# Bayesian inference and gaussian processes for PDF determination
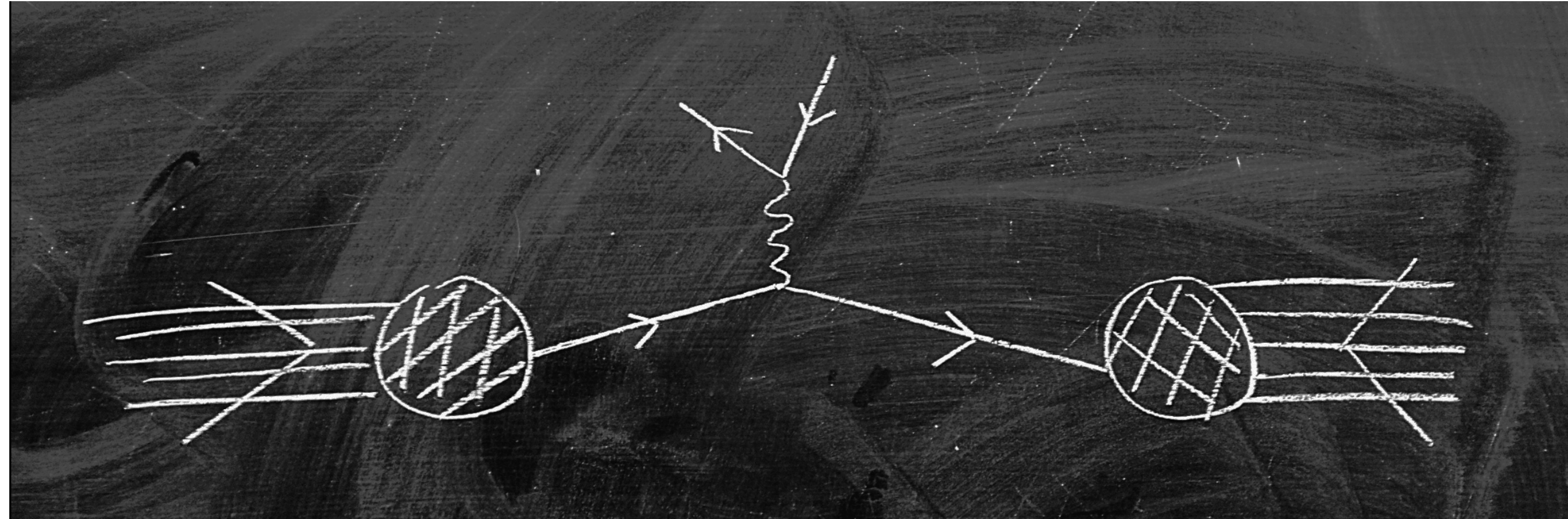
## Tommaso Giani

**Inverse Problems and Uncertainty Quantification in Nuclear Physics**

**8/07/2024**

# Parton Distribution Functions (PDFs)



Hard matrix element: accessible in perturbation theory

$$\sigma = \sum_{i,j} \int dx_1 dx_2\, f_i\left(x_1, \mu\right) f_j\left(x_2, \mu\right) \hat{\sigma}\left(x_1, x_2, \frac{Q}{\mu}\right) \times \left(1 + \mathcal{O}\left(\Lambda/M\right)^p\right)$$

- PDFs: non perturbative objects, extracted from experimental data

$$f_i\left(x, \mu\right)$$

- $x$ : momentum fraction

  $p_{\text{parton}}/p_{\text{proton}}$

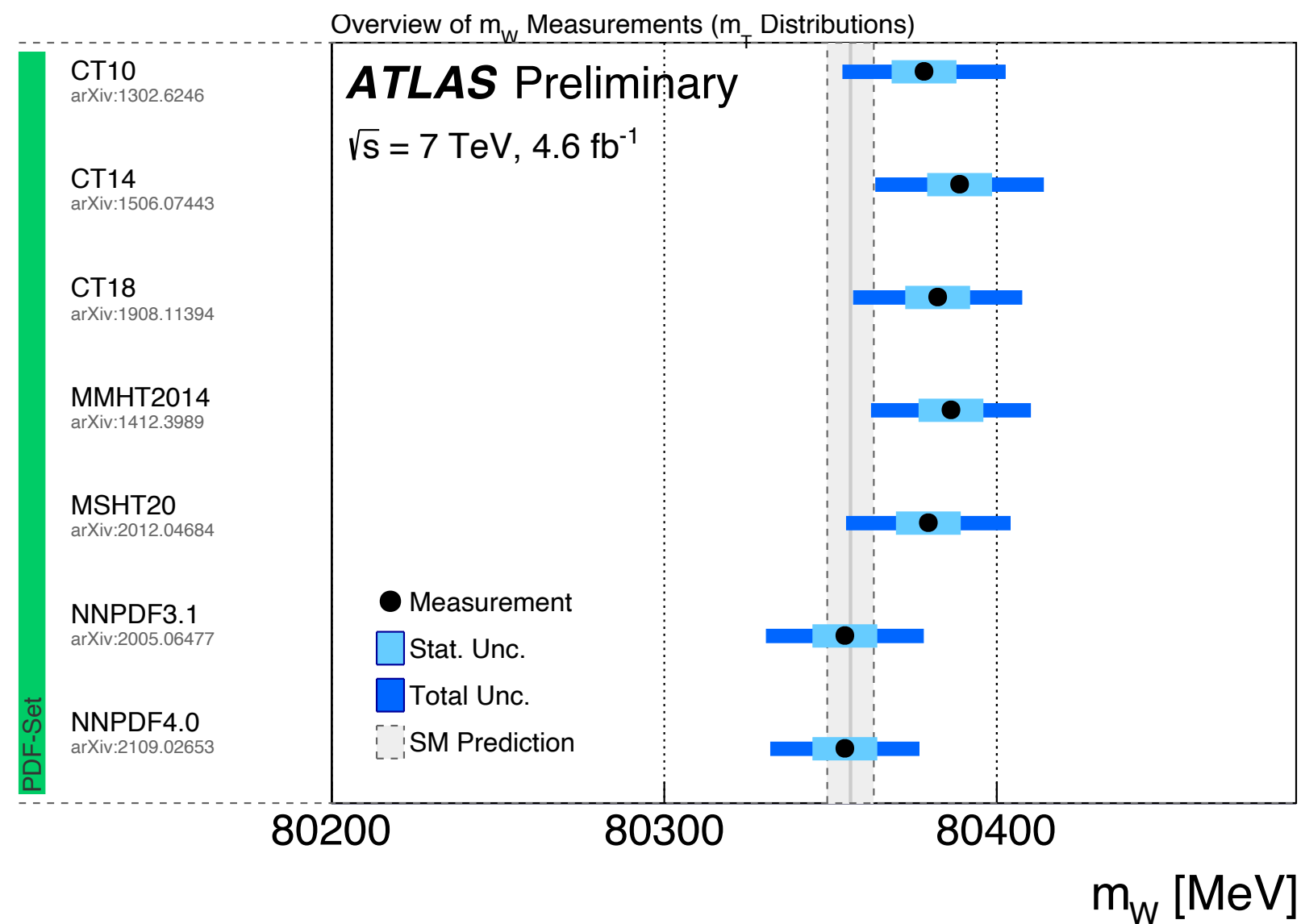- $\mu$ : energy scale, computable perturbation theory

