



# REGRESSION NETWORKS: PRECISION AND UNCERTAINTY ESTIMATION

STEFANO FORTE  
UNIVERSITÀ DI MILANO & INFN

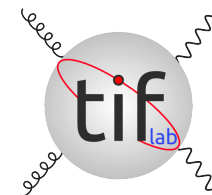
WITH TUTORIALS BY

TOMMASO GIANI

NIKHEF



UNIVERSITÀ DEGLI STUDI DI MILANO  
DIPARTIMENTO DI FISICA



## IV: VALIDATION AND TESTING

- GAUSSIANTY
  - THE “HESSIAN” PROJECTION
  - ASSESSING GAUSSIANTY
- CLOSURE TESTING
  - THE CLOSURE TEST AND ITS METRICS
  - THE NATURE OF UNCERTAINTIES
- FUTURE TESTING
  - THE IDEA
  - THE RESULTS
- CORRELATIONS
  - THE NATURE OF PDF CORRELATIONS
  - DATA-INDUCED VS. METHDOLOGY-INDUCED CORRELATIONS

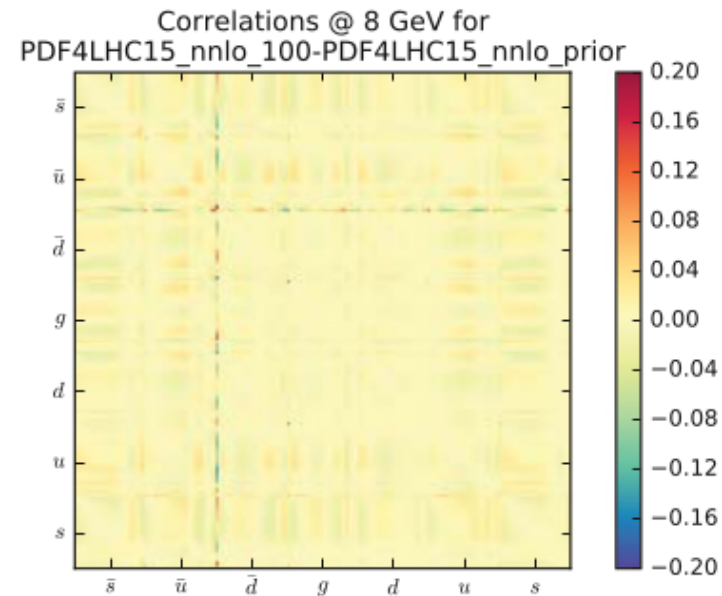
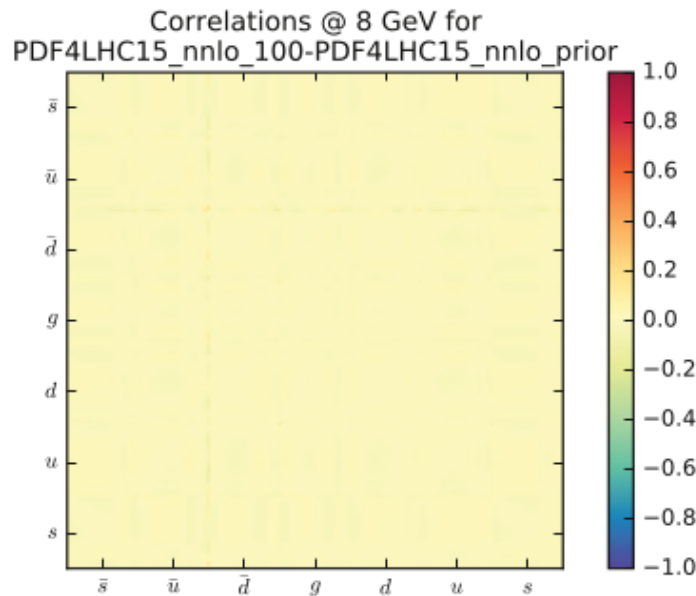
GAUSSIANITY

## MULTIGAUSSIAN REPRESENTATION

- **PARAMETRIC** REGRESSION  $\Rightarrow$  **MAP MULTIGAUSSIAN** IN PARAMETER SPACE  
“HESSIAN”  $\Rightarrow C_{ij}^{-1} = \partial_i \partial_j \chi^2$
- **HESSIAN REPRESENTATION** OF MC POSTERIOR:
  - **SAMPLE**  $k$ -TH PDF **REPLICA** OVER SET OF  $N_p$  POINTS  $f_i^{(k)}(x_j)$   
 $i$  runs over PDF flavors;  $\{ij\} = \{p\}$ ,  $p = 1, \dots, N_p \times N_f$   
 $X_{pk} = f_i^{(k)}(x_j) - f_i^{(0)}(x_j)$ ;  $f_i^{(0)}(x_j) \equiv \langle f_i^{(k)}(x_j) \rangle$  REPLICA AVERAGE
  - $C_{pp'} = \frac{1}{N_{\text{rep}}} X X^t$  (**Cholesky**)
  - $X = U S V^t$ ;  $U \Rightarrow$  EIGENVECTORS OF  $C$   $N_p \times N_{\text{rep}}$ ;  $S \Rightarrow$  DIAGONAL NONZERO EIGENVALUE SQRT MATRIX;  $V \Rightarrow$  ORTHOGONAL  $N_{\text{rep}} \times N_{\text{rep}}$  (**SVD**)
  - $C = \frac{1}{N_{\text{rep}}} X X^t = \frac{1}{N_{\text{rep}}} (U S)(U S)^t \Rightarrow$  **KEEP LARGEST** EIGENVALUES

