



Parton Distributions based on a Maximally Consistent Dataset

Juan Rojo

Rudolf Peierls Centre for Theoretical Physics, 1 Keble Road,
University of Oxford, OX1 3NP Oxford, United Kingdom

Abstract

The choice of data that enters a global QCD analysis can have a substantial impact on the resulting parton distributions and their predictions for collider observables. One of the main reasons for this has to do with the possible presence of inconsistencies, either internal within an experiment or external between different experiments. In order to assess the robustness of the global fit, different definitions of a conservative PDF set, that is, a PDF set based on a maximally consistent dataset, have been introduced. However, these approaches are typically affected by theory biases in the selection of the dataset. In this contribution, after a brief overview of recent NNPDF developments, we propose a new, fully objective, definition of a *conservative* PDF set, based on the Bayesian reweighting approach. Using the new NNPDF3.0 framework, we produce various conservative sets, which turn out to be mutually in agreement within the respective PDF uncertainties, as well as with the global fit. We explore some of their implications for LHC phenomenology, finding also good consistency with the global fit result. These results provide a non-trivial validation test of the new NNPDF3.0 fitting methodology, and indicate that possible inconsistencies in the fitted dataset do not affect substantially the global fit PDFs.

1. Overview of NNPDF developments

The accurate determination of the parton distribution functions (PDFs) of the proton is one of the most important tasks for precision phenomenology at the LHC [1]. PDFs are one of the limiting factors for the precision of our theoretical predictions for Higgs boson production [2], since their uncertainties degrade the accuracy of the Higgs characterization in terms of its couplings; they induce large uncertainties in high-mass New Physics particle production [3]; and they affect Standard Model precision measurements such as the mass of the W boson [4].

Until recently, the most updated set from the NNPDF Collaboration was NNPDF2.3 [5], the first PDF to ever include LHC data from ATLAS, CMS and LHCb. The NNPDF2.3 sets have been used in a large number of phenomenological and experimental studies. In the benchmarking exercise of [6], NNPDF2.3 was com-

pared in great detail to other recent sets including CT10 and MSTW08.

The NNPDF2.3 sets can be accessed through the LHAPDF library, and they are also available as internal sets in various widely used codes. For instance, NNPDF2.3 is one of the internal PDF sets in the MADGRAPH5_AMC@NLO program for the automated computation of NLO cross-sections matched to parton showers [7]. The leading order version of NNPDF2.3 [8, 9] has also been implemented as internal set in the PYTHIA8 Monte Carlo event generator [10], where it has been used as the basis of the recent Monash 2013 Tune [11] of PYTHIA8. As compared to older tunes, the new Monash 2013 tune achieves an improved description of a wide range of collider data, including the recent forward data from the LHC, thanks partly to the steeper small- x gluon in NNPDF2.3LO.

In addition to the standard QCD PDF sets, we also recently produced the NNPDF2.3QED sets [12], which

supplement the QCD DGLAP evolution equations with the corresponding QED contributions (see [13] and references therein). In addition, we provided for the first time an unbiased determination of the photon PDF $\gamma(x, Q^2)$ from experimental data without any theory model assumptions. A precision determination of the photon PDF is relevant since photon-induced contributions can be comparable or even dominant with respect to standard quark-induced production in a number of crucial LHC processes like high-mass dilepton production [14] and WW production [15].

On the polarized side, the NNPDFpol1.1 set [16] was recently released, which is the first global polarized fit using the NNPDF methodology. NNPDFpol1.1 supplements all relevant inclusive polarized DIS data [17] with polarized hadron collider data from the STAR and PHENIX experiments at RHIC on jet and W boson production. Remarkably, we are able to find evidence for the first time a non-zero and positive polarization of the gluon in the proton.¹ The next steps will be using polarized semi-inclusive data to constrain the quark flavor separation, for which a new set of fragmentation functions using the NNPDF methodology is required.

