

Statistical Modelling

Francesca Capel

Max Planck Institute for Physics

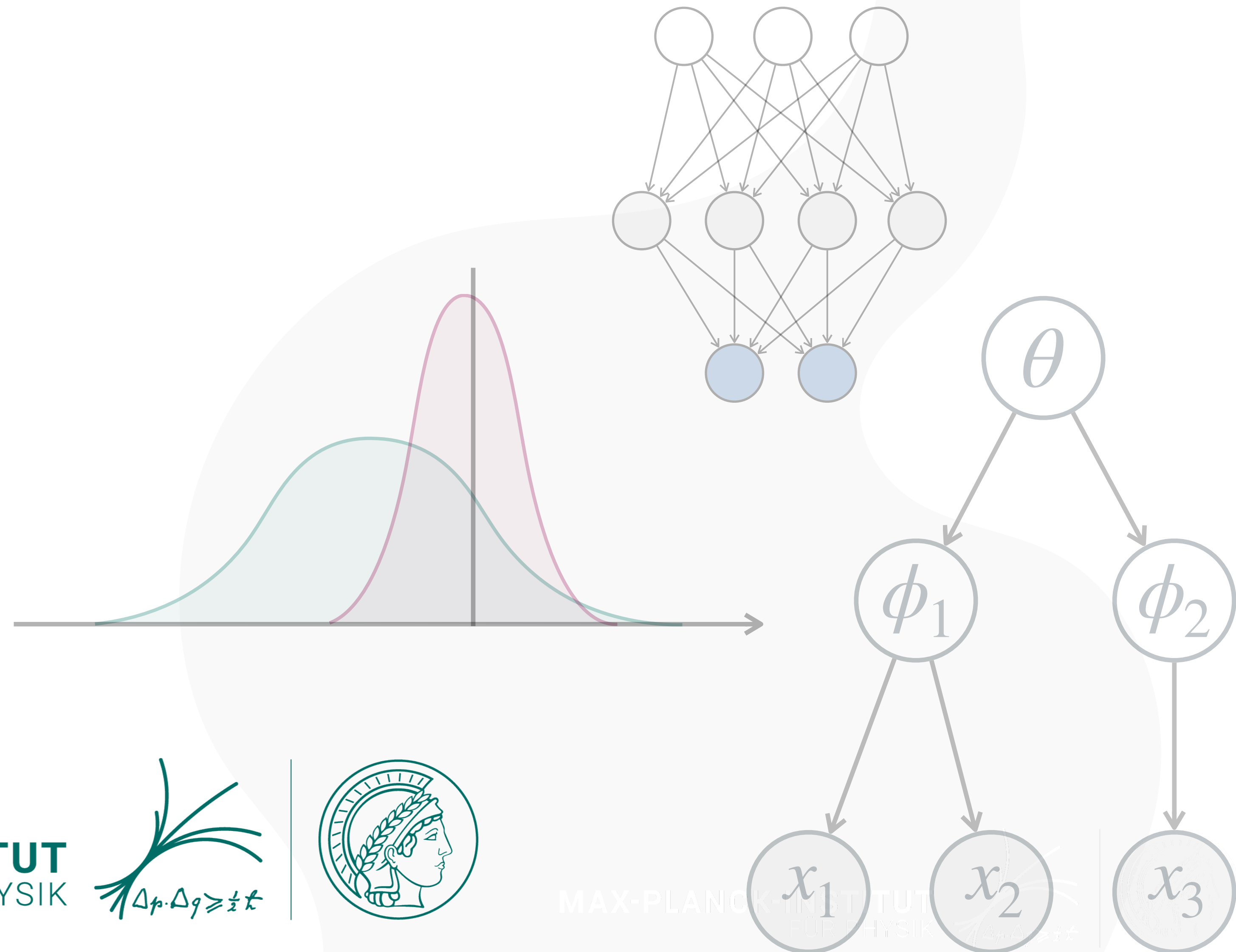
29th March 2025 - TAML Workshop @ ICRR



MAX-PLANCK-INSTITUT
FÜR PHYSIK

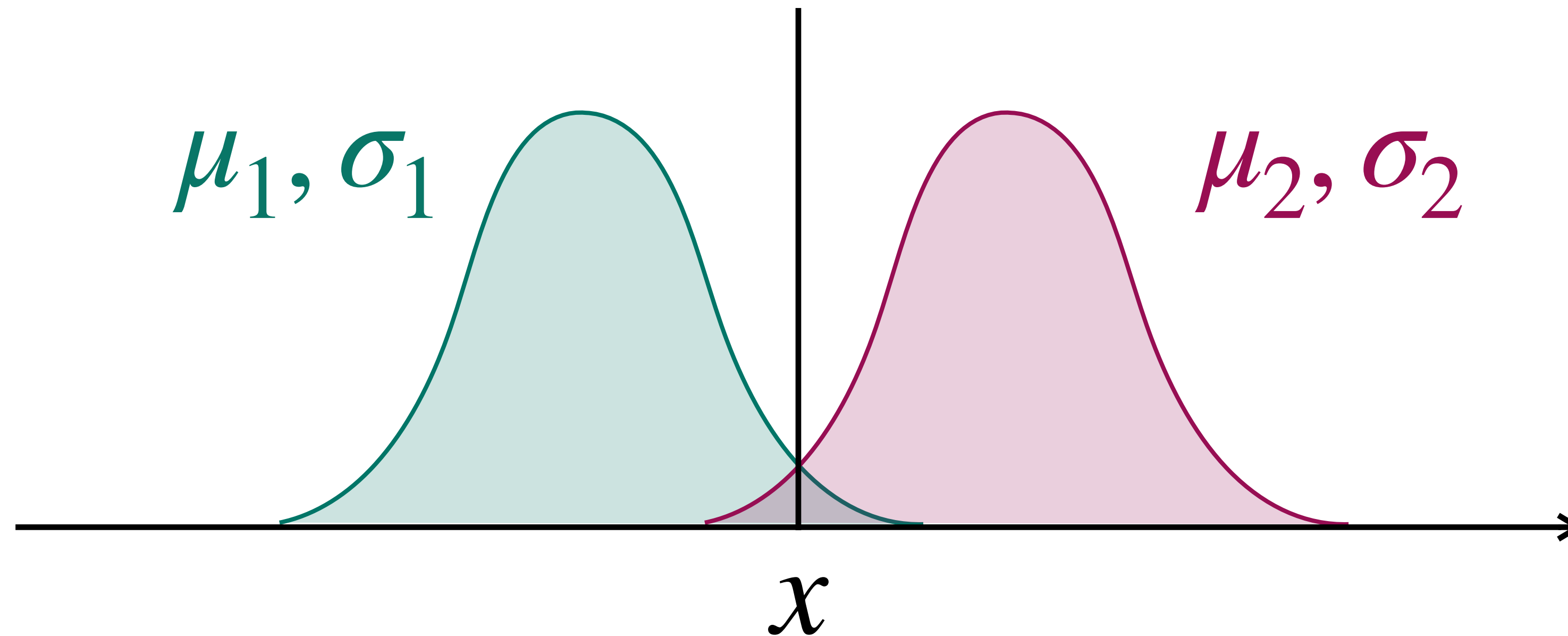


MAX-PLANCK-INSTITUT
FÜR PHYSIK



Statistics and machine learning

A textbook statistics example:



$$\mathcal{L}(x | \mu, \sigma) = w \mathcal{N}(x | \mu_1, \sigma_1) + (1 - w) \mathcal{N}(x | \mu_2, \sigma_2)$$

