



# ML/AI activities at Nikhef Theory

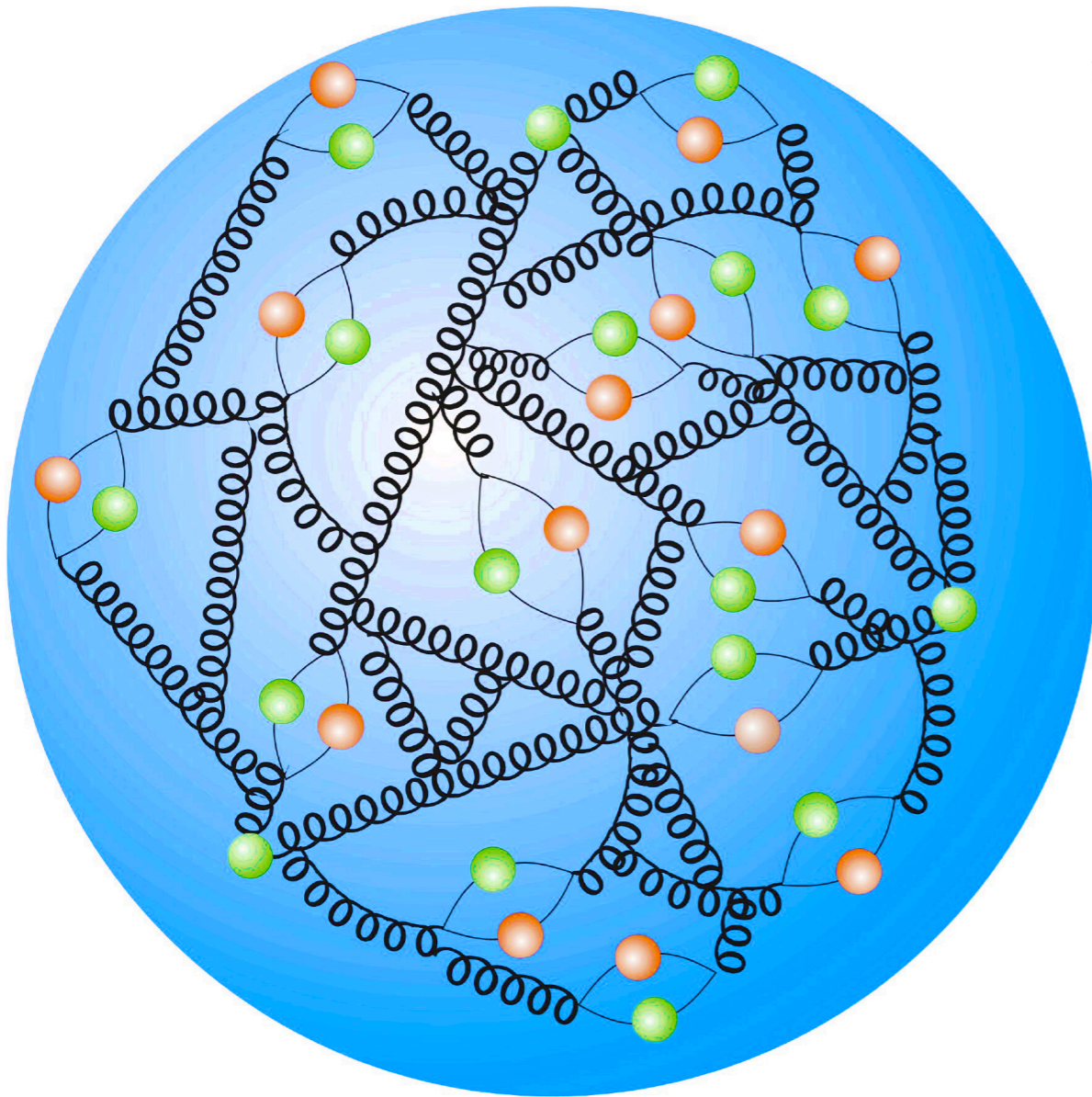
**Juan Rojo, VU Amsterdam & Nikhef**

**Nikhef ML/AI Working Group - Kick-off Meeting**

**Nikhef, 6th September 2024**

# The many faces of the proton

**QCD** bound state of **quarks** and **gluons**



- ☑ **Valence quarks (up and down)** give the proton its quantum numbers (e.g. electric charge)

$$|\Psi\rangle \approx |uud\rangle$$

$$Q_p = +1 \quad \begin{array}{l} Q_u = +2/3 \\ Q_d = -1/3 \end{array}$$

- ☑ **Sea quarks (antiup, antidown, strange, ...)** arise from quantum fluctuations
- ☑ Tightly held together by **gluons**, can only be broken in extremely energetic collisions

# Parton Distributions

$$g(x, Q)$$

**Probability of finding a gluon inside a proton**, carrying a fraction  $x$  of the proton momentum, when probed with energy  $Q$

**Energy** of hard-scattering reaction:  
inverse of resolution length

$x$ : fraction of proton momentum carried by gluon

Dependence on  $x$  fixed by **non-perturbative QCD dynamics**: extract from experimental data

$$g(x, Q_0, \{a_g\}) = f_g(x, a_g^{(1)}, a_g^{(2)}, \dots)$$

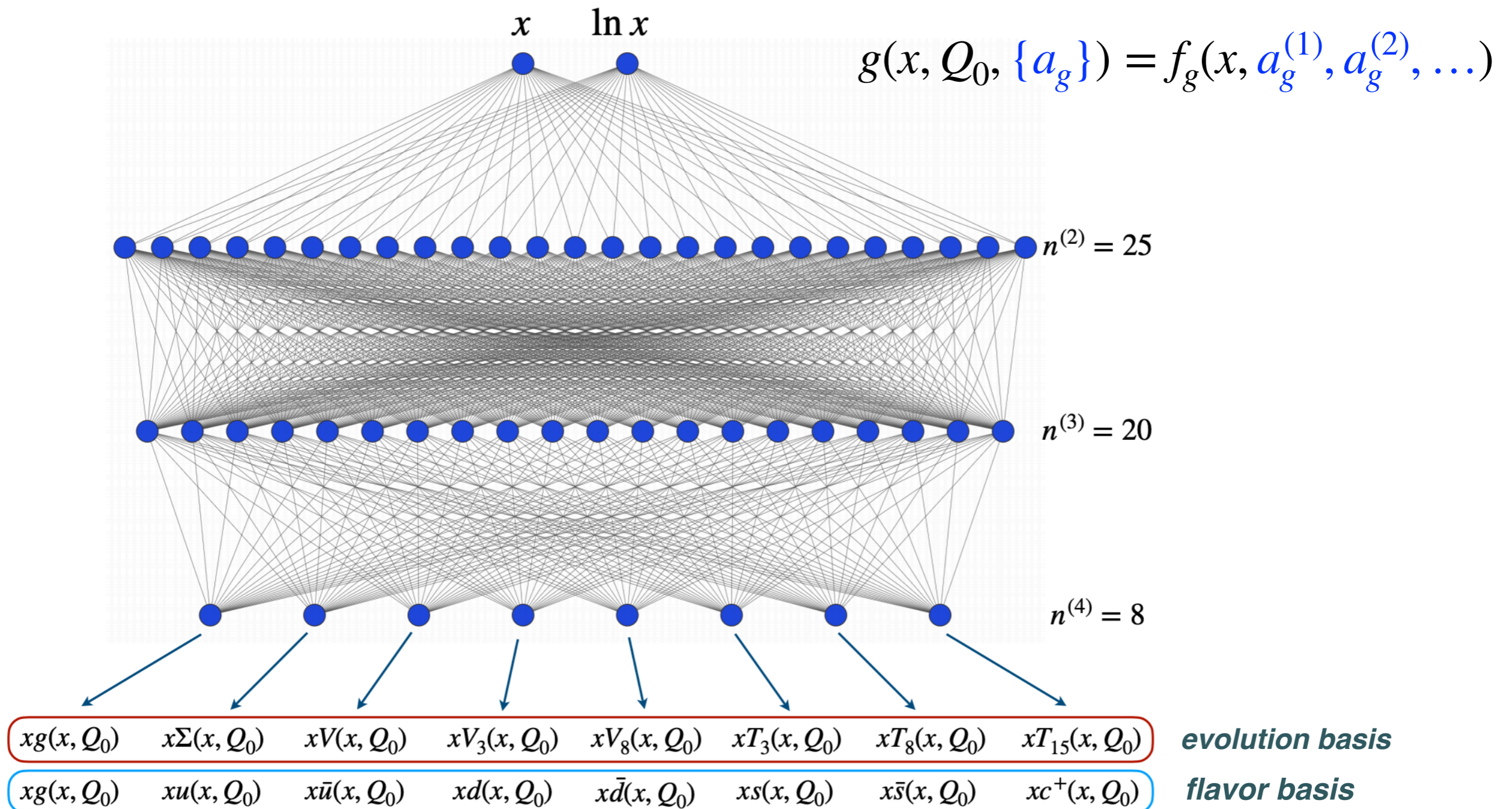
*constrain from data*

📍 Dependence with **resolution scale  $Q$** : DGLAP evolution, computable from first principles