2. The NNPDF3.0 sets

NNPDF3.0 is the new release from the NNPDF Collaboration. It is available in LHAPDF6 since version 6.1.4. As will be discussed in detail in an forthcoming publication, NNPDF3.0 is the first PDF set fully determined from a methodology validated by closure tests. These closure tests ensure that PDFs determined from pseudo-data generated from a known underlying law reproduce the statistical distribution of results expected on the basis of the assumed experimental uncertainties. An important consequence is that it can be demonstrated that methodological uncertainties are rather smaller than the standard theoretical and experimental uncertainties. The updated NNPDF3.0 fitting strategy allows to produce PDFs with different theory inputs and with widely different datasets with a unique consistent methodology, a robustness that has never been achieved in the traditional PDF fitting approach.

NNPDF3.0 is based on a global dataset, that includes all the relevant available experimental constraints on parton distributions. The NNPDF2.3 dataset has been supplemented with the complete HERA-II deep-inelastic cross-sections from H1 and ZEUS, the combined charm production data from HERA, jet production from ATLAS and CMS, vector boson rapidity and

p_T distributions from ATLAS, CMS and LHCb, $W + c$ data from CMS and top quark total cross sections from ATLAS and CMS. Some of these new LHC observables provide precious information on poorly-known PDFs. For instance, the top quark data provide useful constraints on the large- x gluon PDF [19, 20] and $W + c$ data allows to pin down strangeness [21, 22].

NNPDF3.0 uses state-of-the-art theoretical calculations for all collider processes included. At NLO, all LHC observables are computed without any approximation using suitable fast interfaces for NLO calculations: APPLGRID [23], FASTNLO [24] and AMCFast [25]. NLO calculations are then supplemented with NNLO K -factors and electroweak corrections when necessary, using TOP++ for top data [26] and FEWZ for electroweak production [27]. Jet data is included in the NNLO fits using the improved threshold approximation [28], validated with the exact NNLO calculation of the gluon-gluon channel [29], which allows to carefully select only those data points with kinematics for which the threshold approximation is close enough to the exact calculation [30].

In this contribution, I want to focus on a particular aspect of the NNPDF3.0 analysis, namely a new proposal for the definition, in a fully objective way, a set of parton distributions based on a maximally consistent dataset, in order to explore the possible impact of dataset inconsistencies in the global fit.

3. Parton distributions based on a maximally consistent dataset

The choice of data that enters a global QCD analysis has an important impact on the resulting parton distributions. One of the reasons for this is the potential presence of inconsistencies, either between the various datasets or internal within a given experiment. In order to bypass this problem, and to assess the robustness of the global fit, various definitions of a *conservative* PDF set have been introduced. However, these approaches are typically affected by theory bias since the selection of a maximally consistent datasets is done following an expectation of which experiments are more reliable. For example, the NNPDF2.3 collider-only fit [5] is based on the expectation that collider data are presumably more robust than fixed-target data, and the MRST2004 conservative partons [31] excluded various datasets which could be affected by large perturbative corrections.

Here we propose a new alternative definition of a conservative set of parton distributions. The main novelty is removing any theoretical bias that might affect the data selection of such maximally consistent dataset,

¹Consistent results are found by the DSSV collaboration [18].

	$\alpha_{\max} = 1.1$	$\alpha_{\max} = 1.2$	$\alpha_{\max} = 1.3$
NMC d/p	y	y	y
NMC	n	n	n
SLAC	n	n	y
BCDMS	n	y	y
CHORUS	n	y	y
NuTeV	y	y	y
HERA-I	y	y	y
ZEUS HERA-II	n	n	y
H1 HERA-II	n	n	n
HERA σ_{NC}^c	n	y	y
E886 d/p	y	y	y
E886 p	n	y	y
E605	y	y	y
CDF Z rapidity	n	n	n
CDF Run-II k_T jets	n	y	y
D0 Z rapidity	y	y	y
ATLAS W, Z 2010	n	y	y
ATLAS 7 TeV jets	y	y	y
ATLAS 2.76 TeV jets	y	y	y
ATLAS high-mass DY	n	n	n
ATLAS $W p_T$	y	y	y
CMS W electron asy	y	y	y
CMS W muon asy	n	n	n
CMS jets 2011	y	y	y
CMS $W + c$ total	n	n	n
CMS $W + c$ ratio	n	n	n
CMS 2D DY 2011	n	n	y
LHCb W, Z rapidity	n	y	y
$\sigma(t\bar{t})$	n	n	n

Table 1: Datasets included in the NNPDF3.0 global fit, indicating which of these experiments are included in the conservative partons fits for the different values of the threshold α_{\max} that have been used.