Statistics and machine learning

A realistic ML example:



$$\mathcal{L}(x | \text{muffin, dog}) = w p(x | \text{muffin}) + (1 - w) p(x | \text{dog})$$

Statistics and machine learning

We can think about neural networks as a bunch of weights that figure out how to describe the likelihood

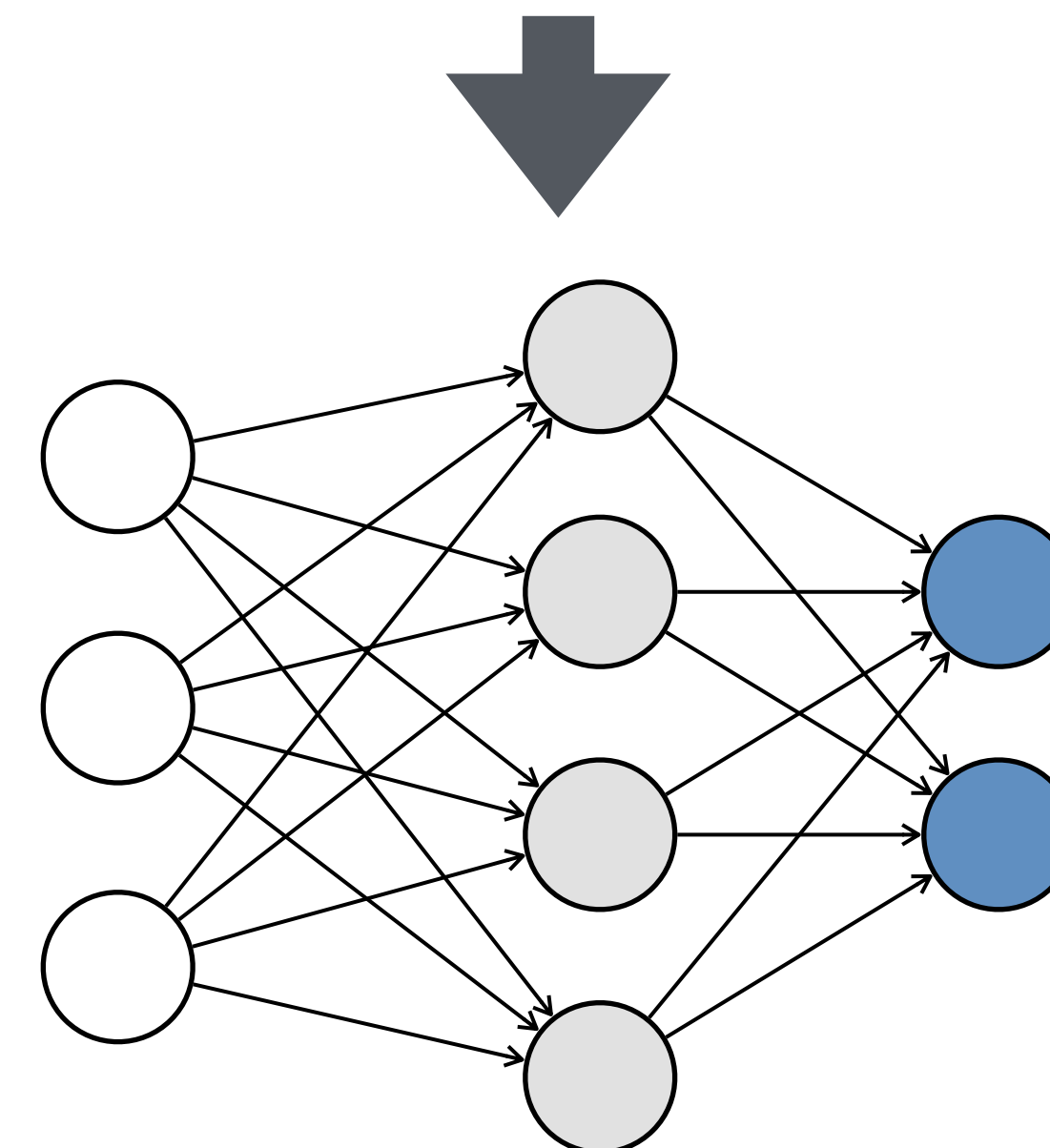
Training: Figure out the likelihood (connection between data and model)

