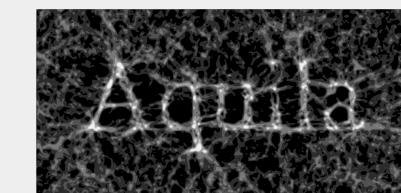


# Parameter inference using neural networks

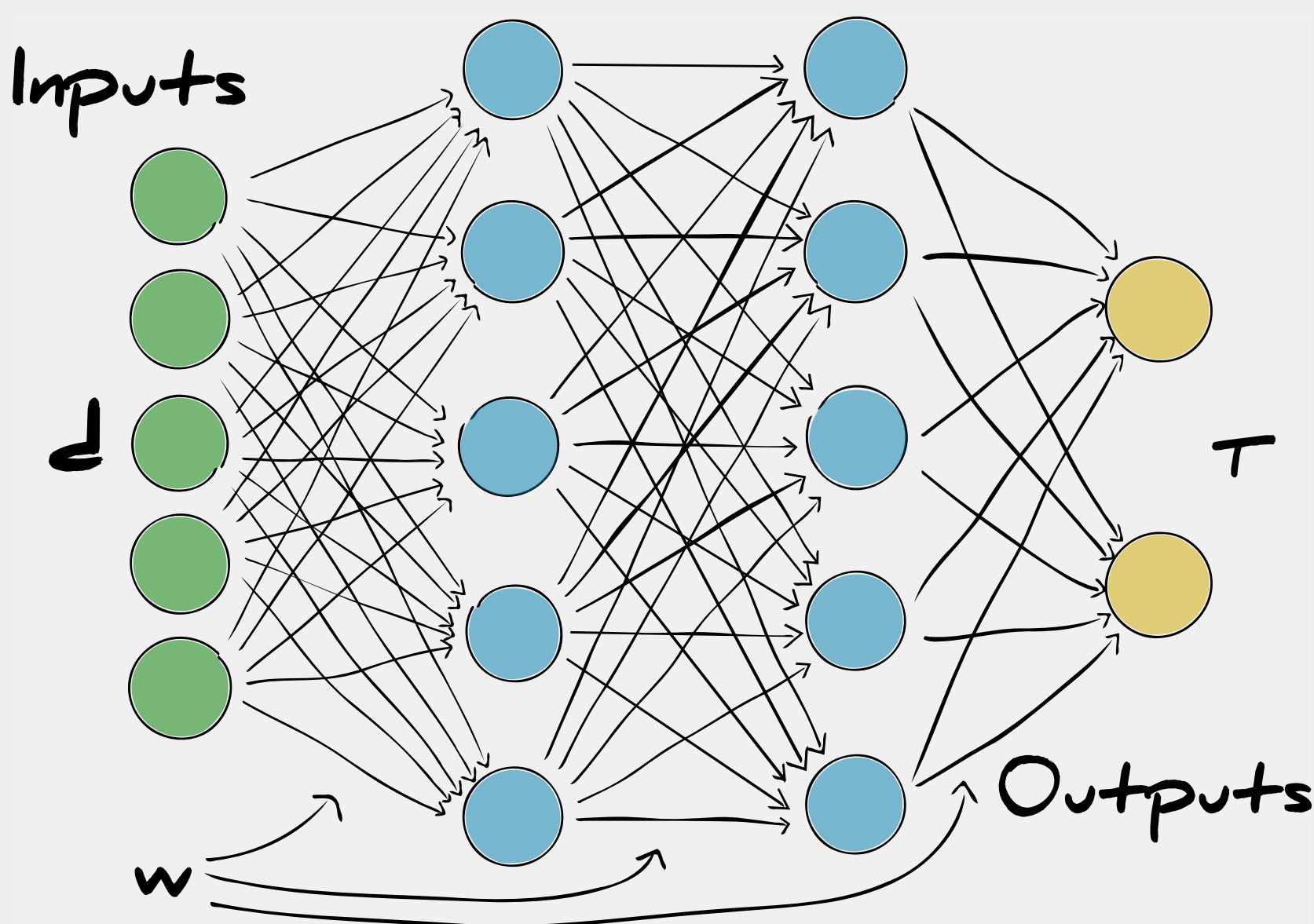
Tom Charnock

Institut d'Astrophysique de Paris



# What is a neural network?

# What is a neural network?



$$\text{NN}(w, \alpha) : \mathbf{d} \rightarrow \tau$$

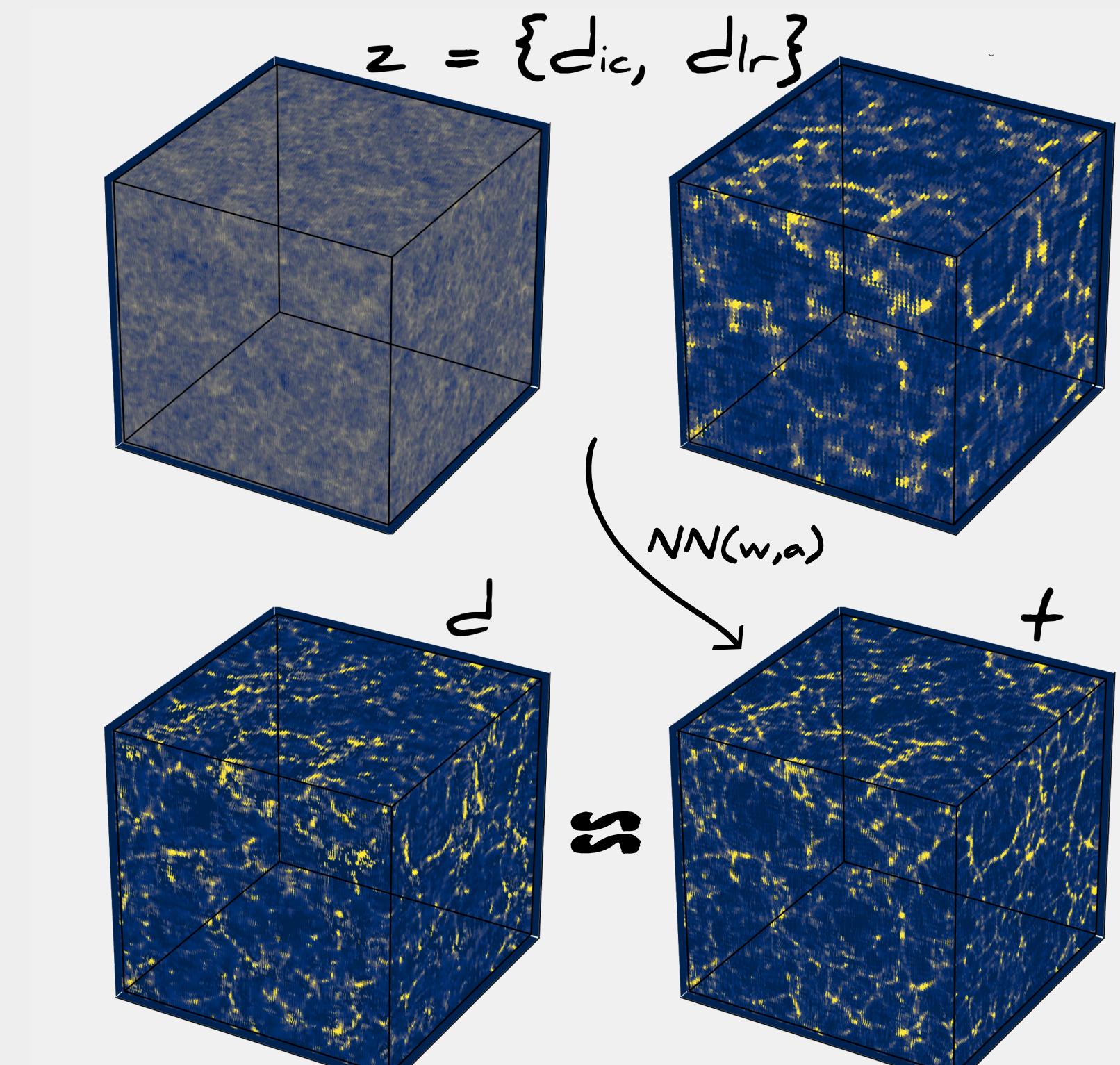
An approximation to a model,  $\mathcal{M} : \mathbf{d} \rightarrow \mathbf{t}$

**What are they used for?**

# Emulation (generative networks)

$$\text{NN}(w, \alpha) : z \rightarrow t$$

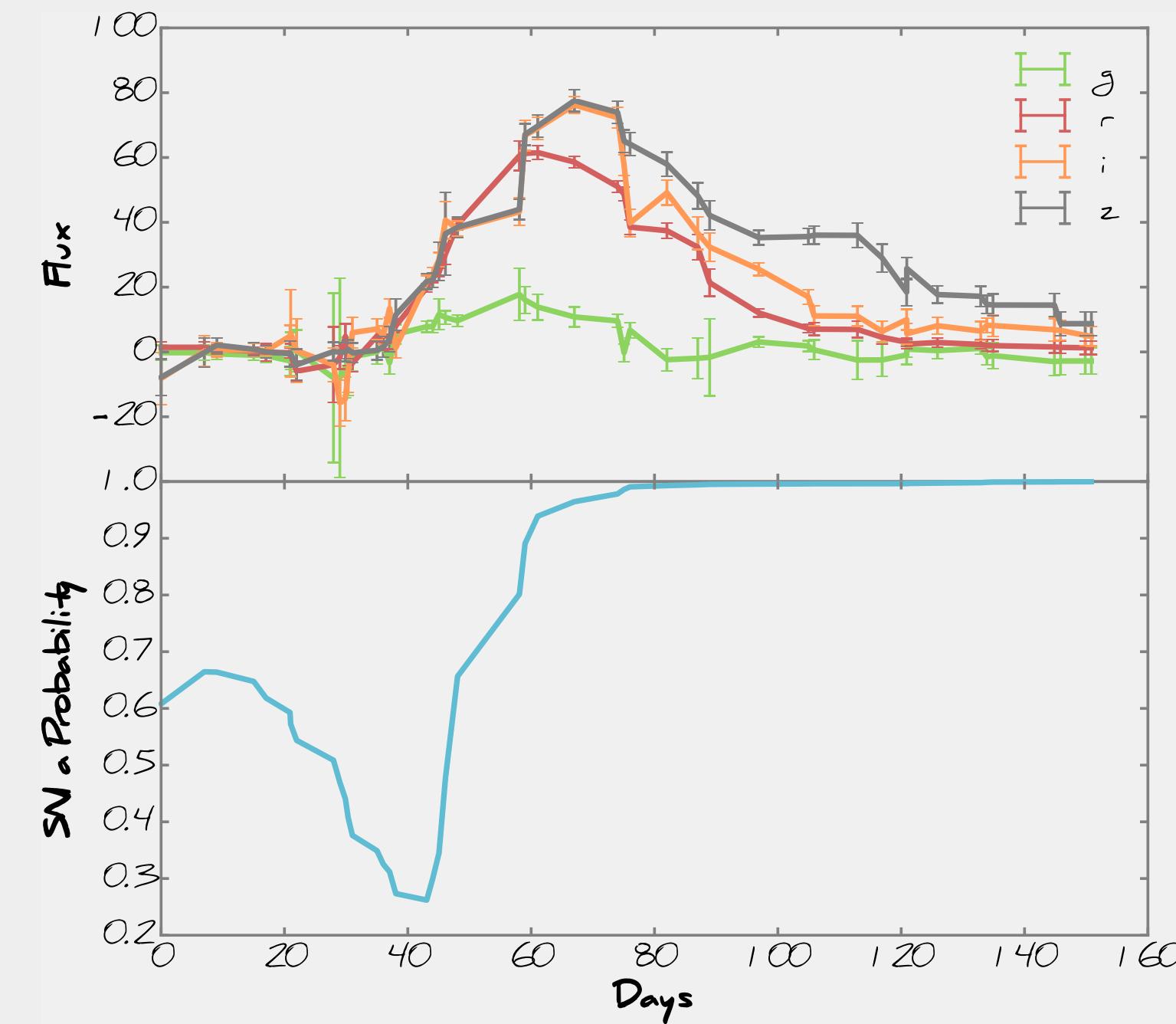
$$t \in \mathcal{P}(t|z) \simeq d \in \mathcal{P}(d)$$



# Classification

$$\text{NN}(w, \alpha) : \mathbf{d} \rightarrow \mathbf{t}$$

$$\mathbf{t} \in \mathcal{P}(\mathbf{t}) \simeq \mathbf{d} \in \mathcal{P}(\text{type}|\mathbf{d})$$



**Extremely popular for cosmological parameter estimation**

**Unfortunately...**

**Unfortunately...**

**None of these work in a scientifically rigourous way!**

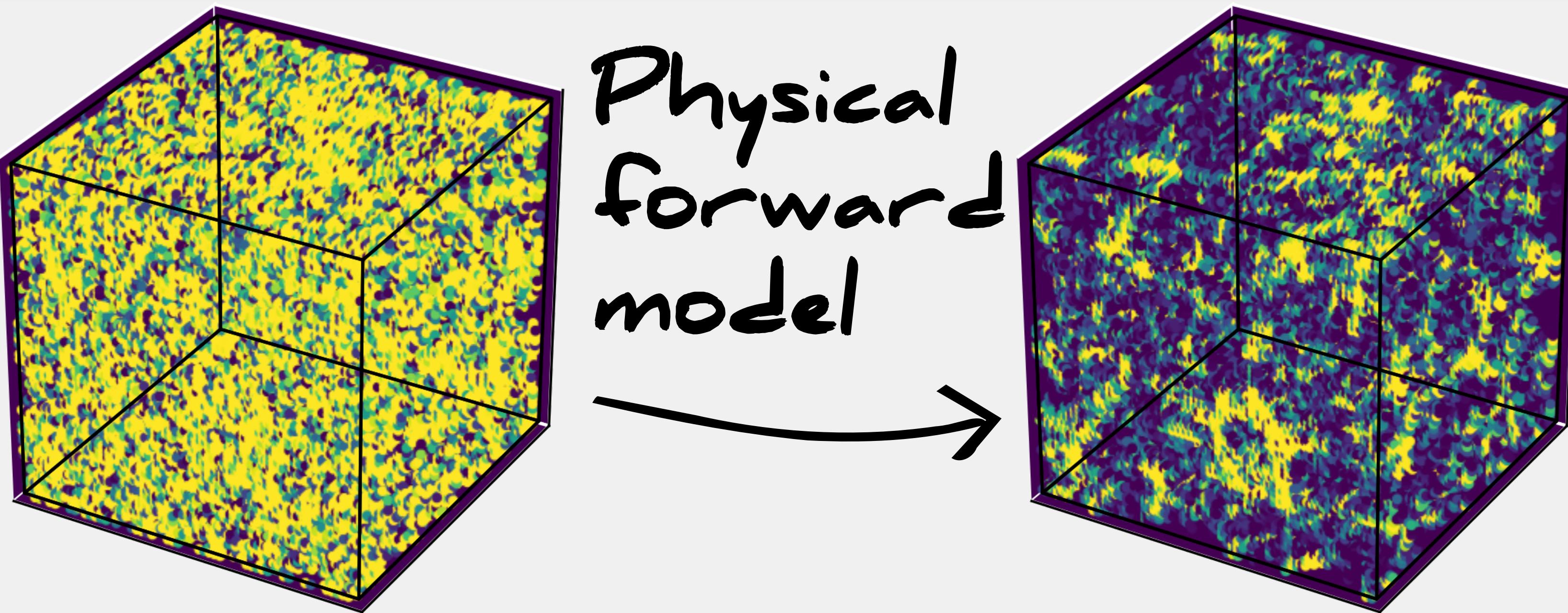
**To see why we need to do some stats**

A physical model  $\mathcal{M} : \theta \rightarrow \mathbf{d}$

The generator of data ( $\mathbf{d}$ ) with parameters ( $\theta$ )

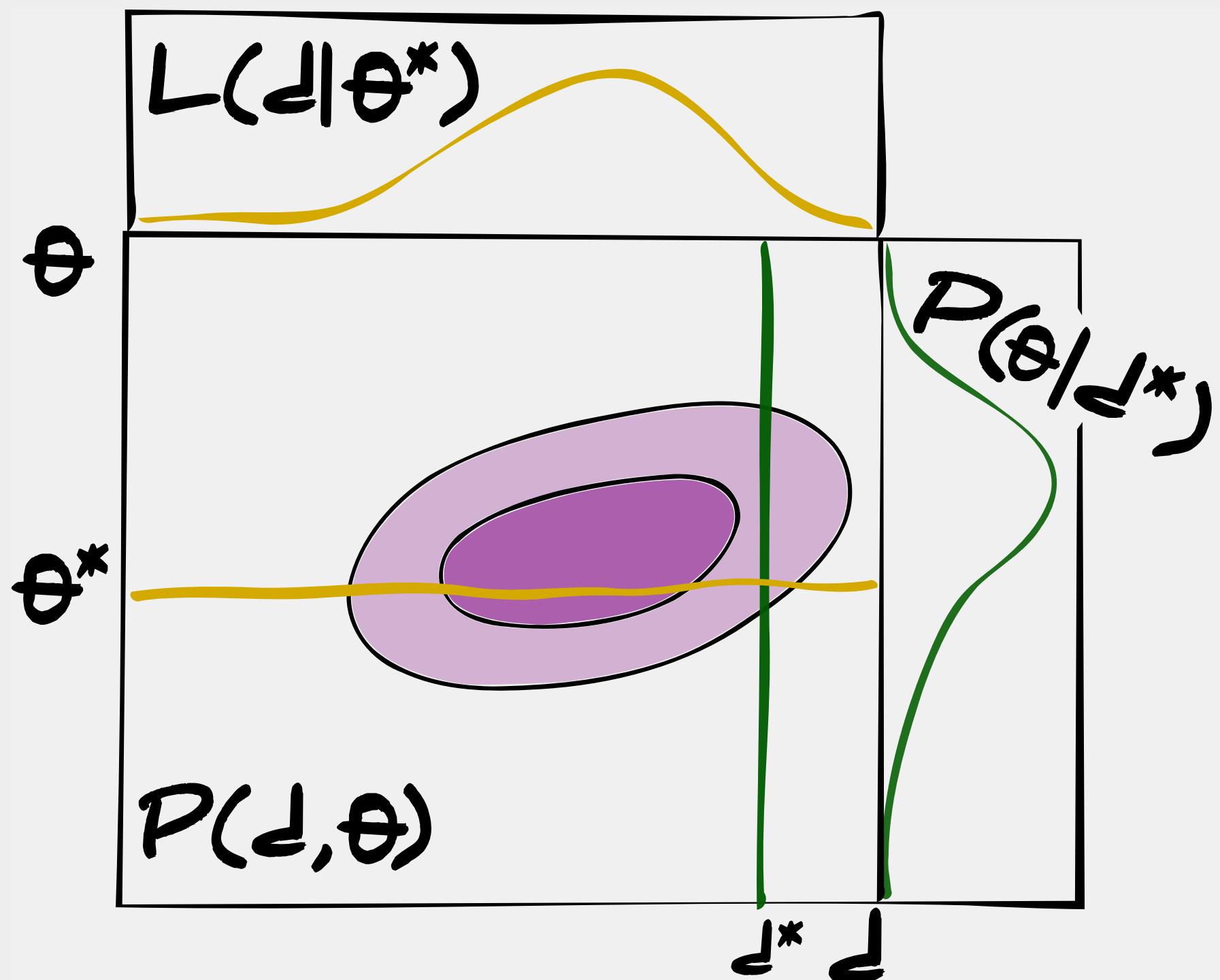
A physical model  $\mathcal{M} : \theta \rightarrow \mathbf{d}$

