

Overcoming limitations to parameter inference using *Neural Ratio Estimation*

Gert Kluge

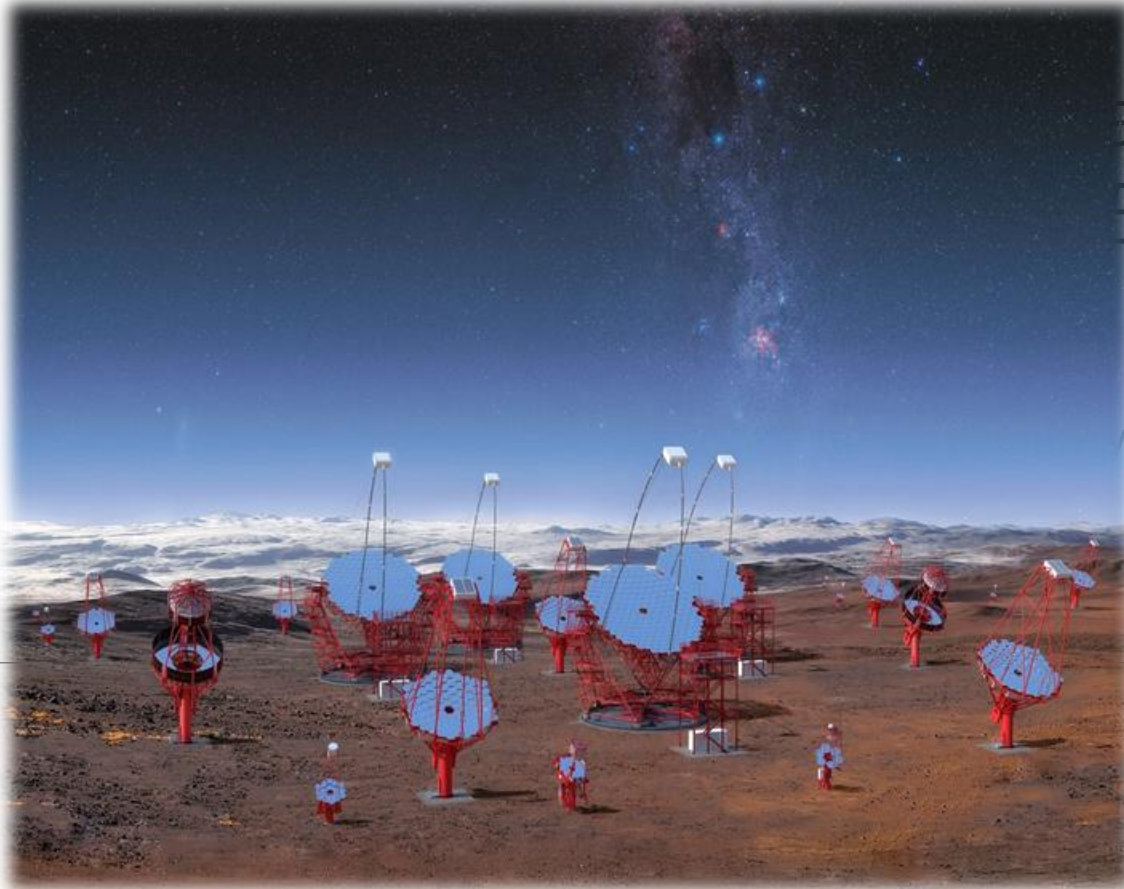
PhD candidate at the University of Oslo

In collaboration with

Giacomo D'Amico, Julia Djuvsland, and Heidi Sandaker

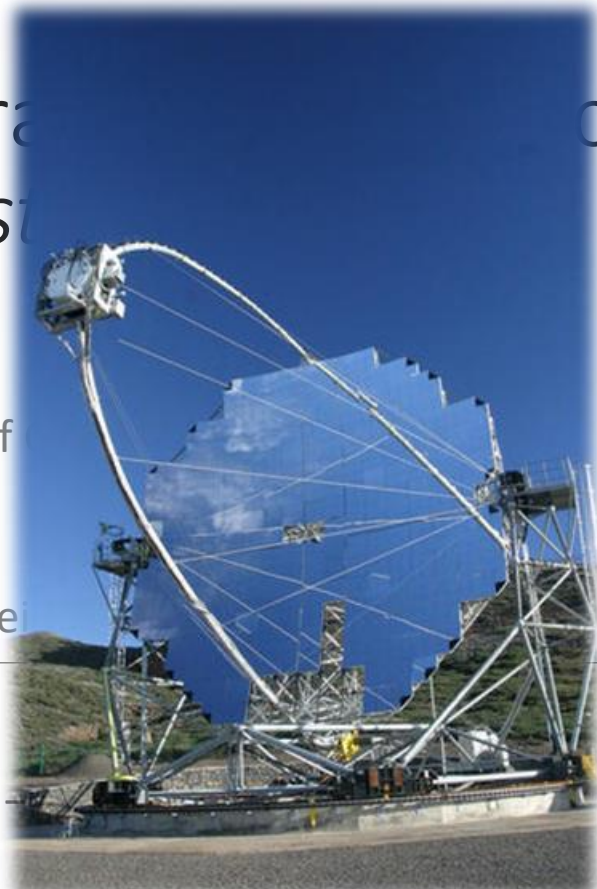
Nordic Conference on Particle Physics – Oppland
7th of January 2023

Cherenkov Telescope Array



Credit: ESO

MAGIC



Credit: Robert Wagner, [CC BY-SA 2.0](https://creativecommons.org/licenses/by-sa/2.0/)

Overcoming limitations to parameter inference using *Neural Ratio Estimation*

Gert Kluge

PhD candidate at the University of Oslo

In collaboration with

Giacomo D'Amico, Julia Djuvslund, and Heidi Sandaker

Nordic Conference on Particle Physics – Oppland
7th of January 2023

Neural Ratio Estimation: A magic trick to do parameter inference

Neural Ratio Estimation: A magic trick to do parameter inference

- Some advantages:
 - General applicability

Neural Ratio Estimation: A magic trick to do parameter inference

- Some advantages:
 - General applicability
 - Assuming that experimental outcomes can be simulated efficiently (*simulation-based inference*)

Neural Ratio Estimation: A magic trick to do parameter inference

- Some advantages:
 - General applicability
 - Assuming that experimental outcomes can be simulated efficiently (*simulation-based inference*)
 - Avoids over-simplifications
(Takes into account the uncertainty on most, or all, nuisance parameters)

Neural Ratio Estimation: A magic trick to do parameter inference

- Some advantages:
 - General applicability
 - Assuming that experimental outcomes can be simulated efficiently (*simulation-based inference*)
 - Avoids over-simplifications
(Takes into account the uncertainty on most, or all, nuisance parameters)
 - Faster than MCMC

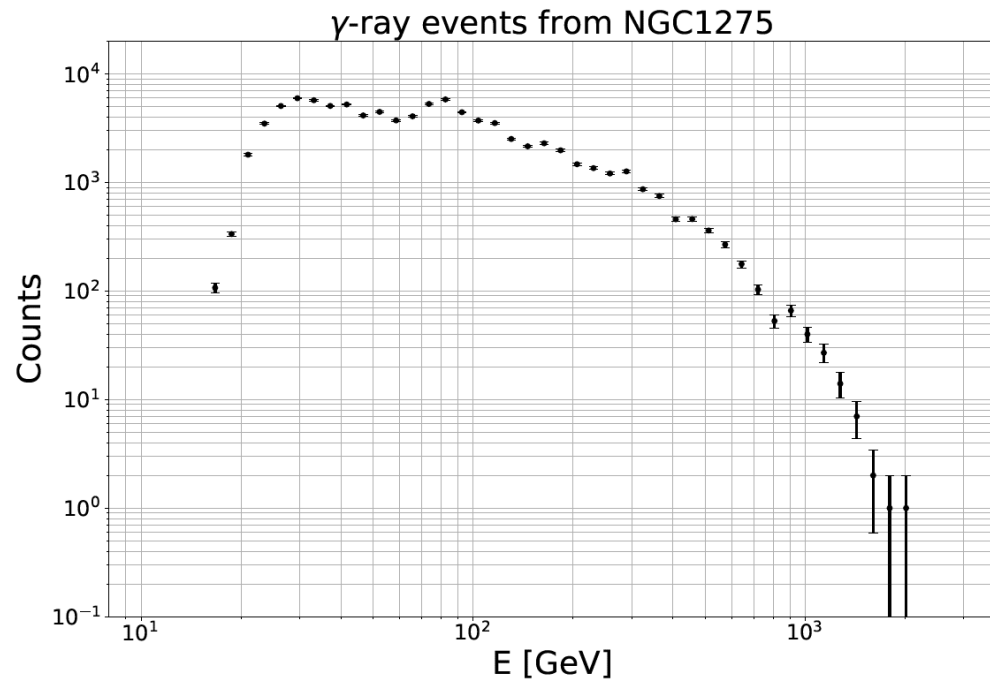
Neural Ratio Estimation: A magic trick to do parameter inference

- Some advantages:
 - General applicability
 - Assuming that experimental outcomes can be simulated efficiently (*simulation-based inference*)
 - Avoids over-simplifications
(Takes into account the uncertainty on most, or all, nuisance parameters)
 - Faster than MCMC



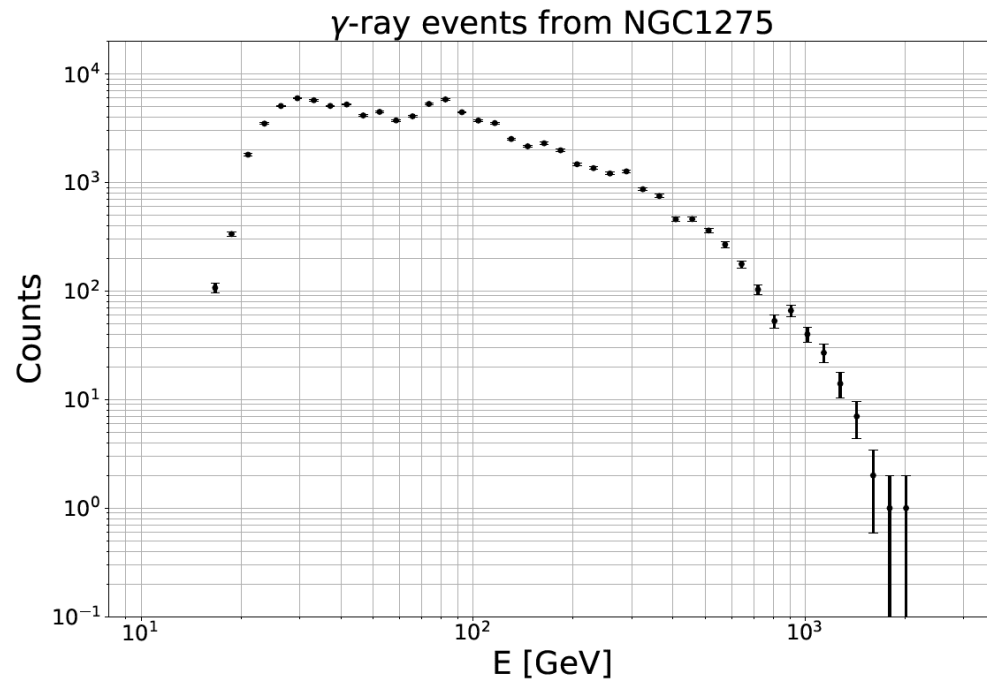
Working example for this talk: fitting the spectrum of NGC1275

Working example for this talk: fitting the spectrum of NGC1275



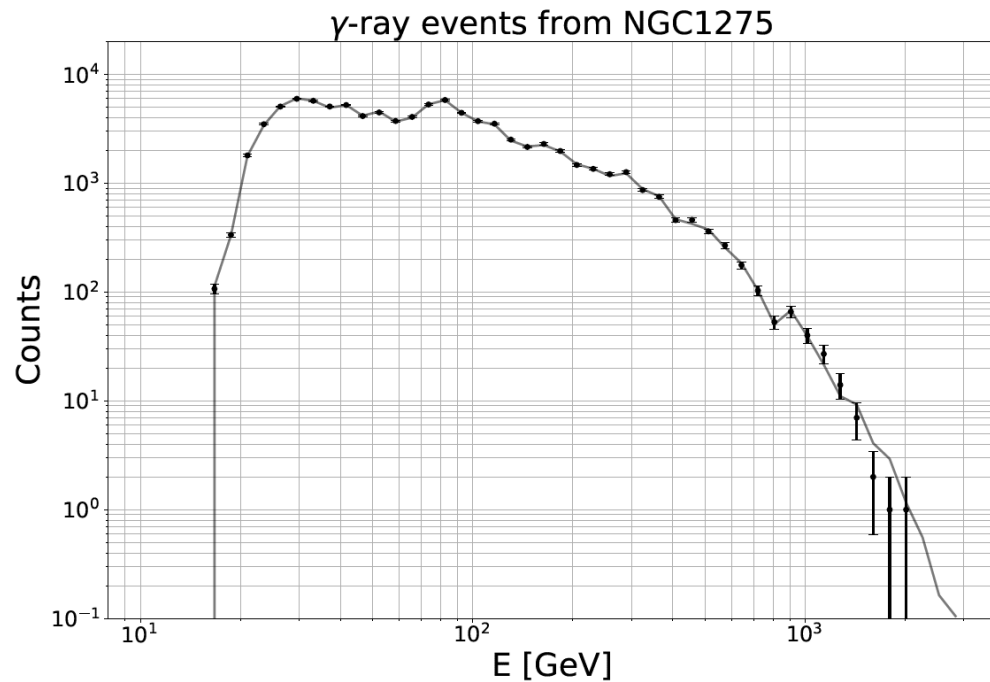
Working example for this talk: fitting the spectrum of NGC1275

$$\varphi(E) = \varphi_0 \left(\frac{E}{E_0} \right)^\gamma e^{-E/E_{cut}} + \text{poisson noise}$$