# MULTIGAUSSIAN REPRESENTATION vs. MONTECARLO

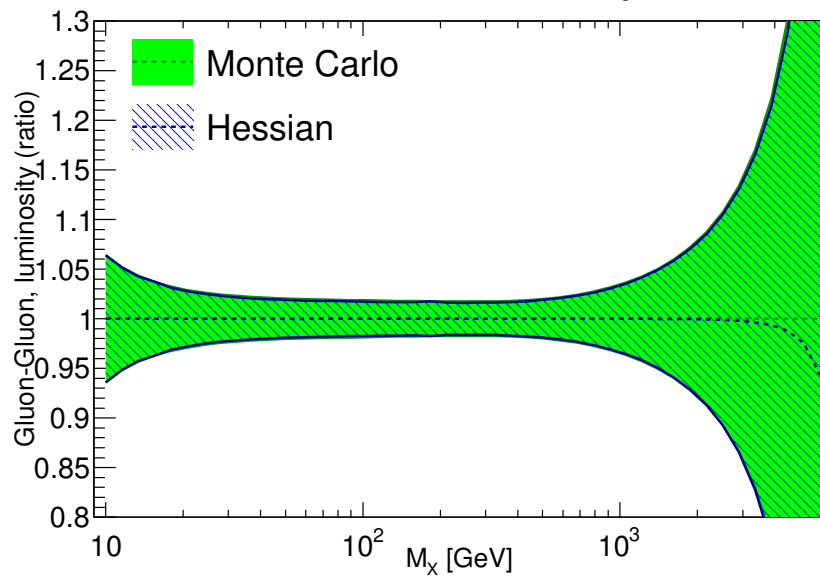
## PDF CORRELATIONS



## LUMINOSITIES

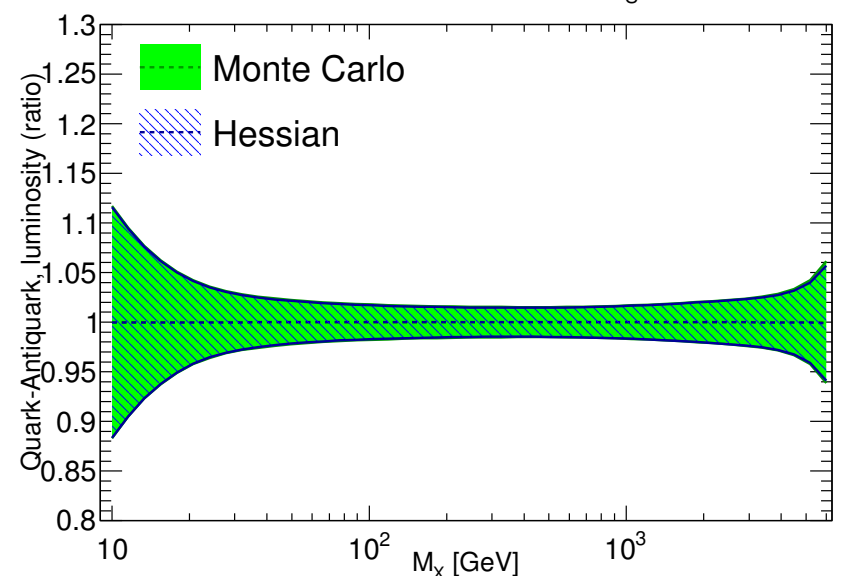
### GLUON-GLUON

LHC 13 TeV - NNPDF3.0 NLO  $\alpha_s = 0.118$



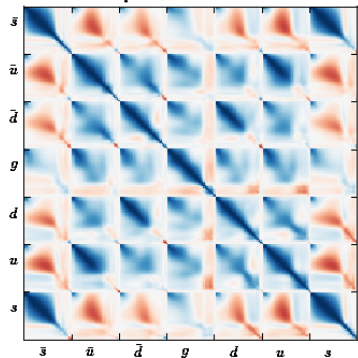
### QUARK-QUARK

LHC 13 TeV - NNPDF3.0 NLO  $\alpha_s = 0.118$

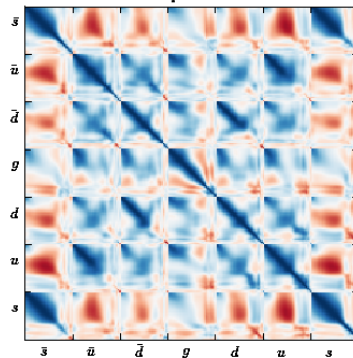


# COMPRESSION MONTECARLO

Correlations for NNPDF3.0 NLO  
Prior  $N_{rep}=1000$  @  $Q=100$  GeV



Correlations for NNPDF3.0 NLO  
Compressed  $N_{rep}=50$  @  $Q=100$  GeV

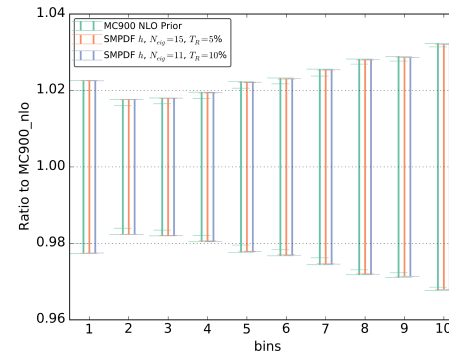


- CONSTRUCT A **VERY LARGE REPLICAS SAMPLE**
- SELECT BY **GENETIC ALGORITHM** A **SUBSET OF REPLICAS** WHOSE STATISTICAL FEATURES ARE AS CLOSE AS POSSIBLE TO THOSE OF THE PRIOR
- $\Rightarrow$  FOR ALL PDFs ON A GRID OF POINTS // **MINIMIZE DIFFERENCE** OF: FIRST FOUR MOMENTS, CORRELATIONS; OUTPUT OF KOLMOGOROV-SMIRNOV TEST (NUMBER OF REPLICAS BETWEEN MEAN AND  $\sigma$ ,  $2\sigma$ , INFINITY)
- 50 COMPRESSED REPLICAS REPRODUCE 1000 REPLICAS SET TO PRESENT ACCURACY

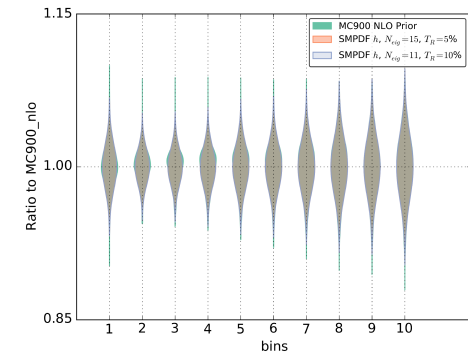
## MULTIGAUSSIAN

- SELECT **SUBSET OF THE COVARIANCE MATRIX CORRELATED** TO A GIVEN SET OF PROCESSES
- PERFORM **SVD ON THE REDUCED COVARIANCE MATRIX**, SELECT DOMINANT EIGENVECTOR, PROJECT OUT ORTHOGONAL SUBSPACE
- ITERATE UNTIL DESIRED ACCURACY REACHED
- 15 EIGENVECTORS DESCRIBE ALL HIGGS MODES + JETS +  $W$ ,  $Z$  PRODUCTION

ggh\_pt\_13tev\_larger(NLO)

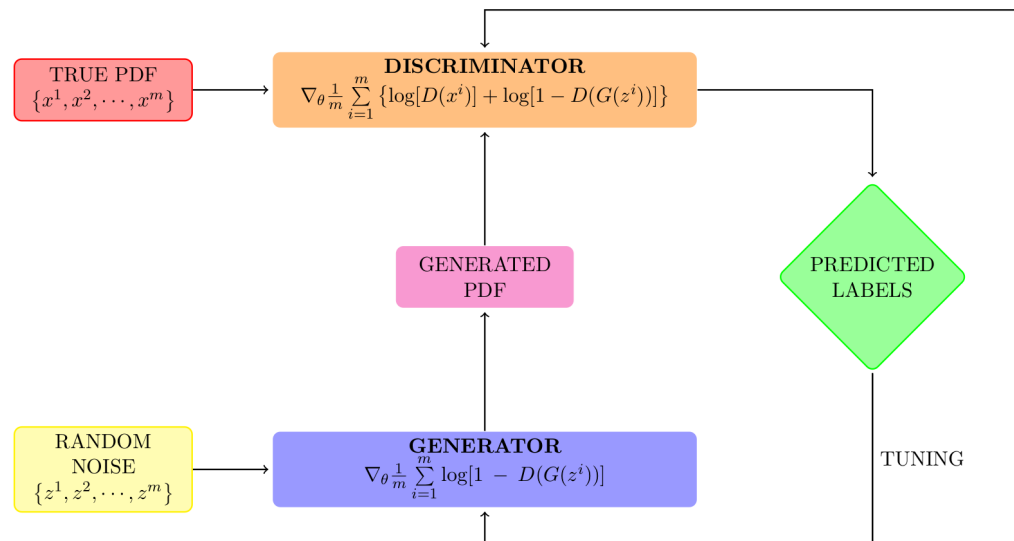


ggh\_pt\_13tev\_larger(NLO)