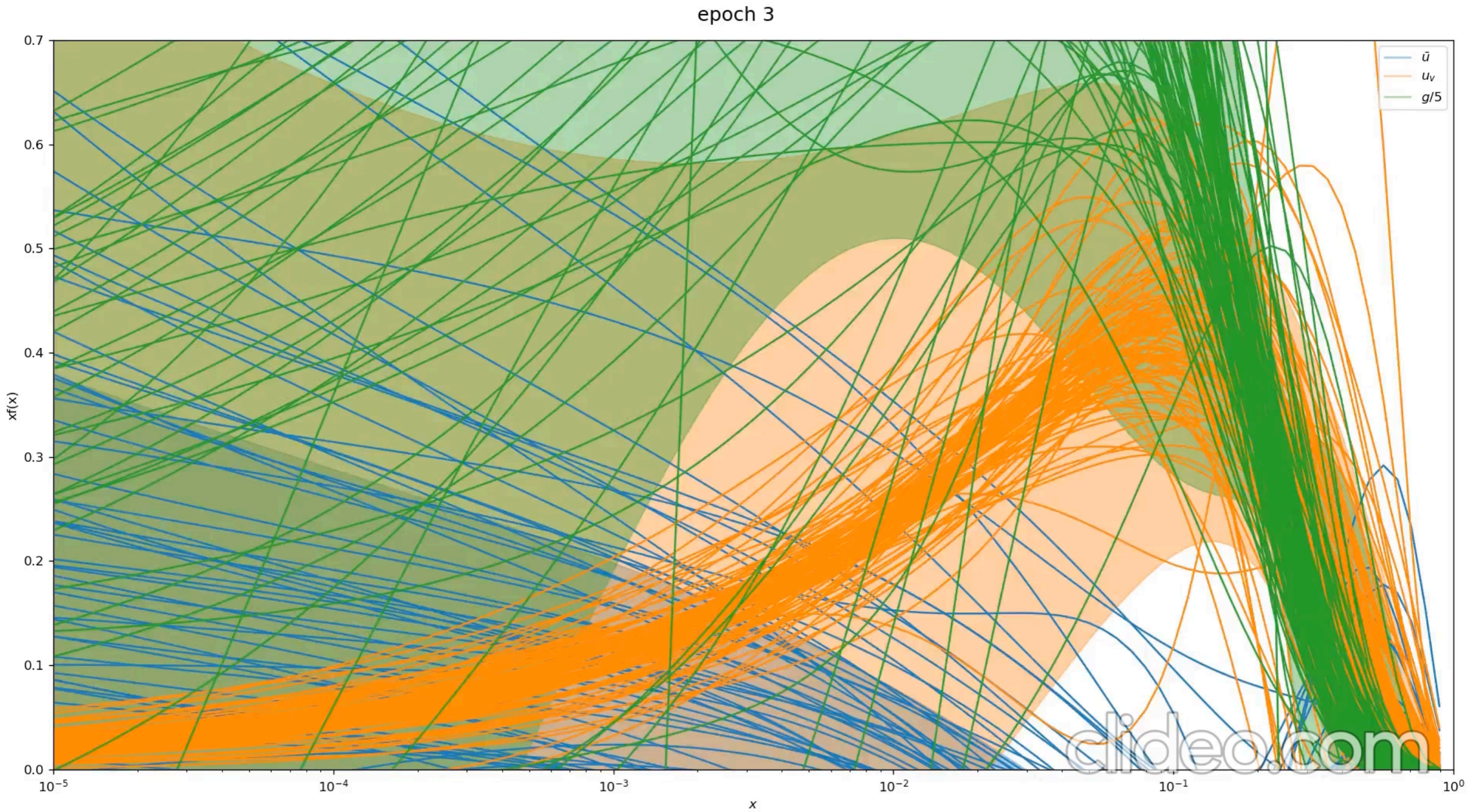
📍 **Energy conservation** and **quark number conservation** are fixed boundary conditions

# ML Proton Structure

- ☑ Model-independent PDF parametrisation with neural networks as **universal unbiased interpolants**
- ☑ **Stochastic Gradient Descent** via TensorFlow for neural network training
- ☑ Automated model **hyperparameter optimisation**: NN architecture, minimiser, learning rates ...



# Machine Learning PDFs



Error estimate based on **Monte Carlo replica method** (band: standard deviation over the MC replicas)

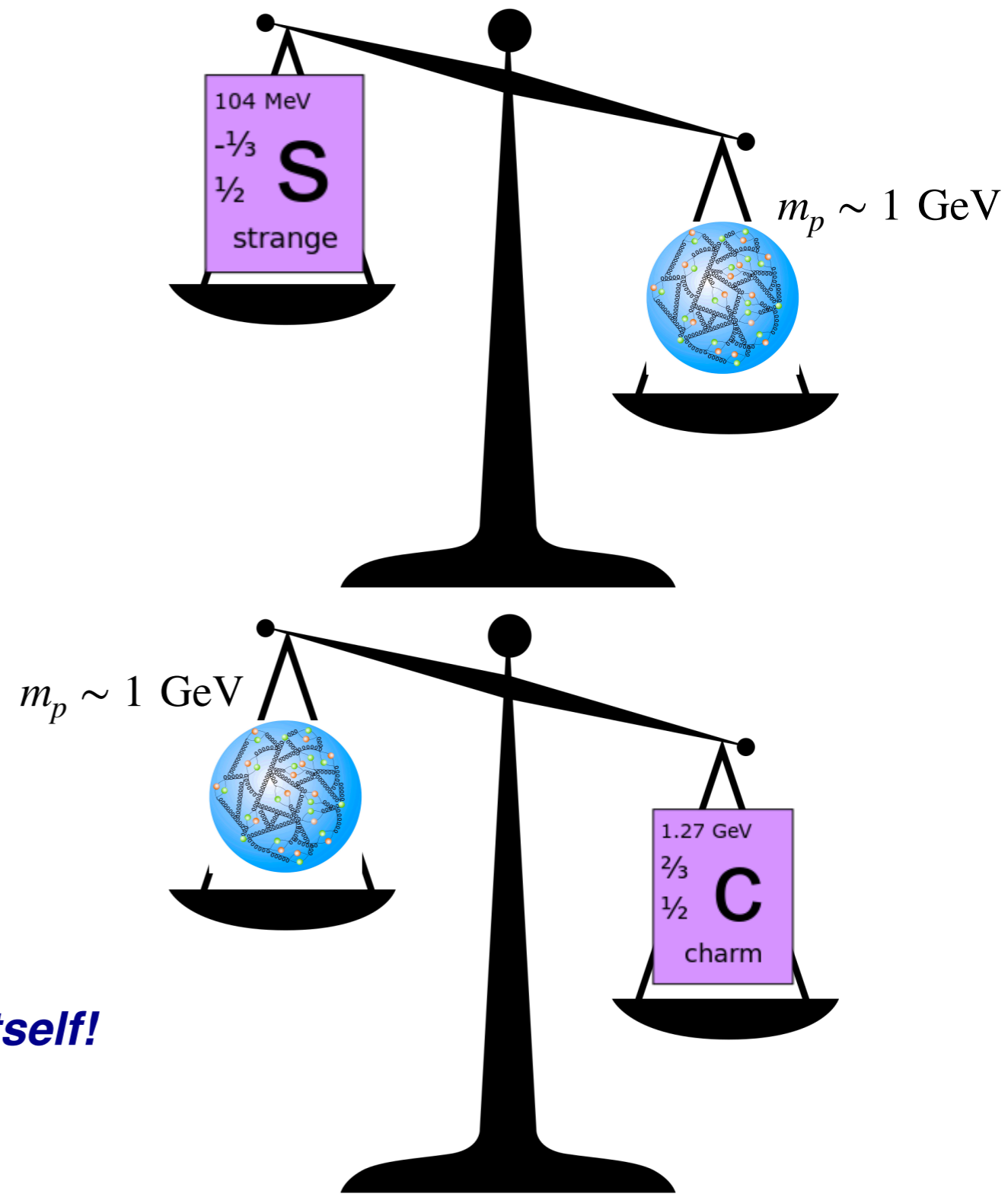
*each curve is a separately trained neural network*

