

# Bayesian inference and gaussian processes for PDF determination

Tommaso Giani

In collaboration with A. Candido, L.Del Debbio and G.Petrillo

Nikhef



DIS2024, Grenoble,  
10/04/2024

# $\alpha_s$ from Z pT

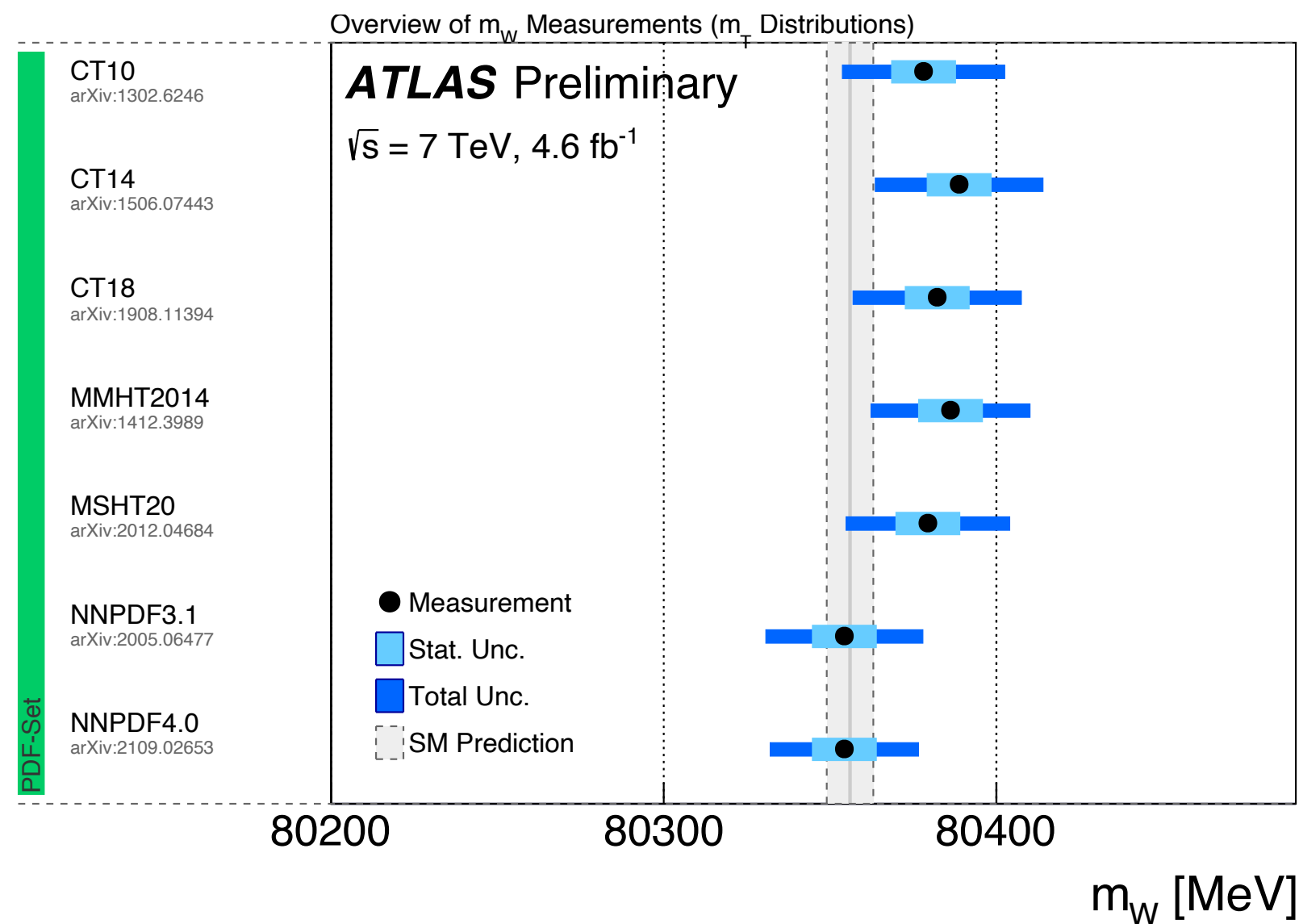
[arXiv:2309.12986](https://arxiv.org/abs/2309.12986)

$$\alpha_s(m_Z) = 0.11847 + 0.00091 - 0.00088$$

$\sim 0.76\%$

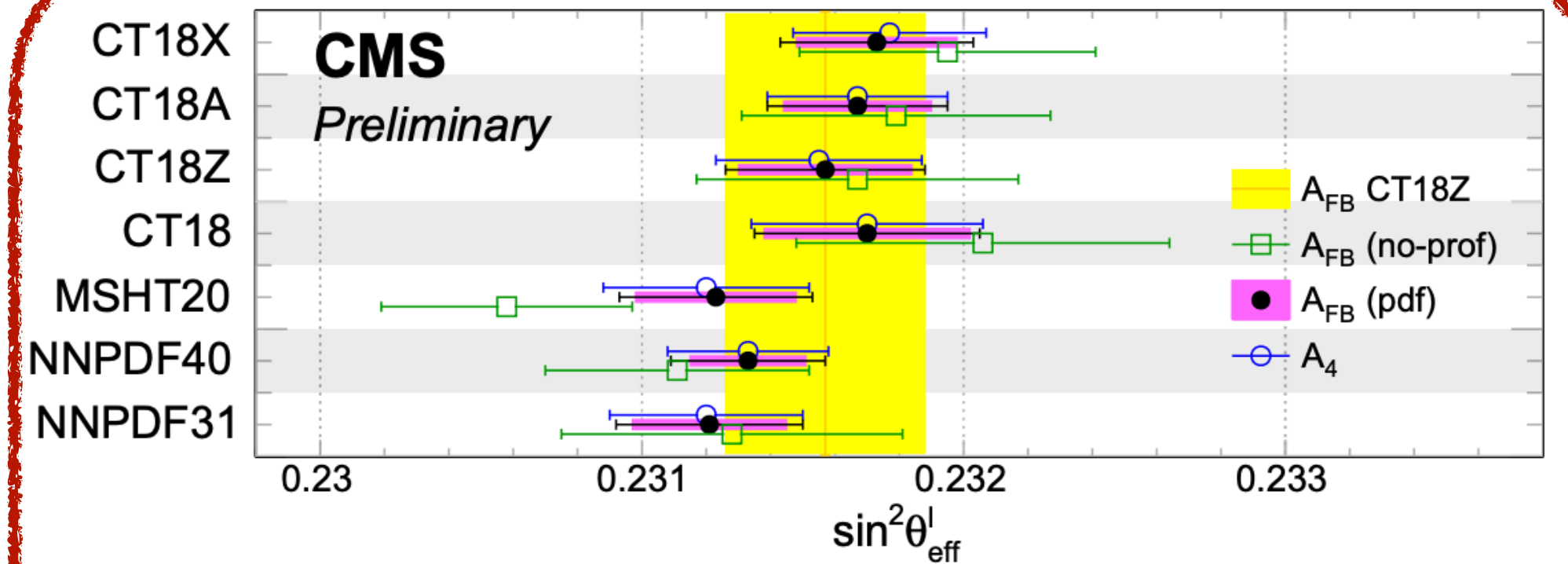
PDF set	$\alpha_s(m_Z)$	PDF uncertainty	$g [GeV^2]$	$q [GeV^4]$
MSHT20 [37]	0.11839	0.00040	0.44	-0.07
NNPDF4.0 [84]	0.11779	0.00024	0.50	-0.08
CT18A [29]	0.11982	0.00050	0.36	-0.03
HERAPDF2.0 [65]	0.11890	0.00027	0.40	-0.04

$\sim 1.7\%$



## W mass determination

[ATLAS-CONF-2023-004](#)



## weak mixing angle at 13 TeV

[CMS-PAS-SMP-22-010](#)

# $\alpha_s$ from Z pT

[arXiv:2309.12986](https://arxiv.org/abs/2309.12986)

$$\alpha_s(m_Z) = 0.11847 + 0.00091 - 0.00088$$

PDF set	$\alpha_s(m_Z)$	PDF uncertainty	$g [GeV^2]$	$q [GeV^4]$
MSHT20 [37]	0.11839	0.00040	0.44	-0.07
NNPDF4.0 [84]	0.11779	0.00024	0.50	-0.08
CT18A [29]	0.11982	0.00050	0.36	-0.03
HERAPDF2.0 [65]	0.11890	0.00027	0.40	-0.04

$\sim 1.7\%$

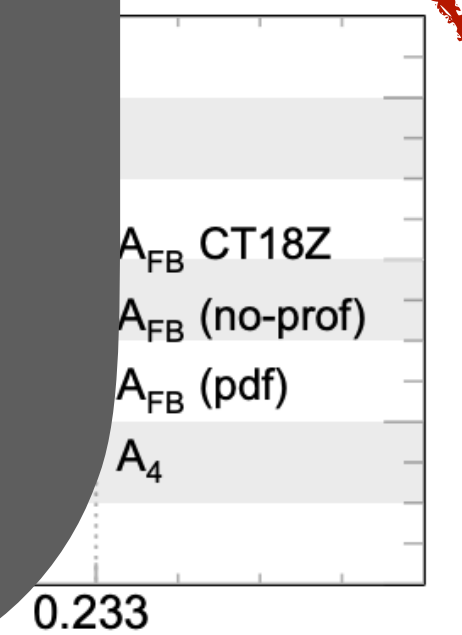
## Uncertainties in PDF determinations

- Experimental (from the data)
- Theoretical (missing higher order for theory predictions, ...)
- Input SM parameters
- Methodological

PDF-Set

- CT10  
arXiv:1302.6246
- CT14  
arXiv:1506.07443
- CT18  
arXiv:1908.11394
- MMHT2014  
arXiv:1412.3989
- MSHT20  
arXiv:2012.04684
- NNPDF3.1  
arXiv:2005.06477
- NNPDF4.0  
arXiv:2109.02653

802

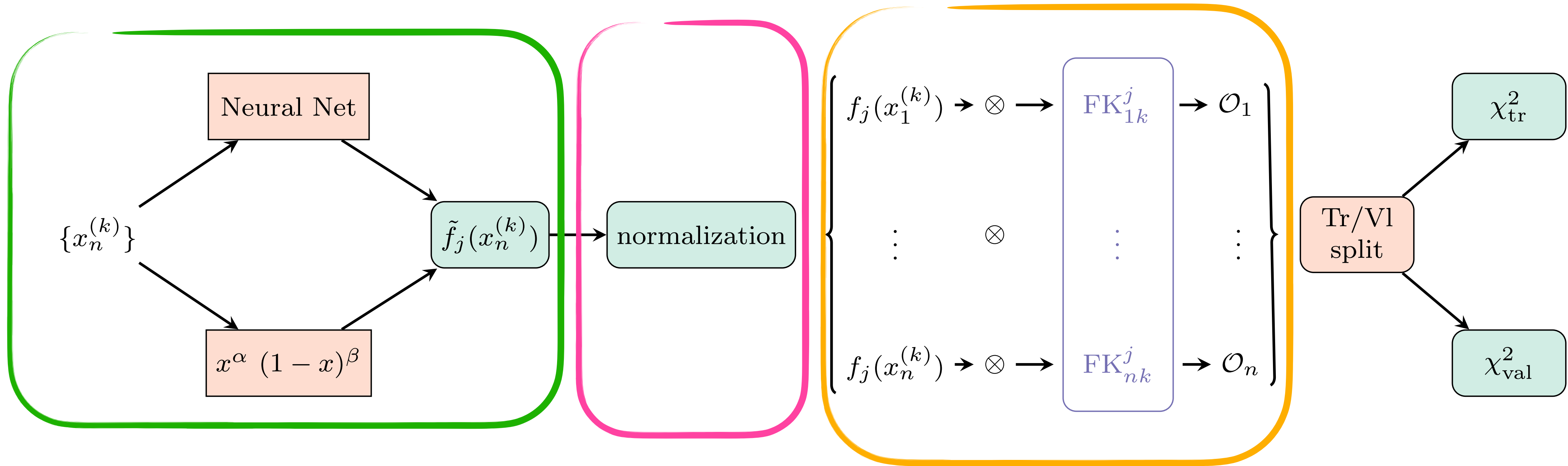


## weak mixing angle at 13 TeV

CMS-PAS-SMP-22-010

# Parametric regression

Build theory predictions for observables entering the fit



PDFs are parametrised at some initial scale  $Q_0 = 1.65$  GeV. Sum rules are imposed with suitable normalisation

Use data to build  $\chi^2$  and minimise

# Bayesian approach

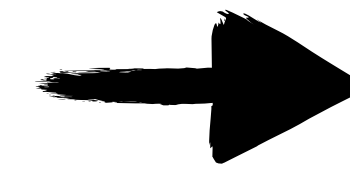
- Start from a prior on the model  $p(f)$
- Look at the data
- Get the posterior  $p(f|D)$

$$p(f|D) = \frac{p(D|f)p(f)}{p(D)}$$

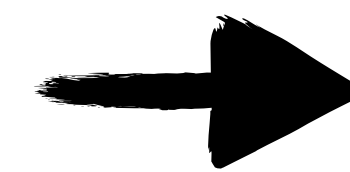
The equation is enclosed in a yellow rounded rectangle. The term  $p(f|D)$  is circled in green, and the term  $p(f)$  in the numerator is circled in pink.

