modifications, a number of determinations were performed to different datasets. Firstly the global fit was performed to the entire 2.1 dataset with the addition of the LHC data. This was followed by a ‘noLHC’ fit which applied the methodological improvements to the same dataset as NNPDF2.1, both in order to understand the impact of the methodological modifications upon the fit and to provide a set for applications where the inclusion of an LHC dataset is undesirable. Finally a collider-only determination was performed, which excluded the older low scale fixed-target data in an attempt to reduce the effect of nuclear, higher twist and non-perturbative corrections.

	$N_{\text{gen}}^{\text{mut}}$	$N_{\text{gen}}^{\max}$	$N_{\text{mut}}^a$	$N_{\text{mut}}^b$
2.1 NLO	2500	30000	80	10
2.1 NNLO	2500	30000	80	30
2.3 NLO	2500	50000	80	30
2.3 NNLO	2500	50000	80	30

PDF	2.1 NLO		2.1 NNLO and 2.3	
	$N_{\text{mut}}$	$\eta^k$	$N_{\text{mut}}$	$\eta^k$
$\Sigma(x)$	2	10, 1	2	10, 1
$g(x)$	2	10, 1	3	10, 3, 0.4
$T_3(x)$	2	1, 0.1	2	1, 0.1
$V(x)$	2	1, 0.1	3	8, 1, 0.1
$\Delta_S(x)$	2	1, 0.1	3	5, 1, 0.1
$s^+(x)$	2	5, 0.5	2	5, 0.5
$s^-(x)$	2	1, 0.1	2	1, 0.1

Table 5.5: Summary of modifications to the genetic algorithm minimisation between NNPDF2.1 and NNPDF2.3. The table on top describes the number of mutations ( $N_{\text{mut}}$ ) and the number of generations ( $N_{\text{gen}}$ ) in the different training epochs, while the lower table shows the number of mutations per PDF and the corresponding mutation sizes. Table from [?].

### 5.2.2 NNPDF2.3 results

Here we shall discuss some of the results obtained in the NNPDF2.3 family of PDF determinations. Assisted by the developments in the fitting methodology, all of the 2.3 determinations were able to provide an excellent description of their included datasets. Table ?? details the agreement through the  $\chi^2$  measure to each experiment in the analysis, for every variation of the dataset.

Experiment	NNPDF2.1				NNPDF2.3					
	Global		Global Fit		Global RW		noLHC		Collider	
	NLO	NNLO	NLO	NNLO	NLO	NNLO	NLO	NNLO	NLO	NNLO
Total	1.145	1.162	1.101	1.139	1.105	1.139	1.101	1.142	0.971	0.993
NMC-pd	0.97	0.93	0.95	0.95	0.93	0.93	0.93	0.94	[5.33]	[5.13]
NMC	1.68	1.58	1.61	1.59	1.62	1.57	1.59	1.56	[1.89]	[1.83]
SLAC	1.34	1.04	1.24	1.00	1.27	1.01	1.28	1.04	[1.72]	[1.41]
BCDMS	1.21	1.29	1.20	1.28	1.20	1.28	1.20	1.28	[1.85]	[2.15]
CHORUS	1.10	1.08	1.10	1.07	1.10	1.06	1.09	1.07	[1.73]	[1.70]
NTVDMN	0.70	0.50	0.43	0.56	0.42	0.51	0.42	0.48	[26.69]	[21.13]
HERA1-AV	1.04	1.04	1.00	1.01	1.00	1.02	1.01	1.03	0.97	0.99
FLH108	1.34	1.23	1.29	1.20	1.29	1.20	1.29	1.21	1.35	1.25
ZEUS-H2	1.21	1.21	1.20	1.22	1.20	1.22	1.20	1.22	1.29	1.32
ZEUS $F_2^c$	0.75	0.81	0.82	0.90	0.80	0.90	0.81	0.86	0.71	0.77
H1 $F_2^c$	1.50	1.44	1.59	1.53	1.57	1.52	1.58	1.49	1.33	1.30
DYE605	0.94	1.08	0.86	1.04	0.88	1.04	0.85	1.06	[3.58]	[1.02]
DYE886	1.42	1.69	1.27	1.58	1.27	1.55	1.24	1.55	[5.65]	[5.14]
CDF $W$ asy	1.88	1.63	1.57	1.64	1.57	1.72	1.45	1.67	1.05	1.21
CDF $Z$ rap	1.77	2.38	1.80	2.03	1.77	2.17	1.76	2.13	1.32	1.37
D0 $Z$ rap	0.57	0.67	0.56	0.61	0.57	0.63	0.57	0.63	0.56	0.58
ATLAS $W, Z$	[1.57]	[2.21]	1.26	1.43	1.31	1.65	[1.37]	[1.94]	1.02	1.05
CMS $W$ el asy	[2.02]	[1.27]	0.82	0.81	1.09	0.99	[1.32]	[1.20]	0.87	0.85
LHCb $W$	[0.89]	[1.13]	0.67	0.83	0.77	0.98	[0.76]	[1.03]	0.74	0.72
CDF RII $k_T$	0.68	0.65	0.60	0.68	0.61	0.67	0.60	0.67	0.60	0.59
D0 RII cone	0.90	0.98	0.84	0.94	0.84	0.93	0.84	0.94	0.85	0.92
ATLAS jets	[1.06]	[0.95]	1.00	0.94	1.00	0.92	[1.01]	[0.94]	0.98	0.93

Table 5.6: The fit quality to each individual dataset in the global NNPDF2.3 determination provided by various NNPDF sets. The global, noLHC and collider only 2.3 determinations are shown along with the NNPDF2.1 values for comparison. Additionally the values for a reweighting of 2.1 with LHC data is shown in order to test the efficacy of the fitting procedure. The figures in square brackets are for datasets that were not included in the associated PDF set.

The total  $\chi^2$  values achieved by the global fits were 1.101 at NLO and 1.139 at NNLO, both indicating fine agreement with the experimental data and demonstrating improvement over the fit quality obtained in the NNPDF2.1 series. The noLHC fits obtained similar levels of fit quality, while the collider only determinations demonstrated the excellent consistency in the dataset with  $\chi^2$

values of 0.971 and 0.993 for the NLO and NNLO fits respectively.

Notably the collider only dataset fails to describe the older fixed-target data, particularly the NuTeV dimuon measurements, by a large margin. A  $\chi^2$  value of 26.69 at NLO to the NuTeV data suggests that the collider-only dataset may be in some tension with the older, low scale measurements. Despite this the global fit is able to provide a good description of both the collider only and fixed target data simultaneously, therefore any tension present between the datasets remains at the moment compatible within experimental errors.

The average training length at NLO is predictably extended in 2.3 over 2.1. The more stringent stopping condition leading to more replicas running for the extended maximum  $N_{\text{gen}} = 50,000$ . The training length comparison is shown in Figure ??.

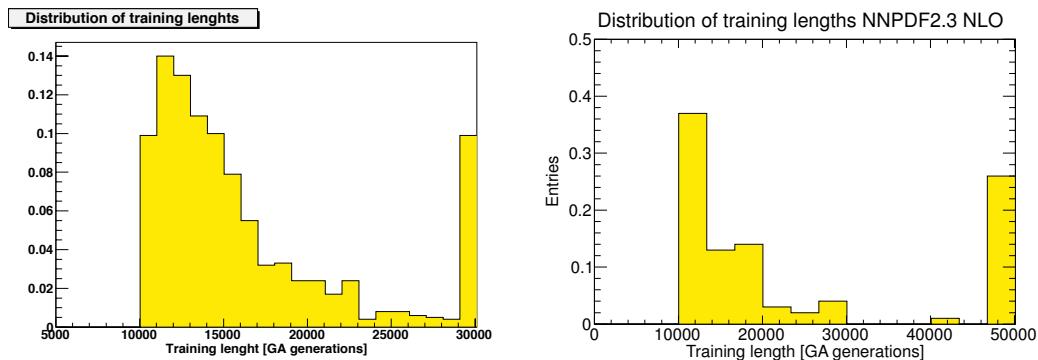


Figure 5.5: Replica training lengths in NNPDF2.1 NLO (left) and NNPDF2.3 NLO (right). These histograms display the relative frequency of replicas stopping in bins along the full training length. In NNPDF2.3 both the maximum number of generations was increased to 50,000, and the criteria governing the replica stopping was tightened, causing more replicas to stop later.

We shall now examine the changes between the NNPDF2.1 and NNPDF2.3 distributions at the level of PDFs. Firstly discussing the impact of the methodological changes to the NNPDF determination by examining the NNPDF2.3 noLHC fits, before moving on to look at the direct impact of the LHC dataset by performing comparisons of the noLHC and full 2.3 datasets. Finally we shall

discuss the impact of the LHC data upon a collider only determination. The issue of the strange content of the proton is a particularly delicate one, and therefore will be discussed separately from the other five light quark distributions.

### NNPDF2.3 noLHC

The NNPDF2.3 noLHC set has two primary uses. To understand the improvements made in the NNPDF methodology by applying the updated procedure to the older dataset, and for applications such as BSM searches at the LHC where a dataset without the influence of LHC data may be desirable. Here we shall directly compare the 2.3 noLHC results with NNPDF2.1 to see the methodological improvement. These improvements were expected to be clearer in the NLO PDF sets, as for NNPDF2.1 NNLO several of the improvements in the minimisation were already implemented. Aside from the strange sector (which will be discussed later), the methodological changes largely only impact the gluon and singlet distributions, with other distributions undergoing small changes, so we shall restrict ourselves here to comparisons of the gluon and singlet PDFs. The upper section of Figure ?? compares NNPDF2.1 and NNPDF2.3 noLHC at NLO, for those PDFs most affected by the improvements; the singlet and gluon. The clearest improvements are in the low- $x$  region, where the combination of more aggressive minimisation and tighter stopping criteria lead to substantially smaller uncertainty in the singlet, and a moderate shift in the gluon. These improvements suggests that there was potentially a degree of underlearning present in the small- $x$  region of NNPDF2.1 generated by stopping too early.

The lower part of Figure ?? demonstrates the same comparison for the NNLO determination. From this figure it is clear that the degree of underlearning present in the NLO fit was avoided by the use of the updated fit settings, leading to slightly narrower uncertainty bands. The relatively insignificant differences remaining due to the presence of more data in the training sets, although the

difference remains statistically insignificant at the level of PDFs.

### NNPDF2.3 global

The NNPDF2.3 global set includes all of the methodological improvements along with the constraints from the new LHC dataset. There are therefore comprehensive improvements available in the 2.3 set over 2.1, both at NLO and NNLO in QCD. Figure ?? shows the same comparison as in Figure ??, but including the impact of the LHC data by comparing NNPDF2.1 to the full global NNPDF2.3 set. As much of the improvements are driven by methodology, the largest modifications in the global comparison can also be found in the gluon and singlet distributions. To obtain a clearer view of the impact of the LHC data upon the PDFs we can compare the 2.3 noLHC fit with the global determination, with the only differences in the two sets due to the LHC data.

In comparing the 2.3 global and noLHC sets, the clearest improvements can be found in the singlet, gluon and valence sectors as would be expected from the expanded dataset. Figure ?? compares these distributions at NLO and NNLO to study the influence of the new data. In the singlet sector, the LHC data prefers a rather higher value for the PDF in the small- $x$  region, with the central value being systematically higher below  $x \sim 0.1$ , an effect which is clearer at NNLO. In the NLO fit there is a broadening of uncertainties for the extrapolation region  $x < 10^{-4}$ , but a moderate degree of uncertainty reduction in the data region. For the NNLO singlet the uncertainties are larger over a broad kinematic range, generated by the larger upwards shift preferred by the LHC data at NNLO.

The gluon distribution at NLO enjoys a great deal of consistency between the 2.3 noLHC and 2.3 global fits. With the additional LHC data contributing to a broad reduction of uncertainties in the data region. The NNLO fit, while demonstrating a good deal of consistency, does not make any significant reduction in uncertainty outside the region of  $x \sim 10^{-2}$ .

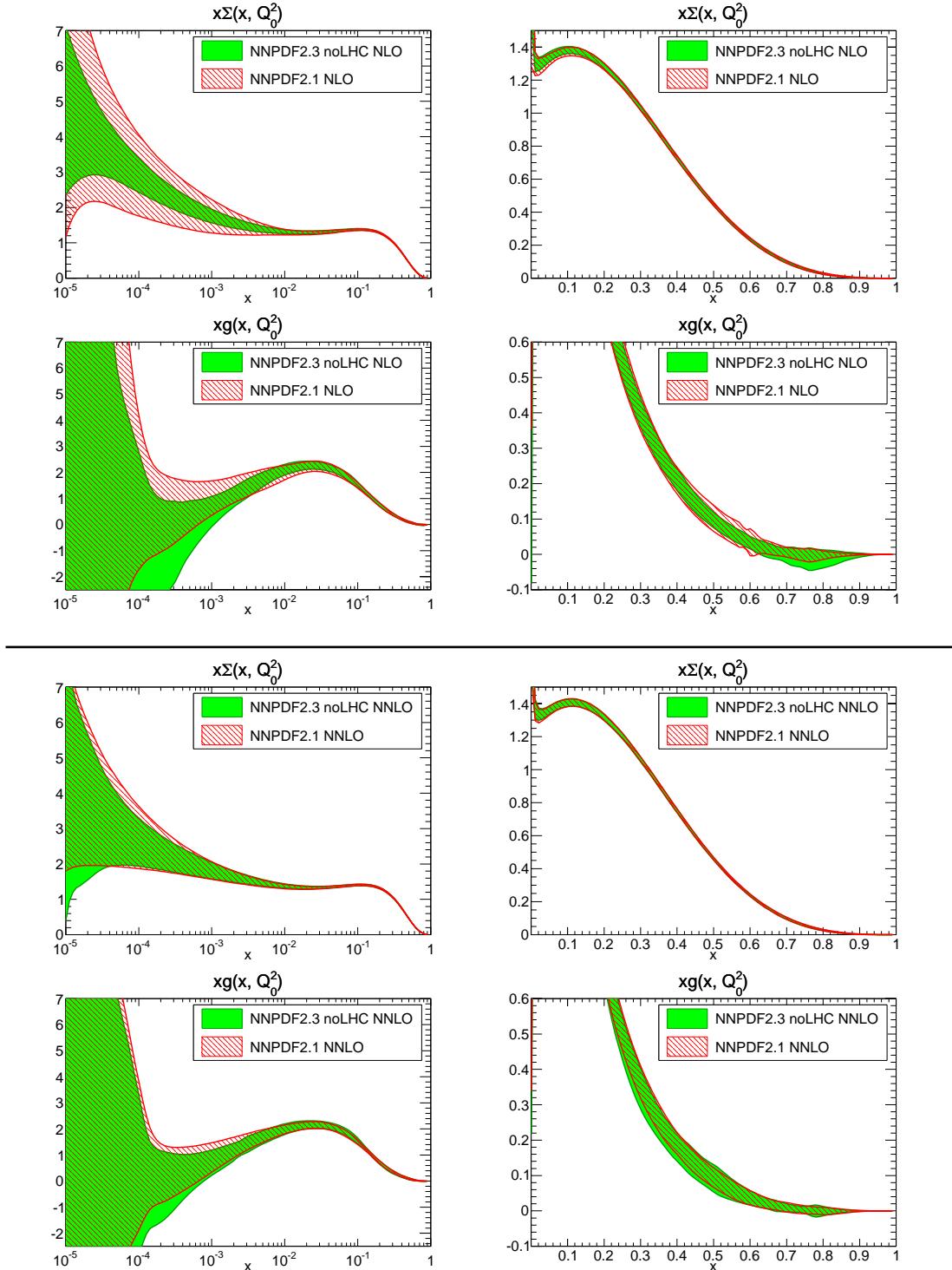


Figure 5.6: The impact of the improved methodology in NNPDF2.3 against NNPDF2.1 in the gluon and singlet sectors for the NLO (top) and NNLO (bottom) distributions. The red curves show the results of NNPDF2.1 while the green curves show NNPDF2.3 noLHC. Figures on the left are shown in a logarithmic scale in  $x$ .

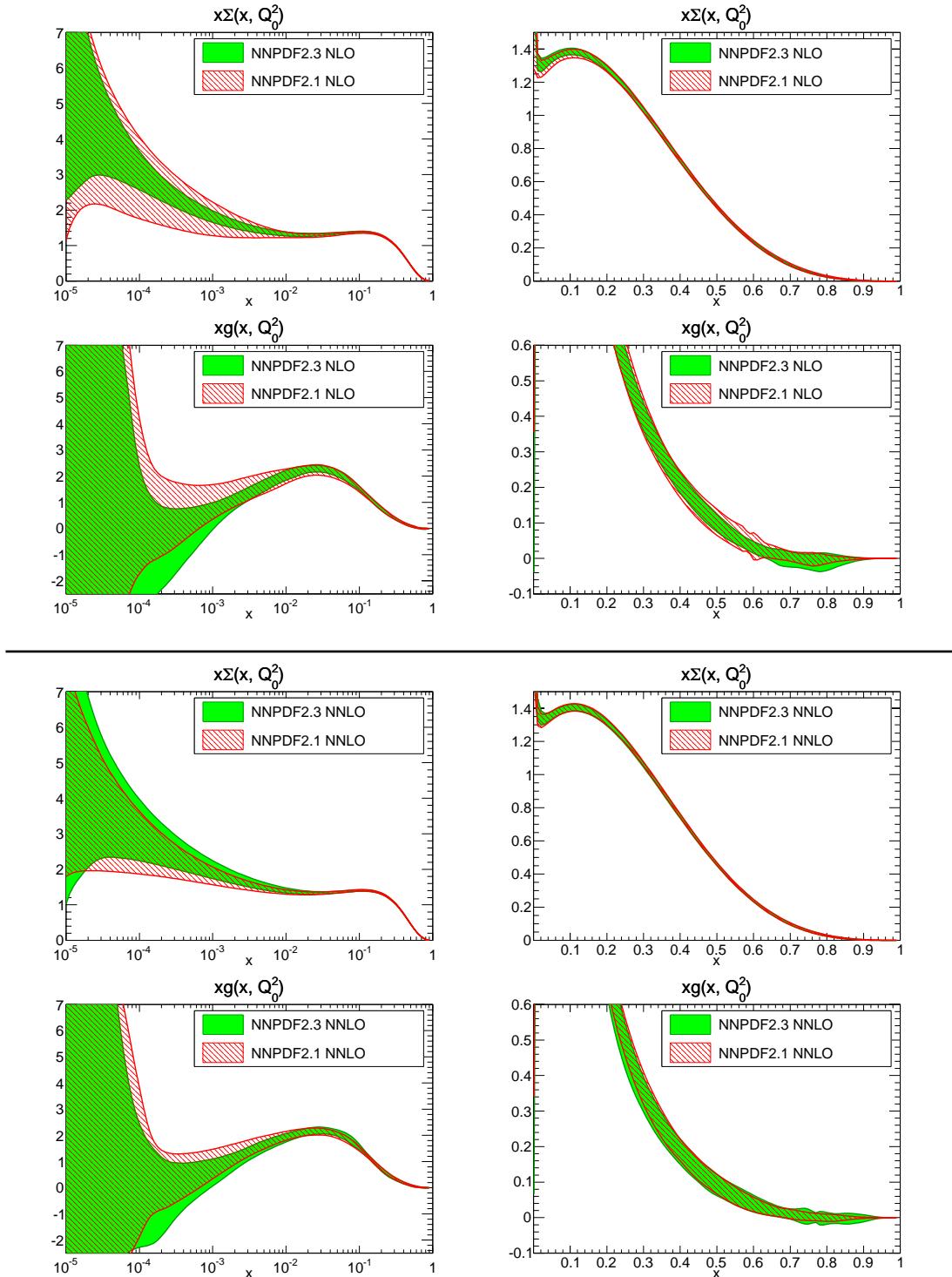


Figure 5.7: A comparison of the NNPDF2.1 and NNPDF2.3 global determinations at NLO (top) and NNLO (bottom). The gluon and singlet distributions are shown, with a logarithmic  $x$  scale on the left, and linear on the right. Red distributions are those given by the NNPDF2.1 set, while green represent NNPDF2.3.

The PDF benefiting the most from the inclusion of the LHC data is the NLO valence distribution, where significant reductions in uncertainty are achieved across a wide kinematic range. Despite this improvement, the NNLO fit is not able to make such significant gains on the basis of the new data, with improvements constrained to the moderate to large- $x$  region.

Figure ?? specifically demonstrates the changes in uncertainties upon the addition of the new data. While uncertainty reduction has been achieved for some PDF combinations, several areas undergo an increase in their uncertainties due to central value shifts.

### NNPDF2.3 collider only

In order to investigate whether an NNPDF determination of a collider only PDF determination is viable with the current collider dataset, we now compare the resulting distributions from the NNPDF2.3 global and collider-only fits. In Figure ?? the distributions for the singlet, gluon, sea-asymmetry and triplet are shown for the two fits at NNLO. The combination of HERA DIS data along with Tevatron and LHC inclusive jet data ensure that the gluon and singlet, although deviating not insignificantly from the global fit, are well constrained by data. The preference for a higher singlet distribution by the LHC data seen in the comparison between the global and noLHC fits is very clear in this comparison, with the collider only dataset preferring a significantly higher singlet also. The gluon distribution demonstrates also rather different shape in the medium- $x$  region. One may therefore at first be tempted to prefer the collider only determination for phenomenology, given its theoretically cleaner underpinnings. However the ability of the fit to obtain a good handle on PDF combinations involving flavour separation is substantially reduced in the collider only fit. The lower two plots in Figure ?? demonstrate that the collider dataset is not able to provide sufficient constraints for these combinations even after the addition of

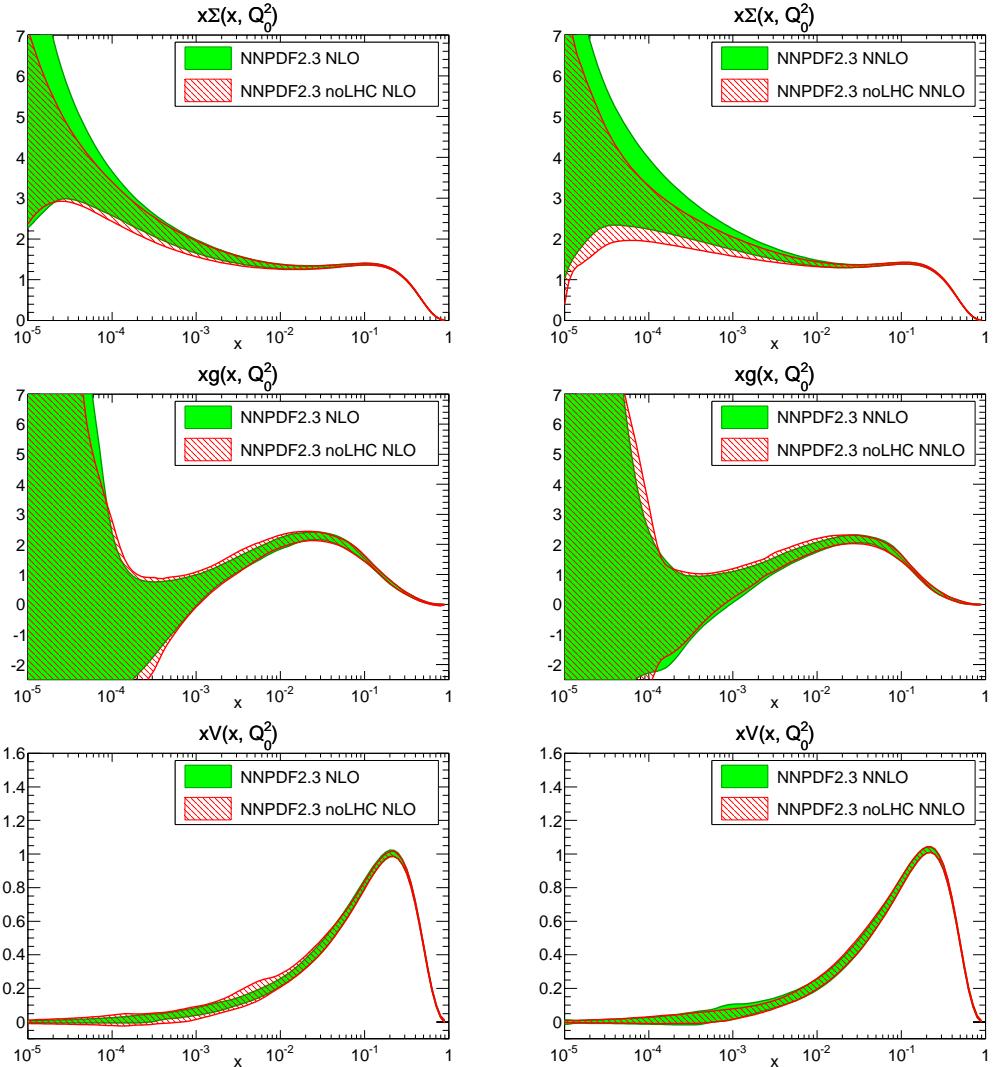


Figure 5.8: Comparison of NNPDF2.3 and NNPDF2.3 noLHC at NLO (left) and NNLO (right) for the singlet (top), gluon (middle) and valence (bottom) distributions. The figures therefore show directly the influence of the LHC data in the fit.