# The charm content of the proton

common assumption: the proton wave function does not contain charm quarks

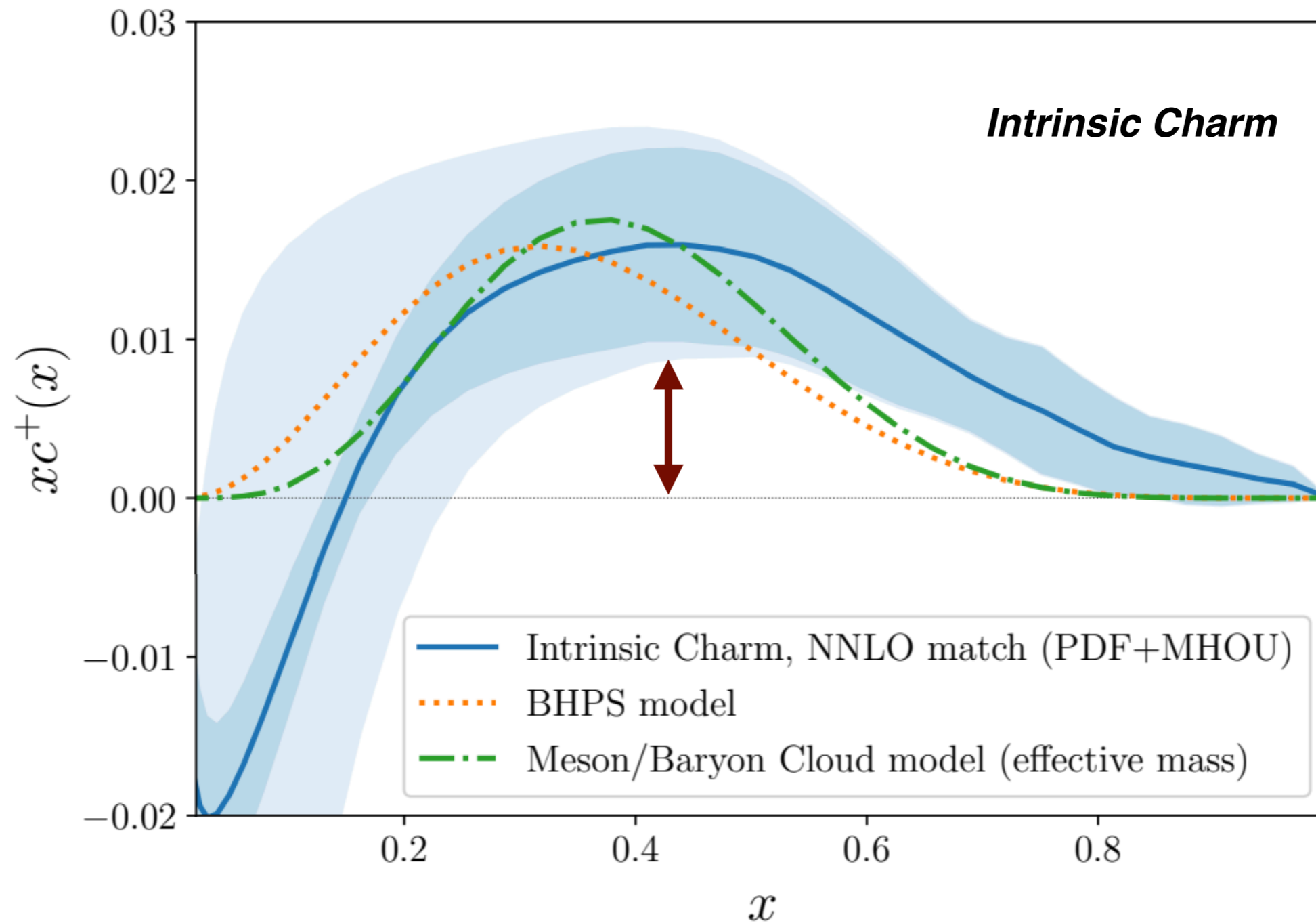
the proton contains **intrinsic up, down, strange (anti-)quarks** but **no intrinsic charm quarks**

mass →	2.4 MeV	1.27 GeV	171.2 GeV
charge →	$\frac{2}{3}$	$\frac{2}{3}$	$\frac{2}{3}$
spin →	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{2}$
name →	<b>u</b> up	<b>c</b> charm	<b>t</b> top
Quarks	4.8 MeV	104 MeV	4.2 GeV
	$-\frac{1}{3}$	$-\frac{1}{3}$	$-\frac{1}{3}$
	<b>d</b> down	<b>s</b> strange	<b>b</b> bottom