Prior on the model

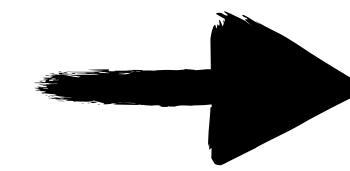
Posterior of model given the data



Introduce probability distribution on a space of functions



Build a suitable prior



Use Bayes' theorem

# Gaussian Processes

$$\mathbf{f} = \begin{pmatrix} f(x_1) \\ \vdots \\ f(x_N) \end{pmatrix} \in \mathbb{R}^N$$

**Parameters  $\mathbf{f}$** : stochastic variables representing values of the PDF on a grid of points

**Kernel  $\mathbf{K}$  and mean function  $\mathbf{m}$** : functions modelling the correlation between parameters

$$m(x_i; \theta) = \mathbb{E} \left( f(x_i) \right)$$

$$k(x_i, x_j; \theta) = \text{COV} \left( f(x_i), f(x_j) \right)$$

**Hyperparameters  $\theta$** : set of parameters entering the definition of the kernel (they control some specific feature of the prior)

**Joint probability distribution of  $\mathbf{f}$  and  $\theta$** : target of the analysis

$$p(\mathbf{f}, \theta | \text{data})$$

## Some examples of application of GPs in physics

### Gaussian process models—I. A framework for probabilistic continuous inverse theory FREE

Andrew P Valentine ✉, Malcolm Sambridge

*Geophysical Journal International*, Volume 220, Issue 3, March 2020, Pages 1632–1647,  
<https://doi.org/10.1093/gji/ggz520>

### Reconstructing QCD spectral functions with Gaussian processes

Jan Horak, Jan M. Pawłowski, José Rodríguez-Quintero, Jonas Turnwald, Julian M. Urban, Nicolas Wink, and Savvas Zafeiropoulos

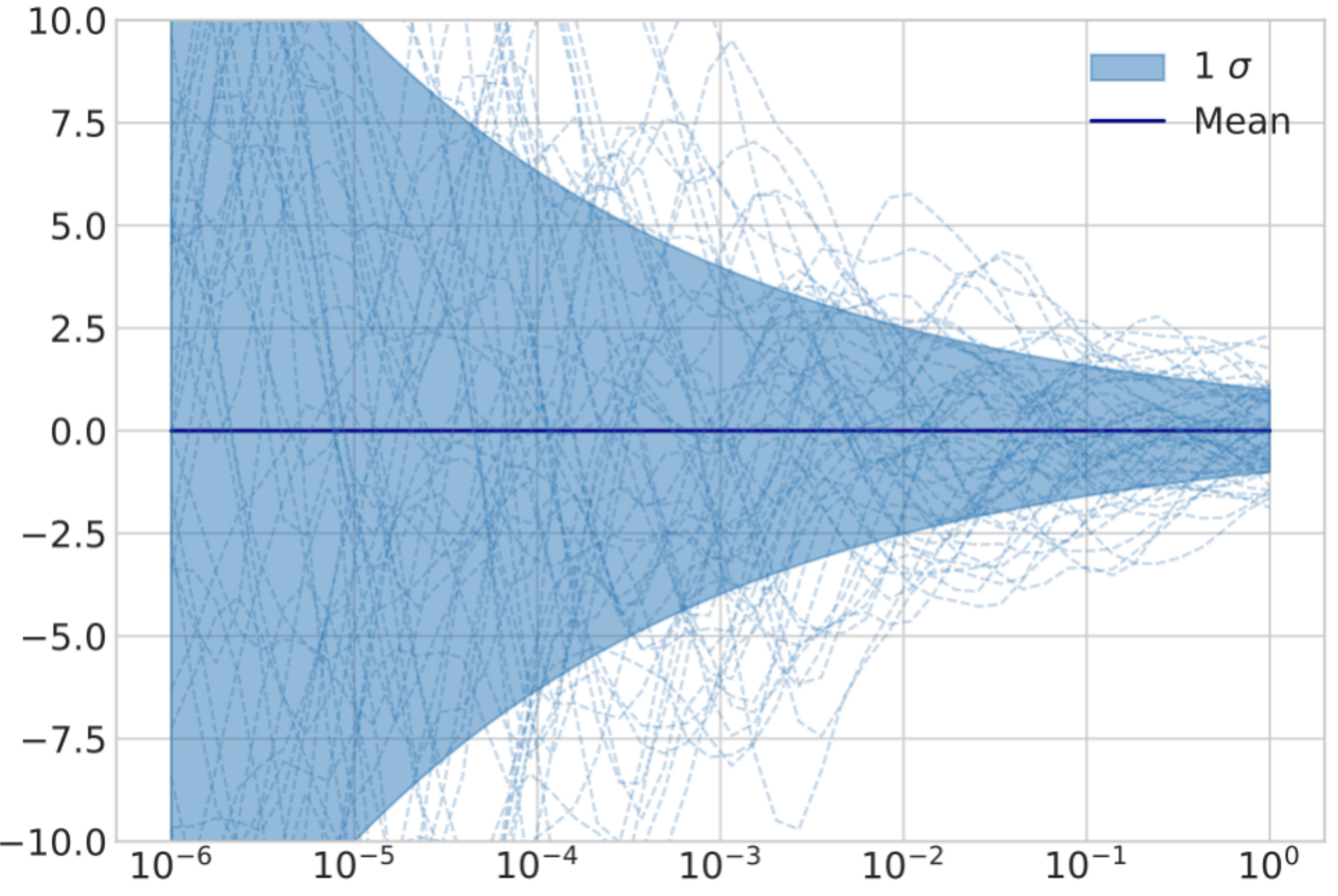
Phys. Rev. D **105**, 036014 – Published 23 February 2022

## What about PDFs?

# Prior for PDF

## Gibbs Kernel

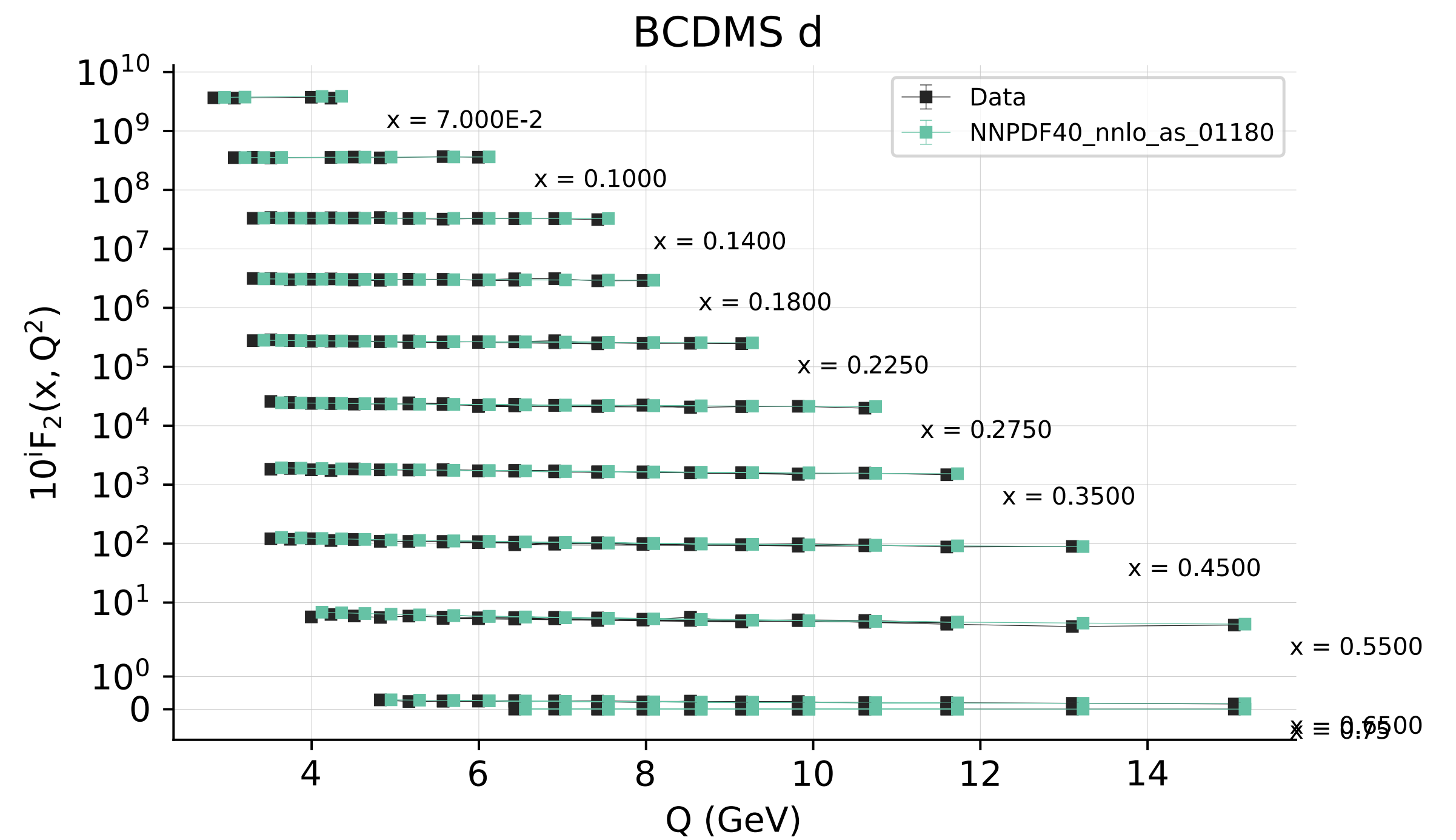
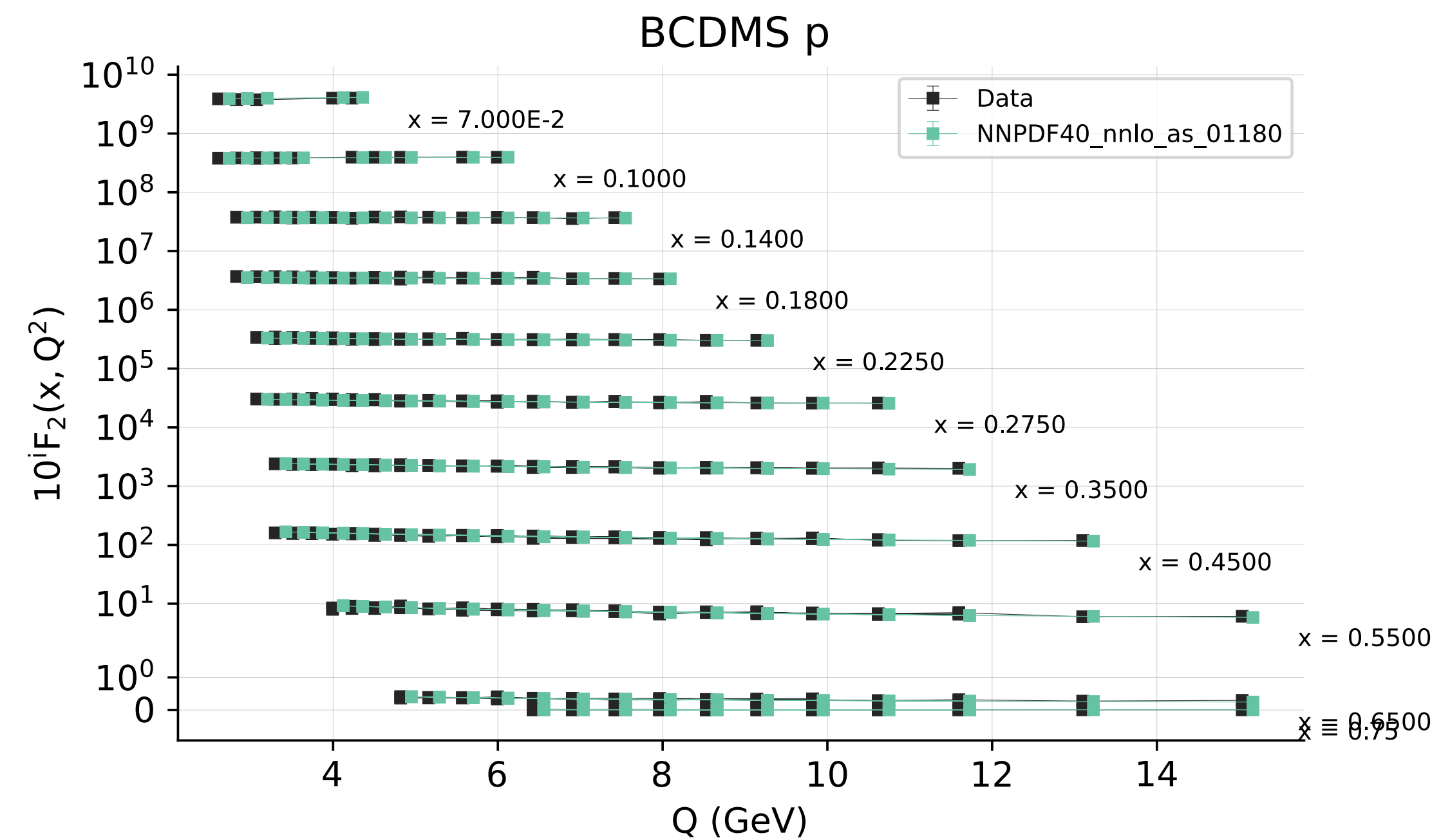
$$\tilde{k}(x, y) = x^\alpha y^\alpha \sigma^2 \sqrt{\frac{2l(x)l(y)}{l^2(x) + l^2(y)}} \exp\left[-\frac{(x-y)^2}{l^2(x) + l^2(y)}\right] \quad \text{with } l(x) = (x + \epsilon) \times l_0$$



