

Toward precise and robust unpolarized PDFs

[Parton distributions need a representative sampling, in progress]

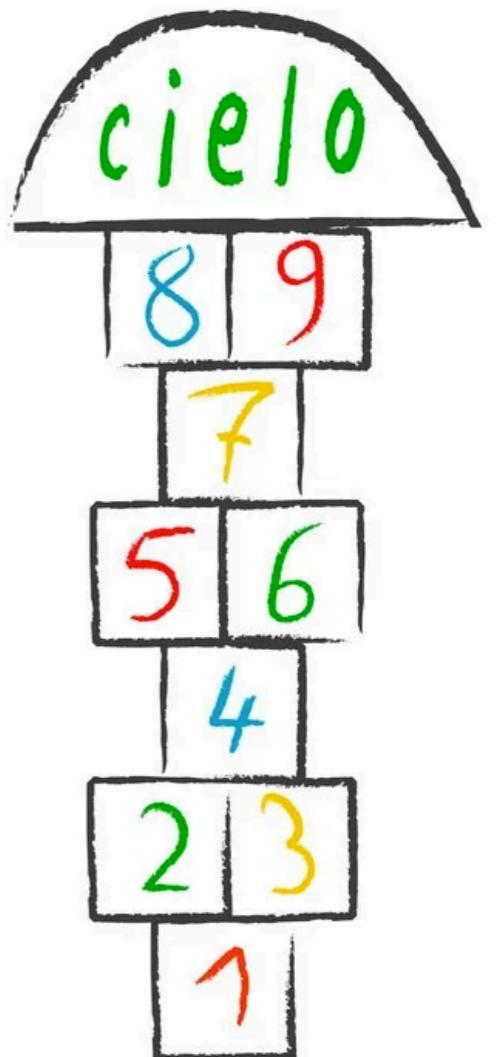
Aurore Courtoy

Instituto de Física

National Autonomous University of Mexico (UNAM)

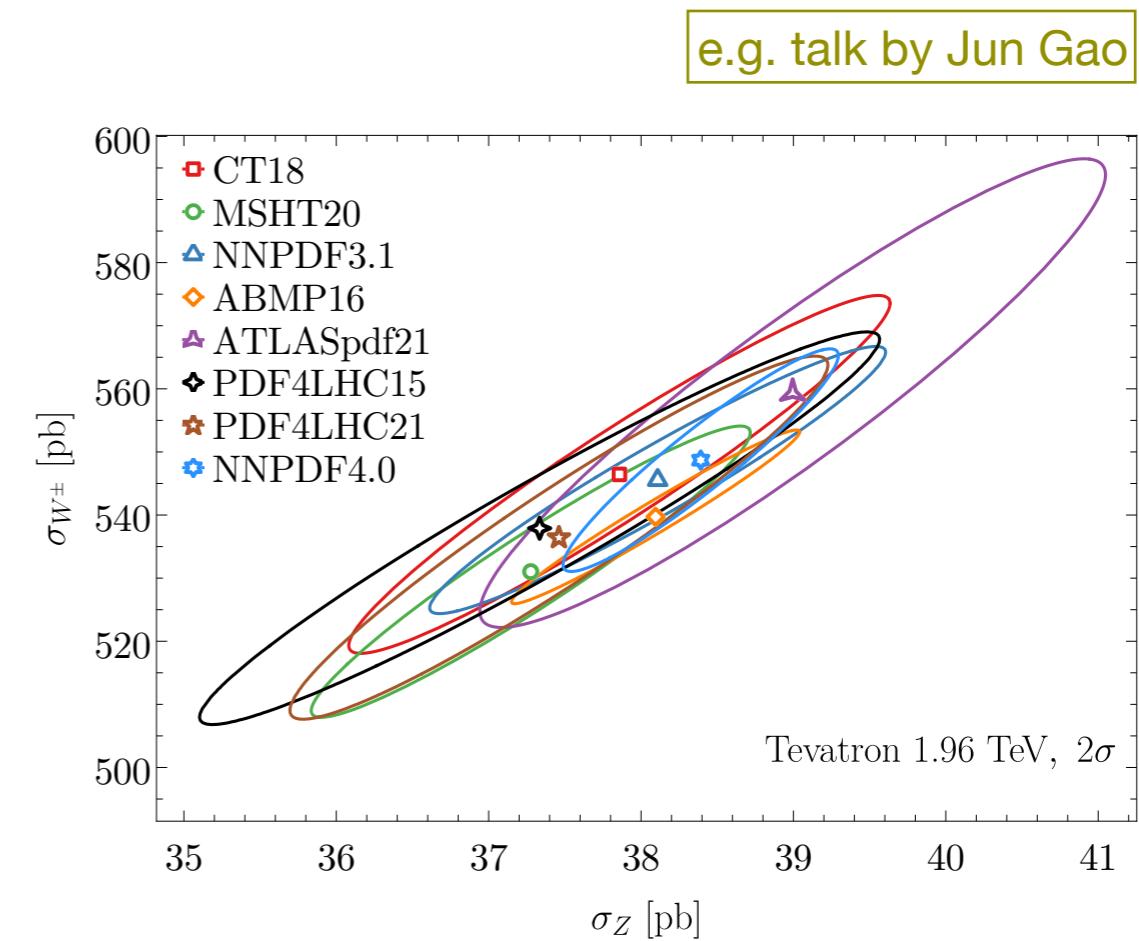
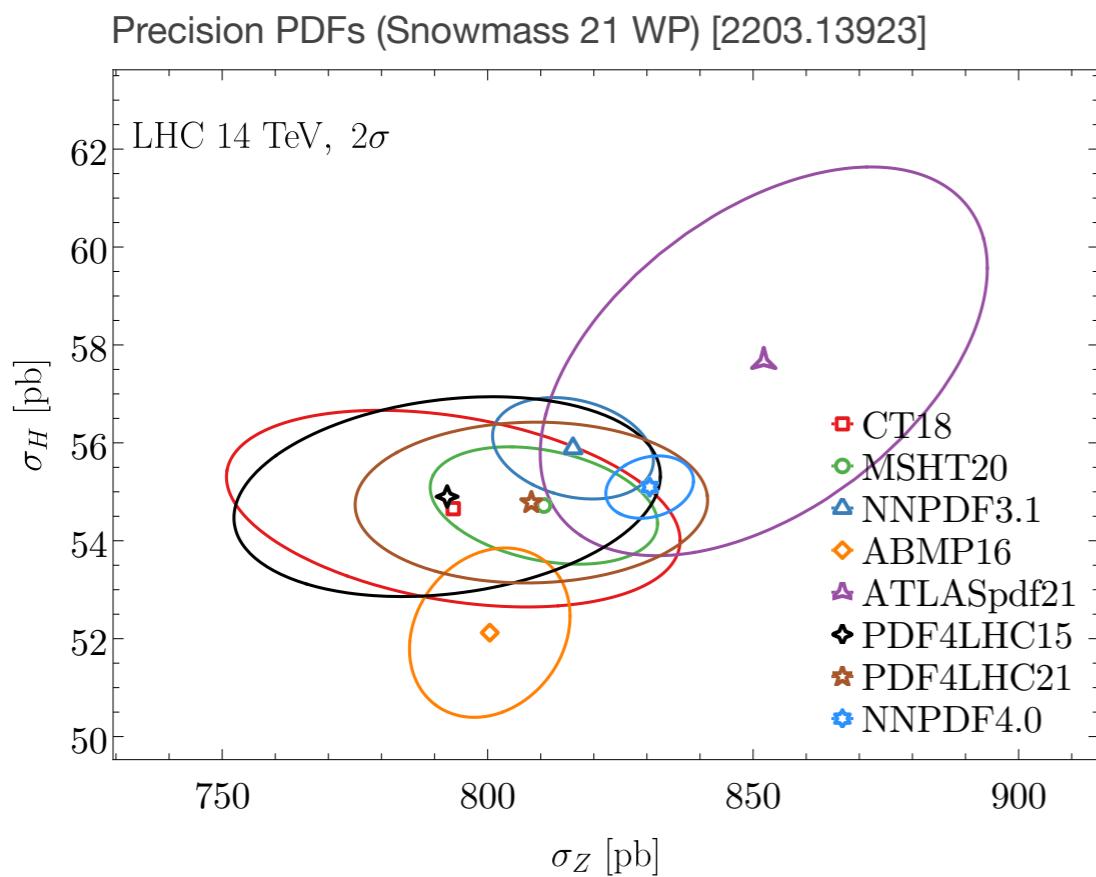
with J. Huston, P. Nadolsky, K. Xie, M. Yan, C.-P. Yuan

DIS 2022 – Santiago de Compostela



Uncertainties from global analyses of proton structure

Recent advancements in the determination of unpolarized PDFs:
CT18, MSHT20, NNPDF4.0, ATLASpdf21 as well as PDF4LHC21.



What is a faithful uncertainty coming from PDFs on those cross sections?