# PDFs and precision studies

## $\alpha_s$ from Z pT
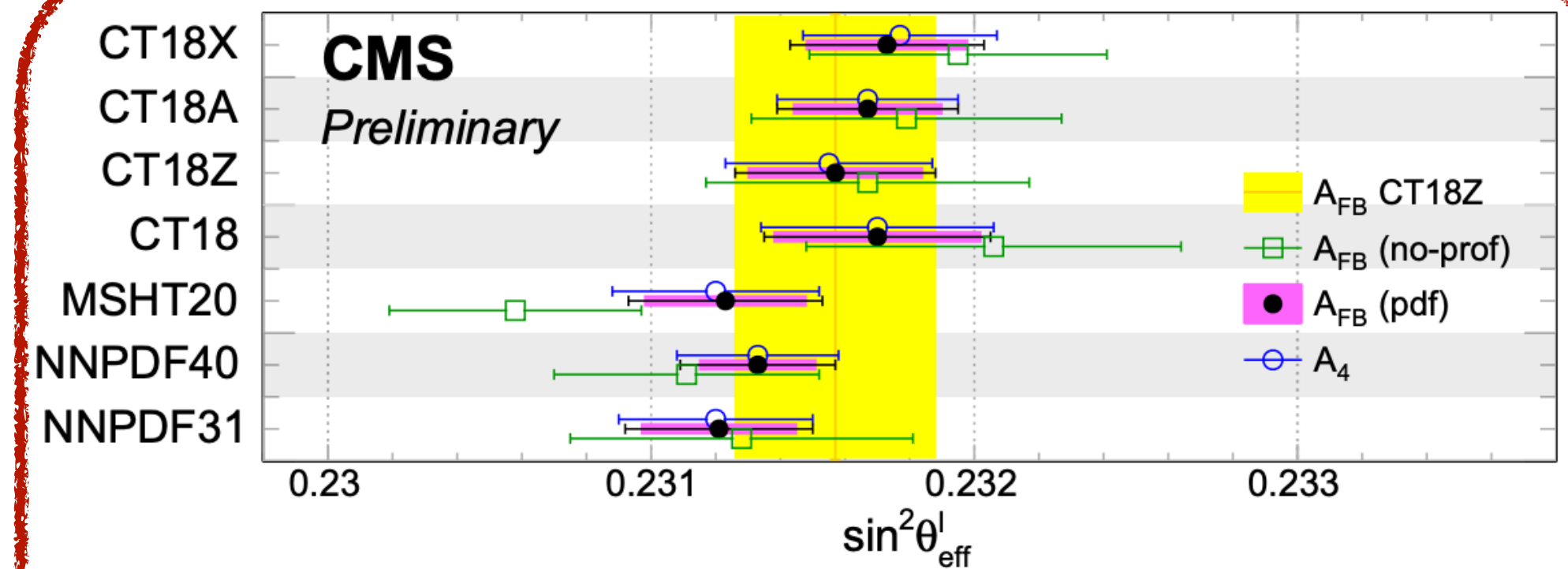
$$\alpha_s(m_Z) = 0.11847 + 0.00091 - 0.00088$$

~ 0.76 %

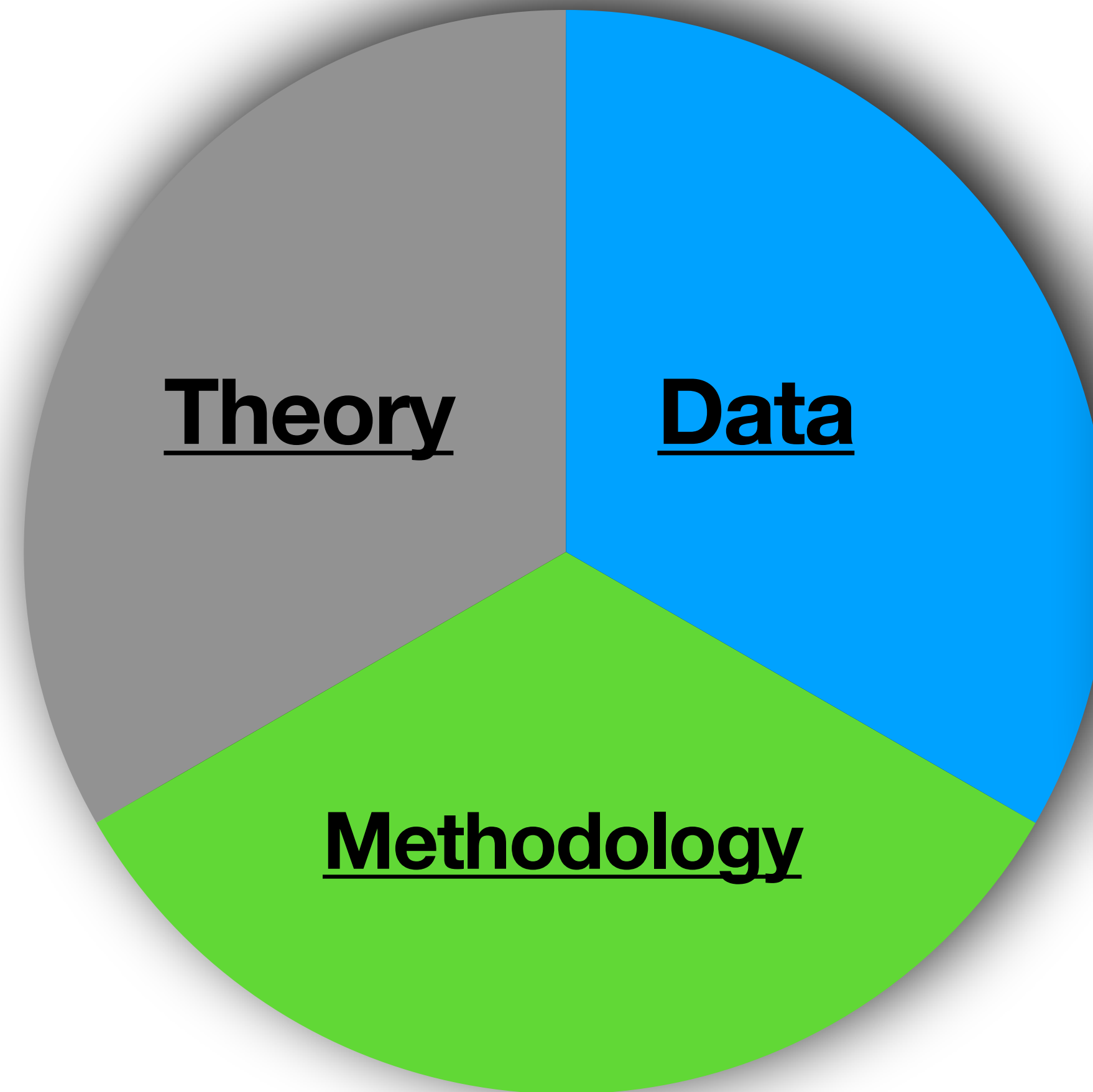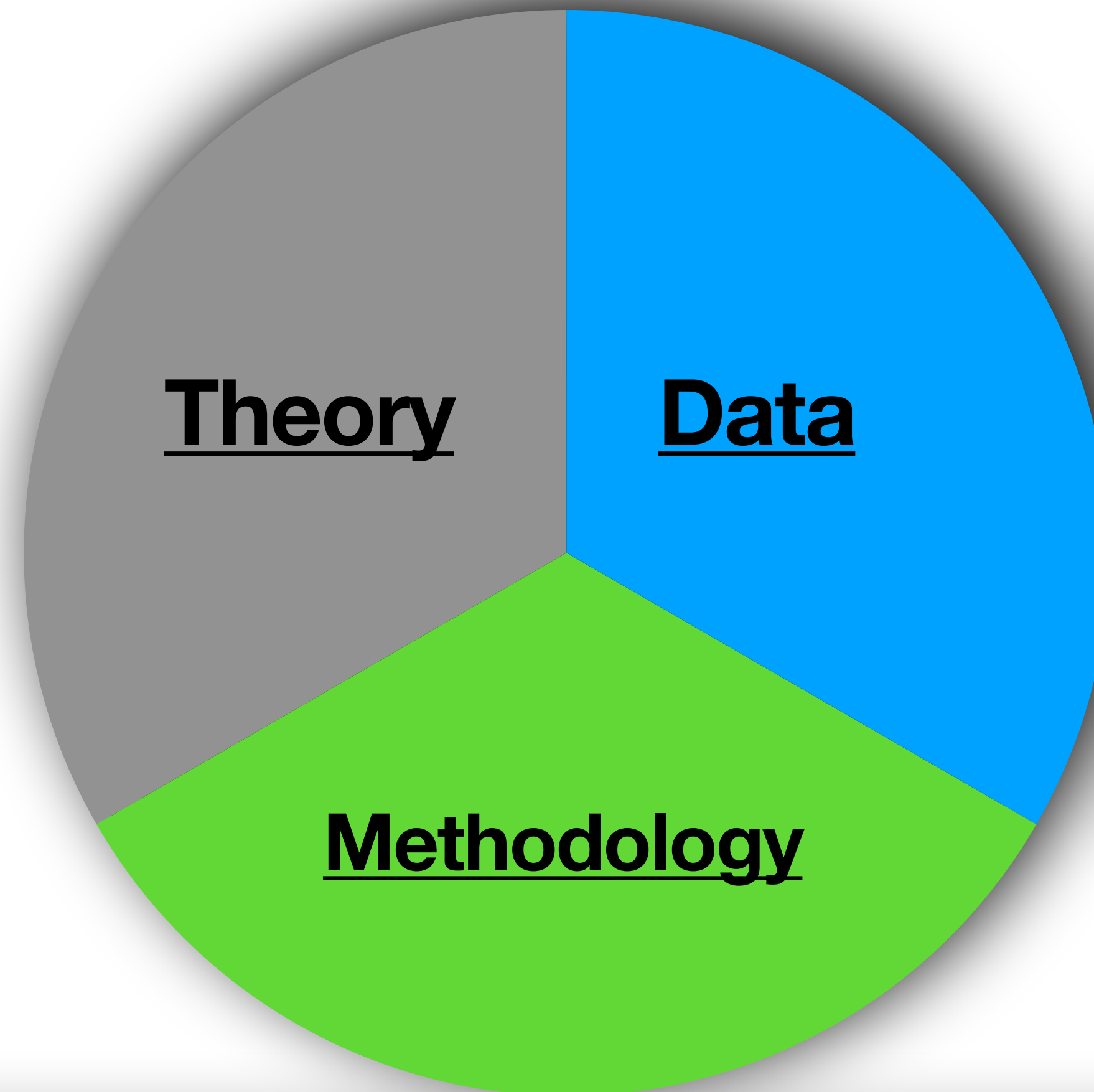| PDF set | $\alpha_s(m_Z)$ | PDF uncertainty | $g\ [GeV^2]$ | $q\ [GeV^4]$ |
|---------|-----------------|-----------------|--------------|--------------|
| MSHT20 [37] | 0.11839 | 0.00040 | 0.44 | −0.07 |
| NNPDF4.0 [84] | 0.11779 | 0.00024 | 0.50 | −0.08 |
| CT18A [29] | 0.11982 | 0.00050 | 0.36 |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

~ 1.7 %



## W mass determination



## weak mixing angle at 13 TeV

- Impact of jets vs diets at N3LO [arXiv:2312.12505]

- Impact of 13 TeV $t\bar{t}$ data [PRD 109 (2024)]

- Impact of future data (HL-LHC [Eur. Phys. J. C (2018) 78], EIC [PRD 103 (2021) 096005], FPF [arXiv:2309.09581])