The generator of data ( $\mathbf{d}$ ) with parameters ( $\theta$ )



# The likelihood and the posterior

# The likelihood and the posterior



**Likelihood  $\mathcal{L}(d|\theta^*)$**

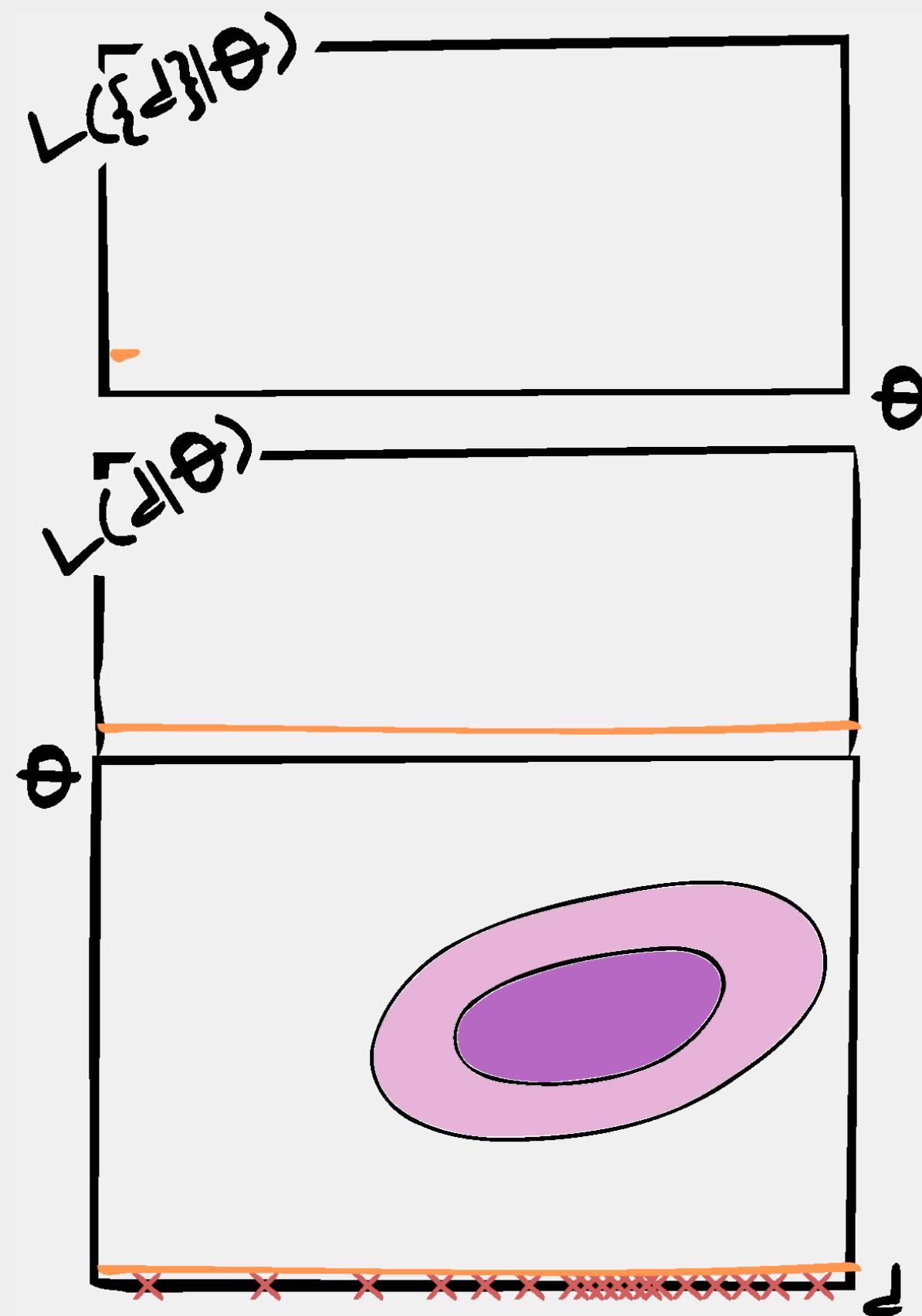
What is the probability that the model generates data  $d$  given a set of model parameters  $\theta^*$ ?

**Posterior  $\mathcal{P}(\theta|d^*)$**

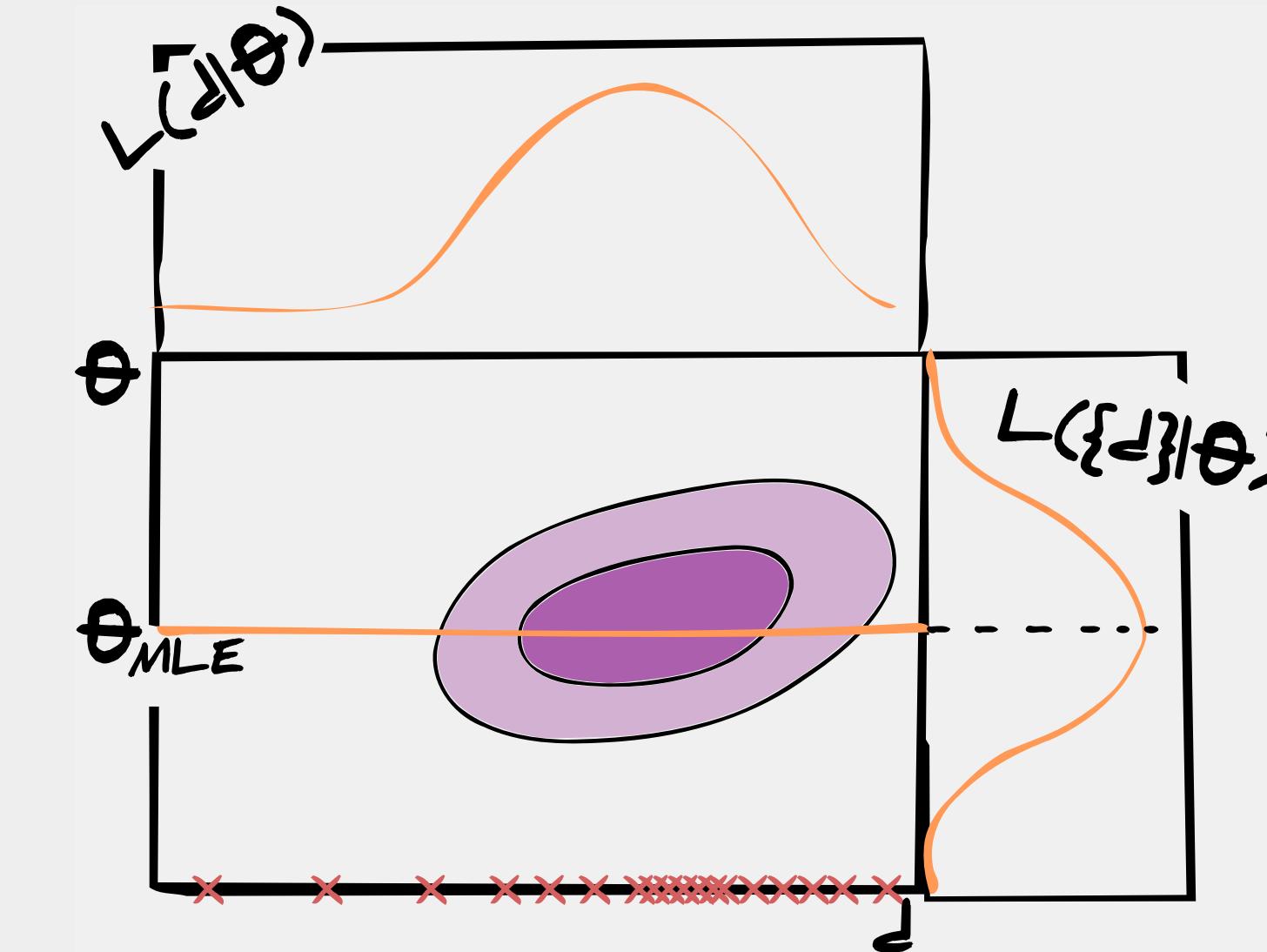
What is the probability that model parameters  $\theta$  generate the data  $d^*$ ?

# Maximum likelihood estimates

# Maximum likelihood estimates



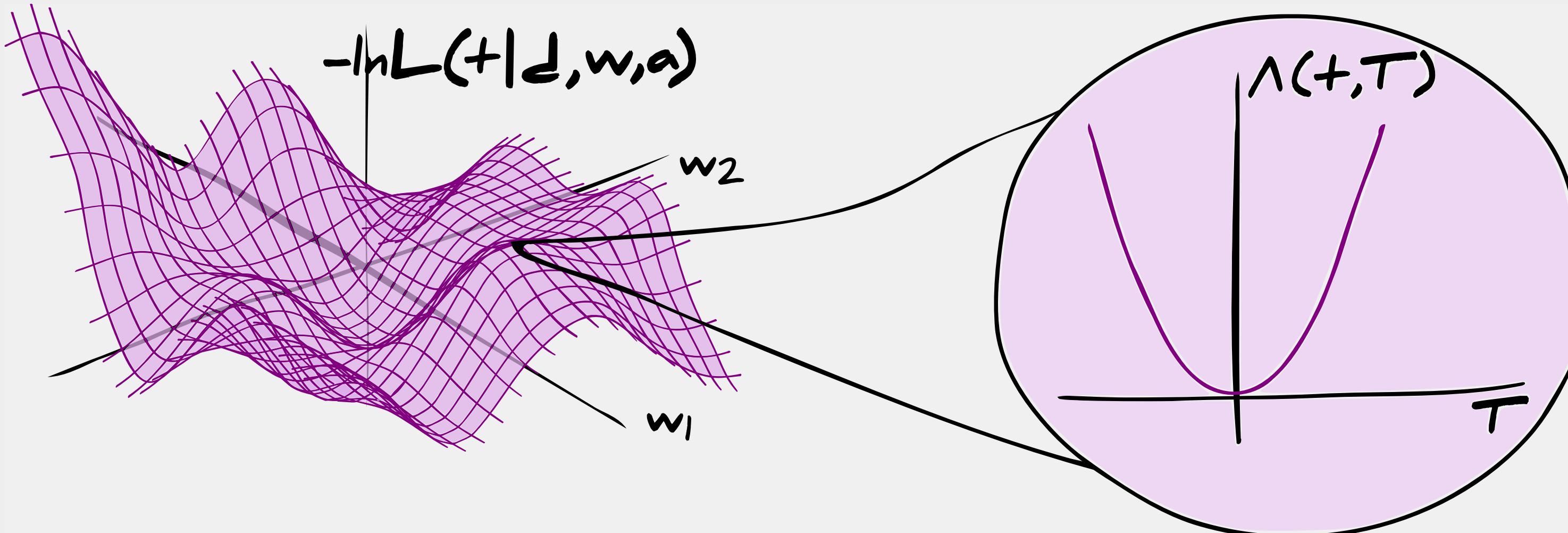
Which value of the model parameters  $\theta$  is most likely to generate the observed population of data  $\{d\}$ ?



$$L(\{d\}|\theta_{MLE}) = \underset{\theta \in \Theta}{\text{Sup}} L(\{d\}|\theta)$$

# Training a network

# Training a network



$\mathcal{M} : \mathbf{d} \rightarrow \mathbf{t}$  approximator

$\Lambda(\mathbf{t}, \tau)$   
smooth and convex

$\text{NN}(w, \alpha) : \mathbf{d} \rightarrow \tau$

$\mathcal{L}(\mathbf{t}|\mathbf{d}, w, \alpha)$

Cost function and likelihood

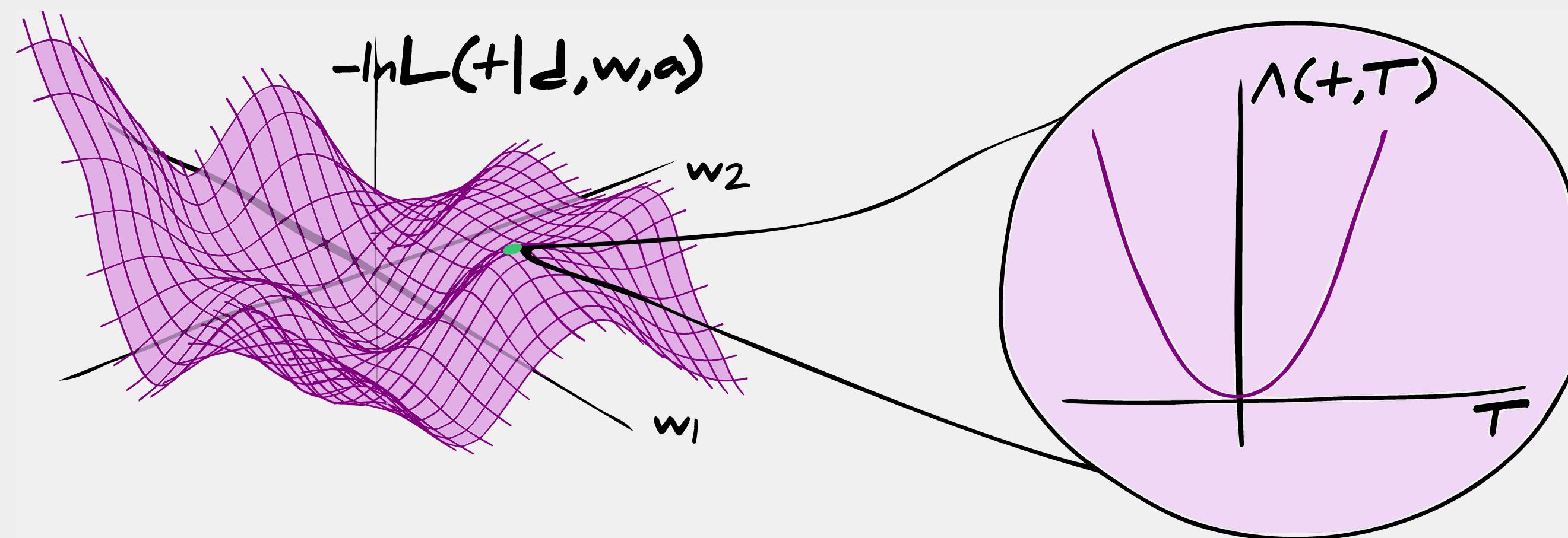
complex, discrete and  
possibly non-convex in  $w$   
and  $\alpha$

$$\Lambda(\mathbf{t}, \tau) = -\ln \mathcal{L}(\mathbf{t}|\mathbf{d}, w^*, \alpha^*)$$

# Optimising a network

# Optimising a network

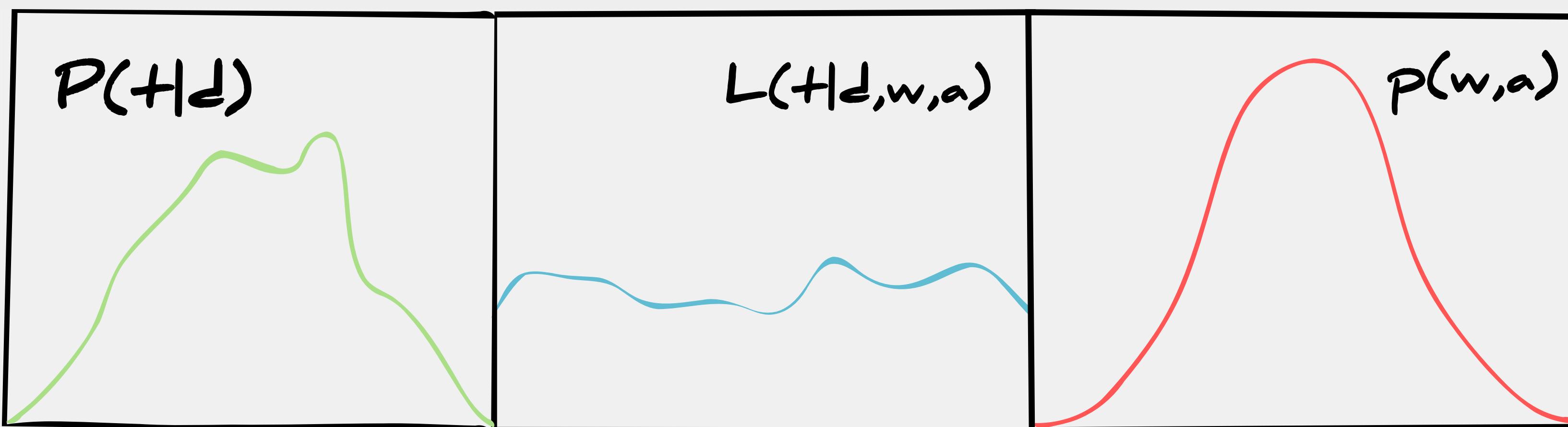
Given a set of training data  $\{\mathbf{d}_i^{\text{train}}, \mathbf{t}_i^{\text{train}} | i \in [1, n_{\text{train}}]\}$  what is the maximum likelihood estimate of the weights?



$$w^{\text{MLE}} = \underset{w}{\operatorname{argmax}} \left[ \mathcal{L}(\{\mathbf{t}_i^{\text{train}} | i \in [1, n_{\text{train}}]\}) | \{\mathbf{d}_i^{\text{train}} | i \in [1, n_{\text{train}}]\} \right]$$