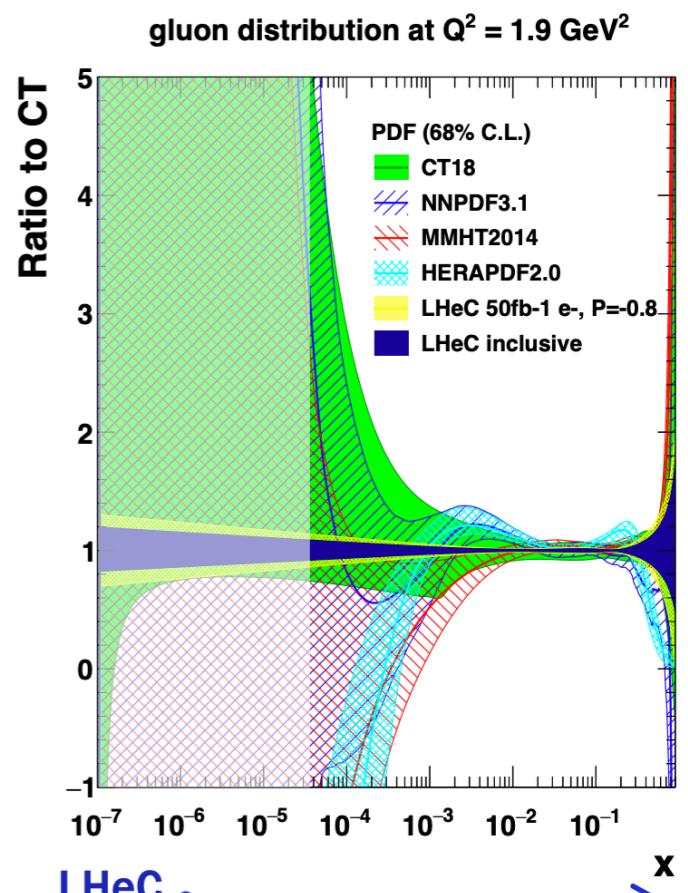
Reducing PDFs and α_s uncertainties for EW and BSM physics

Theoretical progress elevates precision on pQCD predictions.

Measurements of several SM parameters depend on PDF uncertainties.

Future experiments will potentially increase the precision of PDFs:
LHeC, EIC, HL-LHC,...

Future global analyses will require **thorough understanding** of **various sources of uncertainties** in the PDF determination.



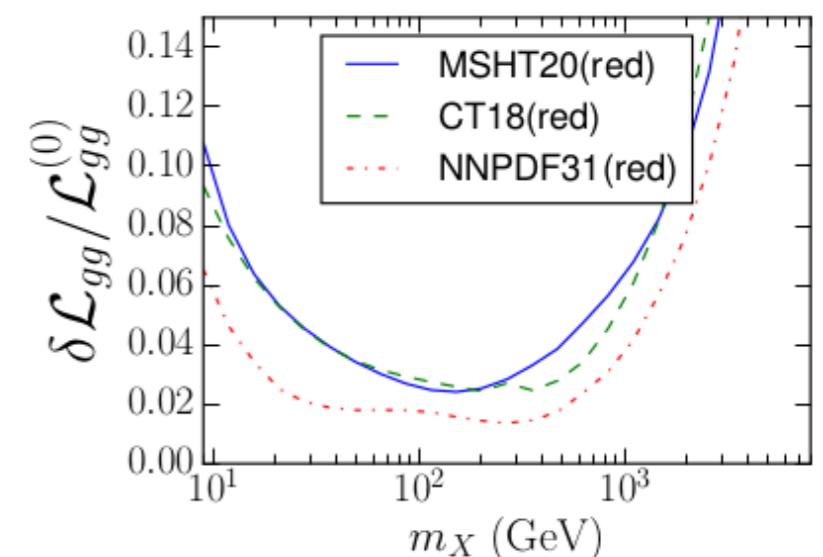
Plot from C. Gwenlan ICHEP 2020

PDF4LHC21 benchmarking exercise:
comparison of uncertainties for same sets of data and QCD settings.

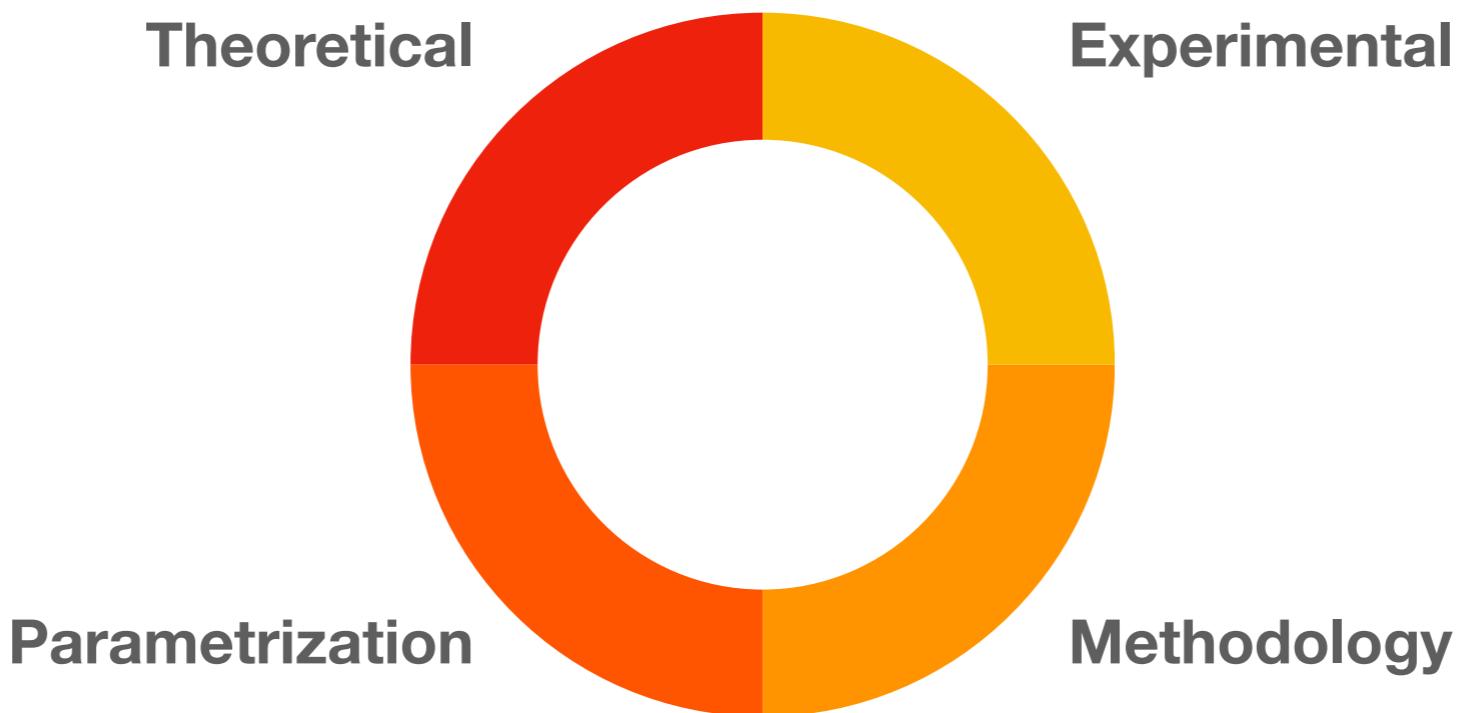
talk by Tom Cridge

The uncertainties for CT18, MSHT20 and NNPDF3.1 reduced sets are still different. Key role played by methodology.

PDF4LHC21 [2203.05506]



Sampling biases contribute to PDF uncertainties



In all four categories of uncertainties, we can further distinguish

PDF fitting accuracy and PDF sampling accuracy.

Accuracy in theo., exp.,... inputs
— commonly integrated in global analyses.

Adequacy of sampling the space of all
solutions — traditionally ignored.