**Theory**

**Data**

**Methodology**

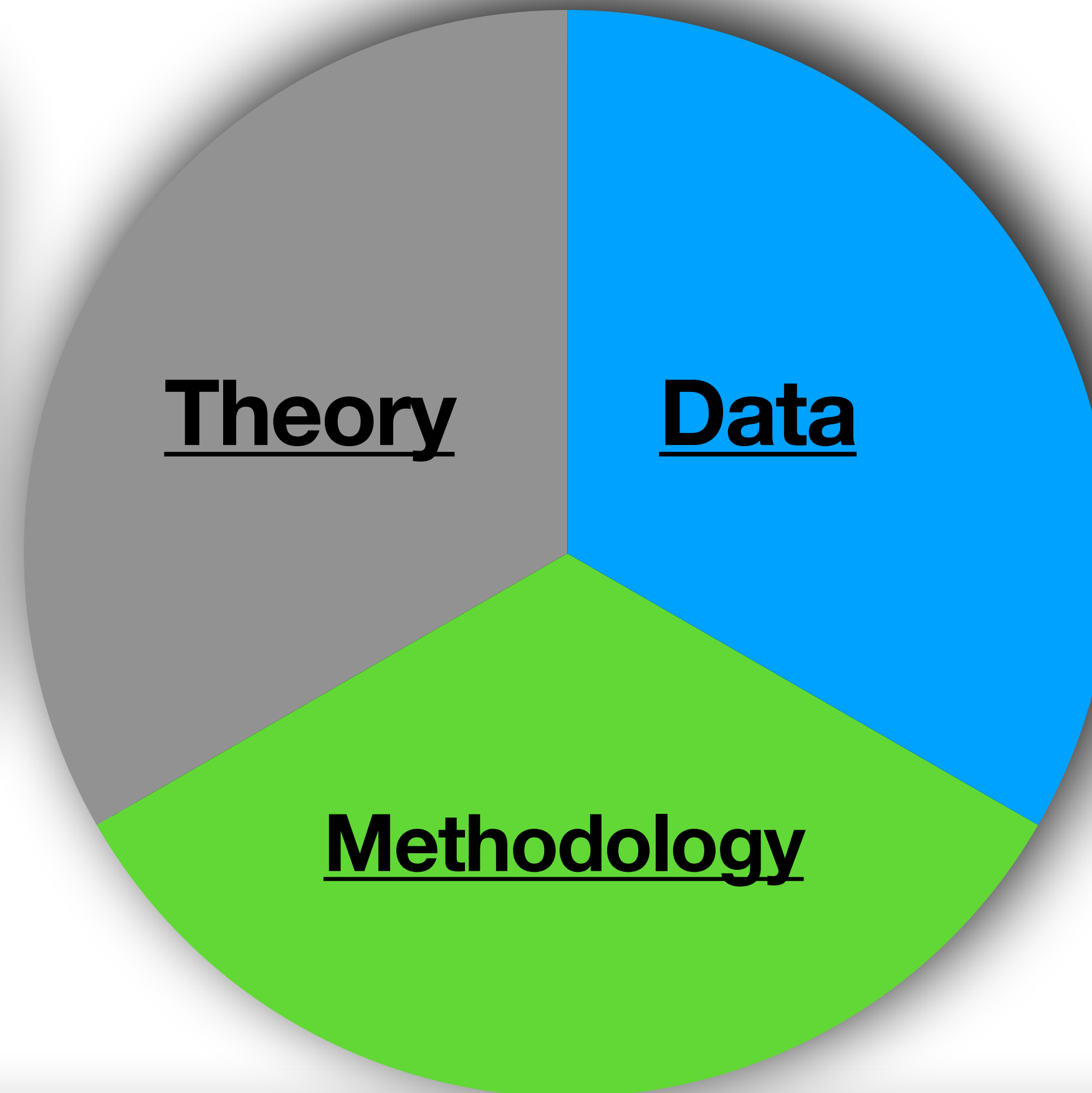- Impact of jets vs diets at N3LO [arXiv:2312.12505]

- Impact of 13 TeV $t\bar{t}$ data [PRD 109 (2024)]

- Impact of future data (HL-LHC [Eur. Phys. J. C (2018) 78], EIC [PRD 103 (2021) 096005], FPF [arXiv:2309.09581])

- Nonparametric regression [arXiv:2404.02964]

- Closure test [EPJC 82 (2022) 4, Talk by Lucian Harland-Lang, DIS2024]

**Theory**
- aN3LO [EPJC 83, arXiv:2402.18635]
- MHOU [arXiv:2401.10319]
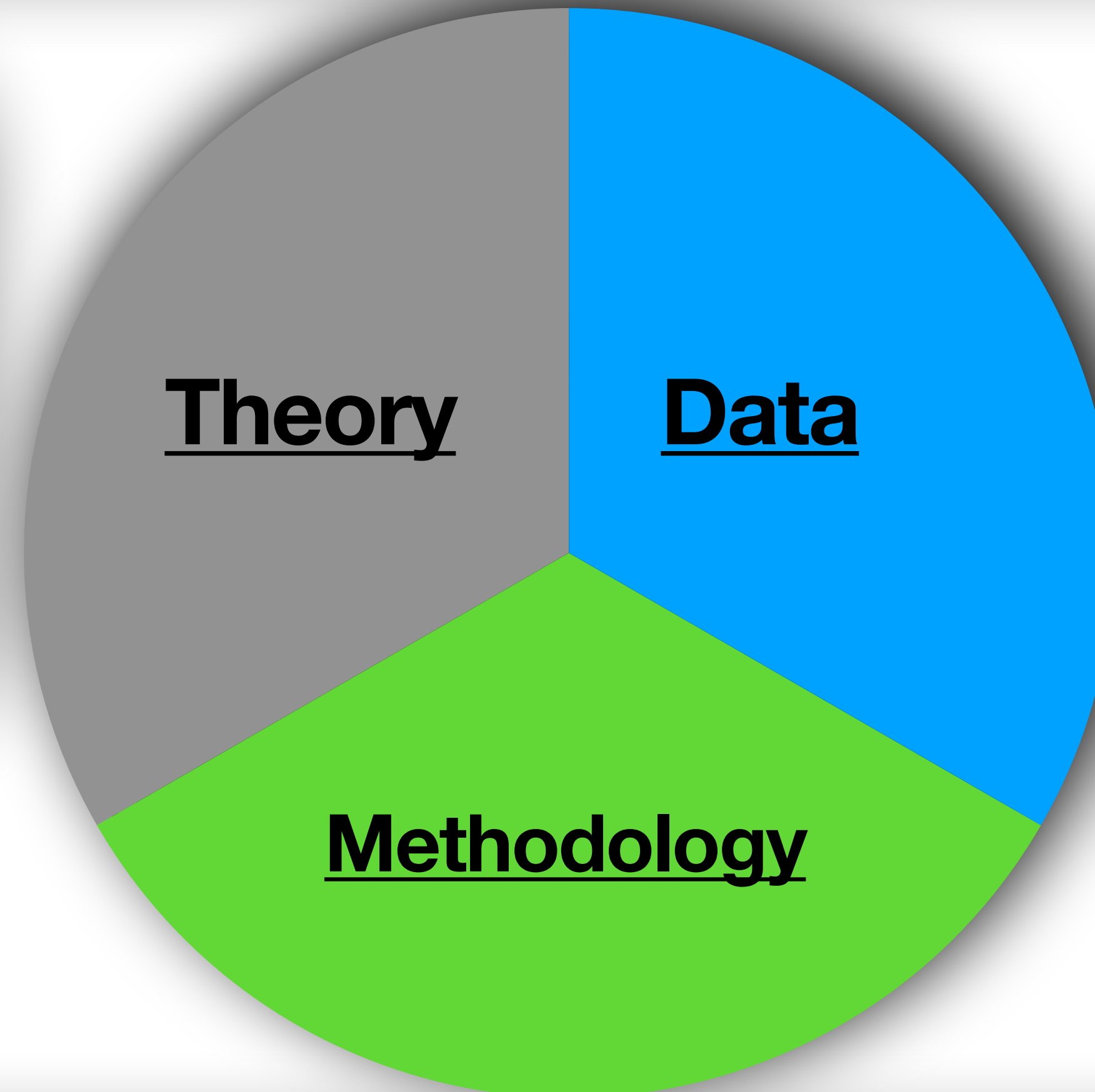- QED [arXiv:2401.08749]
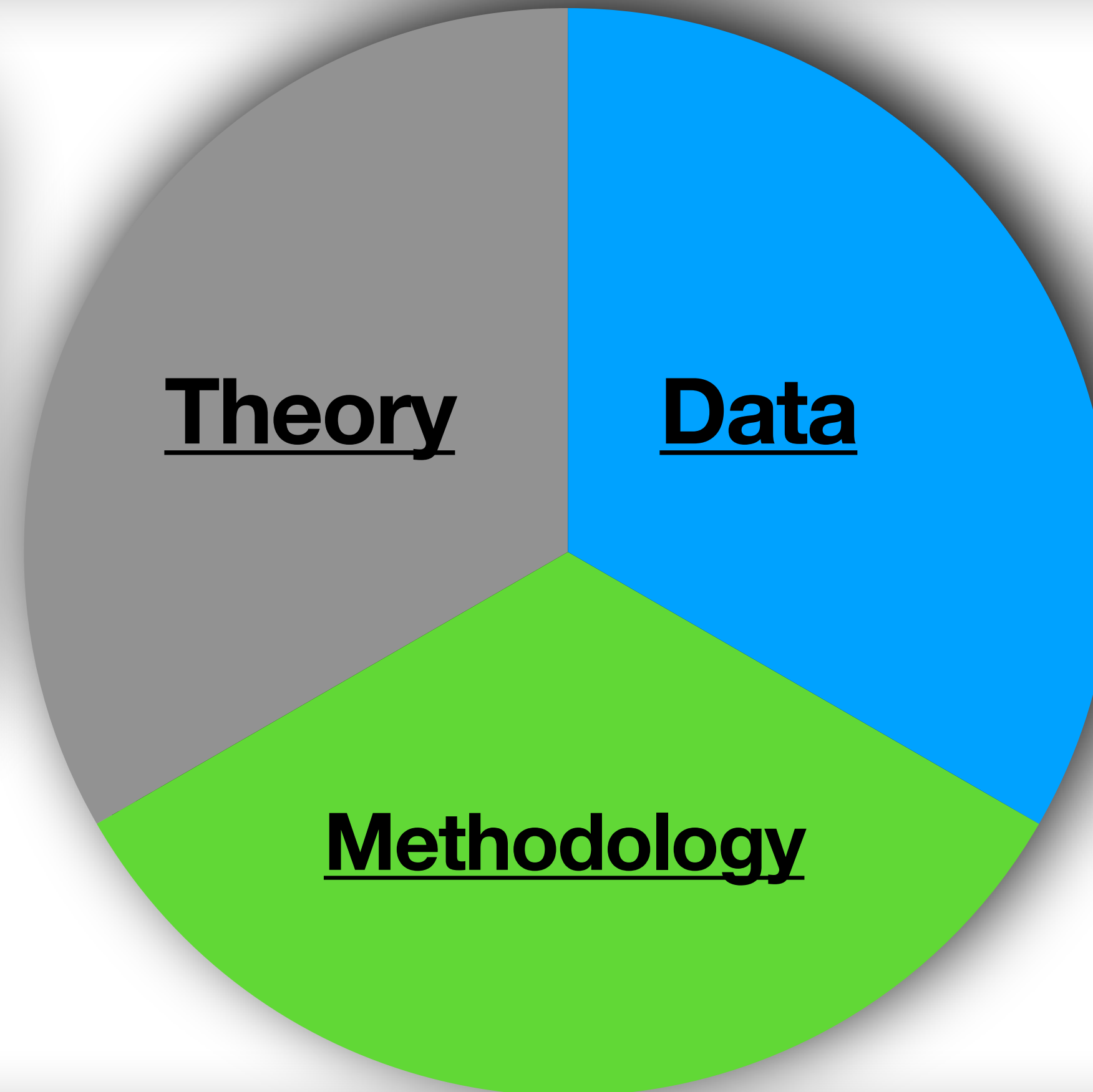- QED + aN3LO [arXiv:2404.02964]