**Where the problem starts**

# The posterior predictive density $\mathcal{P}(\mathbf{t}|\mathbf{d})$

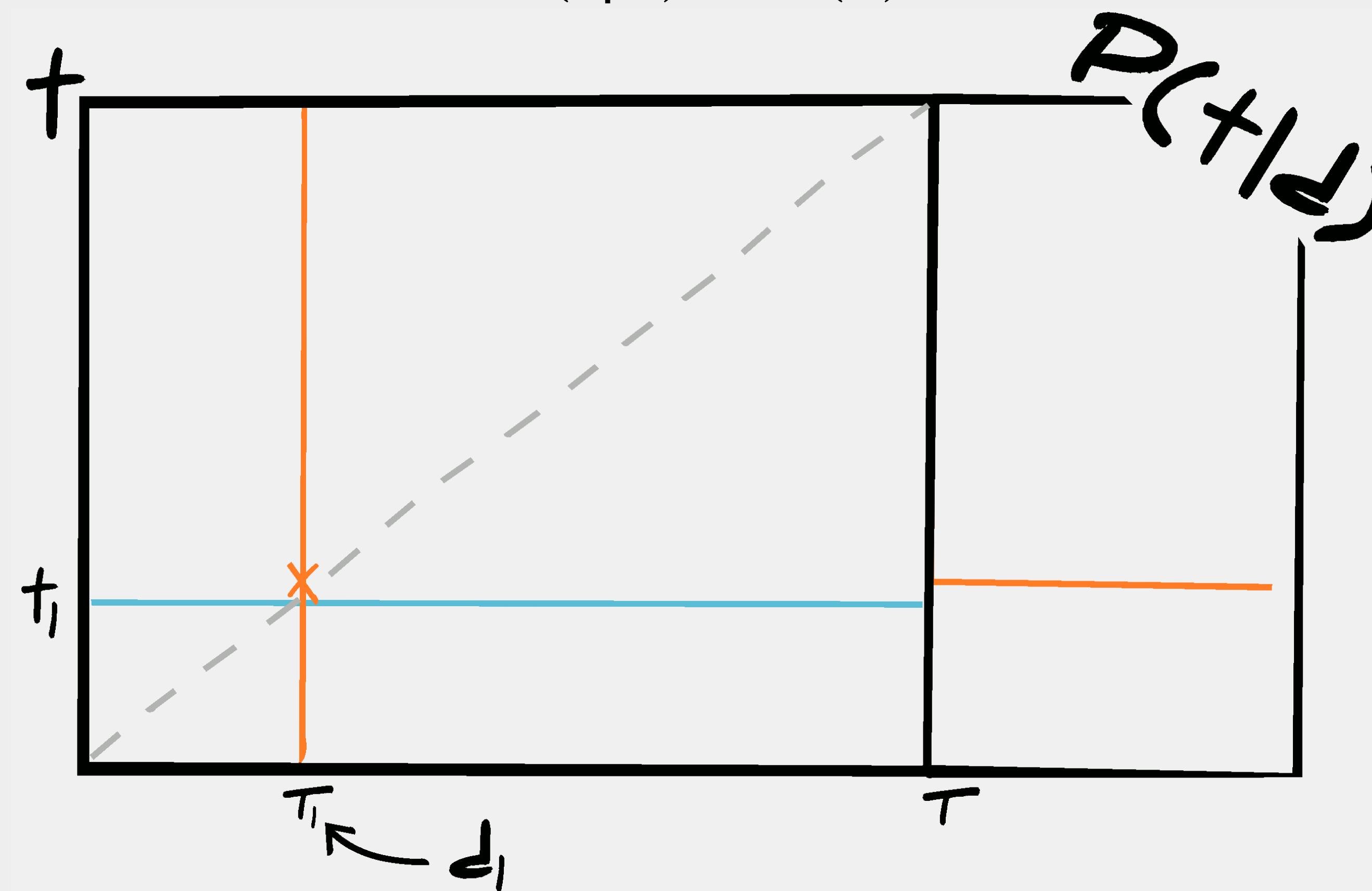


$$\mathcal{P}(\mathbf{t}|\mathbf{d}) = \int d\mathbf{w}\alpha \mathcal{L}(\mathbf{t}|\mathbf{d}, \mathbf{w}, \alpha) p(\mathbf{w}, \alpha)$$

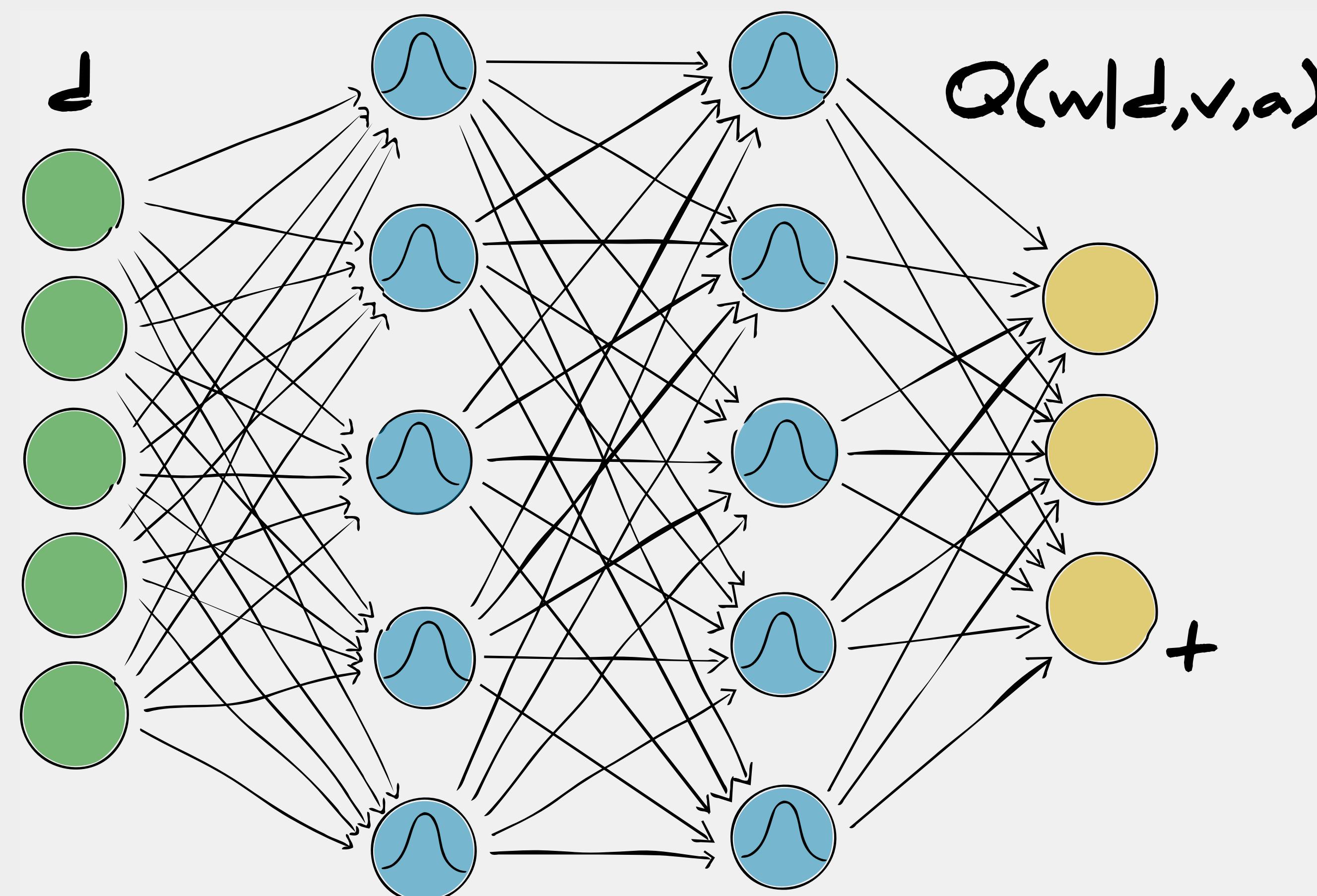
# Failure of traditional networks

$$p(w, \alpha) = \delta(w - w^{\text{local MLE}}, \alpha - \alpha^*)$$

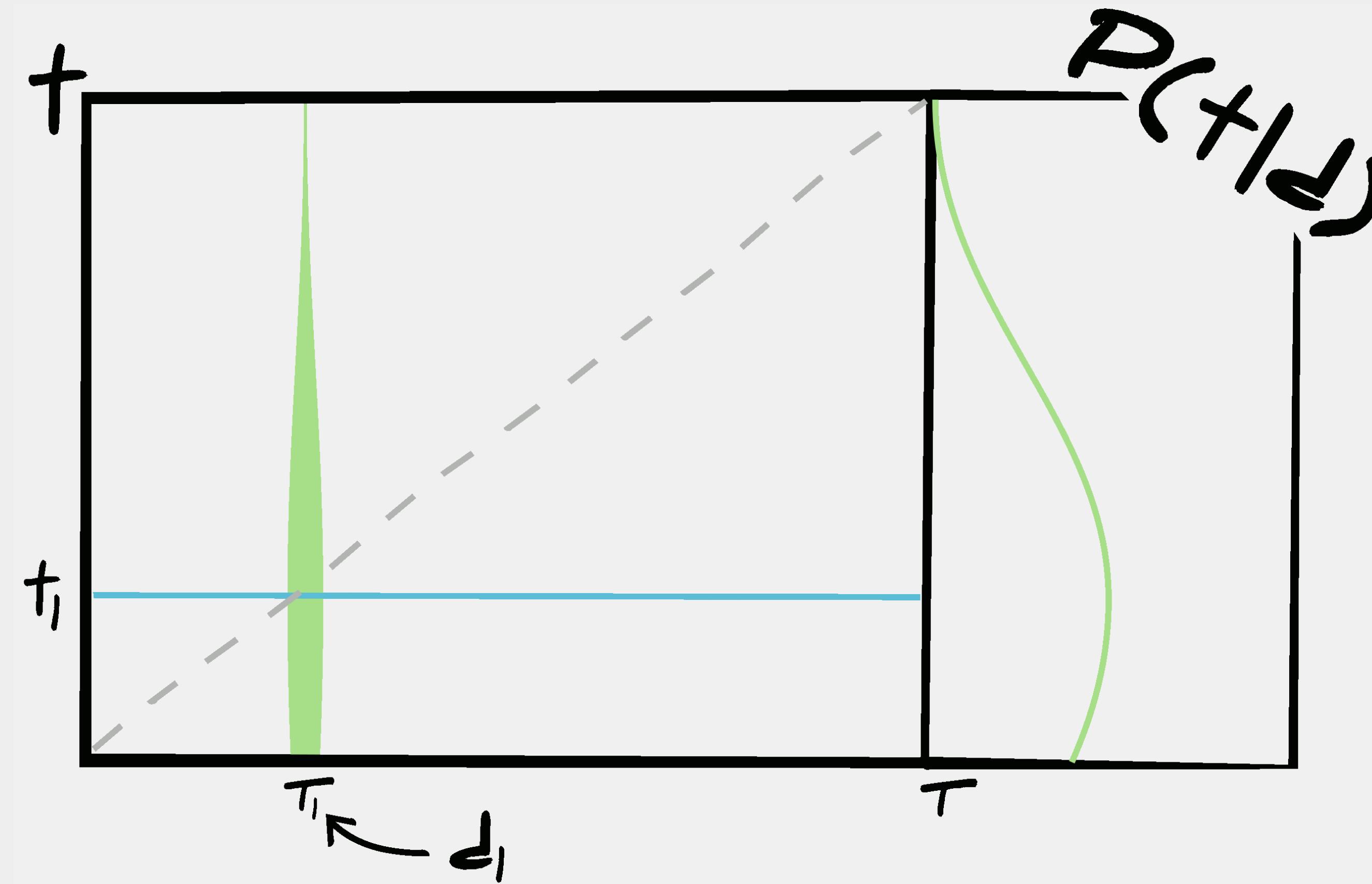
$$\mathcal{P}(t|d) = \delta(\tau)$$



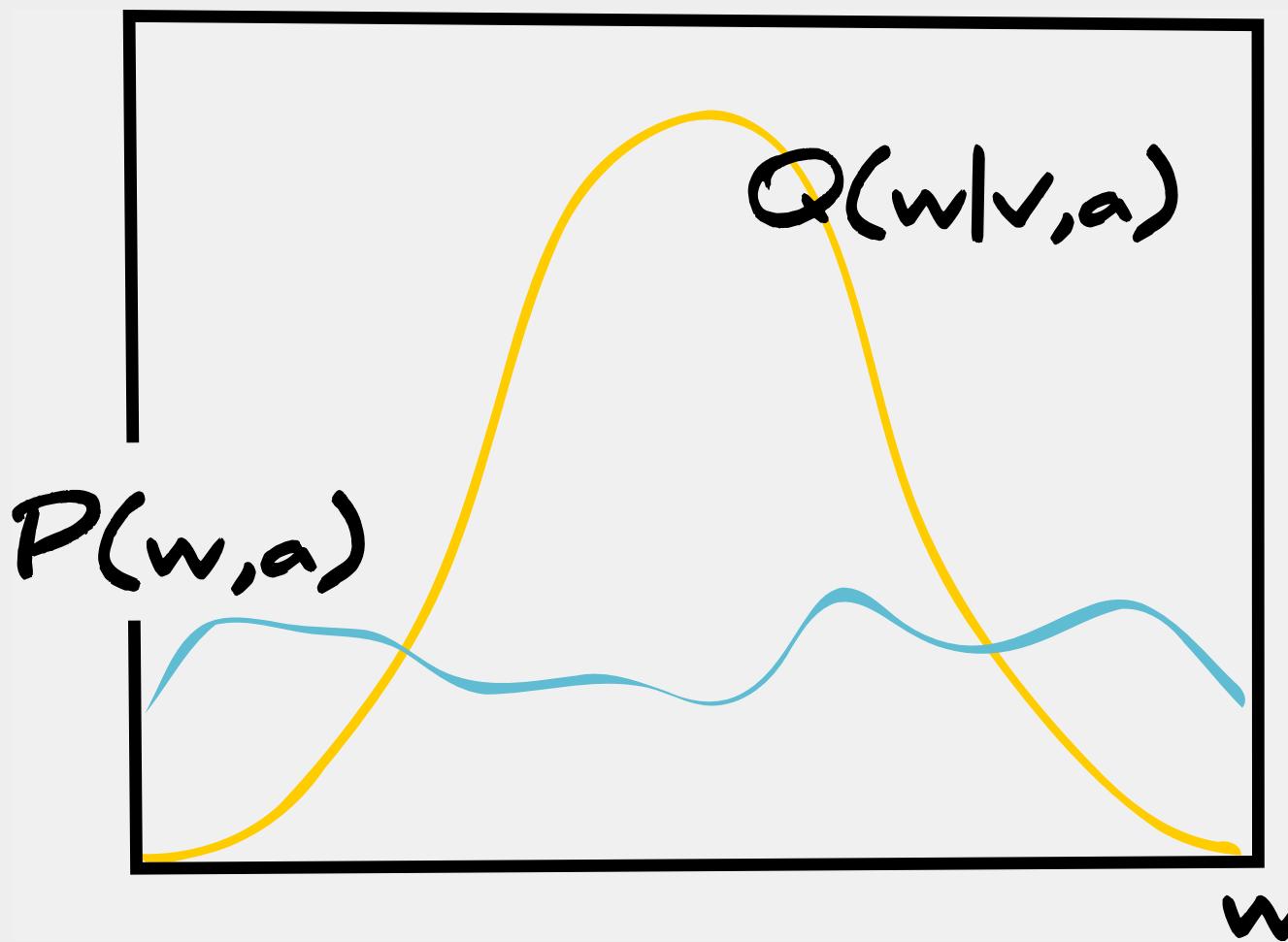
# Variational inference



$$\mathcal{P}(\mathbf{t}|\mathbf{d}) = \int dw dv d\alpha \mathcal{L}(\mathbf{t}|\mathbf{d}, w, \alpha) Q(w|v, \alpha, \{\mathbf{d}_i^{\text{train}}, \mathbf{t}_i^{\text{train}} | i \in [1, n_{\text{train}}]\}) p(v, \alpha)$$