3 hyperparameters controlling different features of the prior:  $\alpha$ ,  $l_0$ ,  $\sigma$



# Example: $u^+ - d^+$ from BCDMS



$$\mathcal{O} = F_2^p - F_2^d = C \otimes f$$

$$f = u^+ - d^+$$

Introduce an interpolation basis for  $f$

$$\mathcal{O}_i = \sum_{\alpha} (FK)_{i\alpha} f(x_{\alpha}) = FK \mathbf{f}$$

# Gaussian inference

<p><b>f</b> Gaussian variable representing PDF on interpolation points <b>x</b></p> <p><math>\mathcal{O} = FK\mathbf{f}</math></p>	<p><b>f*</b> Gaussian variable representing PDF on any set of points <b>x*</b></p>	<p><math>K(x, y; \theta)</math></p> <p>Function modelling correlation</p>	<p><math>y, \epsilon \sim N(0, C_y)</math></p> <p>Data and corresponding experimental error</p>
--	--	---	---

$$\begin{pmatrix} \mathbf{f}^* \\ FK\mathbf{f} \end{pmatrix} \sim \mathcal{N} \left( 0, \begin{pmatrix} K_{\mathbf{x}^*\mathbf{x}^*} & K_{\mathbf{x}^*\mathbf{x}} FK^T \\ FK K_{\mathbf{xx}^*} & FK K_{\mathbf{xx}} FK^T \end{pmatrix} \right)$$

$$p(\mathbf{f}^* | FK\mathbf{f} + \epsilon = y, \theta)$$

This is a gaussian distribution. Its mean and covariance can be computed analytically

# Inference on the hyperparameters

Joint probability distribution  
of  $\mathbf{f}^*$  and  $\theta$

$$p(\mathbf{f}^*, \theta | \text{data}) = p(\mathbf{f}^* | \theta, \text{data}) p(\theta | \text{data})$$

Posterior on the  
hyperparameters given  
the data

$$\propto p(\text{data} | \theta) p_{\theta}(\theta)$$

We can sample from  $p(\theta | \text{data})$  running a MCMC algorithm

# Workflow

Build the prior as a function of hyperparameters:

- Choose kernel
- Encode theory constraints

Collect data and FK tables

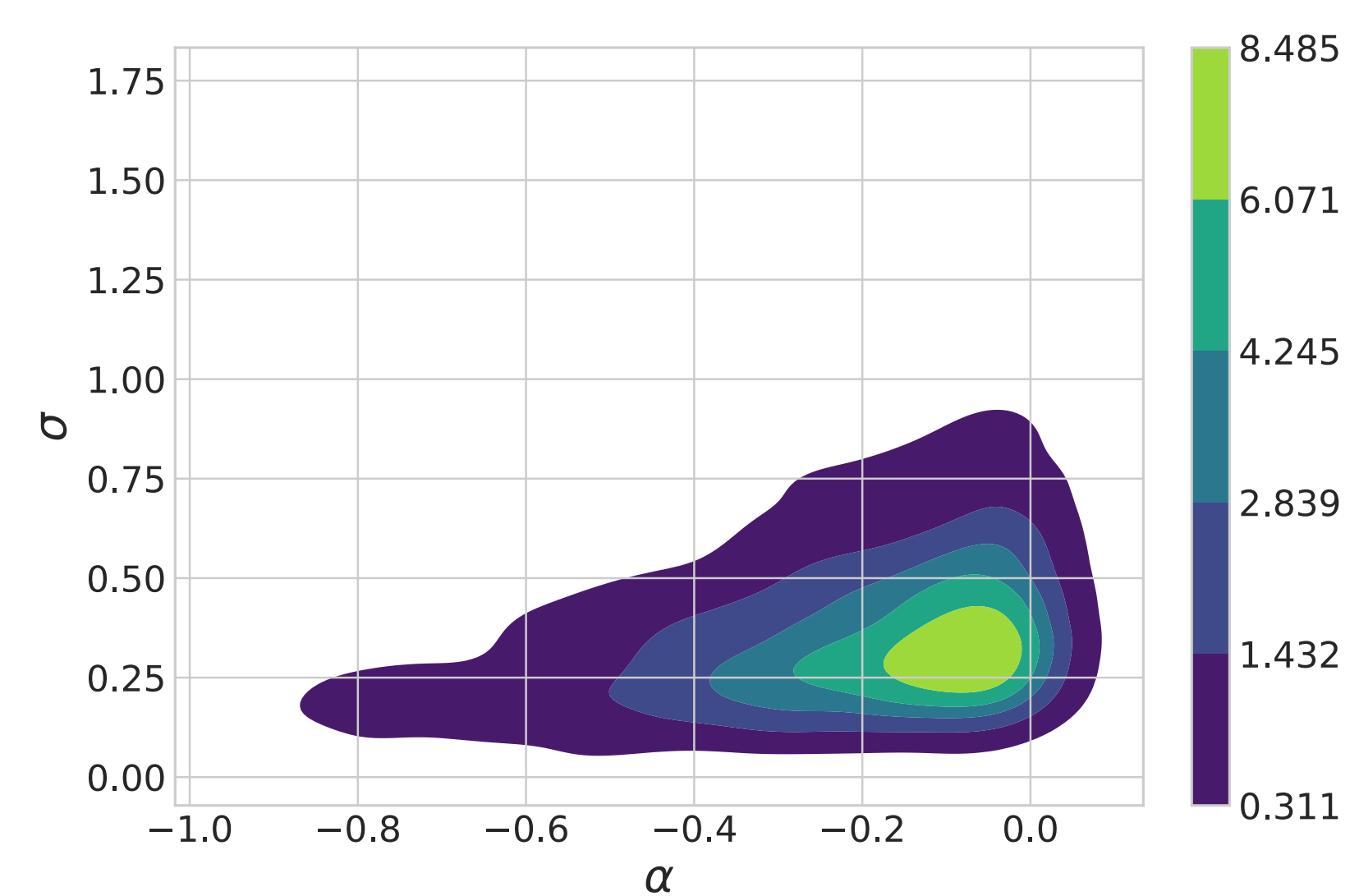
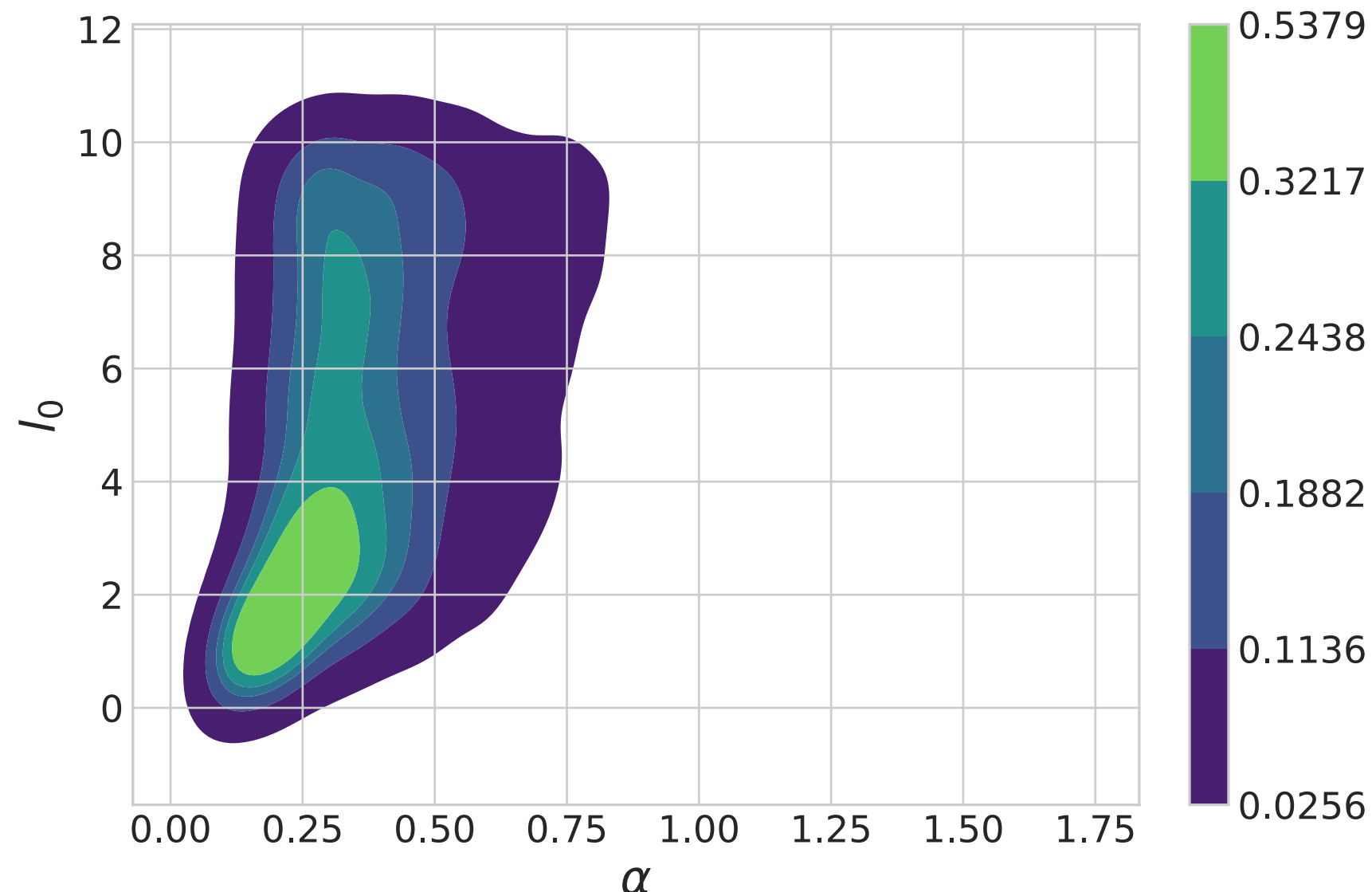
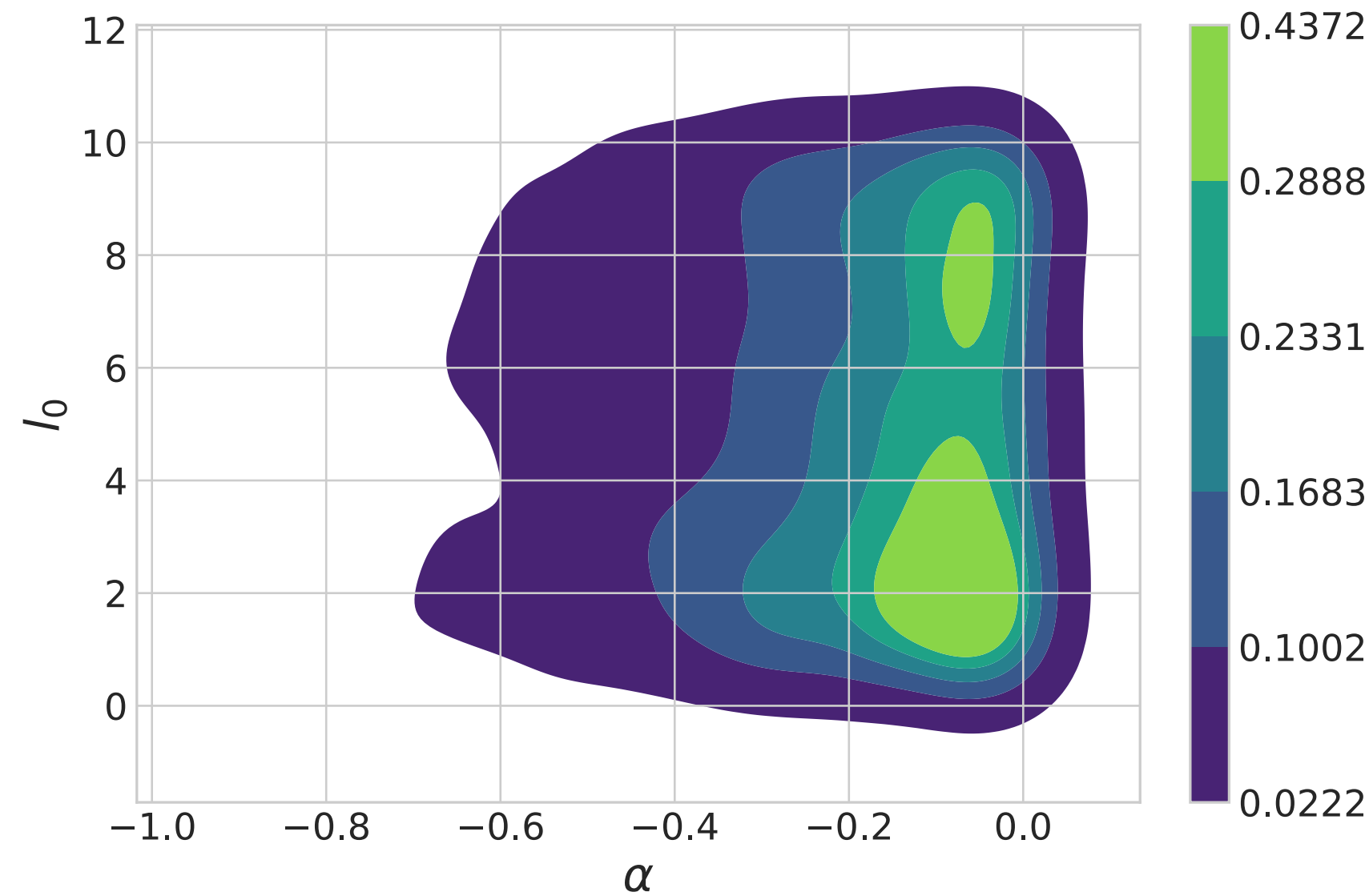
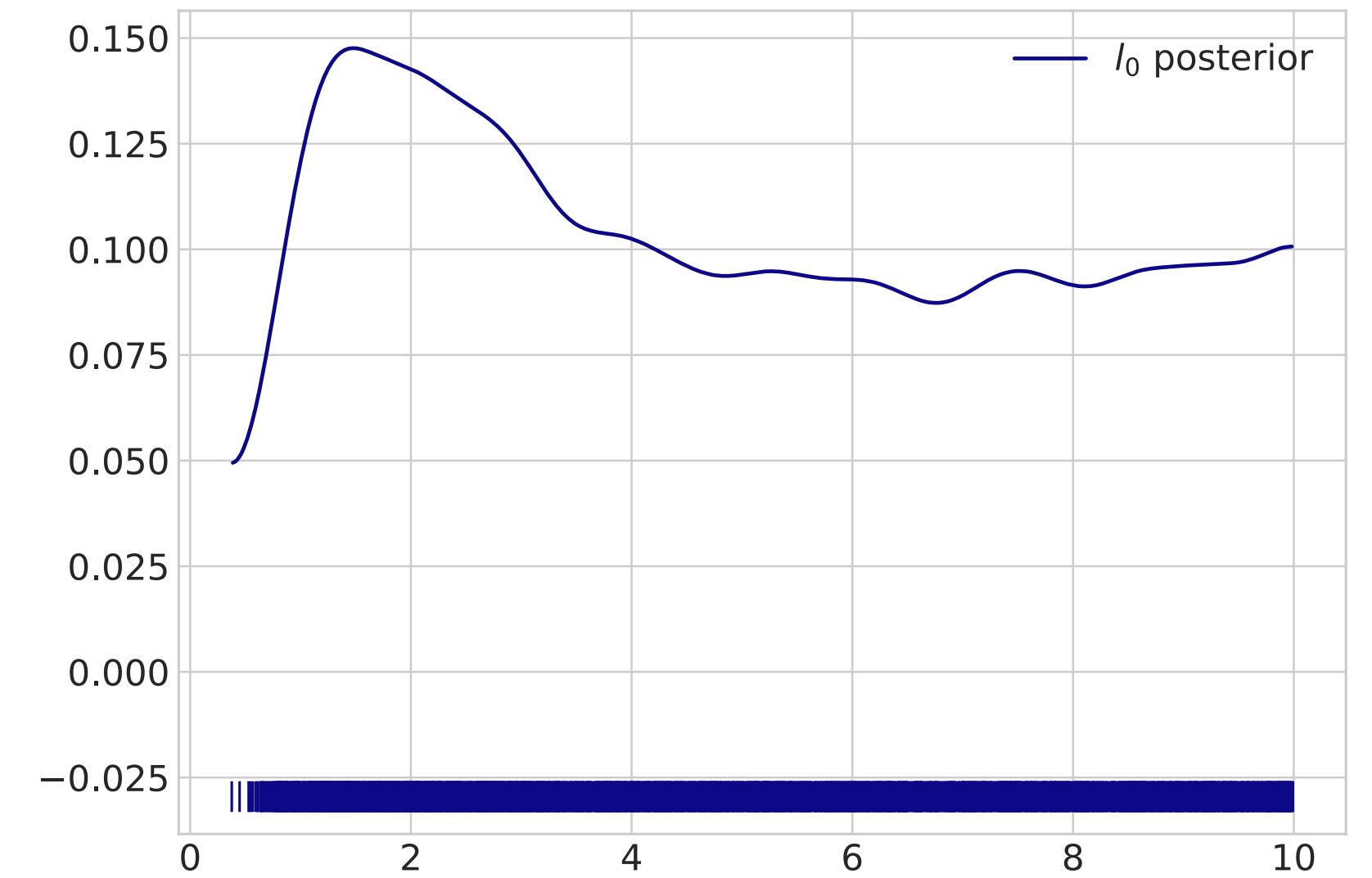
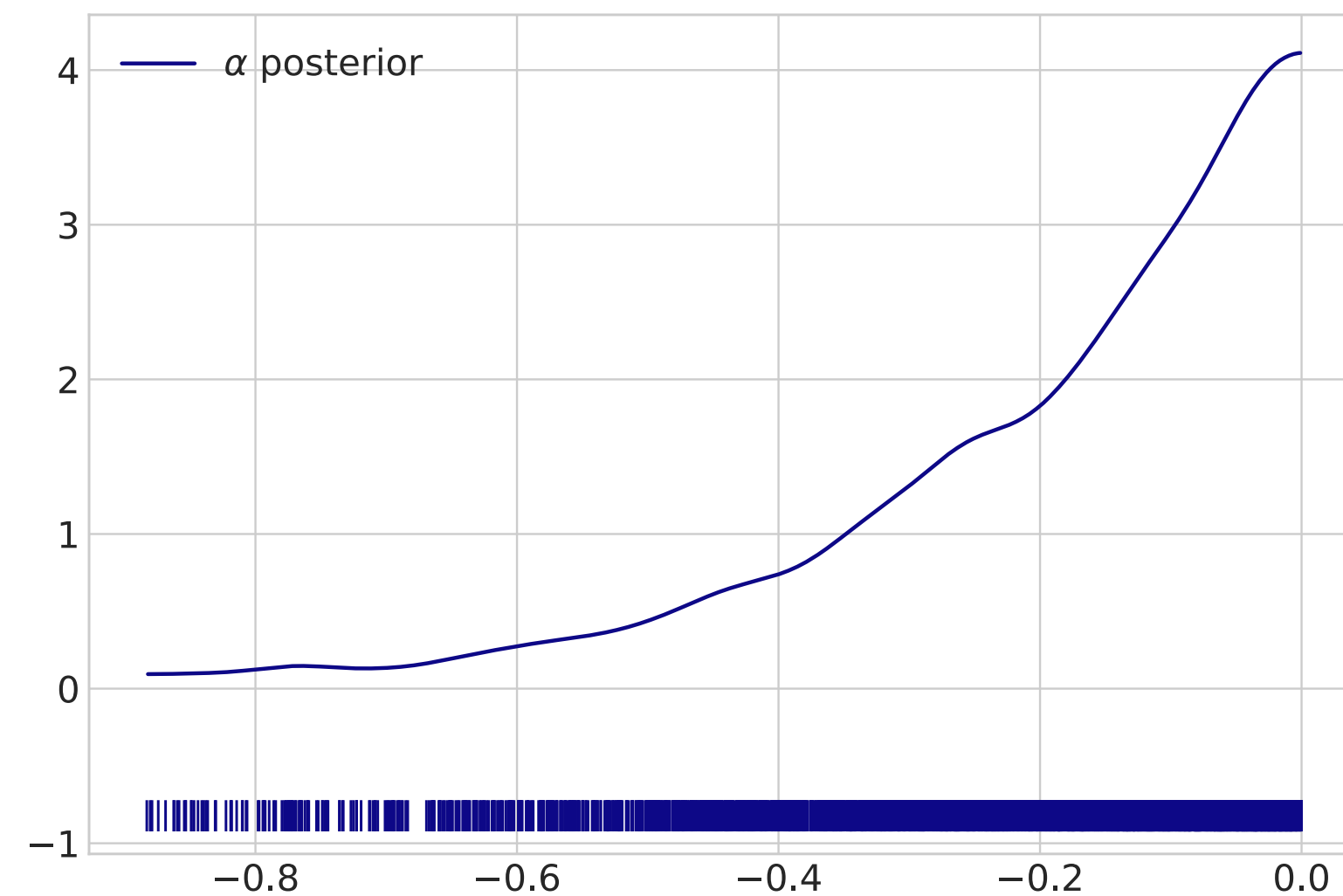
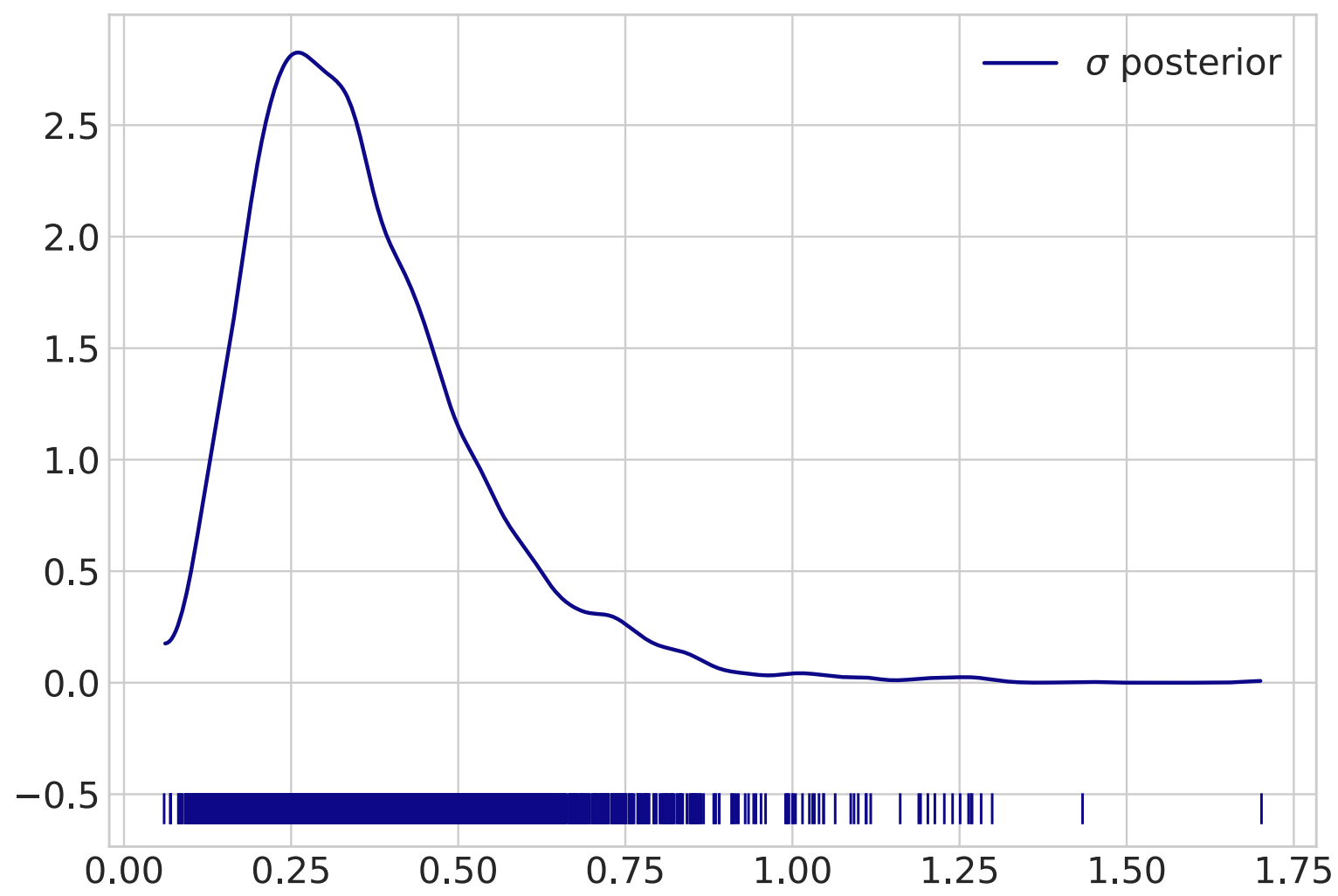
Inference on hyperparameters