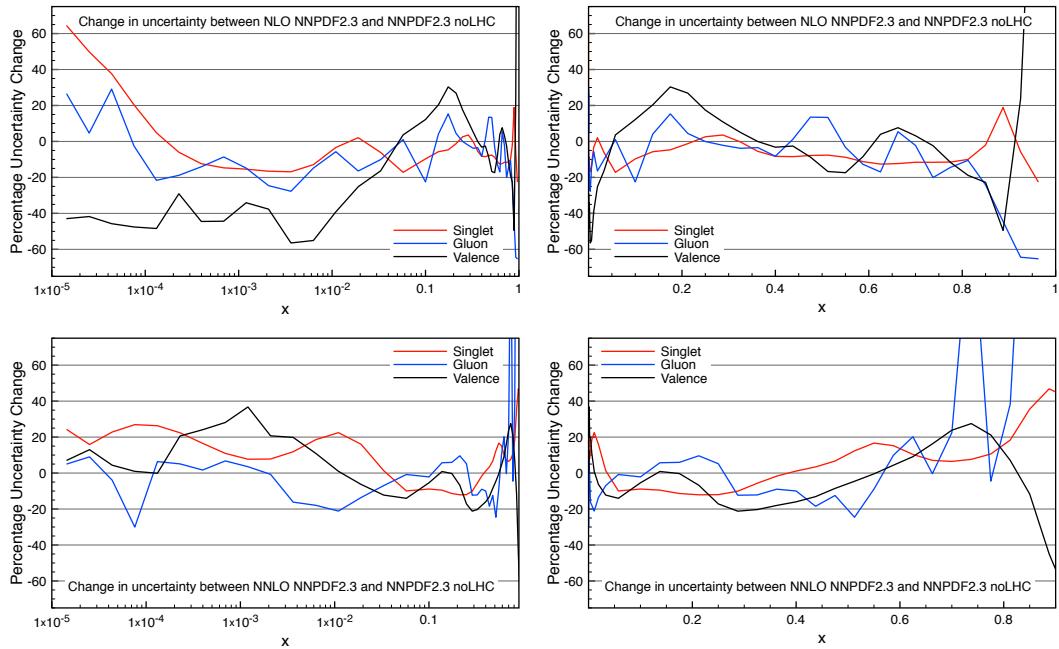


Figure 5.9: Uncertainty change in NNPDF2.3 under the addition of LHC data. The top figures represent the percentage change in uncertainty between NNPDF2.3 global and NNPDF2.3 noLHC at NLO, while the bottom plots show the equivalent comparison at NNLO.

the LHC dataset.

The NNPDF2.3 collider only dataset therefore remains too imprecise to provide an accurate determination of flavour-separation, and is therefore unsuitable for applications sensitive to such parton combinations. Nevertheless, a great deal of progress is evident when we compare PDFs obtained via a collider-only fit to the pre-LHC NNPDF2.1 dataset, and those obtained with the new LHC data. Examining the NNLO PDFs where methodological differences are slight, Figure ?? compares the light quark distributions  $u$  and  $d$  between NNPDF2.1 collider only and NNPDF2.3 collider only, with the only significant difference being the presence of the LHC dataset in the 2.3 determination.

From the figure it is clear that the LHC data provides very large constraints upon the collider only distributions when we compare to the version available in the NNPDF2.1 series. Examining the gluon and singlet PDFs in Figure ?? we

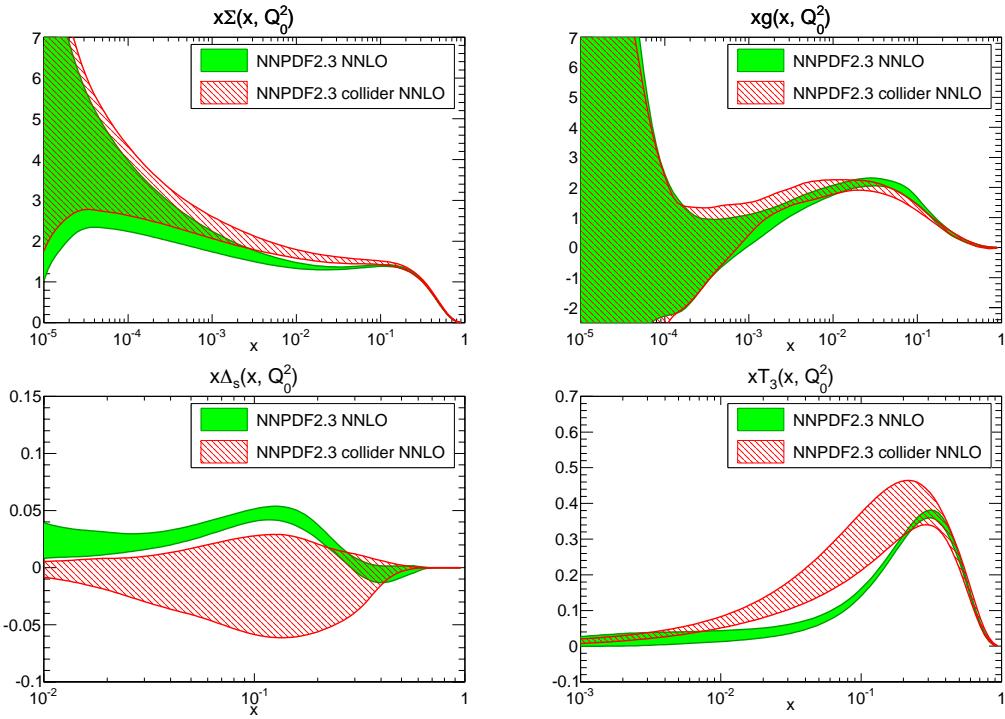


Figure 5.10: NNPDF2.3 collider only compared to the global determination at NNLO. The red distributions shown are those determined via a fit to a collider-only dataset, while the green curves show the results of the NNPDF2.3 global fit.

see that the improvements made in the light quarks carry through to the quark singlet. The gluon distribution was already relatively well determined in the 2.1 series due to the Tevatron jet data, therefore it does not experience such a large improvement.

In Figure ?? we can see explicitly the impact the new data has upon collider only uncertainties. Across nearly the whole kinematic range, very substantial improvements can be seen in both the singlet and valence distributions. The results from the LHC are therefore vital in providing a handle on the collider only distributions, and updated measurements may be able to bring the accuracy of such determinations to near the level of the global fits.

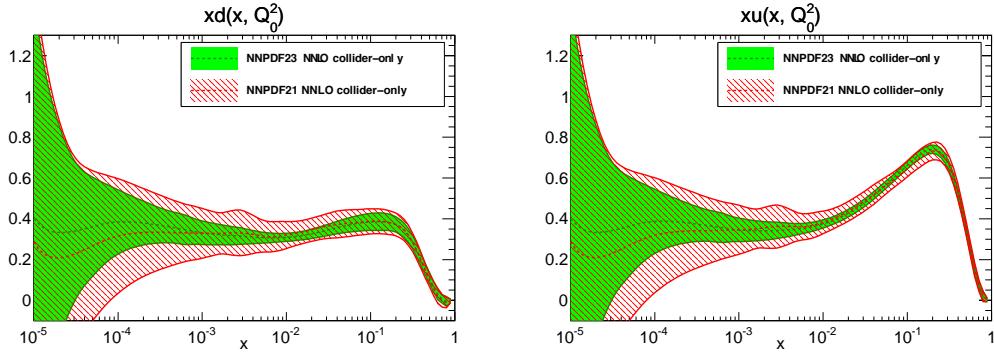


Figure 5.11: Impact of NNPDF2.3 LHC data upon collider only determinations for light flavour PDFs. The green solid curves show the collider only results including the LHC dataset and the red dashed curves show the results of the NNPDF2.1 collider only fits.

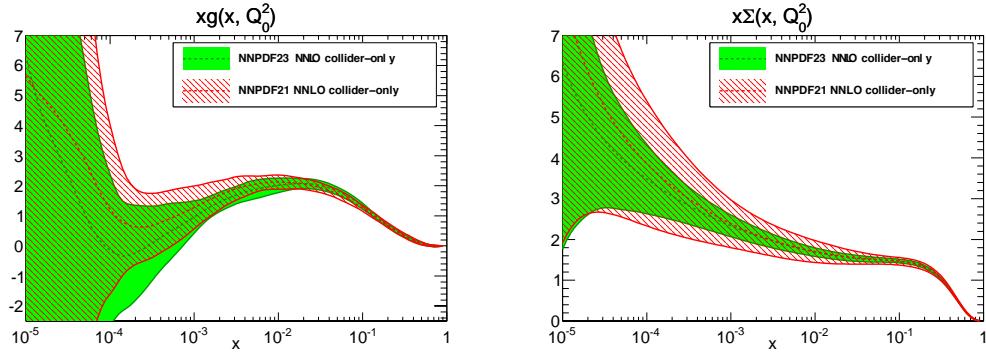


Figure 5.12: Impact of NNPDF2.3 LHC data upon collider only determinations for singlet and gluon PDFs. The green curves show the collider only results including the LHC dataset and the red curves show the results of the NNPDF2.1 collider only fits.

### 5.2.3 Proton strangeness

The issue of the strange content of the proton is a particularly interesting one, and has been the source of discussion due to new results and analyses arising from the LHC experiments. A particular complication lies in the treatment of the NuTeV dimuon data. In the NNPDF2.1 series of fits the expression used for the dimuon data suffered from an error originating from the heavy quark mass handling. Specifically, Eqn. 34 of Ref. [?] presented an incorrect expression for the charm production reduced cross-section in neutrino charged current DIS, where

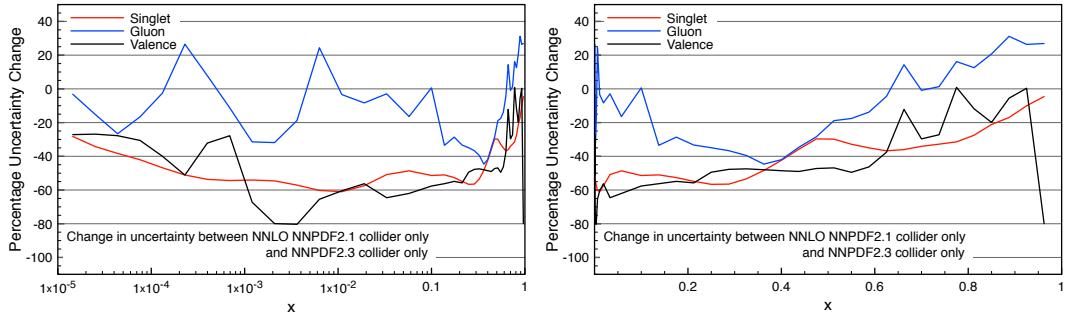


Figure 5.13: Impact of LHC data upon collider only fit uncertainties. Figures show the percentage improvement in the NNPDF2.3 collider only uncertainty compared to the NNPDF2.1 collider only fit, for the singlet, gluon and valence distributions at NNLO.

the correct expression reads

$$\begin{aligned} \tilde{\sigma}^{\nu(\bar{\nu}),c}(x, y, Q^2) &\equiv \frac{1}{E_\nu} \frac{d^2 \sigma^{\nu(\bar{\nu}),c}}{dx dy}(x, y, Q^2) \\ &= \frac{G_F^2 M_N}{2\pi(1+Q^2/M_W^2)^2} \left[ \left( \left( Y_+ - \frac{2M_N^2 x^2 y^2}{Q^2} - y^2 \right) + y^2 \right) F_{2,c}^{\nu(\bar{\nu})}(x, Q^2) \right. \\ &\quad \left. - y^2 F_{L,c}^{\nu(\bar{\nu})}(x, Q^2) \pm Y_- x F_{3,c}^{\nu(\bar{\nu})}(x, Q^2) \right], \end{aligned} \quad (5.1)$$

where here  $x$  and  $y$  are the usual DIS kinematic variables,  $Y_\pm = 1 \pm (1-y)^2$  and the momentum transfer is given by  $Q^2 = 2M_N E_\nu x y$ . The expression in Ref. [?] differs from this by a spurious additional  $\left(1 + \frac{m_c^2}{Q^2}\right)$  term which was corrected prior to the NNPDF2.3 determination.

This error affected only the predictions for the NuTeV data, and consequently after the error was corrected the impact upon PDFs themselves was largely restricted to the strange sector. The impact of the error is shown in Figure ??, where it can be clearly seen that the error led to a small suppression of the total strange distribution across most of the kinematic range, peaking at around half a standard deviation.

The shift towards slightly higher total strangeness is continued upon the addition of the LHC dataset. Figure ?? shows how the strange sea distribution changes under the addition of the new data. The electroweak measurements present in the LHC dataset seem to marginally prefer a slightly larger strange sea at small- $x$  for both the NLO and NNLO distributions.

To investigate the relative contribution of the strange sea with respect to the light quark sea, a commonly used measure [?, ?, ?, ?] is the integrated ratio of the two PDF combinations,

$$K_s = \frac{\int_0^1 dx x (s(x, Q^2) + \bar{s}(x, Q^2))}{\int_0^1 dx x (\bar{u}(x, Q^2) + \bar{d}(x, Q^2))}. \quad (5.2)$$

In most global determinations a significant suppression of the strange sea is typically observed at low scales, with  $K_s < 1$ . In Table ?? we see the results for  $K_s$  obtained through the NNPDF2.1, NNPDF2.3 and NNPDF2.3 noLHC sets demonstrating such a suppression at NNLO. The impact of the incorrect dimuon treatment in NNPDF2.1 is manifest in an exaggerated level of suppression, although it remains consistent with the newer determinations within uncertainties. The preference for a larger strange sea by the LHC measurements is also demonstrated in the difference between NNPDF2.3 and the noLHC dataset fit.

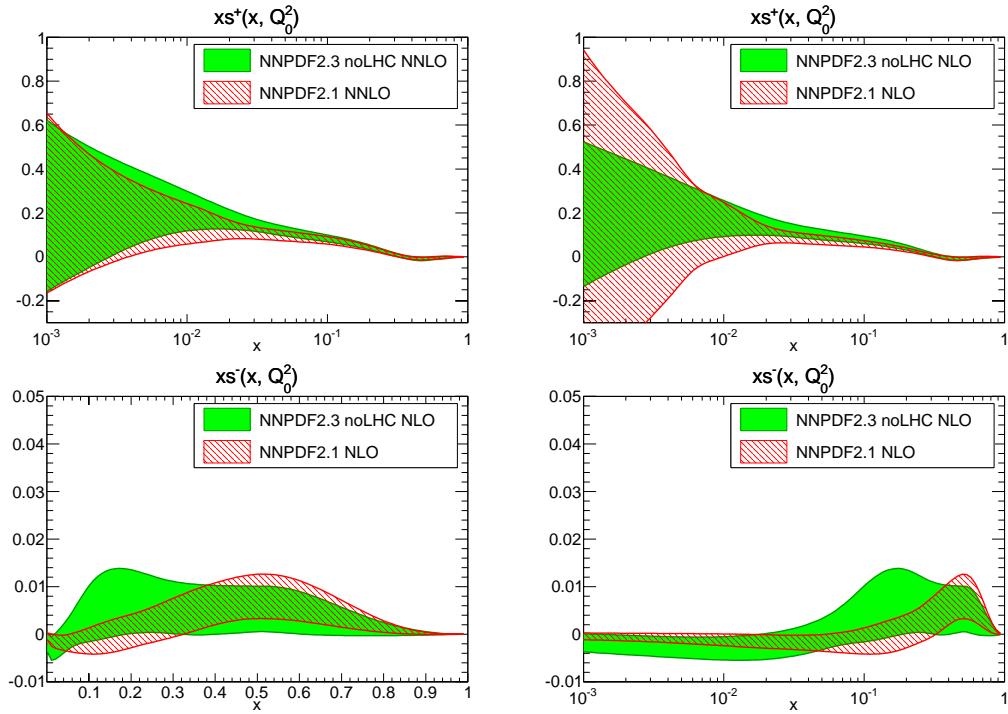


Figure 5.14: Total strangeness and strange valence distributions compared between NNPDF2.1 and NNPDF2.3 noLHC. The NLO bands demonstrate also the improvements due to the more aggressive minimisation, particularly evident at low- $x$ .

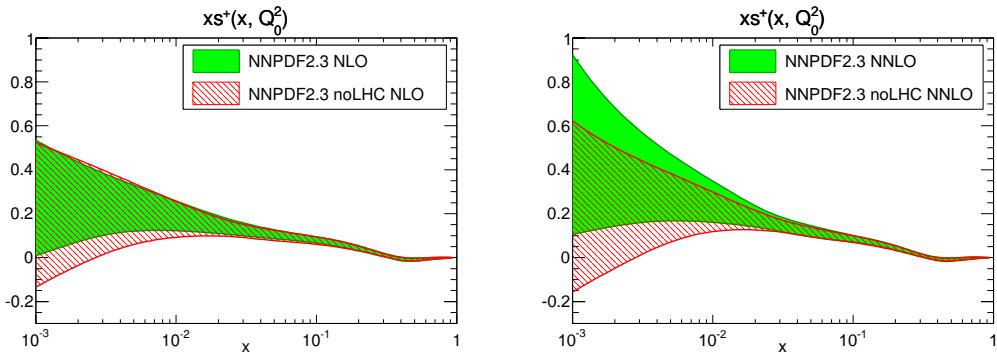


Figure 5.15: Strange sea distributions in NNPDF2.3 and NNPDF2.3 noLHC, for the NLO (left) and NNLO (right) PDF sets. The NNPDF2.3 global set, differing from the noLHC set by the inclusion of LHC measurements, prefers a marginally larger strange distribution.

PDF	$K_s(2 \text{ GeV}^2)$	$K_s(M_Z^2)$
NNPDF2.1	$0.26^{+0.08}_{-0.08}$	$0.63^{+0.04}_{-0.05}$
NNPDF2.3 noLHC	$0.30^{+0.09}_{-0.08}$	$0.65^{+0.05}_{-0.05}$
NNPDF2.3	$0.35^{+0.10}_{-0.08}$	$0.68^{+0.05}_{-0.05}$

Table 5.7: Strange sea suppression in NNPDF2.3 and NNPDF2.1, with the uncertainties given by the 68% confidence interval.

Such a strange sea suppression was challenged by an ATLAS determination of the strange content of the proton [?] based upon a fit to a combined HERA DIS and ATLAS  $W$  and  $Z$  production dataset. Defining a more exclusive measure, the ratio of the strange sea to twice the  $\bar{d}$  distribution at specific points of  $x$  and  $Q^2$ ,

$$r_s(x, Q^2) = \frac{s(x, Q^2) + \bar{s}(x, Q^2)}{2\bar{d}(x, Q^2)}, \quad (5.3)$$

the ATLAS study reported values that significantly differed from the typical results of global fits, with the most extreme disagreement with the NNPDF2.1 set where the two values are separated by more than two sigma. The disagreement is particularly large in the region  $x = 0.023$ , at the initial scales, as shown in the ATLAS plot in Figure ??, where the ATLAS result is consistent with no suppression of the strange sea,  $r_s \sim 1$ .

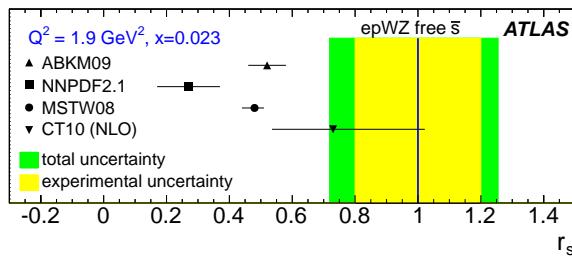


Figure 5.16: ATLAS determination of strange sea suppression at  $x = 0.023$  for a number of PDF sets. The ATLAS result is consistent with no suppression for the strange distributions. Figure from [?].

The impact of the NNPDF2.3 LHC dataset clearly has a preference for a more strange-symmetric sea, as is particularly demonstrated upon the inclusion of the LHC dataset (including the ATLAS data used in their strangeness analysis) to the NNPDF2.1 collider only strange distribution. Figure ?? demonstrates the extensive constraint placed upon the NNPDF2.1 collider only set by the LHC electroweak data in the strange sector, and a clear preference for a larger strange sea. Despite this preference, the results of the global fit remain consistent within the larger uncertainties of the NNPDF collider only determination.

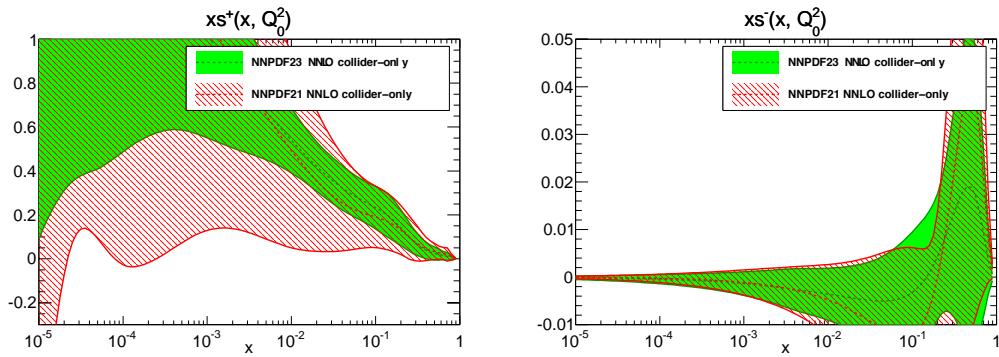


Figure 5.17: Impact of LHC data on collider only strangeness. The strange sea (left) and strange valence (right) distributions are plotted at NNLO, comparing the NNPDF2.1 and NNPDF2.3 collider only results. A very significant impact upon the total strangeness uncertainties can be observed, however little constraint is afforded to the valence distribution.

To investigate the ATLAS result, an NNPDF2.3 fit was performed to the same dataset as in the ATLAS ‘epWZ’ fit. The results of this fit for both the  $r_s$  values quoted by the ATLAS collaboration, and the integrated  $K_s$  values are shown in Figure ???. While the results of the NNPDF2.3 series fits to global datasets remain incompatible with the ATLAS result, the results of all of the fits are perfectly compatible within the very large uncertainties of the NNPDF fit to the restricted ATLAS and HERA dataset used for the ATLAS result.

The much greater uncertainty present in the NNPDF fit to the HERA and ATLAS  $W/Z$  dataset suggests that the uncertainty in the ATLAS result was

underestimated significantly, a conclusion also reached by similar analyses by the MSTW and ABM groups [?, ?].

Measurements of  $W + c$  production, particularly sensitive to the strange fraction can provide additional information for future fits. As an example, the CMS  $W + c$  measurement based upon  $5.0\text{fb}^{-1}$  of 7 TeV data [?] demonstrates good agreement with the results of global PDF sets, as shown in Figure ??.

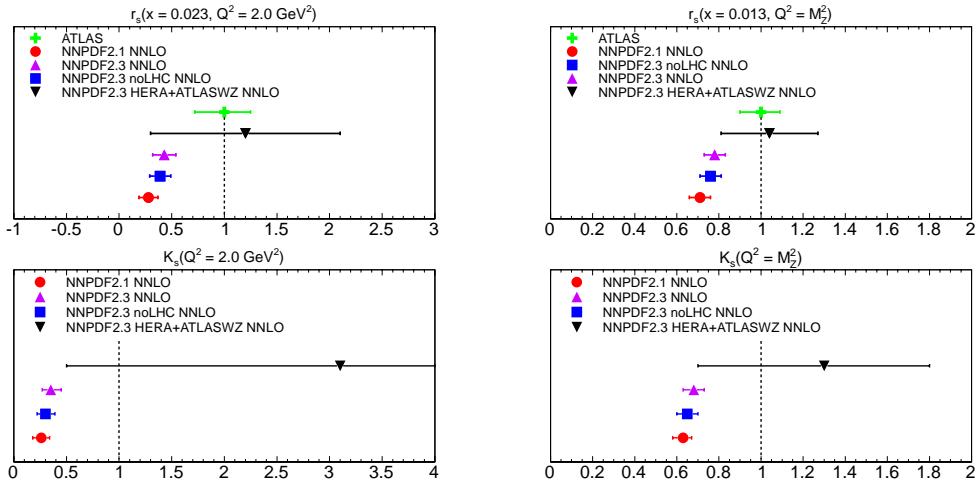
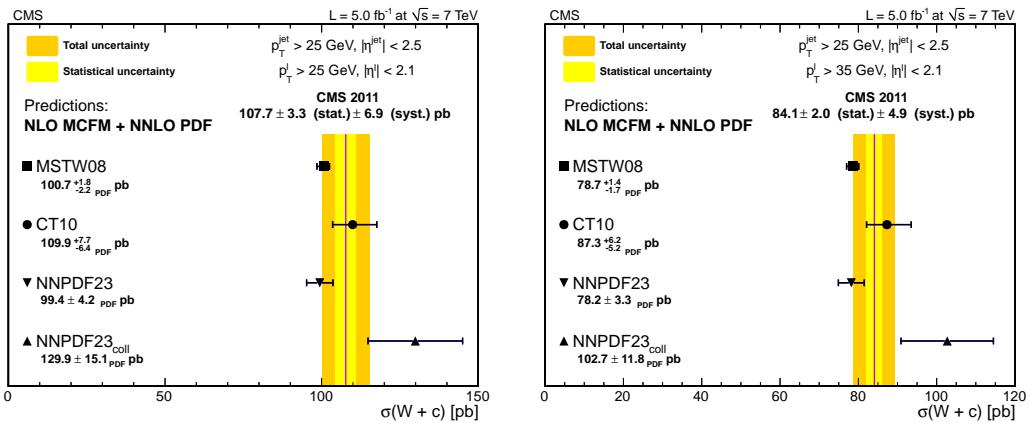


Figure 5.18: Results on the strangeness fraction of the proton from restricted dataset fits. Results are shown for  $r_s$  with the ATLAS kinematics (top plots) and for the integrated strangeness fraction  $K_s$  (below). Values are given for NNPDF2.1, NNPDF2.3, NNPDF2.3 noLHC and the NNPDF2.3 HERA + ATLAS  $W/Z$  dataset.

### 5.2.4 NNPDF2.3 phenomenology

We will begin the study of the phenomenological applications of the NNPDF2.3 set and comparisons to previous determinations, by comparing computations of the LHC measurements included in the 2.3 fit. In this way the improvements in precision available for LHC standard candle predictions can be assessed. We shall follow by looking at some typical total cross-sections of interest.

The impact of the ATLAS inclusive jet measurements upon NNPDF is made clear in Figure ?? where the data is compared to the predictions from the

Figure 5.19: CMS  $W + c$  production data, figure from [?].

NNPDF2.1 and NNPDF2.3 sets. Uncertainties on the predictions are reduced across all datapoints, and there is a general shift to lower values of the differential cross-section. Despite the shift downwards, the theory remains systematically above the experimental datapoints. However the dataset suffers from relatively large systematic uncertainties, within which both NNPDF2.1 and NNPDF2.3 are consistent as demonstrated by the excellent agreement at the level of  $\chi^2$  shown in Table ??.

In the electroweak sector, significant improvements are made across all observables included in the fit. Figure ?? compares the predictions of NNPDF2.1 and NNPDF2.3 to the experimental data for the LHC electroweak measurements, demonstrating the improved agreement between theory and data in the ATLAS and CMS results, while Figure ?? shows the same comparison for the LHCb data, demonstrating the improvements made in the very forward region measured by LHCb. The precise and consistent CMS data provide the clearest reduction of uncertainty of all the datasets, while the ATLAS and LHCb measurements suggest that the previous determination overestimated the electroweak cross-sections, leading to lower distributions with much improved agreement in the new fit.