**Data**
- Impact of jets vs diets at N3LO [arXiv:2312.12505]
- Impact of 13 TeV $t\bar{t}$ data [PRD 109 (2024)]
- Impact of future data (HL-LHC [Eur. Phys. J. C (2018) 78], EIC [PRD 103 (2021) 096005], FPF [arXiv:2309.09581])

**Methodology**
- Nonparametric regression [arXiv:2404.02964]
- Closure test [EPJC 82 (2022) 4, Talk by Lucian Harland-Lang, DIS2024]

# 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 $p(f|D)$

Prior on the model

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

Posterior of model given the data

→ Introduce probability distribution on a space of functions

→ Build a suitable prior

→ Use Bayes' theorem

# Gaussian Processes

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

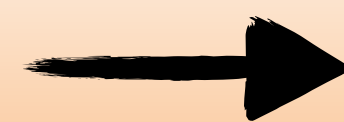**Parameters $\mathbf{f}$**: stochastic variables representing values of the PDF on a grid of points

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

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

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

**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

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

# 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

## Reconstructing QCD spectral functions with Gaussian processes

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

## What about PDFs?

# Prior for PDF: a bad example

$$m(x) = 0$$

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

Exponential quadratic

# Prior for PDF: a possibly better example

Gibbs Kernel

$$m(x) = 0$$

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



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

# Example: PDFs from DIS

Introduce an interpolation basis for $f$

$$F(x, Q^2) = \sum_i \int_x^1 dy\, C_i\left(\frac{x}{y}, \frac{Q}{\mu}, \alpha_s\right) f_i(y, \mu) \longrightarrow F_i = \sum_\alpha (FK)_{i\alpha}\, f(x_\alpha) = FK\,\mathbf{f}$$
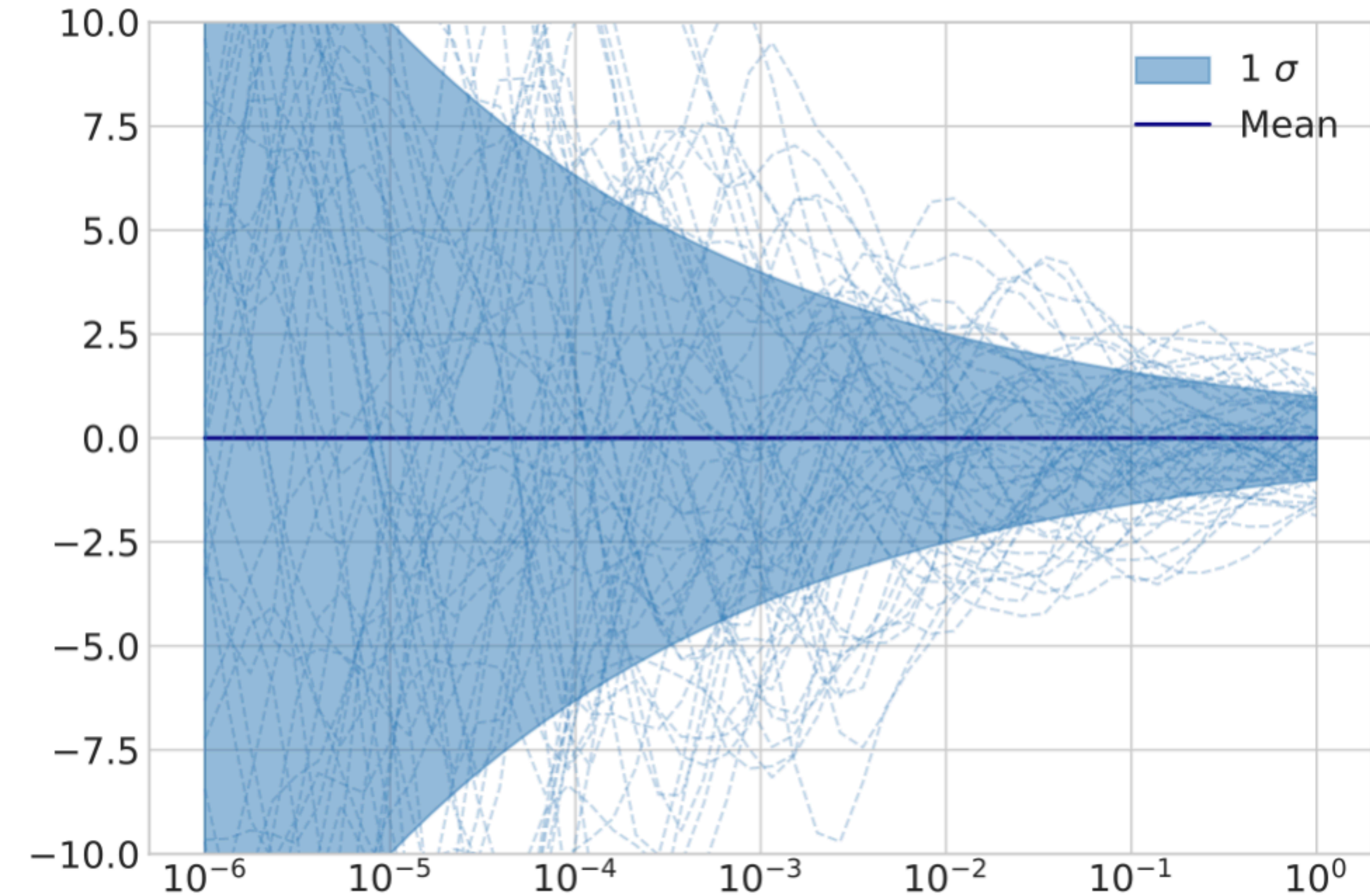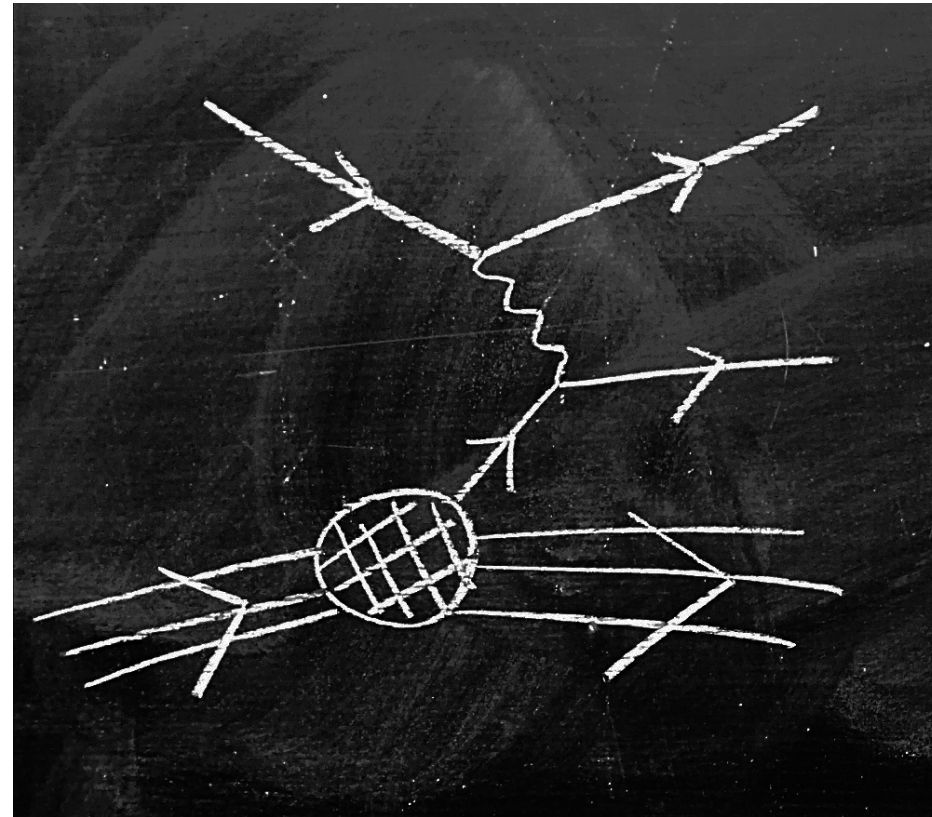


We use a chosen underlying law to generate pseudo-data

$$y = FK\,\mathbf{f_0} + \eta, \quad \eta \sim \mathcal{N}(0, C_y)$$

| 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] |

# Gaussian inference

| $\mathbf{f}$     Gaussian variable representing PDF on interpolation points $\mathbf{x}$ <br><br> $\mathscr{O} = FK\mathbf{f}$ | $\mathbf{f}*$ <br><br> Gaussian variable representing PDF on any set of points $\mathbf{x}*$ | $K\left(x, y; \theta\right)$ <br><br> Function modelling correlation | $y, \quad \epsilon \sim N(0, C_y)$ <br><br> Data and corresponding experimental error |
|---|---|---|---|