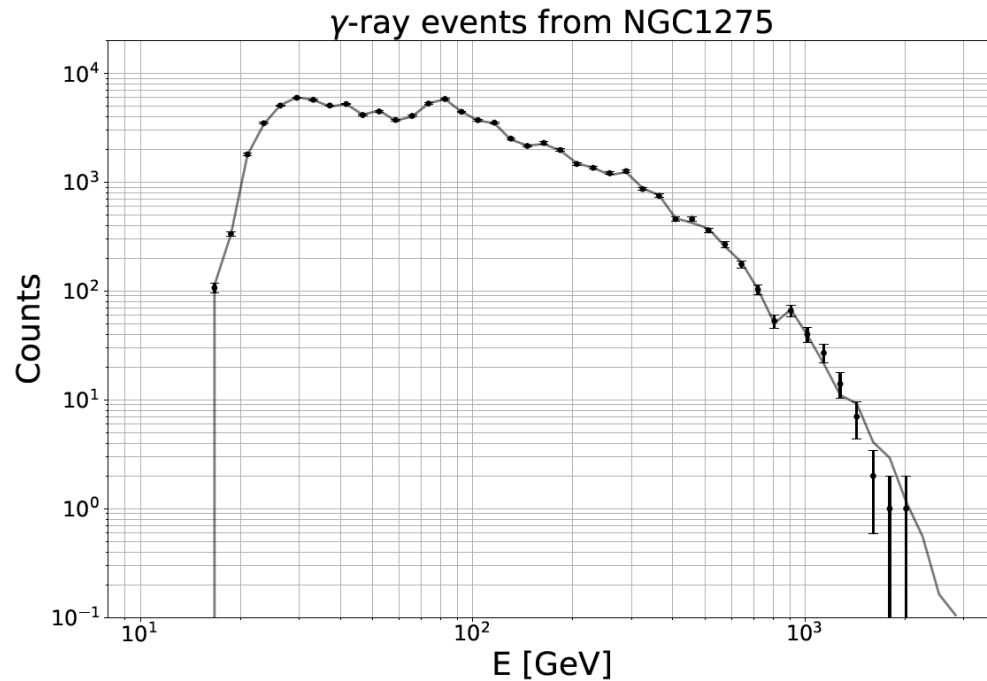
Working example for this talk: fitting the spectrum of NGC1275

$$\varphi(E) = \varphi_0 \left(\frac{E}{E_0} \right)^\gamma e^{-E/E_{cut}} + \text{poisson noise}$$

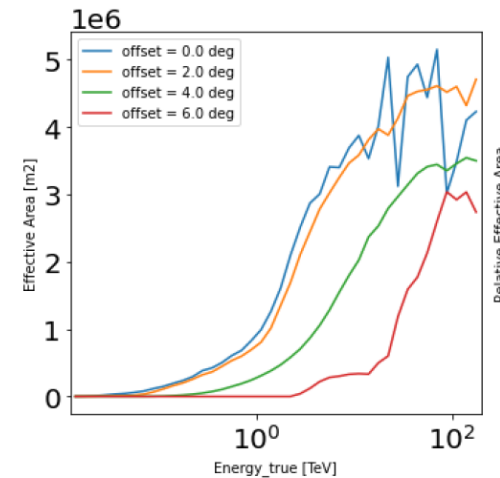


Working example for this talk: fitting the spectrum of NGC1275

$$\varphi(E) = \varphi_0 \left(\frac{E}{E_0} \right)^\gamma e^{-E/E_{cut}} + \text{poisson noise}$$



- IRF: prod3:South_z20_50h

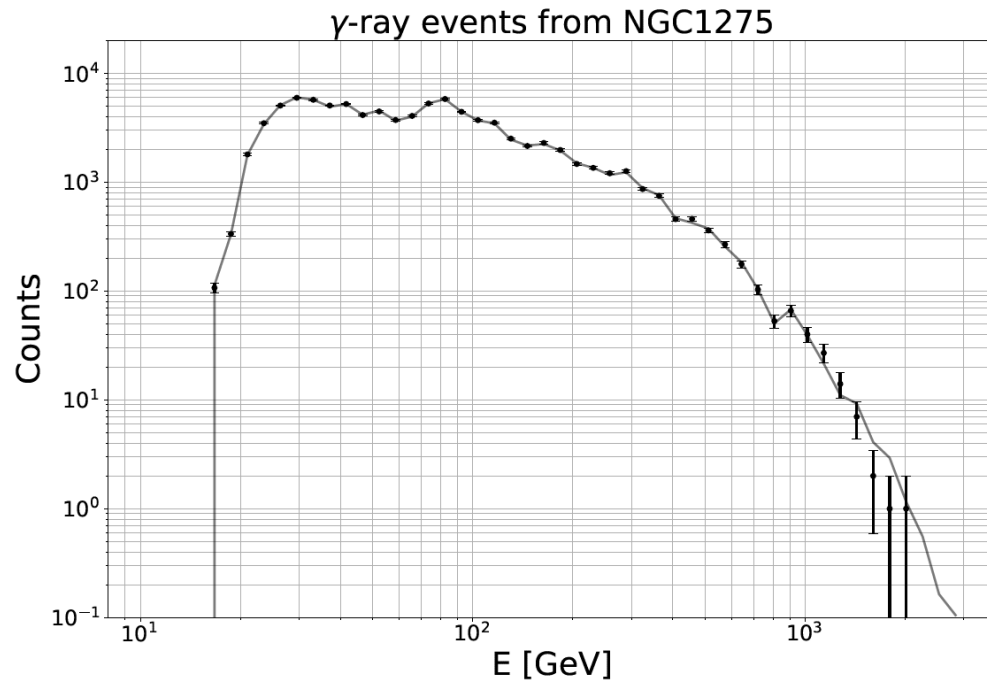


- Livetime = 2 hr
- *Total counts* ~ 90 000
- $E_0 = 153.86$ GeV

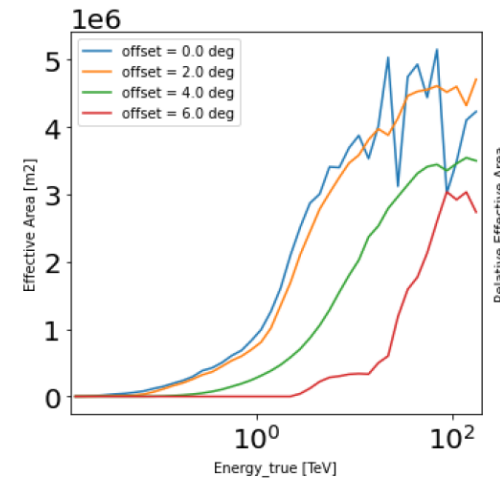
Working example for this talk: fitting the spectrum of NGC1275

$$\varphi(E) = \varphi_0 \left(\frac{E}{E_0} \right)^\gamma e^{-E/E_{cut}} + \text{poisson noise}$$

- Parameters of interest:
 - Amplitude φ_0
 - Spectral index γ
 - E_{cut}



- IRF: prod3:South_z20_50h



- Livetime = 2 hr
- Total counts $\sim 90\,000$
- $E_0 = 153.86$ GeV

We want to use Bayesian inference...

Bayes theorem:

$$p(\boldsymbol{\vartheta}|\boldsymbol{x}) = \frac{p(\boldsymbol{x}|\boldsymbol{\vartheta})}{\int d\boldsymbol{\theta} p(\boldsymbol{x}|\boldsymbol{\theta})} p(\boldsymbol{\vartheta})$$

$\boldsymbol{\vartheta}$ = parameters of interest (POIs)

$\boldsymbol{\theta}$ = POIs with nuisance parameters

\boldsymbol{x} = observed outcome

$p(\boldsymbol{\vartheta})$ = Prior (“a priori” assumption)

$p(\boldsymbol{x}|\boldsymbol{\vartheta})$ = $p(\boldsymbol{x}|\boldsymbol{\theta})$ integrated over
uncertainty in nuisance parameters

We want to use Bayesian inference...

Bayes theorem:

$$p(\boldsymbol{\vartheta}|\boldsymbol{x}) = \frac{p(\boldsymbol{x}|\boldsymbol{\vartheta})}{\int d\boldsymbol{\theta} p(\boldsymbol{x}|\boldsymbol{\theta})} p(\boldsymbol{\vartheta})$$

$\boldsymbol{\vartheta}$ = parameters of interest (POIs)

$\boldsymbol{\theta}$ = POIs with nuisance parameters

\boldsymbol{x} = observed outcome

$p(\boldsymbol{\vartheta})$ = Prior (“a priori” assumption)

$p(\boldsymbol{x}|\boldsymbol{\vartheta}) = p(\boldsymbol{x}|\boldsymbol{\theta})$ integrated over
uncertainty in nuisance parameters

We want to use Bayesian inference...

Bayes theorem:

$$p(\boldsymbol{\vartheta}|\boldsymbol{x}) = \frac{p(\boldsymbol{x}|\boldsymbol{\vartheta})}{\int d\boldsymbol{\theta} p(\boldsymbol{x}|\boldsymbol{\theta})} p(\boldsymbol{\vartheta})$$

$\boldsymbol{\vartheta}$ = parameters of interest (POIs)

$\boldsymbol{\theta}$ = POIs with nuisance parameters

\boldsymbol{x} = observed outcome

$p(\boldsymbol{\vartheta})$ = Prior (“a priori” assumption)

$p(\boldsymbol{x}|\boldsymbol{\vartheta}) = p(\boldsymbol{x}|\boldsymbol{\theta})$ integrated over
uncertainty in nuisance parameters

... but it is often too expensive

We want to use Bayesian inference...

Bayes theorem:

$$p(\boldsymbol{\vartheta}|\boldsymbol{x}) = \frac{p(\boldsymbol{x}|\boldsymbol{\vartheta})}{\int d\boldsymbol{\theta} p(\boldsymbol{x}|\boldsymbol{\theta})} p(\boldsymbol{\vartheta})$$

$\boldsymbol{\vartheta}$ = parameters of interest (POIs)

$\boldsymbol{\theta}$ = POIs with nuisance parameters

\boldsymbol{x} = observed outcome

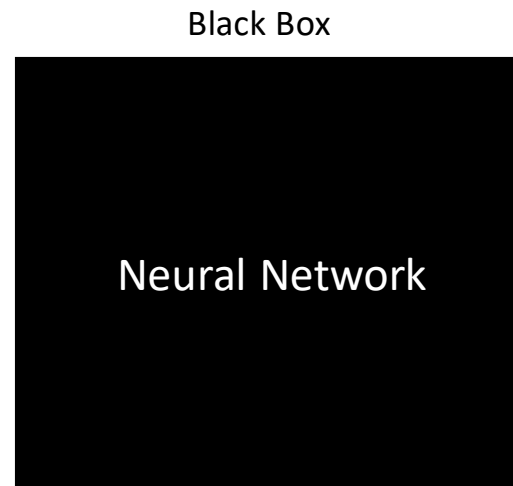
$p(\boldsymbol{\vartheta})$ = Prior (“a priori” assumption)

$p(\boldsymbol{x}|\boldsymbol{\vartheta}) = p(\boldsymbol{x}|\boldsymbol{\theta})$ integrated over
uncertainty in nuisance parameters