Evaluating: Maximum likelihood solution

ML in physics: We care about interpretation and uncertainty → No escape from statistics



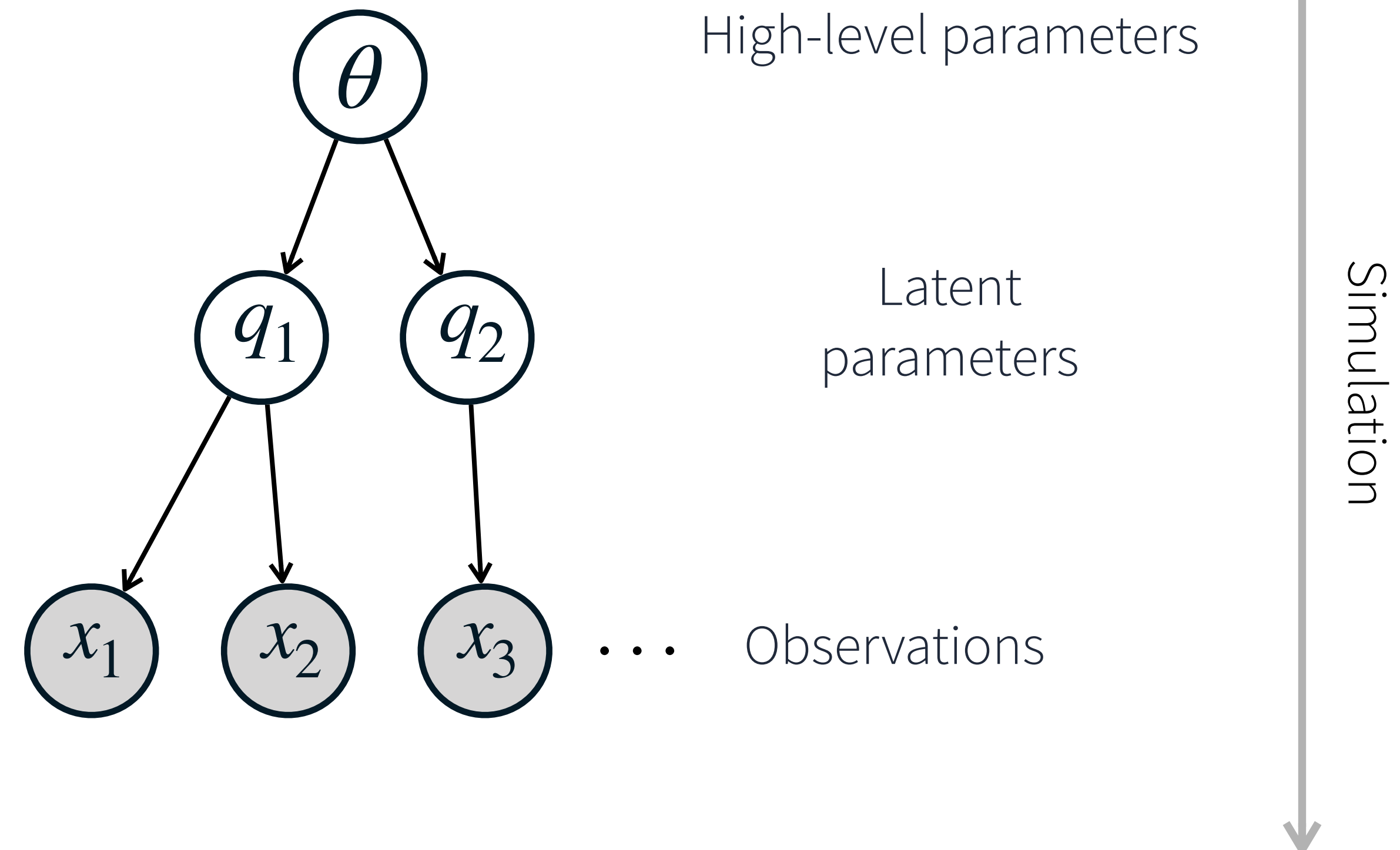
$$\mathcal{L}(x | \text{muffin}, \text{dog}) = w p(x | \text{muffin}) + (1 - w) p(x | \text{dog})$$



Forward modelling

In astroparticle physics we have “models” for our data, typically in the form of Monte Carlo simulations