[Kovarik et al, Rev.Mod.Phys. 92 (2020)]

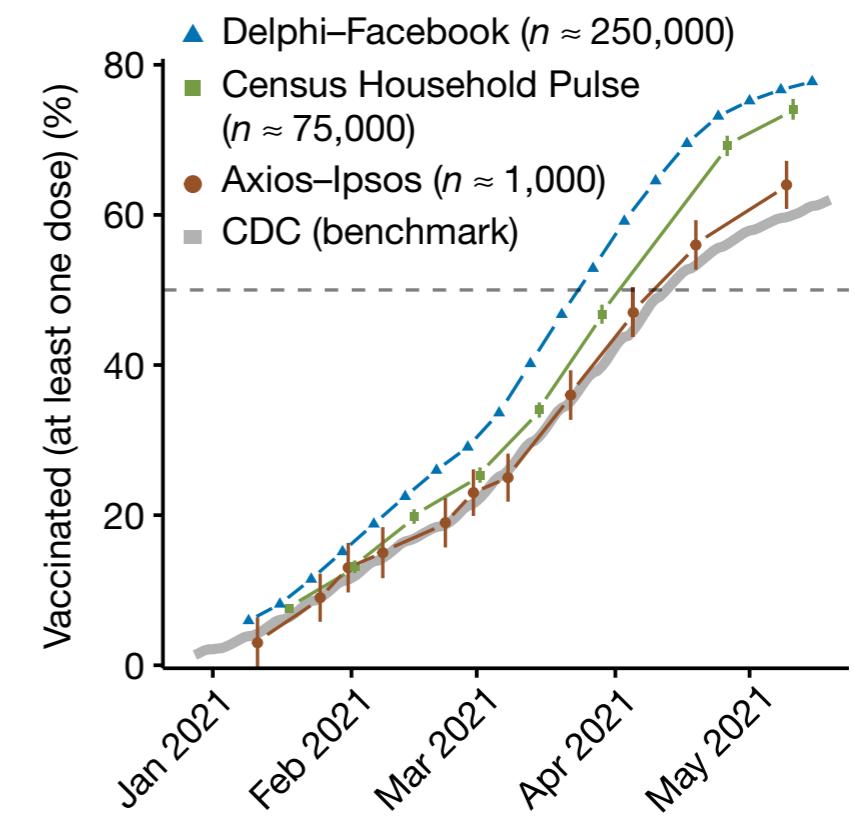
Control for **sampling biases** in determination of PDFs plays a critical role.

Origin of sampling biases — experience with large population surveys

Surveys of the COVID-19 vaccination rate with very large samples of responses and small statistical uncertainties (*Delphi-Facebook*) greatly overestimated the actual vaccination rate published by the Center for Disease Control (*CDC*) after some time delay.

The screenshot shows the **nature** journal website. The main navigation bar includes "Explore content", "About the journal", and "Publish with us". Below the header, the URL "nature > articles > article" is shown. The article title is "Unrepresentative big surveys significantly overestimated US vaccine uptake" by Valerie C. Bradley, Shiro Kuriwaki, Michael Isakov, Dino Sejdinovic, Xiao-Li Meng & Seth Flaxman. It was published on 08 December 2021. The article summary states: "Unrepresentative big surveys significantly overestimated US vaccine uptake". The authors' names are listed, along with the journal issue ("Nature 600, 695–700 (2021)") and a link to "Cite this article".

Based on
[Xiao-Li Meng, The Annals of Applied Statistics, Vol. 12 (2018), p. 685]



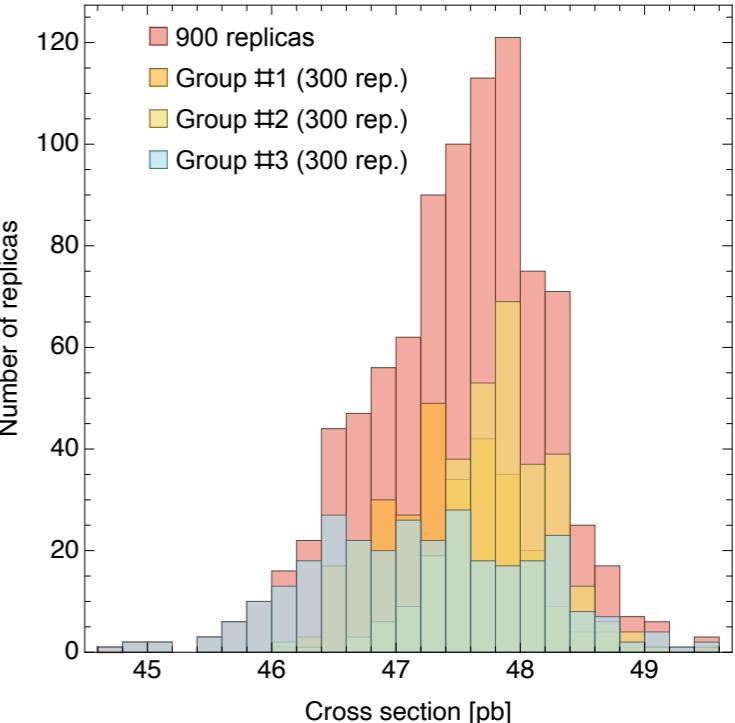
The deviation has been traced to the **sampling bias**.
In contrast to the statistical error, the sampling bias can involve growth with the size of the sample.

Law of large numbers

With an increasing size of sample $n \rightarrow \infty$, under a set of hypotheses, it is usually expected that the deviation on an observable

$$\mu - \hat{\mu} \propto \sigma/\sqrt{n}$$

with σ the standard deviation, μ the true and $\hat{\mu}$ the determined values. *That's the law of large numbers.*

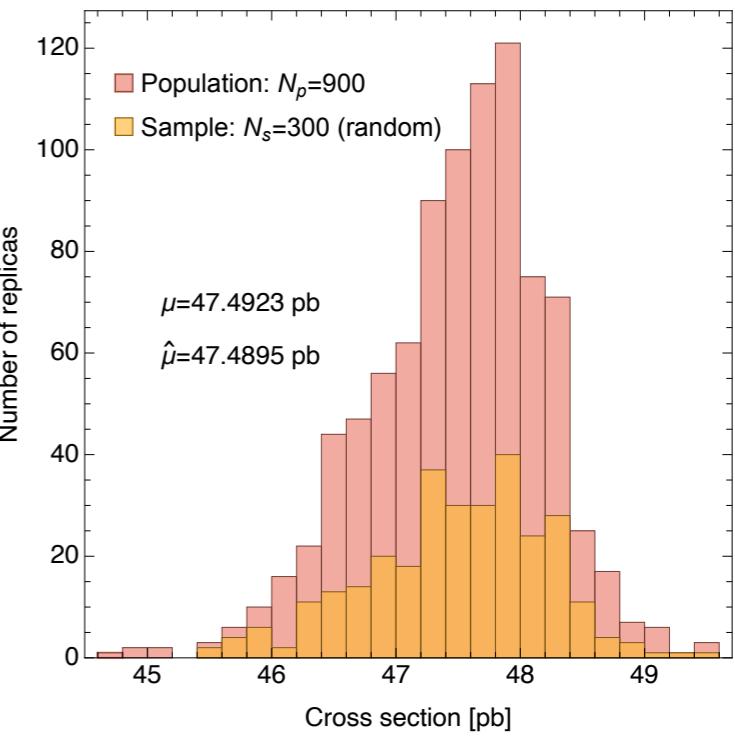


A toy sampling exercise

We take 300×3 groups of **Higgs cross sections** evaluated by 3 different groups.

We **randomly** select 300 out of the 900 cross sections.

The law of large number is fulfilled in this case: there is no bias.

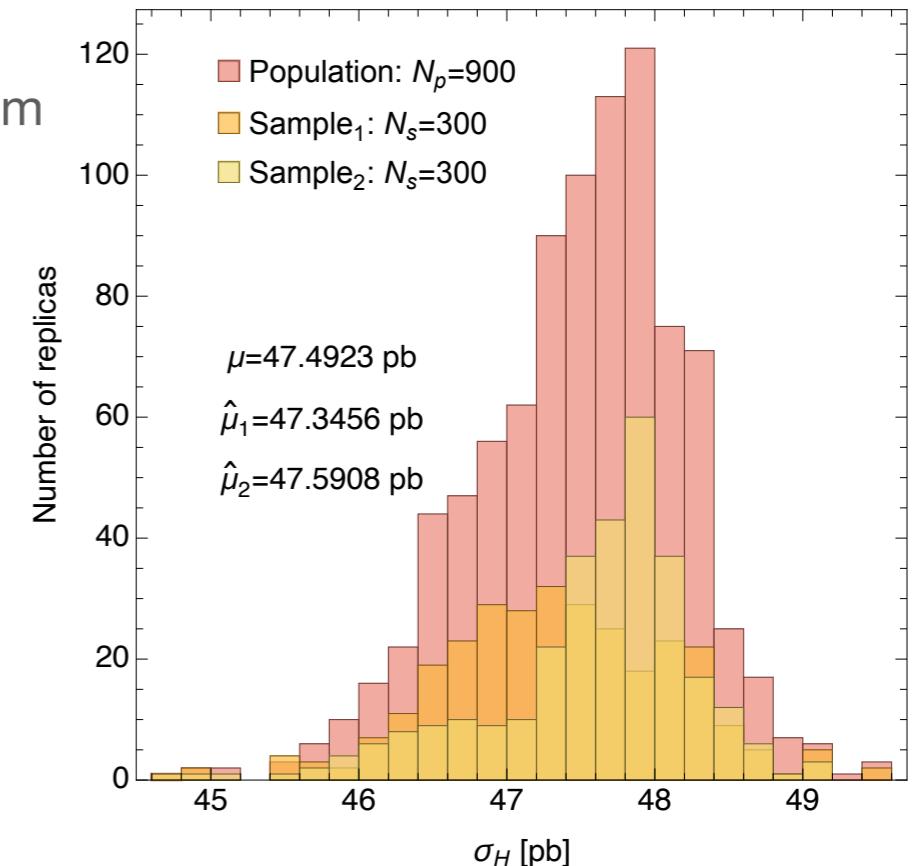


Trio identity

If we **bias** the selection by taking 200 items from one group and 100 from another, the deviation $\mu - \hat{\mu}$ is no longer proportional to σ/\sqrt{n} !