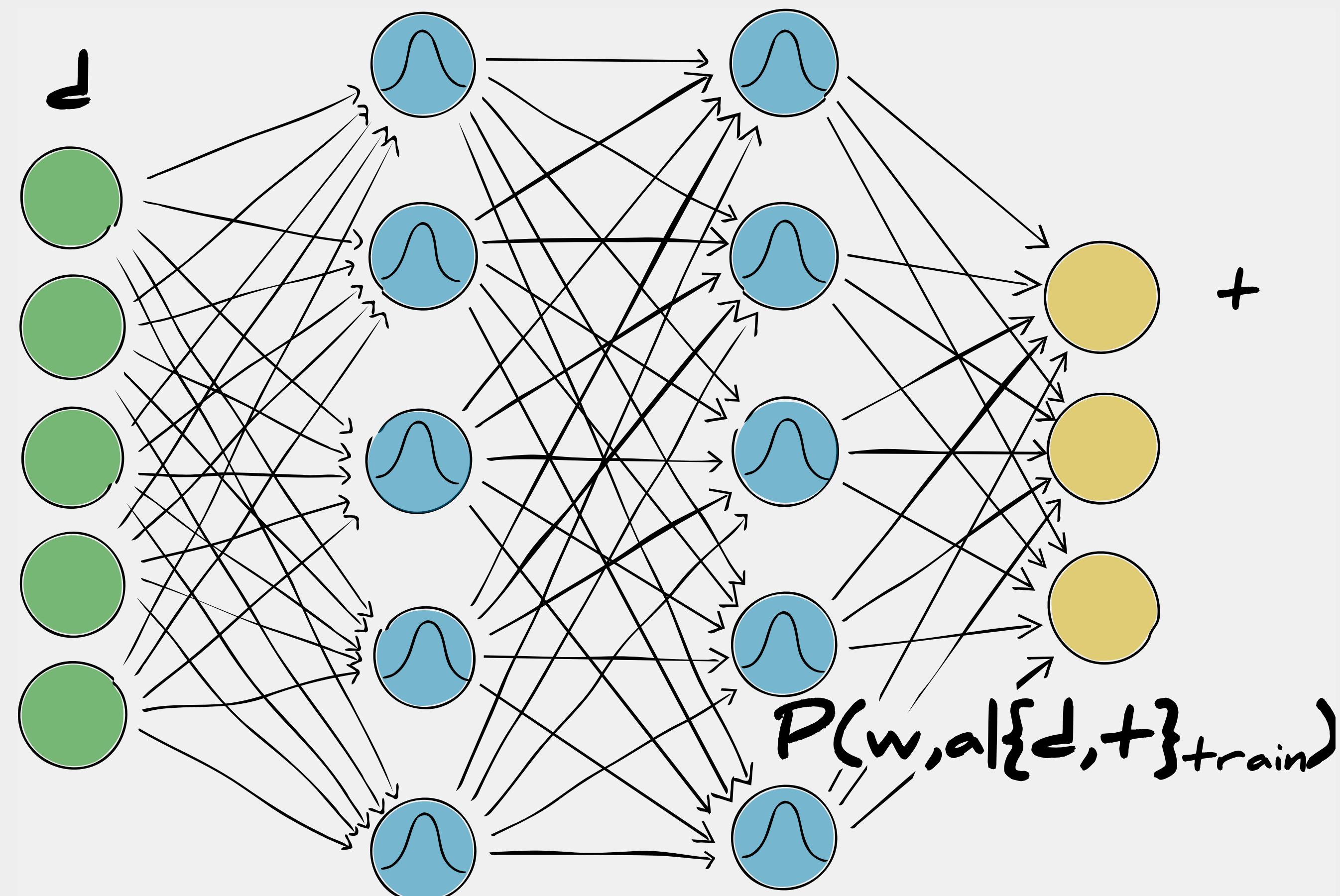
Still depends on training the complex likelihood surface and choice of variational distribution



$$\begin{aligned}
 P(\mathbf{t}|\mathbf{d}) &= \int dw dv d\boldsymbol{\alpha} \mathcal{L}(\mathbf{t}|\mathbf{d}, w, \boldsymbol{\alpha}) Q(w|v, \boldsymbol{\alpha}, \{\mathbf{d}_i^{\text{train}}, \mathbf{t}_i^{\text{train}} | i \in [1, n_{\text{train}}]\}) \\
 &\quad \times \delta(v - v^{\text{local MLE}}, \boldsymbol{\alpha} - \boldsymbol{\alpha}^*) \\
 &= \int dw \mathcal{L}(\mathbf{t}|\mathbf{d}, w, \boldsymbol{\alpha}^*) Q(w|v^{\text{local MLE}}, \boldsymbol{\alpha}^*, \{\mathbf{d}_i^{\text{train}}, \mathbf{t}_i^{\text{train}} | i \in [1, n_{\text{train}}]\}).
 \end{aligned}$$

# Bayesian neural networks

# Bayesian neural networks



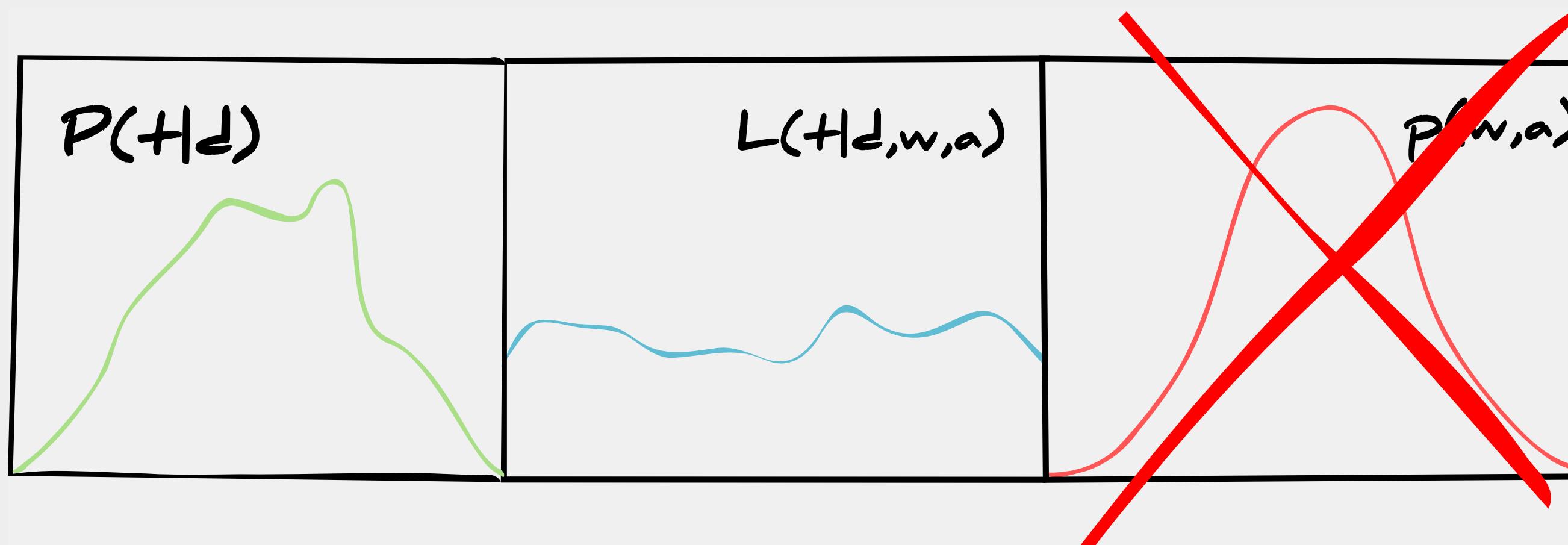
$$\begin{aligned}\mathcal{P}(\mathbf{t}|\mathbf{d}) &= \int d\mathbf{w}d\boldsymbol{\alpha} \mathcal{L}(\mathbf{t}|\mathbf{d}, \mathbf{w}, \boldsymbol{\alpha}) \mathcal{P}(\mathbf{w}, \boldsymbol{\alpha} | \{\mathbf{d}_i^{\text{train}}, \mathbf{t}_i^{\text{train}} | i \in [1, n_{\text{train}}]\}) \\ &= \int d\mathbf{w}d\boldsymbol{\alpha} \mathcal{L}(\mathbf{t}|\mathbf{d}, \mathbf{w}, \boldsymbol{\alpha}) \prod_i^{n_{\text{train}}} \mathcal{L}(\mathbf{t}_i^{\text{train}} | \mathbf{d}_i^{\text{train}}, \mathbf{w}, \boldsymbol{\alpha}) p(\mathbf{w}, \boldsymbol{\alpha}).\end{aligned}$$

# Still dependent on the training data!

Classical network :  $p(w, \alpha) = \delta(w - w^{\text{MLE}}, \alpha - \alpha^*)$

Variational inference :  $p(w, \alpha) = Q(w | v^{\text{local MLE}}, \alpha^*, \{\mathbf{d}_i^{\text{train}}, \mathbf{t}_i^{\text{train}} | i \in [1, n_{\text{train}}]\})$

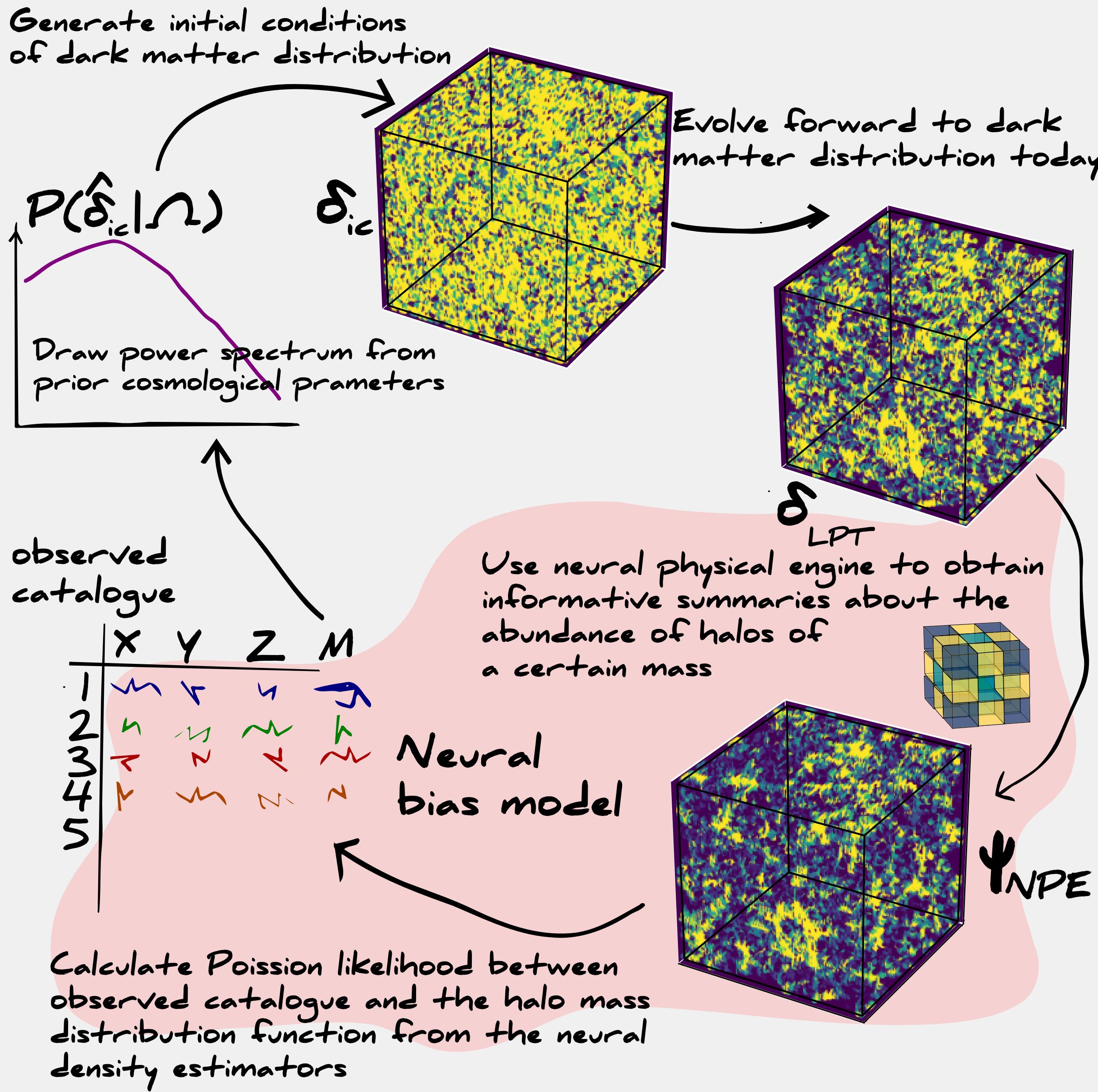
Bayesian neural networks :  $p(w, \alpha) = \prod_i^{n_{\text{train}}} \mathcal{L}(\mathbf{t}_i^{\text{train}} | \mathbf{d}_i^{\text{train}}, w, \alpha) p(w, \alpha)$



**How can we use a neural network  
then?**

**Build it into the physical model**

# **Method 1 : Model extension**



**We infer the neural network**

**HMCLET**

# Hamiltonian Monte Carlo

Introduce momentum  $\mathbf{p} \leftarrow \mathcal{N}(\mathbf{0}, \mathbf{M})$  and solve

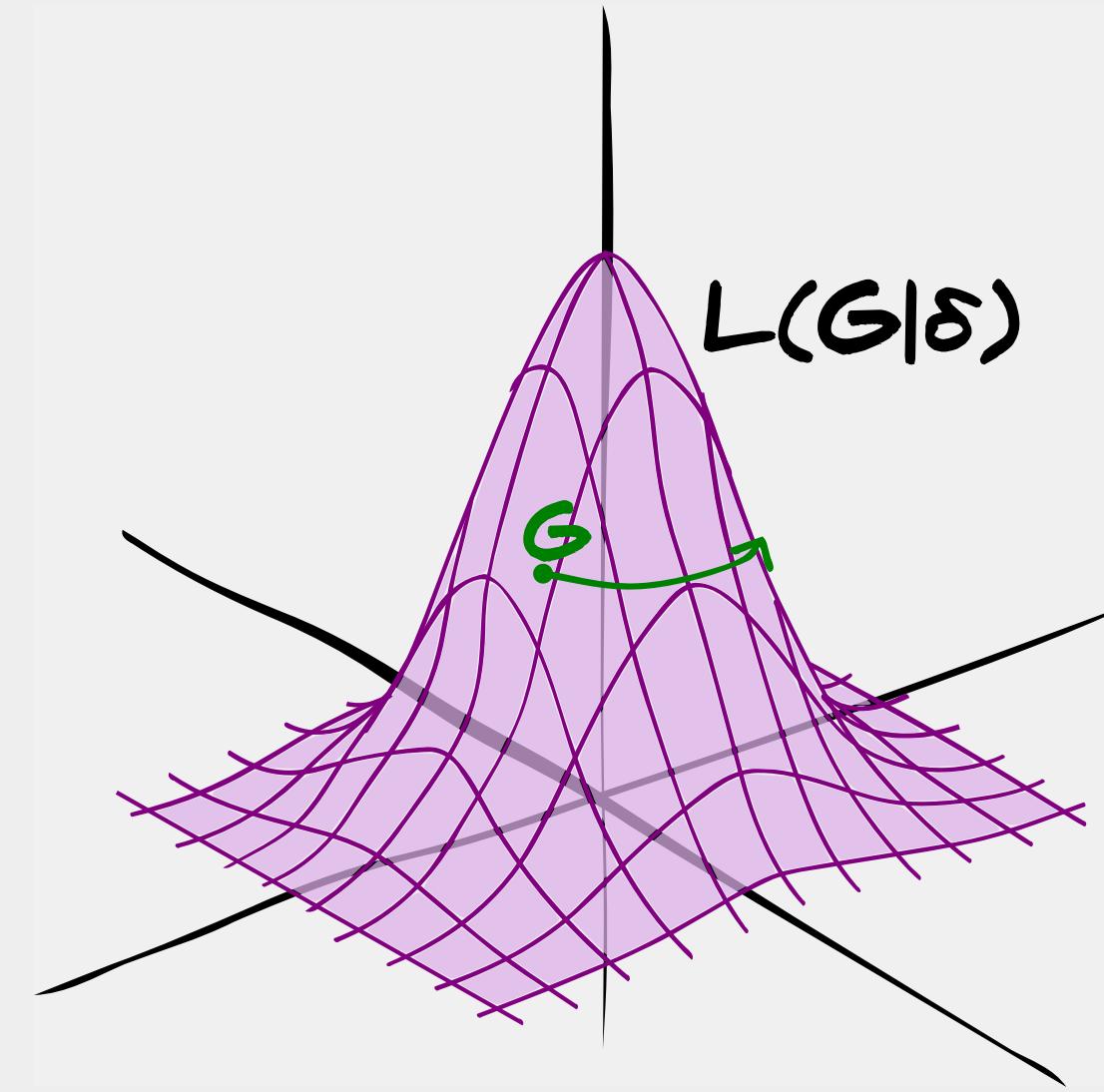
Hamilton's equations

$$\begin{aligned}\mathcal{H}(\theta, \mathbf{p}) &= \mathcal{V}(\theta) + \mathcal{K}(\mathbf{p}) \\ &= \mathcal{L}(\theta|\delta) - \log[\pi(\theta)] + \frac{1}{2} \mathbf{p}^T \mathbf{M}^{-1} \mathbf{p}\end{aligned}$$

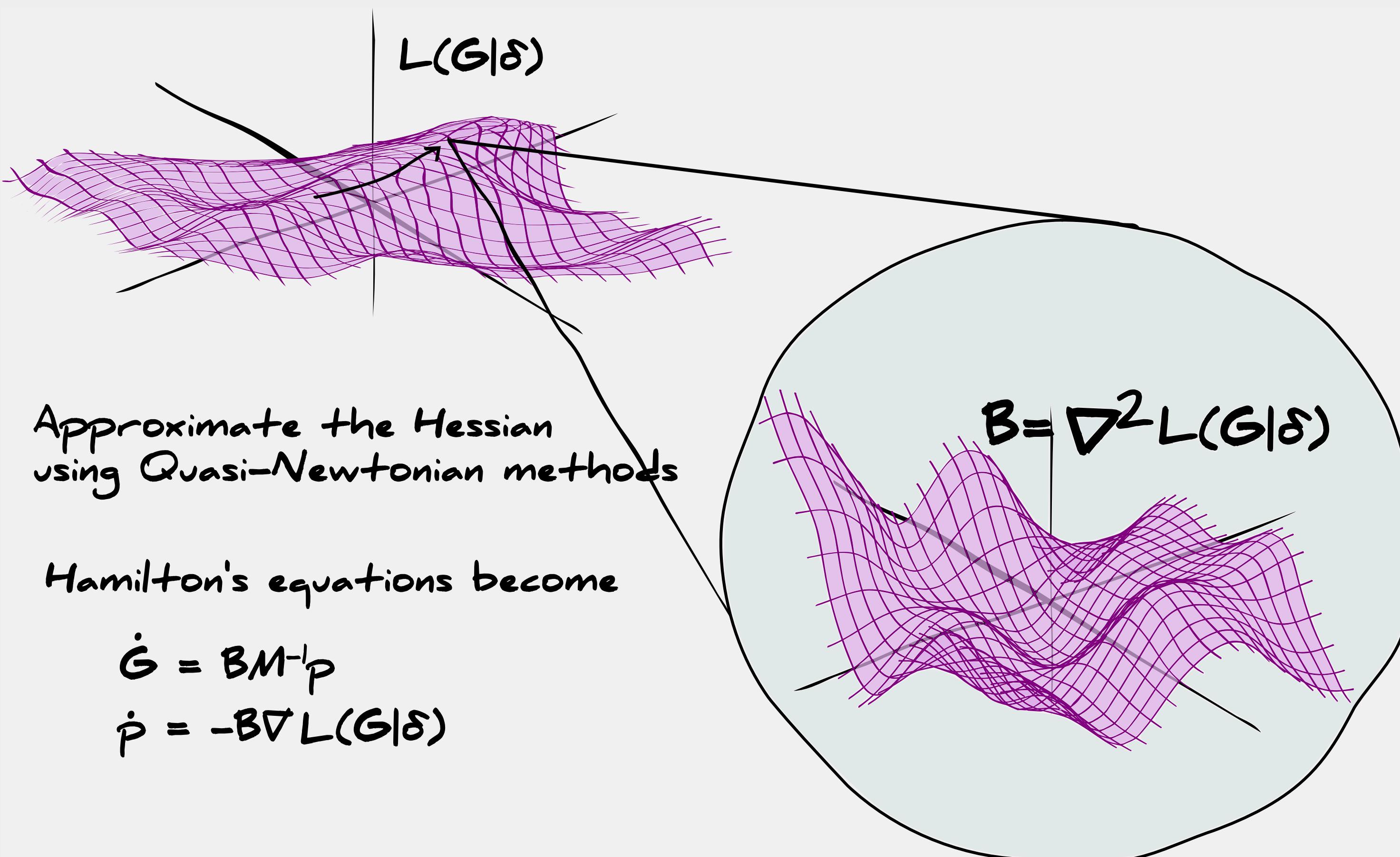
Accept samples according to probability

$$\alpha = \text{Min} [\exp(\Delta\mathcal{H}), 1]$$

Evolve using  $\dot{\theta} = \mathbf{M}^{-1} \mathbf{p}$  and  $\dot{\mathbf{p}} = -\nabla \mathcal{V}(\theta)$