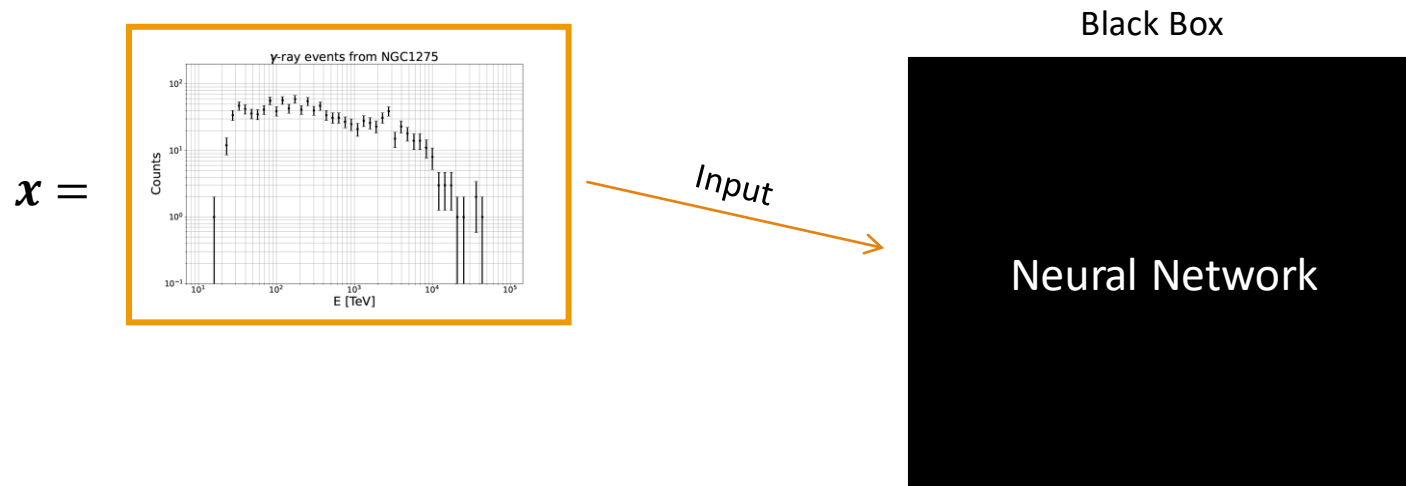
... but it is often too expensive

- Computing cost explodes with number of nuisance parameters
 - Need to make simplifying assumptions

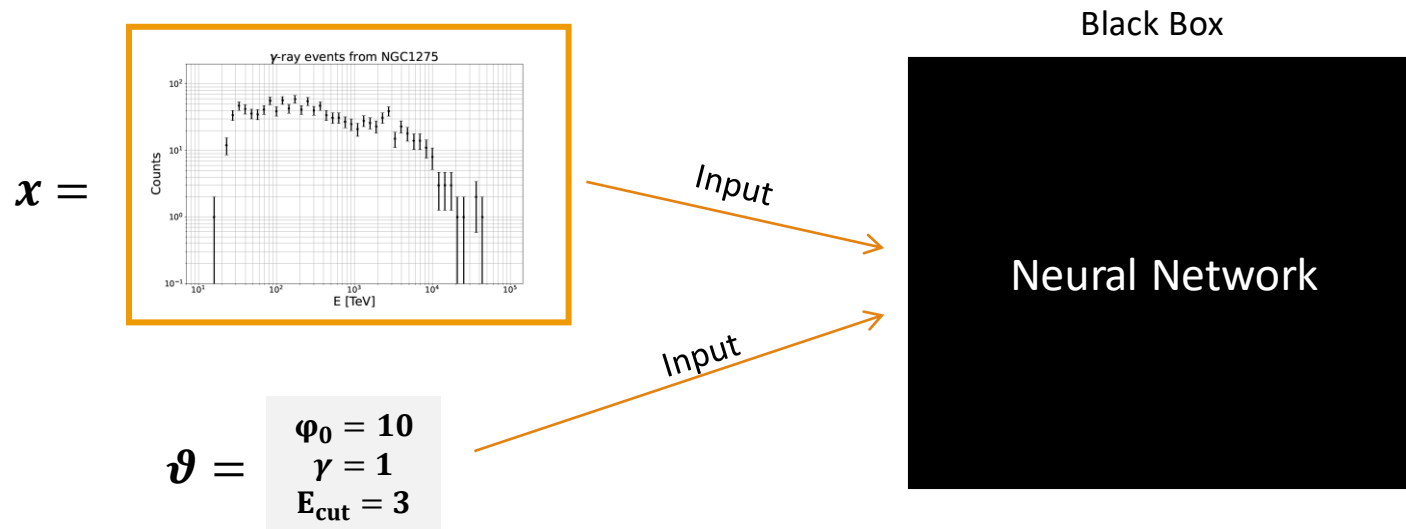
Neural networks can approximate posteriors!



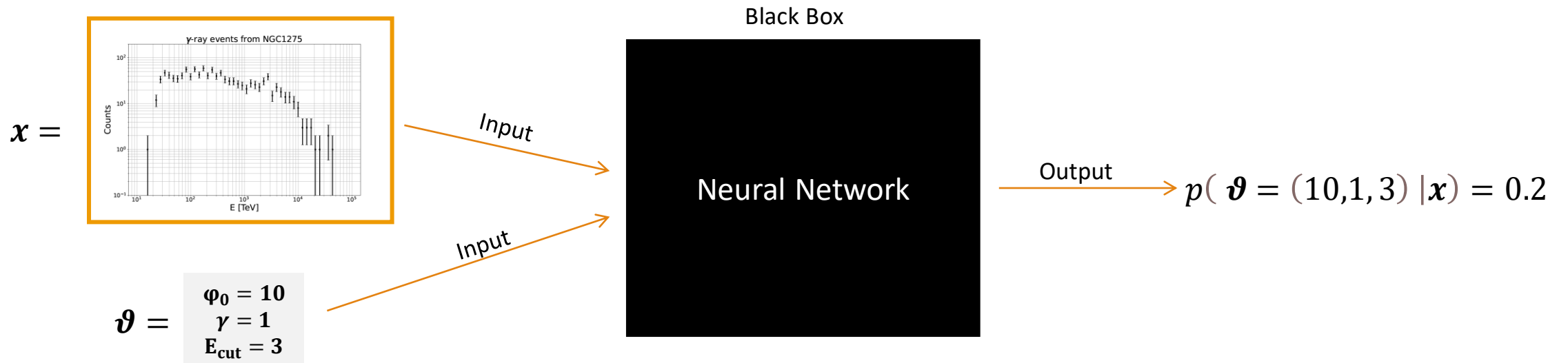
Neural networks can approximate posteriors!



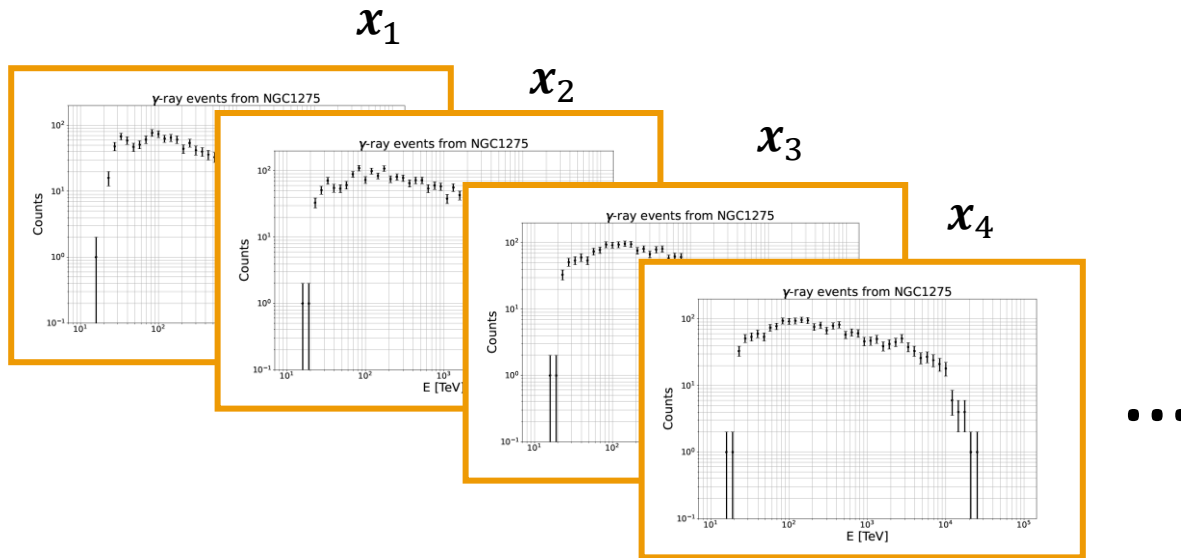
Neural networks can approximate posteriors!



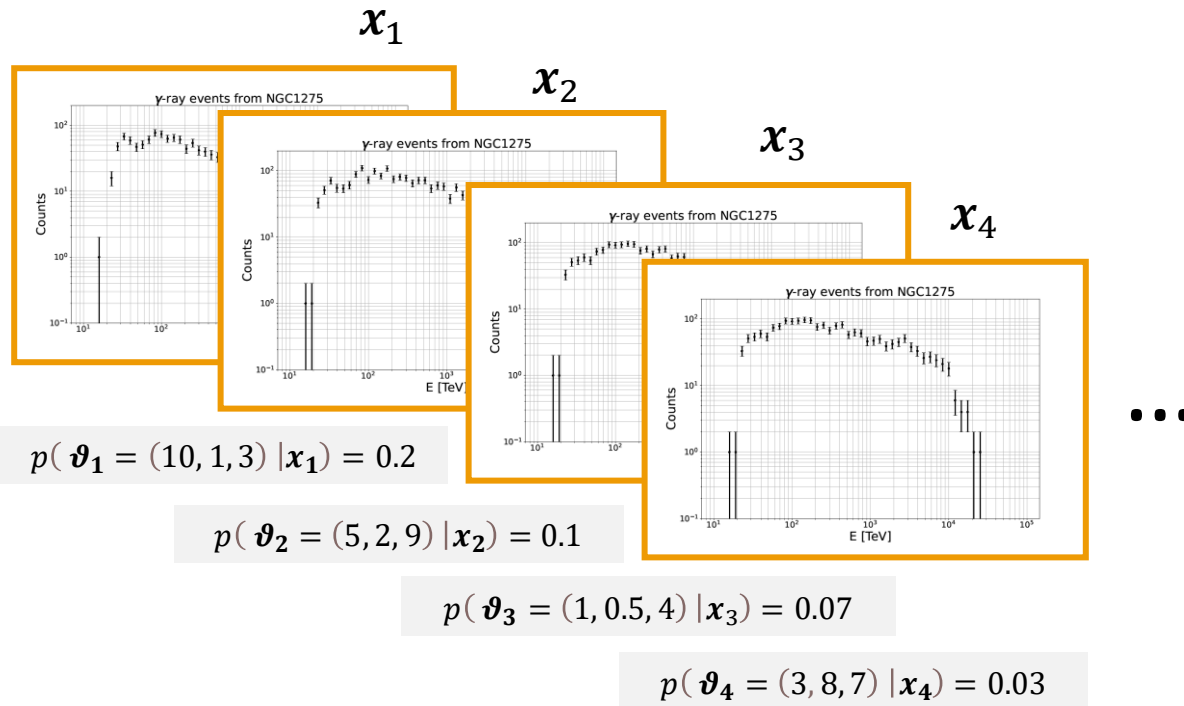
Neural networks can approximate posteriors!



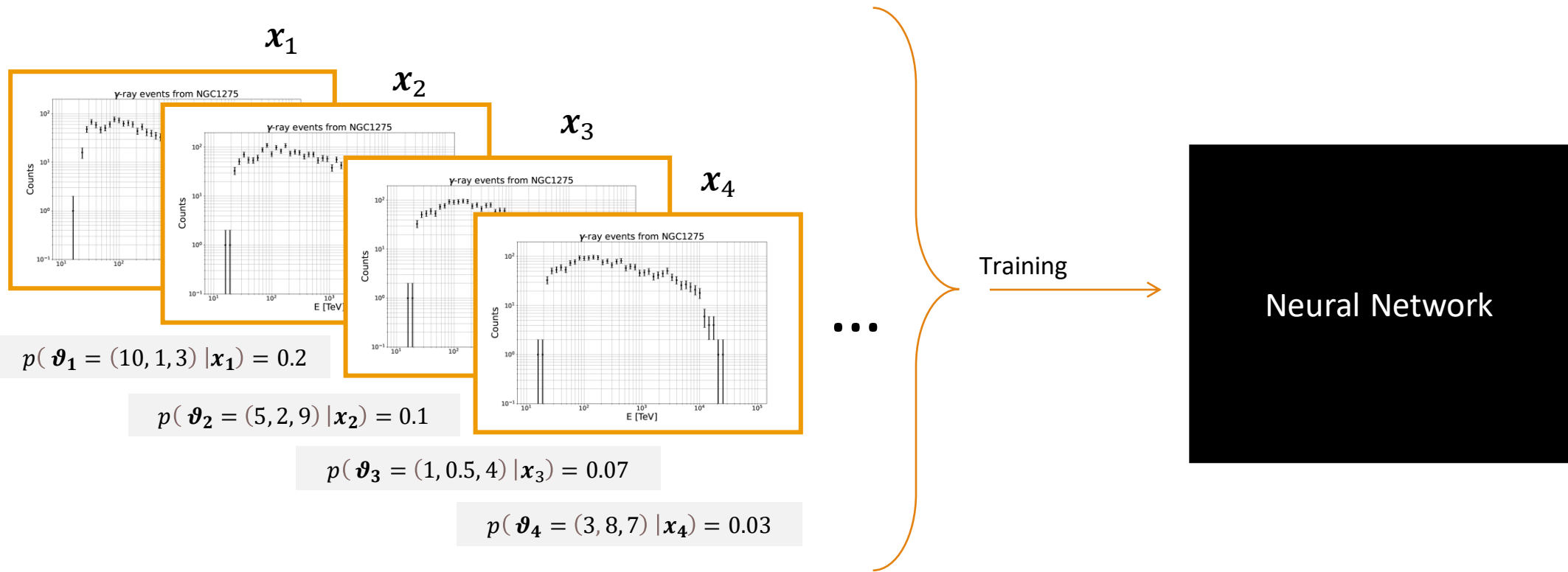
A network learns to do approximations by example:



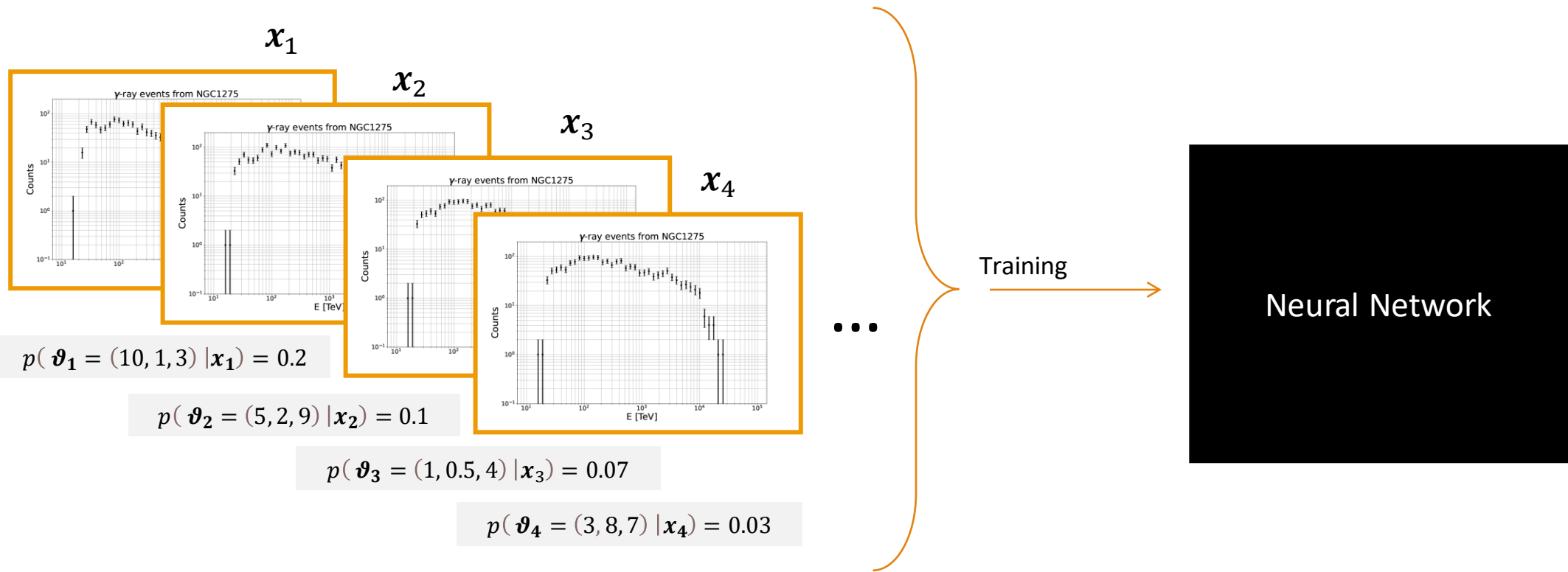
A network learns to do approximations by example:



A network learns to do approximations by example:



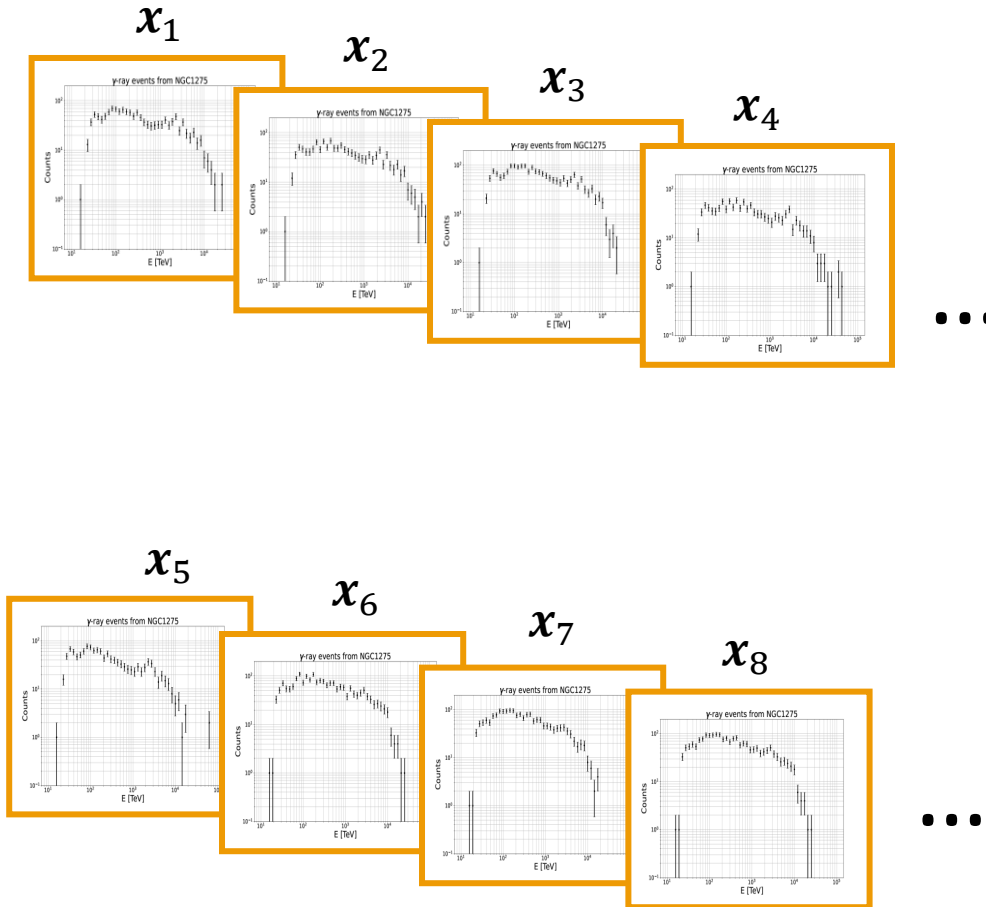
A network learns to do approximations by example:



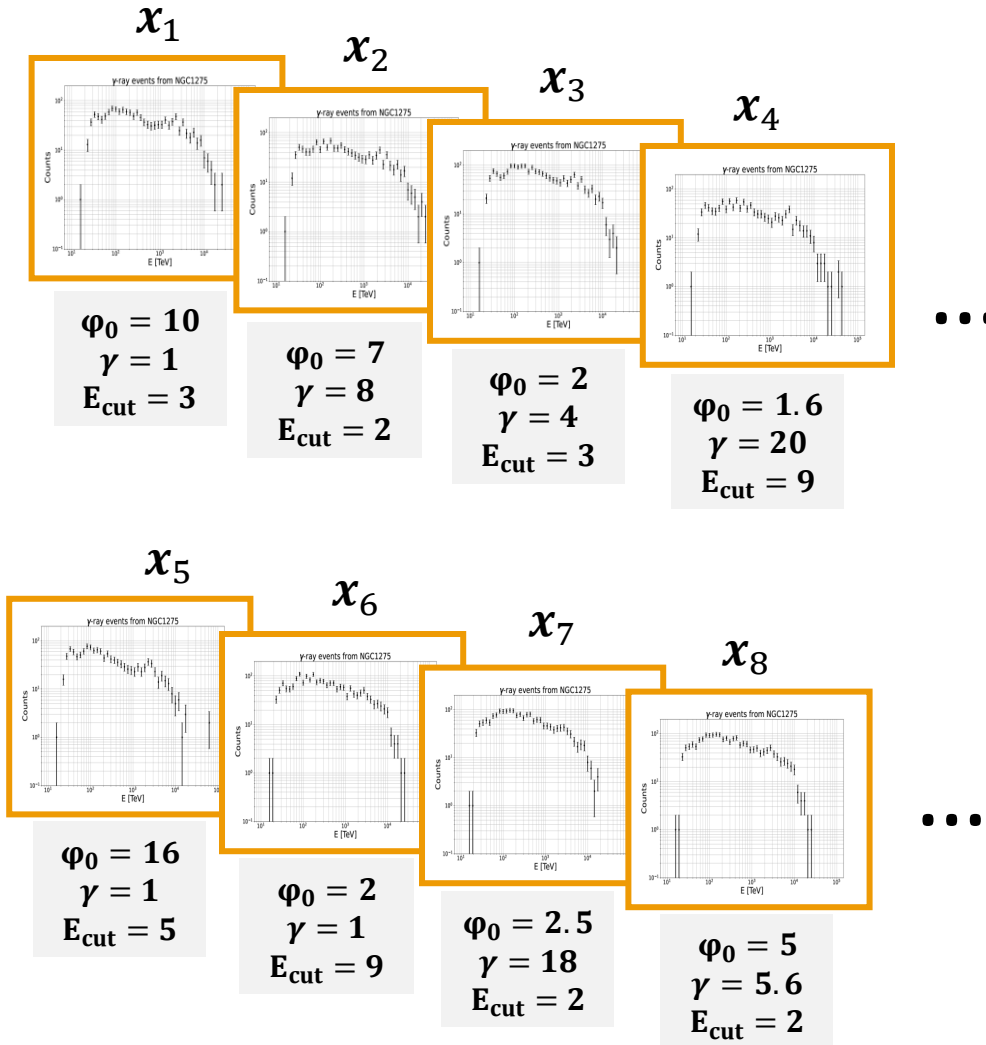
... but this assumes that we can already calculate the posterior!

We can “trick” a network to learn the posterior implicitly

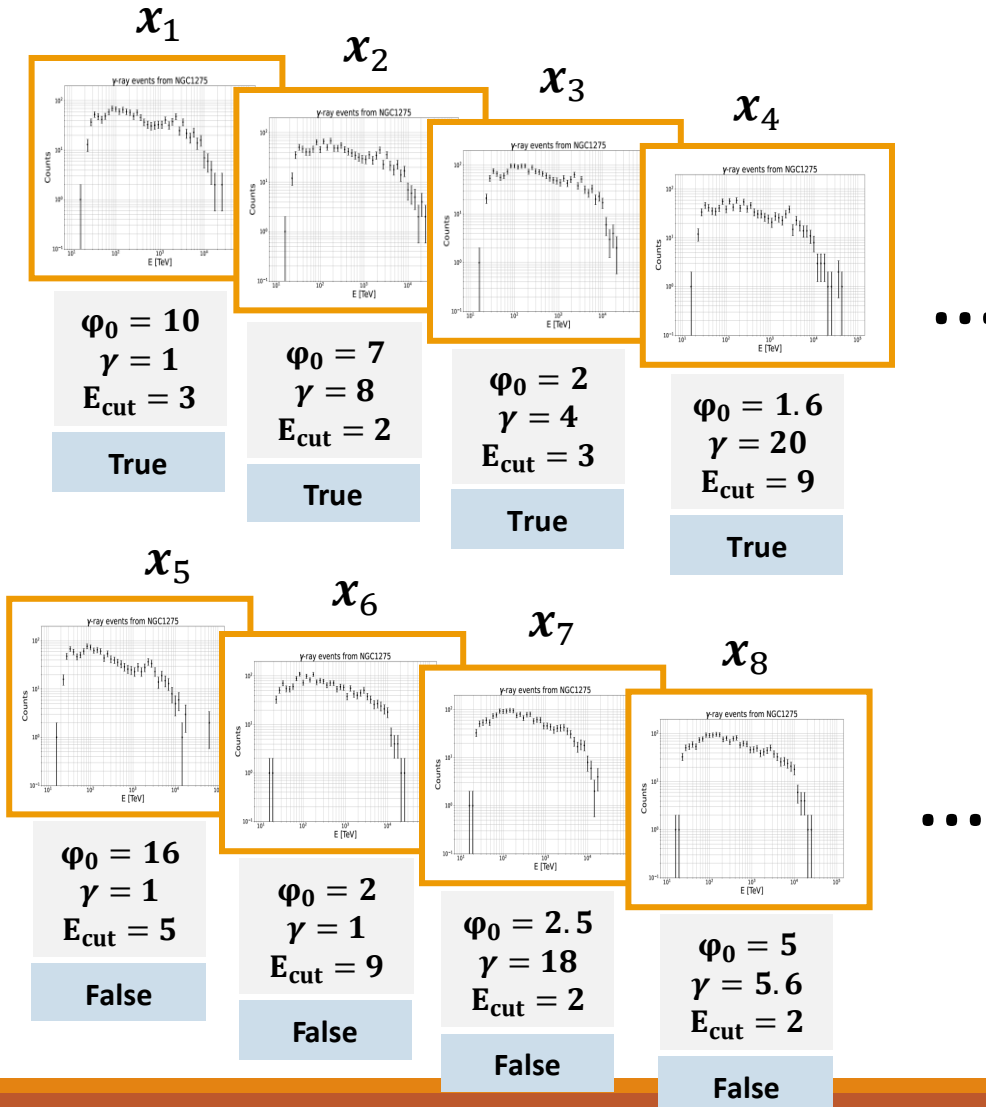
We can “trick” a network to learn the posterior implicitly