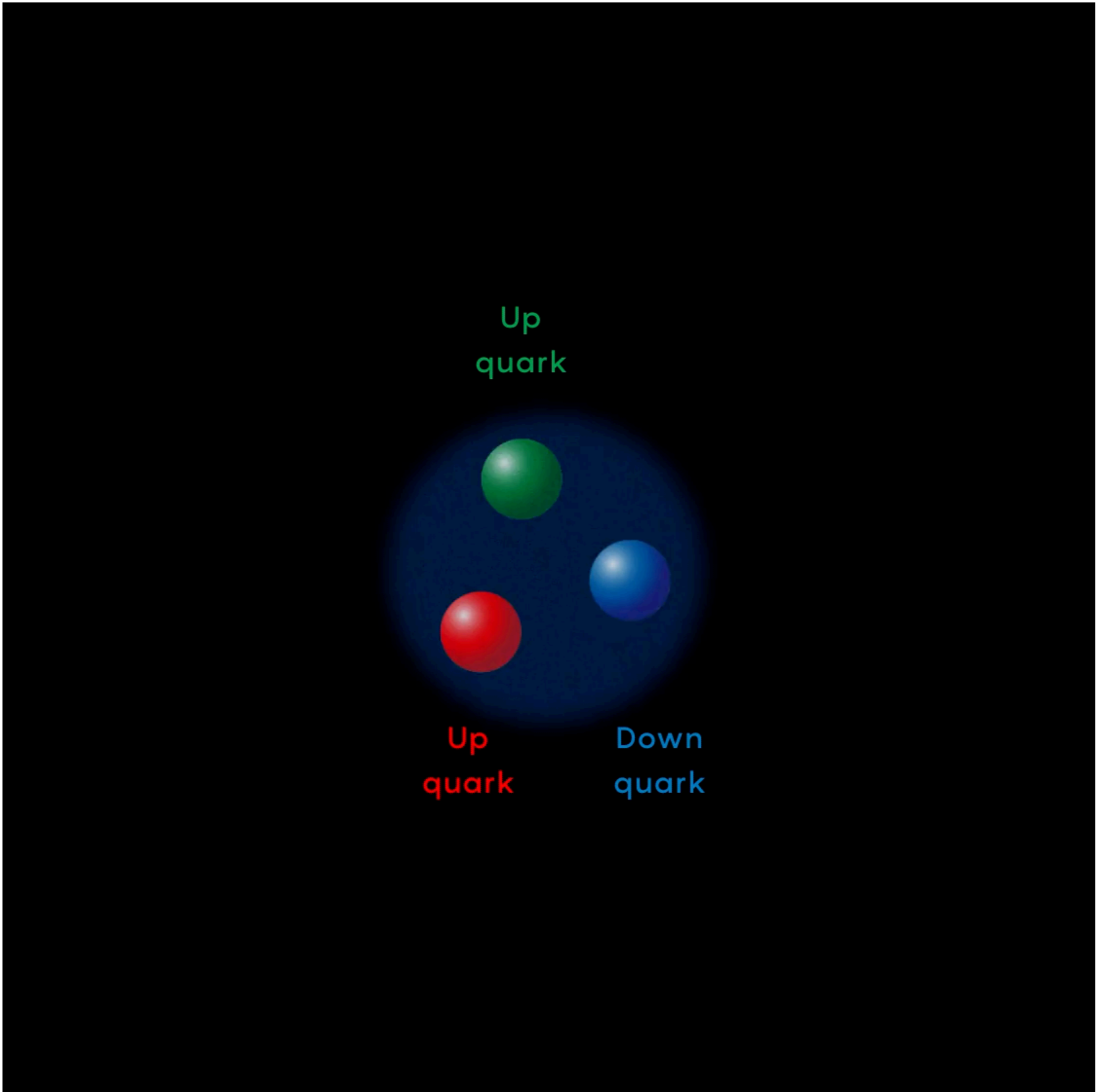


***charm quarks heavier than the proton itself!***

# Intrinsic Charm in the Proton



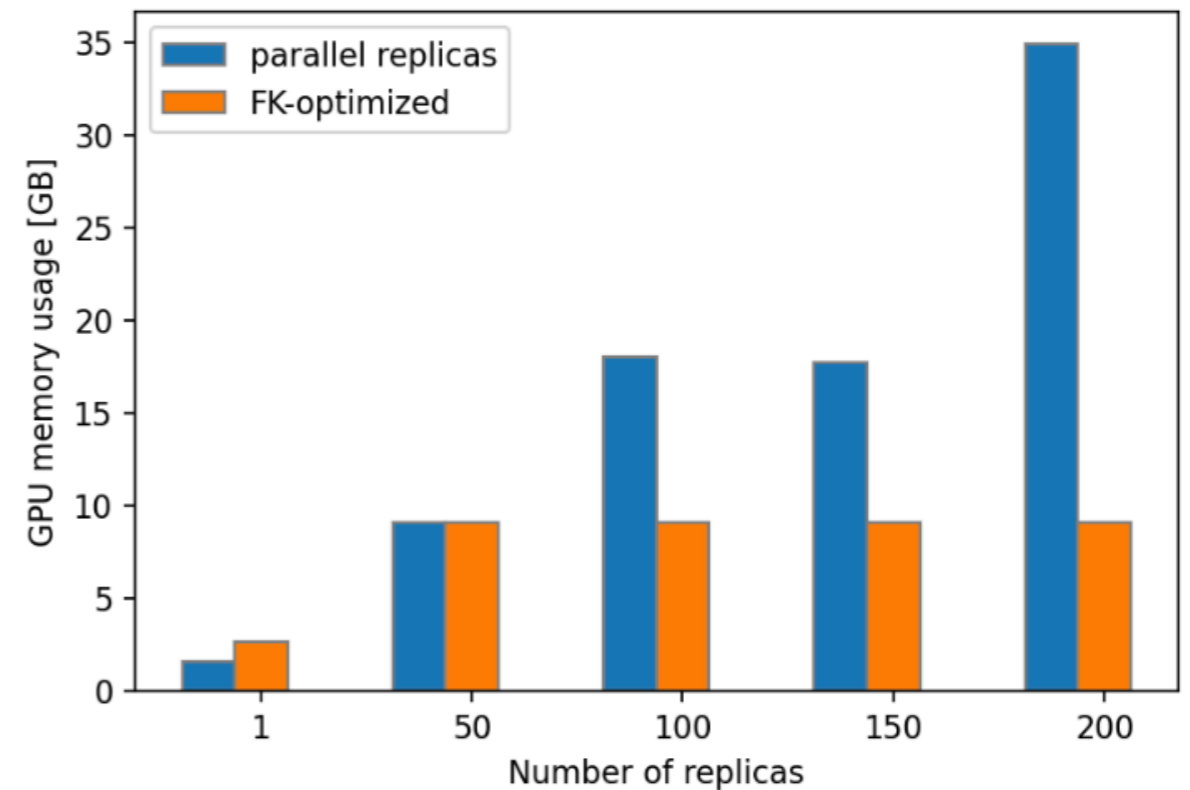
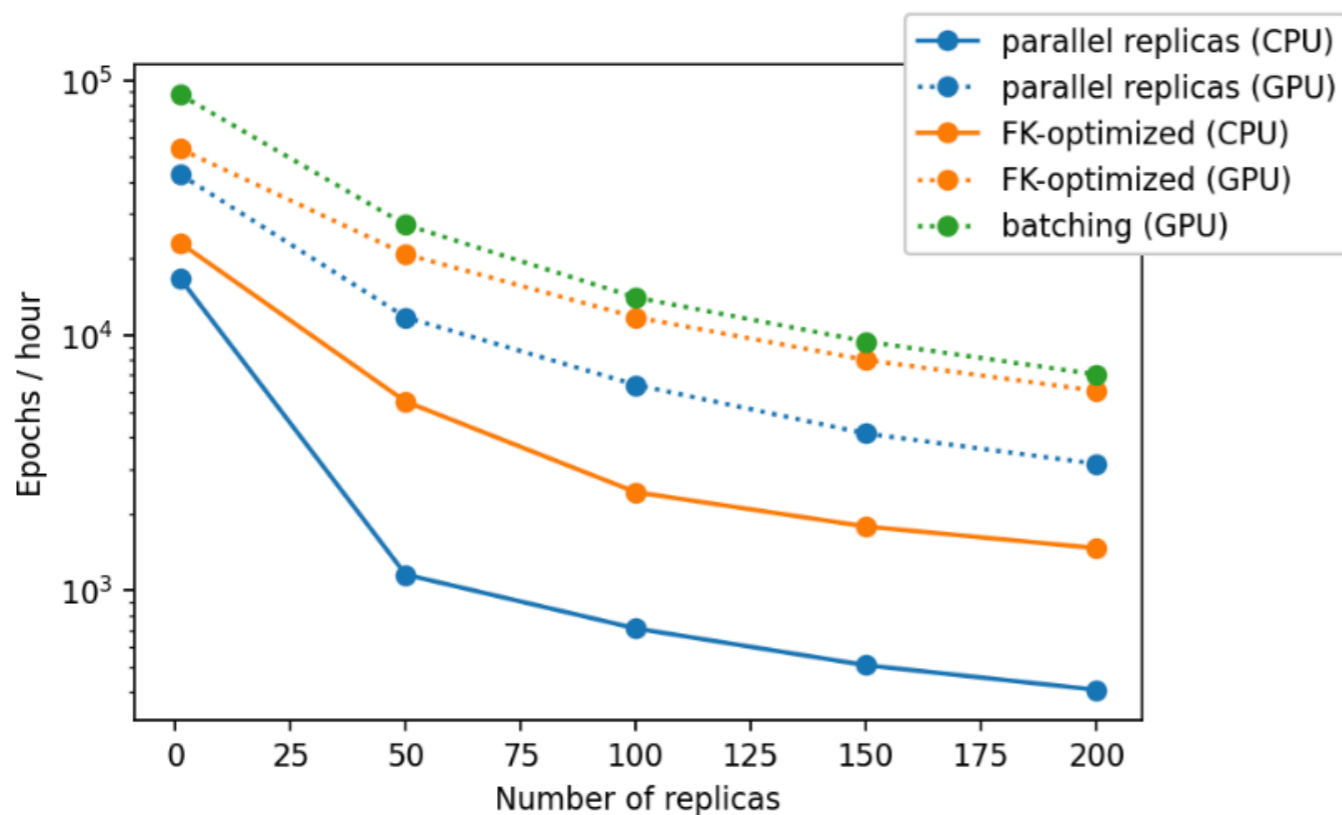
The 3FNS charm PDF displays **non-zero component** peaked at large- $x$  which can be identified with **intrinsic charm**