# GAN ENHANCEMENT

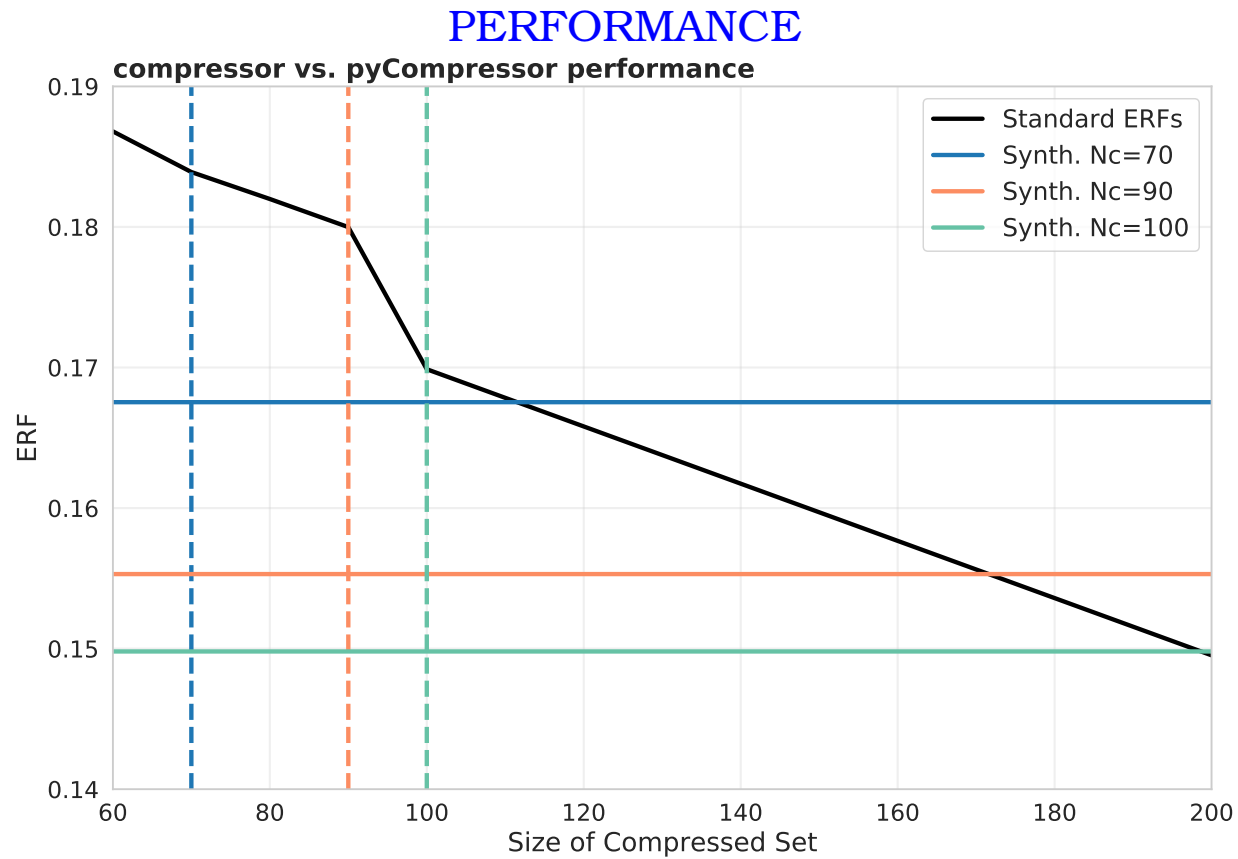
CAN WE FURTHER REDUCE THE NUMBER OF COMPRESSED REPLICAS WITHOUT LOSS OF INFORMATION? GENERATIVE ADVERSARIAL NETWORKS



- TRAIN A NETWORK TO SIMULATE THE TRUE DISTRIBUTION (**GENERATOR**)
- TRAIN A NETWORK TO DISCRIMINATE TRUTH FROM SIMULATION (**DISCRIMINATOR**)
- TRAIN THE **GENERATOR** TO TRICK THE **DISCRIMINATOR**

# GAN ENHANCEMENT

- **ENHANCE** THE STARTING PDF SET BY ADDING GAN-PDFs TO IT
- **PERFORM COMPRESSION** OF THE ENHANCED SET

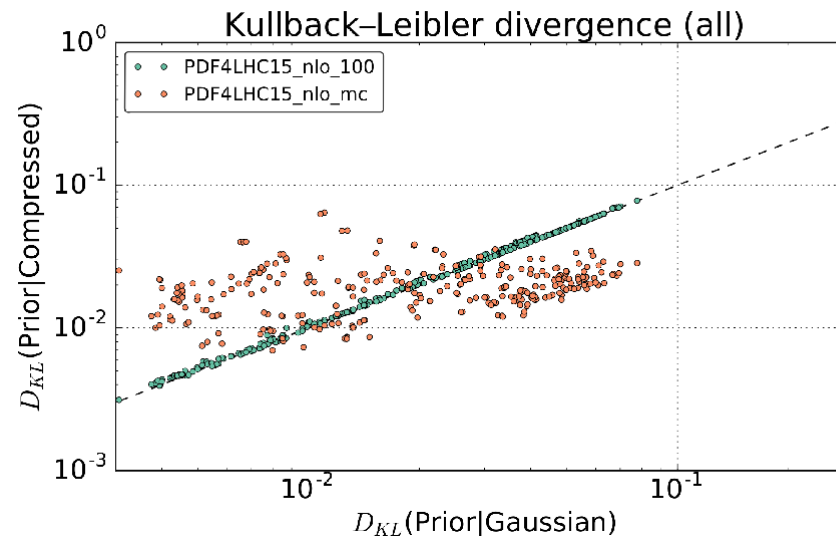


**ENHANCED:** NUMBER OF REPLICAS **CUT IN HALF** FOR SAME TARGET ACCURACY



# ARE UNCERTAINTIES GAUSSIAN?

- **REPLICA HISTOGRAM**  $i$ -TH DATAPOINT  $z_i$  FROM MC  $\Rightarrow$  **CONTINUOUS DISTRIBUTION** WITH KDE
  - POINT  $\Rightarrow$  KERNEL:  $P(z) = \frac{1}{N_{\text{rep}}} \sum_{k=1}^{N_{\text{rep}}} K(z - z_i)$ ;
  - Gaussian kernel  $K(z - z_i) \equiv \frac{1}{h\sqrt{2\pi}} \exp\left(-\frac{(z-z_i)^2}{h}\right)$
  - Silverman bandwidth  $h = \sigma_i \left(\frac{4}{3N_{\text{rep}}}\right)^{\frac{1}{5}} \Rightarrow$  MINIMIZES DIFFERENCE TO GAUSSIAN
- DEFINE **KULLBACK-LEIBLER DIVERGENCE**  
 $D_{\text{KL}} = \int_{-\infty}^{\infty} P(x) \ln \frac{P(x)}{Q(x)} dx$   
BETWEEN A PRIOR  $P$  AND ITS REPRESENTATION  $Q$
- COMPUTE  $D_{\text{KL}}$  **MC PRIOR VS REPRESENTATION** & **MC PRIOR VS GAUSSIAN**
- REPRESENTATIONS: **MULTIGAUSS** OR **MC COMPRESSION**