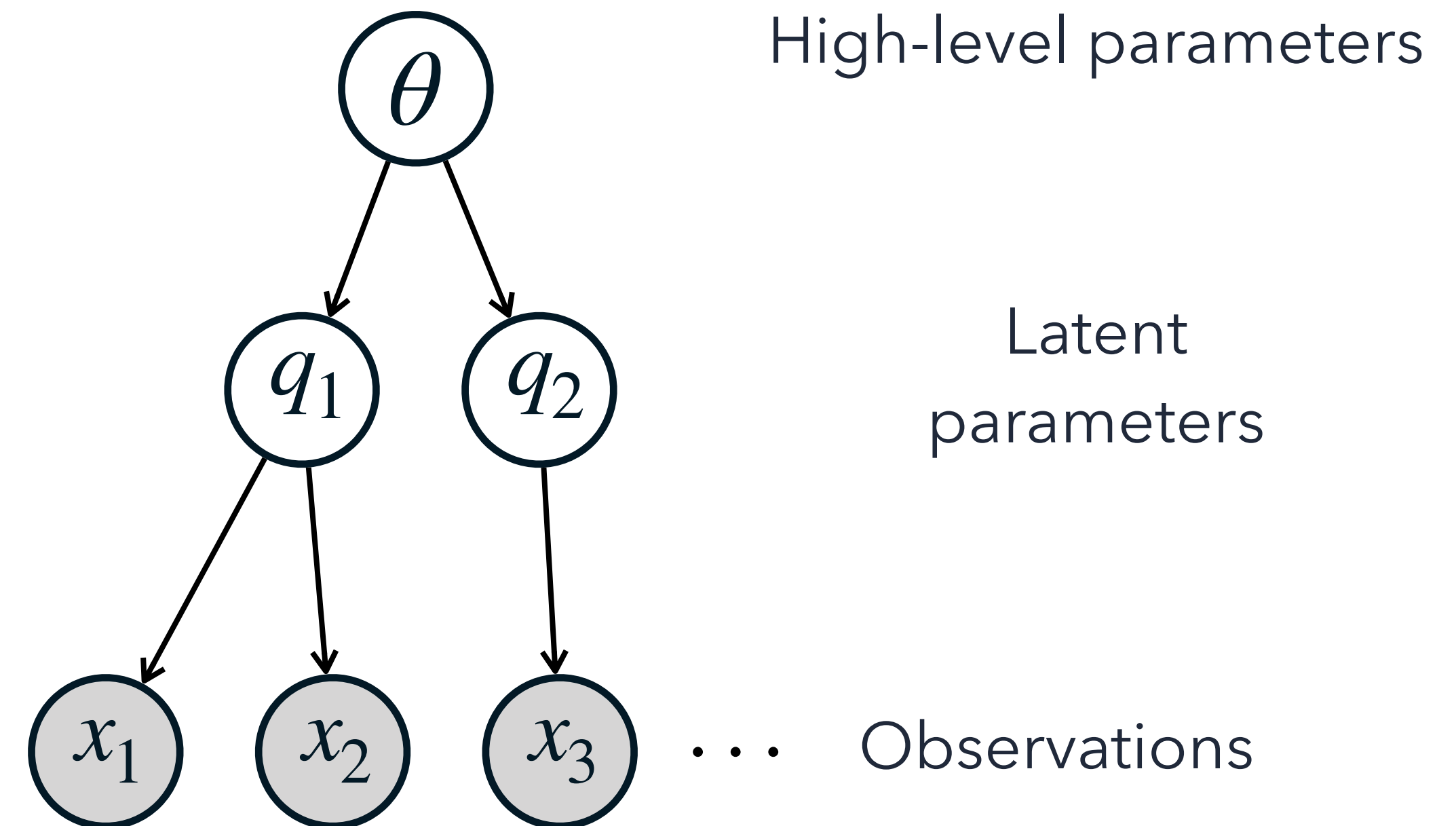
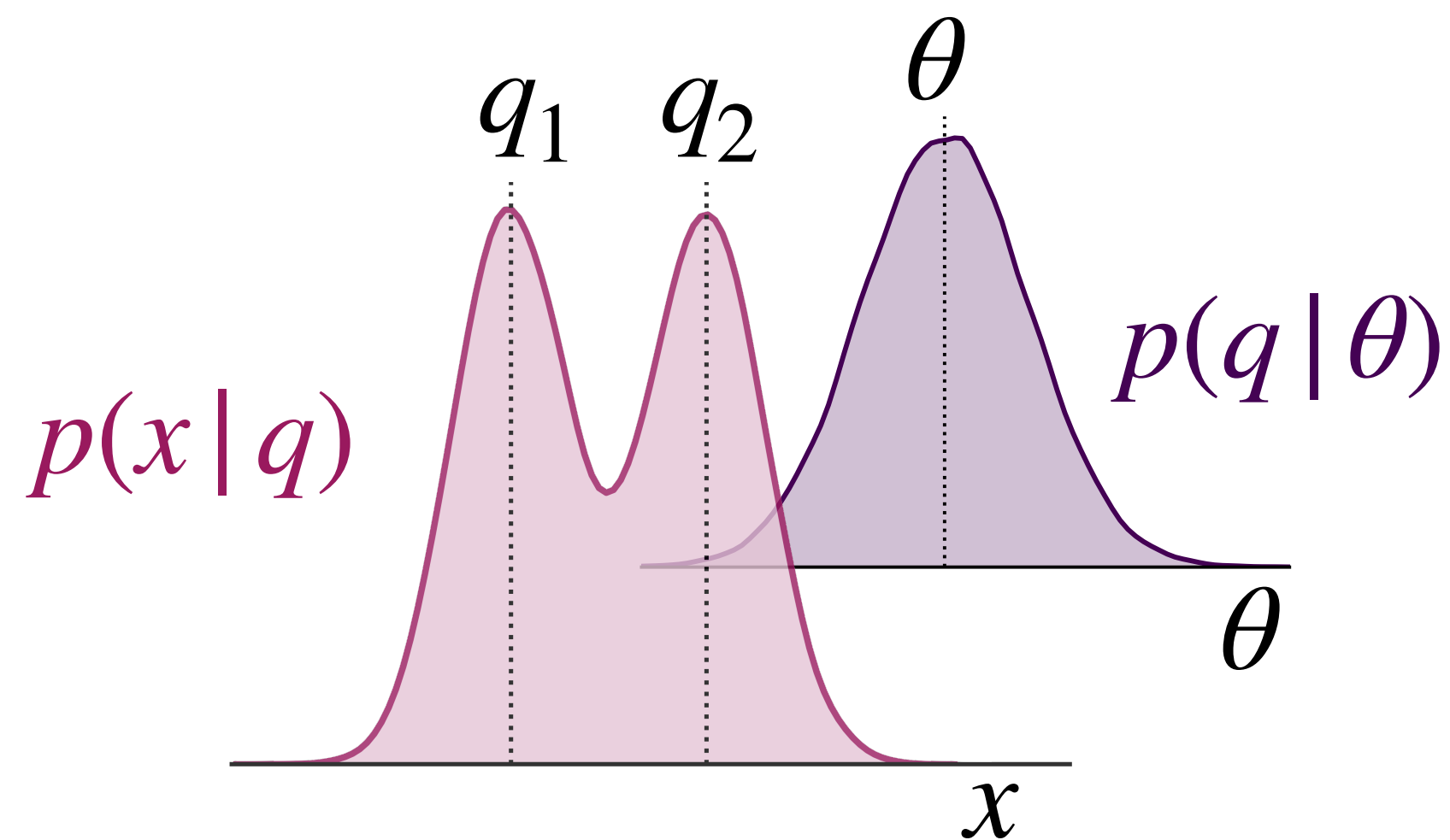
- Quantify expected **uncertainties**
- Map out **complexity**
- Contains all information needed to derive the **likelihood function**
- Testbench to **evaluate possible simplifications** or assumptions