Moving to inclusive cross-sections, predictions for total  $W^\pm$  and  $Z$  boson

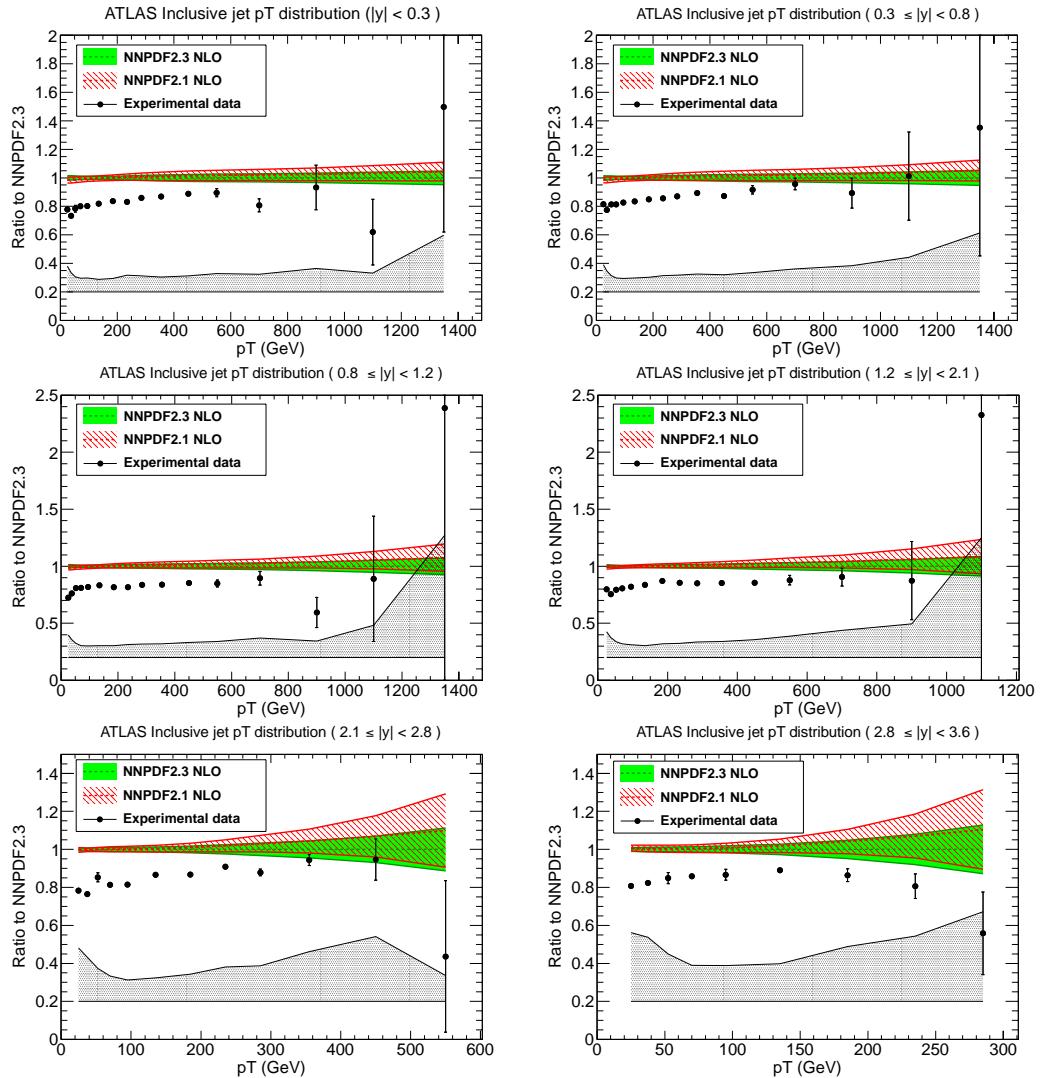


Figure 5.20: Predictions for the ATLAS 2010 inclusive jet data, using NNPDF2.1 (green) and NNPDF2.3 (red). The grey band at the bottom of each figure represents the systematic uncertainty in the data, while the error bars are the statistical error only. Predictions are given for all rapidity bins for the  $R = 0.4$  data as included in NNPDF2.3

production, along with the total  $t\bar{t}$  cross-section are shown in Figure ?? . Predictions for the electroweak observables were calculated using the **VRAP** [?] code, and for the top predictions, **top++** [?, ?] was used. Predictions are provided for the 7 TeV and 8 TeV LHC with  $\alpha_s(M_Z) = 0.119$ . In Figure ?? the total cross-section for Higgs production in gluon fusion is shown with the same settings, predictions provided by **iHixs** [?]. Results across the NNPDF2.1 and NNPDF2.3 sets demonstrate generally good consistency within their errors, with the NNPDF2.3 set providing the most precise predictions. The collider only determination is shown to be reasonably competitive when applied to the electroweak observables, where improved constraint is available from the LHC dataset. A similar pattern can be observed in the top and Higgs production observables, however errors remain systematically larger than for the global set.

The benefits of the NNPDF2.3 PDF set in phenomenological applications to LHC measurements are then clear, with the 2.3 set being the most precise and accurate determination in the NNPDF family.

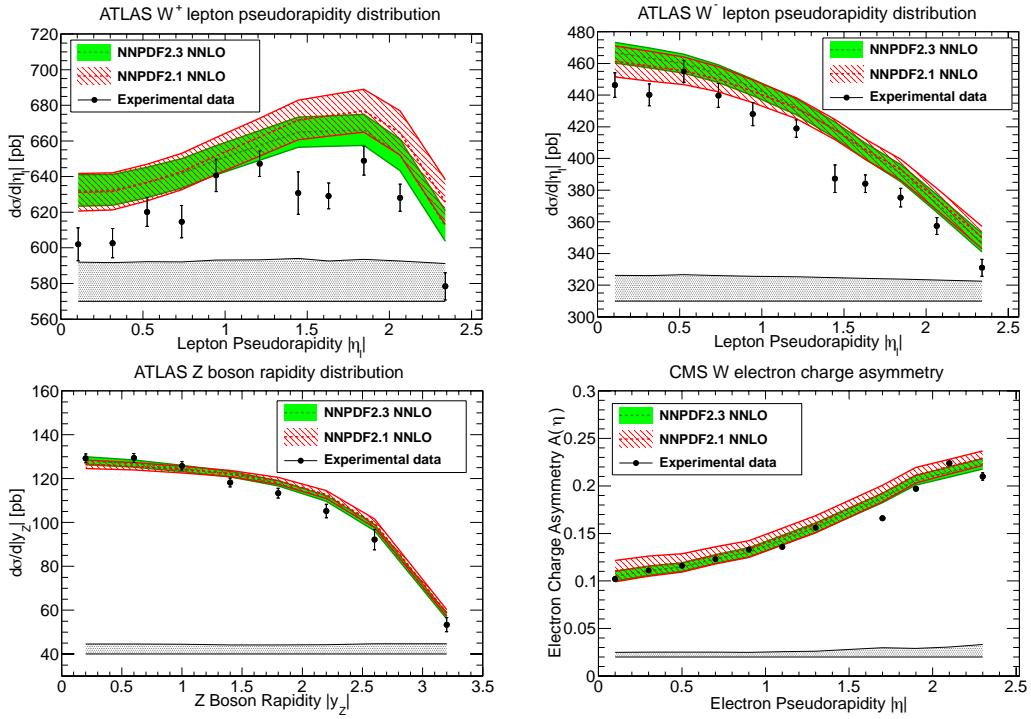


Figure 5.21: Predictions for the ATLAS 2010 electroweak vector boson production and CMS 2011  $W$  electron asymmetry data, using NNPDF2.1 (green) and NNPDF2.3 (red). The grey band at the bottom of each figure represents the systematic uncertainty in the data, while the error bars are the statistical error only.

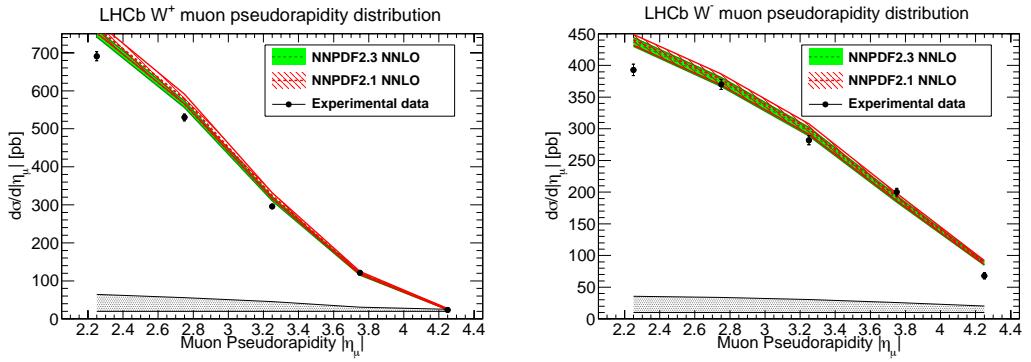


Figure 5.22: Predictions for the LHCb 2010  $W$  boson production data, using NNPDF2.1 (green) and NNPDF2.3 (red). The grey band at the bottom of each figure represents the systematic uncertainty in the data, while the error bars are the statistical error only.

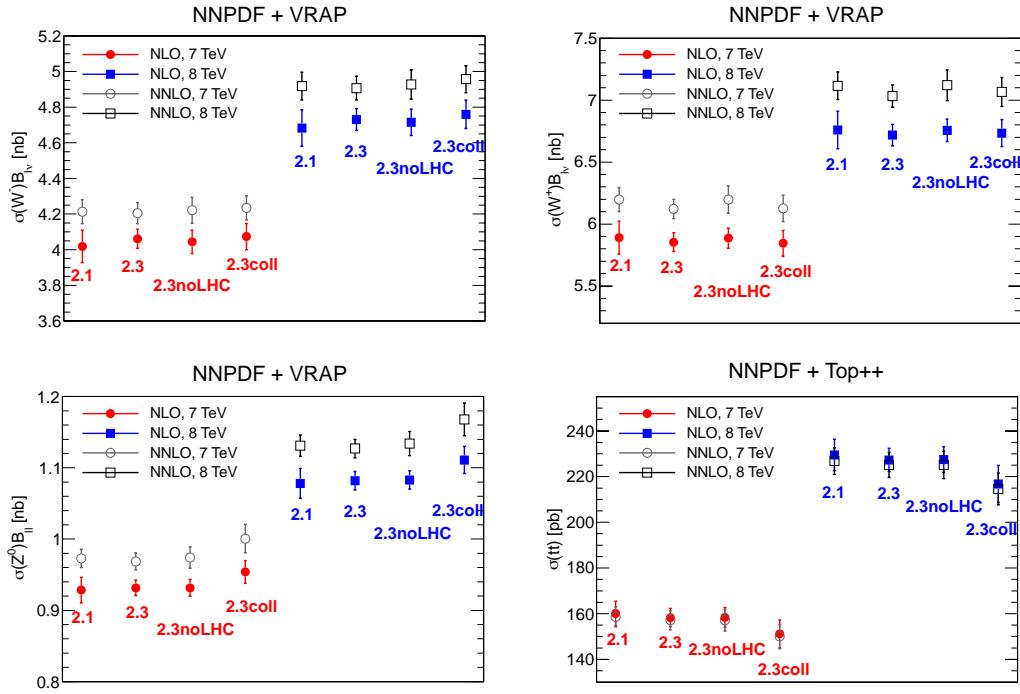


Figure 5.23: Predictions for total cross-sections at the 7 and 8 TeV LHC from NNPDF2.1 and the NNPDF2.3 series at NLO and NNLO. 7 TeV points are shown with circular markers, and 8 TeV with square markers. Theoretical predictions are given for the total  $W^+$  (top left),  $W^-$  (top right),  $Z$  (bottom left) and  $t\bar{t}$  (bottom right) production cross-sections.

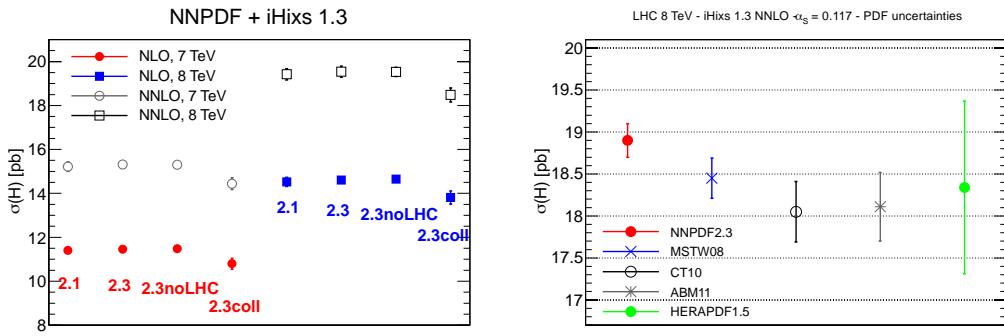


Figure 5.24: Total Higgs production cross-section in the gluon fusion channel at the 7 and 8 TeV LHC. Predictions are given at NLO and NNLO.

# Chapter 6

## Fitting in the light of LHC data

The addition of the first LHC datasets into an NNPDF fit allowed for important gains to be made in the precision of the resulting sets. However the potential dataset available for PDF determination from the LHC is increasing at a considerable rate, and datasets are being rapidly updated with more precise measurements. There is therefore still much more potential in the LHC to provide PDF constraint, especially in collider only fits.

With the ever enlarging dataset comes an important question: whether the fitting methodology applied to the pre-LHC dataset is still the best procedure for the extraction of precise parton densities in the LHC era. In order to accommodate the growing LHC dataset and to be able to efficiently explore methodological options, the toolchain used by the NNPDF collaboration had to be updated. The need for an updated fitting apparatus was recognised near the end of development of the NNPDF2.3 PDF set. Previous NNPDF sets were generated by a **FORTRAN** codebase which grew out of the earliest NNPDF determinations. Consequently the codebase suffered from a great deal of inflexibility with regard to the treatment of data. In particular, performing varying cuts and fits to reduced or special datasets were complicated procedures. Additionally as the fits to the pre-LHC dataset were considerably less computationally intensive the

core fitting apparatus was not designed with computational efficiency as the first priority, meaning that fits with the LHC dataset were rather sluggish. Beyond being a mere technicality, such slow fits actually meant that detailed studies of the methodology applied to an LHC dataset were prohibitively expensive in computer time.

With these issues in mind, the `nnpdf++` project was initiated, whereby the full NNPDF toolchain has been implemented from scratch in C++. The core of the project was built around the efficient FK method described previously, allowing for a much clearer separation in code between theoretical predictions and experimental data along with a much greater efficiency in the convolution. The FK products themselves are accelerated via explicit use of SIMD vectorisation, and OpenMP [?] provides multiprocessor options.

The framework was designed to be as modular as possible, to allow for the simple and safe modification of sections of the NNPDF methodology without requiring major modifications to the remaining codebase. The re-implementation of the whole NNPDF toolchain also provided an extremely thorough cross-check of the two implementations, and allowed for the step-by-step evaluation of several methodological elements. The results of this re-evaluation and investigation of alternative procedures shall be described in this chapter along with the consequences for future determinations.

## 6.1 Closure testing

The central element in the methodological review conducted with the `nnpdf++` code after NNPDF2.3 is the closure testing procedure.

In a closure test, a PDF fitter takes their tools and applies them to a set of pseudo-experimental data generated from a known prior parton distribution set. Provided that the theory used to generate the pseudodata is identical to that

used in the fitting procedure, the results of the fit should reproduce the generating function to within the estimation of PDF error. The test is an extremely sensitive check of a fitting procedure, in that it tests the ability of a methodology to resolve the underlying law when said law is known exactly. The method can also be used to study the effect of data inconsistencies by artificially modifying data uncertainties as is examined in Ref. [?], however here we shall restrict ourselves to examining the quality of reproduction of the underlying law.

Closure tests in the NNPDF methodology can be performed in a number of ways. One possible method is a direct fit to theoretical predictions generated from a known distribution, in this way the pseudo-dataset is free from the statistical noise that would be present in experimental data.  $N_{\text{rep}}$  PDF replicas are then fitted to the theory predictions, without performing the generation of a Monte Carlo artificial data sample. In this type of fit one aims to reproduce as well as possible the generating function at the end of the fit. As no statistical noise is inserted at any point the final fit quality should approach  $\chi^2 = 0$ , we shall therefore denote such a fit a *level zero* closure test.

Alternatively one may perform a fit where statistical noise is introduced to the pseudo-dataset according to the experimental uncertainty present in the real dataset. This can be done in two ways; either the noise is introduced directly to the pseudo-data itself whereby all Monte Carlo replicas fit to the same noisy sample, or noise is introduced on a replica-by-replica basis as in the normal Monte Carlo procedure. These types of fit we denote *level one* closure tests.

Finally one can introduce two levels of noise to the data. The first; applied directly to the pseudo-data, simulates the experimental noise in the distributions. The second level is introduced through the normal Monte Carlo generation of artificial data replicas. This is denoted a *level two* closure test and is the closest to a full fledged PDF fit. The main exception here being the lack of any inconsistency between datasets, as they have all been generated from the

same initial distribution. In the case of a level two fit the PDF fitter wishes to reproduce the underlying law to within their quoted PDF uncertainties, the exact reproduction available at level zero is now unavailable due to the introduced pseudo-experimental noise. The level two fit is therefore the most stringent test of a fitting procedure in that it tests the central claim of a fitting group; that the underlying law should lie within the quoted PDF uncertainty band at the quoted confidence level. The settings used in the different closure tests are summarised in Table ???. As a direct comparison of some example pseudodata, Figure ?? shows example data at closure test levels zero, one and two.

C. Level	Exp. Noise	Art. Data
0	X	X
1a	✓	X
1b	X	✓
2	✓	✓

Table 6.1: Levels available in a closure test fit (C. Level), Exp. Noise corresponds to simulating experimental noise in the pseudodata sample. Art. Data refers to the generation of artificial data replicas in the Monte Carlo uncertainty procedure.

The new structure present in the `nnpdf++` code, particularly the modular treatment of experimental data and theoretical predictions, allows for the straightforward use of predictions in the place of experimental data while keeping the experimental covariance matrices intact. The closure testing method has therefore been extensively applied to the development of the NNPDF methodology, with the procedure used for the NNPDF3.0 determination being guided largely by results from closure testing. Here we shall outline some general results, before demonstrating the application of the procedure to methodological development in the subsequent sections.

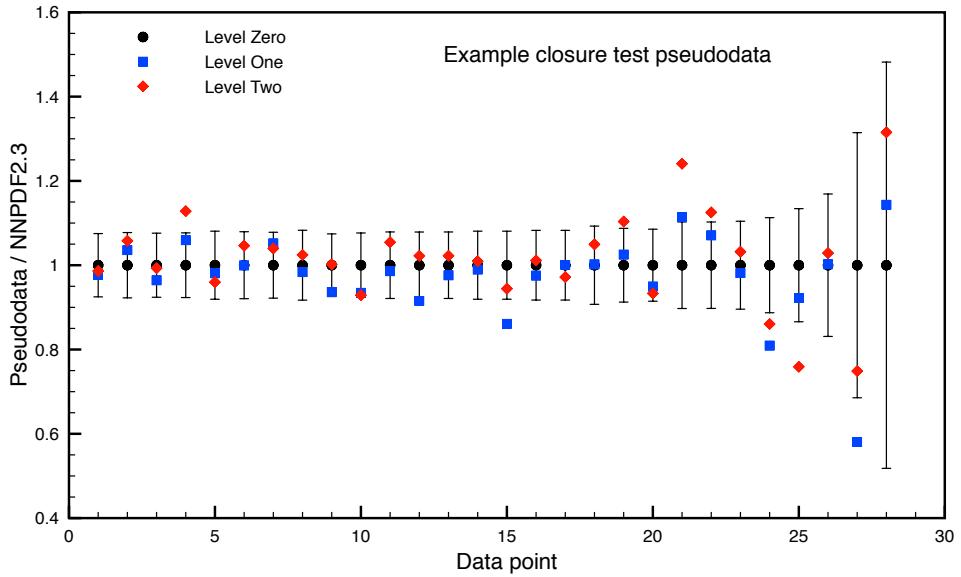


Figure 6.1: Examples of pseudodata used in a closure test for all three levels. The black circles show the level zero pseudodata, and the experimental error bars. The blue squares show the pseudodata after experimental noise has been simulated (level one) and the red diamonds after both statistical noise simulation and Monte Carlo replica generation (level two). All points are normalised to the generating PDF set (NNPDF2.3).

### Early closure tests

The earliest NNPDF closure tests were conducted to assess the usefulness of the procedure, and performed with the full NNPDF2.3 procedure. As an initial test, a fit was performed to the toy PDF parametrisation as used in the Les Houches evolution benchmarks [?], a parametrisation based upon the CTEQ5M determination [?]. In this set, the initial state distributions are given as

$$\begin{aligned}
xu_v(x, \mu_{f,0}^2) &= 5.107200 x^{0.8} (1-x)^3, \\
xd_v(x, \mu_{f,0}^2) &= 3.064320 x^{0.8} (1-x)^4, \\
xg(x, \mu_{f,0}^2) &= 1.700000 x^{-0.1} (1-x)^5, \\
x\bar{d}(x, \mu_{f,0}^2) &= .1939875 x^{-0.1} (1-x)^6, \\
x\bar{u}(x, \mu_{f,0}^2) &= (1-x) x\bar{d}(x, \mu_{f,0}^2), \\
xs(x, \mu_{f,0}^2) &= x\bar{s}(x, \mu_{f,0}^2) = 0.2 x(\bar{u} + \bar{d})(x, \mu_{f,0}^2),
\end{aligned} \tag{6.1}$$

where  $u_v$  and  $d_v$  refer to the up and down valence distributions respectively. Predictions for the NNPDF2.3 dataset were made according to these distributions, and used in the place of experimental data. Experimental noise was simulated in the pseudodata by application of the same procedure used to provide artificial data replicas. The full NNPDF2.3 procedure including Monte Carlo artificial replicas was then applied to the dataset, the resulting PDF set therefore being a level two type closure test where the generating PDF set should be recovered by the fit within the estimated uncertainties.

Figure ?? displays the results of the level two closure test fit with the Les Houches toy PDFs used as a generating function. The result demonstrates impressive agreement, with the NNPDF2.3 methodology able to accommodate the predictions of the Les Houches toy generating function despite it deviating significantly from the standard NNPDF2.3 result. For all four PDF combinations shown, the results of the closure test maintain distances of less than one standard deviation to the generating function across a wide kinematic range. Of additional interest are the strange distributions, relatively poorly constrained by the data included in the pseudo-dataset. The strange valence in particular is set to zero in the Les Houches toy. Figure ?? shows the results from the closure test for both the total strangeness and strange valence distributions, the NNPDF methodology

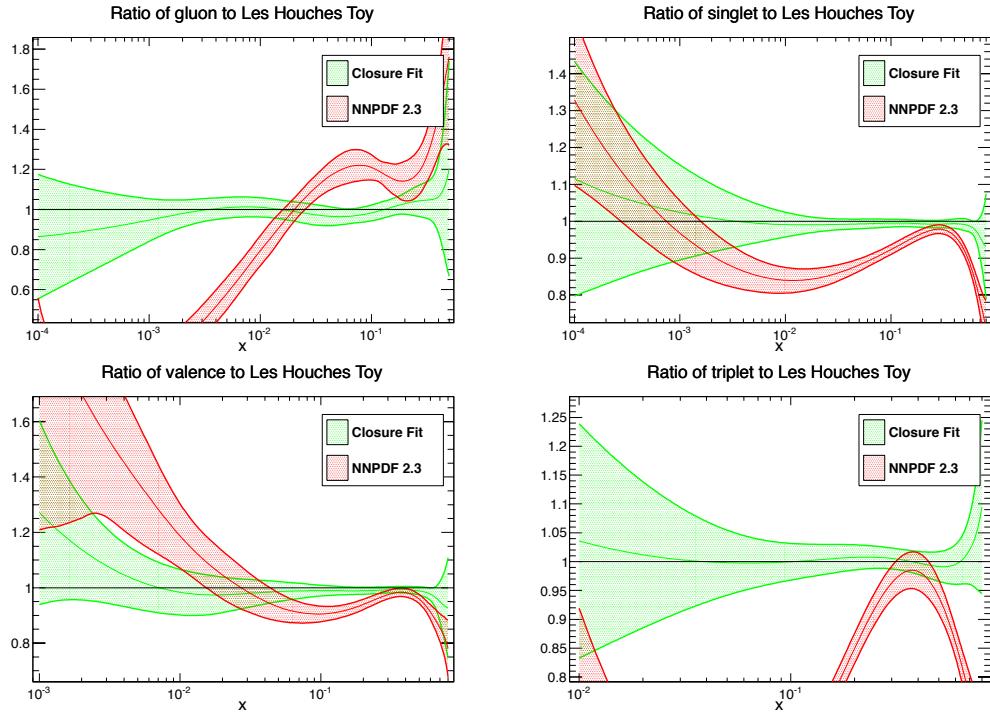


Figure 6.2: PDFs obtained through a closure test fit with Les Houches toy PDFs as a generating function, displayed as a ratio to the generating function. Shown are the distributions for the gluon, singlet, valence and triplet PDFs. In green are the results obtained through the closure test, and the red curves show the standard NNPDF2.3 result.

is able to clearly reproduce the underlying law within uncertainties in both cases, and is able to comfortably resolve a zero strange valence contribution.

The results are particularly impressive considering that this is a test of a methodology that has not been previously verified by closure test. The example case of a pseudo-dataset generated according to the Les Houches toy PDF is however a rather simplified case, and methodological refinements can be made by examining closure tests with greater structure in the generating function.