- $D_{\text{KL}}$  TO GAUSSIAN **GENERALLY SMALL**
- ONLY FOR **FEW POINTS** COMPRESSION **MORE EFFICIENT** THAN MULTIGAUSS CONVERSION

# PDF UNCERTAINTIES: DATA

# THE CLOSURE TEST: THE BASIC IDEA

- **POSTERIOR** REPLICA DISTRIBUTION  $\Rightarrow$  APPROXIMATELY **GAUSSIAN**
- CAN DETERMINE CONFIDENCE LEVEL OF TRUTH ABOUT PREDICTION
  - DATA SPACE
    - \* **IN SAMPLE** (USED FOR TRAINING)
    - \* **OUT OF SAMPLE** (PREDICTIONS)
  - PDF SPACE
- **MEASURABLE** BASED ON **ASSUMED UNDERLYING TRUTH** = **RUNS OF THE UNIVERSE**
- $n\sigma$  **GAUSSIAN** CONFIDENCE INTERVAL:  $\xi_{n\sigma} = \text{erf}\left(\frac{n\sigma}{\sqrt{2}\sigma^0}\right)$   
 $n\sigma \rightarrow$  CONFIDENCE INTERVAL;  $\sigma^0 \Rightarrow$  WIDTH OF GAUSSIAN

## STATISTICAL INDICATORS

- **BIAS**  $b = \frac{1}{N_{\text{points}}} \sum_{i=1}^{N_{\text{points}}} (\mathcal{G}_i(f) - z_i)^2$   
 $(\mathcal{G}_i(f) = \langle \mathcal{G}_i(f) \rangle \equiv \frac{1}{N_{\text{replicas}}} \sum_{j=1}^{N_{\text{replicas}}} \mathcal{G}_i(f_j)$  prediction,  $z_i$  true);  
**NORMALIZED:**  $\frac{1}{N_{\text{points}}} |\mathcal{G}(f) - z|_C^2$ ,  
**C COVARIANCE MATRIX:**
  - DATA  $\Rightarrow$  FROM EXPERIMENT
  - PDF  $\Rightarrow$  FROM REPLICAS
- **VARIANCE**  $v = \frac{1}{N_{\text{points}}} \sum_{i=1}^{N_{\text{points}}} \sigma_i^2$ ;  $\sigma_i = \langle (\mathcal{G}_i(f) - \langle \mathcal{G}_i(f) \rangle)^2 \rangle$ ; **NORMALIZED**  
 $v = |\mathcal{G}(f) - \langle \mathcal{G}(f) \rangle|_C^2$
- **BIAS-VARIANCE RATIO**  $R_{bv} = \sqrt{\frac{b}{v}}$ : **AVERAGED** OVER RUNS OF THE UNIVERSE (**RUS**)
- **EMPIRICAL CONFIDENCE LVL**  $\xi_{n\sigma} = \frac{1}{N_{\text{points}}} \sum_{i=1}^{N_{\text{points}}} I_{[-n\sigma, n\sigma]} (\langle \mathcal{G}_i(f) \rangle - z_i)$  **OVER RUS**

# CLOSURE TEST IMPLEMENTATION

- ASSUME UNDERLYING “TRUTH” PDF (SAY A RANDOM PDF REPLICA)
- GENERATE DATA ACCORDING TO STATISTICAL AND CORRELATED SYSTEMATICS (SAY FOR NNPDF4.0 DATASET)
- DETERMINE PDFs & COMPARED TO “TRUTH” BASED ON INDICATORS

## THE NATURE OF UNCERTAINTIES

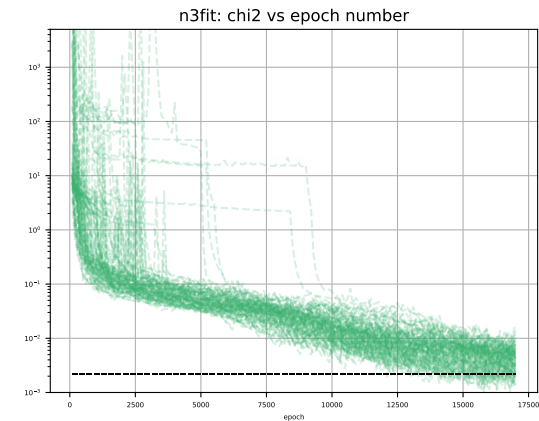
- LEVEL 0:
  - EACH DATAPOINT EQUAL TO THE “TRUTH VALUE”; ZERO UNCERTAINTY
  - FIT → MUST FIND  $\chi^2 = 0$  (GET BACK “TRUTH”)
  - $\chi^2 \approx 0$  BOTH REPLICA TO REPLICA AND AVERAGE TO TRUTH
  - INTERPOLATION/EXTRAPOLATION UNCERTAINTY
- LEVEL 1:
  - EACH PSEUDO- DATAPOINT IS OBTAINED AS A RANDOM FLUCTUATION WITH GIVEN COVARIANCE MATRIX ABOUT “TRUTH”  
⇒ “RUN OF THE UNIVERSE”
  - FIT DATA OVER AND OVER AGAIN
  - $\chi^2 \approx 1$  BOTH REPLICA TO REPLICA AND AVERAGE TO TRUTH
  - FUNCTIONAL UNCERTAINTY
- LEVEL 2:
  - DATA AS IN LEVEL 1
  - GENERATE DATA REPLICAS OF THESE “DATA”
  - FIT PDF REPLICAS TO DATA REPLICAS
  - $\chi^2 \approx 2$  REPLICA TO REPLICA;  $\chi^2 \approx 1$  AVERAGE TO TRUTH
  - DATA UNCERTAINTY