# GPU & Hyperparameter Optimisation

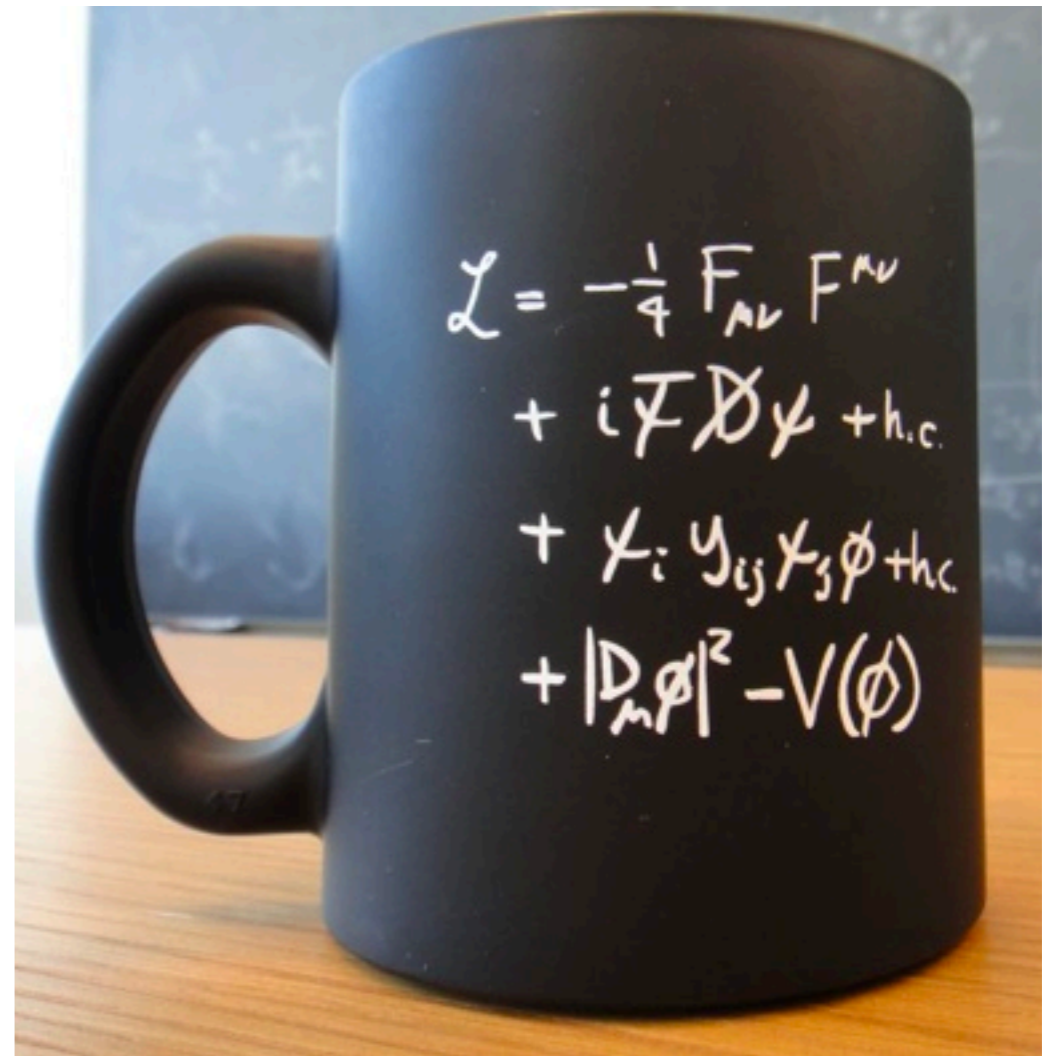
- Deploy NNPDF machinery on GPUs & optimise performance: **factor speed 200 improvement!**
- Also ensure CPU memory consumption keep reasonable
- Develop new strategies for **hyperparameter optimisation** based on the full posterior probability distribution, not only on first moment as most approaches



# The Standard Model as an Effective Theory

The Standard Model EFT is defined by:

- **Particle (matter) content:** quarks and leptons
- **Gauge** (local) symmetries and their eventual breaking mechanisms
- **Lorentz** invariance and other global symmetries
- Linearly realised  $SU(2)_L$  EW symmetry breaking
- **Validity only up to certain energy scale  $\Lambda$**



$$\mathcal{L}_{\text{SMEFT}}(\{c_i\}, \Lambda) = \mathcal{L}_{\text{SM}} + \sum_{d=5}^{\infty} \sum_{i=1}^{N_d} c_i^{(d)} \frac{\mathcal{O}_i^{(d)}}{\Lambda^{d-4}}$$

EFT coupling constants,  
to be determined from **data**

All possible operators of **mass-dimension  $d$**  consistent with  
above requirements

# Statistically optimal observables from ML

Optimal observables depend on **all kinematic variables** and **all EFT coefficients**

$$r_{\sigma}(\mathbf{x}, \mathbf{c}) \equiv \frac{f_{\sigma}(\mathbf{x}, \mathbf{c})}{f_{\sigma}(\mathbf{x}, \mathbf{0})} = 1 + \sum_{j=1}^{n_{\text{eft}}} r_{\sigma}^{(j)}(\mathbf{x}) c_j + \sum_{j=1}^{n_{\text{eft}}} \sum_{k \geq j}^{n_{\text{eft}}} r_{\sigma}^{(j,k)}(\mathbf{x}) c_j c_k$$

parametrised with **neural networks** trained to Monte Carlo simulations & benchmarked with exact calculations

$$\hat{r}_{\sigma}(\mathbf{x}, \mathbf{c}) = 1 + \sum_{j=1}^{n_{\text{eft}}} \text{NN}^{(j)}(\mathbf{x}) c_j + \sum_{j=1}^{n_{\text{eft}}} \sum_{k \geq j}^{n_{\text{eft}}} \text{NN}^{(j,k)}(\mathbf{x}) c_j c_k$$

extendable to **arbitrary number** of kinematic variables and EFT coefficients: training can be parallelised

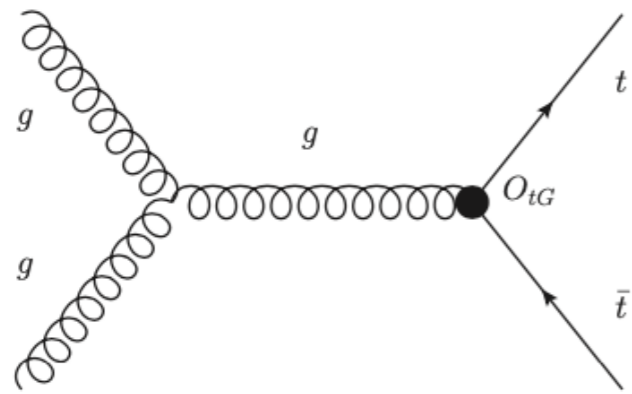
methodological uncertainties (e.g. finite training samples) assess with the **replica method**

$$\hat{r}_{\sigma}^{(i)}(\mathbf{x}, \mathbf{c}) \equiv 1 + \sum_{j=1}^{n_{\text{eft}}} \text{NN}_i^{(j)}(\mathbf{x}) c_j + \sum_{j=1}^{n_{\text{eft}}} \sum_{k \geq j}^{n_{\text{eft}}} \text{NN}_i^{(j,k)}(\mathbf{x}) c_j c_k, \quad i = 1, \dots, N_{\text{rep}}$$