$$\begin{pmatrix} \mathbf{f}* \\ FK\mathbf{f} \end{pmatrix} \sim \mathscr{N}\left( 0, \begin{pmatrix} K_{\mathbf{x}*\mathbf{x}*} & K_{\mathbf{x}*\mathbf{x}}FK^T \\ FKK_{\mathbf{x}\mathbf{x}*} & FKK_{\mathbf{x}\mathbf{x}}FK^T \end{pmatrix} \right)$$
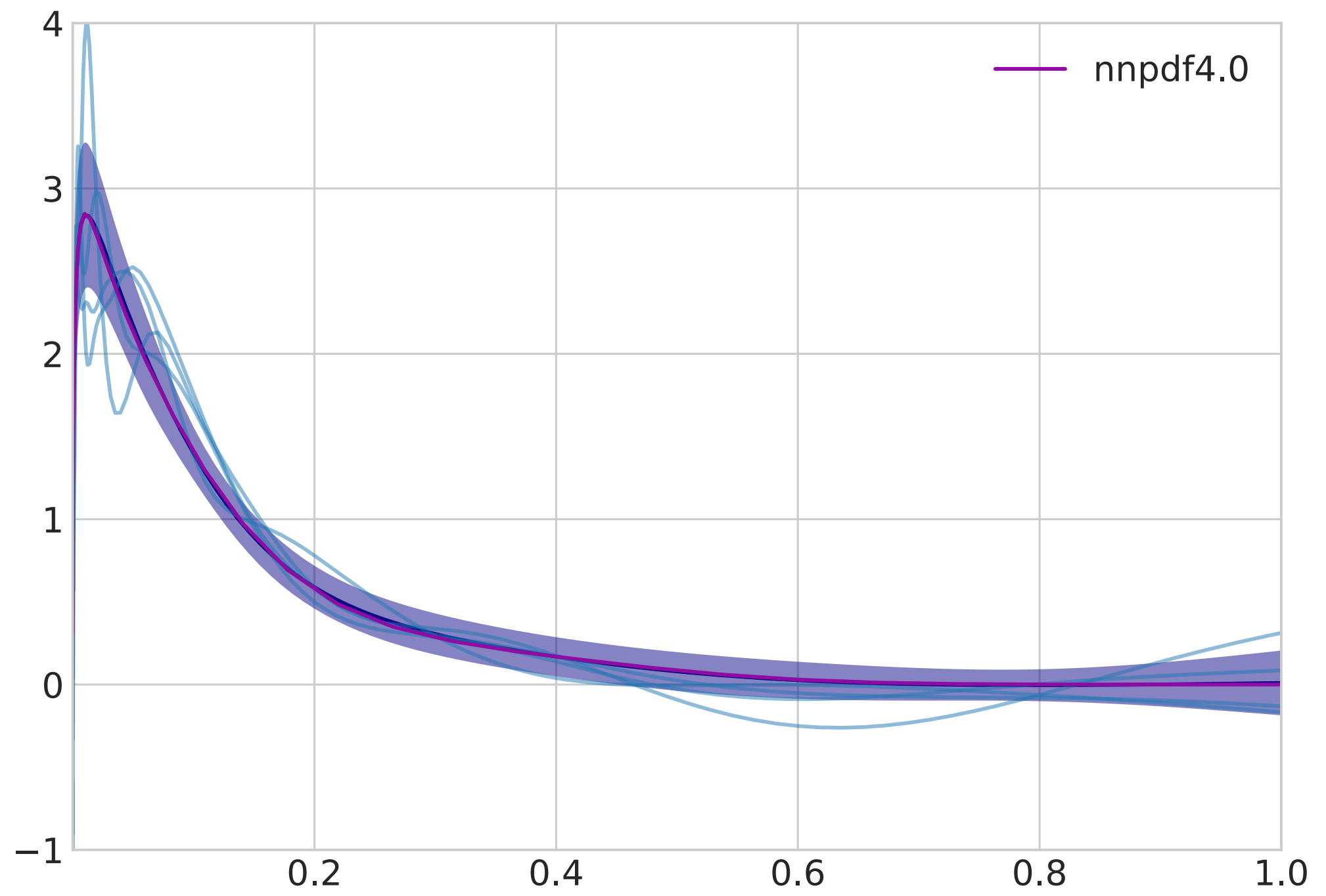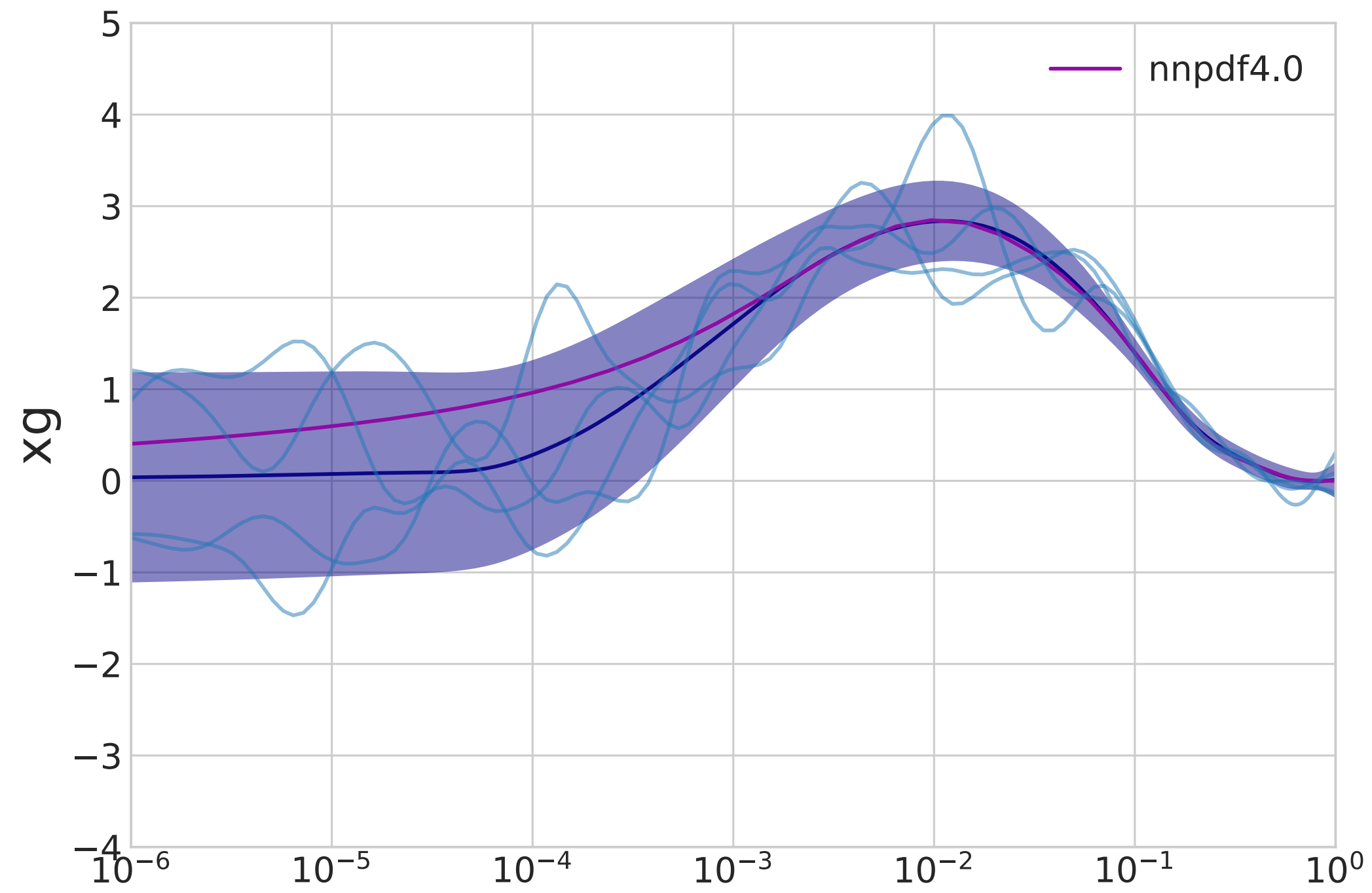
$$p\left(\mathbf{f}* \mid FK\mathbf{f} + \epsilon = y, \theta\right)$$