Forward model = likelihood

Organise the free parameters into a statistical model that describes the **data generating process**
Lots of free parameters \rightarrow machine learning is useful (but not only way)

$$\mathcal{L}(x, \theta) = p(x | q) p(q | \theta) \dots$$



Inference in many dimensions

Markov chain Monte Carlo

Numerically approximate **high-dimensional integrals**
(e.g. expectation values, variances of parameters)

Demo: <https://chi-feng.github.io/mcmc-demo/app.html>

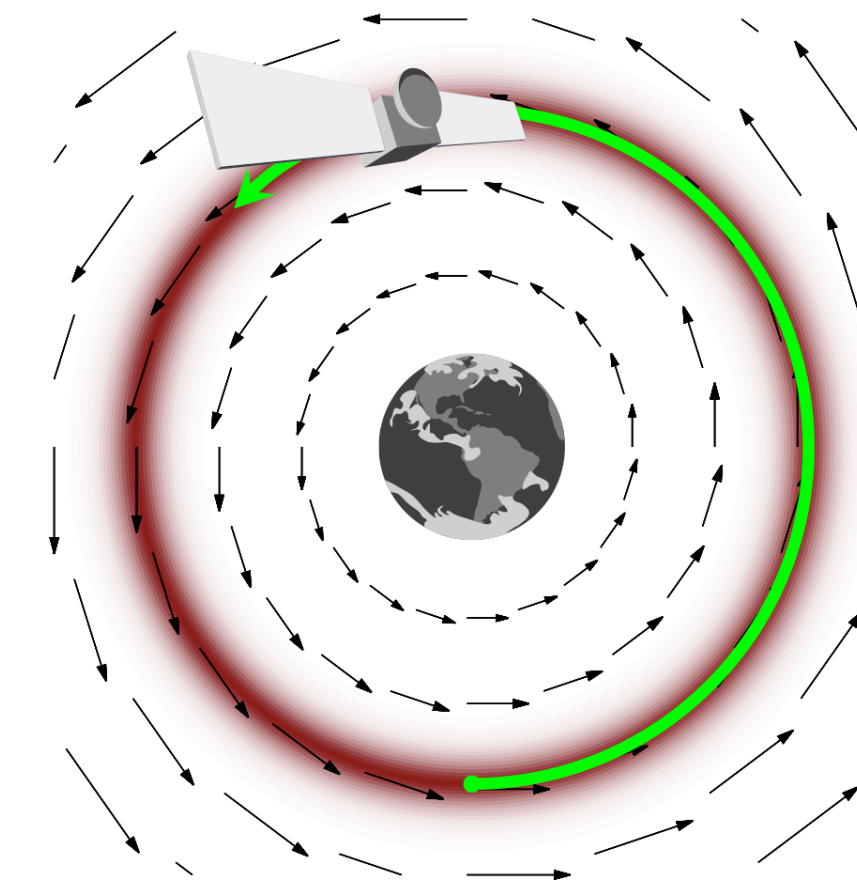
Hamiltonian Monte Carlo

Markov chain Monte Carlo that uses Hamiltonian dynamics to move efficiently through high-dimensional parameter spaces

$$\theta \longrightarrow (\theta, p) \quad \frac{d\theta}{dt} = \frac{\partial H}{\partial p}$$

$$H(\theta, p) \equiv \log P(\theta, p) \quad \frac{dp}{dt} = -\frac{\partial H}{\partial \theta}$$

$$\int_{\Theta} d\theta f(\theta)p(\theta) = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=0}^N f(\theta_n)$$