# UNCERTAINTIES: TYPE AND SIZE

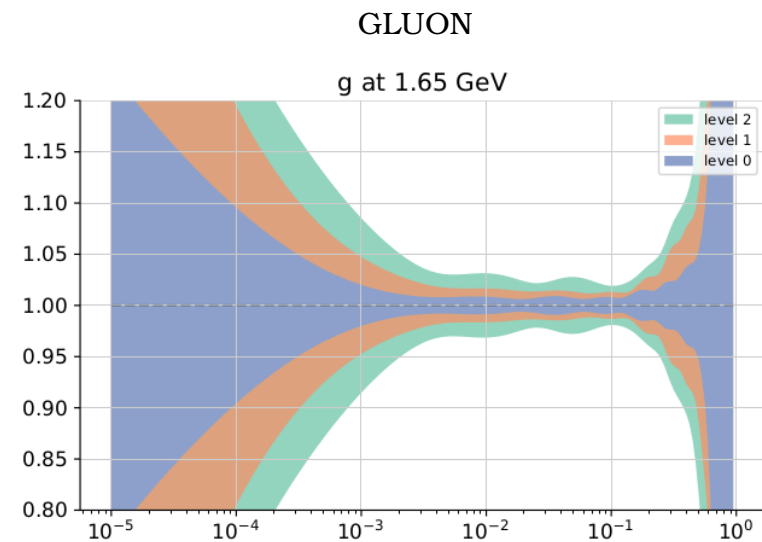
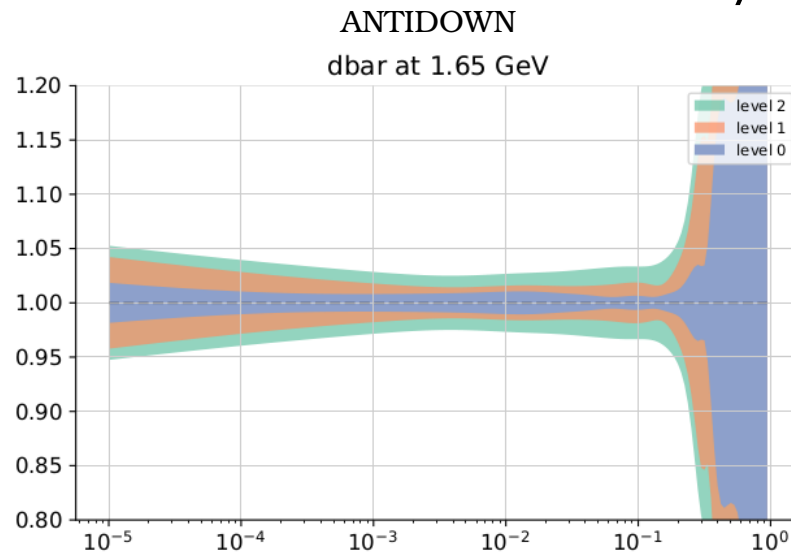
## CLOSURE TEST RESULTS (NNPDF4.0)

- **LEVEL 0** (TRUTH DATA)  $\Rightarrow \chi^2 \approx 0$ , YET **UNCERTAINTY NONZERO**  
 $\Rightarrow$  NEURAL NETS  $\Leftrightarrow$  **MANY FUNCTIONAL FORMS**
- **LEVEL 1** (RUNS OF UNIVERSE)  $\Rightarrow$  REPLICAS ALL FITTED TO SAME DATA, YET **UNCERTAINTY NONZERO**  
 $\Rightarrow$  **DITTO**
- **LEVEL 0, 1 AND 2 UNCERTAINTIES COMPARABLE IN SIZE**

LEVEL 0  $\chi^2$  VS TRAINING



## LEVEL 0/1/2 UNCERTAINTIES



# TESTING: THE INDICATORS

## BIAS/VARIANCE RATIO AND ONE- $\sigma$ QUANTILE

DATA-SPACE, DATA COVARIANCE MATRIX, OUT-OF-SAMPLE

PDF-SPACE & COV MATRIX

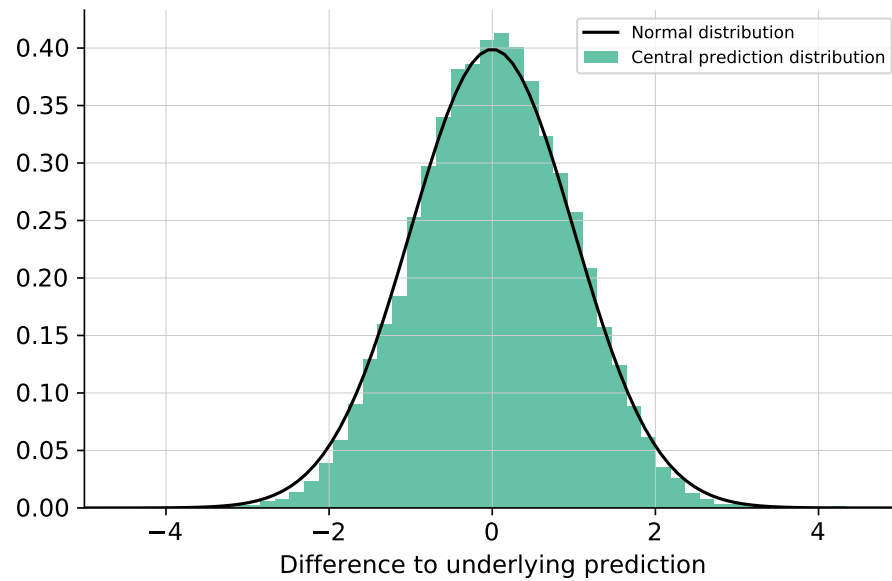
Dataset	$\sqrt{b/v}$	$\xi_{1\sigma}^{(\text{data})}$	$\text{erf}(R_{bv}/\sqrt{2})$	flavour	$\xi_{1\sigma}^{(\text{pdf})}$
DY	$0.99 \pm 0.08$	$0.69 \pm 0.02$	$0.69 \pm 0.04$	$\Sigma$	$0.82 \pm 0.04$
Top-pair	$0.75 \pm 0.06$	$0.75 \pm 0.03$	$0.82 \pm 0.03$	$g$	$0.70 \pm 0.05$
Jets	$1.14 \pm 0.05$	$0.63 \pm 0.03$	$0.62 \pm 0.02$	$V$	$0.65 \pm 0.05$
Dijets	$0.99 \pm 0.07$	$0.70 \pm 0.03$	$0.69 \pm 0.04$	$V_3$	$0.63 \pm 0.05$
Direct photon	$0.71 \pm 0.06$	$0.81 \pm 0.03$	$0.84 \pm 0.03$	$V_8$	$0.72 \pm 0.04$
Single top	$0.87 \pm 0.07$	$0.69 \pm 0.04$	$0.75 \pm 0.04$	$T_3$	$0.71 \pm 0.05$
Total	$1.03 \pm 0.05$	$0.68 \pm 0.02$	$0.67 \pm 0.03$	$T_8$	$0.71 \pm 0.05$
				Total	$0.71 \pm 0.02$

- 25 “UNIVERSE RUNS”, 45 REPLICAS EACH
- IN-SAMPLE DATA: PRE 2015
- OUT OF SAMPLE DATA: 2015-2020, MOSTLY LHC
- PDFs HIGHLY CORRELATED  $\Rightarrow$  SAMPLED AT 4 POINTS EACH