using the tools that the Bayesian reweighting framework provides [32, 33]. In this new approach, we start from the NNPDF3.0 NLO and NNLO global fits with $N_{\text{rep}} = 1000$ replicas, and compute for each replica the weight associated with each individual experiment included in the fit:

$$w_k = \frac{(\chi_k^2)^{\frac{1}{2}(n-1)} e^{-\frac{1}{2}\chi_k^2}}{\frac{1}{N_{\text{rep}}} \sum_{k=1}^{N_{\text{rep}}} (\chi_k^2)^{\frac{1}{2}(n-1)} e^{-\frac{1}{2}\chi_k^2}}, \quad (1)$$

with n the number of points of this experiment and χ_k^2 is the $t_0 \chi^2$ for replica k .

Next, we compute the $P(\alpha)$ distributions for each of the individual datasets. This distribution is a measure of the consistency of the various experiments with the global fit: the parameter α measures by how much experimental uncertainties should be rescaled to achieve perfect consistency. This distribution is defined by the fact that when we rescale the uncertainties of the data by a factor α , we can use inverse probability to calculate the probability density for the rescaling parameter

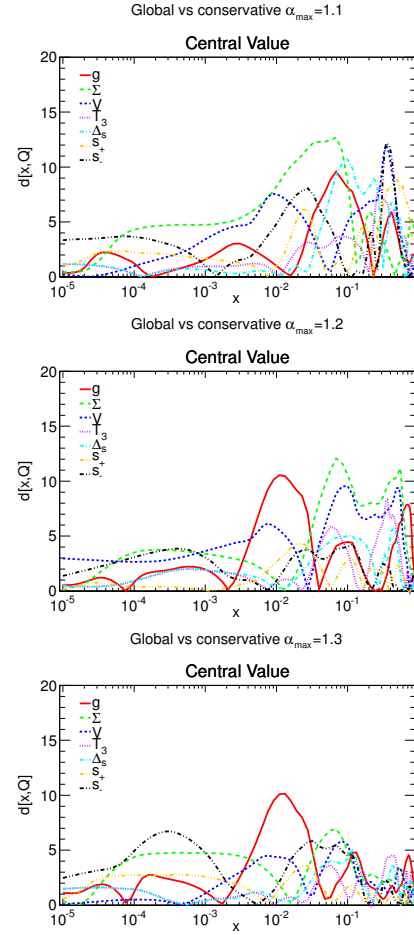


Figure 1: Distances between the global and conservative NNPDF3.0 NNLO partons for the three values of α_{\max} used. In all these comparisons the sets with $N_{\text{rep}} = 100$ replicas have been used.

α :

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

where $w_k(\alpha)$ are the weights Eq. (1) evaluated by replacing χ_k^2 with χ_k^2/α^2 . When $P(\alpha)$ peaks close to one, this particular dataset is consistent with the global fit, while if it peaks far above one, then it is likely that the errors in the data have been underestimated, or that the theoretical calculations are not accurate enough, for instance in the presence of large perturbative corrections.

This information suggest an objective criterion to select a maximally consistent dataset. For each of the experiment included in the global fit, we compute the mean, the median and the mode of the corresponding $P(\alpha)$ distribution. This experiment will be included in

the conservative set if at least two out of these three estimators are below some fixed threshold, denoted by α_{\max} . In order to separate effects from different datasets from perturbative uncertainties, we use the same data at NLO and NNLO, and we only include in the conservative partons experiments which satisfy the above criterion both at NLO and NNLO. There is some degree of arbitrariness in this criterion, in particular in the choice of α_{\max} , or by the requirements that NLO and NNLO datasets are the same in the conservative fits, which is compensated by its objective character.

The main difference between this and other criteria is that while previous conservative sets were based on an expectation of which data are more reliable, the new criterion is based on a measure of which data are consistent or inconsistent with the rest of the experiments in the global fit. Again, let me emphasize that this consistency can mean either consistency of the new data with the other datasets or internal self-consistency.

We have produced various sets of conservative partons, obtained for three different values of α_{\max} , namely 1.1, 1.2 and 1.3, corresponding to different degrees of tolerance about the inconsistencies that we admit in the fitted dataset. In Table 1 we list all the datasets included in the NNPDF3.0 global fit, and indicate which of these experiments are included in the conservative partons fit for the values of the threshold α_{\max} that have been used.