[Betancourt (2014)]

Large numbers
of free
parameters
possible

Uncertainty
quantification for
free

Bayesian Inference

POSTERIOR

LIKELIHOOD

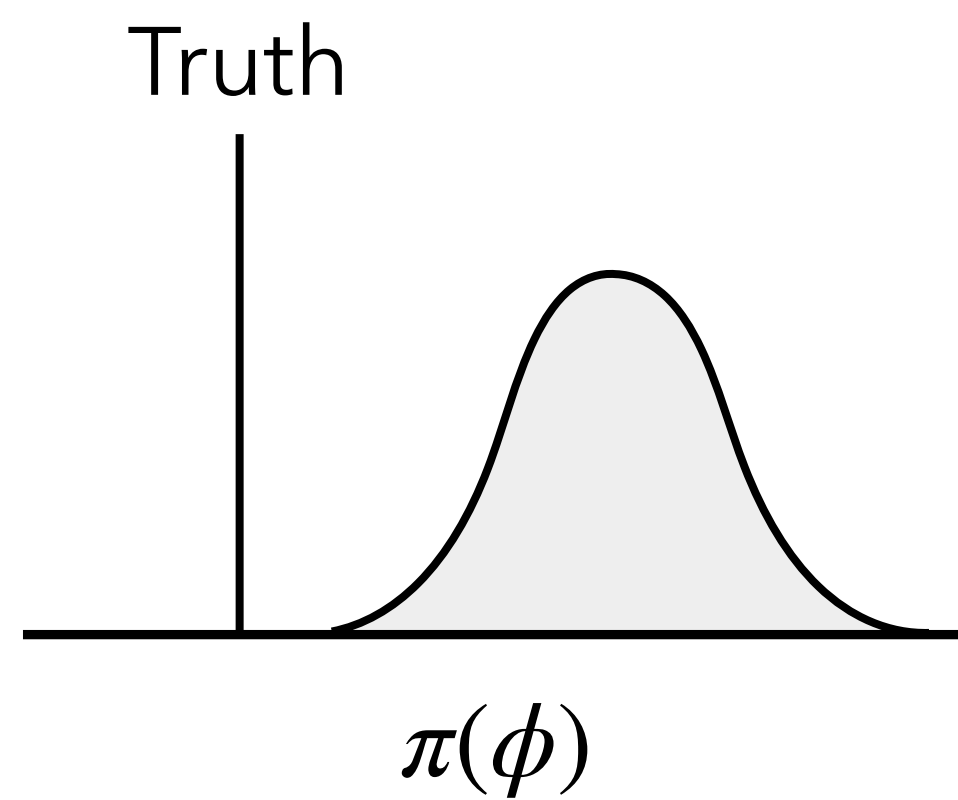
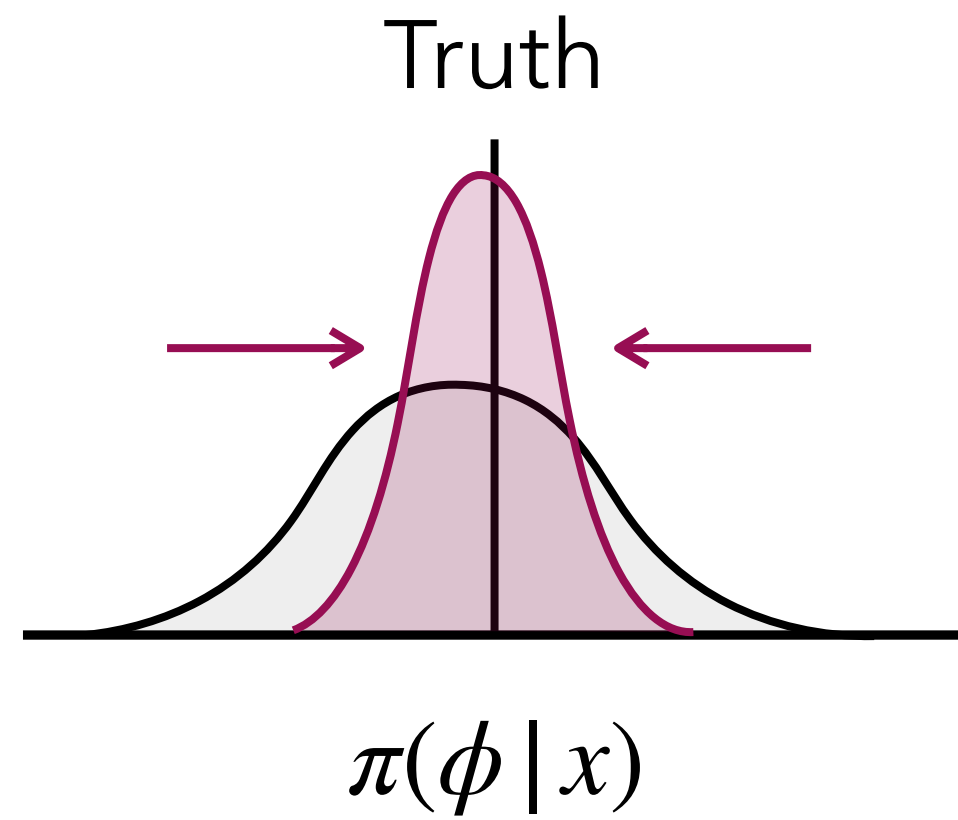
PRIOR