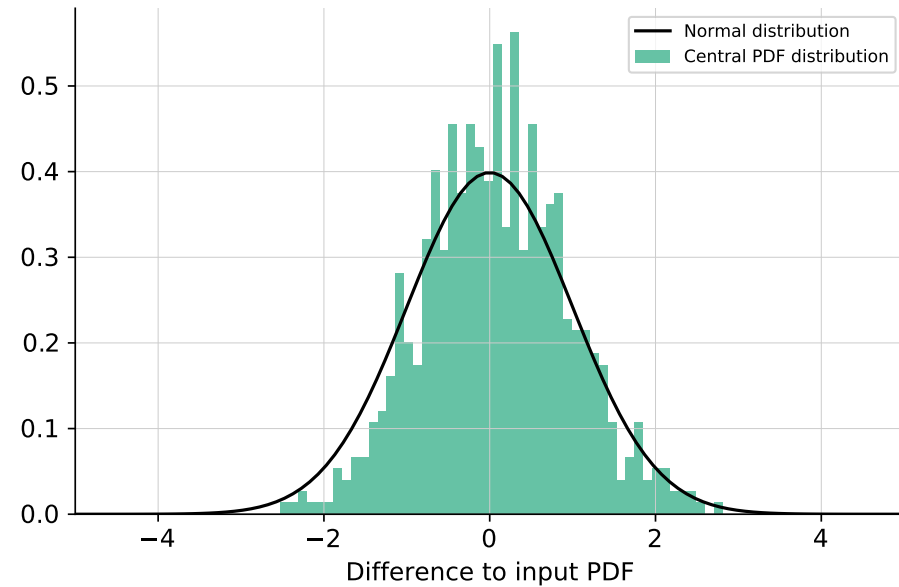
# TESTING: THE INDICATORS

## DISTRIBUTION OF DEVIATIONS FROM TRUTH

DATA SPACE (OUT OF SAMPLE)



PDF SPACE



- PDF-SPACE MORE NOISY THAN DATA SPACE

## ASIDE: ERRORS IN MC ESTIMATES

### THE JACKKNIFE/BOOTSTRAP METHOD

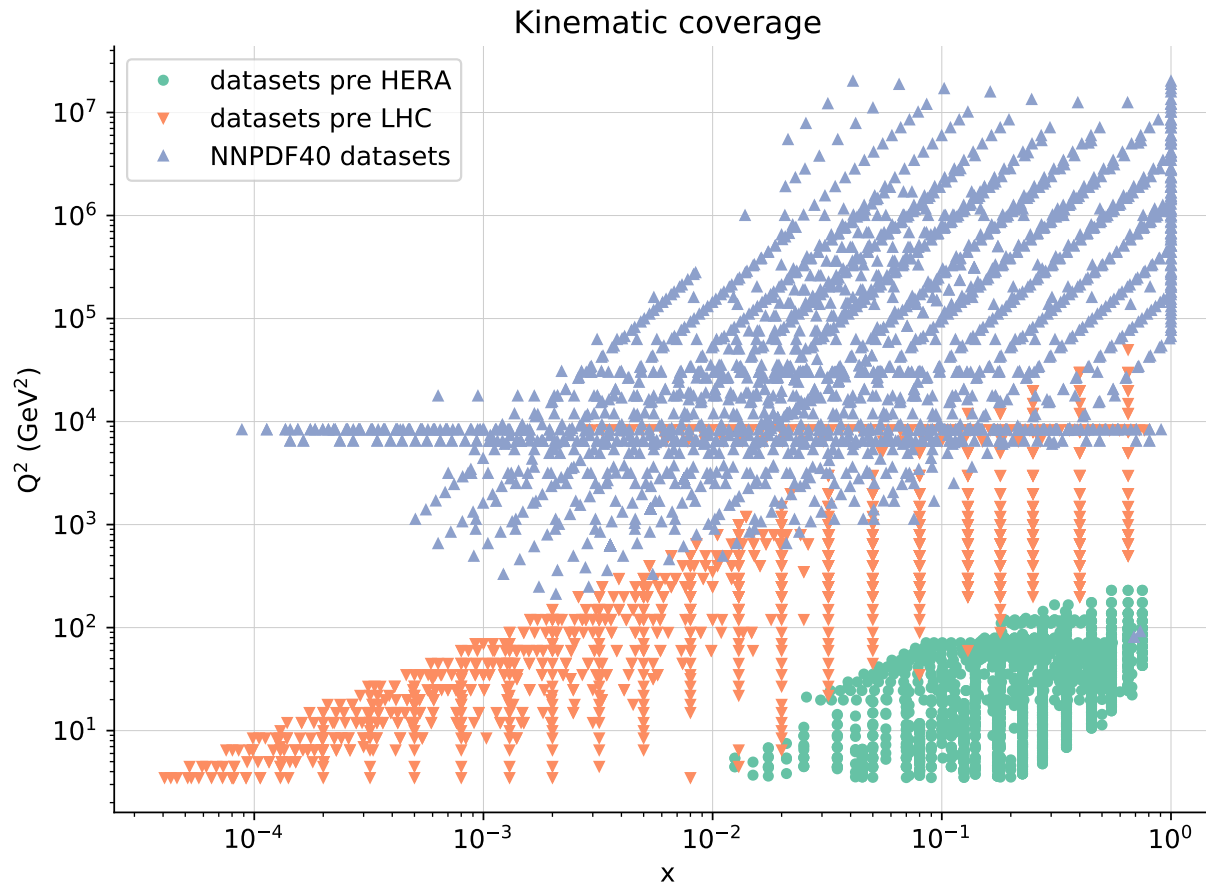
- GIVEN  $N_{\text{est}}$  ESTIMATES  $x^i$  OF  $x$ , COMBINED ESTIMATE  $x = \langle x \rangle \pm \sigma$ ,  
 $\langle x \rangle = \frac{1}{N_{\text{est}}} \sum_{i=1}^{N_{\text{est}}} x^i$ ;  $\sigma^2 = \langle (x - \langle x \rangle)^2 \rangle$ .
- GIVEN A MC SAMPLE OF ESTIMATES, **EXTRACT RANDOMLY**  $n < N_{\text{est}}$  FROM IT WITH REPLACEMENT  $\Rightarrow$  EXTRACTIONS INDEPENDENT, REPETITIONS ALLOWED  
**JACKKNIFE**:  $n = N_{\text{est}} - 1$
- **REPEAT** EXTRACTION  $N_B$  TIMES  $\Rightarrow N_b$  SAMPLES OF  $n$  REPLICAS
- **COMPUTE AVERAGE** FOR EACH EXTRACTED  $n$ -REPLICA **SAMPLE**
- AVERAGE OF THESE EQUAL TO THE STARTING SAMPLE AVERAGE
- COMPUTE  $\langle x \rangle$  FROM FULL SAMPLE, ESTIMATE UNCERTAINTY ON IT FROM **VARIANCE OF BOOTSTRAP** EXTRACTIONS