In order to gauge the fit quality of the various conservative PDF sets, in Table 2 we show the NLO and NNLO experimental χ^2 (see [6] for the precise definition) for both the global and the three conservative fits. The results are as expected: by increasing α_{\max} we interpolate between the maximally consistent dataset with $\chi^2 \sim 1$ and the global fit result. These results suggest that $\alpha_{\max} = 1.1$ could be used to define the best approximation to a conservative parton set.

	χ^2 NLO	χ^2 NNLO
$\alpha_{\max} = 1.1$	0.97	1.01
$\alpha_{\max} = 1.2$	1.06	1.09
$\alpha_{\max} = 1.3$	1.12	1.15
Global	1.23	1.28

Table 2: The total χ^2 per data point for the global and conservative fits for different values of the threshold α_{\max} .

To quantify how the various conservative partons differ from the global fit, we show in Fig. 1 the corresponding distances for $\alpha_{\max} = 1.1, 1.2$ and 1.3 . Let us recall that for $N_{\text{rep}} = 100$, a distance of order 10 correspond to PDFs that agree at the one-sigma level. As can be seen, for the three fits there is a nice consistency

with the global dataset, with PDFs differing at most at the one-sigma level. Of course PDF uncertainties are larger in the fits to reduced datasets, but the statistical compatibility with the global fit is very similar for three values of α_{\max} .

Next, we compare various NNLO PDFs in the global fit and in the two conservative fits with $\alpha_{\max} = 1.1$ and 1.2 , in Fig. 2, at a scale of $Q^2 = 2 \text{ GeV}^2$. As we expected from the distance comparison, there is a nice agreement between the different fits, typically at the one-sigma level or better. This is a non-trivial consistency check of both the global fit approach and of the definition of conservative partons that we propose here. The small- x gluon is similar in all cases because is driven by the HERA-I data. Larger differences are found in the quark sector, at medium and large- x and for the most conservative fit with $\alpha_{\max} = 1.1$. The region around $x \sim 0.01$ for the gluon, relevant for Higgs production in gluon fusion, is stable at the one-sigma level, consistent with the studies based on NNPDF2.3 in the Les Houches 2014 proceedings [34].

4. Implications for LHC phenomenology

As we have found, the conservative partons are nicely consistent with the global fit results for all values of the threshold α_{\max} , with of course rather larger PDF uncertainties, specially for small α_{\max} . To study what are the implications of these conservative partons for LHC phenomenology, in the following we have used MADGRAPH5_AMC@NLO [7] to compute NLO predictions for a variety of LHC observables both in the global fit and in the different conservative fit. Cross-sections are computed at 13 TeV with typical experimental analysis cuts in the final states. The results are summarized in Fig. 3, where we show the predictions for the various processes as ratios with respect to the NNPDF2.3 NLO global fit predictions.

Is clear that all the fits are consistent at the one-sigma level. The fits with smaller values of α_{\max} have larger PDF uncertainties and therefore larger fluctuations of the central values, as expected, while the fit with $\alpha_{\max} = 1.3$ is quite close to the global fit. The increase in PDF errors is quite noticeable, for example in cross-sections that depend on the large- x gluon such as $t\bar{t}$ and $h\bar{t}t$, as well as for those that depend on strangeness. On the other hand, the predictions for Higgs production in gluon fusion are rather stable with respect to the choice of dataset, as had already been observed in [34].