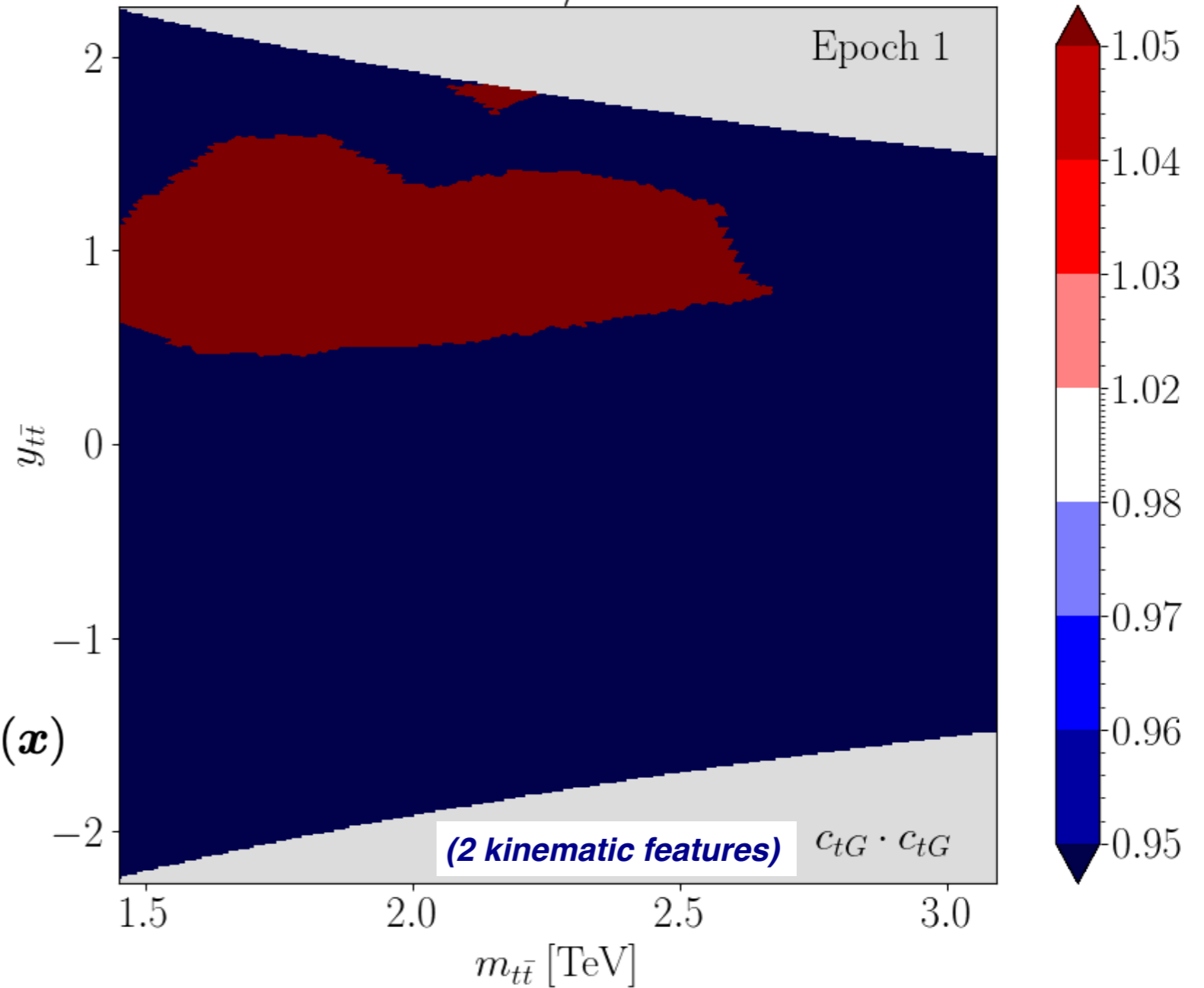
*each replica trained to an independent set of MC events*

*representation of the probability distribution in the space of ML models*

# Neural network training



Unbinned exact/Unbinned ML



$$r_\sigma(\mathbf{x}, c_j^{(\text{tr})}, c_k^{(\text{tr})}) = 1 + c_j^{(\text{tr})} c_k^{(\text{tr})} \text{NN}^{(j,k)}(\mathbf{x})$$

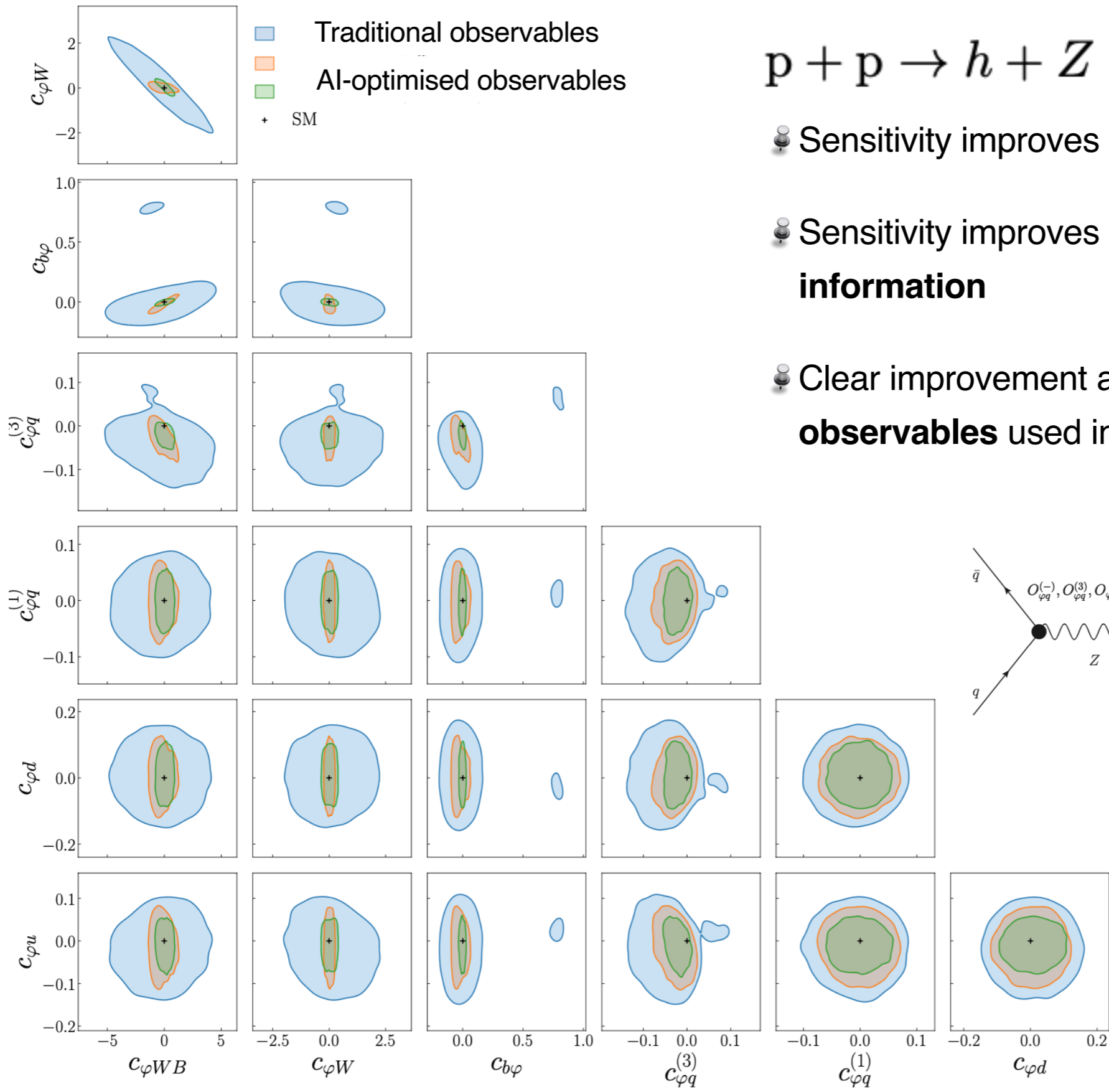
$$\mathbf{x} = (m_{t\bar{t}}, y_{t\bar{t}})$$

**NN training by minimising cross-entropy loss function**

$$L[g(\mathbf{x}, \mathbf{c})] = -\sigma_{\text{fid}}(\mathbf{c}) \sum_{i=1}^{N_{\text{ev}}} \log(1 - g(\mathbf{x}_i, \mathbf{c})) - \sigma_{\text{fid}}(\mathbf{0}) \sum_{j=1}^{N_{\text{ev}}} \log g(\mathbf{x}_j, \mathbf{c}) \quad g = (1 + r_\sigma)^{-1}$$

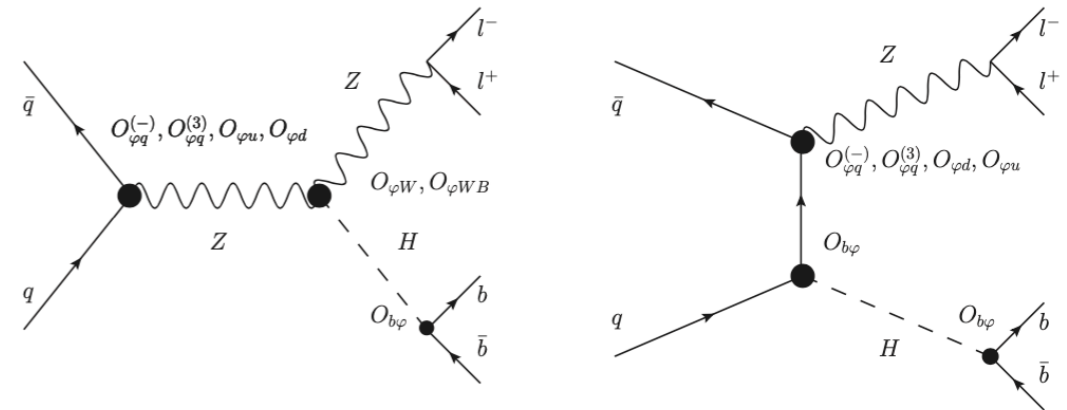
# Results: Higgs+Z production

Marginalised 95 % C.L. intervals,  $\mathcal{O}(\Lambda^{-4})$  at  $\mathcal{L} = 300 \text{ fb}^{-1}$