# PDF UNCERTAINTIES: EXTRAPOLATION

# FUTURE TESTS

IDEA: USE (REAL) HIERARCHICAL DATASETS



- DEFINE “PRE-HERA”, “PRE-LHC” AND “CURRENT” DATASETS  
EACH LATER DATASET IS EXTRAPOLATION OF PREVIOUS
- DETERMINE PDFs & COMPARE TO “FUTURE” DATA
- COMPUTE  $\chi^2$  TO FUTURE DATA:
  - WITHOUT PDF UNCERTAINTIES  $\Rightarrow$  IF  $\gg 1$ , MISSING INFORMATION
  - WITH PDF UNCERTAINTY  $\Rightarrow$  IF  $\sim 1$ , TEST PASSED  
MISSING INFO REPRODUCED BY UNCERTAINTY

# ASSESSING EXTRAPOLATION UNCERTAINTIES

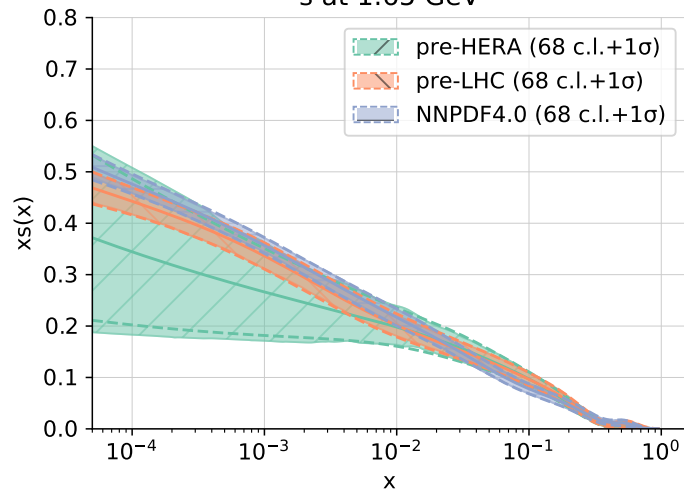
## FUTURE TEST RESULTS (NNPDF4.0)

$\chi^2$ : FITTED VS EXTRAPOLATED: **WITHOUT**/**WITH** PDF UNC.

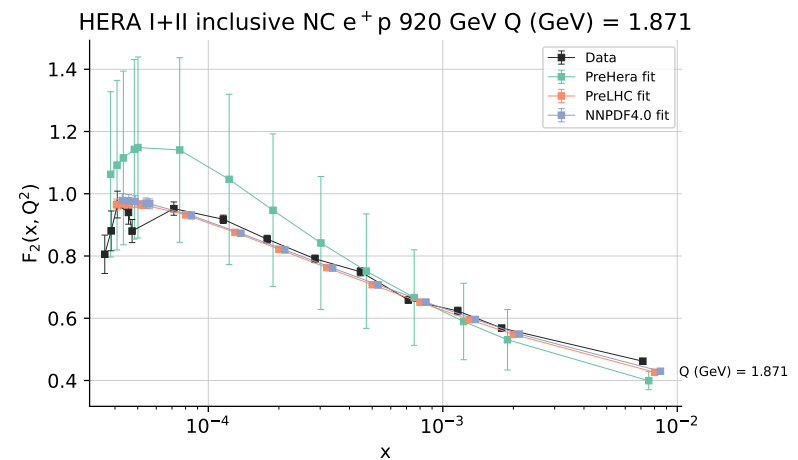
PROCESS	PRE-HERA	PRE-LHC	NNPDF4.0
FT DIS (NC)	1.05	1.18	1.23
FT DIS (CC)	0.80	0.85	0.87
FT DY	0.92	1.27	1.59
HERA	<b>27.20/1.23</b>	1.22	1.20
COLL. DY (TEV.)	<b>5.52/1.02</b>	0.99	1.11
COLL. DY (LHC)	<b>18.91/1.31</b>	<b>2.63/1.58</b>	1.53
TOP QUARK	<b>20.01/1.06</b>	<b>1.30/0.87</b>	1.01
JETS	<b>2.69/0.98</b>	<b>2.12/1.10</b>	1.26
TOTAL OUT OF SAMPLE	<b>19.48/1.16</b>	<b>2.10/1.15</b>	-

### strange PDF

s at 1.65 GeV



### HERA $F_2^p$



PDFs ARE FUTURE-COMPATIBLE!

# PDF CORRELATIONS

# CORRELATION BETWEEN MODEL FEATURES

example: up vs down PDFs

COVARIANCE:  $\text{Cov}[u, d](x, x') = \langle u(x, Q_0^2)d(x', Q_0^2) \rangle - \langle u(x, Q_0^2) \rangle \langle d(x', Q_0^2) \rangle$ ;

CORRELATION:  $\rho[u, d](x, x') = \frac{\text{Cov}[u, d](x, x')}{\sqrt{\text{Var}[u](x)\text{Var}[d](x')}}$

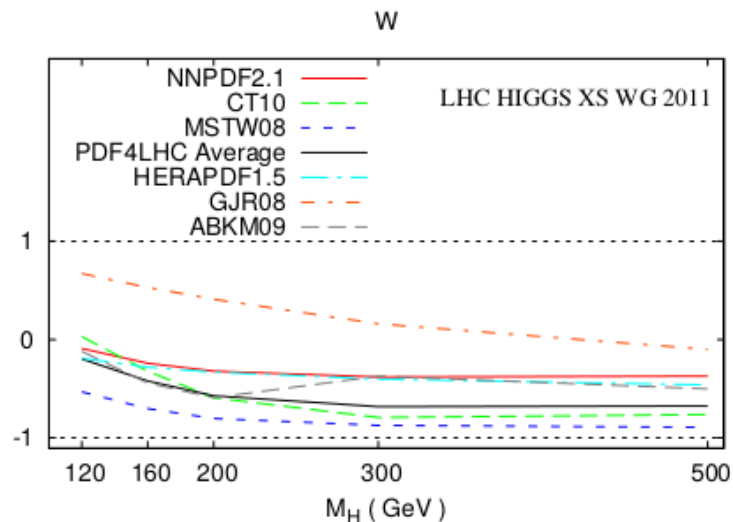
COMPUTATION IN MC APPROACH:  $\langle u(x, Q_0^2)d(x', Q_0^2) \rangle = \frac{1}{N} \sum_{r=1}^N u^{(r)}(x, Q_0^2)d^{(r)}(x', Q_0^2)$ ;  
 $u^{(r)}(x, Q_0^2)$  REPLICAS

- CORRELATION INDUCED BY DATA, THEORY (E.G. SUM RULES), METHODOLOGY (E.G. ASSUMPTIONS ON EXTRAPOLATION)
- USED E.G. TO ASSESS CORRELATION BETWEEN SIGNAL AND BACKGROUND PROCESSES

## PDF-INDUCED CORRELATIONS

BETWEEN HIGGS SIGNAL & BACKGROUND PROCESSES (HXSWG, YR2, 2011)

Higgs in gluon fusion vs. W production



# CORRELATIONS BETWEEN MODELS

CORRELATE PDFS IN **DIFFERENT SETS**

example: up NN model vs down parametric model

$$\text{Cov}[u^N, d^P](x, x') = \langle u^N(x, Q_0^2) d^P(x', Q_0^2) \rangle - \langle u^N(x, Q_0^2) \rangle \langle d^P(x', Q_0^2) \rangle$$