We can “trick” a network to learn the posterior implicitly



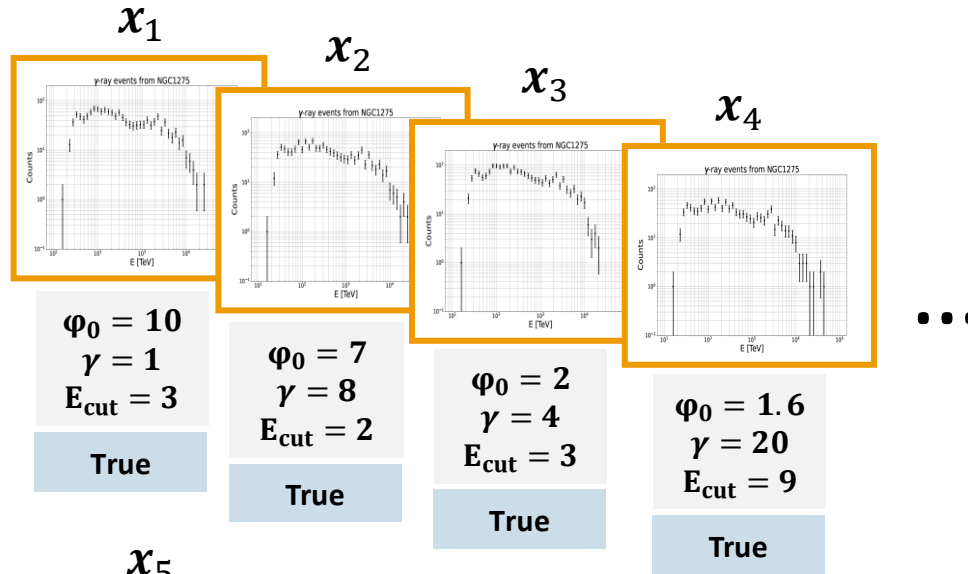
We can “trick” a network to learn the posterior implicitly



We can “trick” a network to learn the posterior implicitly

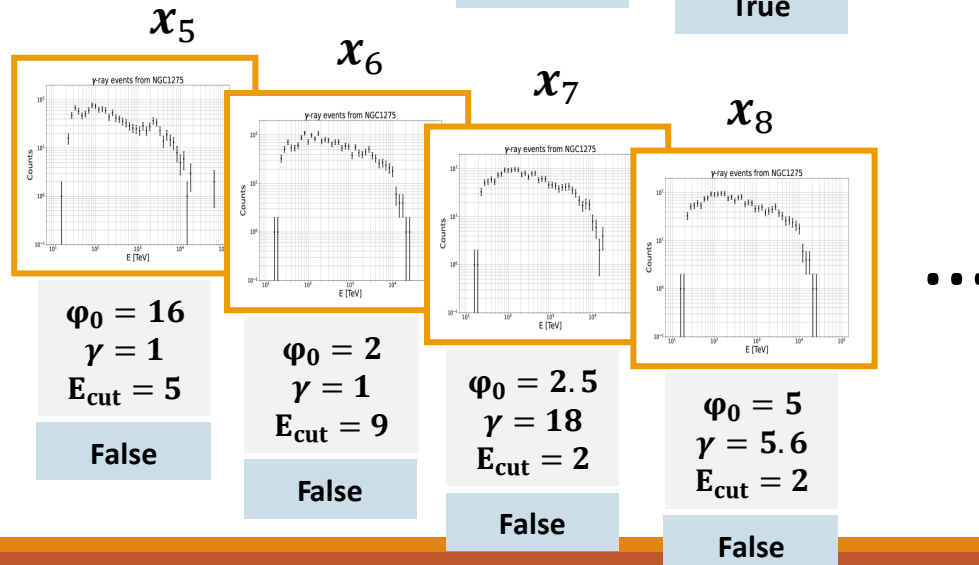
True

The observations are simulated according to the parameter values



False

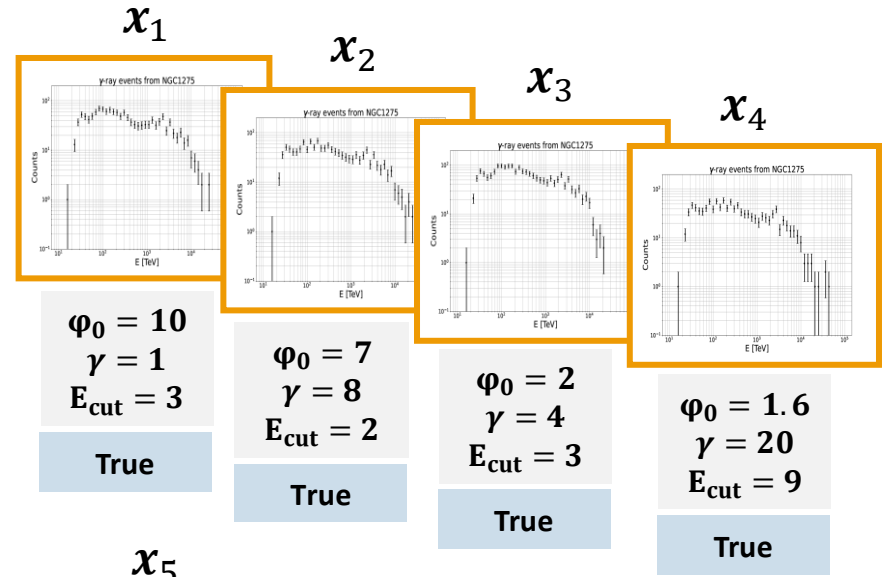
The parameter values are chosen independently from the simulations