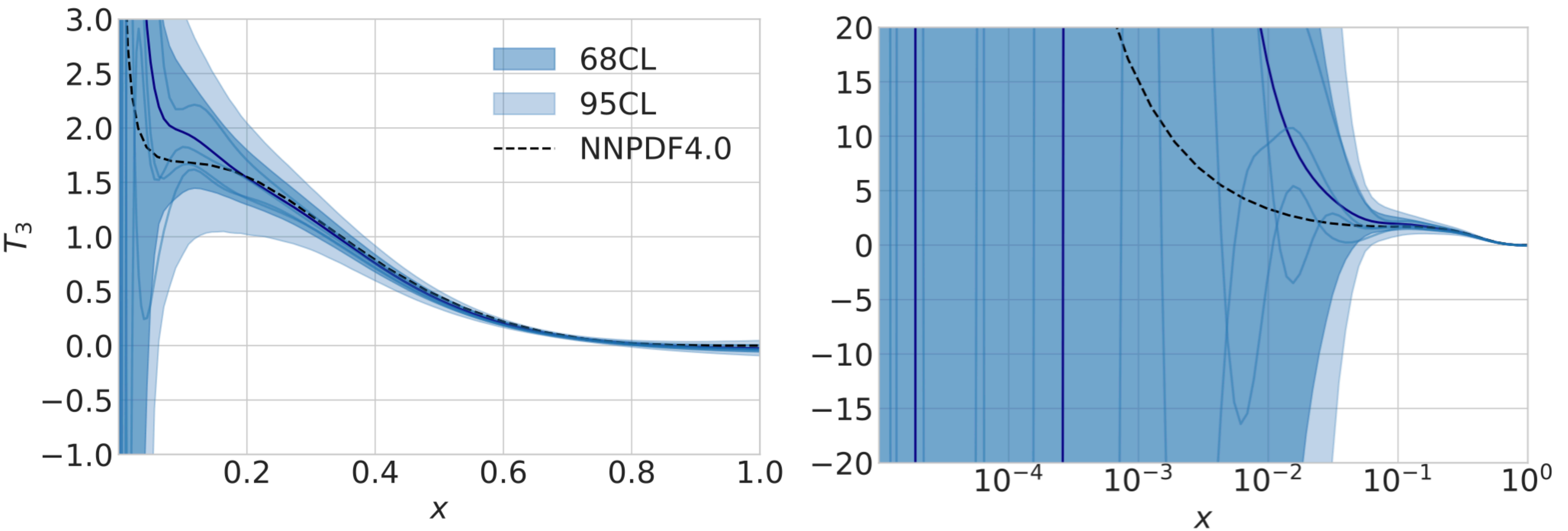


Inference on parameters

For DIS this step is analytical

# Posterior for hyperparameters





Samples from  $p(\mathbf{f}^*, \theta | data)$

uncertainty due to hyperparameter selection is incorporated into the final PDF uncertainty

# Towards a full DIS fit

- 9 flavours
- Extended set of DIS data
- Same formalism and simplifications apply
- Kinetic limit  $f(1) = 0$
- Momentum and valence sumrules

Dataset	References	$N_{\text{dat}}$	$x$	$Q$ [GeV]
NMC $F_2^d/F_2^p$	[33]	260 (121/121)	[0.012, 0.680]	[2.1, 10.]
NMC $\sigma^{\text{NC},p}$	[34]	292 (204/204)	[0.012, 0.500]	[1.8, 7.9]
SLAC $F_2^p$	[35]	211 (33/33)	[0.140, 0.550]	[1.9, 4.4]
SLAC $F_2^d$	[35]	211 (34/34)	[0.140, 0.550]	[1.9, 4.4]
BCDMS $F_2^p$	[36]	351 (333/333)	[0.070, 0.750]	[2.7, 15.]
BCDMS $F_2^d$	[36]	254 (248/248)	[0.070, 0.750]	[2.7, 15.]
CHORUS $\sigma_{CC}^v$	[37]	607 (416/416)	[0.045, 0.650]	[1.9, 9.8]
CHORUS $\sigma_{CC}^{\bar{v}}$	[37]	607 (416/416)	[0.045, 0.650]	[1.9, 9.8]
NuTeV $\sigma_{CC}^v$ (dimuon)	[38,39]	45 (39/39)	[0.020, 0.330]	[2.0, 11.]
NuTeV $\sigma_{CC}^{\bar{v}}$ (dimuon)	[38,39]	45 (36/37)	[0.020, 0.210]	[1.9, 8.3]
[NOMAD $\mathcal{R}_{\mu\mu}(E_\nu)$ ] (*)	[111]	15 (-/15)	[0.030, 0.640]	[1.0, 28.]
[EMC $F_2^c$ ]	[44]	21 (-/16)	[0.014, 0.440]	[2.1, 8.8]
HERA I+II $\sigma_{\text{NC},CC}^p$	[40]	1306 (1011/1145)	$[4 \cdot 10^{-5}, 0.65]$	[1.87, 223]
HERA I+II $\sigma_{\text{NC}}^c$ (*)	[145]	52 (-/37)	$[7 \cdot 10^{-5}, 0.05]$	[2.2, 45]
HERA I+II $\sigma_{\text{NC}}^b$ (*)	[145]	27 (26/26)	$[2 \cdot 10^{-4}, 0.50]$	[2.2, 45]

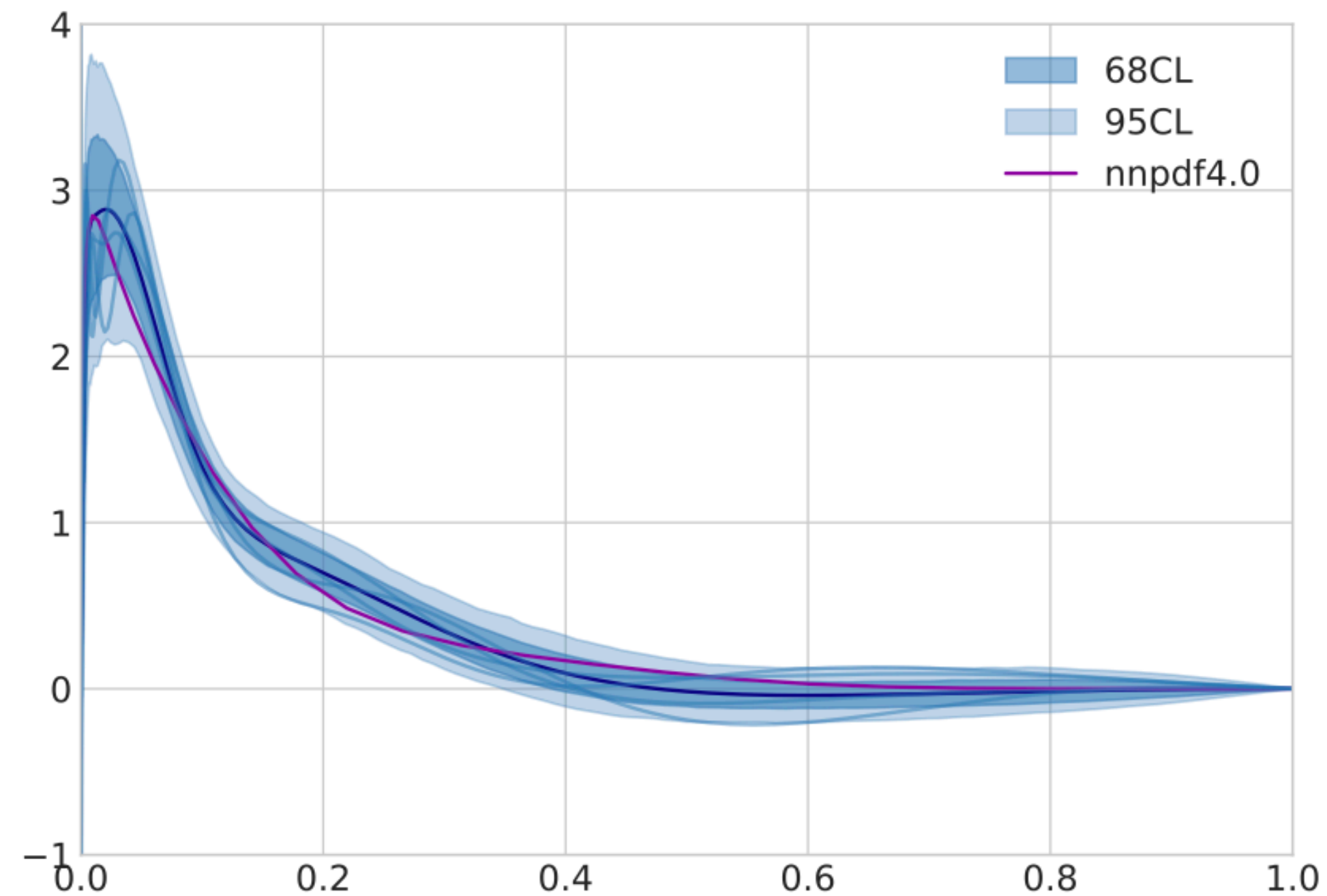
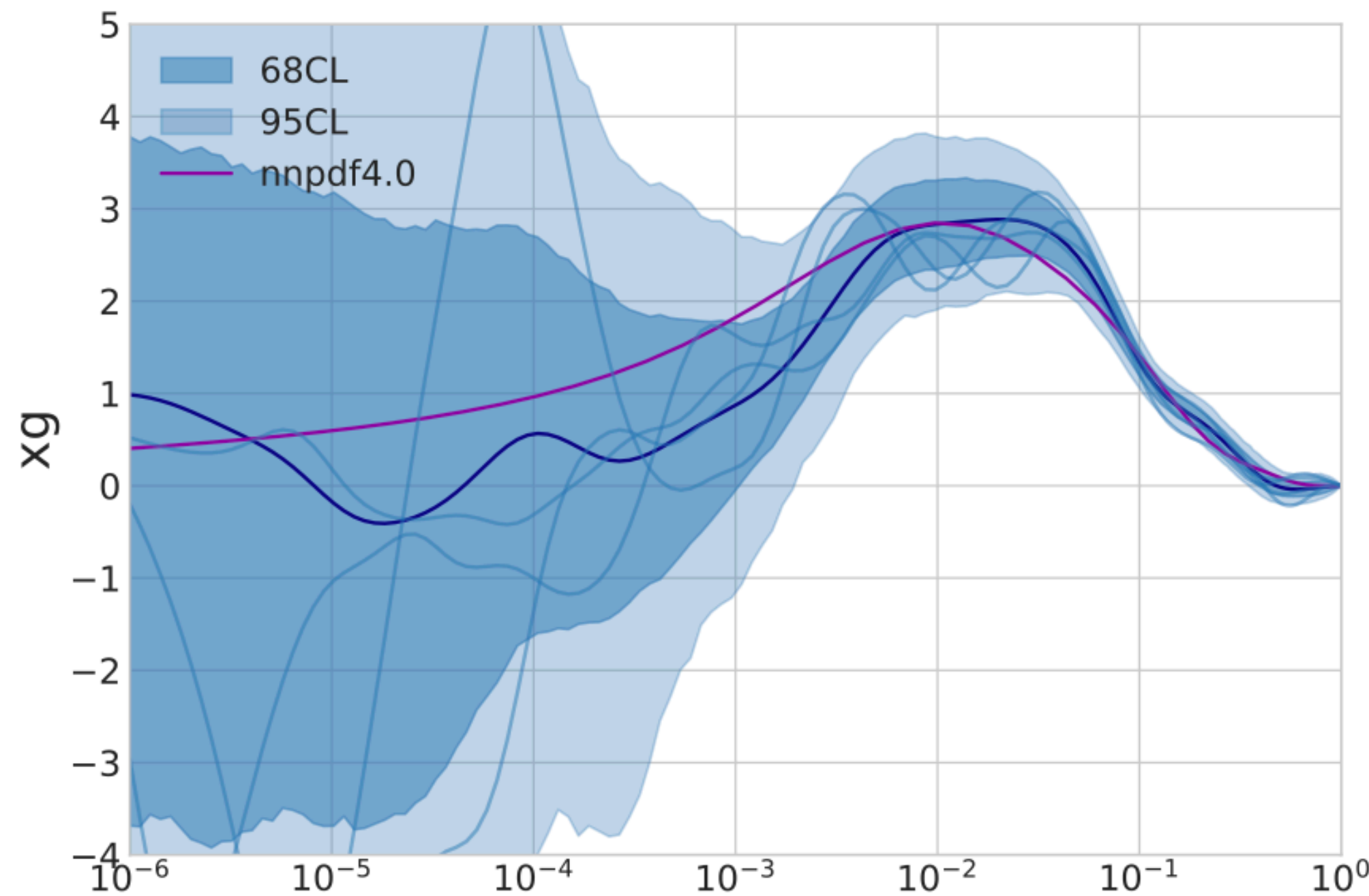