**S-CORRELATION** VS **F-CORRELATION**

$\rho[u^N, u^P]$  **DIFFERENT SETS, SAME PDF** VS.  $\rho[u^N, d^N]$  **SAME SET, DIFFERENT PDFS**

- SAME REPLICA **MUST BE USED** FOR NONZERO CORRELATION:

**IF REPLICAS UNCORRELATED**  $\langle u(x, Q_0^2) d(x, Q_0^2) \rangle \stackrel{?}{=} \frac{1}{N} \sum_{r=1}^N u^{(r)}(x, Q_0^2) d^{(r)}(x, Q_0^2) = \langle u \rangle \langle d \rangle$

**THEN** CORRELATION VANISHES

## REPLICA CORRELATION

- FIT PDF REPLICAS  $f_i^{(r, N)}(x, Q_0^2)$  &  $f_i^{(r, P)}(x, Q_0^2)$  for all  $x, i$  TO **SAME DATA REPLICA**
- COMPUTE COVARIANCE & CORRELATION USING

$$\langle u(x, Q_0^2) d(x, Q_0^2) \rangle = \frac{1}{N} \sum_{r=1}^N u^{(r, N)}(x, Q_0^2) d^{(r, P)}(x, Q_0^2)$$

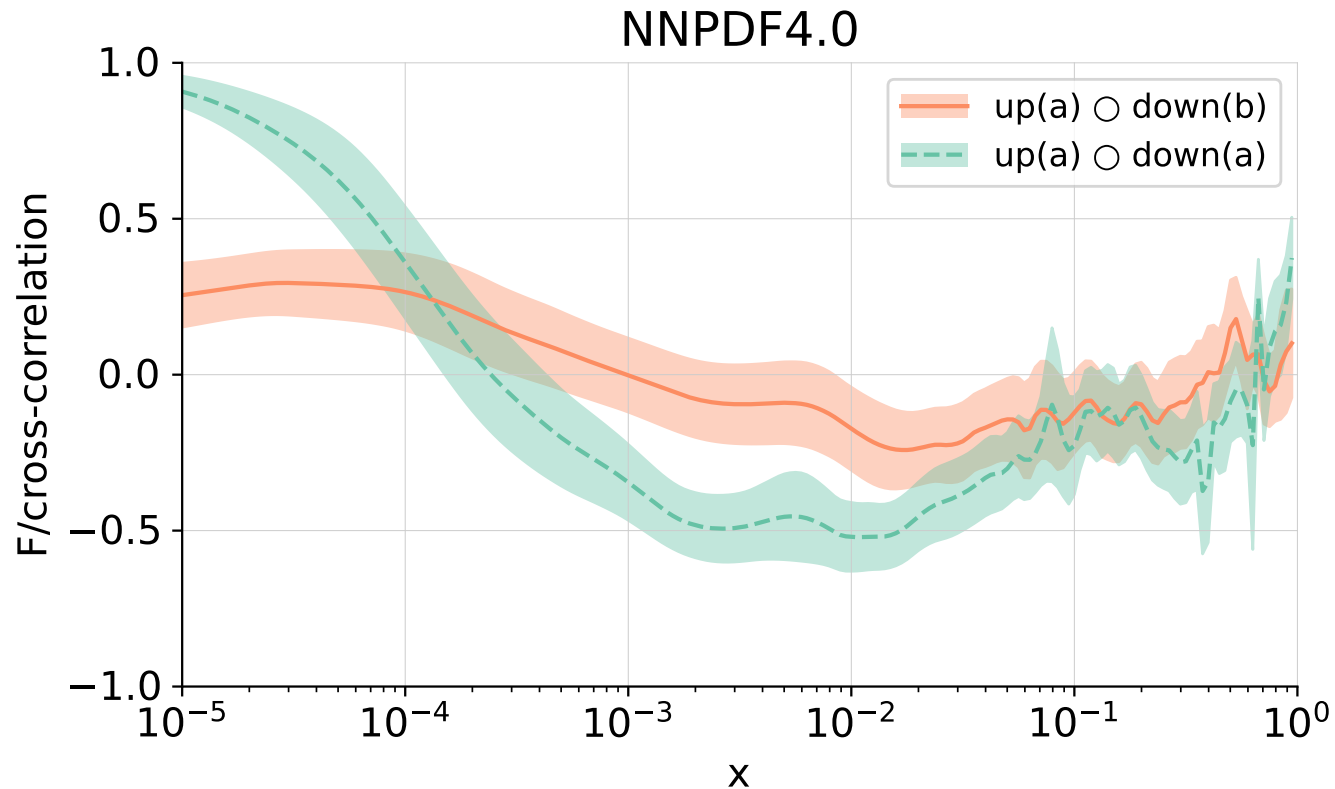
# DATA vs METHODOLOGY CORRELATION

- **NONZERO LEVEL-1** UNCERTAINTY  $\Rightarrow$  **DATA REPLICA DOES NOT DETERMINE** UNIQUELY THE PDF REPLICA
- **IN PRINCIPLE** FULL CORRELATION:  $r \Leftrightarrow$  **DATA REPLICA** AND  $r' \Leftrightarrow$  **LEVEL-1 (METHODOLOGY)** REPLICAS  
 REPLICAS (UP QUARK)  $u^{(r,r')}(x, Q_0^2)$ ;

$$\left| \frac{1}{N} \sum_{r=1}^N u^{(r,r')}(x, Q_0^2) d^{(r,r')}(x, Q_0^2) - \langle u \rangle \langle d \rangle \right| \leq \left| \frac{1}{NM} \sum_{r=1}^N \sum_{r'=1}^M u^{(r,r')}(x, Q_0^2) d^{(r,r')}(x, Q_0^2) - \langle u \rangle \langle d \rangle \right|$$

- **IN PRACTICE METHODOLOGY** CORRELATION **NOT INCLUDED**  $\Rightarrow$  CORRELATION LOSS

**FULL VS DATA-INDUCED**



# MEASURING METHODOLOGY DECORRELATION

- SELF-CORRELATION: **S-CORRELATION OF A PDF SET TO ITSELF**  
= **F-CORRELATION OF A PDF TO ITSELF**
- USE **TWO DIFFERENT SETS** OF PDF REPLICAS FITTED TO  
THE **SAME DATA REPLICAS**

$$\langle u(x, Q_0^2)u(x, Q_0^2) \rangle = \frac{1}{N} \sum_{r=1}^N u^{(r, r')}(x, Q_0^2)u^{(r, r'')}(x, Q_0^2)$$

- **DEVIATION OF CORRELATION FROM 100%** MEASURES THE  
**CORRELATION LOSS**  $\Rightarrow$  UNCORRELATED FUNCTIONAL UNCERTAINTY
- **HIGHER CORRELATION**  $\Rightarrow$  **MORE EFFICIENT METHODOLOGY**