A good level of agreement can also be found at the level of the  $\chi^2$  to both the pseudodata sample, and the real experimental data. In Figure ?? we compare the fit quality of a closure test and its generating PDF dataset by dataset by presenting the  $\chi^2$  to each measurement from both the closure test result and the

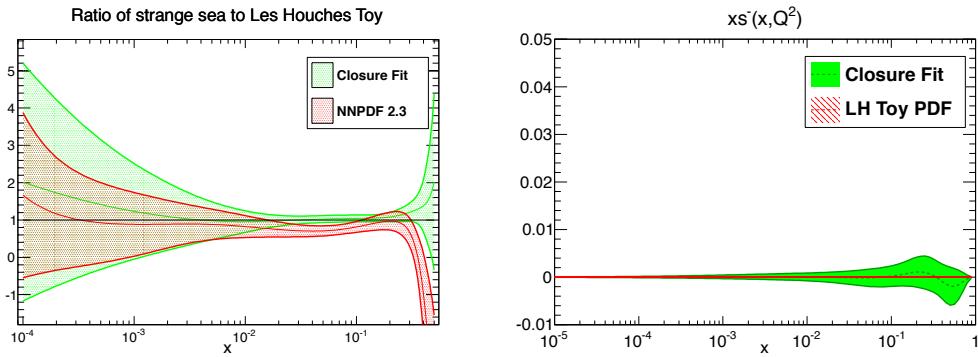


Figure 6.3: Strange sea (left) and valence (right) PDFs obtained through a closure test fit with Les Houches toy PDFs as a generating function. The strange sea is presented as a ratio to the LH toy PDF, and the strange valence is presented directly as the PDF, with the (zero) LH toy line shown.

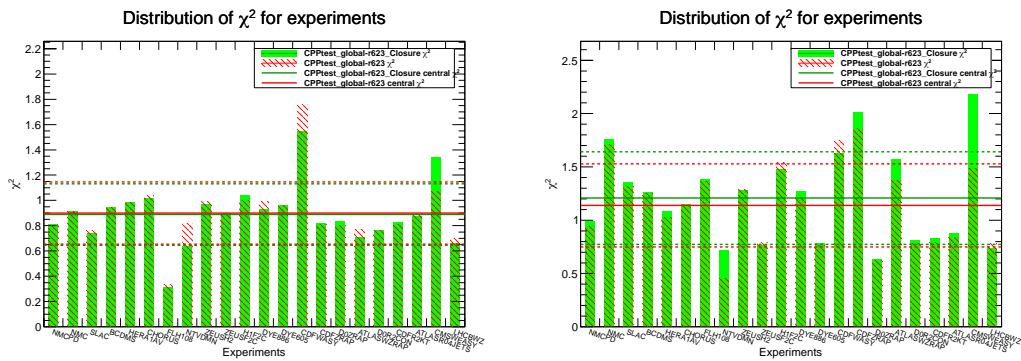


Figure 6.4: Example  $\chi^2$  values to the pseudo- (left) and experimental- (right) datasets of a closure fit and the generating PDF from early `nnpdf++` test fits. The red bars show the fit quality of the generating PDF set while the green bars demonstrate the  $\chi^2$  for the closure test set. The horizontal lines indicate the average and  $1\sigma$  of the fit qualities in their associated colours.

generating PDF. In this case the generating function has considerably greater complexity, being an early `nnpdf++` test fit with most of the NNPDF methodology in place. While agreement is generally very good, especially on the level of total  $\chi^2$ ; we begin to see some elements of discrepancy in datasets sensitive to flavour separation and strangeness such as the NuTeV dataset and electroweak vector boson production data. Such discrepancies can help in pinpointing areas where further development is needed.

## 6.2 Preprocessing

Early closure tests performed with the NNPDF2.3 methodology showed generally very good agreement between the produced PDFs and the underlying functions used to generate the pseudo-dataset. However some PDF combinations demonstrated rather poorer agreement than others, particularly distributions sensitive to flavour separation. Such disagreements became more apparent when considering closure tests to underlying functions with more structure than available in the Les Houches toy set. The disagreements were found to originate in the choice of the preprocessing exponents used in the definition of the NNPDF parametrisation. Recalling Eqn. ??, the structure of the basic NNPDF parametrisation follows

$$f(x) \propto x^{-\alpha} (1-x)^{\beta} \text{NN}(x), \quad (6.2)$$

where NN represents the neural network itself, and  $\alpha$  and  $\beta$  are the preprocessing exponents randomised on a replica-by-replica basis at the start of a fit. The range in which the exponents were randomised has been fixed in the fits up to and including NNPDF2.3, set to a span large enough such that the dependence of the results upon the choice of range was minimised. In such a way the preprocessing was considered to provide a backbone for the neural-network fit and, aside from improving fitting efficiency, to have a minimal impact upon the results.

To study the effect of different preprocessing ranges we can look at estimators for the *effective* asymptotic exponents,

$$\alpha_{\text{eff}} = -\frac{\log(|f(x)|)}{\log(x)}, \quad \beta_{\text{eff}} = \frac{\log(|f(x)|)}{\log(1-x)}, \quad (6.3)$$

such that in the limits of  $x \rightarrow 0, 1$  the exponents  $\alpha, \beta$  are recovered. By examining these effective exponents in the high- and low- $x$  regions, we can ascertain if there

is a data preference for a different preprocessing range than was used in a fit, and if the preprocessing range used was too restrictive.

In Figure ?? an example preprocessing analysis is shown for a closure test based upon an MSTW08 underlying law at NLO. The sea asymmetry  $\bar{u} - \bar{d}$  is shown for two choices of preprocessing range, the NNPDF2.3 standard and a range modified to better accommodate the data preference visible in the effective exponents. From the figure we can see that the choice of exponent randomisation range has a significant effect on the resulting distributions, and that the effective exponents can show a clear data preference for a different range. In Figure ?? we can see the same analysis applied to the triplet PDF where similar conclusions may be drawn.

These analyses demonstrate that the sensitivity to the preprocessing exponent randomisation ranges is somewhat larger than suspected previously, and needs to be studied in detail in order to avoid minimisation difficulties in a fit where the preprocessing ranges are ill-suited to the dataset. Furthermore, the uncertainty bulges visible in both the triplet and sea asymmetry distributions in Figures ?? and ?? are generated by the preprocessing suppressing genuine data uncertainty in the asymptotic regions. These problems may be alleviated by lifting the requirement that such distributions should be preprocessed to zero at low- $x$ , and implementing a procedure for the iterative and data-driven determination of preprocessing exponents.

To improve the minimisation performance, hampered by ill-suited preprocessing, NNPDF fits have now adopted the following iterative procedure for the determination of both high and low- $x$  randomisation ranges:

- **Singlet and gluon PDFs**

Exponent randomisation ranges are set to be twice the  $1\sigma$  interval of the previous iteration's effective exponents at the asymptotic points.

- **Nonsinglet PDF combinations**

The low- $x$  randomisation interval is set to be the maximal extent of two effective exponent ranges; twice the  $1\sigma$  interval at the asymptotic point and twice the  $1\sigma$  interval at the point  $x = 1 \times 10^{-3}$ . The high- $x$  interval is set identically as with the singlet and gluon.

In such a way convergence of the randomisation interval can be established typically in two or three fit iterations, and the preprocessing exponents are obtained from the preference of the experimental dataset. As an example of a fit generated from such an iterative procedure consider Figure ?? which demonstrates the preprocessing analysis for the  $\Delta_s$  and Triplet distributions resulting from the new procedure. In comparison to Figure ?? where the old settings are used, the low- $x$  preprocessing ranges have relaxed considerably and are no longer constrained by the chosen exponent range but driven by the experimental data. Furthermore the agreement with the underlying law is noticeably improved over the previous result shown in Figure ??.

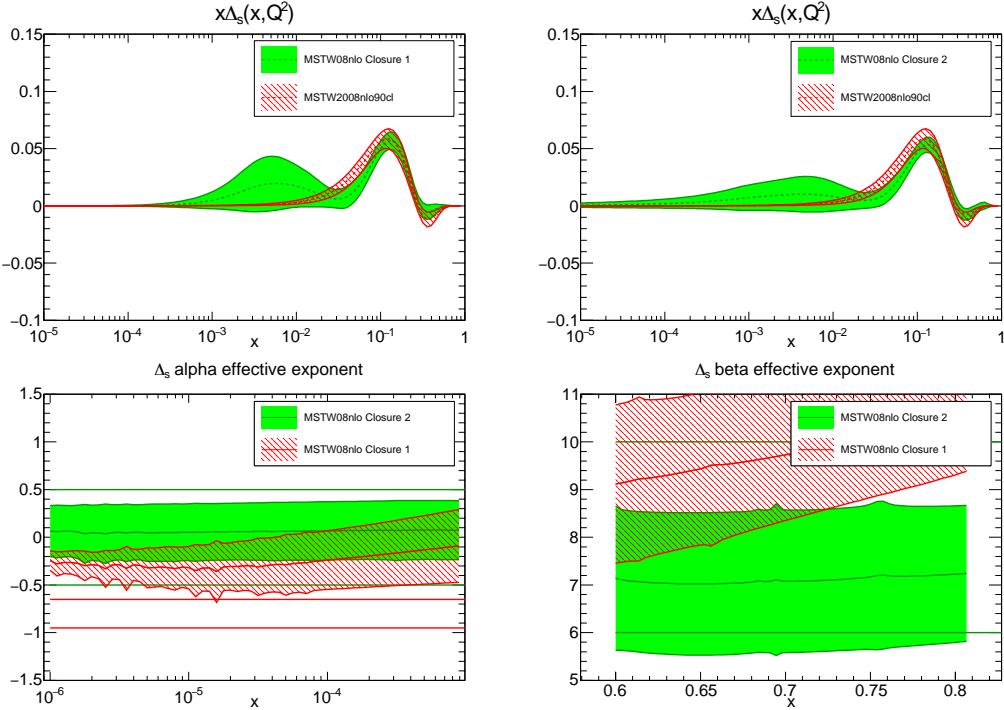


Figure 6.5: Demonstration of the impact made by changes in preprocessing to the sea asymmetry PDF in a closure test fit to an MSTW08 NLO underlying law. The top two figures demonstrate the results for the  $\Delta_s$  distribution for two choices of preprocessing ranges, with the left figure using NNPDF2.3 standard preprocessing. In both cases, the red curve shows the underlying law used in the Closure test. The right figure demonstrates slightly improved agreement, particularly at low- $x$ . The lower figures show the low and high  $x$  effective exponent plots for the two ranges. The solid horizontal lines delineate the regions in which the preprocessing exponents were initialised, and the bands show the  $1\sigma$  contours of the effective exponents.

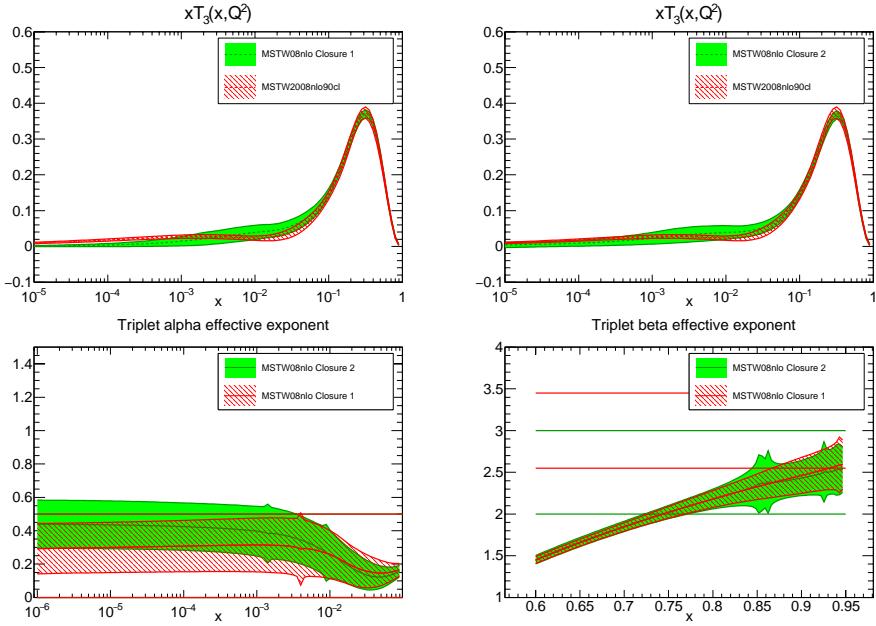


Figure 6.6: A further preprocessing analysis as in Figure ??, performed upon the Triplet PDF combination for the same two closure test fits.

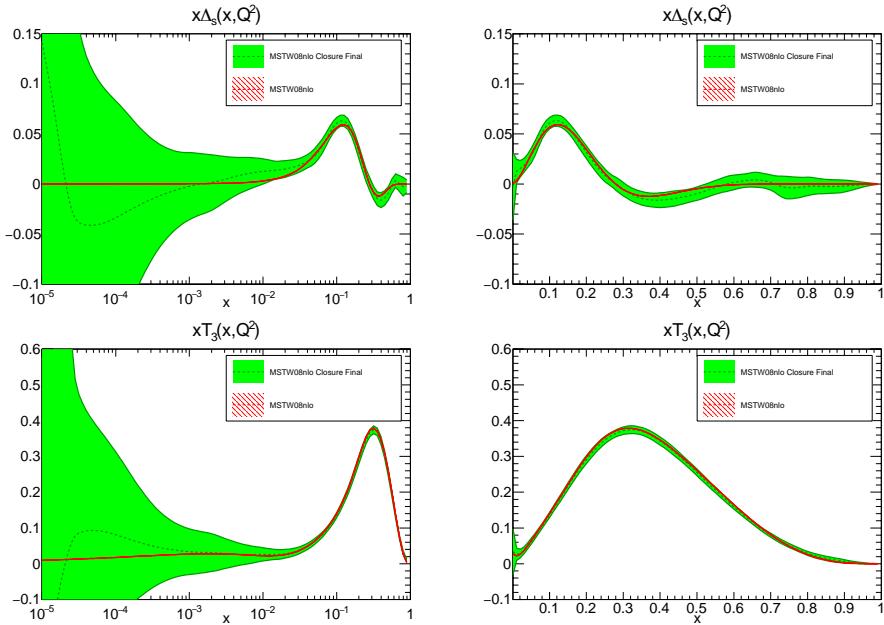


Figure 6.7: Impact of improved preprocessing range selection in the sea asymmetry and triplet PDFs. The top figures demonstrate the  $\Delta_s$  PDF obtained via a closure test to MSTW08 using the improved preprocessing procedure in green, with the underlying law shown in red. The figures below show the equivalent plots for the triplet PDF with the improved preprocessing ranges.

### 6.2.1 Strange valence preprocessing

A special case when considering the preprocessing of the neural networks is that of the strange valence distribution. As specified in Equation ??, the strange valence PDF in the NNPDF2.3 determination had an auxiliary term to encourage the PDF to perform its required sign change in the valence region. Such an additional term has been previously needed due to the lack of specific data constraints upon the strange valence distribution before the LHC, introducing a bias, albeit a physically motivated one. Additionally the auxiliary term provides a mechanism by which the strange valence sum rule may be imposed. In the NNPDF3.0 determination and beyond this auxiliary term has been removed given the enlarged dataset and it's improved sensitivity to the strange PDF.

Figure ?? demonstrates the effect of the removal of the strange auxiliary term upon a closure test fit to the MSTW08 set. While the NNPDF2.3 methodology closure fit struggles to accommodate the MSTW08 strange valence distribution, the updated methodology is able to reproduce the underlying law well, within enlarged uncertainties.

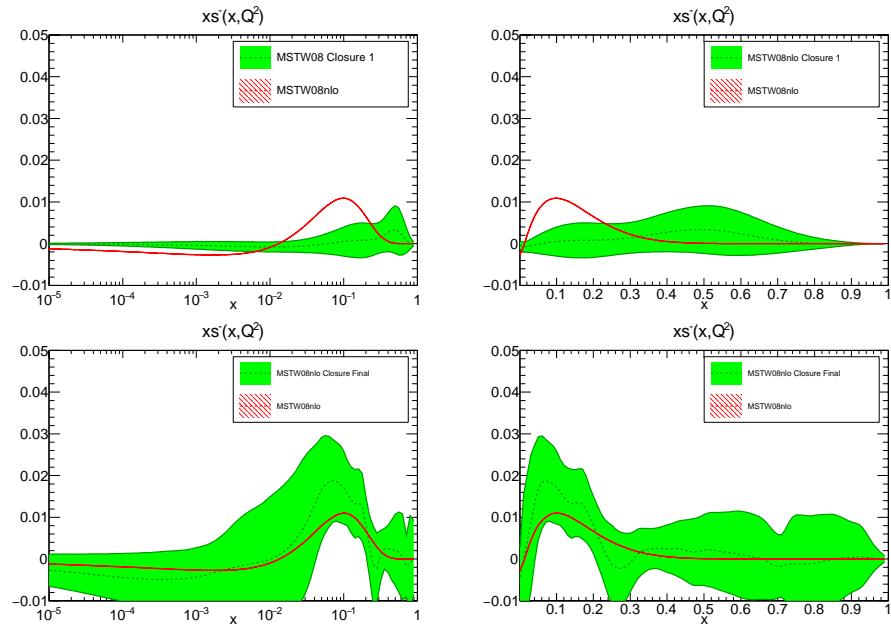


Figure 6.8: Impact of more flexible treatment of strange valence PDF in fits post NNPDF2.3. The top two figures show a comparison of a closure test performed with the NNPDF2.3 preprocessing, and the figures below show the results using the more flexible parametrisation.

### 6.3 PDF parametrisation

The choice of PDF parameterisation and basis used in the fitting procedure has been reassessed with the help of the closure test procedure. In particular, a modification to the choice of fitting basis has been made necessary by the removal of the strange valence sum rule enforcing auxiliary term in the strange valence parametrisation. The most direct choice of fitting basis is the same basis as is used in PDF evolution, and therefore the basis required for PDFs in the FK product. In this basis, the required quantum number sum rules may be applied as normalisations to the total valence,  $V_3$  and  $V_8$  distributions,

$$\begin{aligned} V(x, Q_0^2) &= N_V (u^- + d^- + s^-)(x, Q_0^2), \\ V_3(x, Q_0^2) &= N_{V3} (u^- - d^-)(x, Q_0^2), \\ V_8(x, Q_0^2) &= N_{V8} (u^- + d^- - 2s^-)(x, Q_0^2), \end{aligned} \quad (6.4)$$

where the normalisations  $N$  are set such that

$$\int_0^1 dx V(x, Q_0^2) = 3, \quad (6.5)$$

$$\int_0^1 dx V_3(x, Q_0^2) = 1, \quad (6.6)$$

$$\int_0^1 dx V_8(x, Q_0^2) = 3. \quad (6.7)$$

In such a way, the total valence quantum number is fixed, along with the up, down and strange valence quantum numbers. The evolution basis also has the advantage of being particularly efficient, not requiring any transformation before combination with FK tables to calculate physical observables. We have shown, based upon closure test results, that the fit results show a good degree of stability under such a change in parametrisation basis. While the previous strategy was designed to construct PDF combinations with specific data constraints,

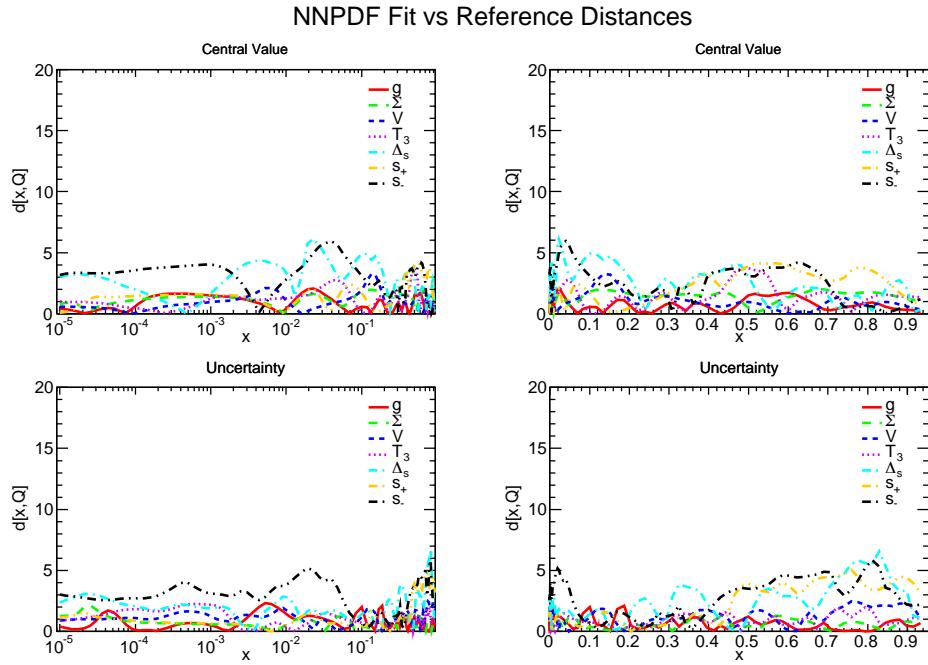


Figure 6.9: Distance comparison of two closure test fits with differing parametrisation bases. Distances are defined through the measure in Appendix ??, whereby a distance of 10 corresponds to  $1\sigma$ .

the flexibility of the fit means that the results do not suffer when moving away from such a basis. Figure ?? shows the statistical distance between a fit with the full evolution basis with a fit based upon the standard NNPDF2.3 parametrisation basis. As expected, any differences are isolated to those PDFs whose parametrisation (and therefore preprocessing) has substantially changed e.g.  $\Delta_S$  and the strange PDFs. Even in these PDFs the differences are typically less than half a standard deviation.

## 6.4 Minimisation and stopping

In addition to examining areas where the choice of parametrisation may lead to some degree of bias, the closure test procedure is particularly useful for assessing the efficacy of a fitting methodology. Furthermore, the substantial gains in

computational efficiency made in the transition to the `nnpdf++` code mean that far more aggressive genetic minimisation strategies may be implemented.

The entirety of the NNPDF minimisation procedure has therefore been re-examined to ensure that it is the most effective methodology in the light of additional constraints coming from the LHC. Here we shall summarise some of the major modifications made since the NNPDF2.3 determination.

#### 6.4.1 Target weighted training

Target Weighted Training (TWT) was a central feature of previous NNPDF determinations. TWT was developed in early NNPDF fits as a method of obtaining a balanced training across datasets, solving a problem with early neural network fits whereby some smaller datasets were largely ignored by the minimisation in favour of larger, more constraining sets. This typically led to a very uneven fit quality profile over the complete experimental dataset. The TWT procedure solved this problem by introducing a training epoch at the beginning of a fit where each dataset had a target  $\chi^2$ . In the event where a fit iteration reached a  $\chi^2$  value higher than the target, a large weight in fit quality was applied to that dataset in order to bring its fit quality down.

While ensuring a relatively even training profile, the TWT procedure had a number of difficulties. The most important being the restriction of the early fit to a  $\chi^2$  fit quality measure applied on a dataset-by-dataset basis, ignoring experimental uncertainty cross-correlations such as luminosity uncertainties, between datasets. Furthermore the TWT procedure introduced a considerable amount of complexity in the fitting procedure. With this in mind, real data fits with target weighted training were compared to fits without in the `nnpdf++` framework with the large experimental dataset of NNPDF2.3 and updated genetic algorithm parameters. Figure ?? compares the dataset-by-dataset fit quality of two such example fits. With these fits we can see clearly that with a larger dataset

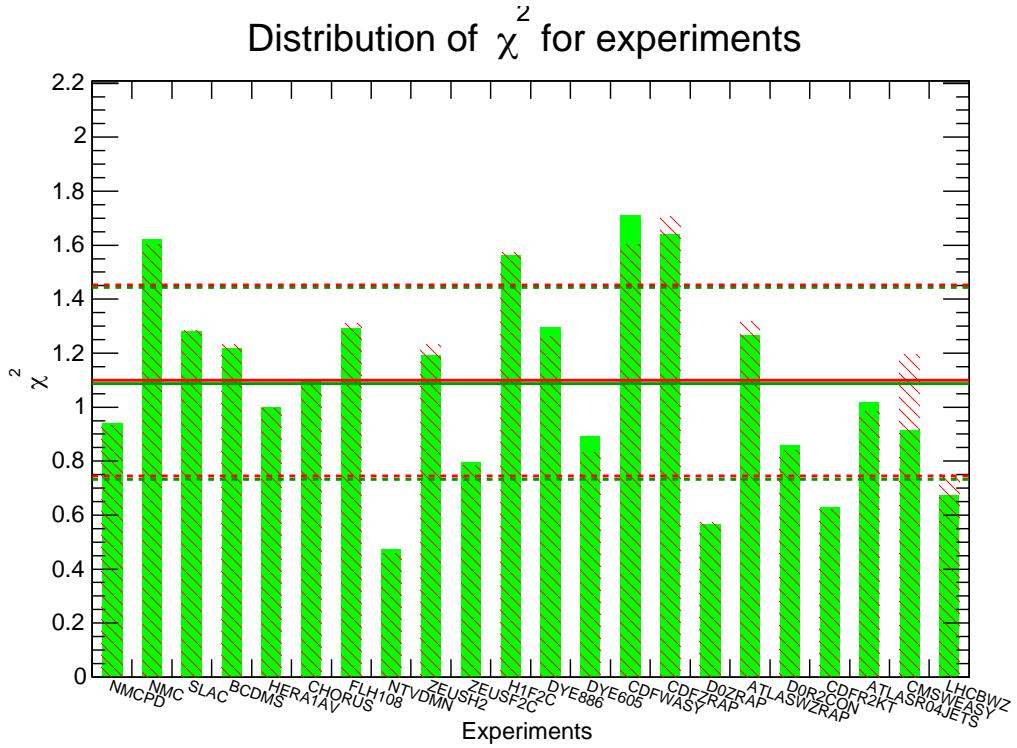


Figure 6.10: Comparison of  $\chi^2$  by dataset between real data fits with (green bars) and without (red bars) Target Weighted Training.

and more efficient GA procedure, no large training imbalance can be seen in the fits even without the TWT procedure applied. Future NNPDF fits will therefore be performed without target weights, allowing for the consistent application of the experimental correlations across datasets throughout the fitting procedure.