The law of large numbers obviates the *quality of the sampling*.



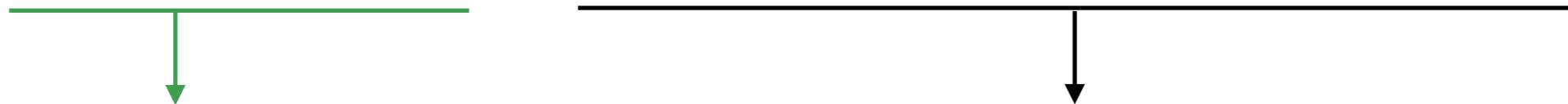
The **trio identity** remedies to that problem by accounting for sampling bias:

$$\mu - \hat{\mu} = (\text{data+sampling defect}) \times (\text{measure discrepancy}) \times (\text{inherent problem difficulty})$$

This identity originates from the statistics of large-scale surveys
[Xiao-Li Meng, The Annals of Applied Statistics, Vol. 12 (2018), p. 685]

Trio identity

$$\mu - \hat{\mu} = (\text{data+sampling defect}) \times (\text{measure discrepancy}) \times (\text{inherent problem difficulty})$$



depends on the sampling algorithm

can tend to σ/\sqrt{n} for random sampling

For a sample of n items from the population of size N , we can consider an array built by the random spanning of the binary responses of the $N - n$ (0) and n (1) items, so that

$$\mu - \hat{\mu} = \text{Corr}[\text{observable, sampling quality}] \times \sqrt{\frac{N}{n} - 1} \times \sigma(\text{observable})$$

[X.-L.Meng, The Annals of Applied Statistics, Vol. 12 (2018), p. 685]
[Hickernell, MCQMC 2016, 1702.01487]

The sample deviation can be large if the sampling is not sufficiently random.

Standard error estimates can be misleadingly small.

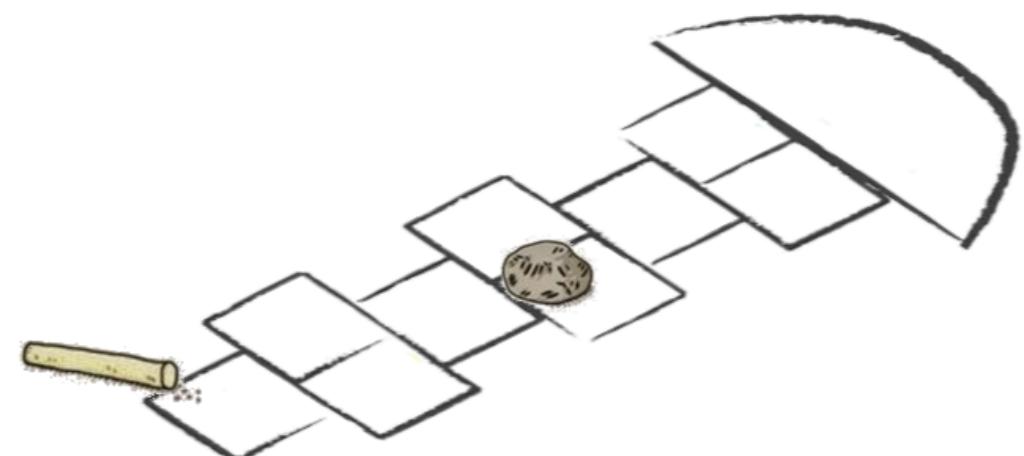
⇒ critical role of controlling for **sampling biases** in determination of PDFs.

Estimation of a representative uncertainty on a cross section σ

To sample the PDF dependence: sample primarily the coordinates with large variations of σ .

We employ:

1. Basis coordinates in the PDF space
2. Knowledge of 4-8 "large dimensions" in PDF space controlling variation of σ
3. A moderate number of MC PDF replicas varying primarily in these directions



Based on the ideas of [Hickernell, MCQMC 2016, 1702.01487]
[Sloan,I.H.,Wo'zniakowski, 1997]

ELEMENTOS PARA LLEGAR AL CIELO

A hopscotch scan of LHC cross sections for NNPDF4.0 PDFs

Step 1

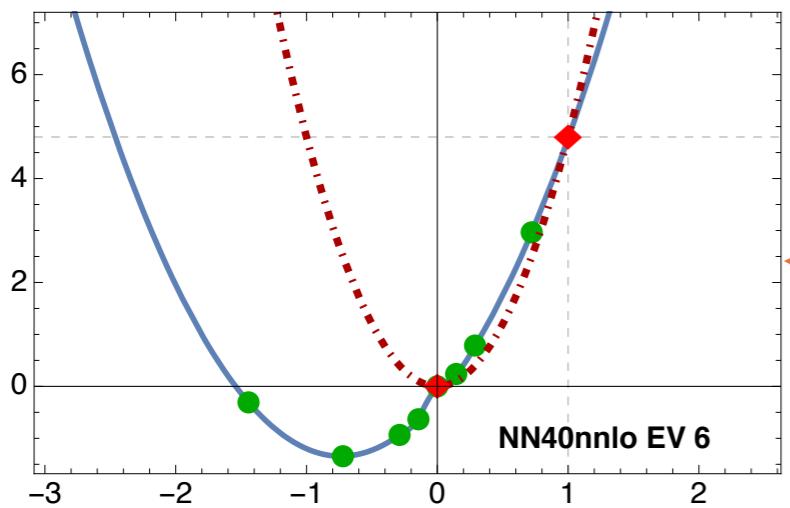
The NNPDF4.0 Hessian set ($n = 50$) defines a coordinate system on a manifold corresponding to the largest variations of the PDF uncertainty — **red dots and curve**.

[NNPDF, 2109.02653]

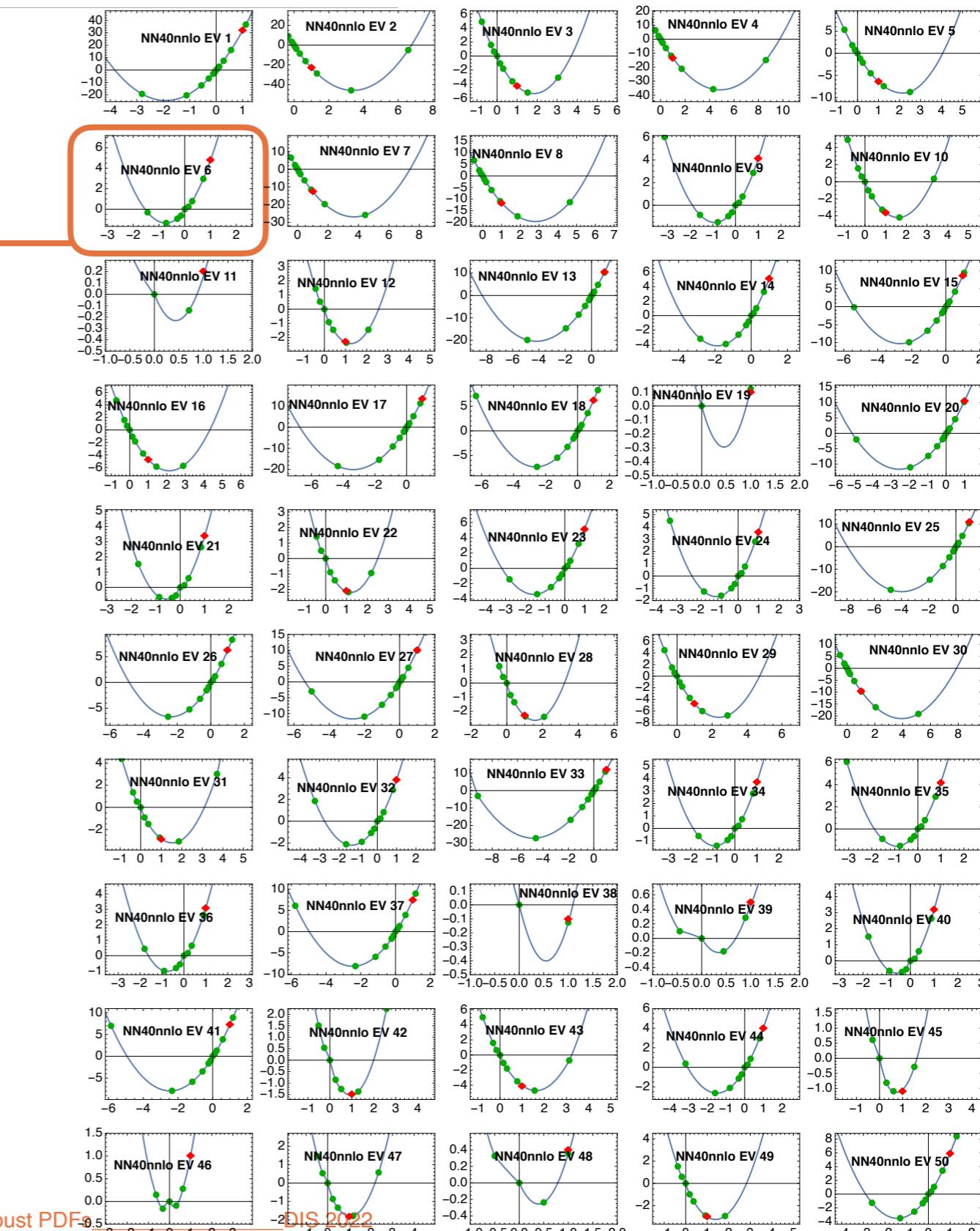
Step 2

Using the public NNPDF code, scan χ^2_{tot} along the 50 EV directions to identify a hypercube corresponding to $\Delta\chi^2 \leq T^2$ (where $T^2 > 0$ is a user-selected value).