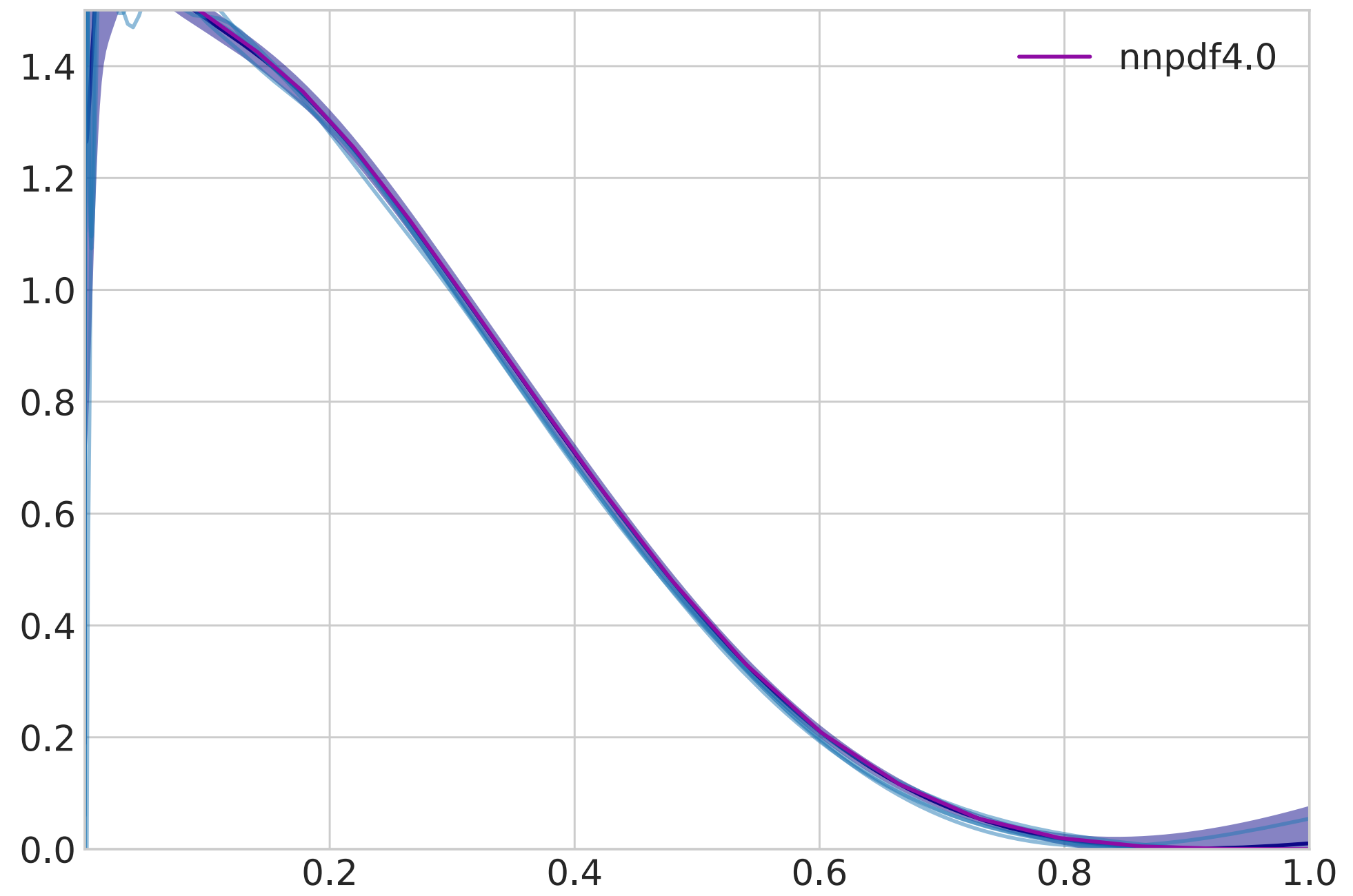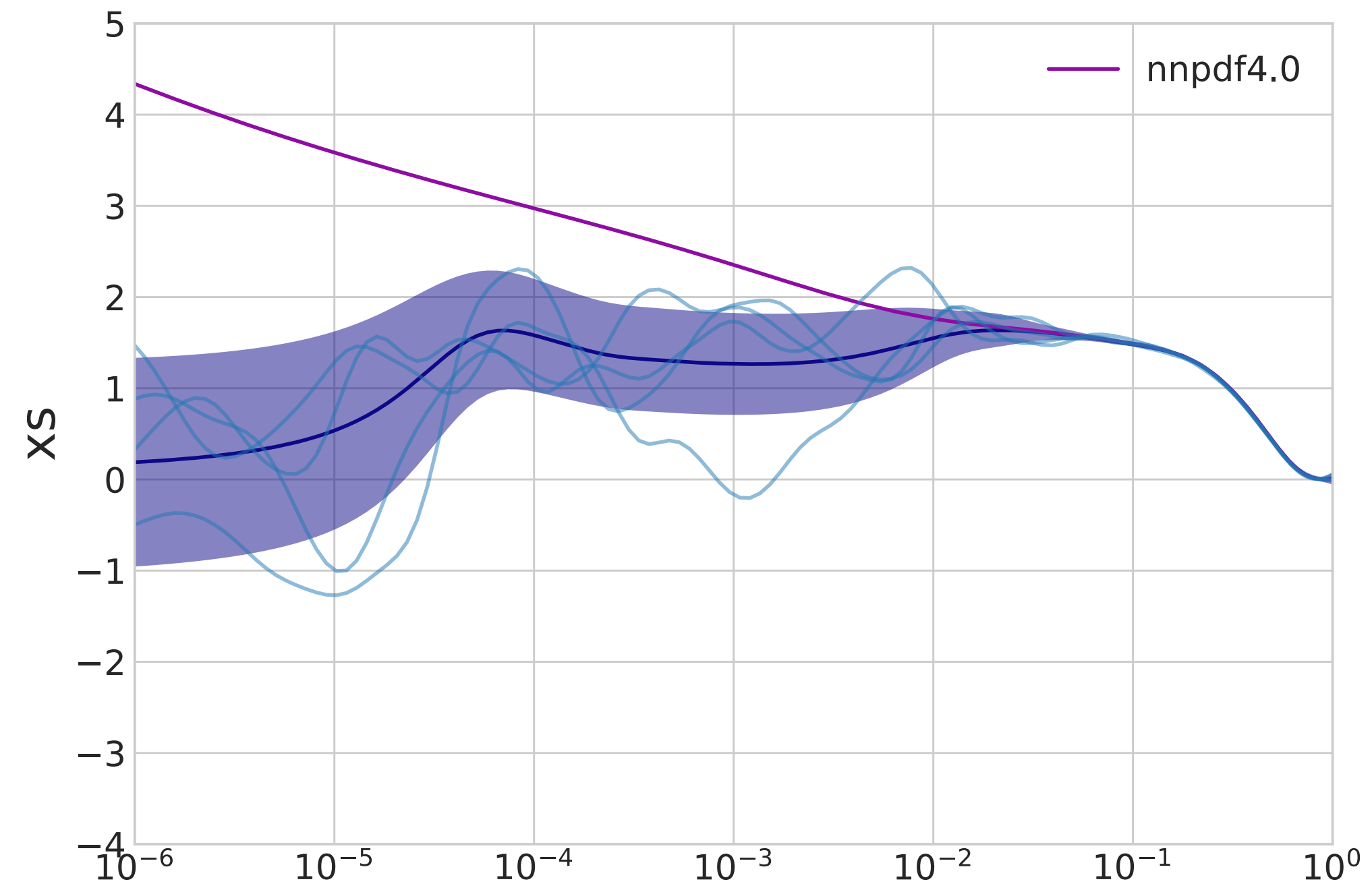This is a gaussian distribution. Its mean and covariance can be computed analytically

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

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

$p\left(\mathbf{f^*} \mid \mathbf{data}, \theta\right)$ **gluon**

# $p\left(\mathbf{f}^{*} \mid \mathbf{data}, \theta\right)$ singlet

# Inference on the hyperparameters

Joint probability distribution
of **f**\* and $\theta$

Posterior on the
hyperparamters given
the data

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

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

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

# Workflow

Build the prior as a function of hyperparameters:

- Choose kernel

- Encode theory constraints
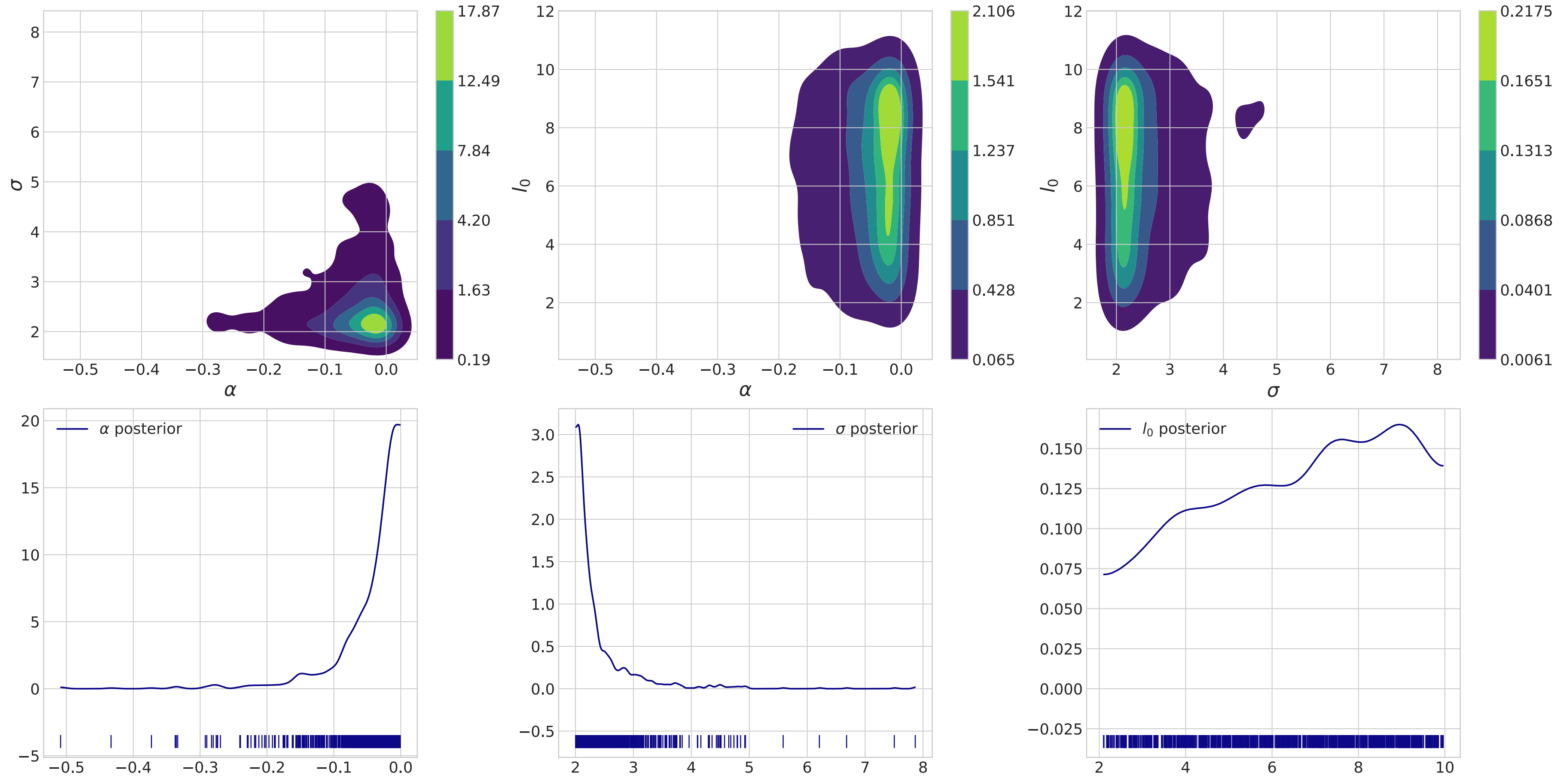
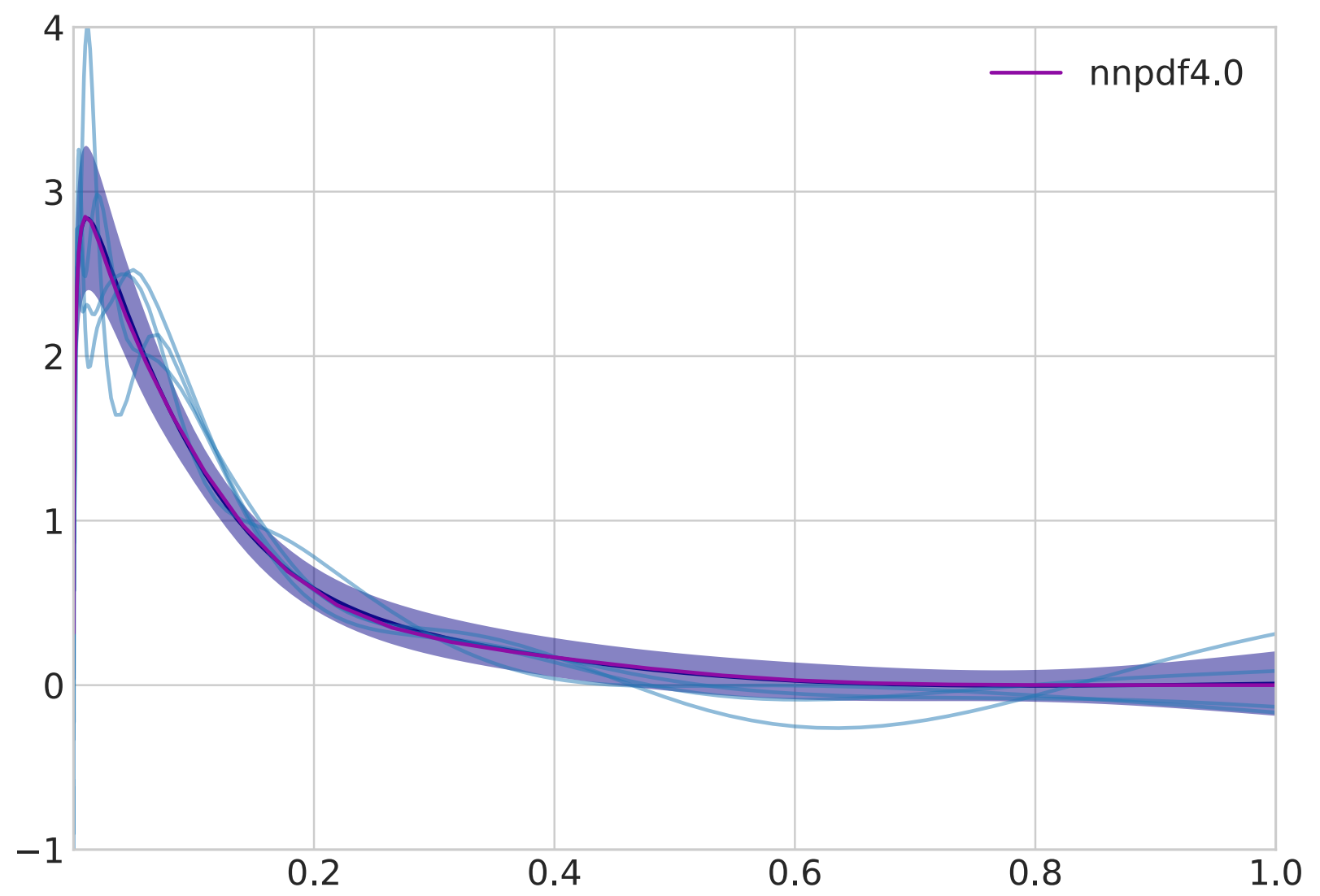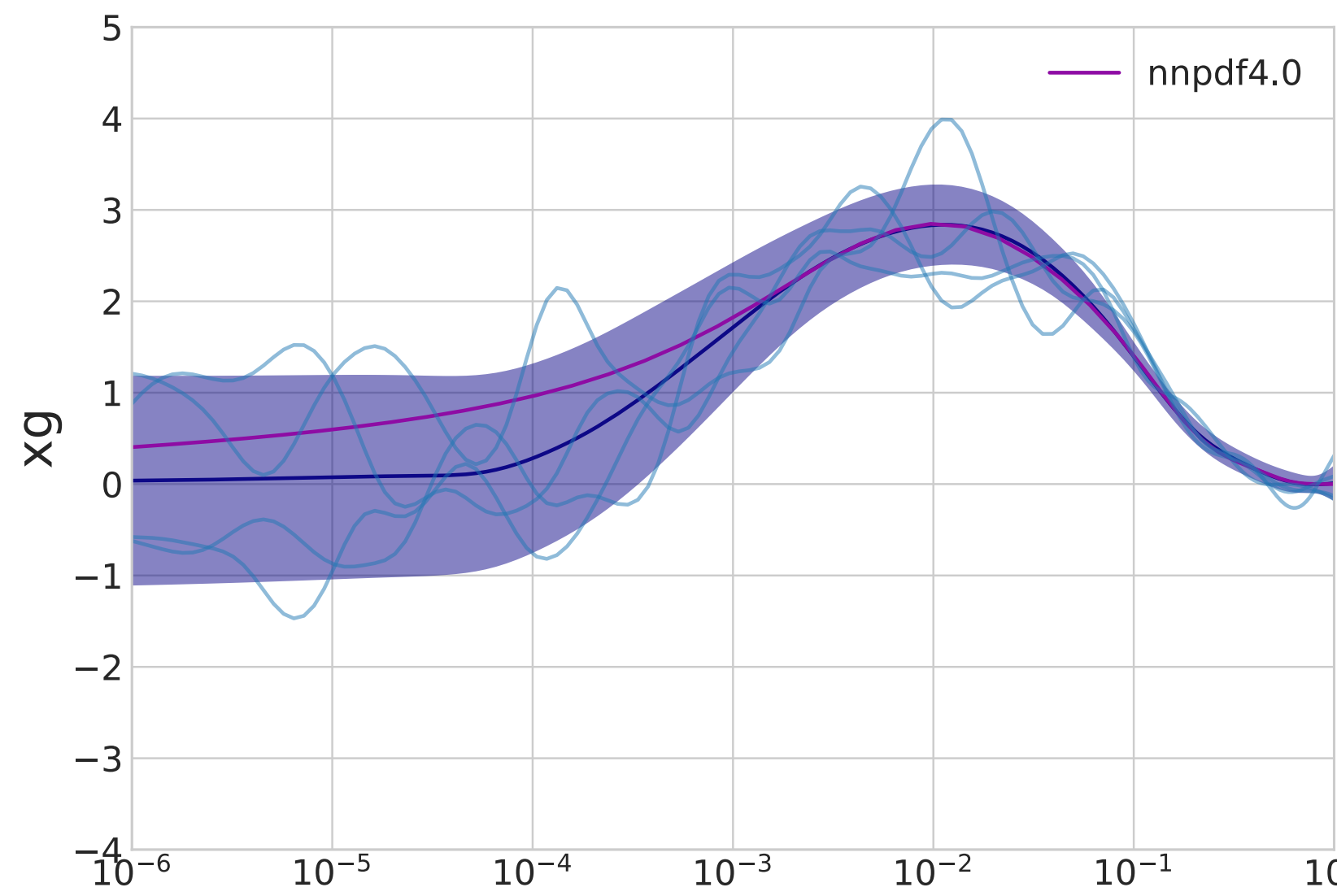Collect data and FK tables

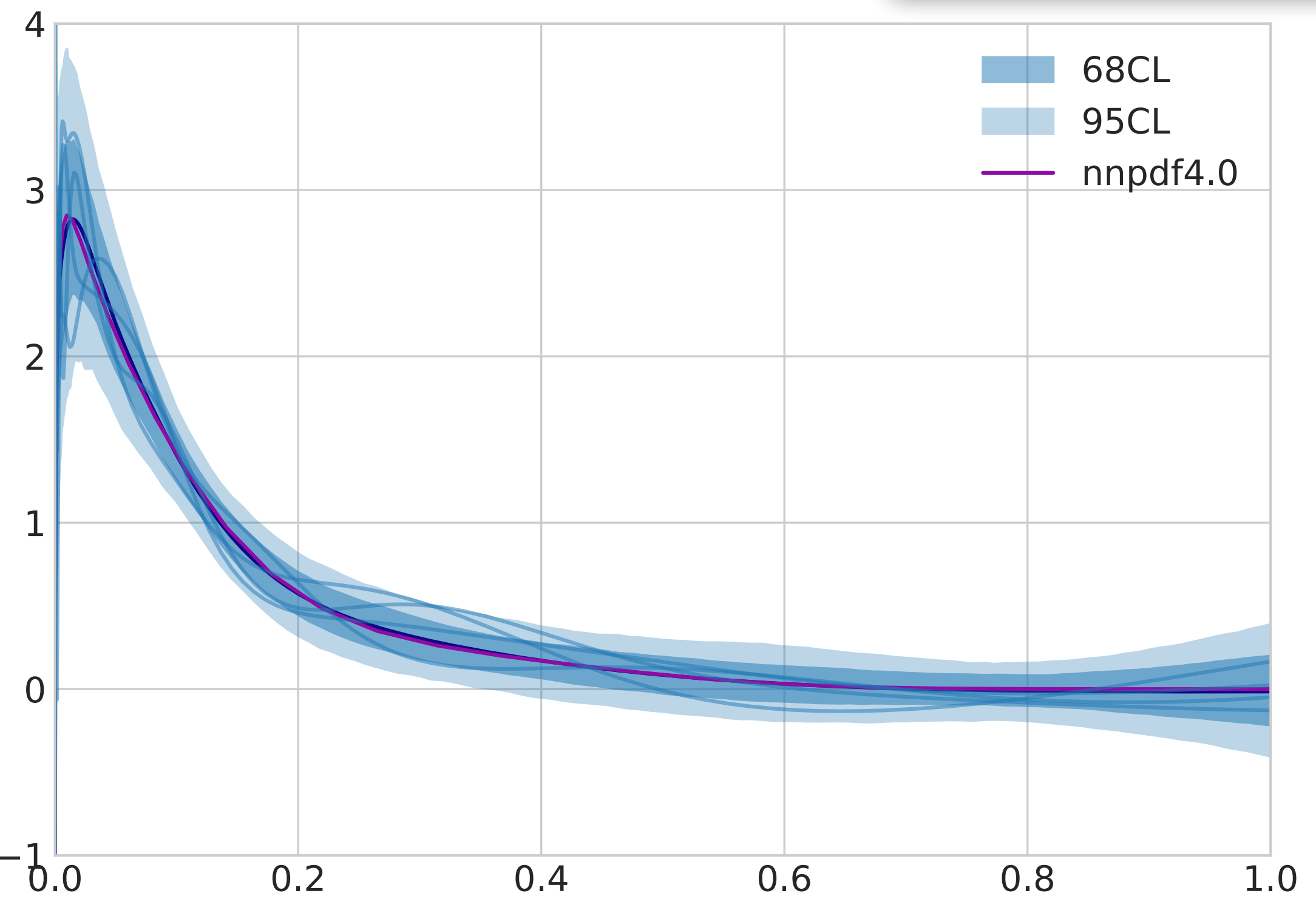Inference on hyperparameters
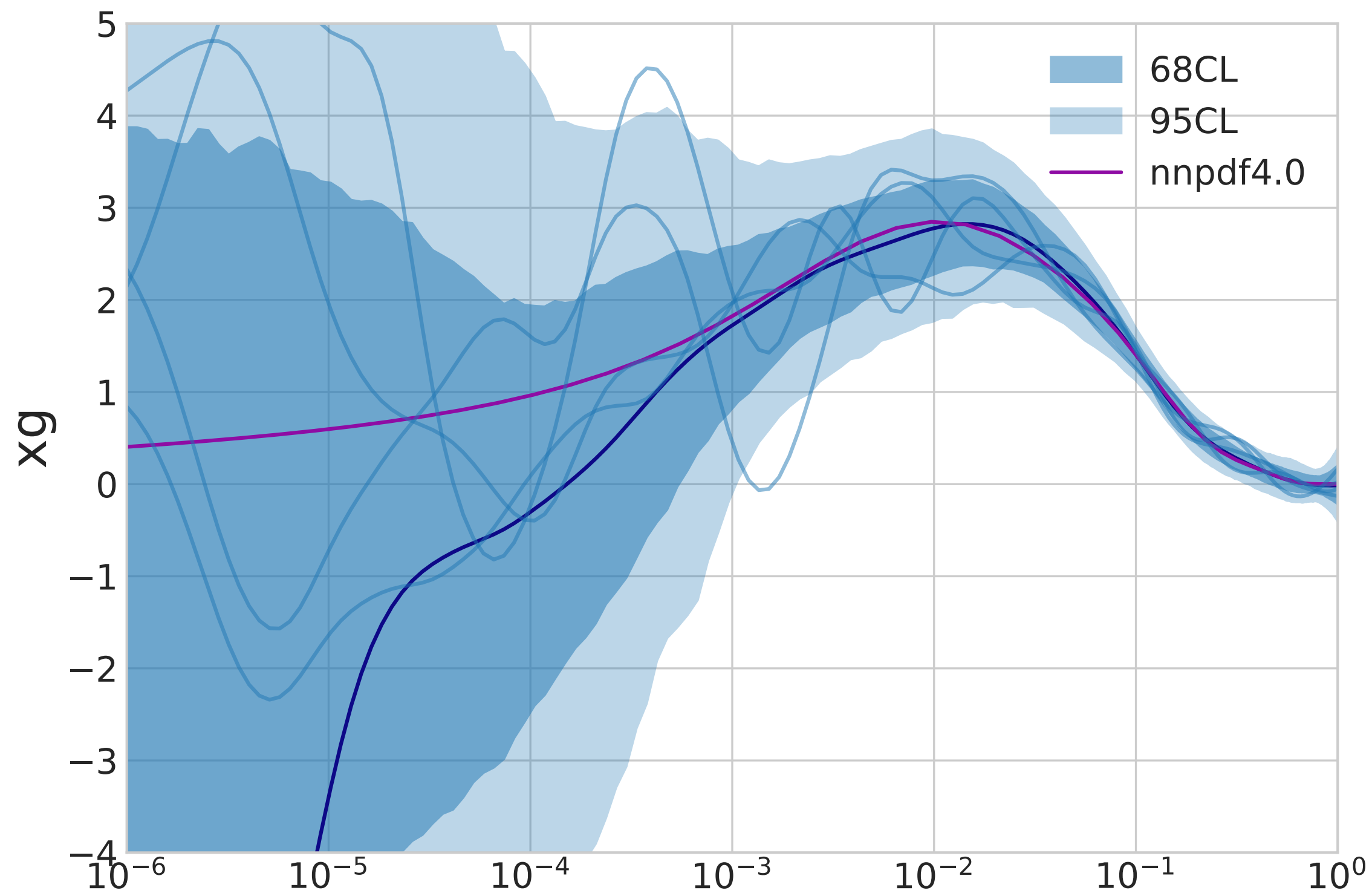
**PyMC**