# Use second order geometric information

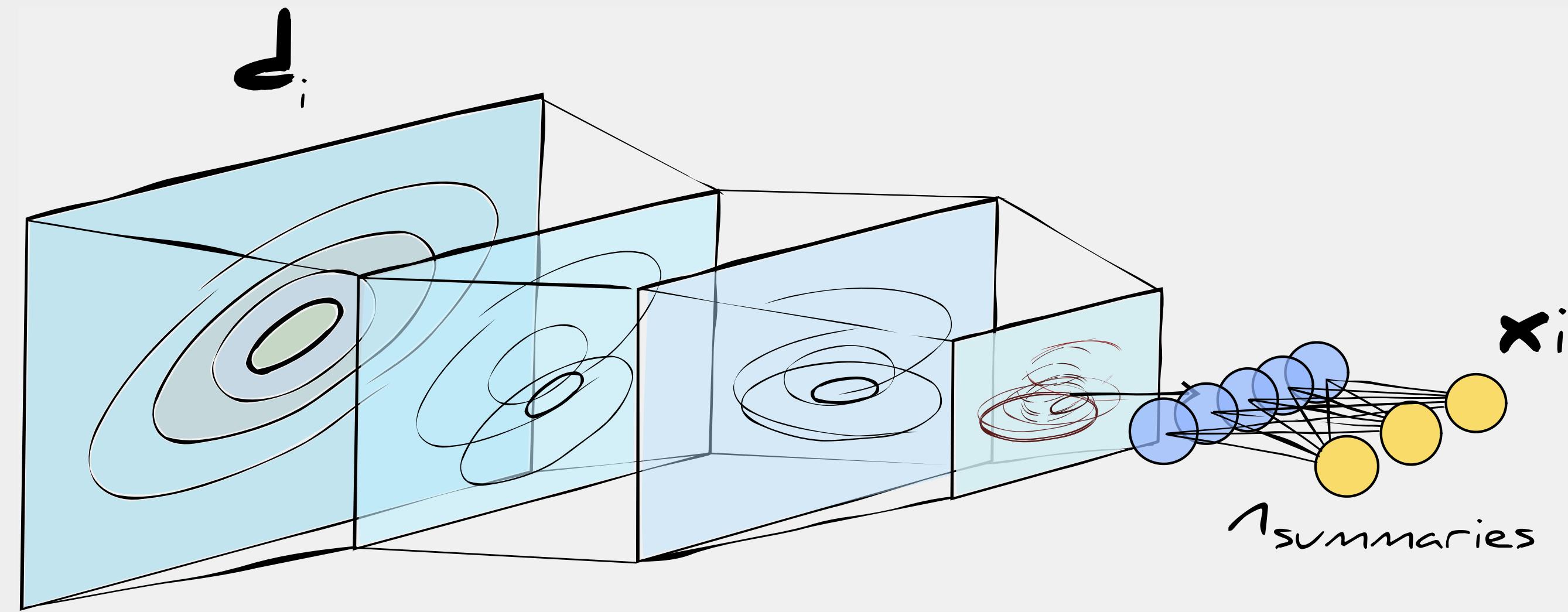


# Neural physical engines

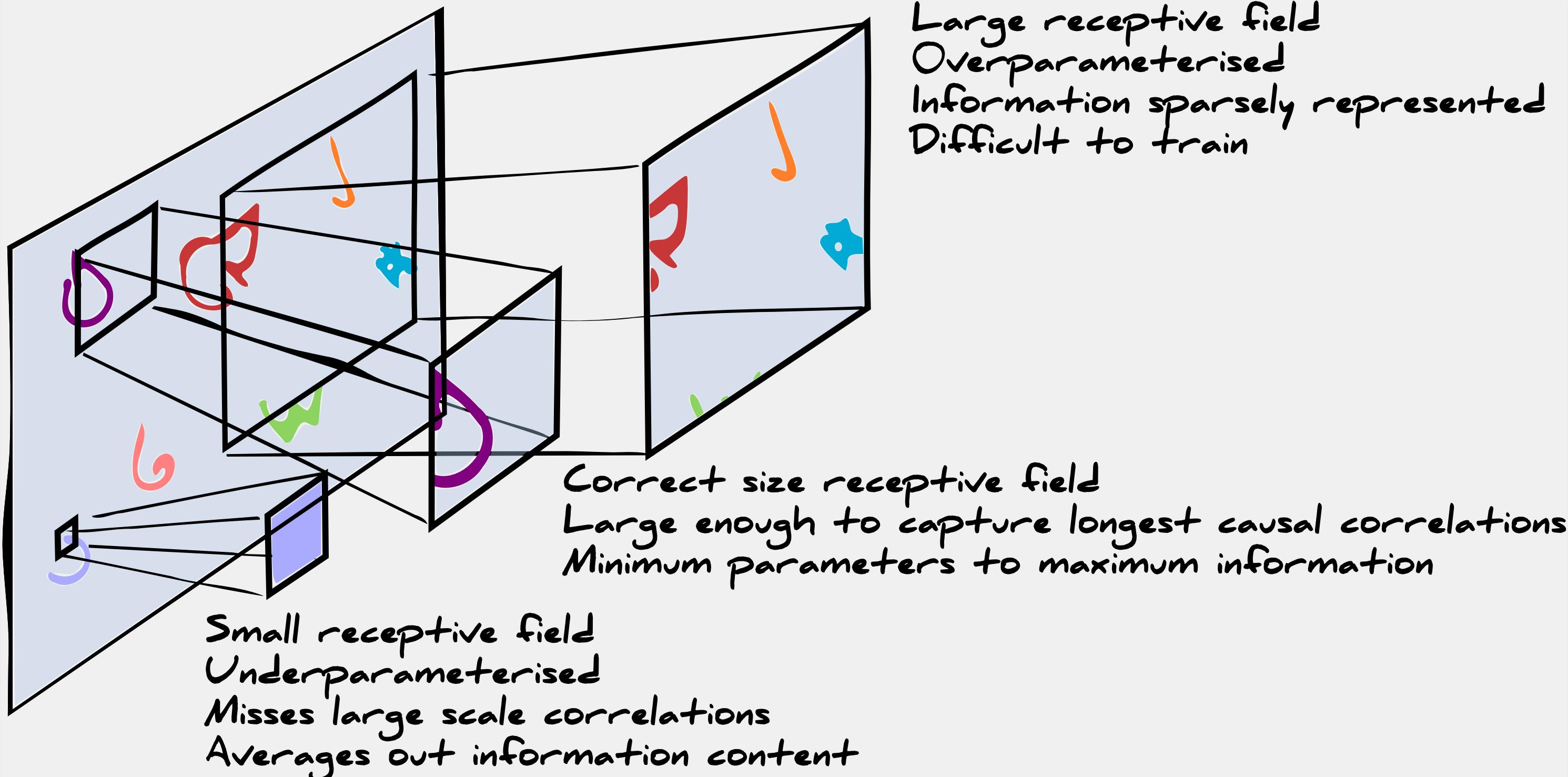
Build networks using physical principles.

- Reduces number of parameters
- Increases computational efficiency
- Decreases overfitting
- Improves interpretability

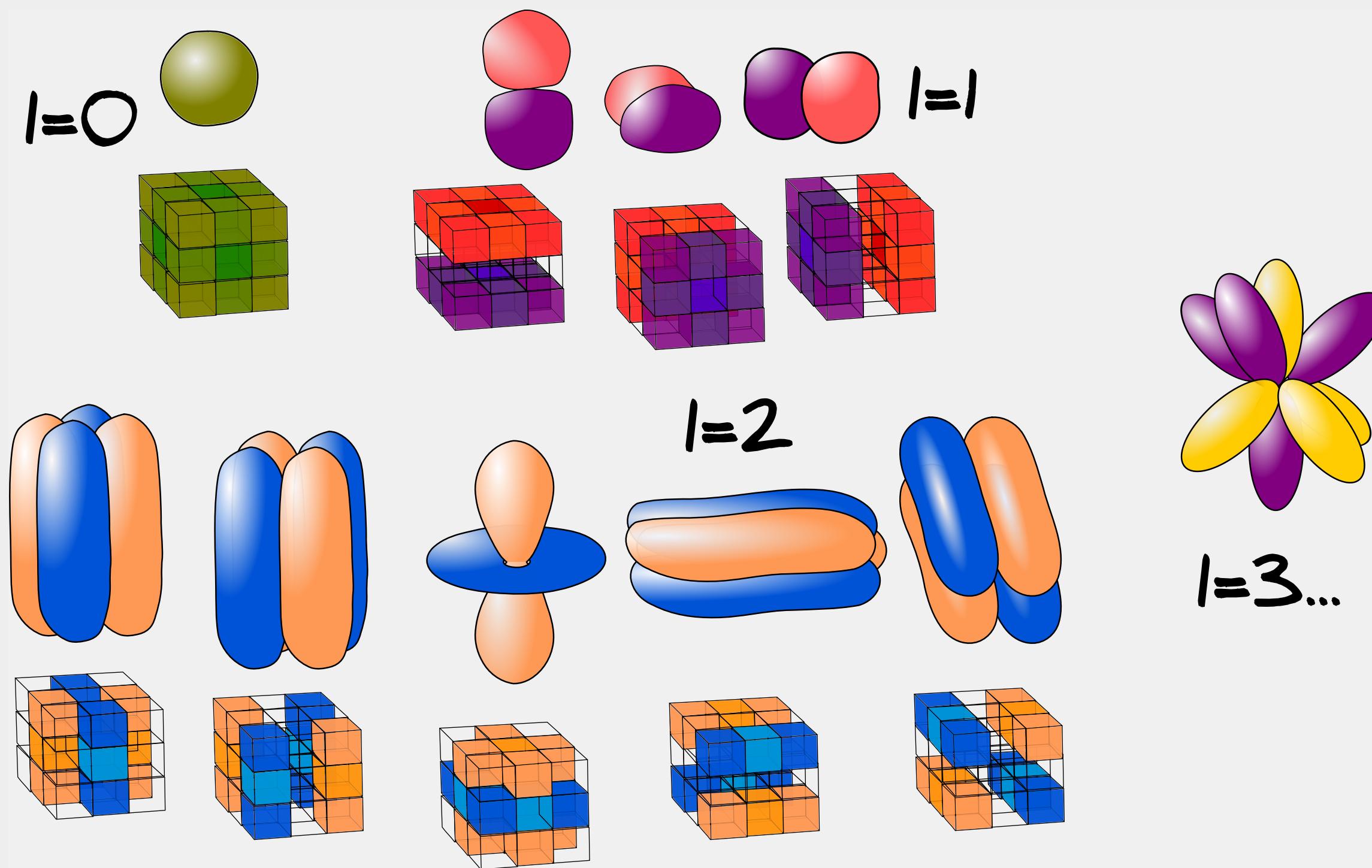
# Convolutions for translational invariance



# Use only causally relevant features

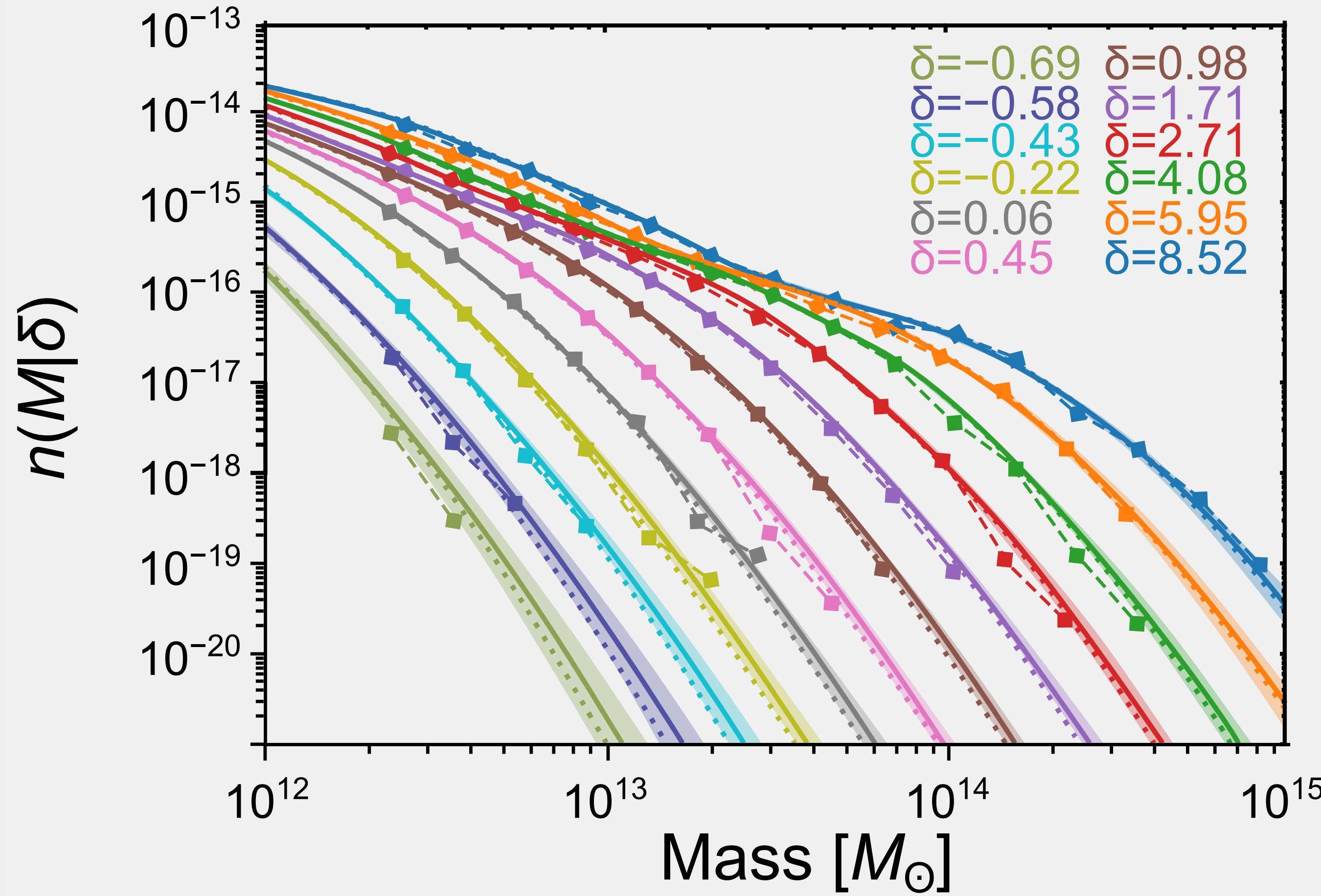


# Use only informative features



[github:multipole kernels](https://github.com/multipole/kernels)