Lagrange multiplier scan confirms the approximate Gaussian profiles, but suggest that there exist solutions with lower χ^2 — **green dots and blue curve**.



A. Courtoy—IFUNAM



Robust PDFs

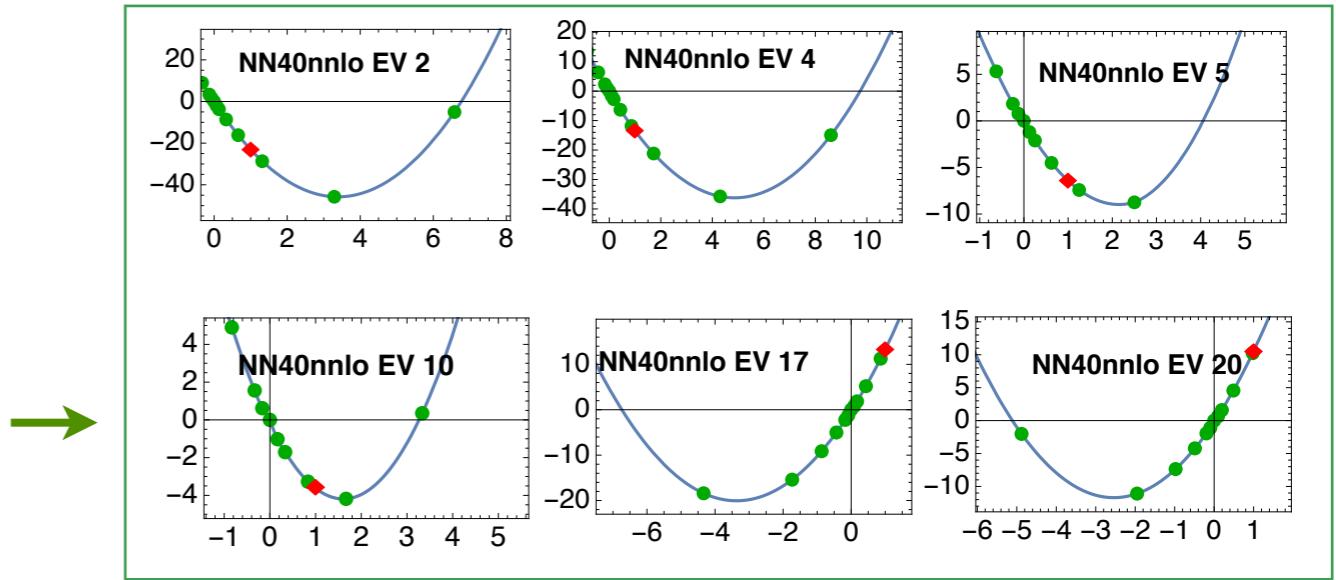
DIS 2022

A hopscotch scan of LHC cross sections for NNPDF4.0 PDFs

Step 3

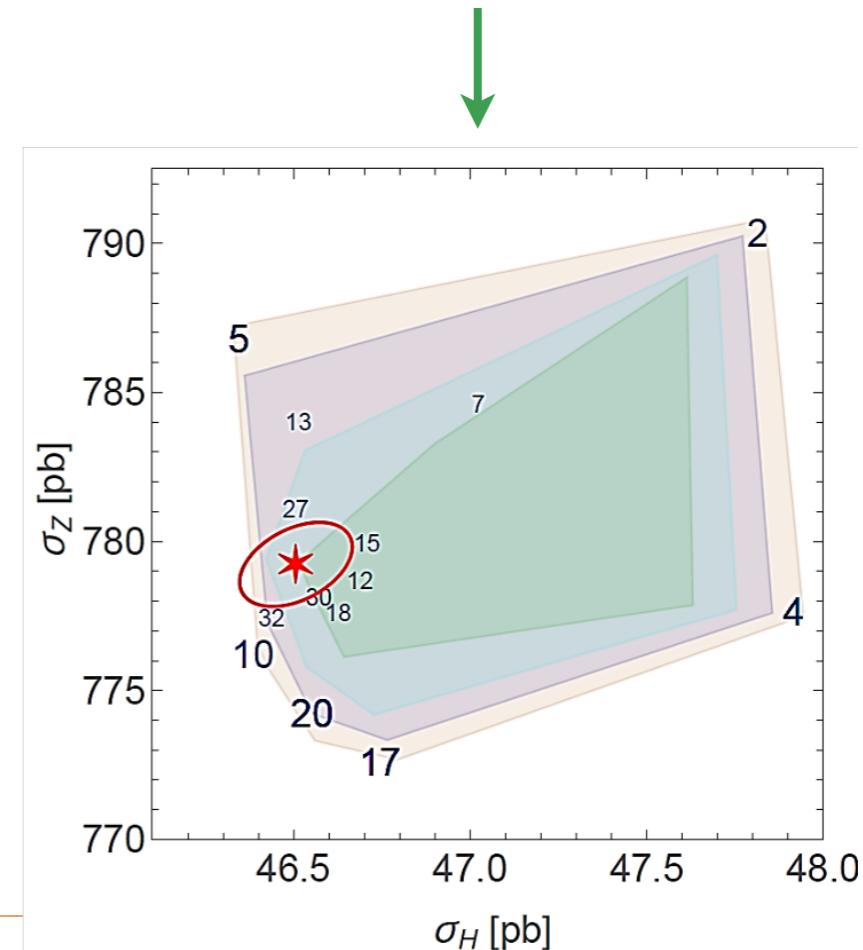
Guidance from specific cross sections:
we identify 4-7 EV directions that give the largest
displacements for a given $\Delta\chi^2$ per pair.

E.g., σ_Z vs. σ_H is represented by the 6 corners of a projected
octahedron, corresponding to “large” EV directions: 2, 4, 5,
10, 17, 20.



Other directions generally give smaller displacements.
Large EV directions are shared among various pairs of
cross sections.

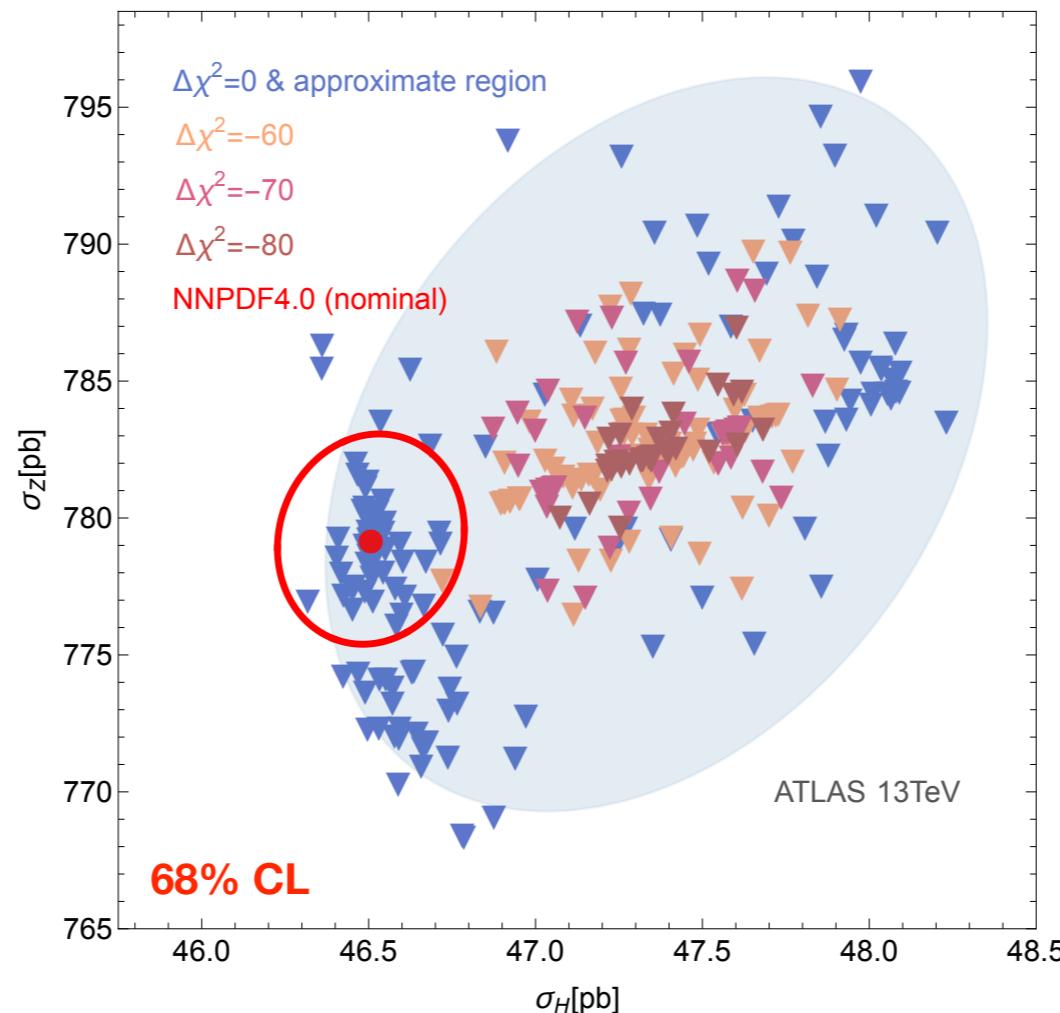
The contours are for $\Delta\chi^2 = +10, 0, -10, -20$ w.r.t.
NNPDF4.0 replica 0 (red).