Table from  
Eur.Phys.J.C 82 (2022)  
5

# Towards a full DIS fit

- 9 flavours
- Extended set of DIS data
- Same formalism and simplifications apply
- Kinetic limit  $f(1) = 0$
- Momentum and valence sumrules

Dataset	References	$N_{\text{dat}}$	$x$	$Q$ [GeV]
NMC $F_2^d/F_2^p$	[33]	260 (121/121)	[0.012, 0.680]	[2.1, 10.]
NMC $\sigma^{\text{NC},p}$	[34]	292 (204/204)	[0.012, 0.500]	[1.8, 7.9]
SLAC $F_2^p$	[35]	211 (33/33)	[0.140, 0.550]	[1.9, 4.4]
SLAC $F_2^d$	[35]	211 (34/34)	[0.140, 0.550]	[1.9, 4.4]
BCDMS $F_2^p$	[36]	351 (333/333)	[0.070, 0.750]	[2.7, 15.]
BCDMS $F_2^d$	[36]	254 (248/248)	[0.070, 0.750]	[2.7, 15.]
CHORUS $\sigma_{CC}^v$	[37]	607 (416/416)	[0.045, 0.650]	[1.9, 9.8]
CHORUS $\sigma_{CC}^{\bar{v}}$	[37]	607 (416/416)	[0.045, 0.650]	[1.9, 9.8]
NuTeV $\sigma_{CC}^v$ (dimuon)	[38,39]	45 (39/39)	[0.020, 0.330]	[2.0, 11.]
NuTeV $\sigma_{CC}^{\bar{v}}$ (dimuon)	[38,39]	45 (36/37)	[0.020, 0.210]	[1.9, 8.3]
[NOMAD $\mathcal{R}_{\mu\mu}(E_\nu)$ ] (*)	[111]	15 (-/15)	[0.030, 0.640]	[1.0, 28.]
[EMC $F_2^c$ ]	[44]	21 (-/16)	[0.014, 0.440]	[2.1, 8.8]
HERA I+II $\sigma_{\text{NC},\text{CC}}^p$	[40]	1306 (1011/1145)	$[4 \cdot 10^{-5}, 0.65]$	[1.87, 223]
HERA I+II $\sigma_{\text{NC}}^c$ (*)	[145]	52 (-/37)	$[7 \cdot 10^{-5}, 0.05]$	[2.2, 45]
HERA I+II $\sigma_{\text{NC}}^b$ (*)	[145]	27 (26/26)	$[2 \cdot 10^{-4}, 0.50]$	[2.2, 45]

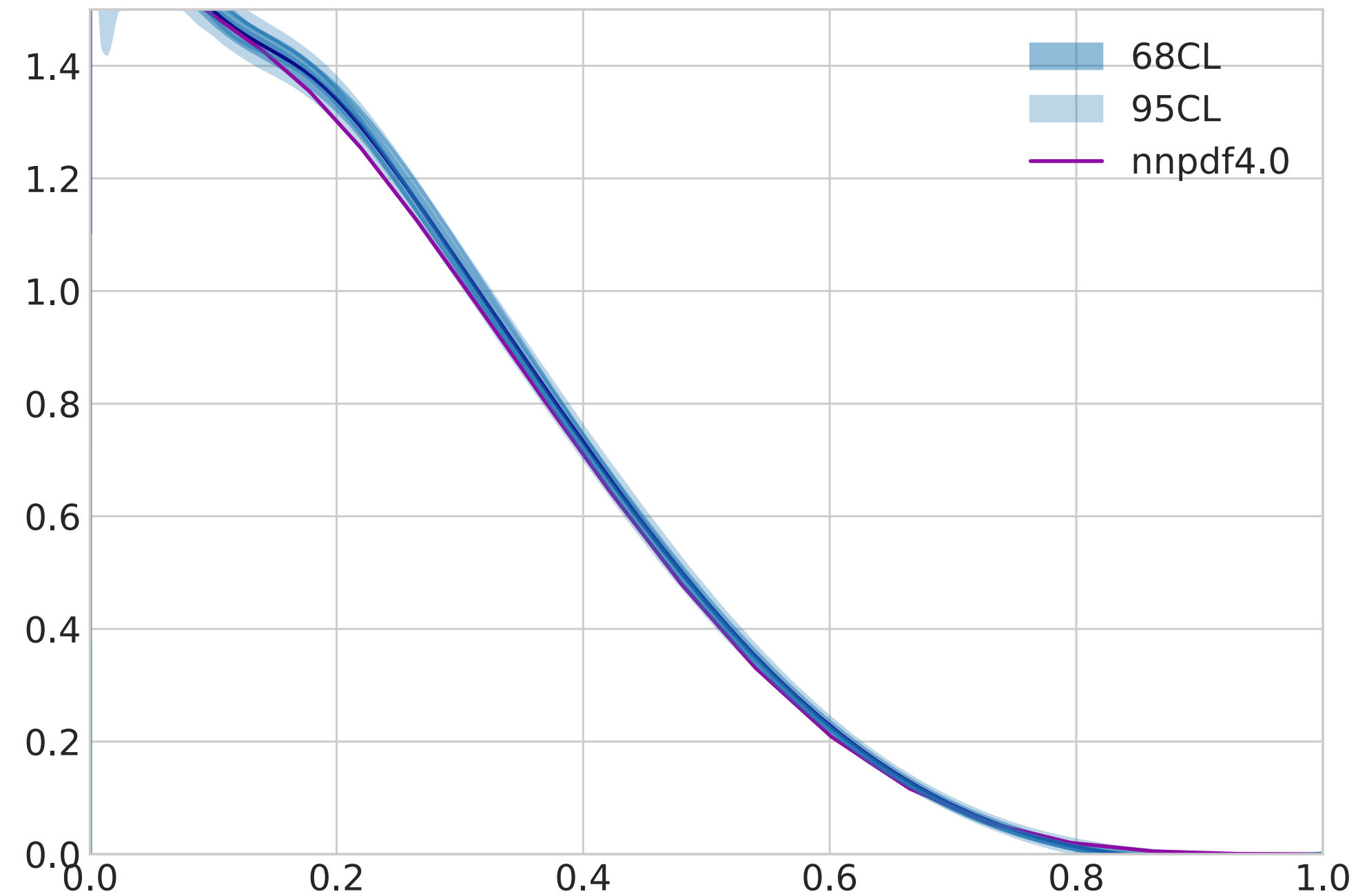
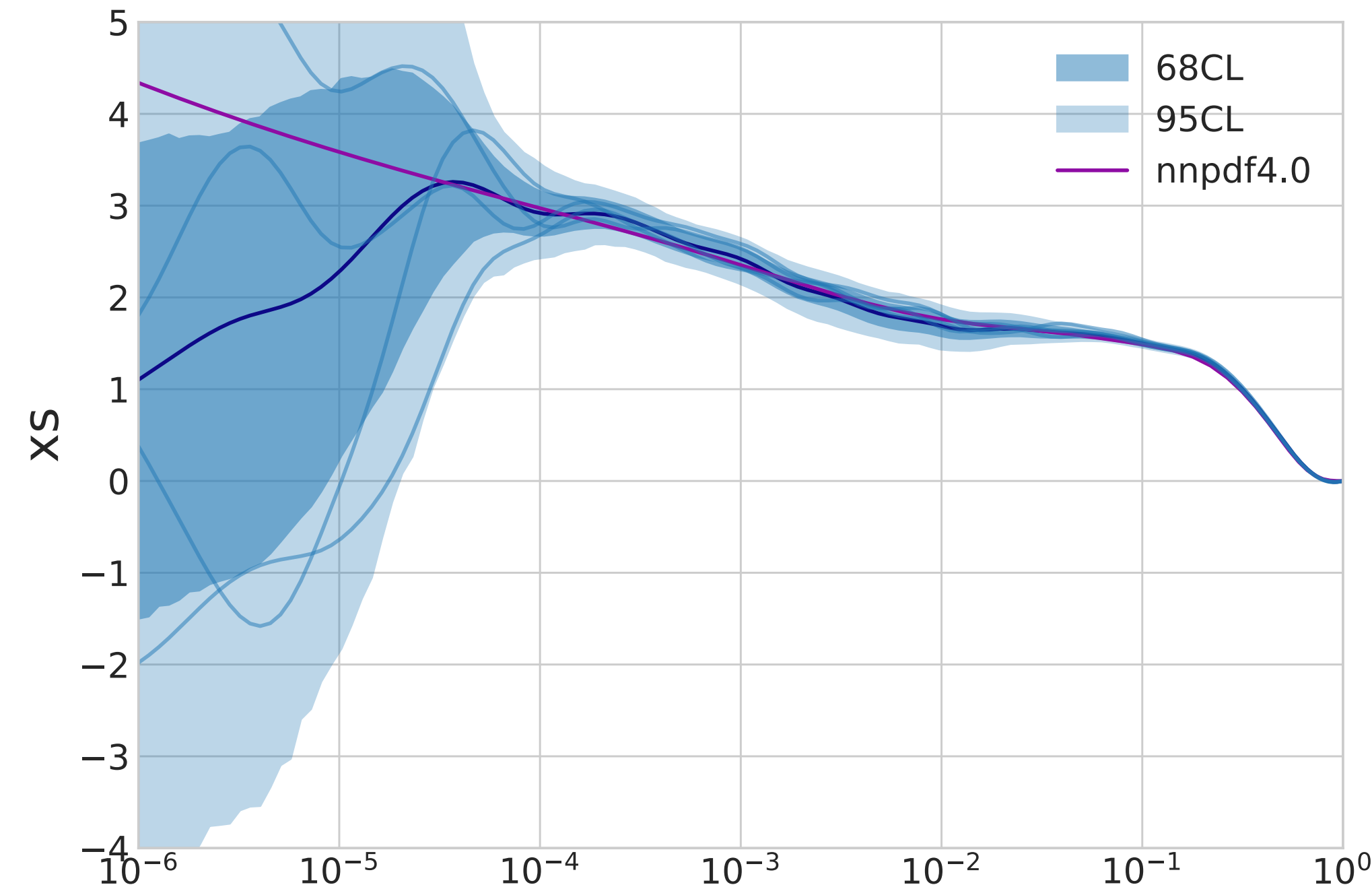


Table from  
Eur.Phys.J.C 82 (2022)  
5



# Towards a full DIS fit

- 9 flavours
- Extended set of DIS data
- Same formalism and simplifications apply
- Kinetic limit  $f(1) = 0$
- Momentum and valence sumrules