# Halo mass distribution function

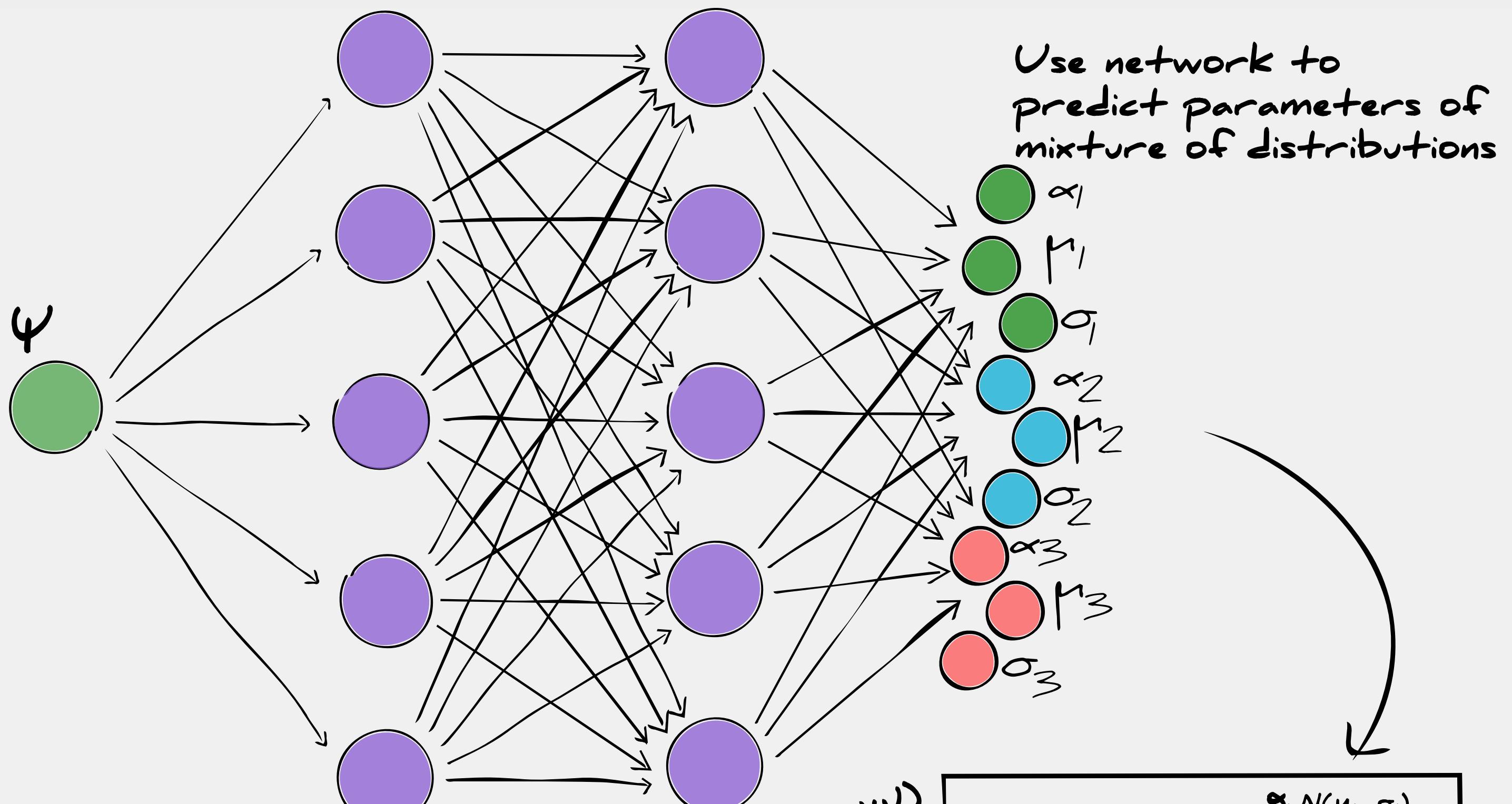


# Neural density estimators

**Halo mass distribution function is a smooth function of mass given a density environment**

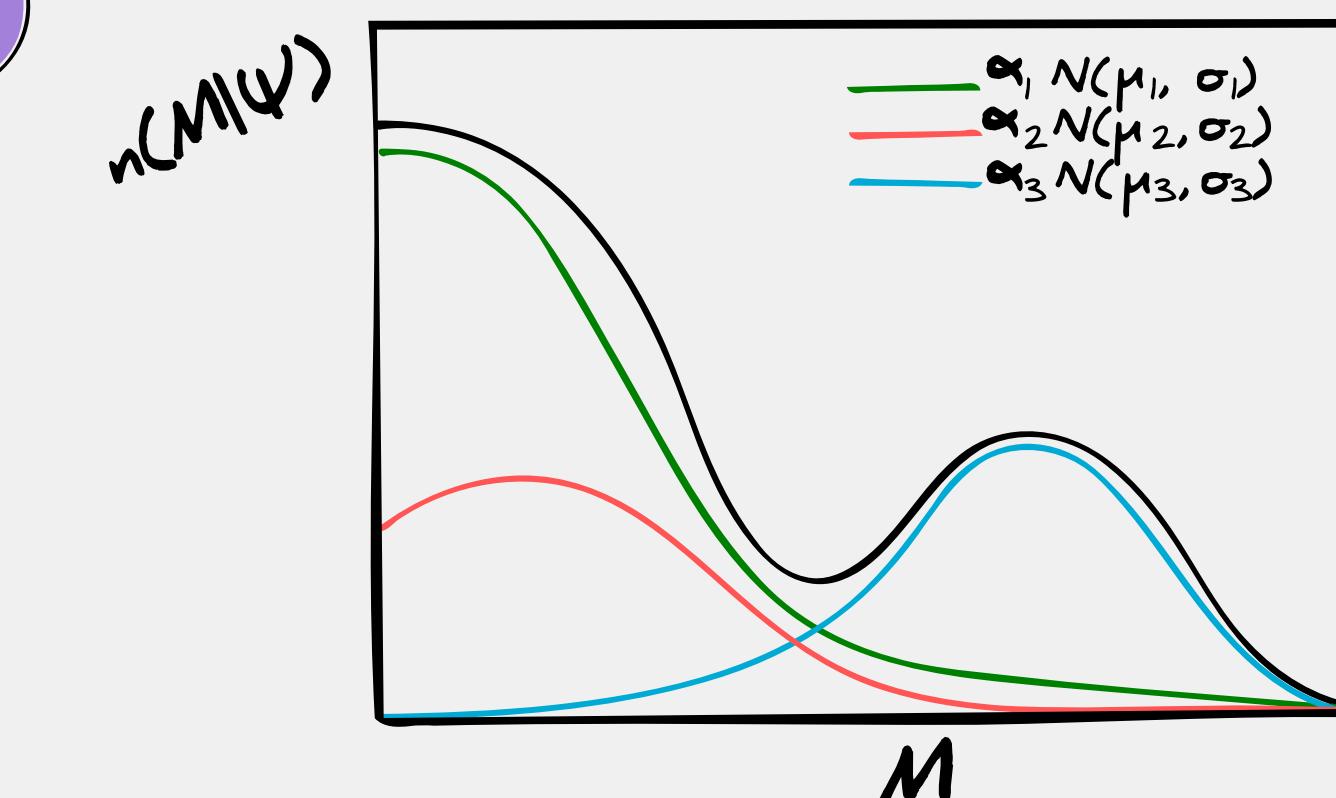
**Halo mass distribution function is a smooth function of mass given a density environment**

**Use a mixture density network**



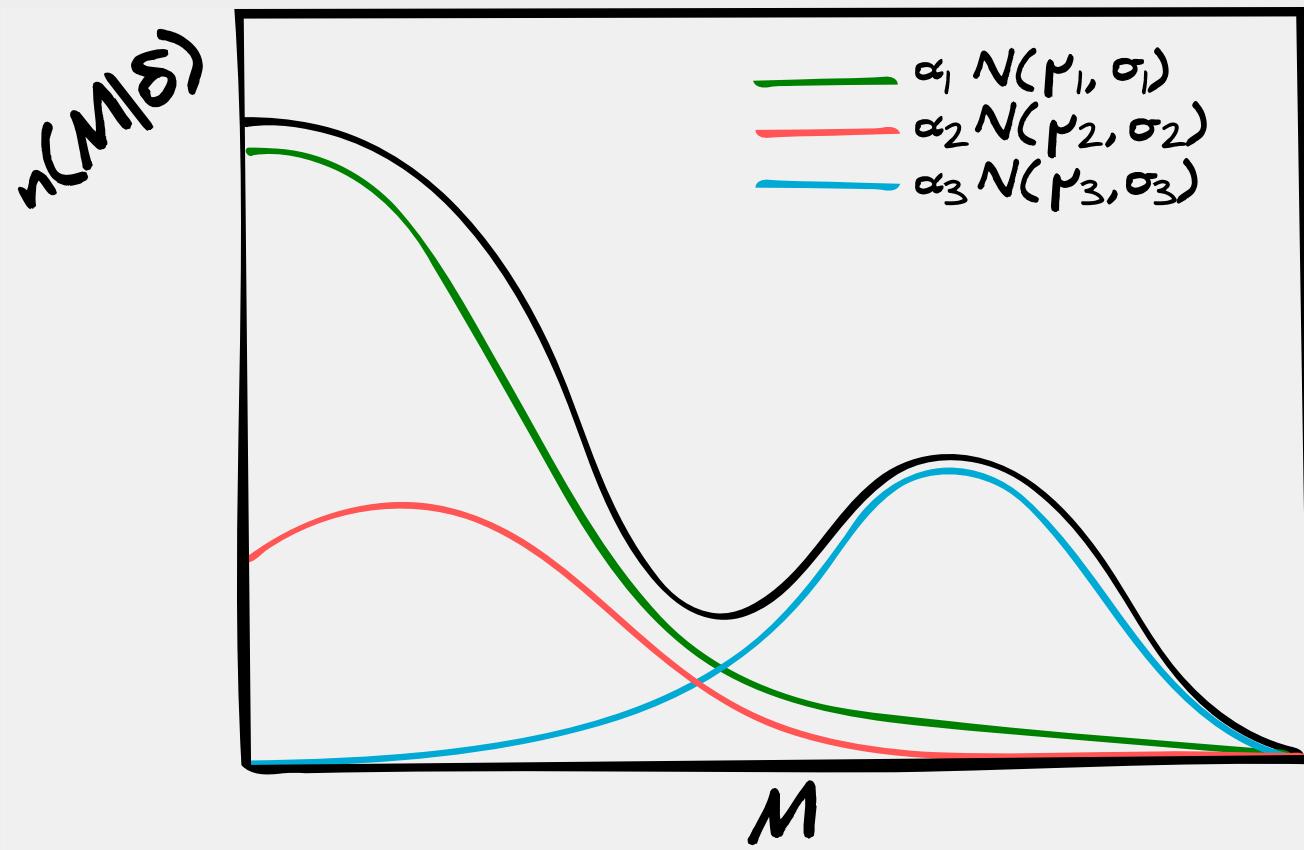
$$\alpha_i N(\mu_i, \sigma_i) = \frac{\alpha_i \exp[-(\log M - \mu_i)^2 / 2\sigma_i^2]}{\sqrt{2\pi\sigma_i^2}}$$

$$L = -\ln \sum_i \alpha_i N(\mu_i, \sigma_i)$$



# Our neural density estimator

$$\begin{aligned}
n(M|\delta) &= \sum_i^N \alpha(\psi, \theta_i^\alpha) \mathcal{N} (\mu(\psi, \theta_i^\mu), \sigma(\psi, \theta_i^\sigma) | M), \\
&= \sum_i^N \frac{\alpha(\psi, \theta_i^\alpha)}{\sqrt{2\pi(\sigma(\psi, \theta_i^\sigma))^2}} \exp \left[ -\frac{(\log(M) - \mu(\psi, \theta_i^\mu))^2}{2(\sigma(\psi, \theta_i^\sigma))^2} \right],
\end{aligned}$$



$$\begin{aligned}
\alpha_i &= \text{softplus}(w_i^\alpha \psi + b_i^\alpha), \\
\mu_i &= \begin{cases} w_i^\mu \psi + b_i^\mu & i = 0 \\ \text{Max} [0, w_i^\mu \psi + b_i^\mu] + \mu_{i-1} & i > 0 \end{cases}, \\
\sigma_i &= \text{softplus}(w_i^\sigma \psi + b_i^\sigma),
\end{aligned}$$

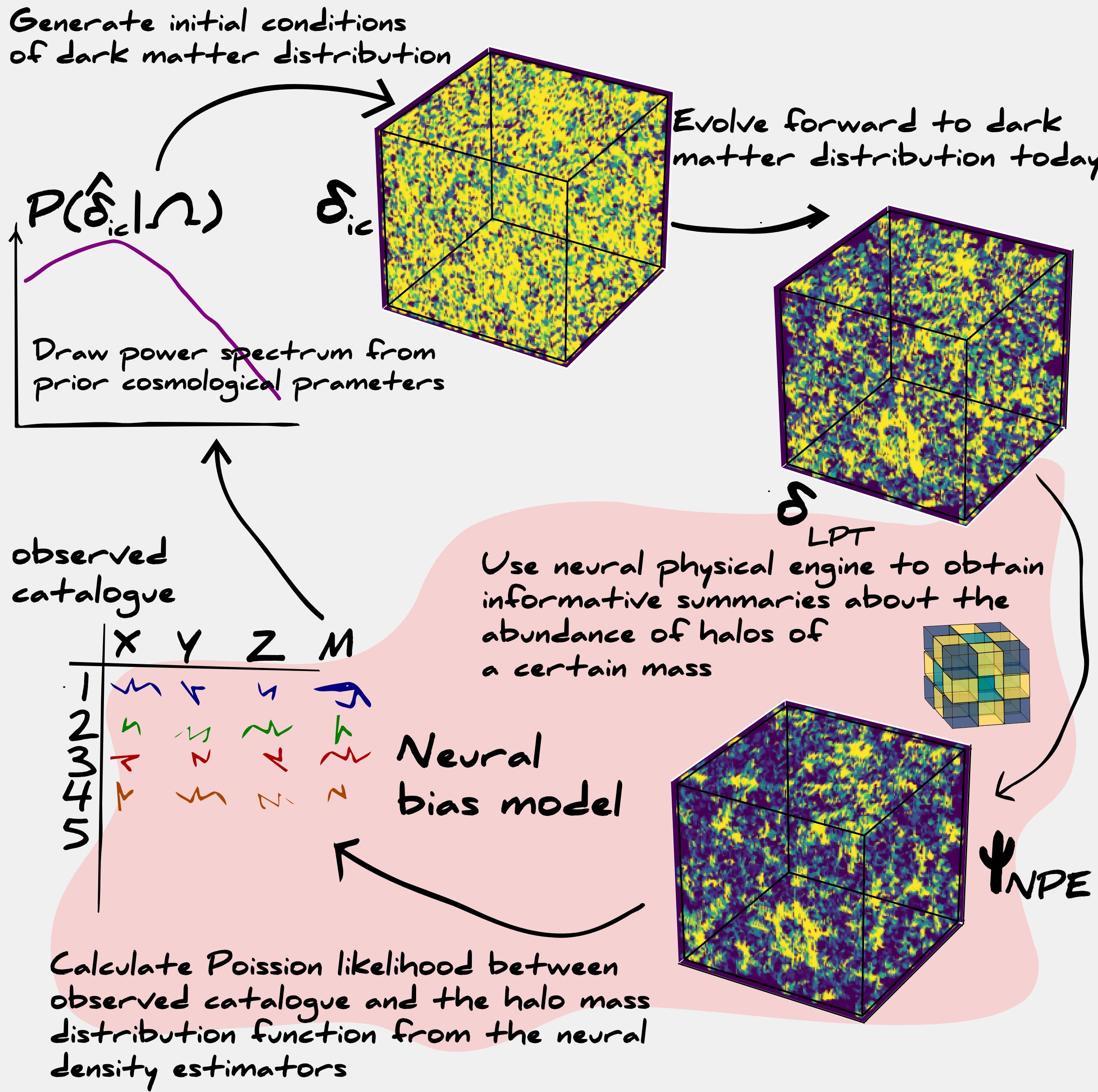
2 Gaussians with

$$\beta^\alpha \rightarrow \beta^\alpha + \log(10^{-3}), b_0^\mu \rightarrow b_0^\mu + \log(2 \times 10^{12}) \text{ and } \beta^\sigma \rightarrow \beta^\alpha + \log(10^3).$$

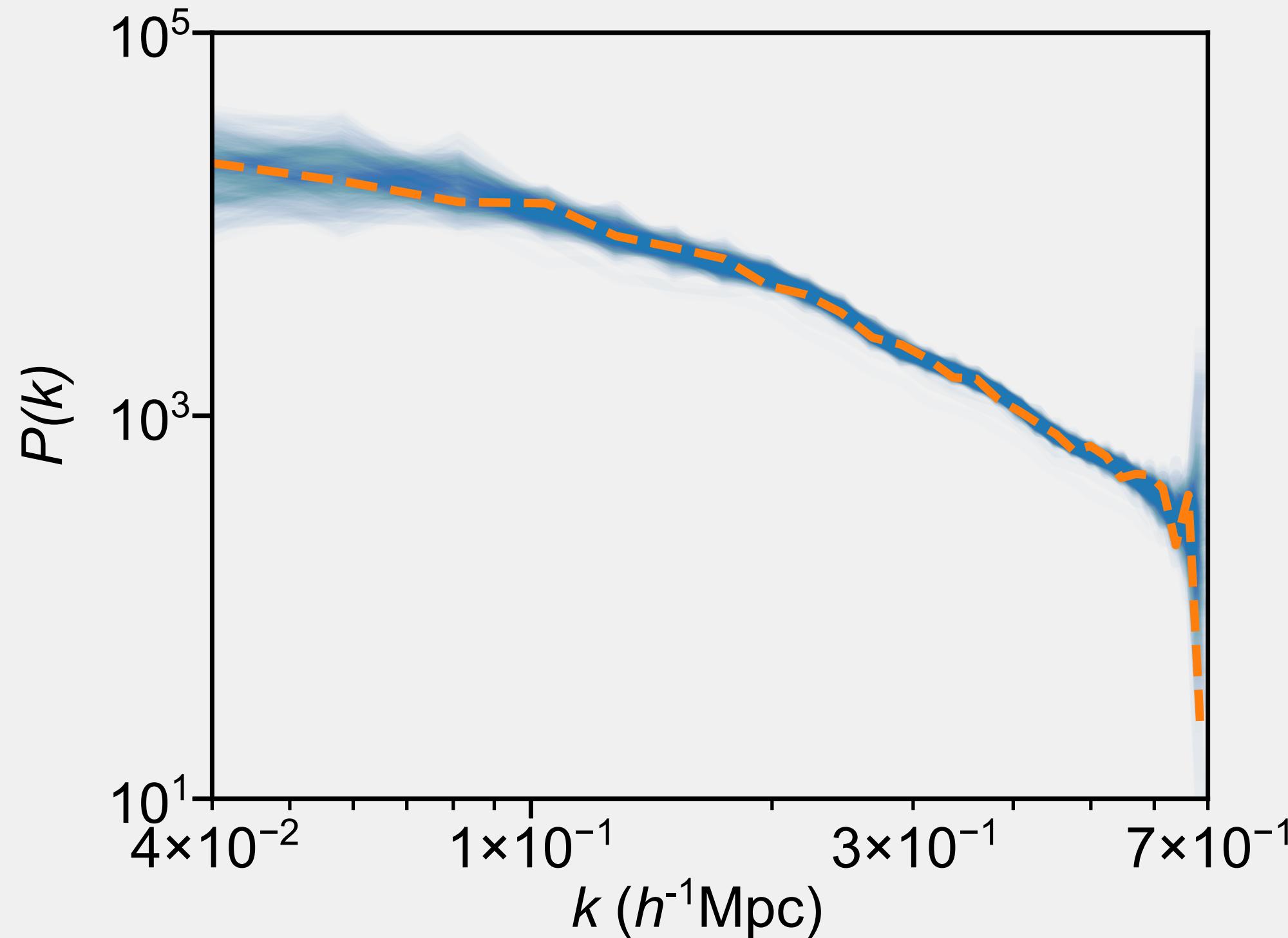
# Likelihood

# Likelihood

$$\begin{aligned}\mathcal{L} = & \sum_{j \in \text{catalogue}} \log \left[ \sum_i^N \frac{\alpha_{i,j}}{\sqrt{2\pi\sigma_{i,j}^2}} \exp \left[ -\frac{(\log(M_j) - \mu_{i,j})^2}{2\sigma_{i,j}^2} \right] \right] \\ & - V \sum_{j \in \text{voxels}, i=1}^N \frac{\alpha_{i,j}}{2} \exp \left[ \frac{\sigma_{i,j}^2}{2} \right] \operatorname{erfc} \left[ \frac{\log(M_{\text{th}}) - \mu_{i,j} - \sigma_{i,j}^2}{\sqrt{2\sigma_{i,j}^2}} \right].\end{aligned}$$