We can “trick” a network to learn the posterior implicitly

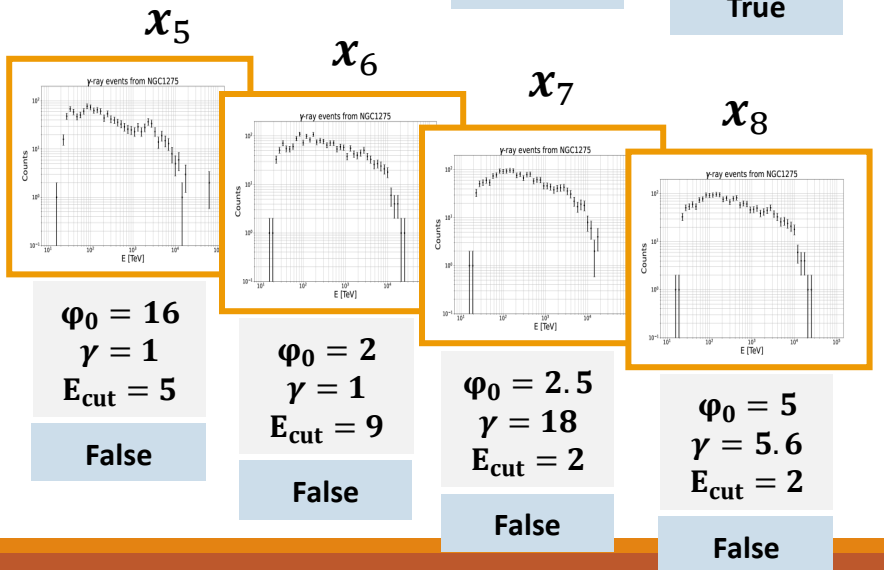
True

The observations are simulated according to the parameter values

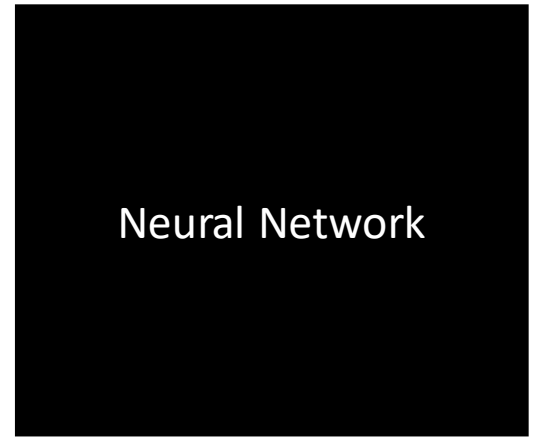


False

The parameter values are chosen independently from the simulations

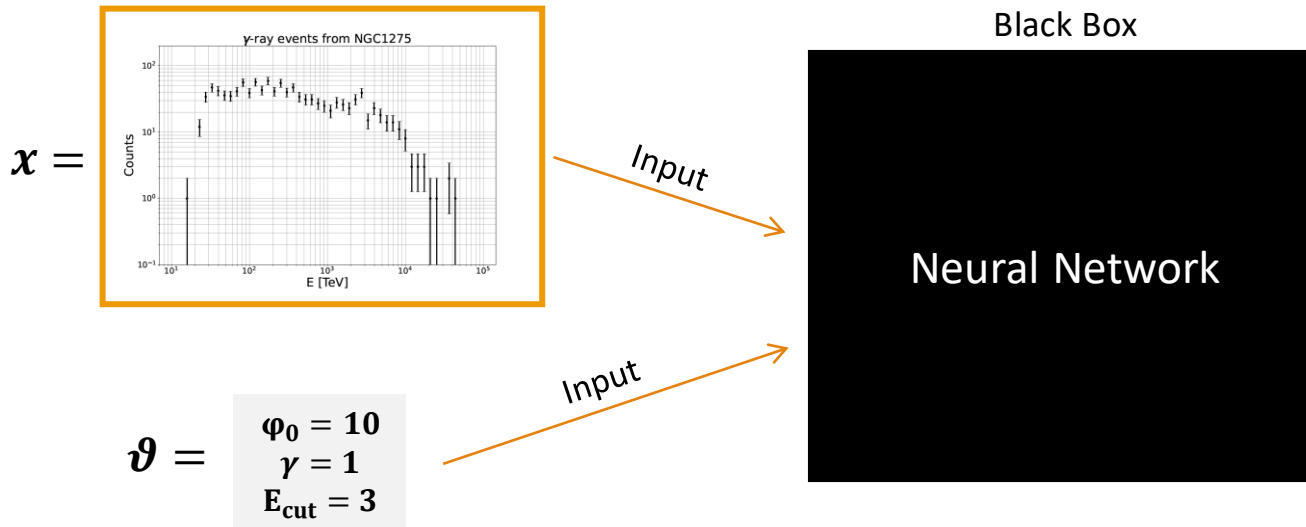


Training
to distinguish between
True and **False**

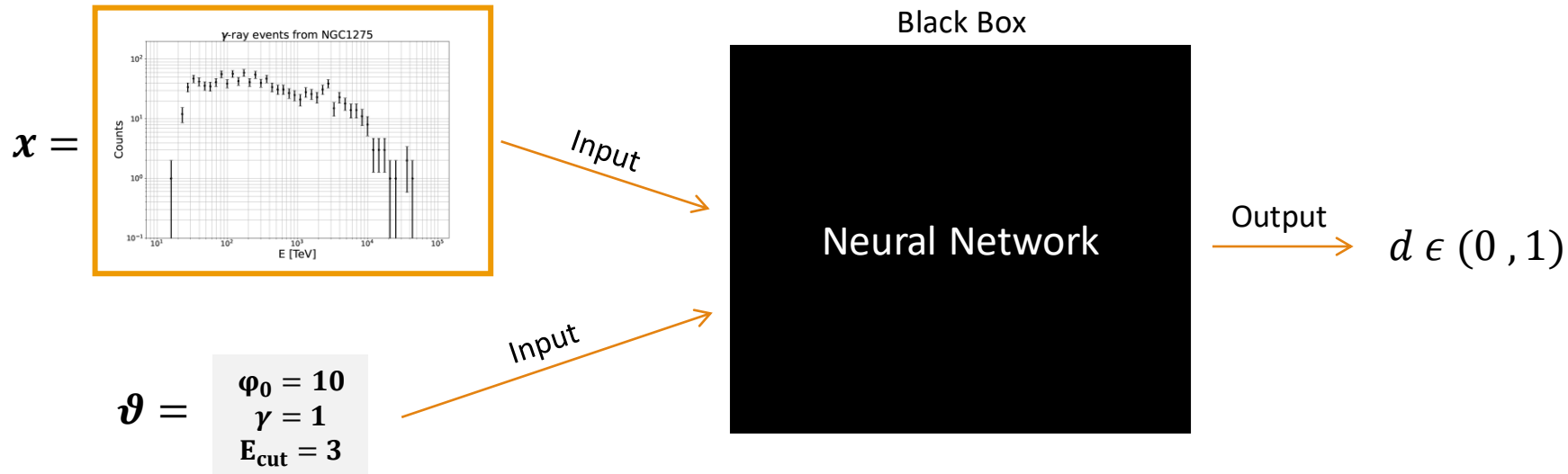


The network now has an implicit understanding of the posterior

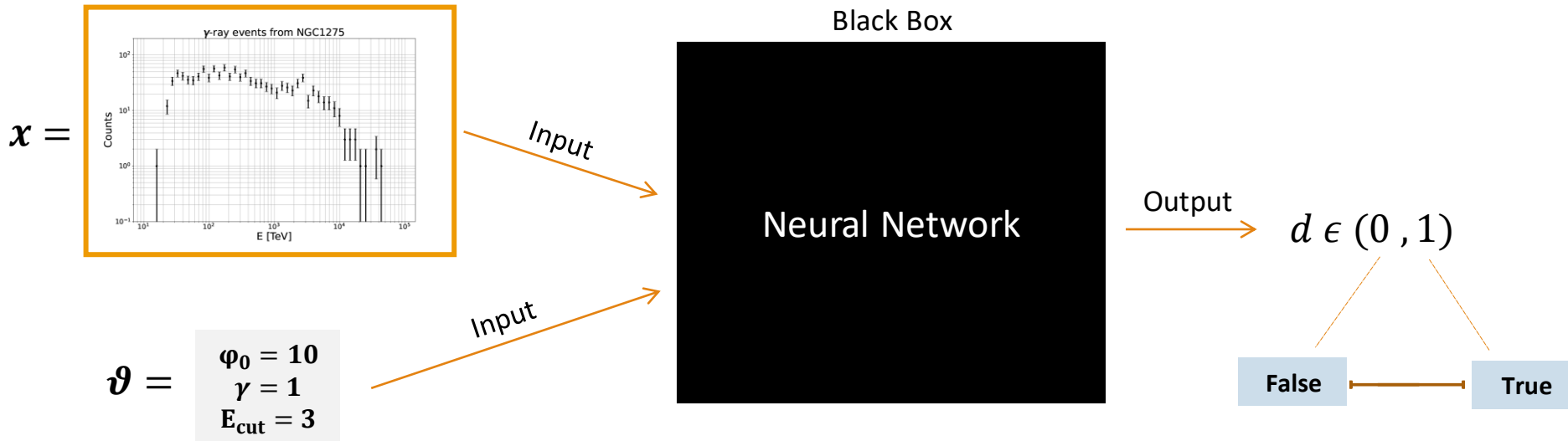
The network now has an implicit understanding of the posterior



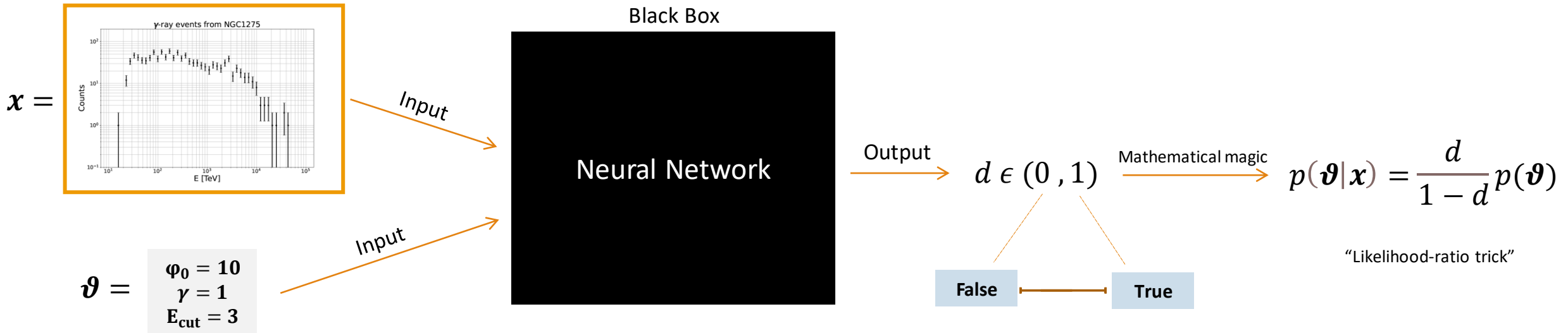
The network now has an implicit understanding of the posterior



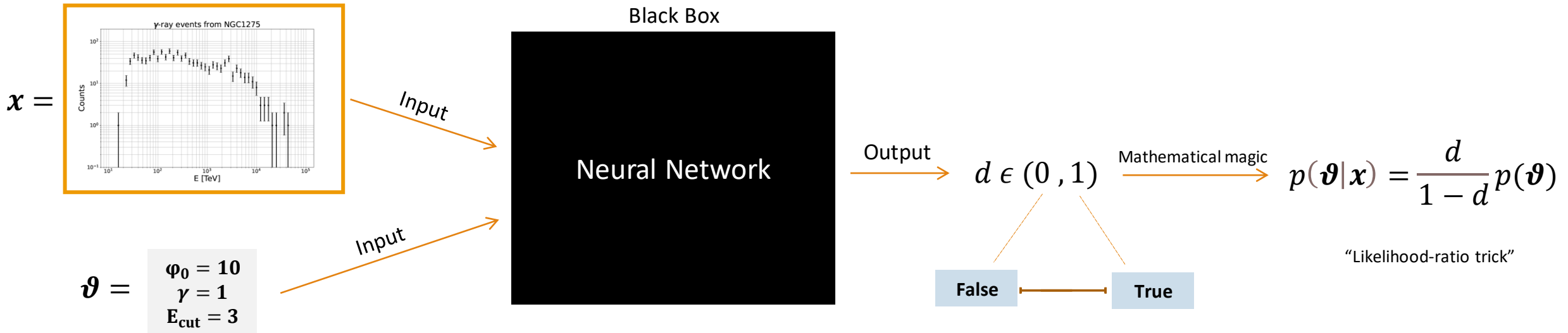
The network now has an implicit understanding of the posterior



The network now has an implicit understanding of the posterior



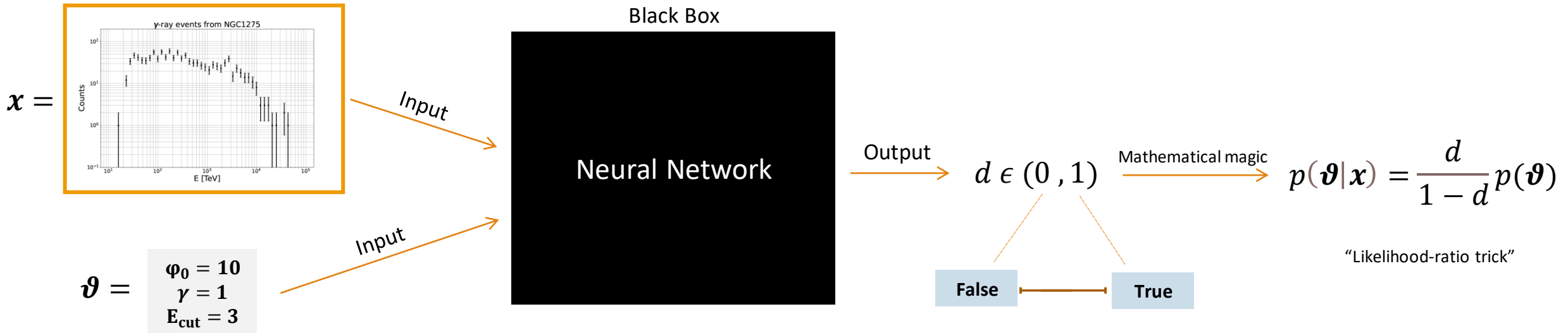
The network now has an implicit understanding of the posterior



Conditions for the trick to work:

- When generating the training set, we have to draw all of the parameters of interest from the prior distribution.
- The network’s loss function must be the binary cross entropy.

The network now has an implicit understanding of the posterior

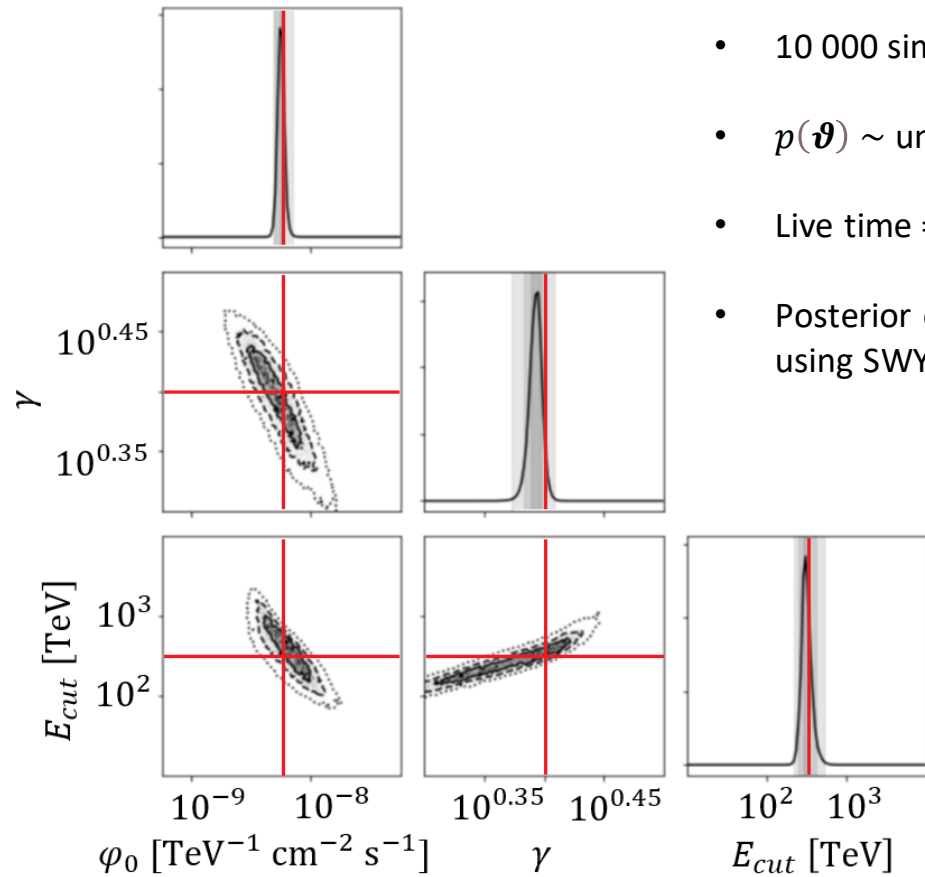


Vary this input
to scan the posterior
over parameter space

Conditions for the trick to work:

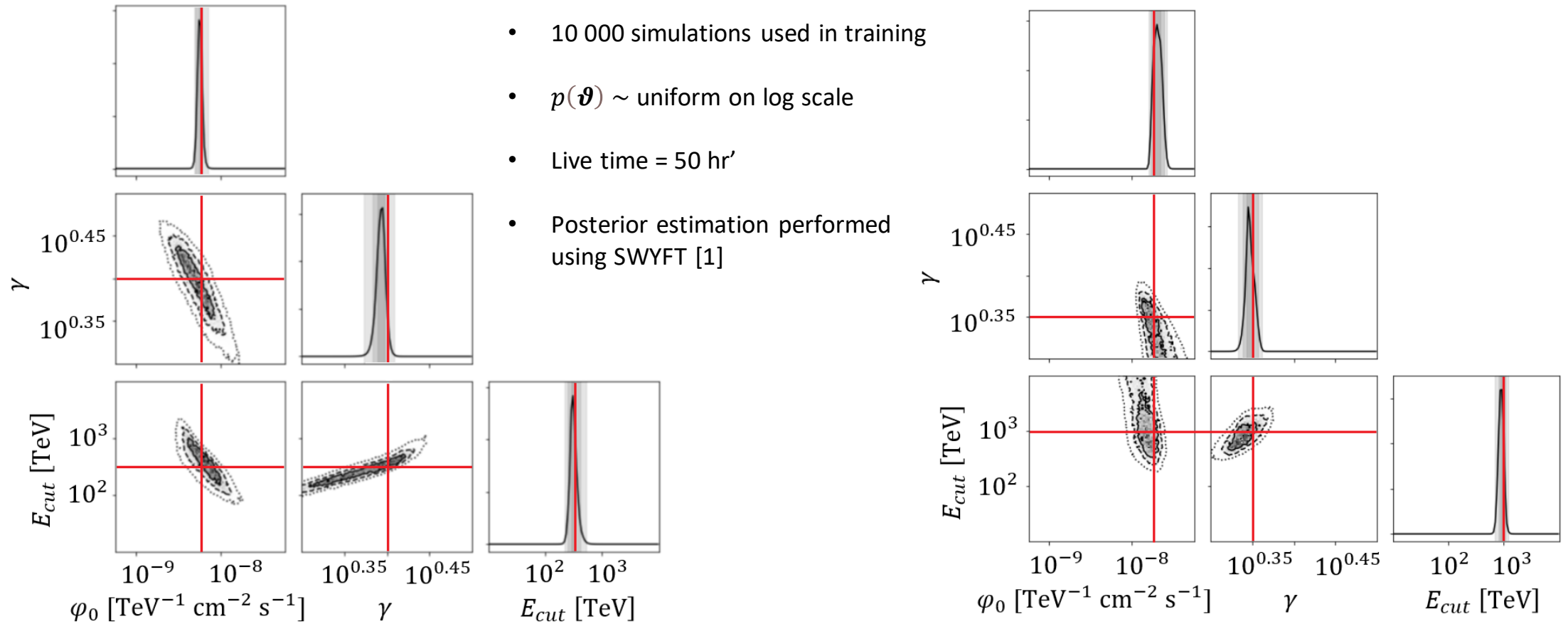
- When generating the training set, we have to draw all of the parameters of interest from the prior distribution.
- The network’s loss function must be the binary cross entropy.