Dataset	References	$N_{\text{dat}}$	$x$	$Q$ [GeV]
NMC $F_2^d/F_2^p$	[33]	260 (121/121)	[0.012, 0.680]	[2.1, 10.]
NMC $\sigma^{\text{NC},p}$	[34]	292 (204/204)	[0.012, 0.500]	[1.8, 7.9]
SLAC $F_2^p$	[35]	211 (33/33)	[0.140, 0.550]	[1.9, 4.4]
SLAC $F_2^d$	[35]	211 (34/34)	[0.140, 0.550]	[1.9, 4.4]
BCDMS $F_2^p$	[36]	351 (333/333)	[0.070, 0.750]	[2.7, 15.]
BCDMS $F_2^d$	[36]	254 (248/248)	[0.070, 0.750]	[2.7, 15.]
CHORUS $\sigma_{CC}^{\nu}$	[37]	607 (416/416)	[0.045, 0.650]	[1.9, 9.8]
CHORUS $\sigma_{CC}^{\bar{\nu}}$	[37]	607 (416/416)	[0.045, 0.650]	[1.9, 9.8]
NuTeV $\sigma_{CC}^{\nu}$ (dimuon)	[38,39]	45 (39/39)	[0.020, 0.330]	[2.0, 11.]
NuTeV $\sigma_{CC}^{\bar{\nu}}$ (dimuon)	[38,39]	45 (36/37)	[0.020, 0.210]	[1.9, 8.3]
[NOMAD $\mathcal{R}_{\mu\mu}(E_{\nu})$ ] (*)	[111]	15 (-/15)	[0.030, 0.640]	[1.0, 28.]
[EMC $F_2^c$ ]	[44]	21 (-/16)	[0.014, 0.440]	[2.1, 8.8]
HERA I+II $\sigma_{\text{NC,CC}}^p$	[40]	1306 (1011/1145)	$[4 \cdot 10^{-5}, 0.65]$	[1.87, 223]
HERA I+II $\sigma_{\text{NC}}^c$ (*)	[145]	52 (-/37)	$[7 \cdot 10^{-5}, 0.05]$	[2.2, 45]
HERA I+II $\sigma_{\text{NC}}^b$ (*)	[145]	27 (26/26)	$[2 \cdot 10^{-4}, 0.50]$	[2.2, 45]

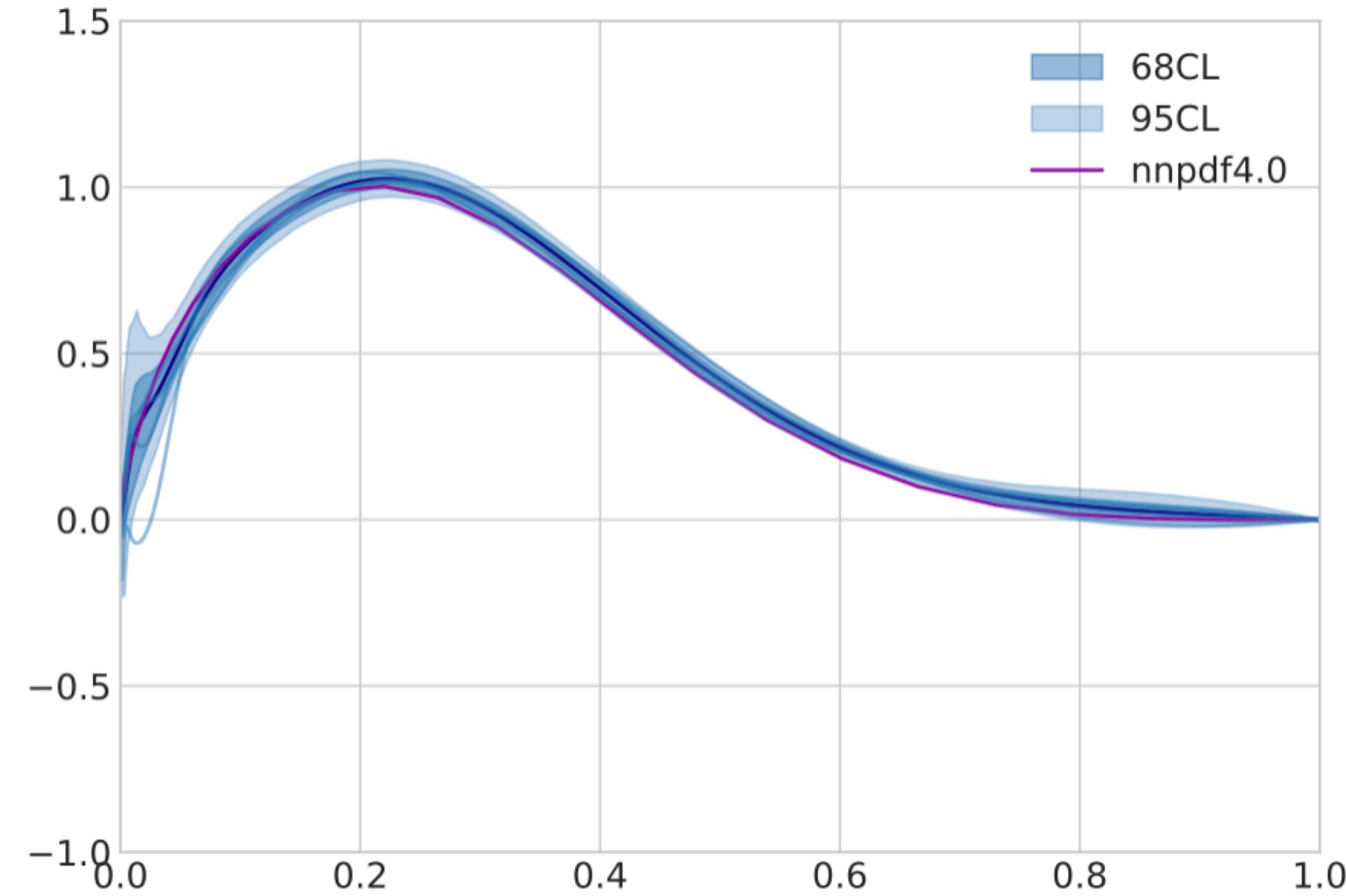
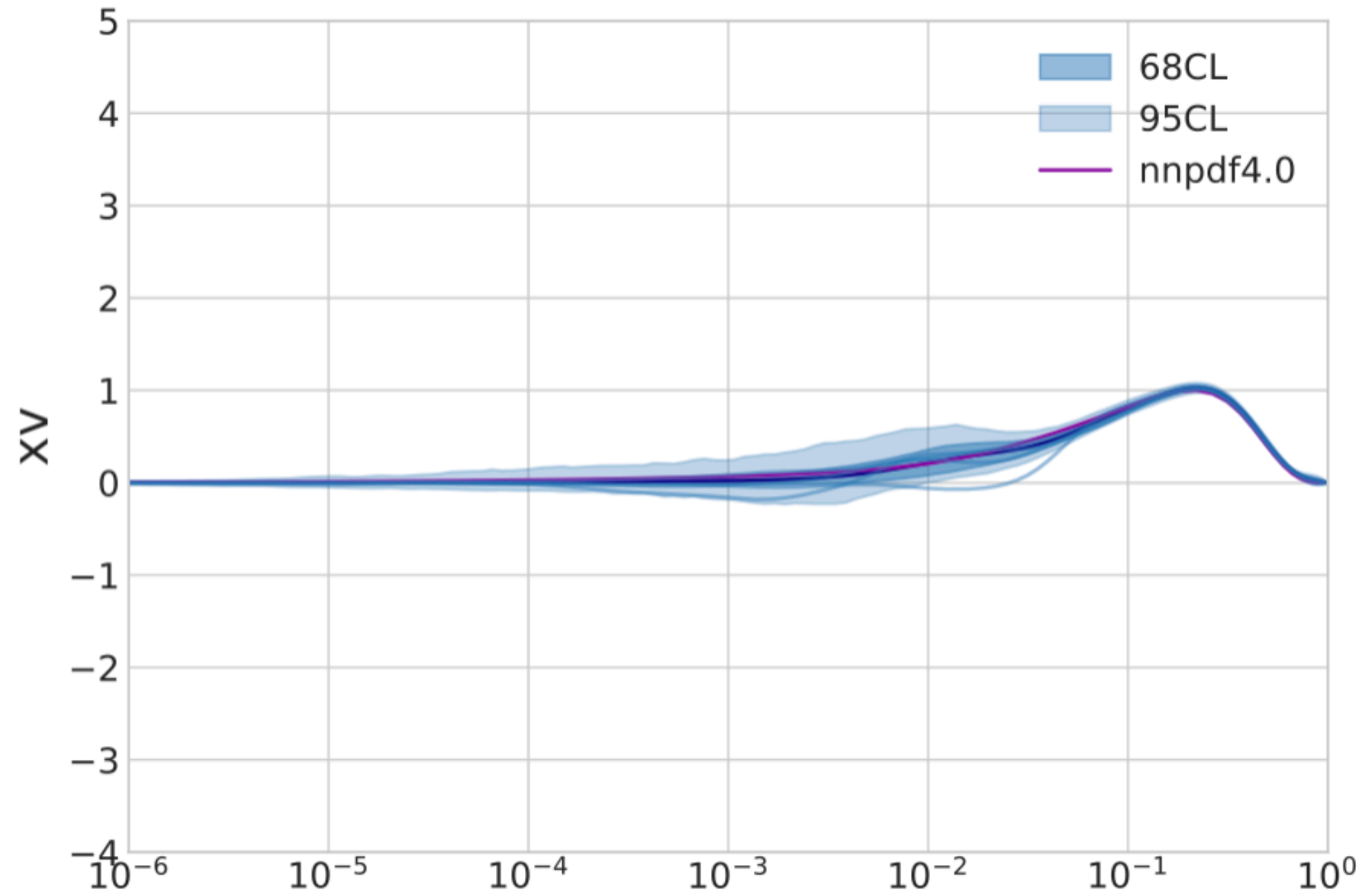


Table from  
Eur.Phys.J.C 82 (2022)  
5

## What changes in a global fit?

$$p(\mathbf{f}, \theta | \text{data}) = p(\mathbf{f} | \text{data}, \theta) p(\theta | \text{data})$$

This bit is not a gaussian distribution any longer

To access the posterior we have to run a MCMC having dimension  $\dim \mathbf{f} + \dim \theta$

---

### TO DO :

- Provide a DIS only PDF fit
- Study dependence on the kernel
- Compare with existing methodology. Are there differences?

## What changes in a global fit?

$$p(\mathbf{f}, \theta | \text{data}) = p(\mathbf{f} | \text{data}, \theta) p(\theta | \text{data})$$

This bit is not a gaussian distribution any longer

To access the posterior we have to run a MCMC having dimension  $\dim \mathbf{f} + \dim \theta$

---

### TO DO :

- Provide a DIS only PDF fit
- Study dependence on the kernel
- Compare with existing methodology. Are there differences?

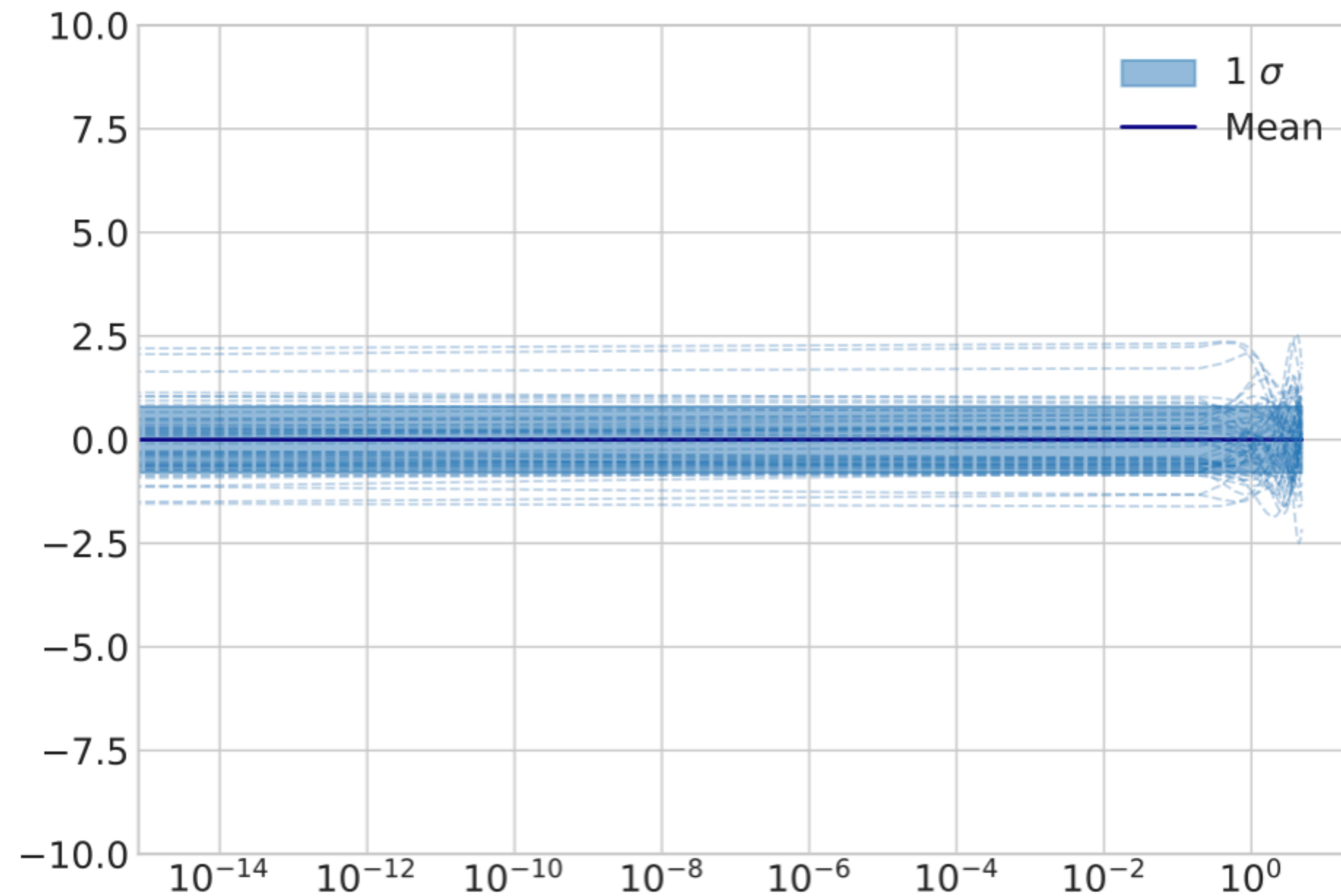
Thanks!