#### 6.4.2 Genetic algorithm

A number of changes have been made to the GA procedure used in NNPDF fits in order to improve fitting efficiency and provide more precise PDF determinations. In the analysis of the efficacy of a GA, level zero closure tests are particularly helpful in that they directly test the ability of a minimisation procedure to reproduce a given function precisely. In these fits the closure test fit should be able to effectively draw a line between datapoints, leading to an ideal  $\chi^2$  of

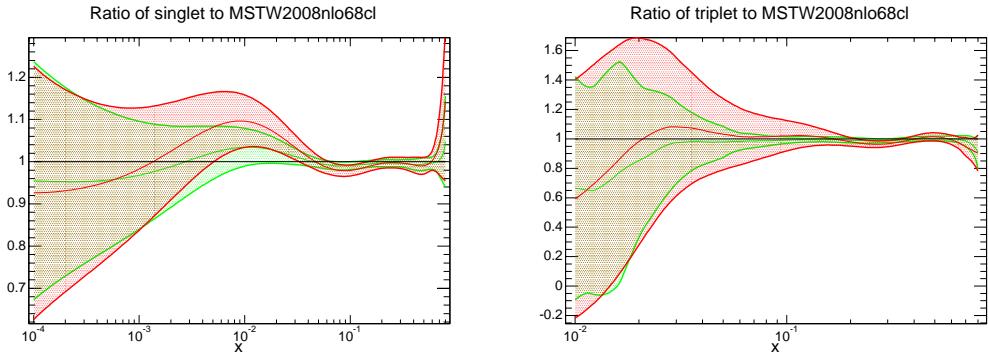


Figure 6.11: Comparison of a conventional NNPDF GA fit (red bands) with a Nodal GA fit (green bands) in a closure test to MSTW08. PDFs are given as a ratio to the generating PDF set for the singlet (left) and triplet (right) distributions.

zero to the pseudo-data. A number of modifications to the procedure have been tested, the most effective of which is the implementation of *Nodal* mutations in the GA [?]. In previous versions of the NNPDF GA, mutations were performed upon individual parameters of each neural network with no consideration as to their position in the network.

The concept of nodal mutations introduces the strategy of mutating all parameters associated with a particular neural network *node* at once. In this procedure a node of the network is chosen at random, then all of its associated weights connected to the earlier layer are mutated along with its threshold parameter. Doing so yields a much more effective genetic algorithm as demonstrated in the comparison in Figure ??, where a standard GA is compared to a nodal mutation GA in their reproduction of the MSTW underlying law. The nodal GA is able to better resolve the underlying law, and to a greater precision. The comparison in Figure ?? is corroborated by the  $\chi^2$  values of the two fits to the perfect pseudo-data in the level zero fit. The standard GA fit shown in the figure obtained a final  $\chi^2$  of 0.0279 compared to 0.0043 for the nodal GA. The nodal GA strategy has therefore been adopted for future NNPDF determinations.

### 6.4.3 Dynamical stopping

The cross-validation dynamical stopping procedure utilised in previous FORTRAN based NNPDF fits was triggered by a slope-detection algorithm applied to the fit quality profiles of each replica to the validation dataset. While providing a reasonable stopping criteria and preventing excessive overfitting, the relative balance between the degree of under- and over-learning was governed by the parameters of the slope-detection algorithm. Such sensitivity to the stopping parameters meant that a re-tune was often necessary upon large modifications to the dataset or minimisation algorithm.

The modular nature of the stopping criteria implemented in the `nnpdf++` framework means that alternative stopping procedures may be quickly and safely implemented to investigate their impact. One such stopping criterion that has demonstrated greater stability than the previous slope-detection based procedure is that of *look-back* cross-validation.

In this procedure all replicas are run for the maximum number of generations  $N_{\text{gen}}^{\max}$ , all the while storing the GA generation that best described the validation dataset. At the end of the fit, the GA generation that minimised the  $\chi^2$  to the validation set is selected as the best-fit stopping point, and that replica is used as a member of the Monte Carlo ensemble. This method yields an extremely clean stopping criterion, having no tuneable parameters aside from the maximum number of generations, and offers a very faithful implementation of the cross-validation method. Furthermore, the look-back procedure is not practically more time-consuming to implement despite running each replica to the maximum number of generations, as even in the previous dynamical stopping procedure the time taken to run a fit is typically given by the time taken by the slowest replica. In Figure ?? the fit quality profile for a single PDF replica can be seen for the training and validation sets alongside the look-back stopping point. In this case, the look-back method can clearly discern an overlearning signal, as the fit quality

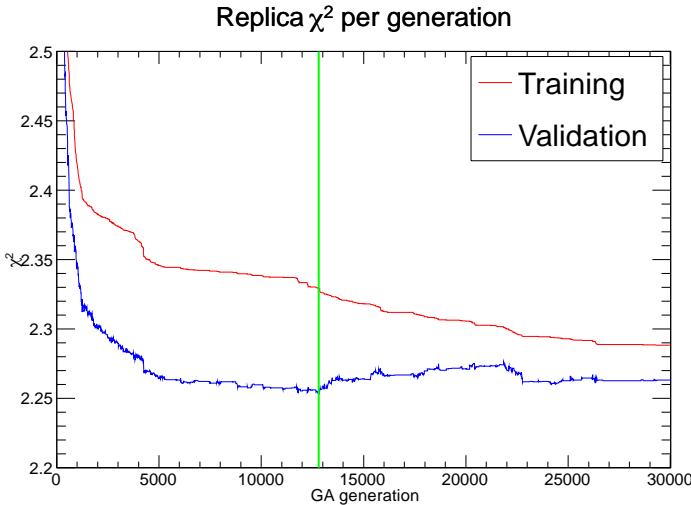


Figure 6.12: Fit quality profiles for the training and validation sets in look-back cross validation. The red curve shows the fit quality to the training set, and the blue curve to the validation set as the number of fit generations goes on. The green line indicates the stopping point selected by the look-back criterion, generation 12813 having the minimum validation  $\chi^2$ .

to the validation set worsens while the training set  $\chi^2$  improves.

In Figure ?? we compare the results for the singlet and gluon PDFs in the case of a look-back fit with  $N_{\text{gen}}^{\max} = 30,000$  generations, and a fit with the NNPDF2.3 standard dynamical stopping. In both instances, the fit performed was a level two closure test using MSTW2008 as the underlying law. While differences are small the look-back fit demonstrates slightly smaller uncertainties, implying a marginal underlearning present in the NNPDF2.3 dynamical stopping procedure. The fits yield essentially equivalent results, although the optimal point determined in the look-back method is typically somewhat later than in the dynamical stopping as can be seen in the comparison of training length histograms in Figure ???. In this figure it is clear also that several PDF replicas in the look-back method stop close to the maximum number of generations available, implying that no significant overlearning can be resolved in their cases over the given GA interval.

In order to examine the effect of increasing the length of the look-back period, we compare the 30,000 generation fit to an extended 60,000 generation fit in

Figure ?? where we use the PDF distance definition in Appendix ???. Distances of effectively zero throughout the PDF combinations and  $x$ -range mean that no change is observed between the two fits, demonstrating the stability of the method once a sufficiently large look-back length is used. The look-back cross-validation method as discussed here will therefore be implemented as the default stopping criterion for the NNPDF3.0 family of fits.

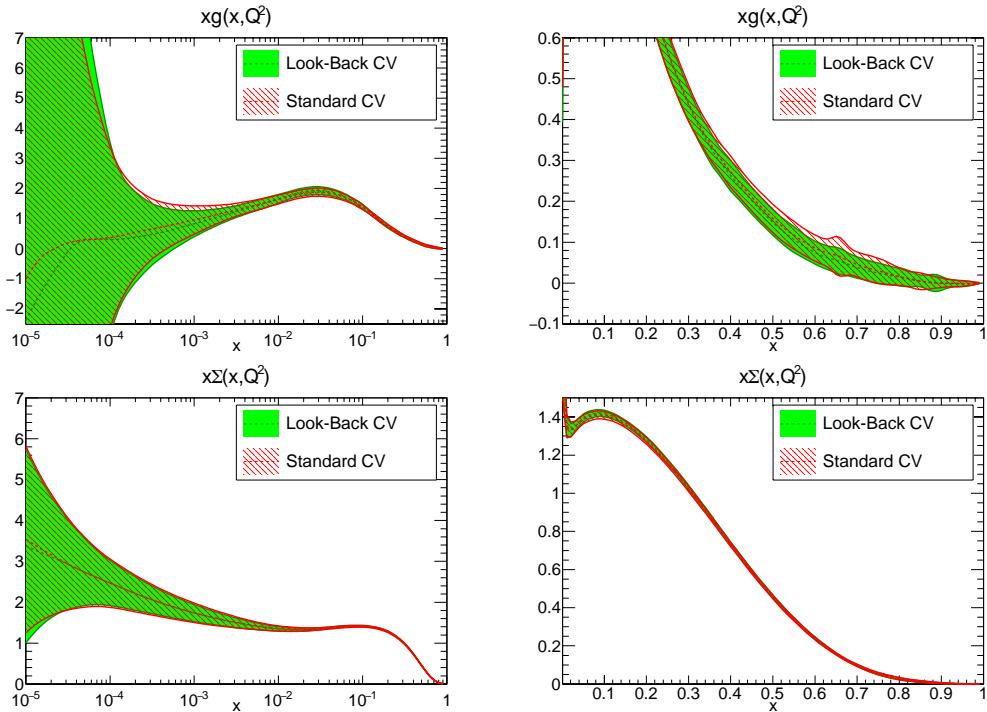


Figure 6.13: Comparison of PDFs obtained through look-back cross validation and NNPDF2.3 standard dynamical stopping. PDFs for the singlet and gluon are shown, with green bands representing fits using the look-back method and red demonstrating those with the slope-detection algorithm used in NNPDF2.3 and earlier.

## 6.5 Methodology for NNPDF3.0

We have performed an overview of the methodological developments made since the release of the NNPDF2.3 PDF set, with an aim to outline the procedure to

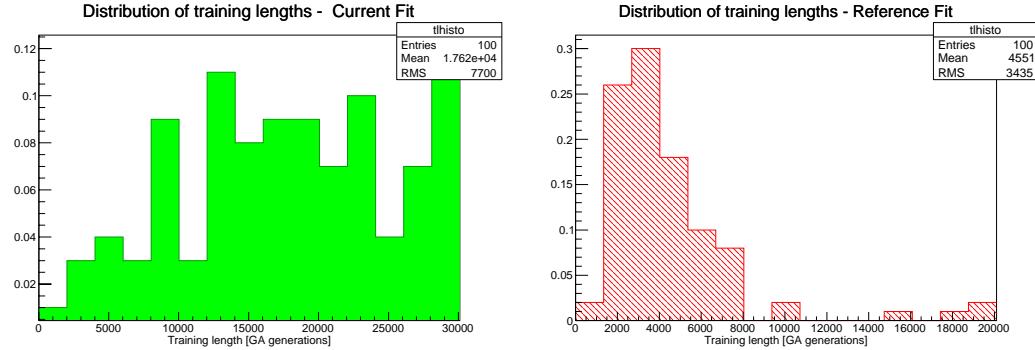


Figure 6.14: Comparison of training lengths in look-back cross-validation and NNPDF2.3 standard dynamical stopping. The left figure demonstrates the ‘optimal point’ determined by looking back over the while GA interval for the minimum validation  $\chi^2$ . The right figure shows the stopping point based upon the slope-detection algorithm.

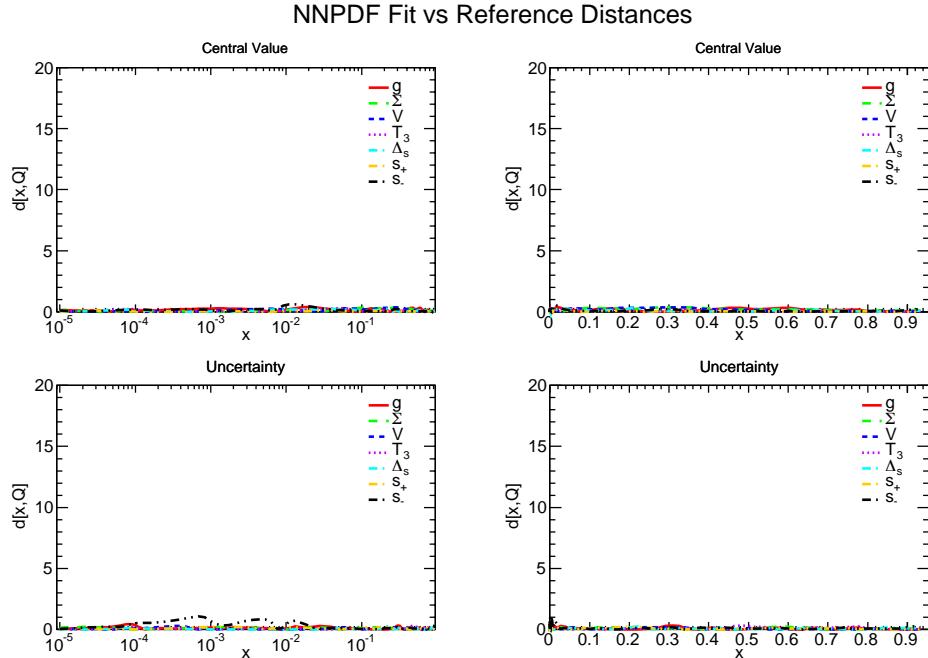


Figure 6.15: Distance comparison of two closure test fits with look-back stopping enabled and different maximum training lengths. Distances are computed between all evolution basis PDFs at the initial scale between  $N_{\text{gen}}^{\max} = 30,000$  and  $N_{\text{gen}}^{\max} = 60,000$  generation look-back fits.

be used in the forthcoming NNPDF3.0 set. To provide a stringent verification of the combined procedure, we shall now examine a set of closure test fits performed at various levels to differing generating PDF sets. In this section we present fits based upon a nodal genetic algorithm minimisation with look-back cross-validation stopping as detailed previously, with the iterative preprocessing procedure and new PDF fitting basis. Therefore the fits represent preliminary closure test results for the NNPDF3.0 methodology, upon a global pseudo-dataset of hadronic and DIS data. Results in this section will be presented using NLO calculations for the observables in the fit, although the conclusions will be very similar for an identical analysis at NNLO, as the closure test procedure is relatively insensitive to theory choices.

### 6.5.1 Closure tests for NNPDF3.0

Firstly let's consider the results obtained when fitting to an MSTW2008 generating PDF, the closure test guiding the methodological choices made so far in this section. In Figure ?? the ratio of the resulting closure test PDFs to the generating MSTW08 distributions are shown for some of the evolution basis PDFs. Here we show results for the kinematic region most constrained by the experimental pseudo-dataset:  $10^{-2} \leq x \leq 1$ . The level zero curves in Figure ?? closely reproduce the MSTW central values, achieving a final total  $\chi^2/N_{\text{dat}} = 0.00182$ . The uncertainty band in the case of the level zero result corresponds directly to the functional freedom available within the fitted pseudo-dataset. The level two fit clearly demonstrates the variations introduced by the simulated experimental noise, with the expected level of deviation clearly visible in the resulting PDFs. Given the simulated noise in the pseudo-dataset, the closure test still tracks the central value to an excellent level of accuracy, achieving an almost statistically ideal fit quality of  $\chi^2/N_{\text{dat}} = 1.00021$ .

As the preliminary NNPDF3.0 methodology has been validated against

closure test fits to the MSTW2008 set, it is important to test the procedure’s ability to reproduce a generating PDF with greater functional complexity. To verify the preliminary methodology in this case we now consider a closure test fit to the NNPDF2.3 PDF set. Figure ?? demonstrates once more the level zero and two closure test fits to NNPDF2.3. Even given the greater functional freedom present in the previous NNPDF determination, the 3.0 closure test provides an excellent reproduction of the generating functions, with fit qualities of  $\chi^2/N_{\text{dat}} = 0.00287$  and 1.01356 respectively. Once again the uncertainty due to parametrisation flexibility is demonstrated in the level zero fit, while the level two fit provides a closer simulation of a full fledged experimental data fit. These figures therefore suggest that the preliminary NNPDF3.0 methodological choices can accurately determine complex functional forms without any modification with respect to fits to much simpler parametrisations.

For a final closure test, we shall now consider a fit using the CT10 PDF set as a set of generating functions. In this way we can verify the NNPDF3.0 method in a way that is independent of the closure PDF set guiding the methodological development (MSTW2008) and previous NNPDF determinations. The results of the test, once more at level zero and two, are shown in Figure ?? . The closure test fit provides once again an excellent description of data, with  $\chi^2/N_{\text{dat}} = 0.00130$  for the level zero fit and 1.01324 for the level one. The procedure detailed here has now been validated against three different generating PDFs in a closure test and is able to convincingly reproduce the generating sets in each of them. We can therefore be confident that when applied to real experimental data the procedure will yield an accurate result up to theoretical uncertainties.

### 6.5.2 Improvements in data fits for NNPDF3.0

While we have now validated much of the methodology to be used in the NNPDF3.0 determination, we shall now finally investigate some of the expected

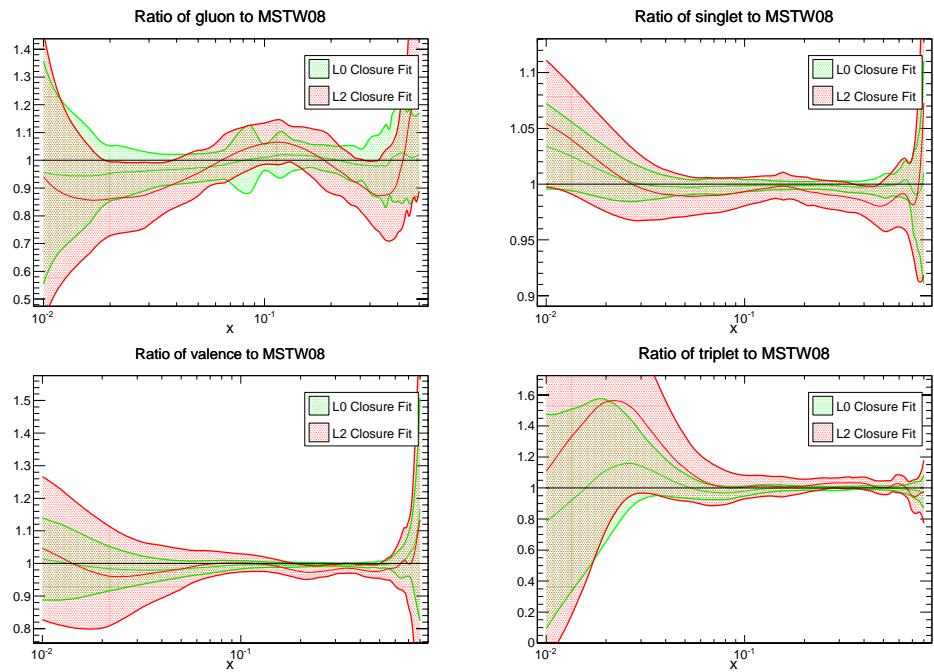


Figure 6.16: NNPDF3.0 methodology closure test fit to MSTW2008 NLO. Curves are shown normalised to the generating PDF for the gluon, singlet, triplet and valence distributions. The green curves show the results of a level zero closure test, while the red curves show the results of a level two test.

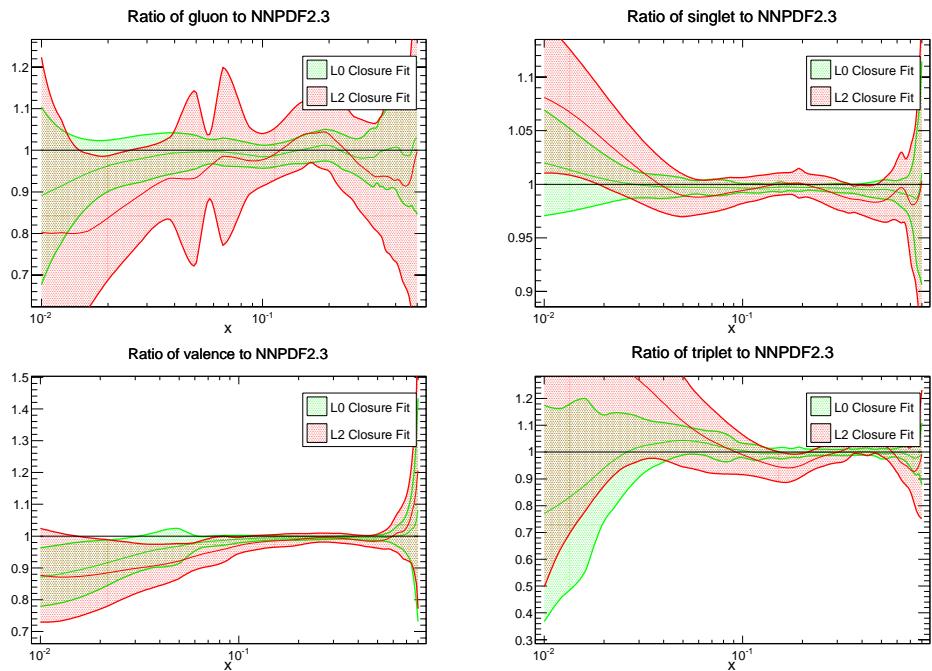


Figure 6.17: NNPDF3.0 methodology closure test fit to NNPDF2.3 NLO. Plots as in Figure ??.

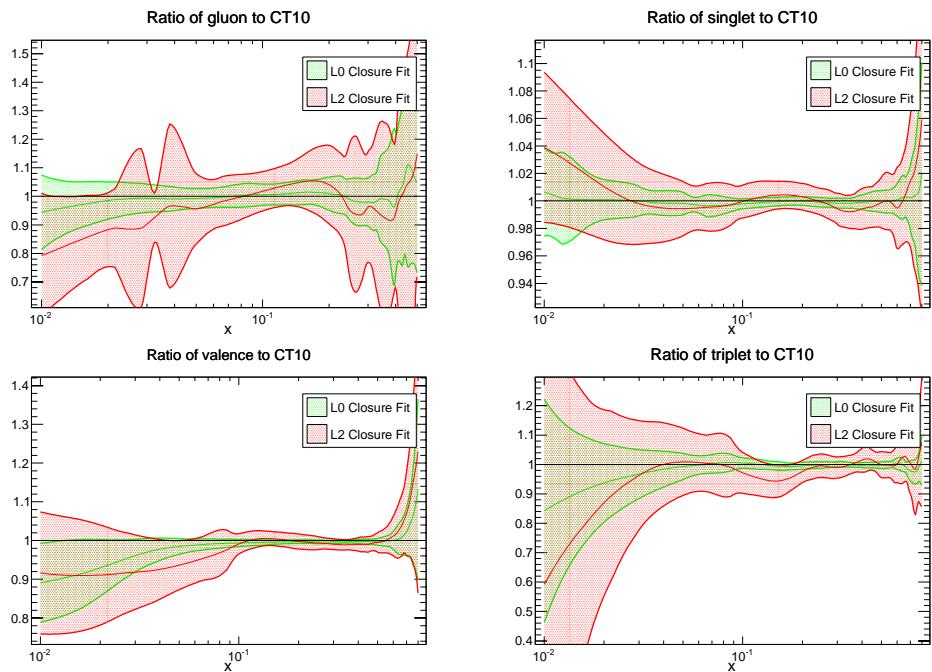


Figure 6.18: NNPDF3.0 methodology closure test fit to CT10 NLO. Plots as in Figure ??.

differences arising with respect to the NNPDF2.3 results in the case of a fit to experimental data. In order to directly assess the changes arising purely from the methodological differences in the two approaches, we shall perform two fits to a small common dataset, one with the full NNPDF2.3 machinery and the second with the improvements implemented in the NNPDF3.0 procedure. It should be noted that these results are of an extremely preliminary nature and as so should only be taken as roughly indicative of the final results. Furthermore, the full NNPDF3.0 set will benefit from a considerably expanded dataset with respect to the NNPDF2.3 determination.

For these test fits, we use a collider-only dataset to ensure a maximally consistent set of experimental data, including the full NNPDF2.3 LHC and Tevatron datasets, and the HERA-1 combined DIS results. Once more, the fits were run with a maximum number of generations of  $N_{\text{gen}} = 30,000$ . The NNPDF2.3-like fit was otherwise performed according to the settings of the central NNPDF2.3 fit. The NNPDF3.0 fits were performed with identical settings to the closure test fits described in the previous section.

Looking firstly at the gluon and singlet sectors, in Figure ?? we see the results of the two methodology test fits compared as a ratio to the NNPDF2.3 methodology fit's central value. The first feature to note is that in the region where data constraints in this test fit are largest, the two methodologies remain very consistent in their results, with the most significant changes occurring in the extrapolation regions and for the large- $x$  singlet. At small- $x$  the NNPDF3.0 methodology fit is more confident in the extrapolation for both singlet and gluon PDFs, resulting in a systematically smaller uncertainty. At large- $x$  there is a moderate shift in the gluon central value in the NNPDF3.0 result, and a broadening of uncertainties. The same pattern can be found in the large- $x$  gluon, where once again uncertainties are slightly larger and there is some change in central value. However both distributions remain in agreement within their

uncertainties, validating that the two methodologies remain compatible within the experimental uncertainty present in the test dataset.

To investigate the impact of the methodological changes to PDFs sensitive to the valence distributions and quark flavour separation, we plot the valence and triplet PDF combinations in Figure ???. In the valence PDF comparison, we see a similar pattern as for the singlet and gluon PDFs, where the low- $x$  result from the NNPDF3.0 methodology fit obtained a narrower distribution, and at high- $x$  the uncertainties are systematically larger. The triplet PDF shows by some way the largest differences between the two methodologies, with PDF uncertainties being significantly larger across the whole range of  $x$ . This effect is largely due to the much more flexible preprocessing used for the triplet PDF, where now there is no requirement that the PDF be preprocessed to zero at low- $x$ , the constraint now being entirely based on experimental data. Such a treatment leads to a rather conservative determination of the low- $x$  triplet, however this effect should be at least partially offset by increased data constraints in the full NNPDF3.0 determination.

Here we have seen that the results of the two methodologies provide consistent results when applied to the same experimental dataset. However there are significant changes in the fit results due to methodological improvements, particularly important in the PDF extrapolation regions at large and small values of parton- $x$ , and for PDF combinations sensitive to light flavour separation. As has been shown in the validation with closure tests, the methodological modifications, particularly in allowing for greater preprocessing flexibility, result in an improved reproduction of a test PDF distribution. The upgraded methodology should therefore provide a more reliable estimate of the parton densities in the proton.