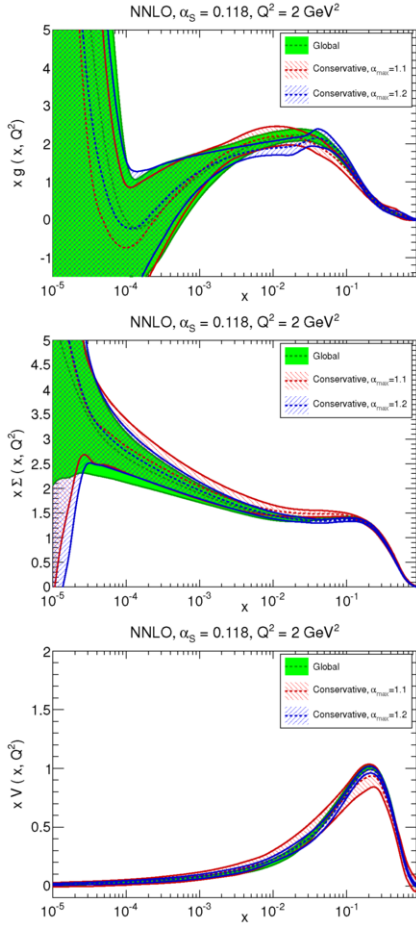


Figure 2: Comparison of various NNLO PDFs in the global fit and in the two conservative fits with $\alpha_{\max} = 1.1$ and 1.2 . The comparison is performed at a scale of $Q^2 = 2 \text{ GeV}^2$. From top to bottom, we have gluon, the singlet PDF and the total valence PDFs.

5. Outlook

In this contribution I have presented a new strategy to select a maximally consistent dataset for a PDF analysis. We have found that in this new approach there is a good consistency between fits to reduced datasets and the global fit, at the level of one-sigma PDF uncertainties or better. The corresponding predictions for LHC observables are also in reasonable agreement. Therefore, these conservative partons can be used in phenomenological analyses that aim to study how fits based on small but maximally consistent datasets affect LHC observables.

In addition to the release of conservative partons, it is conceivable that this method could be used in fu-

ture NNPDF releases in order to systematically decide which data enters into the global analysis. The idea would be to systematically remove experiments from the global fit, compute their $P(\alpha)$ distributions as discussed above and upon the results decide whether or not it should be kept in the global dataset. As more and more data from HERA and the LHC are becoming available, the question of which data use in the global fits will become more and more pressing, and the strategy outlined here could provide a possible way forward.

References

- [1] S. Forte, G. Watt, Progress in the Determination of the Partonic Structure of the Proton, arXiv:1301.6754.
- [2] S. Dittmaier, S. Dittmaier, C. Mariotti, G. Passarino, R. Tanaka, et al., Handbook of LHC Higgs Cross Sections: 2. Differential Distributions, arXiv:1201.3084, doi:10.5170/CERN-2012-002.
- [3] C. Borschensky, M. Kramer, A. Kulesza, M. Mangano, S. Padhi, et al., Squark and gluino production cross sections in pp collisions at $\sqrt{s} = 13, 14, 33$ and 100 TeV , arXiv:1407.5066.
- [4] G. Bozzi, J. Rojo, A. Vicini, The Impact of PDF uncertainties on the measurement of the W boson mass at the Tevatron and the LHC, Phys.Rev. D83 (2011) 113008. arXiv:1104.2056, doi:10.1103/PhysRevD.83.113008.
- [5] R. D. Ball, V. Bertone, S. Carrazza, C. S. Deans, L. Del Debbio, et al., Parton distributions with LHC data, Nucl.Phys. B867 (2013) 244–289. arXiv:1207.1303, doi:10.1016/j.nuclphysb.2012.10.003.
- [6] R. D. Ball, S. Carrazza, L. Del Debbio, S. Forte, J. Gao, et al., Parton Distribution Benchmarking with LHC Data, JHEP 1304 (2013) 125. arXiv:1211.5142, doi:10.1007/JHEP04(2013)125.
- [7] J. Alwall, R. Frederix, S. Frixione, V. Hirschi, F. Maltoni, et al., The automated computation of tree-level and next-to-leading order differential cross sections, and their matching to parton shower simulations, arXiv:1405.0301.
- [8] R. D. Ball, et al., Unbiased global determination of parton distributions and their uncertainties at NNLO and at LO, Nucl.Phys. B855 (2012) 153–221. arXiv:1107.2652.
- [9] S. Carrazza, S. Forte, J. Rojo, Parton Distributions and Event Generators arXiv:1311.5887.
- [10] T. Sjostrand, S. Mrenna, P. Z. Skands, A Brief Introduction to PYTHIA 8.1, Comput. Phys. Commun. 178 (2008) 852–867. arXiv:0710.3820, doi:10.1016/j.cpc.2008.01.036.
- [11] P. Skands, S. Carrazza, J. Rojo, Tuning PYTHIA 8.1: the Monash 2013 Tune, arXiv:1404.5630.
- [12] R. D. Ball, et al., Parton distributions with QED corrections, Nucl.Phys. B877 (2013) 290–320. arXiv:1308.0598, doi:10.1016/j.nuclphysb.2013.10.010.
- [13] V. Bertone, S. Carrazza, J. Rojo, APFEL: A PDF Evolution Library with QED corrections, Comput.Phys.Commun. 185 (2014) 1647–1668. arXiv:1310.1394, doi:10.1016/j.cpc.2014.03.007.
- [14] R. Boughezal, Y. Li, F. Petriello, Disentangling radiative corrections using high-mass Drell-Yan at the LHC, Phys.Rev. D89 (2014) 034030. arXiv:1312.3972, doi:10.1103/PhysRevD.89.034030.
- [15] A. Bierweiler, T. Kasprzik, J. H. Kuhn, Vector-boson pair production at the LHC to $\mathcal{O}(\alpha^3)$ accuracy, JHEP 1312 (2013) 071. arXiv:1305.5402, doi:10.1007/JHEP12(2013)071.
- [16] E. R. Nocera, R. D. Ball, S. Forte, G. Ridolfi, J. Rojo, A first unbiased global determination of po-