A hopscotch scan of LHC cross sections for NNPDF4.0 PDFs

Step 4

For each pair of cross sections, we generate 300 replicas by sampling uniformly along the “large” EV directions. Sort the $n_{pairs} \times 300$ resulting replicas according to their $\Delta\chi^2$ w.r.t. to NN40 replica 0.



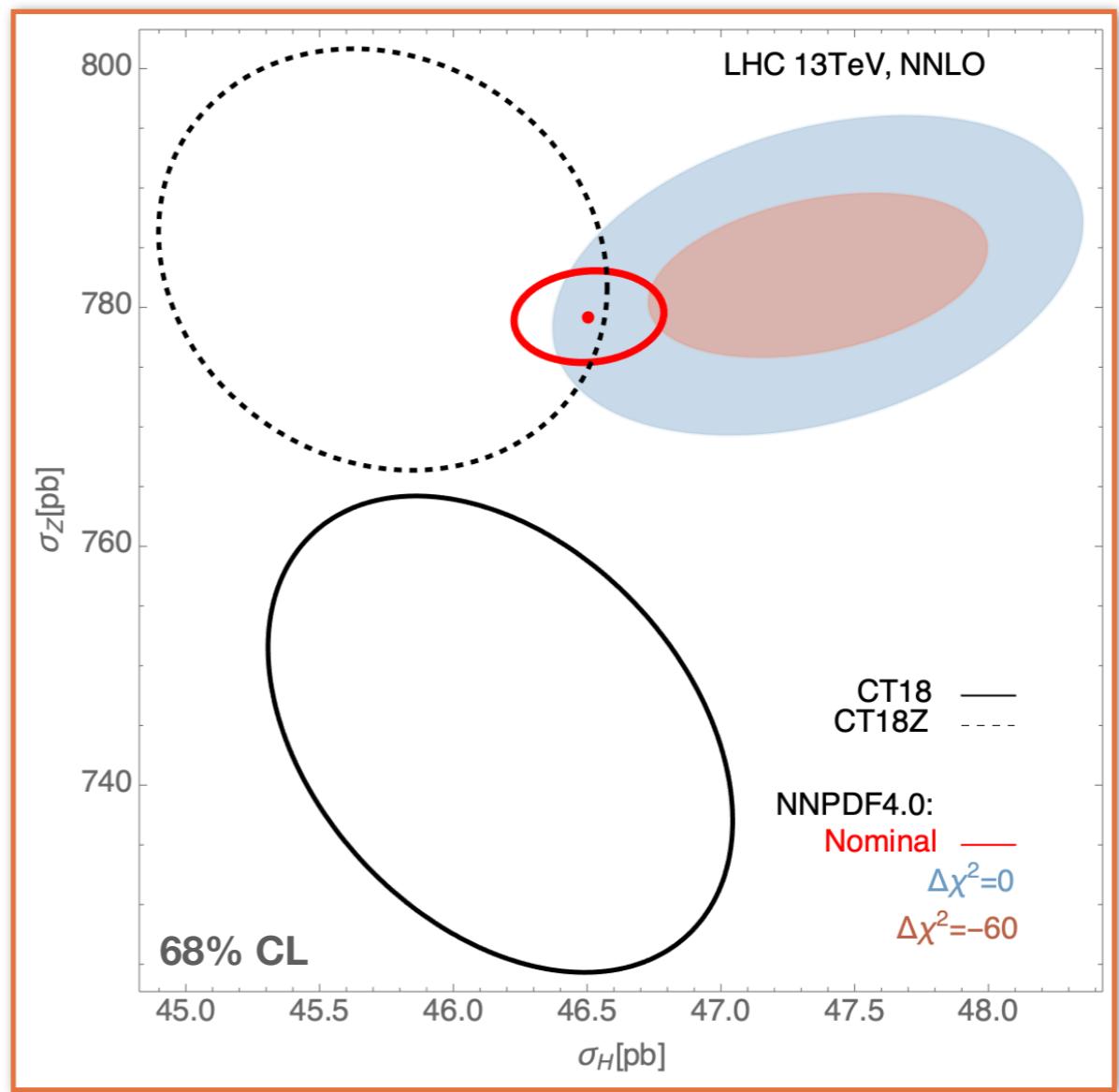
Each of the $\Delta\chi^2 = 0 \pm 3$ replicas is an acceptable PDF set from the NNPDF4.0 fit.

The blue ellipse (constructed using a convex hull method) is an approximate region containing all found replicas with $\Delta\chi^2 = 0 \pm 3$.

[Anwar, Hamilton, Nadolsky, 1901.05511]

The blue area is larger than the nominal NNPDF4.0 uncertainty (red ellipse).

Monte-Carlo sampling for PDF parametrizations: cross sections for LHC



Blue and brown filled ellipses:

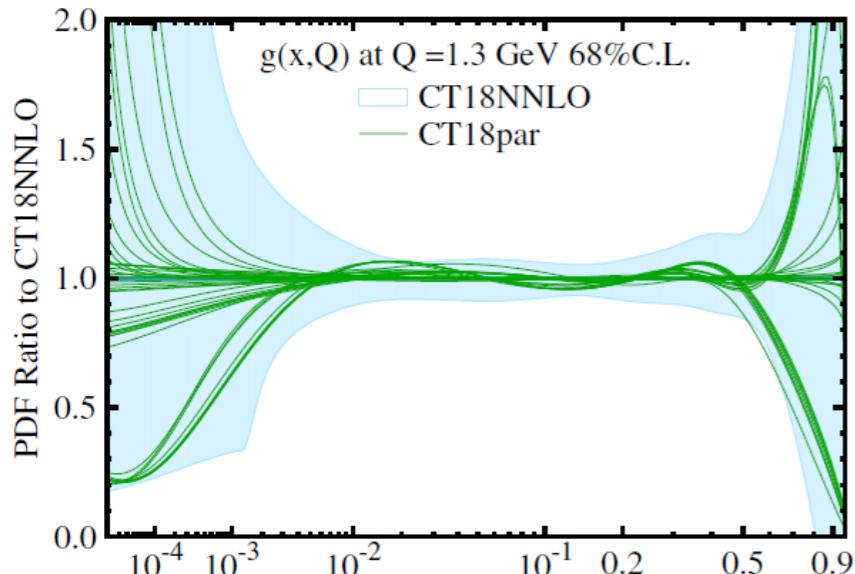
- areas of possible solutions corresponding to an equal ($\Delta\chi^2 = 0$) or lower ($\Delta\chi^2 = -60$) chi square w.r.t. the nominal solution
- found through the hopscotch scan — a dimensionality reduction method.
- size of blue areas comparable to 68% CL CT18 ellipses