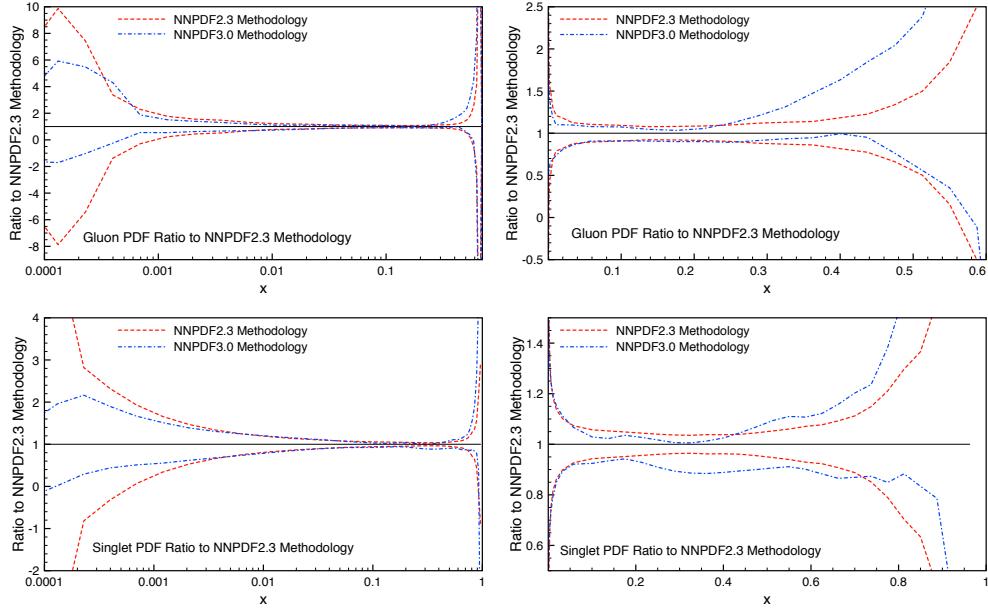


Figure 6.19: Comparison of NNPDF2.3 and NNPDF3.0 fitting methodologies when applied to a common experimental dataset. Here the gluon (top) and singlet (bottom) PDFs are shown, with all values normalised to the result of the NNPDF2.3 methodology fit.

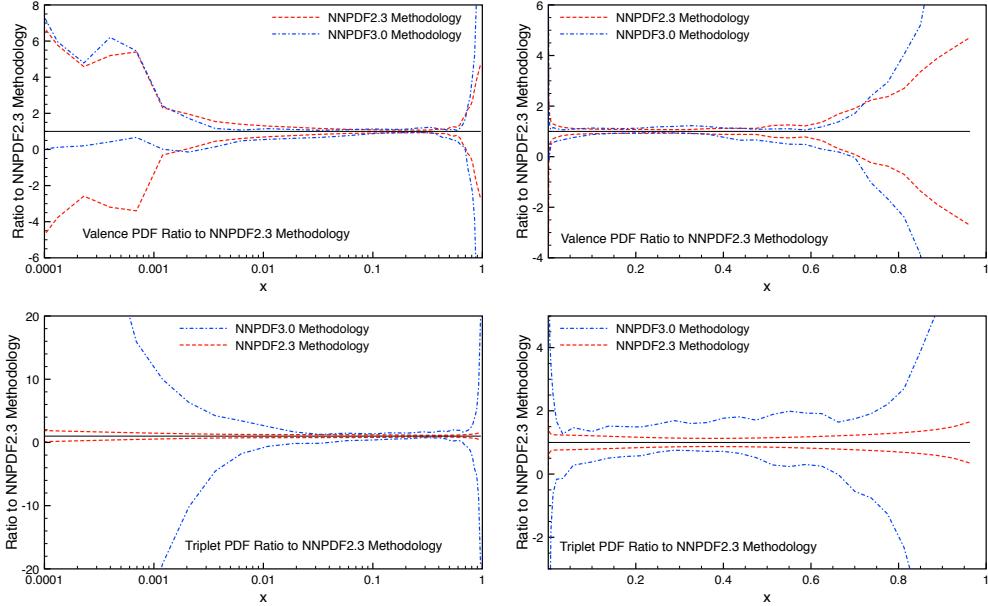


Figure 6.20: Comparison of NNPDF2.3 and NNPDF3.0 fitting methodologies when applied to a common experimental dataset, for the valence (top) and triplet (bottom) PDF combinations. Plots as in Figure ??.

# Chapter 7

## Conclusion

In this thesis we have discussed the impact of LHC measurements upon the extraction of parton distribution functions, and methods for the inclusion of such data into parton fits. With the LHC due to begin collisions at centre of mass energies of 13 TeV in 2015, the need for precise determinations of proton structure is as great as ever. Here we have described the efforts undertaken to provide the particle physics community with PDFs extracted via a methodologically sound procedure to an extensive experimental dataset including measurements from the first data runs of the LHC.

A number of tools have been developed that enable the study of a large collider dataset in the context of the NNPDF procedure. The Bayesian reweighting procedure, first implemented in a PDF determination by the NNPDF collaboration, allows for a rapid assessment of data impact. Bayesian reweighting also provides a deep check of the consistency of the NNPDF methodology, as a verification of the statistical behaviour of the Monte Carlo PDF ensemble. The method is however unsuitable for the inclusion of a large dataset, due to the need for an unpractically large prior Monte Carlo distribution. Therefore in order to enable the inclusion of a large LHC dataset in an NNPDF fit, further tools are required.

---

The FK method was introduced as a method by which the PDF evolution may be combined with the theoretical predictions for experimental observables in a process independent way. The resulting matrices, or FK tables provide an extremely efficient method of computing theoretical predictions based upon a varying input PDF set, reducing the task to a simple scalar product which can be efficiently optimised. This method enabled the inclusion of LHC data into a full NNPDF fit for the first time, without having to compromise on the accuracy of the calculation. While fast, the FK procedure is rather specialised to the task of PDF fitting as it only permits for the variation of input PDF set. For wider studies of the dependence of theoretical predictions upon parameters such as the strong coupling or factorisation and renormalisation scales, more flexible approaches as implemented in packages such as **FastNLO** or **APPLgrid** are more relevant.

A package for the interfacing of automated NLO calculational tools to such fast interpolating codes has been developed. The **MCgrid** package allows for the use of Monte Carlo event generators such as **SHERPA** along with a suitable one-loop generator, for the efficient variation of QCD parameters in theoretical predictions. Such an interface opens up the possibility of using such codes in applications such as  $\alpha_S$  determinations or PDF fits, alongside making PDF and scale variation more accessible in computationally challenging processes. The interface between event generators and interpolating tools remains however at the level of fixed order perturbation theory. In principle an extension for the fast computation of observables with parton shower effects included would be particularly desirable, making this an important avenue for future research.

These tools have been applied to the study of LHC measurements and their impact upon PDF distributions. To date two NNPDF results including LHC constraints have been published. The NNPDF2.2 determination included a set of  $W$  boson production asymmetry measurements at both the LHC and the Tevatron, by the method of Bayesian reweighting. In such a way the

reweighting method was extensively validated, and new constraints were placed upon PDFs from LHC data for the first time. As the dataset available from the LHC expanded, the need for a comprehensive refit including all appropriate measurements increased. The resulting PDF set, NNPDF2.3, utilised the FK procedure for all of the included processes and so was able to include all available LHC measurements of interest to PDF fits at the time. The NNPDF2.3 set provided a precise determination ideally suited for further applications at the LHC.

Following the NNPDF2.3 set, the development of the `nnpdf++` project has allowed for a greater scope in investigating methodological elements, permitting a large scale re-evaluation of the procedure used in the NNPDF2.3 family of fits. The closure testing procedure introduced in Chapter ?? now forming the basis for the development in methodology post-NNPDF2.3. Insights provided by a detailed study of the NNPDF procedure when applied to closure tests have informed a number of new approaches in minimisation and stopping for the next global PDF set produced by the collaboration. The next release, NNPDF3.0, being validated using the closure testing procedure and including an expanded LHC dataset will provide the most precise and methodologically sound determination of parton distribution functions.

# Bibliography

- [1] R.P. Feynman. *Photon hadron interactions.* W.A. Benjamin, New York, 1972.
- [2] R.P. Feynman. The behaviour of hadron collisions at extreme energies. In C. N Yang, editor, *Stony Brook 1969, Proceedings, Conference On High Energy Collisions*, pages 237–258, 1969.
- [3] R.P. Feynman. Very high-energy collision of hadrons. *Phys. Rev. Lett.*, 23:1415–1417, 1969.
- [4] J.D. Bjorken. Asymptotic Sum Rules at Infinite Momentum. *Phys. Rev.*, 179:1547–1553, 1969.
- [5] D Callan, C Gross. *Phys. Rev. Lett.*, 22:156, 1969.
- [6] G Sterman et al. Handbook of perturbative qcd. *Rev. Mod. Phys.*, 67(1):157–248, January 1995.
- [7] S. Aid et al. A Measurement and QCD analysis of the proton structure function  $f_2(x, q^{**2})$  at HERA. *Nucl.Phys.*, B470:3–40, 1996, hep-ex/9603004.
- [8] D.J. Gross and Frank Wilczek. Asymptotically Free Gauge Theories. 1. *Phys. Rev.*, D8:3633–3652, 1973.

- [9] Howard Georgi and H. David Politzer. Electroproduction scaling in an asymptotically free theory of strong interactions. *Phys.Rev.*, D9:416–420, 1974.
- [10] E.G. Floratos, D.A. Ross, and Christopher T. Sachrajda. Higher Order Effects in Asymptotically Free Gauge Theories: The Anomalous Dimensions of Wilson Operators. *Nucl.Phys.*, B129:66–88, 1977.
- [11] Guido Altarelli and G. Parisi. Asymptotic Freedom in Parton Language. *Nucl.Phys.*, B126:298, 1977.
- [12] Antonio Gonzalez-Arroyo, C. Lopez, and F.J. Yndurain. Second Order Contributions to the Structure Functions in Deep Inelastic Scattering. 1. Theoretical Calculations. *Nucl.Phys.*, B153:161–186, 1979.
- [13] E.G. Floratos, D.A. Ross, and Christopher T. Sachrajda. Higher Order Effects in Asymptotically Free Gauge Theories. 2. Flavor Singlet Wilson Operators and Coefficient Functions. *Nucl.Phys.*, B152:493, 1979.
- [14] W. Furmanski and R. Petronzio. Singlet Parton Densities Beyond Leading Order. *Phys.Lett.*, B97:437, 1980.
- [15] G. Curci, W. Furmanski, and R. Petronzio. Evolution of Parton Densities Beyond Leading Order: The Nonsinglet Case. *Nucl.Phys.*, B175:27, 1980.
- [16] Antonio Gonzalez-Arroyo and C. Lopez. Second Order Contributions to the Structure Functions in Deep Inelastic Scattering. 3. The Singlet Case. *Nucl.Phys.*, B166:429, 1980.
- [17] E.G. Floratos, C. Kounnas, and R. Lacaze. Higher Order QCD Effects in Inclusive Annihilation and Deep Inelastic Scattering. *Nucl.Phys.*, B192:417, 1981.

- [18] R. Hamberg and W.L. van Neerven. The Correct renormalization of the gluon operator in a covariant gauge. *Nucl.Phys.*, B379:143–171, 1992.
- [19] S. Moch, J.A.M. Vermaseren, and A. Vogt. The Three loop splitting functions in QCD: The Nonsinglet case. *Nucl.Phys.*, B688:101–134, 2004, hep-ph/0403192.
- [20] A. Vogt, S. Moch, and J.A.M. Vermaseren. The Three-loop splitting functions in QCD: The Singlet case. *Nucl.Phys.*, B691:129–181, 2004, hep-ph/0404111.
- [21] Guido Altarelli, R. Keith Ellis, and G. Martinelli. Leptoproduction and Drell-Yan Processes Beyond the Leading Approximation in Chromodynamics. *Nucl. Phys.*, B143:521, 1978.
- [22] G. Altarelli, G. Parisi. *Nucl. Phys. B*, 20:28, 1977.
- [23] Y.L Dokshitzer. *Sov. Phys.*, 46:641, 1977.
- [24] L.N. Gribov, V.N Lipatov. *Sov. J. Nucl. Phys.*, 15:438, 1972.
- [25] L.N. Lipatov. *Sov. J. Nucl. Phys.*, 20:94, 1975.
- [26] Gavin P. Salam and Juan Rojo. A Higher Order Perturbative Parton Evolution Toolkit (HOPPET). *Comput. Phys. Commun.*, 180:120–156, 2009, 0804.3755.
- [27] M. Botje. QCDNUM: Fast QCD Evolution and Convolution. *Comput.Phys.Commun.*, 182:490–532, 2011, 1005.1481.
- [28] Valerio Bertone, Stefano Carrazza, and Juan Rojo. APFEL: A PDF Evolution Library with QED corrections. 2013, 1310.1394.

- [29] A. Vogt. Efficient evolution of unpolarized and polarized parton distributions with QCD-PEGASUS. *Comput. Phys. Commun.*, 170:65–92, 2005, hep-ph/0408244.
- [30] W.K. Tung, H.L. Lai, A. Belyaev, J. Pumplin, D. Stump, et al. Heavy Quark Mass Effects in Deep Inelastic Scattering and Global QCD Analysis. *JHEP*, 0702:053, 2007, hep-ph/0611254.
- [31] M. Buza, Y. Matiounine, J. Smith, R. Migneron, and W.L. van Neerven. Heavy quark coefficient functions at asymptotic values  $Q^{**2} \gg m^{**2}$ . *Nucl.Phys.*, B472:611–658, 1996, hep-ph/9601302.
- [32] M. Buza, Y. Matiounine, J. Smith, and W.L. van Neerven. Charm electroproduction viewed in the variable flavor number scheme versus fixed order perturbation theory. *Eur.Phys.J.*, C1:301–320, 1998, hep-ph/9612398.
- [33] R.S. Thorne and W.K. Tung. PQCD Formulations with Heavy Quark Masses and Global Analysis. 2008, 0809.0714.
- [34] I Kramer, Michael, Fredrick I. Olness, and Davison E. Soper. Treatment of heavy quarks in deeply inelastic scattering. *Phys.Rev.*, D62:096007, 2000, hep-ph/0003035.
- [35] John C. Collins, Frank Wilczek, and A. Zee. Low-Energy Manifestations of Heavy Particles: Application to the Neutral Current. *Phys.Rev.*, D18:242, 1978.
- [36] John C. Collins. Hard scattering factorization with heavy quarks: A General treatment. *Phys.Rev.*, D58:094002, 1998, hep-ph/9806259.
- [37] R.S. Thorne and R.G. Roberts. A Practical procedure for evolving heavy flavor structure functions. *Phys.Lett.*, B421:303–311, 1998, hep-ph/9711223.

- [38] R.S. Thorne and R.G. Roberts. An Ordered analysis of heavy flavor production in deep inelastic scattering. *Phys.Rev.*, D57:6871–6898, 1998, hep-ph/9709442.
- [39] Matteo Cacciari, Mario Greco, and Paolo Nason. The P(T) spectrum in heavy flavor hadroproduction. *JHEP*, 9805:007, 1998, hep-ph/9803400.
- [40] Stefano Forte, Eric Laenen, Paolo Nason, and Juan Rojo. Heavy quarks in deep-inelastic scattering. *Nucl.Phys.*, B834:116–162, 2010, 1001.2312.
- [41] D. Dolgov, R. Brower, S. Capitani, John W. Negele, A. Pochinsky, et al. Moments of structure functions in full QCD. *Nucl.Phys.Proc.Suppl.*, 94:303–306, 2001, hep-lat/0011010.
- [42] R. Horsley. The Lattice calculation of moments of structure functions. pages 313–322, 2004, hep-lat/0412007.
- [43] M. Gockeler, R. Horsley, D. Pleiter, Paul E.L. Rakow, and G. Schierholz. A Lattice determination of moments of unpolarised nucleon structure functions using improved Wilson fermions. *Phys.Rev.*, D71:114511, 2005, hep-ph/0410187.
- [44] Wolfram Schroers. Parton distributions from the lattice. *Nucl.Phys.*, A755:333–336, 2005, hep-ph/0501156.
- [45] S. Alekhin, J. Bluemlein, and S. Moch. The ABM parton distributions tuned to LHC data. *Phys.Rev.*, D89:054028, 2014, 1310.3059.
- [46] S. Alekhin, J. Blumlein, and S. Moch. Parton Distribution Functions and Benchmark Cross Sections at NNLO. *Phys.Rev.*, D86:054009, 2012, 1202.2281.

- [47] S. Alekhin, J. Blumlein, S. Klein, and S. Moch. The 3, 4, and 5-flavor NNLO Parton from Deep-Inelastic-Scattering Data and at Hadron Colliders. *Phys.Rev.*, D81:014032, 2010, 0908.2766.
- [48] Jun Gao, Marco Guzzi, Joey Huston, Hung-Liang Lai, Zhao Li, et al. The CT10 NNLO Global Analysis of QCD. *Phys.Rev.*, D89:033009, 2014, 1302.6246.
- [49] Hung-Liang Lai, Marco Guzzi, Joey Huston, Zhao Li, Pavel M. Nadolsky, et al. New parton distributions for collider physics. *Phys.Rev.*, D82:074024, 2010, 1007.2241.
- [50] Pavel M. Nadolsky, Hung-Liang Lai, Qing-Hong Cao, Joey Huston, Jon Pumplin, et al. Implications of CTEQ global analysis for collider observables. *Phys.Rev.*, D78:013004, 2008, 0802.0007.
- [51] J.F. Owens, A. Accardi, and W. Melnitchouk. Global parton distributions with nuclear and finite- $Q^2$  corrections. *Phys.Rev.*, D87(9):094012, 2013, 1212.1702.
- [52] P. Jimenez-Delgado and E. Reya. Dynamical NNLO parton distributions. *Phys.Rev.*, D79:074023, 2009, 0810.4274.
- [53] M. Gluck, P. Jimenez-Delgado, and E. Reya. Dynamical parton distributions of the nucleon and very small-x physics. *Eur.Phys.J.*, C53:355–366, 2008, 0709.0614.
- [54] F.D. Aaron et al. Combined Measurement and QCD Analysis of the Inclusive e+- p Scattering Cross Sections at HERA. *JHEP*, 1001:109, 2010, 0911.0884.
- [55] : et al. Parton distribution functions at LO, NLO and NNLO with correlated uncertainties between orders. 2014, 1404.4234.

- [56] A. D. Martin, W. J. Stirling, R. S. Thorne, and G. Watt. Parton distributions for the LHC. *Eur. Phys. J.*, C63:189–285, 2009, 0901.0002.
- [57] A.D. Martin, A.J. Th.M. Mathijssen, W.J. Stirling, R.S. Thorne, B.J.A. Watt, et al. Extended Parameterisations for MSTW PDFs and their effect on Lepton Charge Asymmetry from W Decays. *Eur.Phys.J.*, C73:2318, 2013, 1211.1215.
- [58] Alan D. Martin, R.G. Roberts, W. James Stirling, and R.S. Thorne. Parton distributions: A New global analysis. *Eur.Phys.J.*, C4:463–496, 1998, hep-ph/9803445.
- [59] Alan D. Martin, R.G. Roberts, W.J. Stirling, and R.S. Thorne. MRST2001: Partons and  $\alpha_s$  from precise deep inelastic scattering and Tevatron jet data. *Eur.Phys.J.*, C23:73–87, 2002, hep-ph/0110215.
- [60] A.D. Martin, R.G. Roberts, W.J. Stirling, and R.S. Thorne. NNLO global parton analysis. *Phys.Lett.*, B531:216–224, 2002, hep-ph/0201127.
- [61] A.D. Martin, R.G. Roberts, W.J. Stirling, and R.S. Thorne. Parton distributions incorporating QED contributions. *Eur.Phys.J.*, C39:155–161, 2005, hep-ph/0411040.
- [62] Richard D. Ball, Valerio Bertone, Stefano Carrazza, Christopher S. Deans, Luigi Del Debbio, et al. Parton distributions with LHC data. *Nucl.Phys.*, B867:244–289, 2013, 1207.1303.
- [63] Richard D. Ball, Valerio Bertone, Francesco Cerutti, Luigi Del Debbio, Stefano Forte, et al. Reweighting and Unweighting of Parton Distributions and the LHC W lepton asymmetry data. *Nucl.Phys.*, B855:608–638, 2012, 1108.1758.

- [64] Richard D.Ball et al. Unbiased global determination of parton distributions and their uncertainties at nnlo and at lo. *Nucl. Phys.*, B855:153–221, 2012, 1107.2652.
- [65] Richard D. Ball, Valerio Bertone, Francesco Cerutti, Luigi Del Debbio, Stefano Forte, et al. Impact of Heavy Quark Masses on Parton Distributions and LHC Phenomenology. *Nucl.Phys.*, B849:296–363, 2011, 1101.1300.
- [66] Richard D. Ball, Luigi Del Debbio, Stefano Forte, Alberto Guffanti, Jose I. Latorre, et al. A first unbiased global NLO determination of parton distributions and their uncertainties. *Nucl.Phys.*, B838:136–206, 2010, 1002.4407.
- [67] Richard D. Ball et al. A Determination of parton distributions with faithful uncertainty estimation. *Nucl.Phys.*, B809:1–63, 2009, 0808.1231.
- [68] Michiel Botje et al. The PDF4LHC Working Group Interim Recommendations. 2011, 1101.0538.
- [69] J. A. M. Vermaseren, A. Vogt, and S. Moch. The third-order QCD corrections to deep-inelastic scattering by photon exchange. *Nucl. Phys.*, B724:3–182, 2005, hep-ph/0504242.
- [70] S. Moch, J.A.M. Vermaseren, and A. Vogt. Third-order QCD corrections to the charged-current structure function F(3). *Nucl.Phys.*, B813:220–258, 2009, 0812.4168.
- [71] M. Buza and W.L. van Neerven. O ( $\alpha_s^{**2}$ ) contributions to charm production in charged current deep inelastic lepton - hadron scattering. *Nucl.Phys.*, B500:301–324, 1997, hep-ph/9702242.

- [72] M. Buza, Y. Matiounine, J. Smith, and W.L. van Neerven.  $O(\alpha_s^{**2})$  corrections to polarized heavy flavor production at  $Q^{**2} \gg m^{**2}$ . *Nucl.Phys.*, B485:420–456, 1997, hep-ph/9608342.
- [73] J. Blumlein, A. De Freitas, W.L. van Neerven, and S. Klein. The Longitudinal Heavy Quark Structure Function  $F^{**Q}$  anti- $Q(L)$  in the Region  $Q^{**2} \gg m^{**2}$  at  $O(\alpha_s^{**3}(s))$ . *Nucl.Phys.*, B755:272–285, 2006, hep-ph/0608024.
- [74] Isabella Bierenbaum, Johannes Blumlein, and Sebastian Klein. Two-Loop Massive Operator Matrix Elements and Unpolarized Heavy Flavor Production at Asymptotic Values  $Q^{**2} \gg m^{**2}$ . *Nucl.Phys.*, B780:40–75, 2007, hep-ph/0703285.
- [75] I. Bierenbaum, J. Blumlein, and S. Klein. Calculation of massive 2-loop operator matrix elements with outer gluon lines. *Phys.Lett.*, B648:195–200, 2007, hep-ph/0702265.
- [76] Isabella Bierenbaum, Johannes Blumlein, and Sebastian Klein. The Gluonic Operator Matrix Elements at  $O(\alpha_s(s)^{**2})$  for DIS Heavy Flavor Production. *Phys.Lett.*, B672:401–406, 2009, 0901.0669.
- [77] Johannes Blümlein, Alexander Hasselhuhn, and Torsten Pfoh. The  $O(\alpha_s^2)$  heavy quark corrections to charged current deep-inelastic scattering at large virtualities. *Nucl.Phys.*, B881:1–41, 2014, 1401.4352.
- [78] B. Badelek and J. Kwiecinski. Shadowing in the deuteron and the new  $f_2(n) / f_2(p)$  measurements. *Phys.Rev.*, D50:4–8, 1994, hep-ph/9401314.
- [79] A.C. Benvenuti et al. A High Statistics Measurement of the Proton Structure Functions  $F(2)$  ( $x, Q^{**2}$ ) and  $R$  from Deep Inelastic Muon Scattering at High  $Q^{**2}$ . *Phys.Lett.*, B223:485, 1989.

- [80] A.C. Benvenuti et al. A High Statistics Measurement of the Deuteron Structure Functions  $F_2(X, Q^2)$  and  $R$  From Deep Inelastic Muon Scattering at High  $Q^2$ . *Phys.Lett.*, B237:592, 1990.
- [81] M. Arneodo et al. Measurement of the proton and deuteron structure functions,  $F_2(p)$  and  $F_2(d)$ , and of the ratio  $\sigma_L / \sigma_T$ . *Nucl.Phys.*, B483:3–43, 1997, hep-ph/9610231.
- [82] M. Arneodo et al. Accurate measurement of  $F_2(d) / F_2(p)$  and  $R^{**d} - R^{**p}$ . *Nucl.Phys.*, B487:3–26, 1997, hep-ex/9611022.
- [83] M.R. Adams et al. Proton and deuteron structure functions in muon scattering at 470-GeV. *Phys.Rev.*, D54:3006–3056, 1996.
- [84] L.W. Whitlow, E.M. Riordan, S. Dasu, Stephen Rock, and A. Bodek. Precise measurements of the proton and deuteron structure functions from a global analysis of the SLAC deep inelastic electron scattering cross-sections. *Phys.Lett.*, B282:475–482, 1992.
- [85] L.W. Whitlow, Stephen Rock, A. Bodek, E.M. Riordan, and S. Dasu. A Precise extraction of  $R = \sigma_L / \sigma_T$  from a global analysis of the SLAC deep inelastic e p and e d scattering cross-sections. *Phys.Lett.*, B250:193–198, 1990.
- [86] J. Breitweg et al. ZEUS results on the measurement and phenomenology of  $F(2)$  at low  $x$  and low  $Q^{**2}$ . *Eur.Phys.J.*, C7:609–630, 1999, hep-ex/9809005.
- [87] S. Chekanov et al. Measurement of the neutral current cross-section and  $F(2)$  structure function for deep inelastic e + p scattering at HERA. *Eur.Phys.J.*, C21:443–471, 2001, hep-ex/0105090.