For DIS this step is analytical

Inference on parameters

# Posterior for hyperparameters (gluon)

Fixed hyperparameters $p\left(\mathbf{f^*}\,|\,\mathrm{data},\theta\right)$

Full posterior $p\left(\mathbf{f^*},\theta\,|\,\mathrm{data}\right)$

# Posterior for hyperparameters (singlet)

Fixed hyperparameters $p\left(\mathbf{f}^*\mid\mathrm{data},\theta\right)$

Full posterior $p\left(\mathbf{f}^*,\theta\mid\mathrm{data}\right)$

# Theory constraints

## Kinetic limit

$$f(1) = 0$$

Additional linear constrain on the PDF: can be implemented directly in the FK table

## Sum rules

$$g \sim GP\left(0, K\left(x, y\right)\right) \quad \longrightarrow \quad \frac{dg}{dx} \sim GP\left(0, \partial_x \partial_y K\left(x, y\right)\right)$$

Sum rules can be implemented as additional linear constrains on the primitive of the PDF

# Global fits

$$\sigma = \sum_{i,j} \int dx_1 dx_2 \, f_i\left(x_1, \mu\right) f_j\left(x_2, \mu\right) \hat{\sigma}\left(x_1, x_2, \frac{Q}{\mu}\right) \times \left(1 + \mathcal{O}\left(\Lambda/M\right)^p\right)$$

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

This bit is not a gaussian distribution any longer

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

# Summary and future work

- A Bayesian methodology based on GP has been presented

-  Might be helpful to further understand and quantify PDF uncertainties

- It is possible to implement physical constraints such as sum rules and kinetic limit

- Preliminary study on an extended DIS dataset is being finalised

- Systematic study of different possible kernels

- Comparison with non Bayesian methodologies. Are there any differences?

- Implementation of a full global analysis

Backup slides

# Fit quality

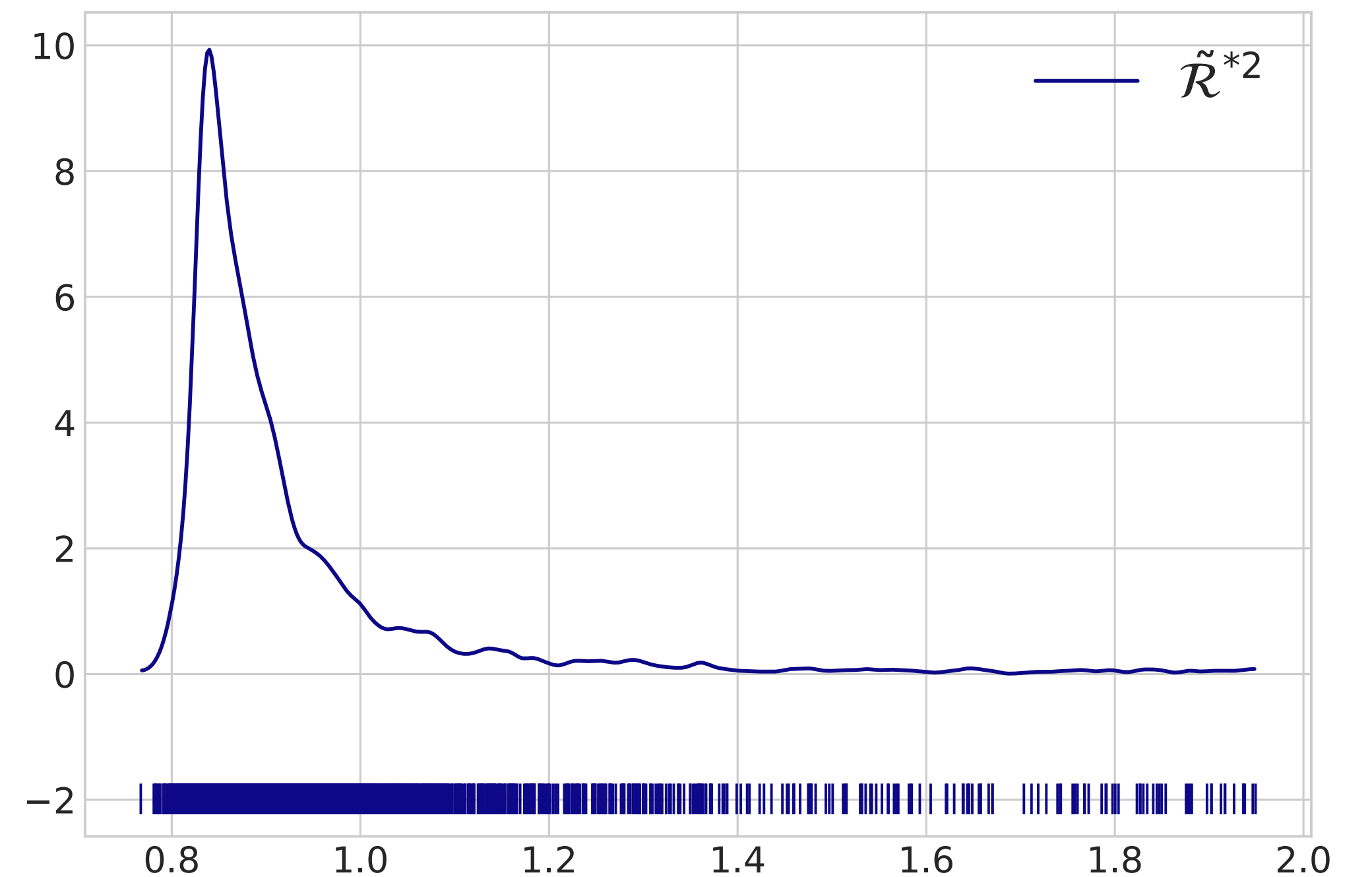$$\frac{S}{dof} = \frac{1}{N_{\text{data}}} \Big( (\mathbf{m} - \tilde{\mathbf{m}})^T K_{xx}^{-1} (\mathbf{m} - \tilde{\mathbf{m}}) + (y - FK\tilde{\mathbf{m}})^T C_Y^{-1} (y - FK\tilde{\mathbf{m}}) \Big)$$

# Generalisation on unseen data

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

# Further possible applications

- simultaneous fits of PDFs and Wilson coefficients

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

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

Zahari Kassabov,[a] Maeve Madigan,[a] Luca Mantani,[a] James Moore,[a] Manuel Morales Alvarado,[a] Juan Rojo[b,c] and Maria Ubiali[a]

*JHEP* 05 (2023) 205

- Inverse problems relevant for the lattice community

**Reconstructing QCD Spectral Functions with Gaussian Processes**

Jan Horak,[1] Jan M. Pawlowski,[1,2] José Rodríguez-Quintero,[3] Jonas Turnwald,[1] Julian M. Urban,[1,*] Nicolas Wink,[1] and Savvas Zafeiropoulos[4]

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

# Decomposition of PDF uncertainty

$$\tilde{K} = \boxed{\left(I - R_{xx}\right) K_{\mathbf{xx}} \left(I - R_{xx}\right)^T} + \boxed{a_{xx}^T C_y a_{xx}}$$

<span style="color:magenta">Methodology</span>　　　<span style="color:gray">Experimental error</span>

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

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