Monte-Carlo sampling for PDF parametrizations: cross sections for LHC

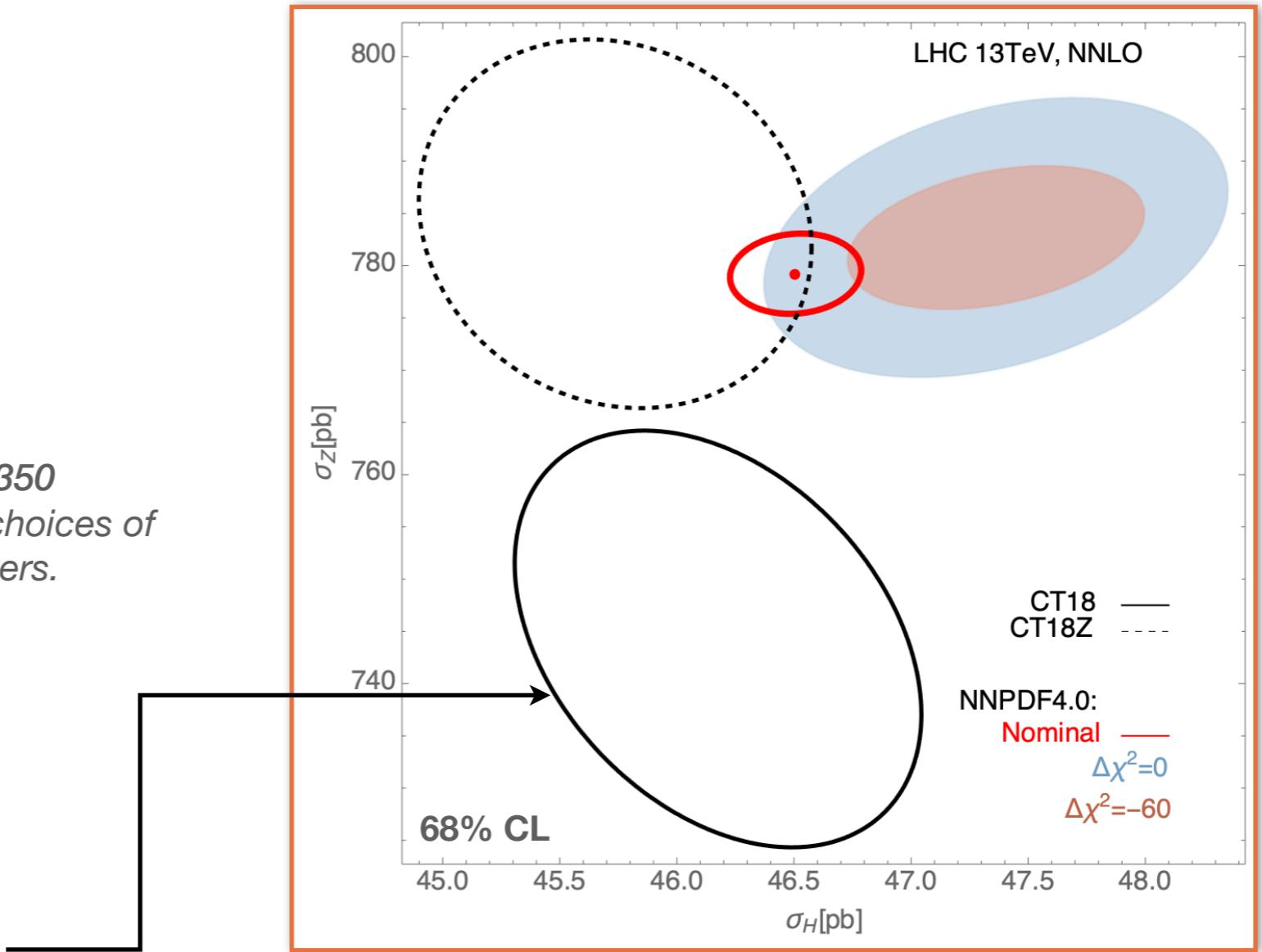
CT18 PDF uncertainty:

Accounts for the sampling over 250-350 parametrization forms and possible choices of fitted experiments and fitting parameters.

Reflected in choice of tolerance.



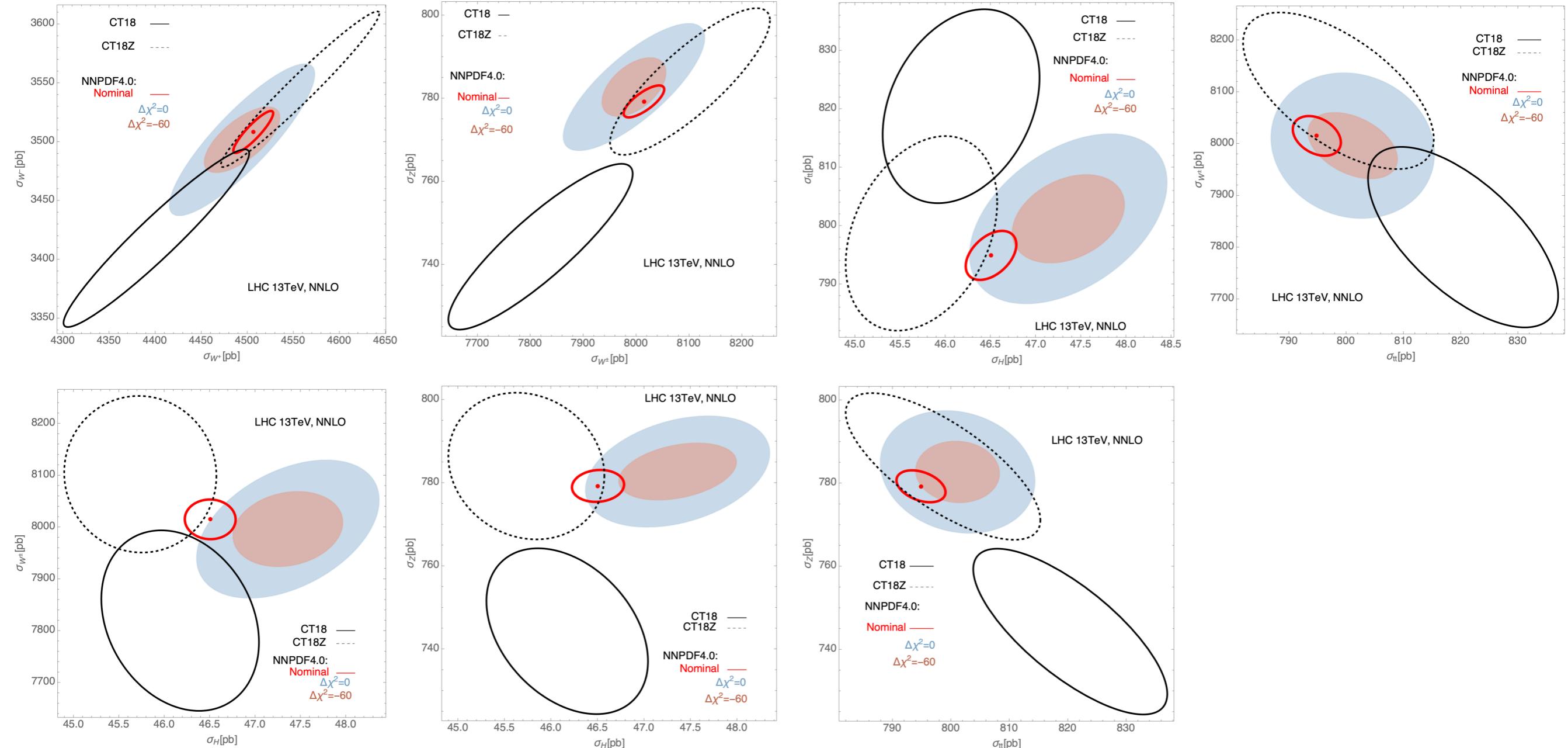
[Hou et al, Phys.Rev.D 103 (2021)]



Blue and brown filled ellipses:

- areas of possible solutions corresponding to an equal ($\Delta\chi^2 = 0$) or lower ($\Delta\chi^2 = -60$) chi square w.r.t. the nominal solution
- found through the hopscotch scan — a dimensionality reduction method.
- size of blue areas comparable to 68% CL CT18 ellipses

Monte-Carlo sampling for PDF parametrizations: cross sections for LHC



Ellipses at 68% CL

Conclusions

What is a faithful uncertainty coming from PDFs on those cross sections?

PDF uncertainties in high-stake measurements (Higgs cross sections, W mass...) should be examined for *robustness of sampling*.

Sampling biases: may arise in PDF fits operating with large samples of data or multiparametric functional forms.
The trio identity may take over the law of large numbers.

An undetected sampling bias may result in a wrong prediction with a low nominal uncertainty.

Sample deviation may limit reduction of the PDF uncertainties and may explain some differences between the uncertainties of the PDF sets.

Experience with big surveys and Monte-Carlo integration shows how to quantify such deviations for QCD parameters or cross sections.

⇒ possible framework for systematic study of parametrization within CT.

Hopscotch scans illustrated for the NNPDF4.0 –thanks to the publicly available code.

Applicable to other analyses using similar methodology and a large enough parameter space – e.g. for polarized PDFs.