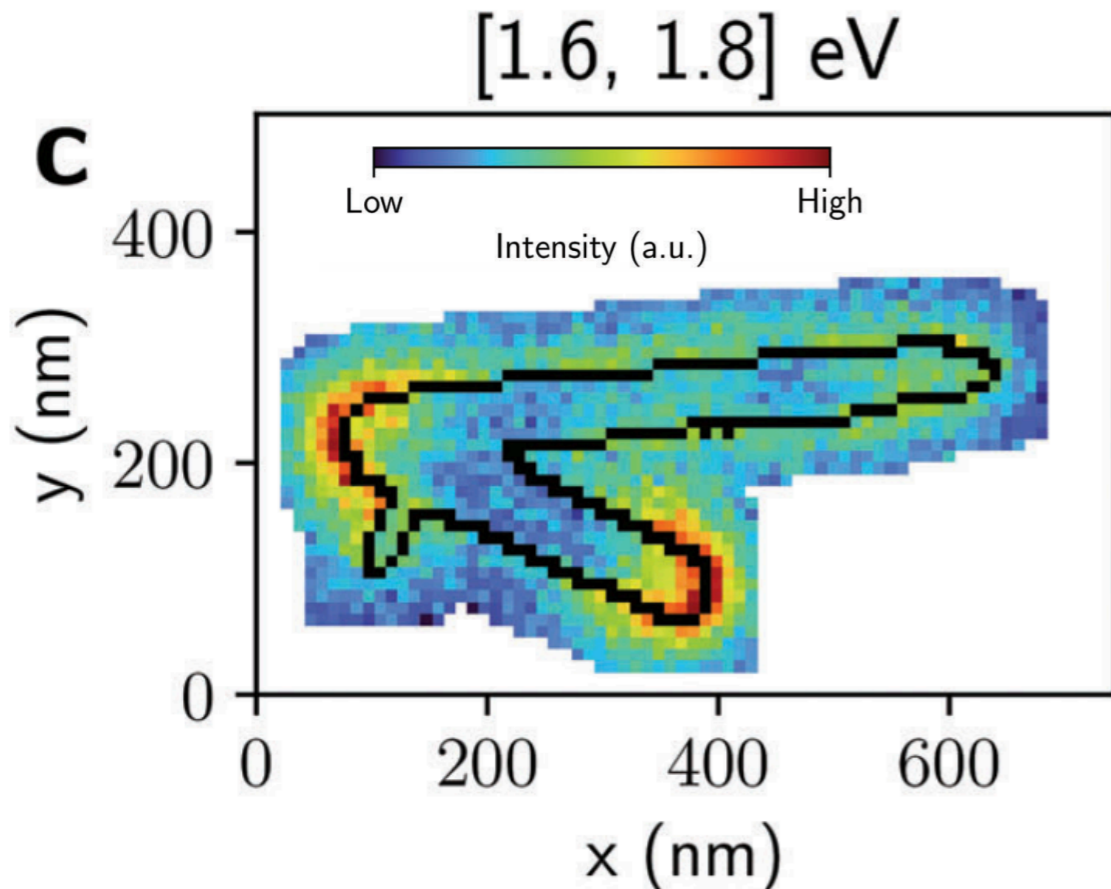
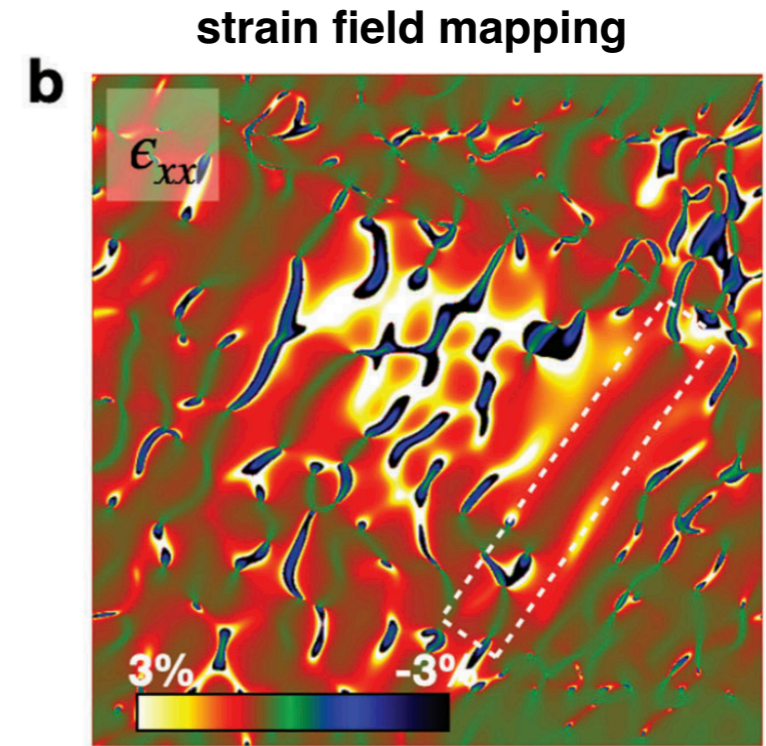
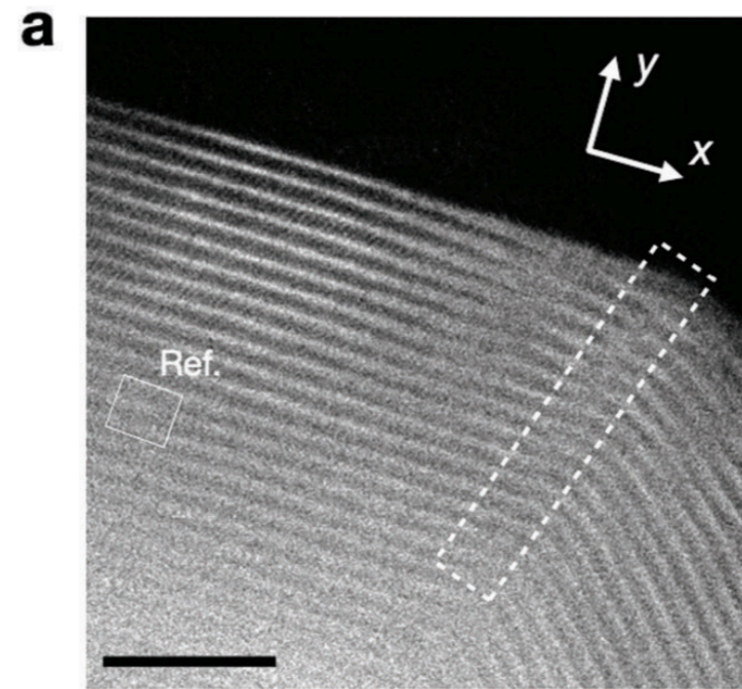
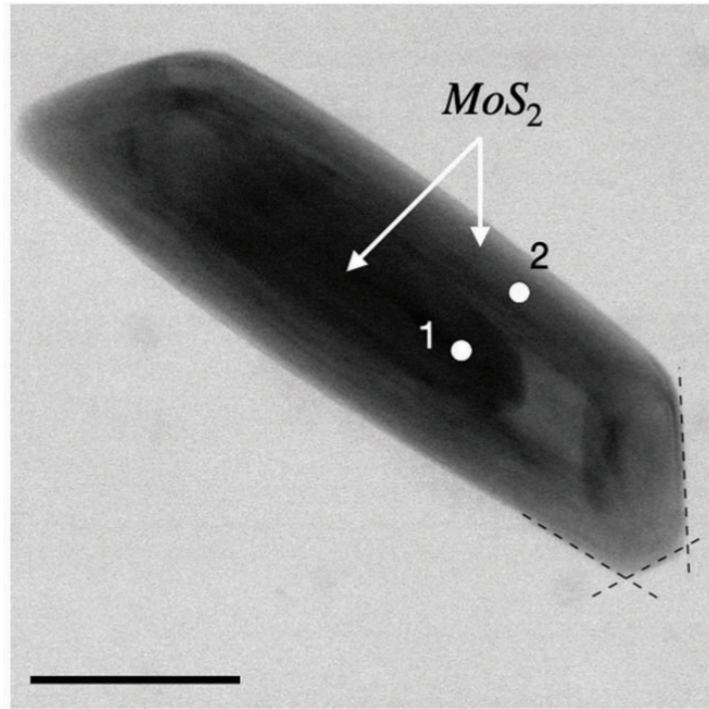


$$p + p \rightarrow h + Z \rightarrow b + \bar{b} + \ell^+ + \ell^-$$

- 📌 Sensitivity improves in **unbinned analysis**
- 📌 Sensitivity improves when **using all kinematic information**
- 📌 Clear improvement as compared to **traditional observables** used in EFT fits