- larized PDFs and their uncertainties, arXiv:1406.5539, doi:10.1016/j.nuclphysb.2014.08.008.
- [17] R. D. Ball, et al., Unbiased determination of polarized parton distributions and their uncertainties, Nucl.Phys. B874 (2013) 36–84. arXiv:1303.7236, doi:10.1016/j.nuclphysb.2013.05.007.
- [18] D. de Florian, R. Sassot, M. Stratmann, W. Vogelsang, Evidence for polarization of gluons in the proton, Phys.Rev.Lett. 113 (2014) 012001. arXiv:1404.4293, doi:10.1103/PhysRevLett.113.012001.
- [19] M. Czakon, M. L. Mangano, A. Mitov, J. Rojo, Constraints on the gluon PDF from top quark pair production at hadron colliders, JHEP 1307 (2013) 167. arXiv:1303.7215, doi:10.1007/JHEP07(2013)167.
- [20] S. Alekhin, J. Bluemlein, S. Moch, The ABM parton distributions tuned to LHC data, Phys.Rev. D89 (2014) 054028. arXiv:1310.3059, doi:10.1103/PhysRevD.89.054028.
- [21] W. Stirling, E. Vryonidou, Charm production in association with an electroweak gauge boson at the LHC, Phys.Rev.Lett. 109 (2012) 082002. arXiv:1203.6781, doi:10.1103/PhysRevLett.109.082002.
- [22] S. 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, arXiv:1312.6283.
- [23] T. Carli, D. Clements, A. Cooper-Sarkar, C. Gwenlan, G. P. Salam, 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 (2010) 503–524. arXiv:0911.2985, doi:10.1140/epjc/s10052-010-1255-0.
- [24] M. Wobisch, D. Britzger, T. Kluge, K. Rabbertz, F. Stober, Theory-Data Comparisons for Jet Measurements in Hadron-Induced Processes, arXiv:1109.1310.
- [25] V. Bertone, R. Frederix, S. Frixione, J. Rojo, M. Sutton, aMCfast: automation of fast NLO computations for PDF fits, JHEP 1408 (2014) 166. arXiv:1406.7693, doi:10.1007/JHEP08(2014)166.
- [26] M. Czakon, A. Mitov, Top++: a program for the calculation of the top-pair cross-section at hadron colliders, arXiv:1112.5675.
- [27] R. Gavin, Y. Li, F. Petriello, S. Quackenbush, W. Physics at the LHC with FEWZ 2.1, Comput.Phys.Commun. 184 (2013) 208–214. arXiv:1201.5896, doi:10.1016/j.cpc.2012.09.005.
- [28] D. de Florian, P. Hinderer, A. Mukherjee, F. Ringer, W. Vogelsang, Approximate next-to-next-to-leading order corrections to hadronic jet production, arXiv:1310.7192.
- [29] A. Gehrmann-De Ridder, T. Gehrmann, E. Glover, J. Pires, Second order QCD corrections to jet production at hadron colliders: the all-gluon contribution, Phys.Rev.Lett. 110 (2013) 162003. arXiv:1301.7310, doi:10.1103/PhysRevLett.110.162003.
- [30] S. Carrazza, J. Pires, Perturbative QCD description of jet data from LHC Run-I and Tevatron Run-II, arXiv:1407.7031.
- [31] A. D. Martin, R. G. Roberts, W. J. Stirling, R. S. Thorne, Uncertainties of predictions from parton distributions. II: Theoretical errors, Eur. Phys. J. C35 (2004) 325–348. arXiv:hep-ph/0308087.
- [32] R. D. Ball, et al., Reweighting NNPDFs: the W lepton asymmetry, Nucl. Phys. B849 (2011) 112–143. arXiv:1012.0836, doi:10.1016/j.nuclphysb.2011.03.017.
- [33] R. D. Ball, V. Bertone, F. Cerutti, L. Del Debbio, S. Forte, et al., Reweighting and Unweighting of Parton Distributions and the LHC W lepton asymmetry data, Nucl.Phys. B855 (2012) 608–638. arXiv:1108.1758, doi:10.1016/j.nuclphysb.2011.10.018.
- [34] J. Butterworth, G. Dissertori, S. Dittmaier, D. de Florian, N. Glover, et al., Les Houches 2013: Physics at TeV Colliders: Standard Model Working Group Report, arXiv:1405.1067.