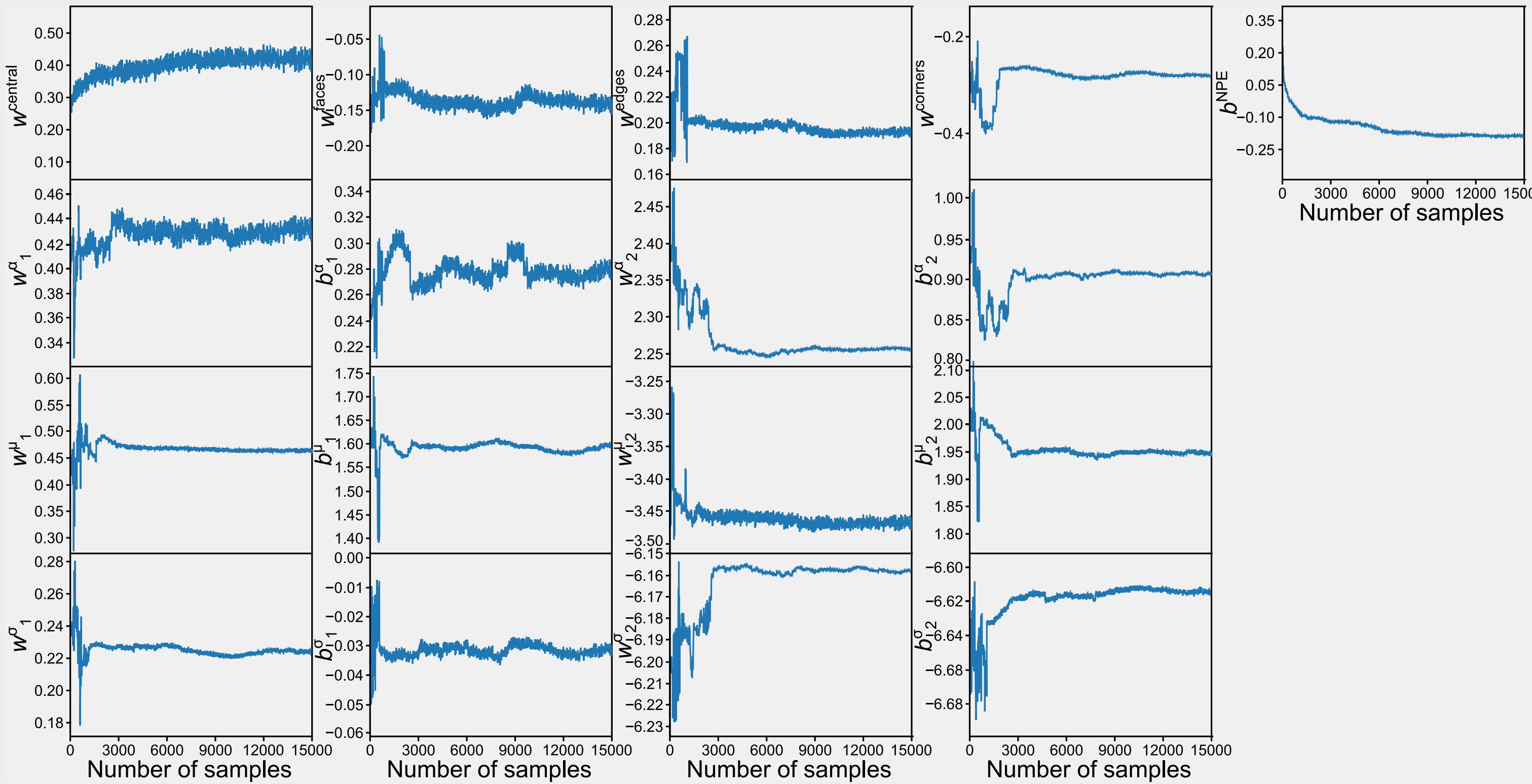


# Sample the power spectrum



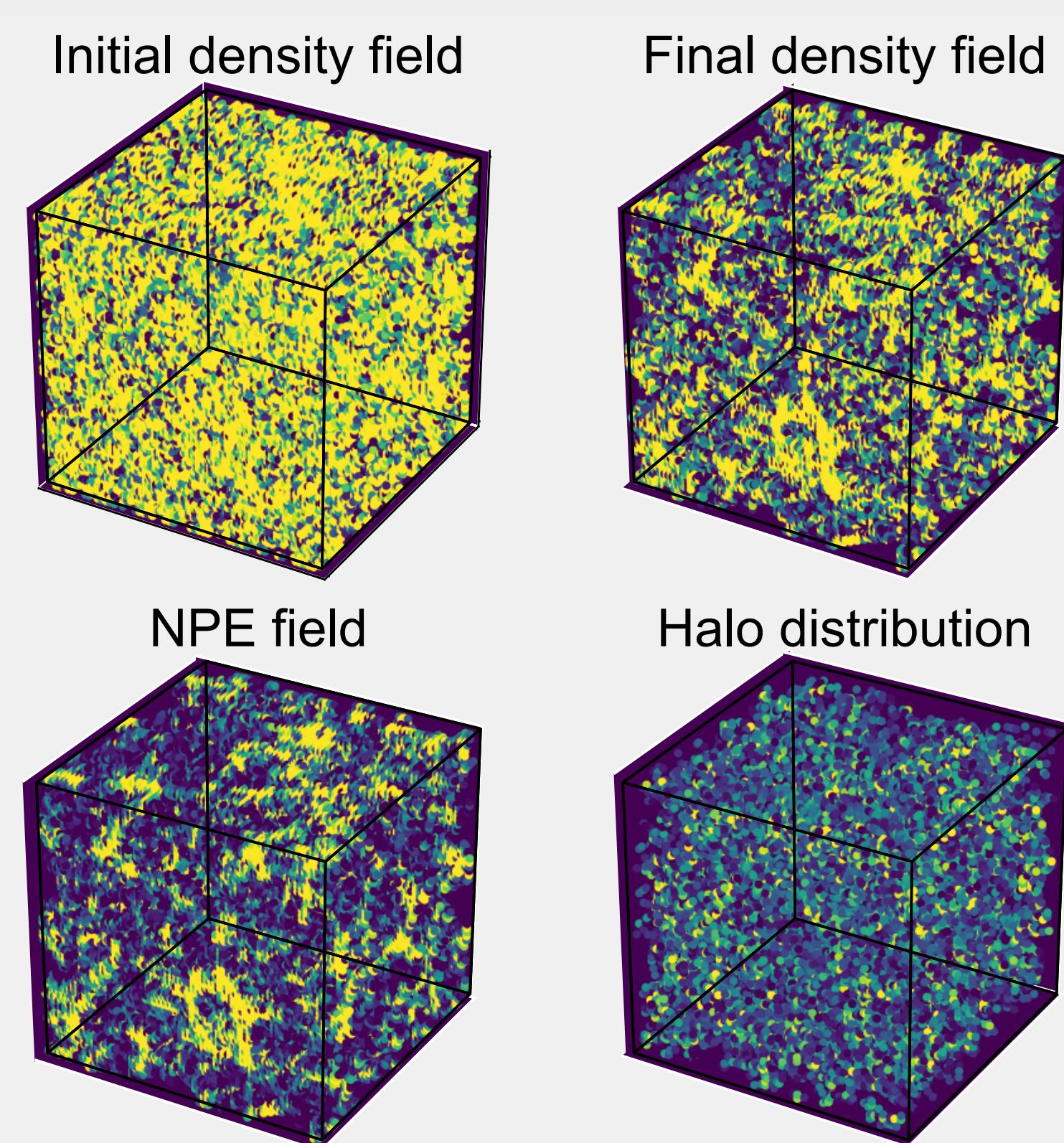
✓ Power spectra remain smooth and close to prior

# Sample the weights of the neural bias model



- ✓ Weight values are properly sampled after burn in
- ✓ NPE acts as a contrast enhancer

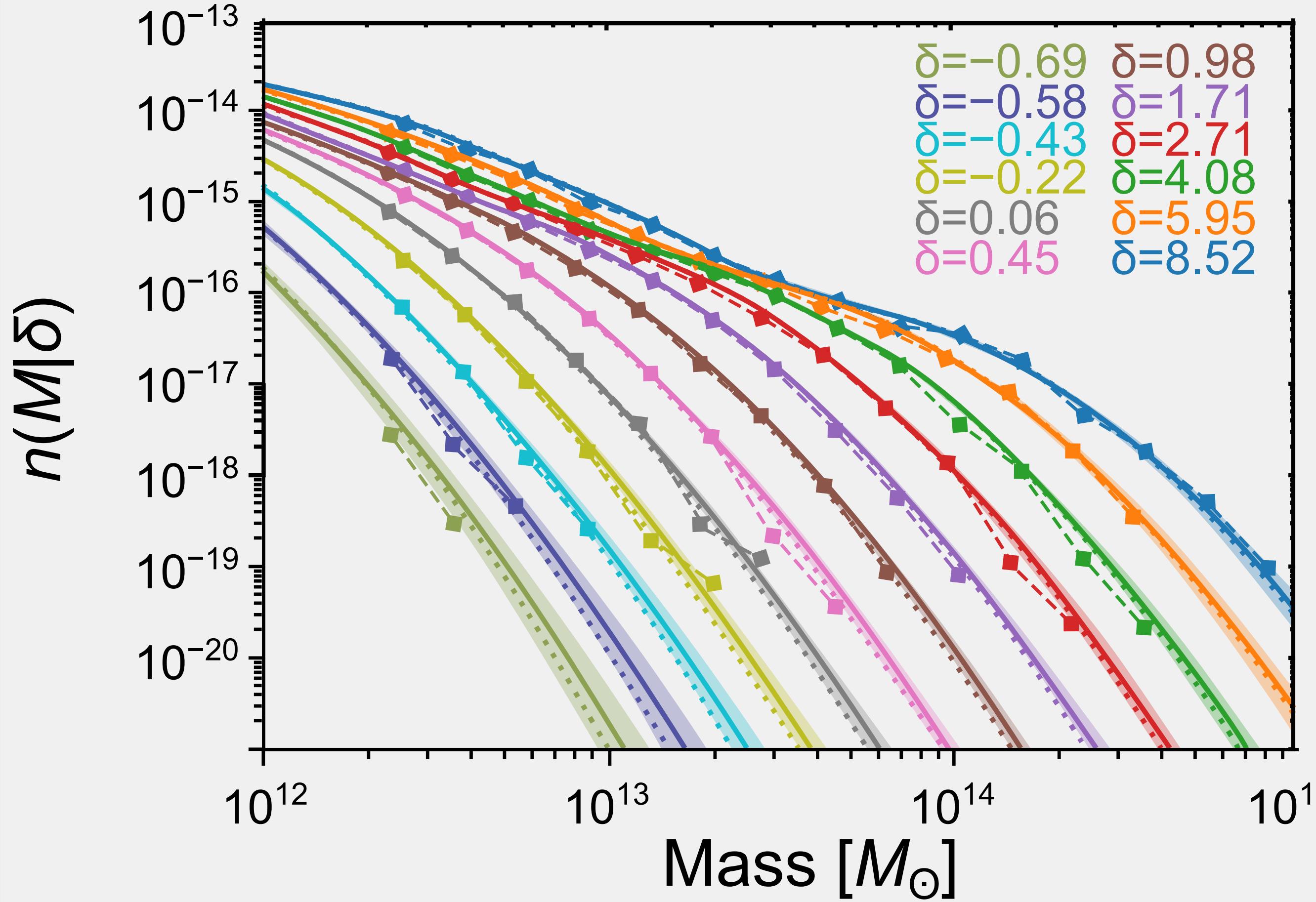
# See the effects of the neural physical engine



✓ Non-local information is used to improve fit

**And what does the halo mass distribution function look like?**

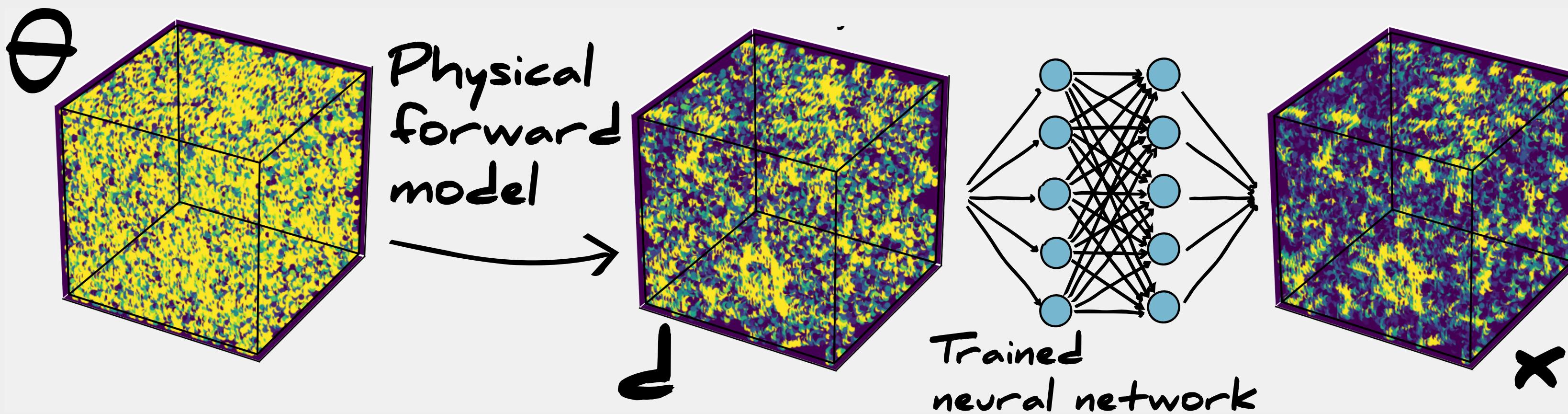
# Halo mass distribution function from neural bias model



✓ Fits data!

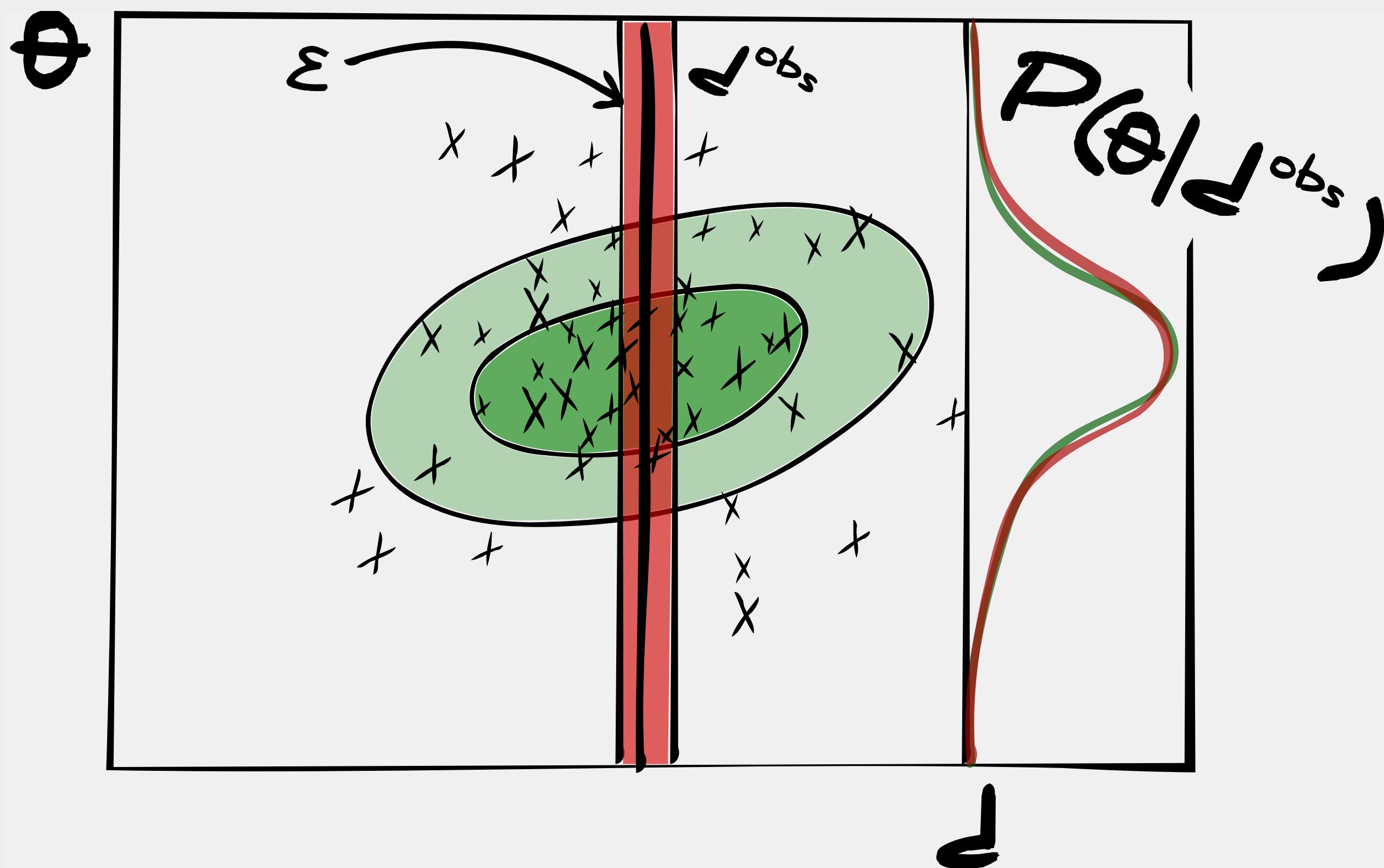
# **Method 2 : Likelihood-free inference**