Back-up slides

Toward robust PDF uncertainties

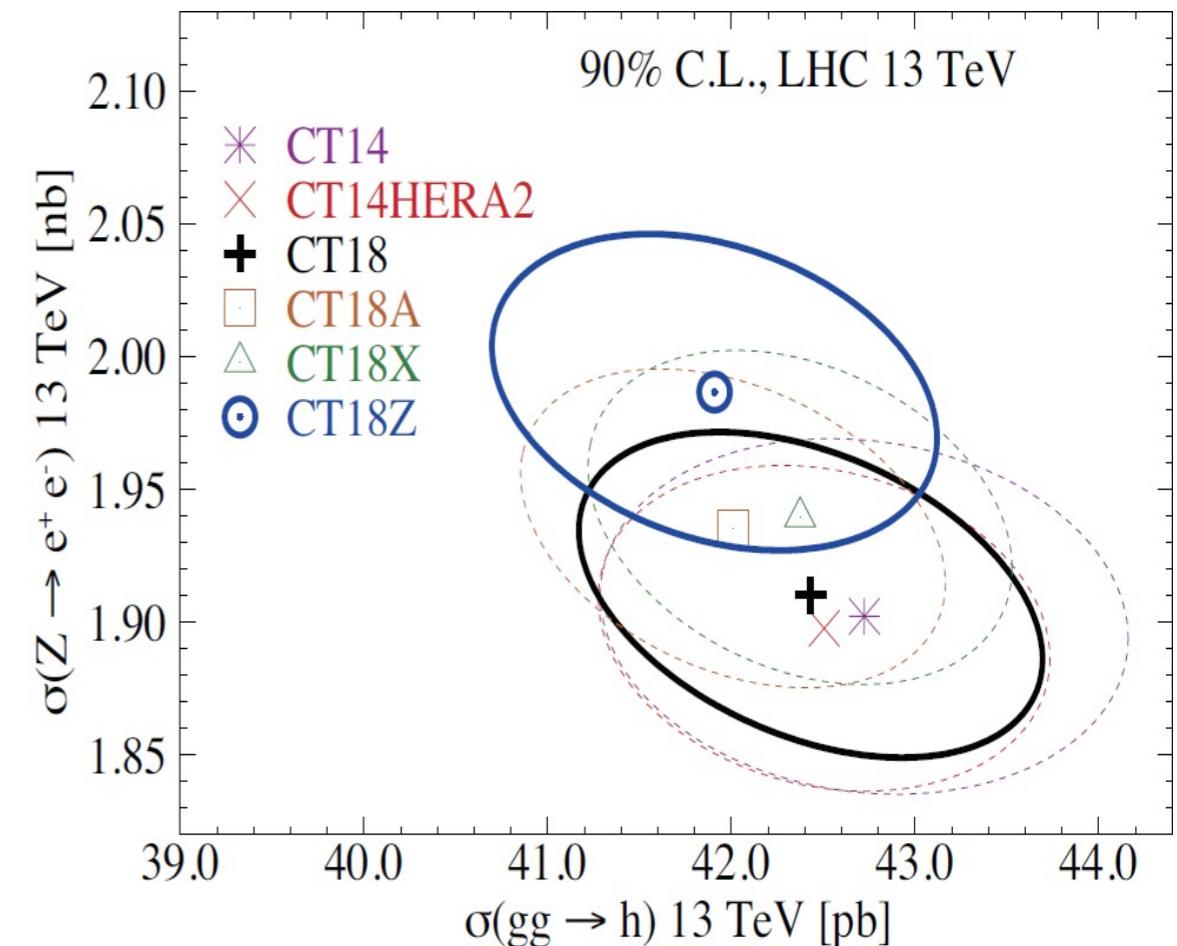
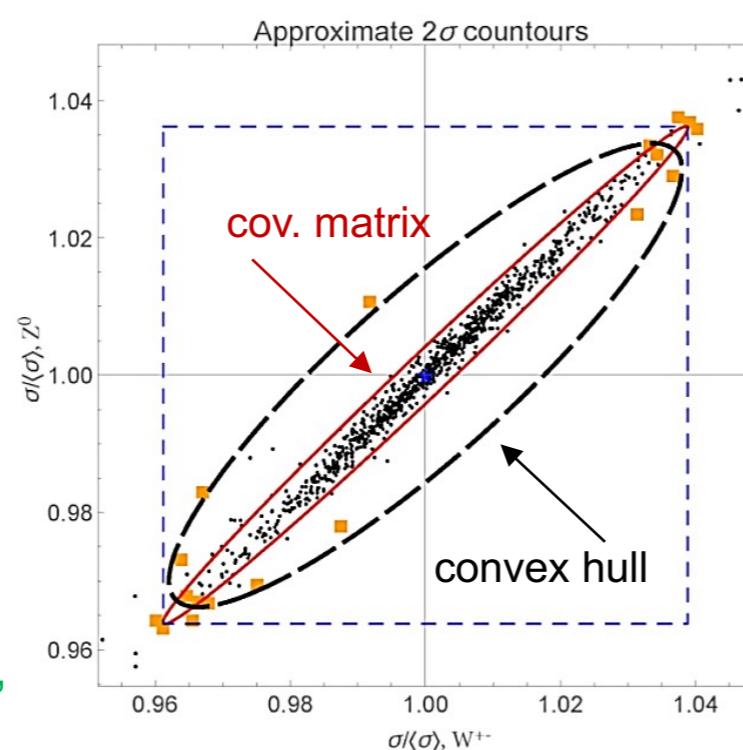
Strong dependence on the definition of corr.
syst. errors would raise a general concern:

**Overreliance on Gaussian distributions
and covariance matrices for poorly
understood effects may produce very
wrong uncertainty estimates**

[N. Taleb, Black Swan & Antifragile]

For instance, the cov. matrix may overestimate the correlation among discrete data points, resulting in a too aggressive error estimate

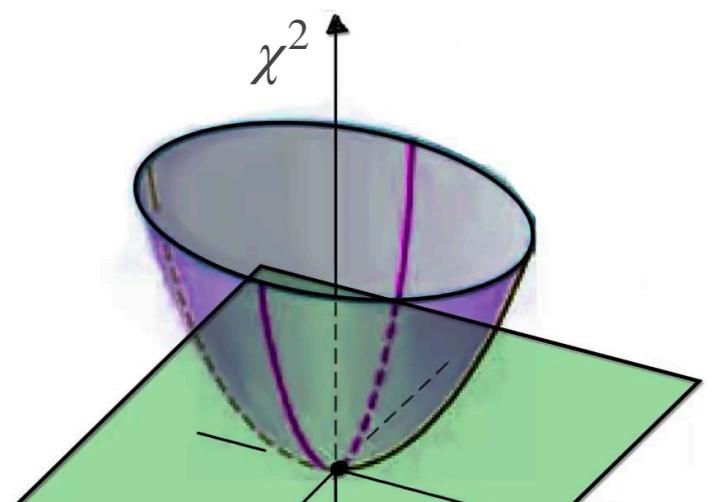
[Anwar, Hamilton, P.N., arXiv:1905.05111]



The CT18/CT18Z uncertainties aim to be **robust**: they largely cover the spread of central predictions obtained with different selections of experiments and assumptions about systematic uncertainties

Uncertainty on QCD observables – the hopscotch

Hessian methods are based on the paraboloid behavior of the χ^2 function
– PDF eigenvector set naturally renders the coordinates giving the largest contribution to a determined value $\hat{\mu}$, with the Principal Component Analysis or a related method.



Uncertainty on QCD observables – the hopscotch

Sampling of multidimensional spaces ($d \gg 20$) can be exponentially inefficient and require $n > 2^d$ replicas to obtain a convergent expectation value.

[Hickernell, MCQMC 2016, 1702.01487]
[Sloan,I.H.,Wo'zniakowski, 1997]

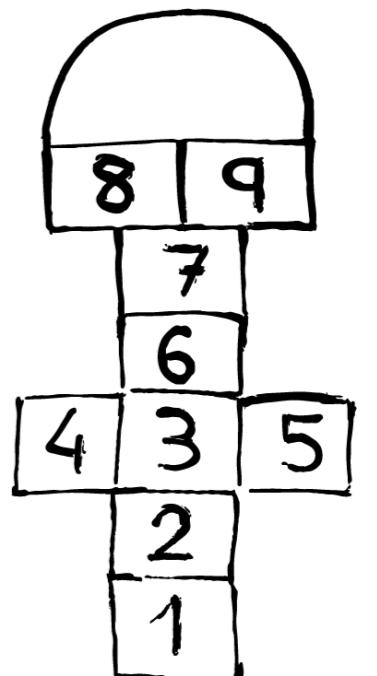
Specific QCD observables: only few effective large dimensions contribute the bulk of the uncertainty.
E.g. compressing MC PDFs into a Hessian set: we construct a basis to identify such large dimensions.

Hopscotch scans:

estimation of a representative uncertainty on a cross section.

The release of a public code for NNPDF4.0's new methodology provide a perfect playground to explore the role of sampling.

[NNPDF, EPJC 81]



Setting for NNPDF4.0 code

The evaluation of χ^2 for NNPDF4.0 nnlo replicas is done by the public NNPDF code [NNPDF, EPJC 81], with its default setting.

χ^2 is computed by the `perreplica_chi2_table` function of `validphys` program of the public NNPDF code.

The kinematics cuts for the correlated uncertainties are fixed as the same of the NNPDF4.0 global analysis.

The minimum value of Q^2 and W^2 for DIS measurements are hence chosen to be 3.49 GeV and 12.5 GeV respectively.