Inference with NRE seems to be precise for the spectral fit



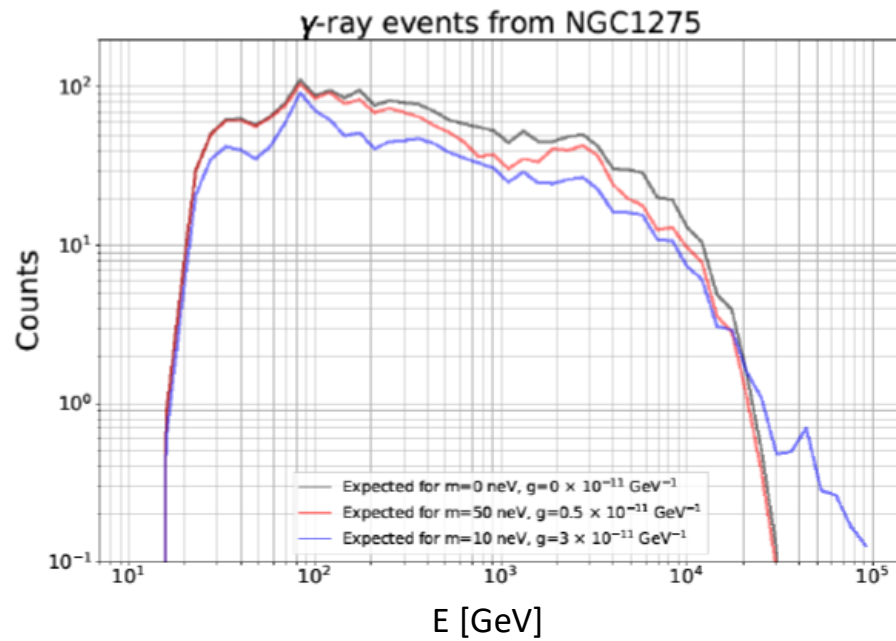
- 10 000 simulations used in training
- $p(\boldsymbol{\vartheta}) \sim$ uniform on log scale
- Live time = 50 hr'
- Posterior estimation performed using SWYFT [1]

Inference with NRE seems to be precise for the spectral fit



NRE seems particularly useful for ALP searches with gamma-telescopes

Expected result assuming ALPs with given mass m and coupling g (using `gammaALPs`):

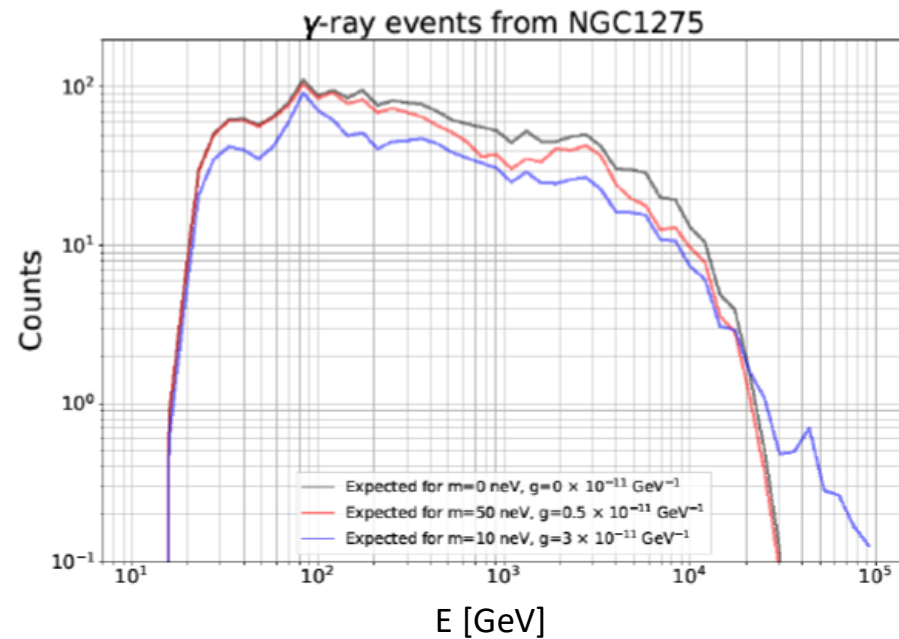


The expected spectrum can be simulated using `gammapy` [2] and `gammaALPs` [3] (by M. Meyer)

NRE seems particularly useful for ALP searches with gamma-telescopes

Expected result assuming ALPs with given mass m and coupling g (using `gammaALPs`):

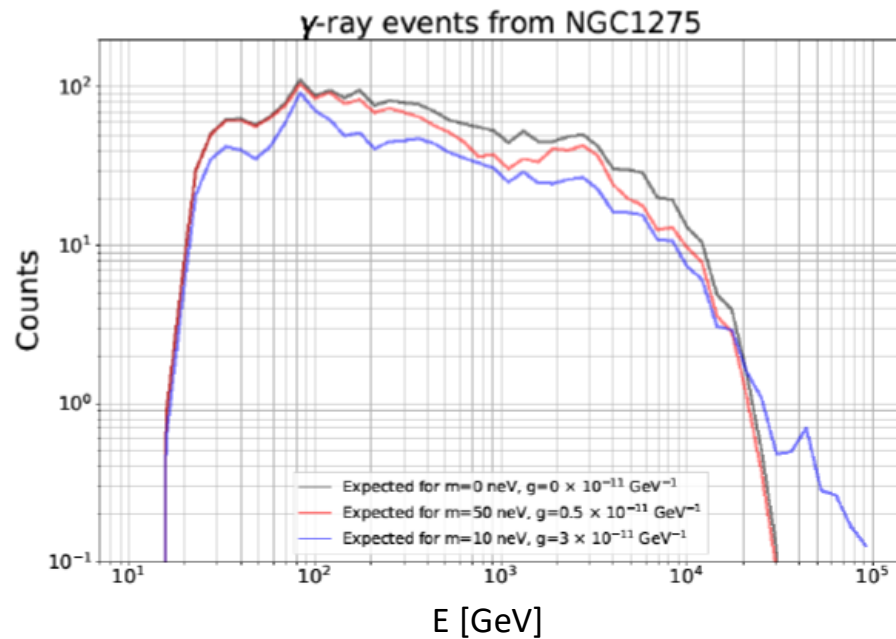
- Parameters of interest:
 - Mass of ALPs, m
 - ALP-photon coupling, g



The expected spectrum can be simulated using `gammapy` [2] and `gammaALPs` [3] (by M. Meyer)

NRE seems particularly useful for ALP searches with gamma-telescopes

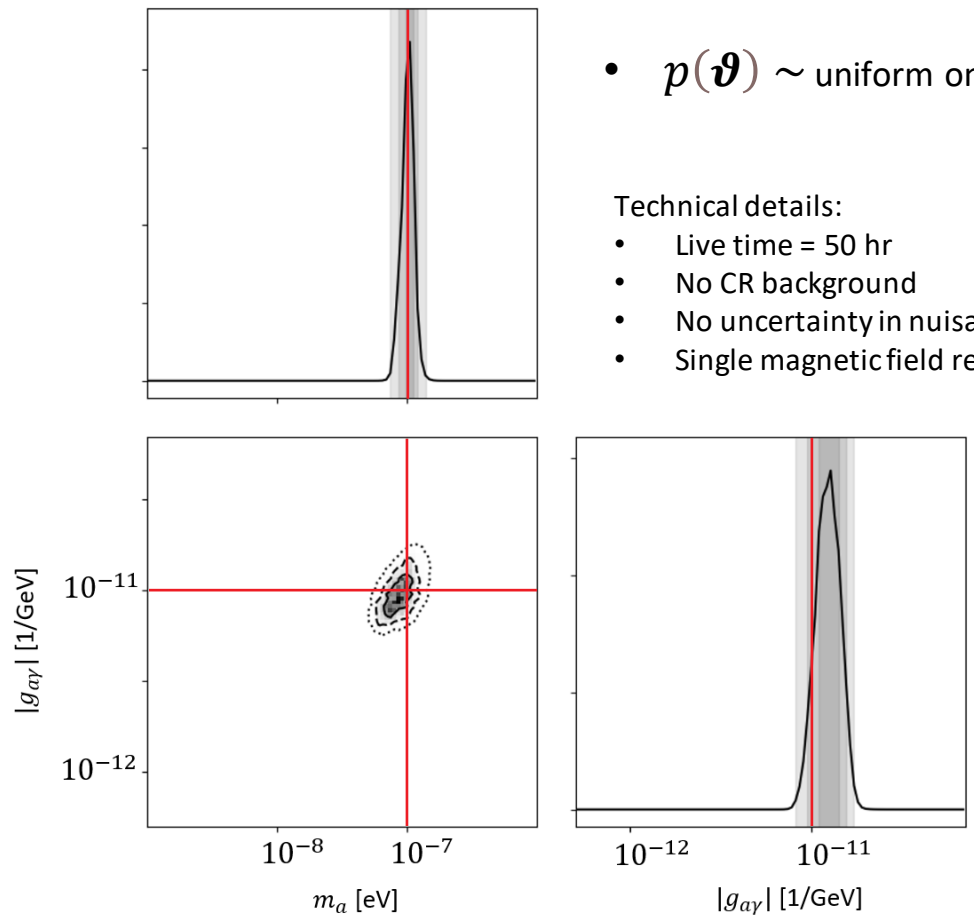
Expected result assuming ALPs with given mass m and coupling g (using `gammaALPs`):



The expected spectrum can be simulated using `gammapy` [2] and `gammaALPs` [3] (by M. Meyer)

- Parameters of interest:
 - Mass of ALPs, m
 - ALP-photon coupling, g
- Nuisance parameters:
 - Amplitude
 - Spectral index
 - Cut-off energy
 - Magnetic field configuration
 - + 12 more related to configuration of NGC1275

Preliminary results indicate the method is suitable for ALP searches

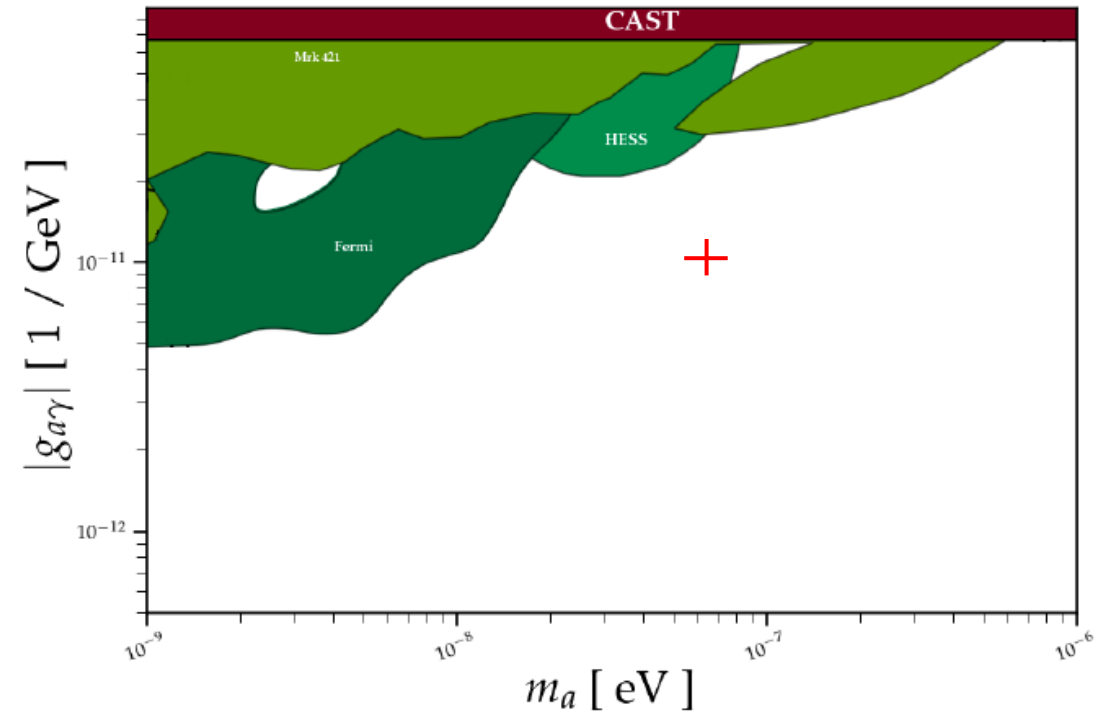


- $p(\boldsymbol{\vartheta}) \sim$ uniform on log scale

Technical details:

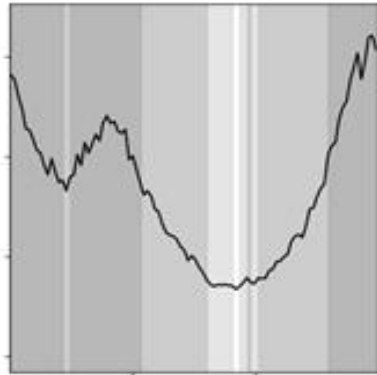
- Live time = 50 hr
- No CR background
- No uncertainty in nuisance parameters
- Single magnetic field realization

Very unrealistic!