- [88] S. Chekanov et al. Measurement of high  $Q^{**2}$  e- p neutral current cross-sections at HERA and the extraction of  $xF(3)$ . *Eur.Phys.J.*, C28:175–201, 2003, hep-ex/0208040.
- [89] S. Chekanov et al. High  $Q^{**2}$  neutral current cross-sections in e+ p deep inelastic scattering at  $s^{**}(1/2) = 318\text{-GeV}$ . *Phys.Rev.*, D70:052001, 2004, hep-ex/0401003.
- [90] C. Adloff et al. Deep inelastic inclusive e p scattering at low x and a determination of alpha(s). *Eur.Phys.J.*, C21:33–61, 2001, hep-ex/0012053.
- [91] C. Adloff et al. Measurement of neutral and charged current cross-sections in electron - proton collisions at high  $Q^2$ . *Eur.Phys.J.*, C19:269–288, 2001, hep-ex/0012052.
- [92] C. Adloff et al. Measurement and QCD analysis of neutral and charged current cross-sections at HERA. *Eur.Phys.J.*, C30:1–32, 2003, hep-ex/0304003.
- [93] S. Chekanov et al. Measurement of high  $Q^{**2}$  charged current cross-sections in e+ p deep inelastic scattering at HERA. *Eur.Phys.J.*, C32:1–16, 2003, hep-ex/0307043.
- [94] V. Andreev et al. Measurement of inclusive  $ep$  cross sections at high  $Q^2$  at  $\sqrt{s} = 225$  and 252 GeV and of the longitudinal proton structure function  $F_L$  at HERA. *Eur.Phys.J.*, C74:2814, 2014, 1312.4821.
- [95] S. Chekanov et al. Measurement of the Longitudinal Proton Structure Function at HERA. *Phys.Lett.*, B682:8–22, 2009, 0904.1092.
- [96] C. Adloff et al. Inclusive D0 and  $D^*{}^{+-}$  production in deep inelastic e p scattering at HERA. *Z.Phys.*, C72:593–605, 1996, hep-ex/9607012.

- [97] C. Adloff et al. Measurement of  $D^{\ast+}$ - meson production and  $F_2(c)$  in deep inelastic scattering at HERA. *Phys.Lett.*, B528:199–214, 2002, hep-ex/0108039.
- [98] A. Aktas et al. Measurement of  $F(2)^{**c}$  anti- $c$  and  $F(2)^{**b}$  anti- $b$  at low  $Q^2$  and  $x$  using the H1 vertex detector at HERA. *Eur.Phys.J.*, C45:23–33, 2006, hep-ex/0507081.
- [99] A. Aktas et al. Measurement of  $F_2(c\bar{c})$  and  $F_2(b\bar{b})$  at high  $Q^2$  using the H1 vertex detector at HERA. *Eur.Phys.J.*, C40:349–359, 2005, hep-ex/0411046.
- [100] J. Breitweg et al. Measurement of  $D^{\ast+}$  production and the charm contribution to  $F(2)$  in deep inelastic scattering at HERA. *Eur.Phys.J.*, C12:35–52, 2000, hep-ex/9908012.
- [101] S. Chekanov et al. Measurement of  $D^{\ast+}$  production in deep inelastic e+p scattering at HERA. *Phys.Rev.*, D69:012004, 2004, hep-ex/0308068.
- [102] S. Chekanov et al. Measurement of D mesons production in deep inelastic scattering at HERA. *JHEP*, 0707:074, 2007, 0704.3562.
- [103] F. D. Aaron et al. Combined Measurement and QCD Analysis of the Inclusive ep Scattering Cross Sections at HERA. *JHEP*, 01:109, 2010, 0911.0884.
- [104] H. Abramowicz et al. Combination and QCD Analysis of Charm Production Cross Section Measurements in Deep-Inelastic ep Scattering at HERA. *Eur.Phys.J.*, C73:2311, 2013, 1211.1182.
- [105] M. Tzanov et al. Precise measurement of neutrino and anti-neutrino differential cross sections. *Phys.Rev.*, D74:012008, 2006, hep-ex/0509010.

- [106] G. Onengut et al. Measurement of nucleon structure functions in neutrino scattering. *Phys.Lett.*, B632:65–75, 2006.
- [107] Stefano Forte and Graeme Watt. Progress in the Determination of the Partonic Structure of the Proton. *Ann.Rev.Nucl.Part.Sci.*, 63:291–328, 2013, 1301.6754.
- [108] M. Goncharov et al. Precise measurement of dimuon production cross-sections in muon neutrino Fe and muon anti-neutrino Fe deep inelastic scattering at the Tevatron. *Phys.Rev.*, D64:112006, 2001, hep-ex/0102049.
- [109] Charalampos Anastasiou, Lance J. Dixon, Kirill Melnikov, and Frank Petriello. High precision QCD at hadron colliders: Electroweak gauge boson rapidity distributions at NNLO. *Phys.Rev.*, D69:094008, 2004, hep-ph/0312266.
- [110] Stefano Catani, Leandro Cieri, Giancarlo Ferrera, Daniel de Florian, and Massimiliano Grazzini. Vector boson production at hadron colliders: a fully exclusive QCD calculation at NNLO. *Phys.Rev.Lett.*, 103:082001, 2009, 0903.2120.
- [111] Stefano Catani, Giancarlo Ferrera, and Massimiliano Grazzini. W Boson Production at Hadron Colliders: The Lepton Charge Asymmetry in NNLO QCD. *JHEP*, 1005:006, 2010, 1002.3115.
- [112] G. Moreno et al. Dimuon production in proton - copper collisions at  $s^{**}(1/2) = 38.8\text{-GeV}$ . *Phys. Rev.*, D43:2815–2836, 1991.
- [113] D. de Florian and R. Sassot. Nuclear parton distributions at next-to-leading order. *Phys.Rev.*, D69:074028, 2004, hep-ph/0311227.

- [114] M. Hirai, S. Kumano, and T.-H. Nagai. Determination of nuclear parton distribution functions and their uncertainties in next-to-leading order. *Phys.Rev.*, C76:065207, 2007, 0709.3038.
- [115] Sergey A. Kulagin and R. Petti. Neutrino inelastic scattering off nuclei. *Phys.Rev.*, D76:094023, 2007, hep-ph/0703033.
- [116] K.J. Eskola, H. Paukkunen, and C.A. Salgado. EPS09: A New Generation of NLO and LO Nuclear Parton Distribution Functions. *JHEP*, 0904:065, 2009, 0902.4154.
- [117] Richard D. Ball et al. Precision determination of electroweak parameters and the strange content of the proton from neutrino deep-inelastic scattering. *Nucl.Phys.*, B823:195–233, 2009, 0906.1958.
- [118] A. Accardi, W. Melnitchouk, J.F. Owens, M.E. Christy, C.E. Keppel, et al. Uncertainties in determining parton distributions at large x. *Phys.Rev.*, D84:014008, 2011, 1102.3686.
- [119] L.T. Brady, A. Accardi, W. Melnitchouk, and J.F. Owens. Impact of PDF uncertainties at large x on heavy boson production. *JHEP*, 1206:019, 2012, 1110.5398.
- [120] Jason C. Webb. Measurement of continuum dimuon production in 800-GeV/C proton nucleon collisions. 2003, hep-ex/0301031.
- [121] R.S. Towell et al. Improved measurement of the anti-d / anti-u asymmetry in the nucleon sea. *Phys.Rev.*, D64:052002, 2001, hep-ex/0103030.
- [122] V.M. Abazov et al. Measurement of the shape of the boson rapidity distribution for  $p\bar{p} \rightarrow Z/\gamma^* \rightarrow e^+e^- + X$  events produced at  $\sqrt{s}$  of 1.96-TeV. *Phys.Rev.*, D76:012003, 2007, hep-ex/0702025.

- [123] D. Acosta et al. Measurement of the forward-backward charge asymmetry from  $W \rightarrow e\nu$  production in  $p\bar{p}$  collisions at  $\sqrt{s} = 1.96$  TeV. *Phys.Rev.*, D71:051104, 2005, hep-ex/0501023.
- [124] V.M. Abazov et al. Measurement of the muon charge asymmetry from  $W$  boson decays. *Phys.Rev.*, D77:011106, 2008, 0709.4254.
- [125] V.M. Abazov et al. Measurement of the electron charge asymmetry in  $p\bar{p} \rightarrow W + X \rightarrow e\nu + X$  events at  $\sqrt{s} = 1.96$ -TeV. *Phys.Rev.Lett.*, 101:211801, 2008, 0807.3367.
- [126] F. Abe et al. Measurement of the lepton charge asymmetry in  $W$  boson decays produced in  $p\bar{p}$  collisions. *Phys.Rev.Lett.*, 81:5754–5759, 1998, hep-ex/9809001.
- [127] T. Aaltonen et al. First measurement of the production of a  $W$  boson in association with a single charm quark in  $p\bar{p}$  collisions at  $\sqrt{s} = 1.96$ -TeV. *Phys.Rev.Lett.*, 100:091803, 2008, 0711.2901.
- [128] V.M. Abazov et al. Measurement of the ratio of the  $p\bar{p} \rightarrow W^+c^-$  jet cross section to the inclusive  $p\bar{p} \rightarrow W +$  jets cross section. *Phys.Lett.*, B666:23–30, 2008, 0803.2259.
- [129] W.J. Stirling and E. Vryonidou. Charm production in association with an electroweak gauge boson at the LHC. *Phys.Rev.Lett.*, 109:082002, 2012, 1203.6781.
- [130] Stephen D. Ellis, Zoltan Kunszt, and Davison E. Soper. Two jet production in hadron collisions at order  $\alpha_s^{**3}$  in QCD. *Phys.Rev.Lett.*, 69:1496–1499, 1992.

- [131] W.T. Giele, E.W. Nigel Glover, and David A. Kosower. The Two-Jet Differential Cross Section at  $\mathcal{O}(\alpha_s^3)$  in Hadron Collisions. *Phys.Rev.Lett.*, 73:2019–2022, 1994, hep-ph/9403347.
- [132] Zoltan Nagy. Three jet cross-sections in hadron hadron collisions at next-to-leading order. *Phys.Rev.Lett.*, 88:122003, 2002, hep-ph/0110315.
- [133] Zoltan Nagy. Next-to-leading order calculation of three jet observables in hadron hadron collision. *Phys.Rev.*, D68:094002, 2003, hep-ph/0307268.
- [134] James Currie, Aude Gehrmann-De Ridder, E.W.N. Glover, and Joao Pires. NNLO QCD corrections to jet production at hadron colliders from gluon scattering. *JHEP*, 1401:110, 2014, 1310.3993.
- [135] E.W. Nigel Glover, C. Oleari, and M.E. Tejeda-Yeomans. Two loop QCD corrections to gluon-gluon scattering. *Nucl.Phys.*, B605:467–485, 2001, hep-ph/0102201.
- [136] E.W. Nigel Glover and M.E. Tejeda-Yeomans. One loop QCD corrections to gluon-gluon scattering at NNLO. *JHEP*, 0105:010, 2001, hep-ph/0104178.
- [137] James Currie, Aude Gehrmann-De Ridder, Thomas Gehrmann, E.W. Nigel Glover, and Joao Pires. NNLO QCD corrections to dijet production at hadron colliders. *PoS*, RADCOR2013:004, 2014, 1312.5608.
- [138] Daniel de Florian, Patriz Hinderer, Asmita Mukherjee, Felix Ringer, and Werner Vogelsang. Approximate next-to-next-to-leading order corrections to hadronic jet production. *Phys.Rev.Lett.*, 112:082001, 2014, 1310.7192.
- [139] Nikolaos Kidonakis and J.F. Owens. Effects of higher order threshold corrections in high E(T) jet production. *Phys.Rev.*, D63:054019, 2001, hep-ph/0007268.

- [140] Meduri C. Kumar and Sven-Olaf Moch. Phenomenology of threshold corrections for inclusive jet production at hadron colliders. *Phys.Lett.*, B730:122–129, 2014, 1309.5311.
- [141] Yuri L. Dokshitzer, G.D. Leder, S. Moretti, and B.R. Webber. Better jet clustering algorithms. *JHEP*, 9708:001, 1997, hep-ph/9707323.
- [142] M. Wobisch and T. Wengler. Hadronization corrections to jet cross-sections in deep inelastic scattering. 1998, hep-ph/9907280.
- [143] Stephen D. Ellis and Davison E. Soper. Successive combination jet algorithm for hadron collisions. *Phys.Rev.*, D48:3160–3166, 1993, hep-ph/9305266.
- [144] Matteo Cacciari, Gavin P. Salam, and Gregory Soyez. The Anti- $k(t)$  jet clustering algorithm. *JHEP*, 0804:063, 2008, 0802.1189.
- [145] Matteo Cacciari, Gavin P. Salam, and Gregory Soyez. FastJet User Manual. *Eur.Phys.J.*, C72:1896, 2012, 1111.6097.
- [146] A. Abulencia et al. Measurement of the Inclusive Jet Cross Section using the  $\mathbf{k_T}$  algorithm in  $p\bar{p}$  Collisions at  $\sqrt{s} = 1.96$  TeV with the CDF II Detector. *Phys.Rev.*, D75:092006, 2007, hep-ex/0701051.
- [147] T. Aaltonen et al. Measurement of the Inclusive Jet Cross Section at the Fermilab Tevatron p anti-p Collider Using a Cone-Based Jet Algorithm. *Phys.Rev.*, D78:052006, 2008, 0807.2204.
- [148] T. Aaltonen et al. Search for new particles decaying into dijets in proton-antiproton collisions at  $s^{**}(1/2) = 1.96$ -TeV. *Phys.Rev.*, D79:112002, 2009, 0812.4036.
- [149] V.M. Abazov et al. Measurement of the inclusive jet cross-section in  $p\bar{p}$  collisions at  $s^{91/2)} = 1.96$ -TeV. *Phys.Rev.Lett.*, 101:062001, 2008, 0802.2400.

- [150] V.M. Abazov et al. Measurement of the dijet invariant mass cross section in  $p\bar{p}$  collisions at  $\sqrt{s} = 1.96$  TeV. *Phys.Lett.*, B693:531–538, 2010, 1002.4594.
- [151] L. Bourhis, M. Fontannaz, and J.P. Guillet. Quarks and gluon fragmentation functions into photons. *Eur.Phys.J.*, C2:529–537, 1998, hep-ph/9704447.
- [152] M. Gluck, E. Reya, and A. Vogt. Parton fragmentation into photons beyond the leading order. *Phys.Rev.*, D48:116, 1993.
- [153] David d'Enterria and Juan Rojo. Quantitative constraints on the gluon distribution function in the proton from collider isolated-photon data. *Nucl.Phys.*, B860:311–338, 2012.
- [154] Stefano Frixione. Isolated photons in perturbative QCD. *Phys.Lett.*, B429:369–374, 1998, hep-ph/9801442.
- [155] J. Owens". Large-momentum-transfer production of direct photons, jets, and particles. *Reviews of Modern Physics*, 59(2):465–503, 1987.
- [156] P. Aurenche, R. Baier, M. Fontannaz, J.F. Owens, and M. Werlen. The Gluon Contents of the Nucleon Probed with Real and Virtual Photons. *Phys.Rev.*, D39:3275, 1989.
- [157] S. Catani, M. Fontannaz, J.P. Guillet, and E. Pilon. Cross-section of isolated prompt photons in hadron hadron collisions. *JHEP*, 0205:028, 2002, hep-ph/0204023.
- [158] Patrick Aurenche, Michel Fontannaz, Jean-Philippe Guillet, Eric Pilon, and Monique Werlen. A New critical study of photon production in hadronic collisions. *Phys.Rev.*, D73:094007, 2006, hep-ph/0602133.

- [159] Z. Belghobsi, M. Fontannaz, J.-Ph. Guillet, G. Heinrich, E. Pilon, et al. Photon - Jet Correlations and Constraints on Fragmentation Functions. *Phys.Rev.*, D79:114024, 2009, 0903.4834.
- [160] C. Albajar et al. Direct Photon Production at the CERN Proton - anti-Proton Collider. *Phys.Lett.*, B209:385–396, 1988.
- [161] J. Alitti et al. A Measurement of single and double prompt photon production at the CERN  $\bar{p}p$  collider. *Phys.Lett.*, B288:386–394, 1992.
- [162] R. Ansari et al. Direct Photon Production in  $\bar{p}p$  Collisions at  $\sqrt{s} = 630$ -GeV. *Z.Phys.*, C41:395, 1988.
- [163] S.S. Adler et al. Measurement of direct photon production in p + p collisions at  $s^{**}(1/2) = 200$ -GeV. *Phys.Rev.Lett.*, 98:012002, 2007, hep-ex/0609031.
- [164] T. Aaltonen et al. Measurement of the Inclusive Isolated Prompt Photon Cross Section in p anti-p Collisions at  $s^{**}(1/2) = 1.96$ -TeV using the CDF Detector. *Phys.Rev.*, D80:111106, 2009, 0910.3623.
- [165] V.M. Abazov et al. Measurement of the isolated photon cross section in  $p\bar{p}$  collisions at  $\sqrt{s} = 1.96$ -TeV. *Phys.Lett.*, B639:151–158, 2006, hep-ex/0511054.
- [166] F. Abe et al. A Precision measurement of the prompt photon cross-section in  $p\bar{p}$  collisions at  $\sqrt{s} = 1.8$  TeV. *Phys.Rev.Lett.*, 73:2662–2666, 1994.
- [167] D. Acosta et al. Comparison of the isolated direct photon cross sections in  $p\bar{p}$  collisions at  $\sqrt{s} = 1.8$ -TeV and  $\sqrt{s} = 0.63$ -TeV. *Phys.Rev.*, D65:112003, 2002, hep-ex/0201004.
- [168] D. Acosta et al. Direct photon cross section with conversions at CDF. *Phys.Rev.*, D70:074008, 2004, hep-ex/0404022.

- [169] S. Abachi et al. Isolated photon cross-section in the central and forward rapidity regions in  $p\bar{p}$  collisions at  $\sqrt{s} = 1.8$  TeV. *Phys.Rev.Lett.*, 77:5011–5015, 1996, hep-ex/9603006.
- [170] V.M. Abazov et al. The ratio of the isolated photon cross sections at  $\sqrt{s} = 630$  GeV and 1800 GeV. *Phys.Rev.Lett.*, 87:251805, 2001, hep-ex/0106026.
- [171] B. Abbott et al. The isolated photon cross-section in  $p\bar{p}$  collisions at  $\sqrt{s} = 1.8$  TeV. *Phys.Rev.Lett.*, 84:2786–2791, 2000, hep-ex/9912017.
- [172] Michał Czakon, Paul Fiedler, and Alexander Mitov. Total Top-Quark Pair-Production Cross Section at Hadron Colliders Through  $O(\frac{4}{S})$ . *Phys.Rev.Lett.*, 110(25):252004, 2013, 1303.6254.
- [173] Peter Bärnreuther, Michal Czakon, and Alexander Mitov. Percent Level Precision Physics at the Tevatron: First Genuine NNLO QCD Corrections to  $q\bar{q} \rightarrow t\bar{t} + X$ . *Phys.Rev.Lett.*, 109:132001, 2012, 1204.5201.
- [174] Michal Czakon and Alexander Mitov. NNLO corrections to top-pair production at hadron colliders: the all-fermionic scattering channels. *JHEP*, 1212:054, 2012, 1207.0236.
- [175] Michal Czakon and Alexander Mitov. NNLO corrections to top pair production at hadron colliders: the quark-gluon reaction. *JHEP*, 1301:080, 2013, 1210.6832.
- [176] T. Aaltonen et al. Tevatron electroweak working group. Combination of the  $t\bar{t}$  production cross section measurements from the tevatron collider. D0-Note-6363, 2012.
- [177] A. De Roeck and R. S. Thorne. Structure Functions. 2011, 1103.0555.

- [178] Stefano Forte, Lluis Garrido, Jose I. Latorre, and Andrea Piccione. Neural network parametrization of deep-inelastic structure functions. *JHEP*, 05:062, 2002, hep-ph/0204232.
- [179] Luigi Del Debbio, Stefano Forte, Jose I. Latorre, Andrea Piccione, and Joan Rojo. Unbiased determination of the proton structure function  $F_2(p)$  with faithful uncertainty estimation. *JHEP*, 03:080, 2005, hep-ph/0501067.
- [180] Luigi Del Debbio, Stefano Forte, Jose I. Latorre, Andrea Piccione, and Joan Rojo. Neural network determination of parton distributions: the nonsinglet case. *JHEP*, 03:039, 2007, hep-ph/0701127.
- [181] Richard D. Ball et al. Parton Distributions: Determining Probabilities in a Space of Functions. 2011, 1110.1863.
- [182] D. Stump, J. Pumplin, R. Brock, D. Casey, J. Huston, et al. Uncertainties of predictions from parton distribution functions. 1. The Lagrange multiplier method. *Phys.Rev.*, D65:014012, 2001, hep-ph/0101051.
- [183] J. Pumplin, D.R. Stump, J. Huston, H.L. Lai, Pavel M. Nadolsky, et al. New generation of parton distributions with uncertainties from global QCD analysis. *JHEP*, 0207:012, 2002, hep-ph/0201195.
- [184] Richard D. Ball, Stefano Carrazza, Luigi Del Debbio, Stefano Forte, Jun Gao, et al. Parton Distribution Benchmarking with LHC Data. *JHEP*, 1304:125, 2013, 1211.5142.
- [185] G. D'Agostini. On the use of the covariance matrix to fit correlated data. *Nucl.Instrum.Meth.*, A346:306–311, 1994.
- [186] Richard D. Ball et al. Fitting Parton Distribution Data with Multiplicative Normalization Uncertainties. *JHEP*, 1005:075, 2010, 0912.2276.

- [187] F James and M Roos. Minuit: Function minimization and error analysis, cern program library long writeup d506. 1994.
- [188] K Levenberg. A method for the solution of certain non-linear problems in least squares. *Quart. Appl. Math.*, 2(164), 1944.
- [189] D W Marquardt. An algorithm for least-squares estimation of nonlinear parameters. *SIAM J. Appl. Math.*, 11(431), 1963.
- [190] C.M. Bishop. *Neural networks for Pattern Recognition*. Oxford University Press, 1995.
- [191] D. J.; Luik A. I. Tetko, I. V.; Livingstone. Neural network studies. 1. comparison of overfitting and overtraining. *J. Chem. Inf. Comput. Sci.*, 35:826–833, 1995.
- [192] Sergey Alekhin. Parton distributions from deep-inelastic scattering data. *Phys. Rev.*, D68:014002, 2003, hep-ph/0211096.
- [193] J. Pumplin, D. R. Stump, and W. K. Tung. Multivariate fitting and the error matrix in global analysis of data. *Phys. Rev.*, D65:014011, 2001, hep-ph/0008191.
- [194] J. Pumplin et al. Uncertainties of predictions from parton distribution functions. 2. The Hessian method. *Phys. Rev.*, D65:014013, 2001, hep-ph/0101032.
- [195] A.D. Martin, R.G. Roberts, W.J. Stirling, and R.S. Thorne. Uncertainties of predictions from parton distributions. 1: Experimental errors. *Eur.Phys.J.*, C28:455–473, 2003, hep-ph/0211080.
- [196] Walter T. Giele and Stephane Keller. Implications of hadron collider observables on parton distribution function uncertainties. *Phys. Rev.*, D58:094023, 1998, hep-ph/9803393.

- [197] Walter T. Giele, Stephane A. Keller, and David A. Kosower. Parton distribution function uncertainties. 2001, hep-ph/0104052.
- [198] G. Watt and R.S. Thorne. Study of Monte Carlo approach to experimental uncertainty propagation with MSTW 2008 PDFs. *JHEP*, 1208:052, 2012, 1205.4024.
- [199] G. Watt. Parton distribution function dependence of benchmark Standard Model total cross sections at the 7 TeV LHC. *JHEP*, 1109:069, 2011, 1106.5788.
- [200] M. Dittmar, S. Forte, A. Glazov, S. Moch, G. Altarelli, et al. Parton Distributions. 2009, 0901.2504.
- [201] Richard D. Ball et al. Reweighting NNPDFs: the W lepton asymmetry. *Nucl. Phys.*, B849:112–143, 2011, 1012.0836.
- [202] Tancredi Carli et al. A posteriori inclusion of parton density functions in NLO QCD final-state calculations at hadron colliders: The APPLGRID Project. *Eur. Phys. J.*, C66:503–524, 2010, 0911.2985.
- [203] T. Kluge, K. Rabbertz, and M. Wobisch. fastnlo: Fast pqcd calculations for pdf fits. pages 483–486, 2006, hep-ph/0609285. Tevatron-for-LHC Workshop (2006).
- [204] Intel Corp. Intel® 64 and ia-32 architectures software developer manual.
- [205] John M. Campbell and R. Keith Ellis. Radiative corrections to Z b anti-b production. *Phys.Rev.*, D62:114012, 2000, hep-ph/0006304.
- [206] Charalampos Anastasiou, Ruth Britto, Bo Feng, Zoltan Kunszt, and Pierpaolo Mastrolia. D-dimensional unitarity cut method. *Phys.Lett.*, B645:213–216, 2007, hep-ph/0609191.

- [207] C.F. Berger, Z. Bern, L.J. Dixon, F. Febres Cordero, D. Forde, et al. An Automated Implementation of On-Shell Methods for One-Loop Amplitudes. *Phys.Rev.*, D78:036003, 2008, 0803.4180.
- [208] Ansgar Denner and S. Dittmaier. Reduction schemes for one-loop tensor integrals. *Nucl.Phys.*, B734:62–115, 2006, hep-ph/0509141.
- [209] R. Keith Ellis, W.T. Giele, and Z. Kunszt. A Numerical Unitarity Formalism for Evaluating One-Loop Amplitudes. *JHEP*, 0803:003, 2008, 0708.2398.
- [210] Walter T. Giele, Zoltan Kunszt, and Kirill Melnikov. Full one-loop amplitudes from tree amplitudes. *JHEP*, 0804:049, 2008, 0801.2237.
- [211] Giovanni Ossola, Costas G. Papadopoulos, and Roberto Pittau. Reducing full one-loop amplitudes to scalar integrals at the integrand level. *Nucl.Phys.*, B763:147–169, 2007, hep-ph/0609007.
- [212] Johan Alwall, Michel Herquet, Fabio Maltoni, Olivier Mattelaer, and Tim Stelzer. MadGraph 5 : Going Beyond. *JHEP*, 1106:128, 2011, 1106.0522.
- [213] G. Bevilacqua, M. Czakon, M.V. Garzelli, A. van Hameren, A. Kardos, et al. HELAC-NLO. *Comput.Phys.Commun.*, 184:986–997, 2013, 1110.1499.
- [214] Tanju Gleisberg, Stefan Hoeche, Frank Krauss, Andreas Schalicke, Steffen Schumann, et al. SHERPA 1. alpha: A Proof of concept version. *JHEP*, 0402:056, 2004, hep-ph/0311263.
- [215] T. Gleisberg, Stefan. Hoeche, F. Krauss, M. Schonherr, S. Schumann, et al. Event generation with SHERPA 1.1. *JHEP*, 0902:007, 2009, 0811.4622.
- [216] Luigi Del Debbio, Nathan P. Hartland, and Steffen Schumann. MCgrid: projecting cross section calculations on grids. *Comput.Phys.Commun.*, 185:2115–2126, 2014, 1312.4460.