**Backup slides**

# An example of a bad prior for PDFs

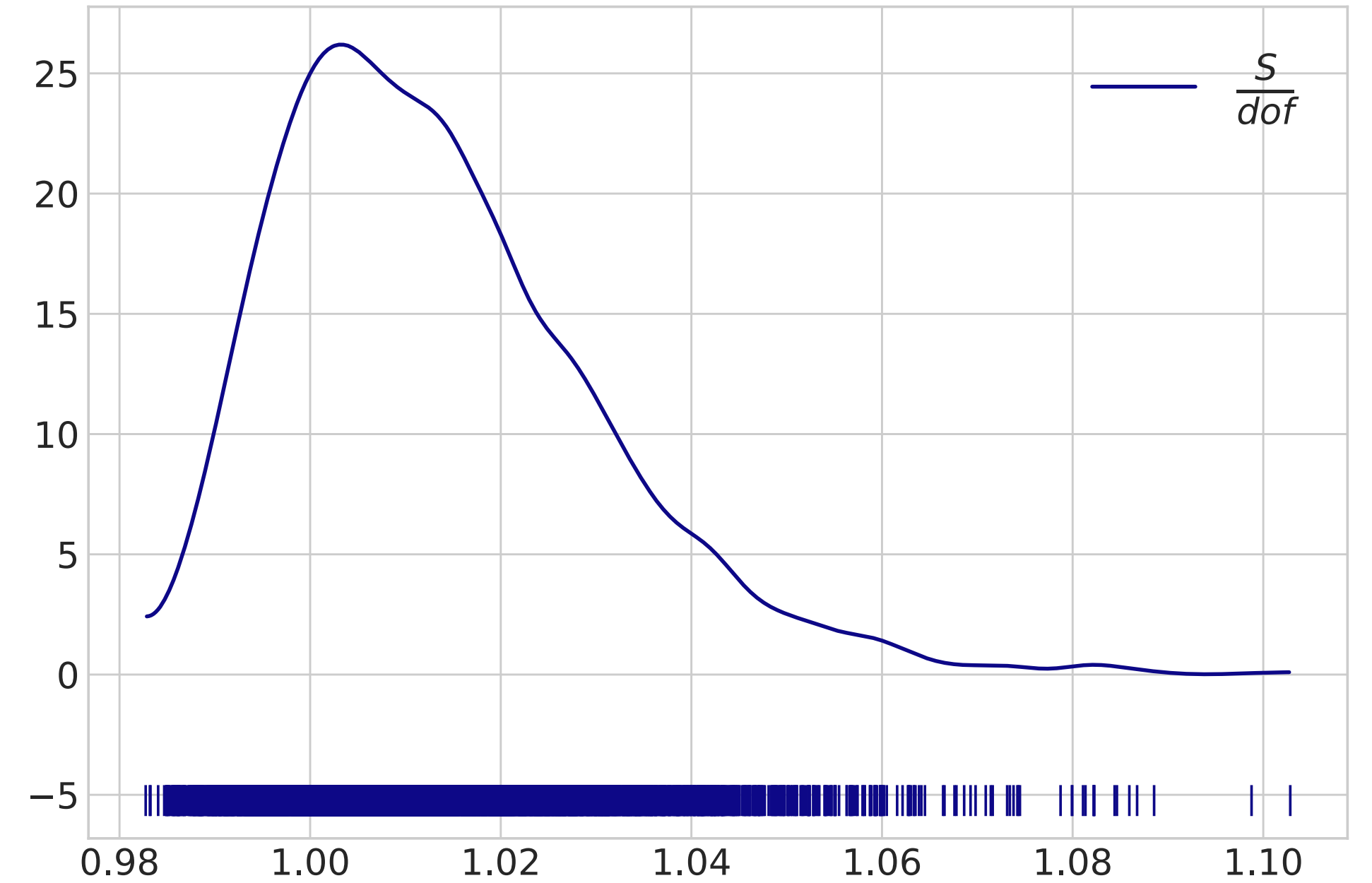
$$k(x, y) = \sigma^2 \exp \left[ -\frac{(x - y)^2}{l^2} \right]$$

Exponential quadratic



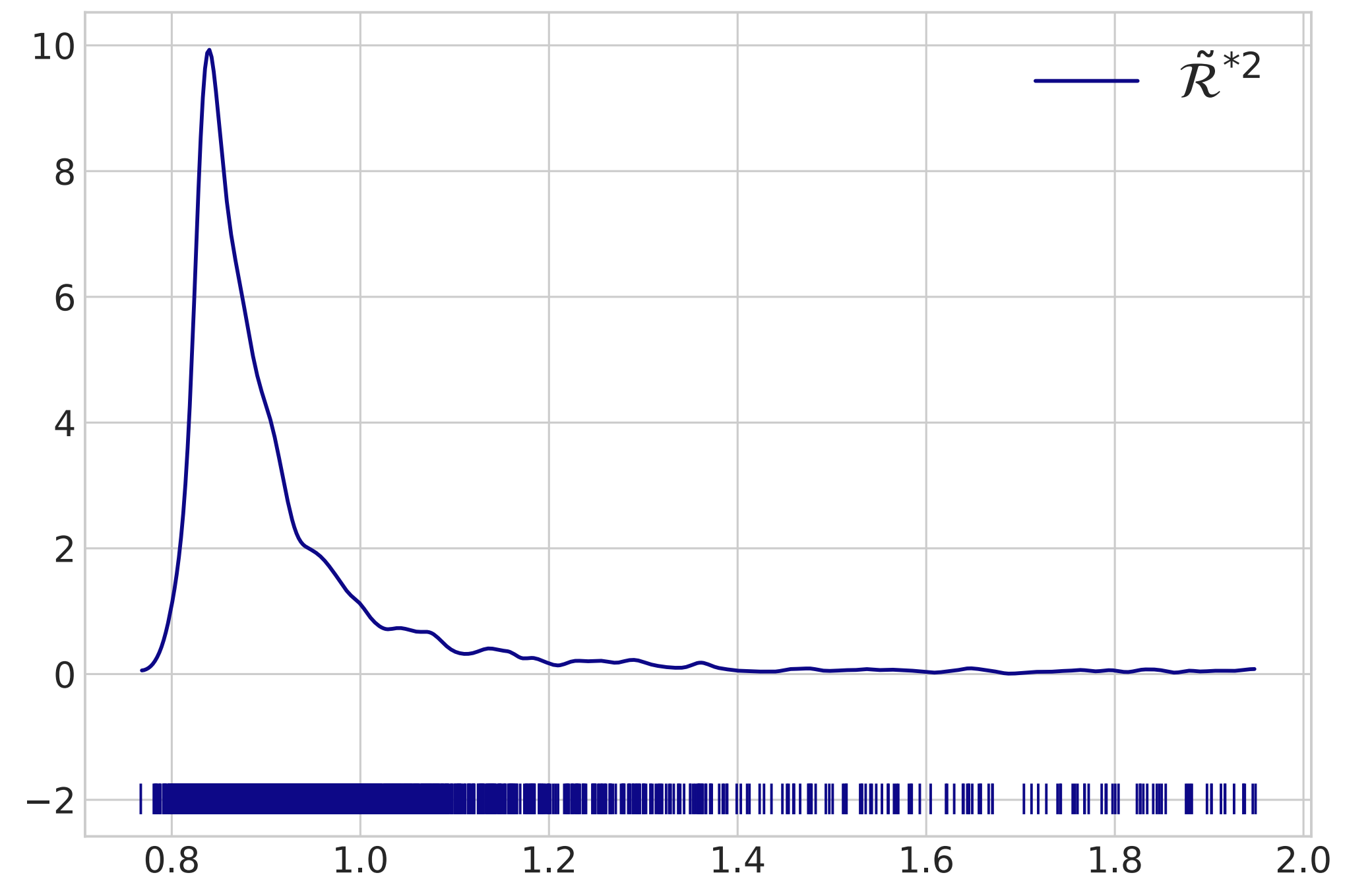
# Fit quality

$$\frac{S}{dof} = \frac{1}{N_{\text{data}}} \left( (\mathbf{m} - \tilde{\mathbf{m}})^T K_{xx}^{-1} (\mathbf{m} - \tilde{\mathbf{m}}) + (y - FK \tilde{\mathbf{m}})^T C_Y^{-1} (y - FK \tilde{\mathbf{m}}) \right)$$



# Generalisation on unseen data

$$\tilde{\mathcal{R}}^{*2} = \frac{1}{\dim(y^* | y)} (FK^* \tilde{\mathbf{m}} - y^*)^T \left( FK^* \tilde{K}_{xx} FK^{*T} + C_Y^* \right)^+ (FK^* \tilde{\mathbf{m}} - y^*)$$



# Gaussian inference

**f**

Gaussian variable  
representing PDF on  
interpolation points **x**

$$\mathcal{O} = FK\mathbf{f}$$

**f\***

Gaussian variable  
representing PDF on any  
set of points **x\***

$K(x, y; \theta)$

Function modelling  
correlation

$y, \epsilon \sim N(0, C_y)$

Data and  
corresponding  
experimental error

$$\tilde{\mathbf{m}}^* = \mathbf{m} + K_{\mathbf{x}^*\mathbf{x}}FK^T \left( FK K_{\mathbf{xx}}FK^T + C_y \right)^+ (\mathbf{y} - \mathbf{m})$$

$$\tilde{K}^* = K_{\mathbf{x}^*\mathbf{x}^*} - K_{\mathbf{x}^*\mathbf{x}}FK^T \left( FK K_{\mathbf{xx}}FK^T + C_y \right)^+ FK K_{\mathbf{xx}}^*$$

inference of the gaussian parameters **f\*** can be done  
analytically

# Further possible applications

- simultaneous fits of PDFs and Wilson coefficients

$$\sigma_{\text{eft}}(\mathbf{c}/\Lambda^2) = \sigma_{\text{SM}} + \sum_i \tilde{\sigma}_i^{\text{LO/NLO}} \frac{c_i}{\Lambda^2} + \sum_{i,j} \tilde{\sigma}_{ij}^{\text{LO/NLO}} \frac{c_i c_j}{\Lambda^4}$$

## The top quark legacy of the LHC Run II for PDF and SMEFT analyses

---

Zahari Kassabov,<sup>a</sup> Maeve Madigan,<sup>a</sup> Luca Mantani,<sup>a</sup> James Moore,<sup>a</sup>  
Manuel Morales Alvarado,<sup>a</sup> Juan Rojo<sup>b,c</sup> and Maria Ubiali<sup>a</sup>

*JHEP* 05 (2023) 205

- Inverse problems relevant for the lattice community

## Reconstructing QCD Spectral Functions with Gaussian Processes

Jan Horak,<sup>1</sup> Jan M. Pawłowski,<sup>1,2</sup> José Rodríguez-Quintero,<sup>3</sup> Jonas Turnwald,<sup>1</sup>  
Julian M. Urban,<sup>1,\*</sup> Nicolas Wink,<sup>1</sup> and Savvas Zafeiropoulos<sup>4</sup>

*Phys.Rev.D* 105 (2022) 3



## Decomposition of PDF uncertainty

$$\tilde{K} = \underbrace{(I - R_{xx}) K_{xx} (I - R_{xx})^T}_{\text{Methodology}} + \underbrace{a_{xx}^T C_y a_{xx}}_{\text{Experimental error}}$$

$$a_{xx}^T = K_{xx} F K^T \left( F K K_{xx} F K^T + C_y \right)^+$$

$$R_{xx} = a_{xx}^T F K$$