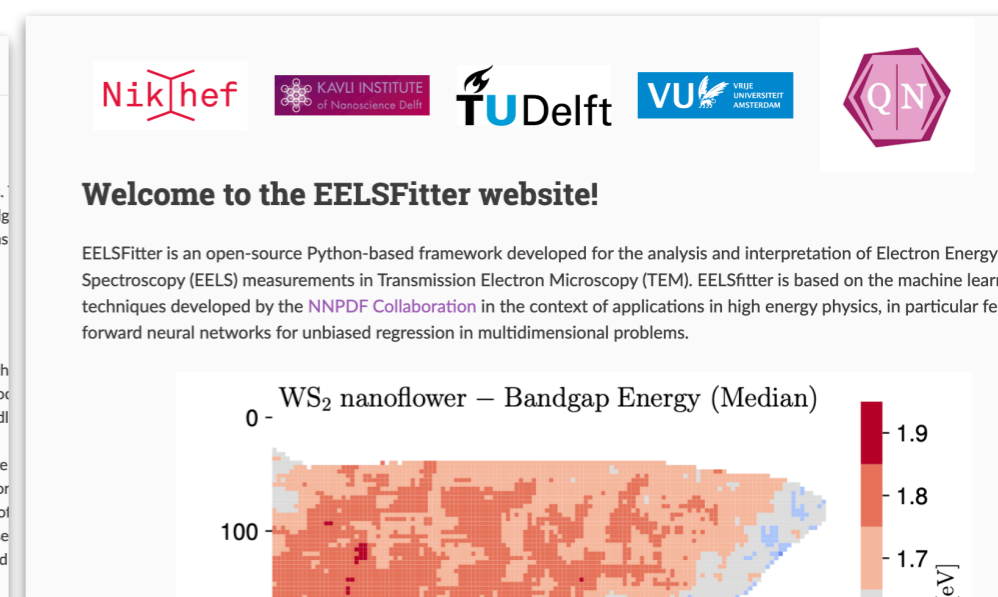
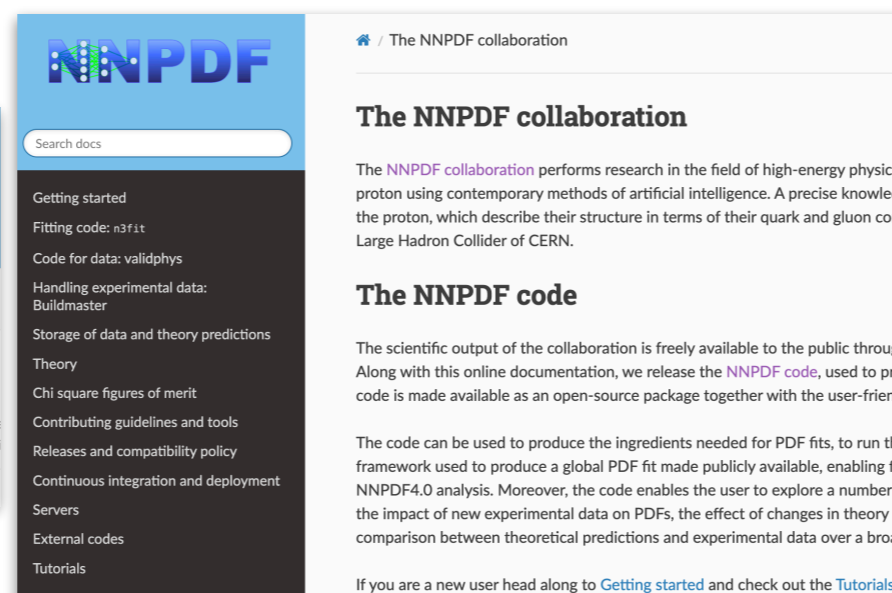
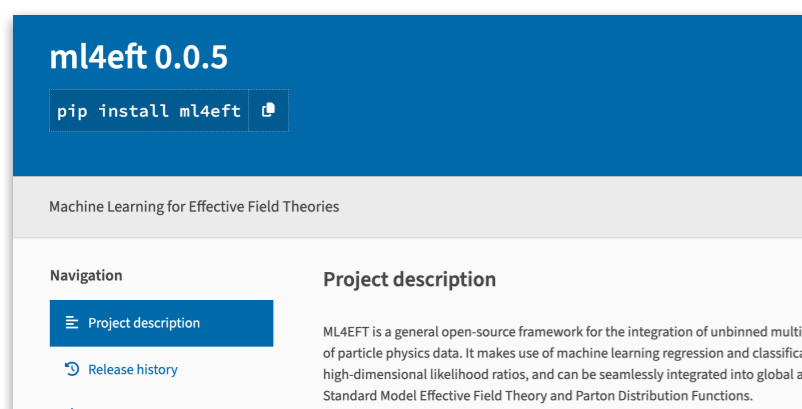
# Applications to Material Science



- ☑ Background subtraction & anomaly detection problems from HEP arise also in material spectroscopy
- ☑ Direct correlation of **strain fields, band gap modulation, and exciton localisation** in 1D- $\text{MoS}_2$  nanostructures with different morphologies
- ☑ Developed EELSfitter ML framework together with TU Delft researchers

# Summary and outlook

- 👤 Machine learning makes possible **identifying patterns in the data** whereby one can efficiently solve problems which are difficult or intractable with traditional approaches
- 👤 **Enable discoveries** such as intrinsic charm quarks in the proton & make possible to **optimise the sensitivity of searches** for interesting phenomena hidden in the data
- 👤 Our technology is portable to many other problems, as demonstrated for their applicability to data analysis in **electron microscopy of quantum materials**
- 👤 Codes are **open source** and extensively documented, and have benefitted from contributions as well from BSc and MSc students in our groups



# Summary and outlook

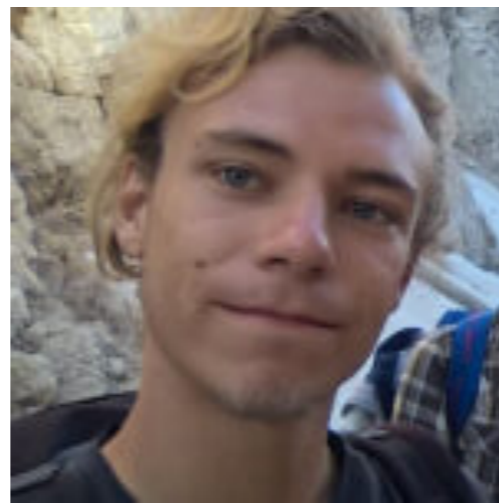
- 🎤 Machine learning makes possible **identifying patterns in the data** whereby one can efficiently solve problems which are difficult or intractable with traditional approaches
- 🎤 **Enable discoveries** such as intrinsic charm quarks in the proton & make possible to **optimise the sensitivity of searches** for interesting phenomena hidden in the data
- 🎤 Our technology is portable to many other problems, as demonstrated for their applicability to data analysis in **electron microscopy of quantum materials**
- 🎤 Codes are **open source** and extensively documented, and have benefitted from contributions as well from BSc and MSc students in our groups



Peter



Jaco



Giacomo



Tanjona



Tommaso