MARGINAL

$$p(\theta|x) = \frac{p(x|\theta)p(\theta)}{p(x)}$$

Systematic uncertainties

[Sinervo PHYSTAT 2003; Heinrich+Lyons Annu. Rev. Nucl. Part. Sci. 2007]



1

Systematics that can be constrained

2

Uncertain assumptions in measurement/analysis

3

Uncertain theoretical framework

Calibration

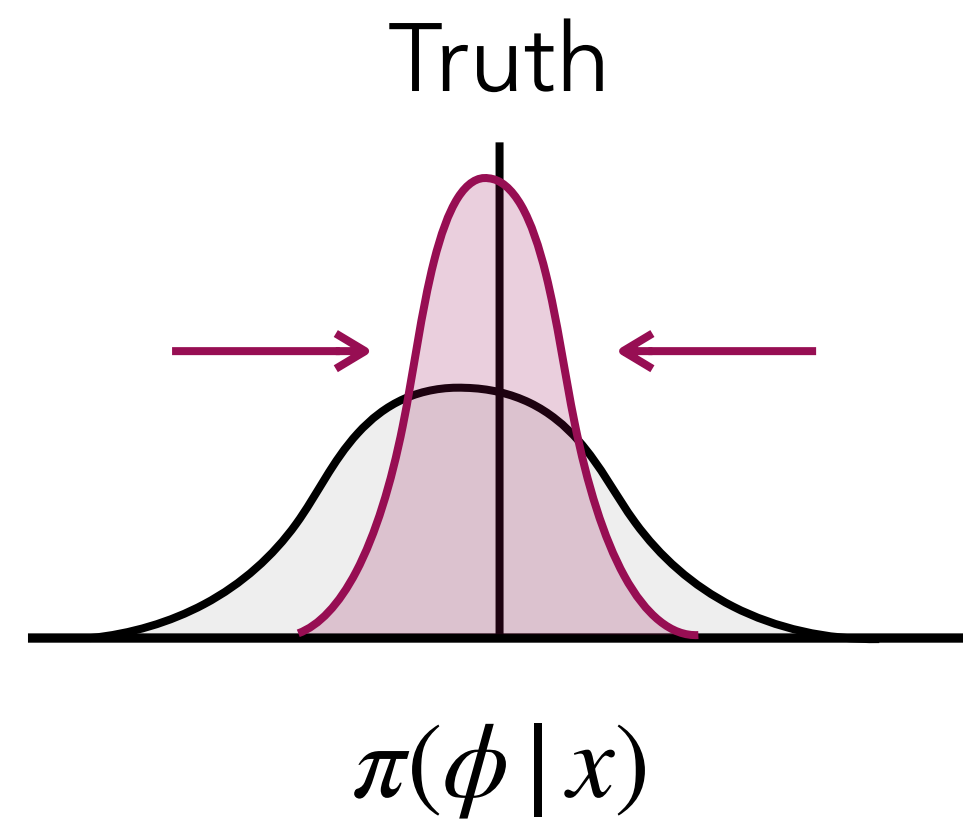
Data

Model misspecification

Theory/
phenomenology

Systematic uncertainties

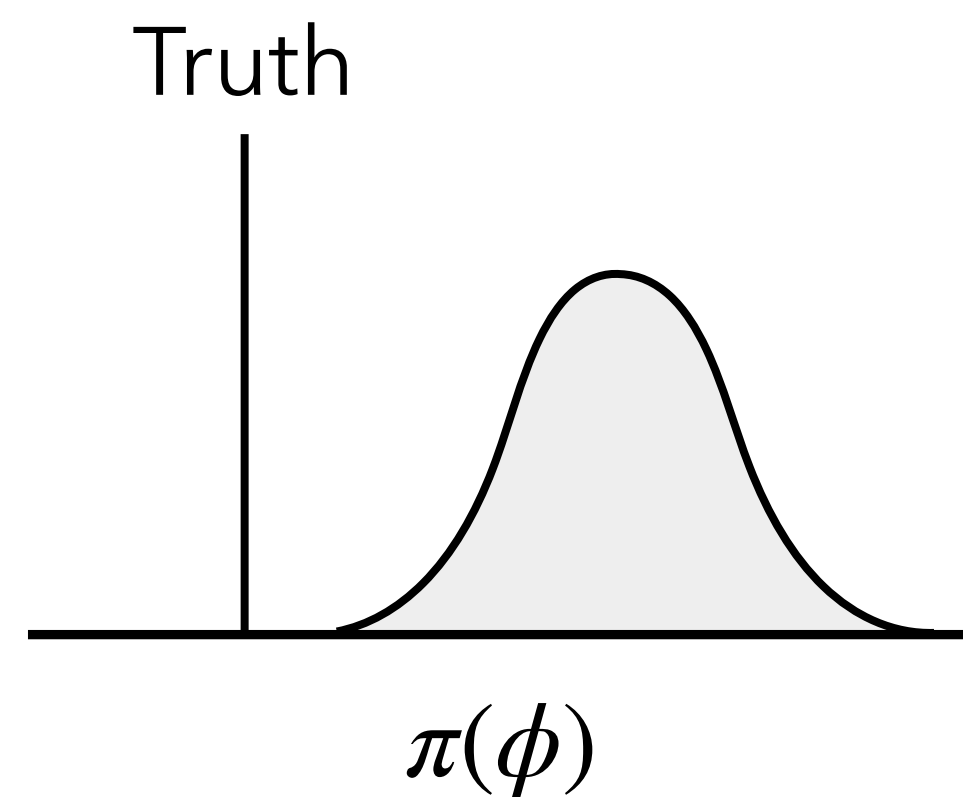
[Sinervo PHYSTAT 2003; Heinrich+Lyons Annu. Rev. Nucl. Part. Sci. 2007]