- [217] Matt Dobbs and Jorgen Beck Hansen. The HepMC C++ Monte Carlo event record for High Energy Physics. *Comput.Phys.Commun.*, 134:41–46, 2001.
- [218] S. Catani and M.H. Seymour. A General algorithm for calculating jet cross-sections in NLO QCD. *Nucl.Phys.*, B485:291–419, 1997, hep-ph/9605323.
- [219] S. Frixione, Z. Kunszt, and A. Signer. Three jet cross-sections to next-to-leading order. *Nucl.Phys.*, B467:399–442, 1996, hep-ph/9512328.
- [220] S. Frixione. A General approach to jet cross-sections in QCD. *Nucl.Phys.*, B507:295–314, 1997, hep-ph/9706545.
- [221] Z. Bern, L.J. Dixon, F. Febres Cordero, S. Höche, H. Ita, et al. Ntuples for NLO Events at Hadron Colliders. *Comput.Phys.Commun.*, 185:1443–1460, 2014, 1310.7439.
- [222] Inclusive jets and dijets in LHCb. *LHCb-CONF-2011-015*, 2011.
- [223] G. Aad et al. Measurement of inclusive jet and dijet cross sections in proton-proton collisions at 7 TeV centre-of-mass energy with the ATLAS detector. *Eur.Phys.J.*, C71:1512, 2011, 1009.5908.
- [224] Georges Aad et al. Measurement of inclusive jet and dijet production in  $pp$  collisions at  $\sqrt{s} = 7$  TeV using the ATLAS detector. *Phys.Rev.*, D86:014022, 2012, 1112.6297.
- [225] Georges Aad et al. Measurement of the inclusive jet cross section in  $pp$  collisions at  $\text{sqrt}(s)=2.76$  TeV and comparison to the inclusive jet cross section at  $\text{sqrt}(s)=7$  TeV using the ATLAS detector. *Eur.Phys.J.*, C73:2509, 2013, 1304.4739.
- [226] Serguei Chatrchyan et al. Measurement of the Inclusive Jet Cross Section in  $pp$  Collisions at  $\sqrt{s} = 7$  TeV. *Phys.Rev.Lett.*, 107:132001, 2011, 1106.0208.

- [227] Serguei Chatrchyan et al. Measurement of the inclusive production cross sections for forward jets and for dijet events with one forward and one central jet in  $pp$  collisions at  $\sqrt{s} = 7$  TeV. *JHEP*, 1206:036, 2012, 1202.0704.
- [228] Serguei Chatrchyan et al. Measurements of differential jet cross sections in proton-proton collisions at  $\sqrt{s} = 7$  TeV with the CMS detector. *Phys.Rev.*, D87(11):112002, 2013, 1212.6660.
- [229] Serguei Chatrchyan et al. Measurement of the Rapidity and Transverse Momentum Distributions of  $Z$  Bosons in  $pp$  Collisions at  $\sqrt{s} = 7$  TeV. *Phys.Rev.*, D85:032002, 2012, 1110.4973.
- [230] Measurement of the transverse momentum distributions of  $Z$  Bosons decaying to dimuons in  $pp$  collisions at  $\text{sqrt}(s)=8$  TeV. Technical Report CMS-PAS-SMP-12-025, CERN, Geneva, 2013.
- [231] Measurement of  $Z$  production as a function of pT, Y. Technical Report CMS-PAS-SMP-13-013, CERN, Geneva, 2014.
- [232] Serguei Chatrchyan et al. Measurement of the lepton charge asymmetry in inclusive  $W$  production in  $pp$  collisions at  $\sqrt{s} = 7$  TeV. *JHEP*, 1104:050, 2011, 1103.3470.
- [233] Serguei Chatrchyan et al. Measurement of the electron charge asymmetry in inclusive  $W$  production in  $pp$  collisions at  $\sqrt{s} = 7$  TeV. *Phys.Rev.Lett.*, 109:111806, 2012, 1206.2598.
- [234] Serguei Chatrchyan et al. Measurement of the muon charge asymmetry in inclusive  $pp \rightarrow W + X$  production at  $\sqrt{s}=7$  TeV and an improved determination of light parton distribution functions. 2013, 1312.6283.

- [235] Georges Aad et al. Measurement of the Muon Charge Asymmetry from W Bosons Produced in pp Collisions at  $\sqrt{s} = 7$  TeV with the ATLAS detector. *Phys.Lett.*, B701:31–49, 2011, 1103.2929.
- [236] Georges Aad et al. Measurement of the transverse momentum distribution of Z/gamma\* bosons in proton-proton collisions at  $\sqrt{s} = 7$  TeV with the ATLAS detector. *Phys.Lett.*, B705:415–434, 2011, 1107.2381.
- [237] Georges Aad et al. Measurement of the Transverse Momentum Distribution of W Bosons in  $pp$  Collisions at  $\sqrt{s} = 7$  TeV with the ATLAS Detector. *Phys.Rev.*, D85:012005, 2012, 1108.6308.
- [238] Georges Aad et al. Measurement of the inclusive  $W^\pm$  and Z/gamma cross sections in the electron and muon decay channels in  $pp$  collisions at  $\sqrt{s} = 7$  TeV with the ATLAS detector. *Phys.Rev.*, D85:072004, 2012, 1109.5141.
- [239] R Aaij et al. Inclusive  $W$  and  $Z$  production in the forward region at  $\sqrt{s} = 7$  TeV. *JHEP*, 1206:058, 2012, 1204.1620.
- [240] Georges Aad et al. Measurement of the inclusive isolated prompt photon cross-section in  $pp$  collisions at  $\sqrt{s} = 7$  TeV using 35 pb-1 of ATLAS data. *Phys.Lett.*, B706:150–167, 2011, 1108.0253.
- [241] Georges Aad et al. Measurement of the production cross section of an isolated photon associated with jets in proton-proton collisions at  $\sqrt{s} = 7$  TeV with the ATLAS detector. *Phys.Rev.*, D85:092014, 2012, 1203.3161.
- [242] Serguei Chatrchyan et al. Measurement of the Differential Cross Section for Isolated Prompt Photon Production in pp Collisions at 7 TeV. *Phys.Rev.*, D84:052011, 2011, 1108.2044.

- [243] Measurement of the top quark pair production cross section in the single-lepton channel with ATLAS in proton-proton collisions at 8 TeV using kinematic fits with b-tagging. 2012.
- [244] Statistical combination of top quark pair production cross-section measurements using dilepton, single-lepton, and all-hadronic final states at  $s = 7$  TeV with the ATLAS detector. 2012.
- [245] Serguei Chatrchyan et al. Measurement of the  $t\bar{t}$  production cross section in the dilepton channel in  $pp$  collisions at  $\sqrt{s} = 7$  TeV. *JHEP*, 1211:067, 2012, 1208.2671.
- [246] CMS Collaboration. Top pair cross section in e/mu+jets at 8 TeV. 2012.
- [247] Michal Czakon and Alexander Mitov. Top++: A Program for the Calculation of the Top-Pair Cross-Section at Hadron Colliders. 2011, 1112.5675.
- [248] Michal Czakon, Michelangelo L. Mangano, Alexander Mitov, and Juan Rojo. Constraints on the gluon PDF from top quark pair production at hadron colliders. *JHEP*, 1307:167, 2013, 1303.7215.
- [249] Stefano Catani and Massimiliano Grazzini. An NNLO subtraction formalism in hadron collisions and its application to Higgs boson production at the LHC. *Phys.Rev.Lett.*, 98:222002, 2007, hep-ph/0703012.
- [250] H.L. Lai, Pavel M. Nadolsky, J. Pumplin, D. Stump, W.K. Tung, et al. The Strange parton distribution of the nucleon: Global analysis and applications. *JHEP*, 0704:089, 2007, hep-ph/0702268.
- [251] S. Alekhin, Sergey A. Kulagin, and R. Petti. Determination of Strange Sea Distributions from Neutrino-Nucleon Deep Inelastic Scattering. *Phys.Lett.*, B675:433–440, 2009, 0812.4448.

- [252] Georges Aad et al. Determination of the strange quark density of the proton from ATLAS measurements of the  $W \rightarrow \ell\nu$  and  $Z \rightarrow \ell\ell$  cross sections. *Phys.Rev.Lett.*, 109:012001, 2012, 1203.4051.
- [253] S. Alekhin, J. Bluemlein, L. Caminadac, K. Lipka, K. Lohwasser, et al. Determination of Strange Sea Quark Distributions from Fixed-target and Collider Data. 2014, 1404.6469.
- [254] Serguei Chatrchyan et al. Measurement of associated  $W +$  charm production in pp collisions at  $\sqrt{s} = 7$  TeV. *JHEP*, 1402:013, 2014, 1310.1138.
- [255] Charalampos Anastasiou, Stephan Buehler, Franz Herzog, and Achilleas Lazopoulos. Total cross-section for Higgs boson hadroproduction with anomalous Standard Model interactions. *JHEP*, 1112:058, 2011, 1107.0683.
- [256] OpenMP Architecture Review Board. OpenMP application program interface version 3.0. <http://www.openmp.org/mp-documents/spec30.pdf>, May 2008.
- [257] W. Giele, E.W. Nigel Glover, I. Hinchliffe, J. Huston, Eric Laenen, et al. The QCD / SM working group: Summary report. pages 275–426, 2002, hep-ph/0204316.
- [258] H.L. Lai et al. Global QCD analysis of parton structure of the nucleon: CTEQ5 parton distributions. *Eur.Phys.J.*, C12:375–392, 2000, hep-ph/9903282.
- [259] David J. Montana and Lawrence Davis. Training feedforward neural networks using genetic algorithms. In *Proceedings of the 11th International Joint Conference on Artificial Intelligence - Volume 1*, IJCAI’89, pages 762–767, San Francisco, CA, USA, 1989. Morgan Kaufmann Publishers Inc.

- [260] J. Beringer et al. Review of Particle Physics (RPP). *Phys.Rev.*, D86:010001, 2012.

# **Appendix A**

## **Summary of experimental data**

In this appendix, the experimental measurements discussed in this thesis are summarised. For each experiment, the underlying process and physical observable measured are specified, along with a brief summary of the PDF flavours and combinations targeted by the data.

Fixed-Target Deep Inelastic Scattering				
Process	Experiment	Obs.	Ref.	Target
$\mu p/d \rightarrow \mu X$	BCDMS NMC Fermilab E665	$F_2^p, F_2^d$	[?]	$q, \bar{q}$
$\mu p/d \rightarrow \mu X$		$F_2^p, F_2^d/F_2^p$	[?]	$q, \bar{q}, d/u$
$\mu p/d \rightarrow \mu X$		$F_2^p, F_2^d$	[?]	$q, \bar{q}$
$\mu p \rightarrow \mu X$	BCDMS NMC SLAC	$F_L$	[?]	$g$
$\mu p/d \rightarrow \mu X$		$F_L$	[?]	$g$
$ep/d \rightarrow eX$		$F_L$	[?]	$g$

Table A.1: Summary of discussed Fixed-Target DIS experiments, arranged as in Table ???. Here deuteron and proton structure function data is summarised.

HERA Deep Inelastic Scattering				
Process	Experiment	Obs.	Ref.	Target
$ep \rightarrow eX$	H1 ZEUS	NC $\sigma$	[?, ?, ?]	$g, q, \bar{q}$
$ep \rightarrow eX$		NC $\sigma$	[?, ?, ?, ?]	$g, q, \bar{q}$
$e^+p \rightarrow \bar{\nu}X$	H1 ZEUS	CC $\sigma$	[?]	$d + s, \bar{u}$
$e^+p \rightarrow \bar{\nu}X$		CC $\sigma$	[?]	$d + s, \bar{u}$
$ep \rightarrow eX + c$	H1 ZEUS	$F_c^2$	[?, ?, ?, ?]	$g$
$ep \rightarrow eX + c$		$F_c^2$	[?, ?, ?]	$g$
$ep \rightarrow eX$	H1 ZEUS	$F_L$	[?]	$g$
$ep \rightarrow eX$		$F_L$	[?]	$g$
$ep \rightarrow eX$	HERA-I HERA-I	NC $\sigma$	[?]	$g, q, \bar{q}$
$ep \rightarrow \nu X$		CC $\sigma$	[?]	$q, \bar{q}$
$ep \rightarrow eX + c$	HERA-I	$F_2^c$	[?]	$g$

Table A.2: Summary of discussed HERA DIS measurements, arranged as in Table ???.

Neutrino Deep Inelastic Scattering				
Process	Experiment	Obs.	Ref.	Target
$\nu(\bar{\nu}) \text{Fe} \rightarrow \mu X$	NuTeV	$F_2/F_3$	[?]	$q, \bar{q}$
$\nu(\bar{\nu}) \text{Pb} \rightarrow \mu X$	CHORUS	$F_2/F_3$	[?]	$q, \bar{q}$
$\nu(\bar{\nu}) \text{Fe} \rightarrow \mu^+ \mu^- X$	NuTeV/CCFR	Dimuon $\sigma$	[?]	$s, \bar{s}$

Table A.3: Summary of discussed Neutrino DIS measurements, arranged as in Table ???.

Drell-Yan				
Process	Experiment	Obs.	Ref.	Target
$p \text{ Cu} \rightarrow \mu^+ \mu^-$	Fermilab E605	$\sigma^{pp}$	[?]	$q + \bar{q}$
$p \text{ H} \rightarrow \mu^+ \mu^-$	NuSea/E866	$\sigma^{pp}$	[?]	$q + \bar{q}$
$p \text{ D/H} \rightarrow \mu^+ \mu^-$	NuSea/E866	$\sigma^{pd}/\sigma^{pp}$	[?]	$\bar{d}/\bar{u}$
$p\bar{p} \rightarrow e^+ e^-$	D0	$Z/\gamma y$	[?]	$u, d$
$p\bar{p} \rightarrow e\nu$	CDF	$W$ asym.	[?]	$u - d$
$p\bar{p} \rightarrow \mu\nu$	D0	$W$ asym.	[?]	$u - d$
$p\bar{p} \rightarrow e\nu$	D0	$W$ asym.	[?]	$u - d$
$pp \rightarrow l^+ l^-$	CMS	$Z$ $p_T/y$	[?]	$u + \bar{u}, d + \bar{d}$
$pp \rightarrow \mu^+ \mu^-$	CMS	$Z$ $p_T/y$	[?, ?]	$u + \bar{u}, d + \bar{d}$
$pp \rightarrow l\nu$	CMS	$W$ asym.	[?]	$u - \bar{d}$
$pp \rightarrow \mu\nu$	CMS	$W$ asym.	[?, ?]	$u - \bar{d}$
$pp \rightarrow \mu\nu$	ATLAS	$W$ asym.	[?]	$u - \bar{d}$
$pp \rightarrow \mu\mu$	ATLAS	$Z/\gamma p_T.$	[?]	$u + \bar{u}, d + \bar{d}$
$pp \rightarrow l\nu$	ATLAS	$W p_T.$	[?]	$u + \bar{d}, \bar{u} + d$
$pp \rightarrow \mu\nu$	LHCb	$W p_T.$	[?]	$u + d, \bar{u} + \bar{d}$
$pp \rightarrow \mu\mu$	LHCb	$Z/\gamma p_T.$	[?]	$u + \bar{u}, d + \bar{d}$

Table A.4: Summary of discussed Drell-Yan measurements, arranged as in Table ???. Here Fixed-Target experiments are shown in the higher segment, and collider experiments in the lower two.

Jet Production				
Process	Experiment	Obs.	Ref.	Target
$p\bar{p} \rightarrow j + X$	CDF	Inclusive Jets	[?, ?]	$g$
$p\bar{p} \rightarrow j + X$	D0	Inclusive Jets	[?]	$g$
$p\bar{p} \rightarrow jj + X$	CDF	Dijets	[?]	$g$
$p\bar{p} \rightarrow jj + X$	D0	Dijets	[?]	$g$
$pp \rightarrow j + X$	LHCb	Inclusive Jets	[?]	$g$
$pp \rightarrow j + X$	ATLAS	Inclusive Jets	[?, ?, ?]	$g$
$pp \rightarrow j + X$	CMS	Inclusive Jets	[?, ?, ?]	$g$
$pp \rightarrow jj + X$	LHCb	Dijets	[?]	$g$
$pp \rightarrow jj + X$	ATLAS	Dijets	[?, ?]	$g$
$pp \rightarrow jj + X$	CMS	Dijets	[?, ?]	$g$

Table A.5: Summary of discussed inclusive jet and dijet measurements, arranged as in Table ??.

Prompt Photon				
Process	Experiment	Obs.	Ref.	Target
$p\bar{p} \rightarrow \gamma X$	UA1/UA2	Photon $E_T$	[?, ?, ?]	$q, g$
$p\bar{p} \rightarrow \gamma X$	CDF	Photon $p_T$	[?, ?, ?, ?, ?]	$q, g$
$p\bar{p} \rightarrow \gamma X$	D0	Photon $E_T$	[?, ?, ?]	$q, g$
$pp \rightarrow \gamma X$	PHENIX	Photon $p_T$	[?]	$g, q + \bar{q}$
$pp \rightarrow \gamma X$	ATLAS	Photon $E_T, \eta$	[?]	$g, q + \bar{q}$
$pp \rightarrow \gamma X + j$	ATLAS	Photon $E_T, \eta$	[?]	$g, q + \bar{q}$
$pp \rightarrow \gamma X$	CMS	Photon $E_T, \eta$	[?]	$g, q + \bar{q}$

Table A.6: Summary of discussed isolated prompt photon measurements. The process column denotes the reaction observed in each experiment, Obs. refers to the physical observable measured and Target illustrates the most relevant partonic channels for the process and observable in question.

Top production				
Process	Experiment	Obs.	Ref.	Target
$p\bar{p} \rightarrow t\bar{t}$	CDF + D0	$\sigma_{t\bar{t}}$	[?]	$q + \bar{q}$
$pp \rightarrow t\bar{t}$	ATLAS	$\sigma_{t\bar{t}}$	[?, ?]	$g$
$pp \rightarrow t\bar{t}$	CMS	$\sigma_{t\bar{t}}$	[?, ?]	$g$

Table A.7: Summary of discussed top production measurements, arranged as in Table ??.

## Appendix B

### Distance Estimators

Here we define a set of useful measures in determining the statistical differences between two sets of parton distributions in a Monte Carlo representation, first introduced in Ref. [?]. Recalling the standard definitions of the central value of a Monte Carlo PDF with  $N_{\text{rep}}$  replicas,

$$\langle f(x, Q^2) \rangle = \frac{1}{N_{\text{rep}}} \sum_i^{N_{\text{rep}}} f_k(x, Q^2), \quad (\text{B.1})$$

and its associated uncertainty,

$$\sigma^2 [f(x, Q^2)] = \frac{1}{N_{\text{rep}} - 1} \sum_i^{N_{\text{rep}}} (f_k(x, Q^2) - \langle f(x, Q^2) \rangle)^2. \quad (\text{B.2})$$

Further estimators are available [?] for the uncertainty upon the central value,

$$\sigma^2 [\langle f(x, Q^2) \rangle] = \frac{1}{N_{\text{rep}}} \sigma^2 [f(x, Q^2)], \quad (\text{B.3})$$

and the uncertainty upon the uncertainty,

$$\sigma^2 [\sigma^2 [f(x, Q^2)]] = \frac{1}{N_{\text{rep}}} \left[ m_4 [f(x, Q^2)] - \frac{N_{\text{rep}} - 3}{N_{\text{rep}} - 1} (\sigma^2 [f(x, Q^2)])^2 \right], \quad (\text{B.4})$$

where  $m_4[f(x, Q^2)]$  refers to the fourth moment of the distribution  $f(x, Q^2)$ .

Given these quantities we can define a distance between the representation of the PDF  $f$  in two PDF sets as the square difference of the PDF central values in units of the uncertainty of the mean,

$$d_{\text{CV}}^2[f^{(1)}, f^{(2)}] = \frac{(\langle f^{(1)} \rangle - \langle f^{(2)} \rangle)^2}{\sigma^2[\langle f^{(1)} \rangle] + \sigma^2[\langle f^{(2)} \rangle]}, \quad (\text{B.5})$$

where the PDF superscripts enumerate the PDF sets being compared and the dependence upon the kinematical variables  $x, Q^2$  is implicit. With this definition of PDF distance, a value of  $d^2 = 1$  corresponds to a discrepancy between PDF sets consistent with one standard deviation of the central values. A similar measure can be defined for the uncertainties of the distribution,

$$d_\sigma^2[f^{(1)}, f^{(2)}] = \frac{(\sigma^2[f^{(1)}] - \sigma^2[f^{(2)}])^2}{\sigma^2[\sigma^2[f^{(1)}]] + \sigma^2[\sigma^2[f^{(2)}]]}. \quad (\text{B.6})$$

These distances quantities are particularly useful in the systematic comparison of all partons in two PDF sets in order to evaluate the size and statistical significance of differences between the two sets. As an example, consider Figure ?? where distances are shown for both estimators  $d_{\text{CV}}$  and  $d_\sigma$  (i.e the square-root of Eqns. ??, ??) between two PDF sets, for seven PDF combinations.

A distinction should be noted between the distances presented in this work and those defined in Ref. [?], where an additional bootstrap sampling of the distributions was used. In this work all distances are presented exactly as in Eqn. ?? and ??.

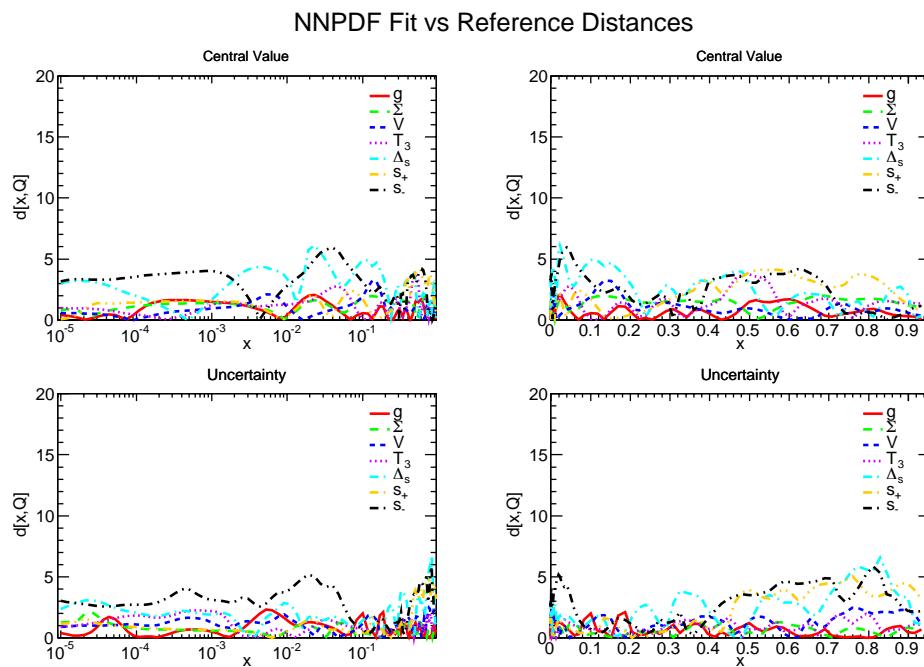


Figure B.1: Example of both PDF central value and uncertainty distances for the seven PDFs parametrised at the NNPDF initial scale.

# Publications

- [1] L. Del Debbio, N. P. Hartland and S. Schumann, “MCgrid: projecting cross section calculations on grids,” *Comput. Phys. Commun.* **185** (2014) 2115 [arXiv:1312.4460 [hep-ph]].
- [2] R. D. Ball *et al.* [NNPDF Collaboration], “Parton distributions with QED corrections,” *Nucl. Phys. B* **877** (2013) 2, 290 arXiv:1308.0598 [hep-ph].
- [3] R. D. Ball, S. Carrazza, L. Del Debbio, S. Forte, J. Gao, N. Hartland, J. Huston and P. Nadolsky *et al.*, “Parton Distribution Benchmarking with LHC Data,” *JHEP* **1304** (2013) 125 arXiv:1211.5142 [hep-ph].
- [4] R. D. Ball, V. Bertone, S. Carrazza, C. S. Deans, L. Del Debbio, S. Forte, A. Guffanti and N. P. Hartland *et al.*, “Parton distributions with LHC data,” *Nucl. Phys. B* **867** (2013) 244 arXiv:1207.1303 [hep-ph].
- [5] R. D. Ball, V. Bertone, F. Cerutti, L. Del Debbio, S. Forte, A. Guffanti, N. P. Hartland and J. I. Latorre *et al.*, “Reweighting and Unweighting of Parton Distributions and the LHC W lepton asymmetry data,” *Nucl. Phys. B* **855** (2012) 608 arXiv:1108.1758 [hep-ph].

# Proceedings

- [1] N. P. Hartland and C. S. Deans, “Towards closure testing of parton determinations,” arXiv:1307.2046 [hep-ph].
- [2] N. P. Hartland and E. R. Nocera, “A Mathematica interface to NNPDFs,” Nucl. Phys. Proc. Suppl. **234** (2013) 54 arXiv:1209.2585 [hep-ph].
- [3] N. Hartland, “LHC data and the proton strangeness,” arXiv:1205.3508 [hep-ph].
- [4] J. Alcaraz Maestre *et al.* [SM AND NLO MULTILEG and SM MC Working Groups Collaboration], arXiv:1203.6803 [hep-ph].
- [5] F. Cerutti and N. P. Hartland, J. Phys. Conf. Ser. **368** (2012) 012063 arXiv:1111.6768 [hep-ph].