Full limits: [4]

Preliminary results indicate the method is suitable for ALP searches

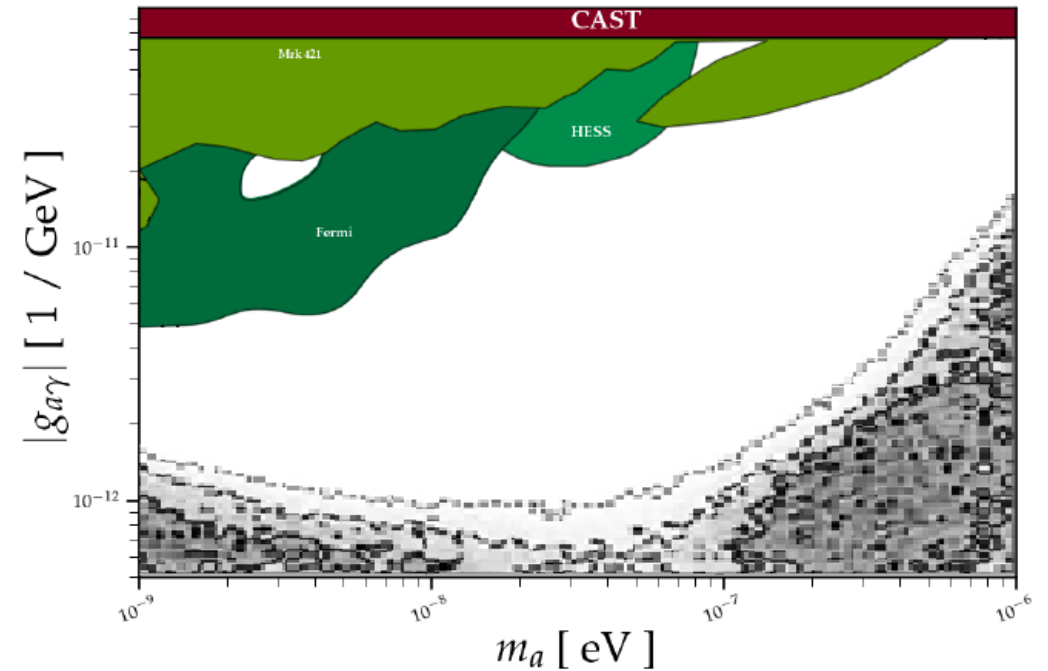
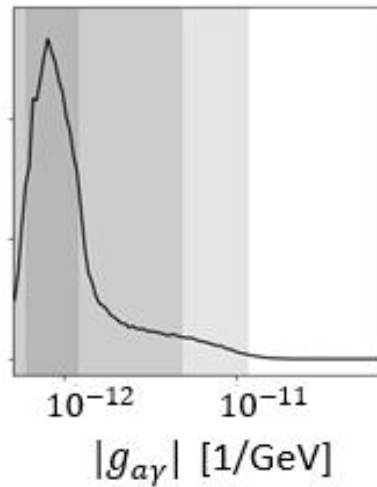
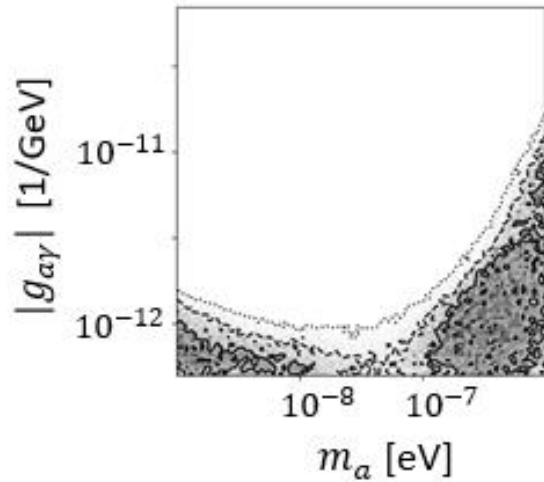


- $(m_{true}, g_{true}) = (0,0)$
- $p(\boldsymbol{\vartheta}) \sim$ uniform on log scale

Technical details:

- Live time = 50 hr
- No CR background
- No uncertainty in nuisance parameters
- Single magnetic field realization

Very unrealistic!



Full limits: [4]

Neural Ratio Estimation: some resources to get you started

- Original paper (to my knowledge) to introduce the concept:

<https://arxiv.org/abs/1903.04057>

- B. K. Miller, A. Cole, P. Forré, G. Louppe, and C. Weniger, “Truncated marginal neural ratio estimation”.

<https://arxiv.org/abs/2107.01214>

- B. K. Miller, A. Cole, G. Louppe, and C. Weniger, “Simulation-efficient marginal posterior estimation with swyft: stop wasting your precious time,”

<https://arxiv.org/abs/2011.13951>

Jun 2020

Likelihood-free MCMC with Amortized Approximate Ratio Estimators

Joeri Hermans¹ Volodimir Begy² Gilles Louppe¹

Abstract

Posterior inference with an intractable likelihood is becoming an increasingly common task in scientific domains which rely on sophisticated computer simulations. Typically, these forward models do not admit tractable densities forcing practitioners to make use of approximations. This work introduces a novel approach to address the intractability of the likelihood and the marginal

ratio of posterior densities between consecutive states in the Markov chain. This allows the posterior to be approximated numerically, provided that the likelihood $p(x|\theta)$ and the prior $p(\theta)$ are tractable. We consider the equally common and more challenging setting, the so-called likelihood-free setup, in which the likelihood cannot be evaluated in a reasonable amount of time or has no tractable closed-form expression. However, drawing samples from the forward model is possible.

- SWYFT (a package under development that does NRE and more): <https://github.com/undark-lab/swyft>

One take away

One take away

If you:

- 1) Want to infer parameters

One take away

If you:

- 1) Want to infer parameters
- 2) Want to avoid simplifications

One take away

If you:

- 1) Want to infer parameters
- 2) Want to avoid simplifications
- 3) Can make many simulations

One take away

If you:

- 1) Want to infer parameters
- 2) Want to avoid simplifications
- 3) Can make many simulations
- 4) Are not a die-hard frequentist

One take away

If you:

- 1) Want to infer parameters
- 2) Want to avoid simplifications
- 3) Can make many simulations
- 4) Are not a die-hard frequentist

Then Neural Ratio Estimation is your friend!

References

- [1] SWYFT: <https://github.com/undark-lab/swyft>
- [2] Gammapy: <https://gammapy.org/>
- [3] gammaALPs: <https://github.com/me-manu/gammaALPs>
- [4] ALP limits: <https://github.com/cajohare/AxionLimits>

Backup

Derivation of the likelihood ratio trick

When the neural network is trained, it is actually trying to minimize a loss function. Different loss functions are possible, but for the likelihood ratio trick to work, it is necessary that we use the *Binary Cross Entropy*, which looks like this:

$$\text{LOSS} = - \sum \left[y \ln d + (1 - y) \ln(1 - d) \right] \quad (\text{Summed over all training samples})$$

- $y = 1$ if training sample is from *True* category, 0 if *False*.
- d = output from neural network for given sample input

Notice that the loss becomes smaller when the network manages to categorize more samples correctly. In the limit of infinite training samples, the loss becomes

$$\text{LOSS} \rightarrow \iint \left[\underbrace{p(\mathbf{x}|\boldsymbol{\vartheta})p(\boldsymbol{\vartheta})}_{\text{Probability of appearance of a sample of the True category (because } \mathbf{x} \text{ is simulated taking } \boldsymbol{\vartheta} \text{ as a premise)}} \ln d + \underbrace{p(\mathbf{x})p(\boldsymbol{\vartheta})}_{\text{Probability of appearance of a sample of the False category (because } \mathbf{x} \text{ and } \boldsymbol{\vartheta} \text{ are independent)}} \ln(1 - d) \right] d\mathbf{x}d\boldsymbol{\vartheta}$$

We assume that the network manages to optimize the loss function perfectly. If this is the case, the derivative of the loss with respect to the neural network's hyperparameters (weights and biases, which are adjusted during training), which we denote by $\boldsymbol{\varphi}$, is 0. Notice that in the integrand, only d is dependent on $\boldsymbol{\varphi}$.

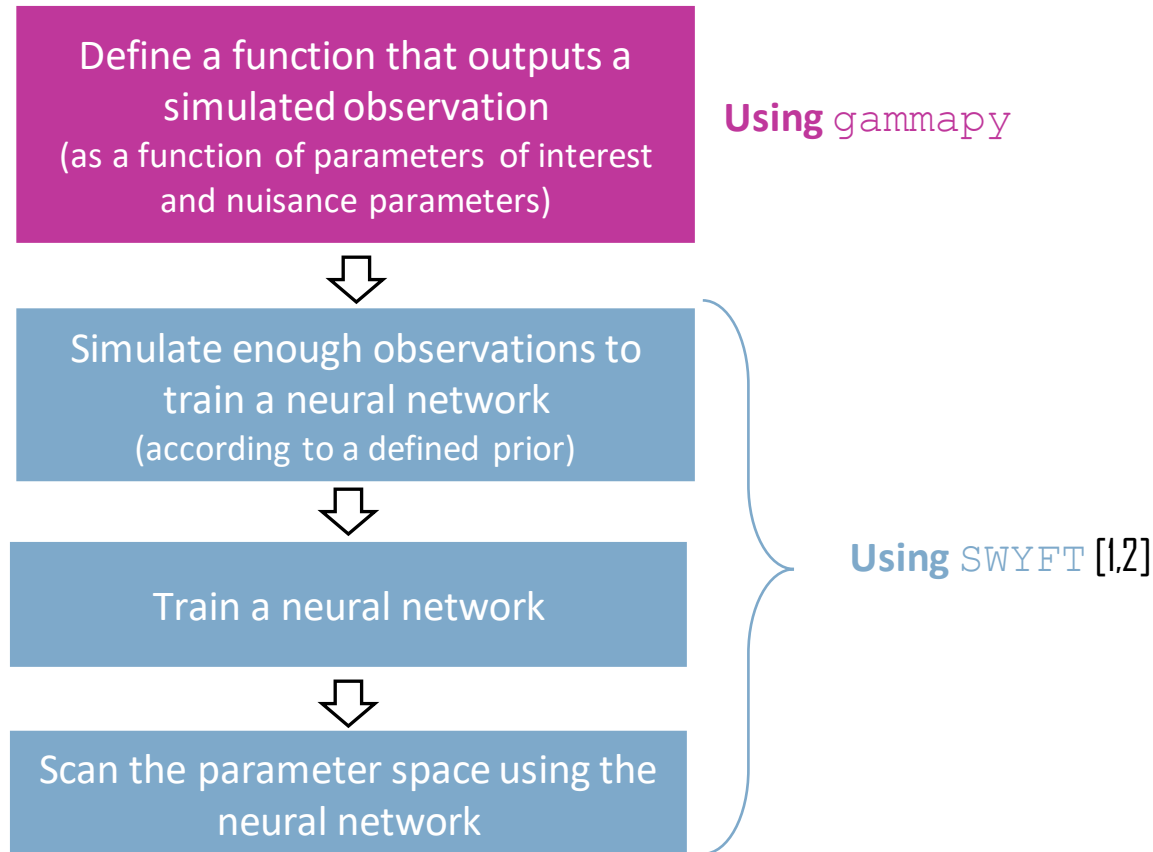
$$\frac{\partial}{\partial \boldsymbol{\varphi}} \text{LOSS} = \iint \left[\frac{p(\mathbf{x}|\boldsymbol{\vartheta})p(\boldsymbol{\vartheta})}{d} - \frac{p(\mathbf{x})p(\boldsymbol{\vartheta})}{1 - d} \right] \frac{\partial d}{\partial \boldsymbol{\varphi}} = 0$$

The likelihood ratio trick follows from setting the square brackets to zero:

$$\Rightarrow \frac{d}{1-d} = \frac{p(\mathbf{x}|\boldsymbol{\vartheta})}{p(\mathbf{x})} = \frac{p(\mathbf{x}|\boldsymbol{\vartheta})}{\int d\boldsymbol{\theta} p(\mathbf{x}|\boldsymbol{\theta})}$$

Workflow

Workflow:



- The simulations implicitly contain the information on the relationship between parameters and observations