1

Systematics that can be constrained

GOOD 😊



2

Uncertain assumptions in measurement/analysis

BAD 😞

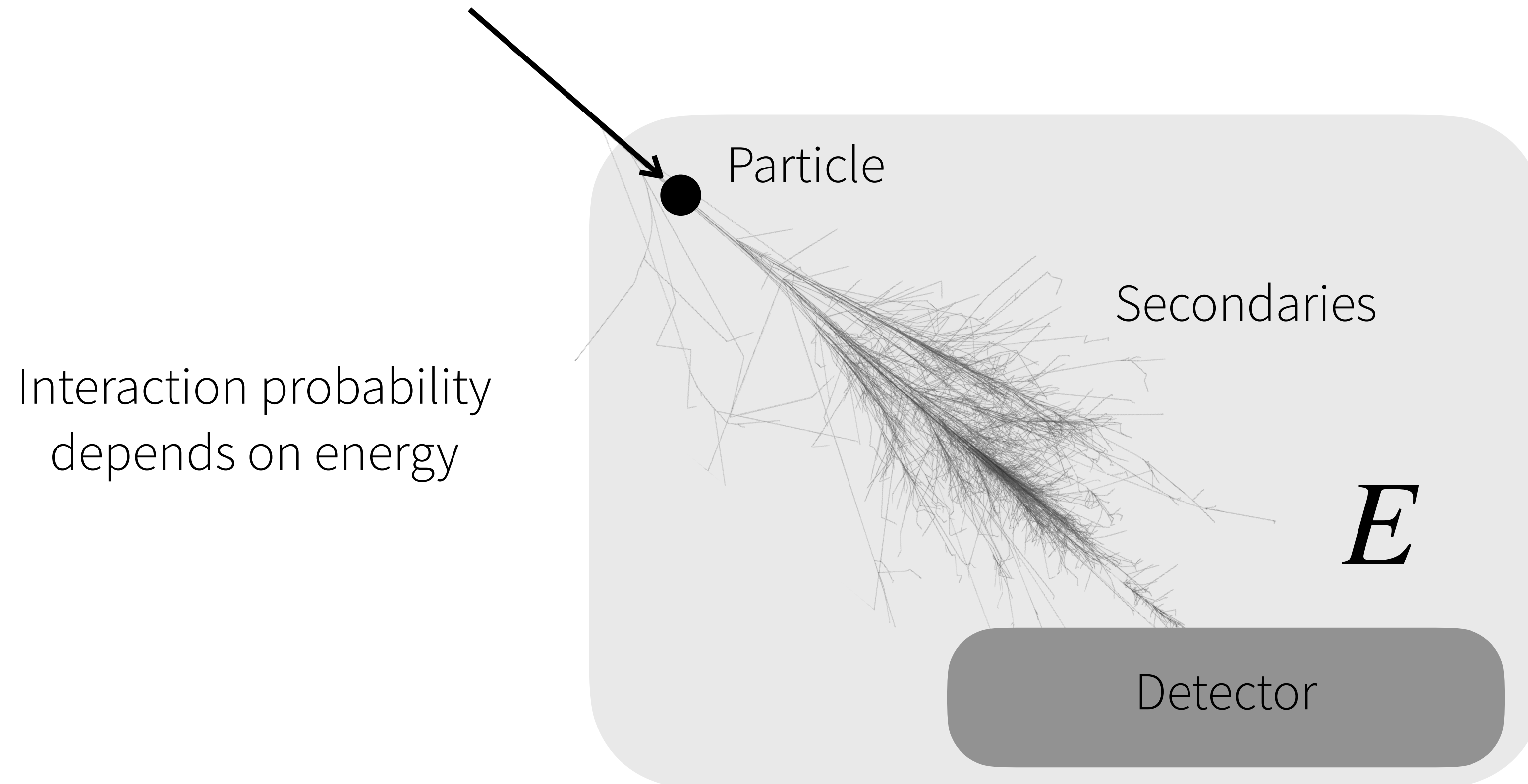
3

Uncertain theoretical framework

UGLY 😱

Today's session

Example: A toy particle detector that measures energy indirectly



Goal:
Measure particle
spectrum

Today's session

- Forward folding
- Goodness-of-fit
- Systematic uncertainties