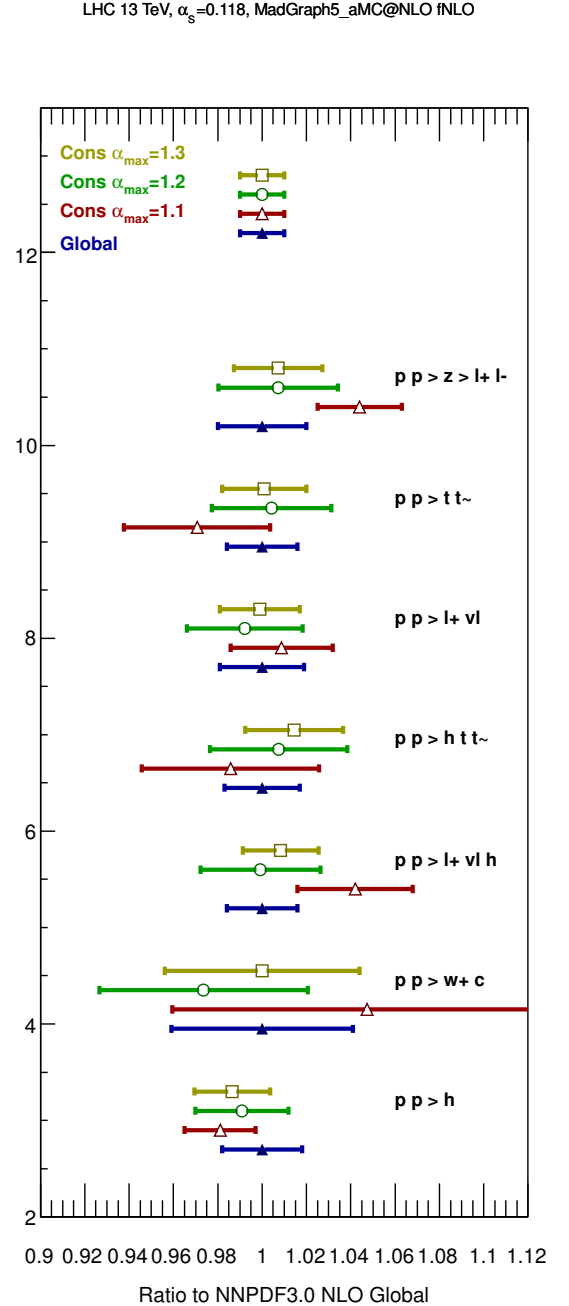


Figure 3: Comparison of the predictions for a number of LHC observables between NNPDF3.0 NLO and the various corresponding conservative fits. Results have been computed with MadGraph5_aMC@NLO in the fNLO mode for the LHC 13 TeV and typical LHC analysis cuts. The cross-sections are shown as ratios with respect to the NNPDF3.0 NLO central value.