# Method 2 : Likelihood-free inference

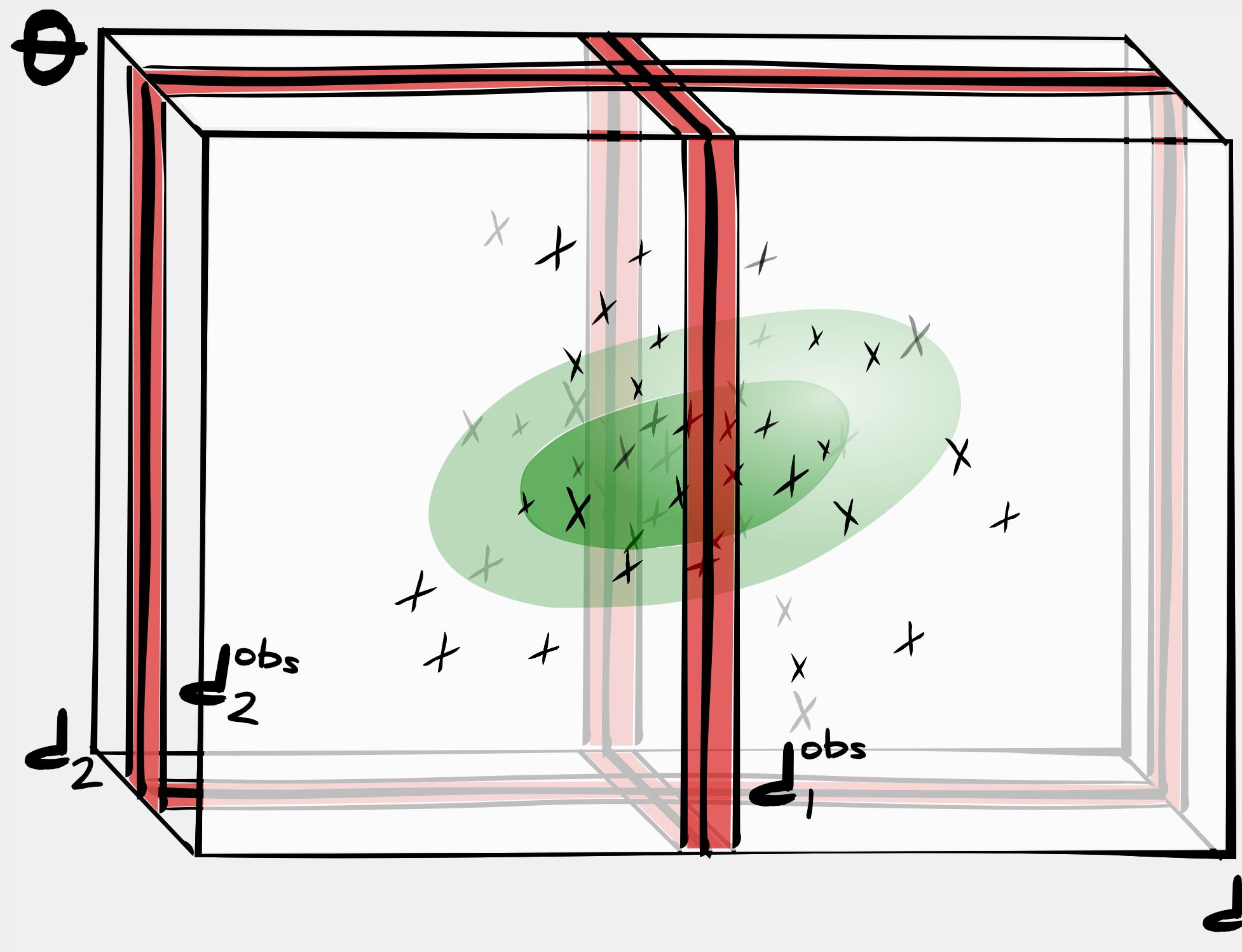


# Likelihood-free inference using machine learning

# Approximate Bayesian computation

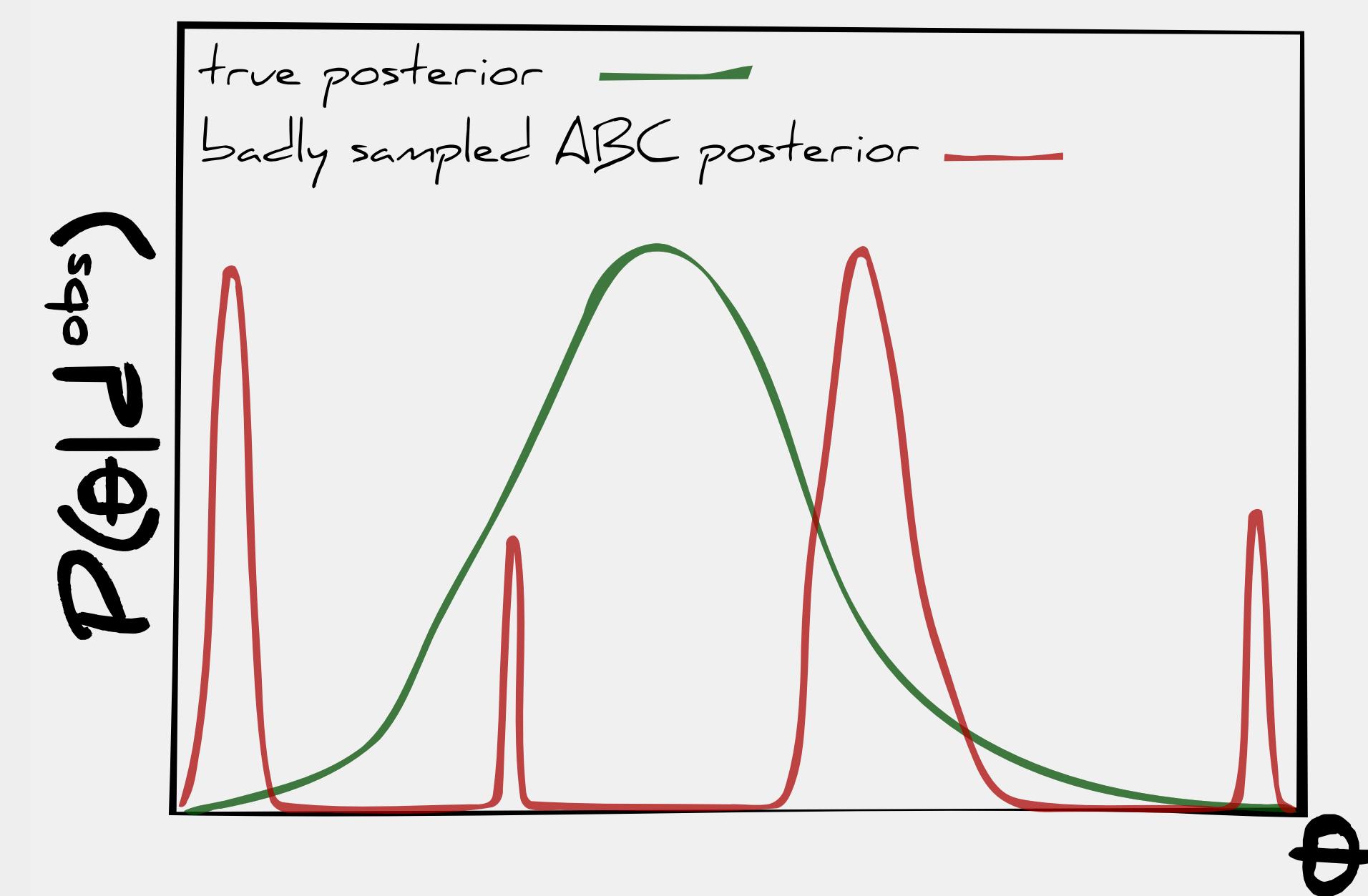


# Why even introduce the neural network?



# Inadequate sampling

Impossibly large numbers of simulations become necessary to correctly sample the posterior



**We need to do some compression**

**Any pretrained neural network *could* be used**

Any pretrained neural network *could* be used

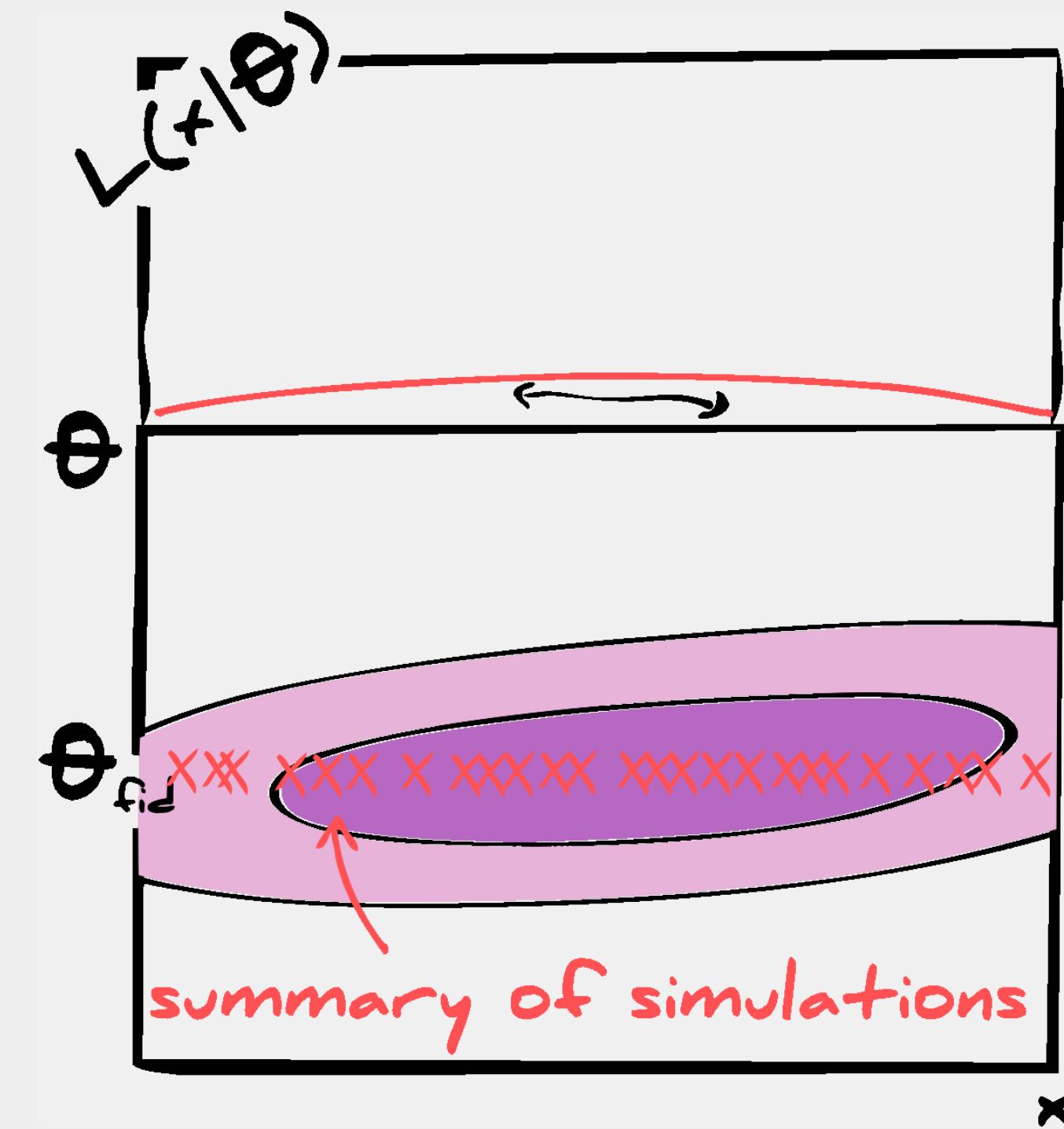
Is there a well motivated choice?



# Information maximising neural networks

Which function  $f : \mathbb{R}^d \rightarrow \mathbb{R}^\theta$  maximises the Fisher information of the summaries  $x$  from that function?

$$\mathcal{L}(d|\theta) \rightarrow \mathcal{L}(x|\theta, d)$$

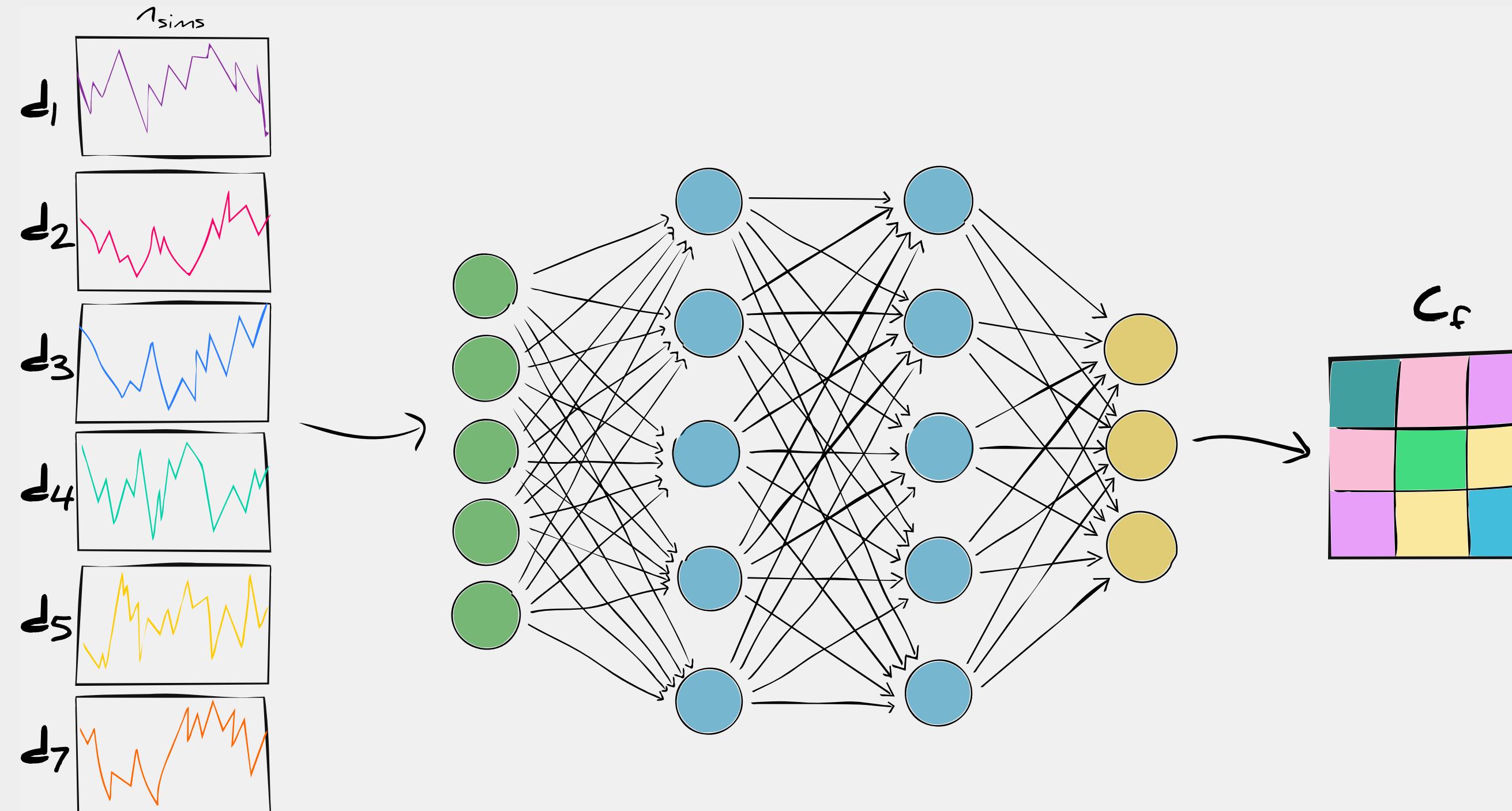


# Fisher information of the network summaries

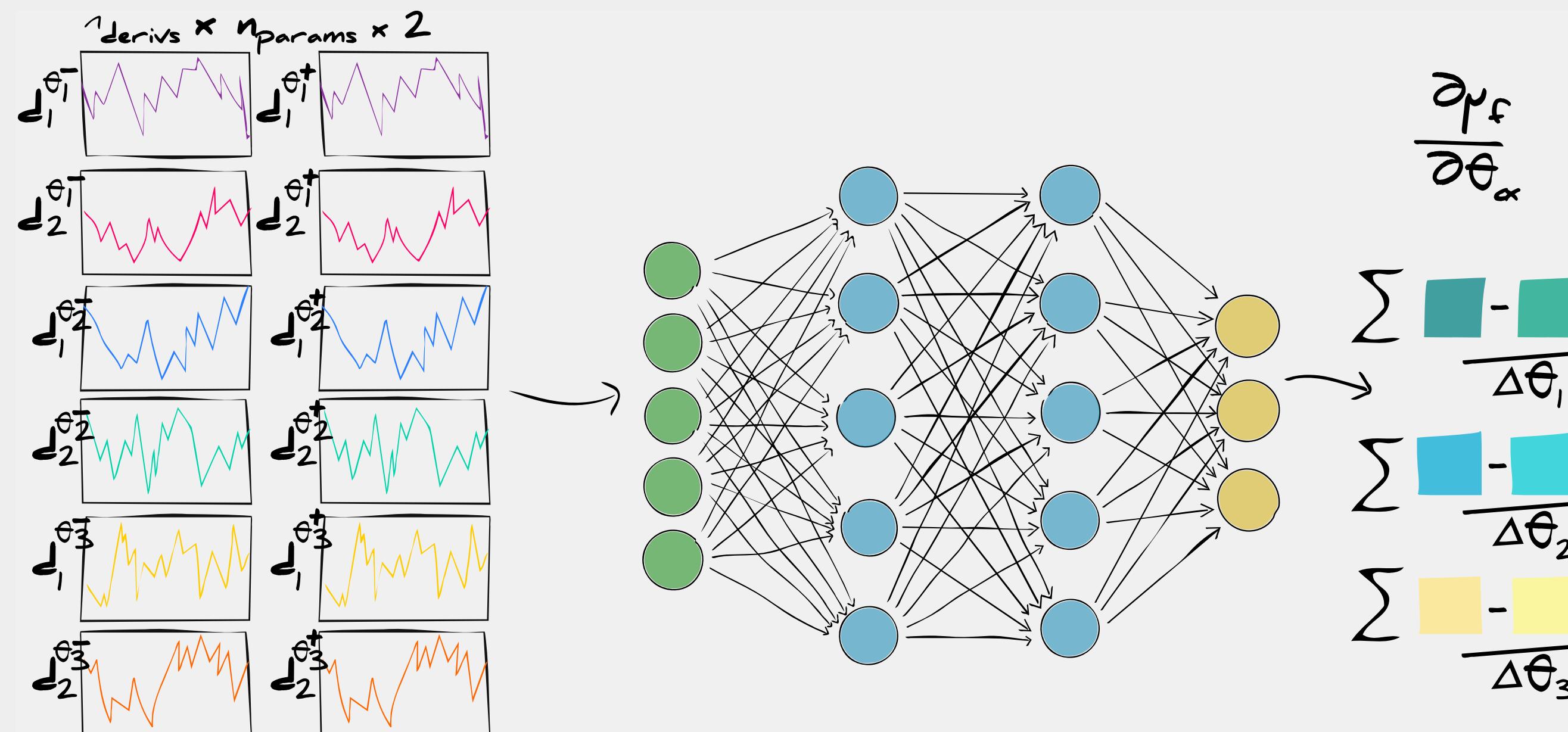
Simulations  $\{\mathbf{d}_i | i \in [1, n_d]\}$  at a single parameter value  $\theta^*$

Seed matched simulations  $\{\mathbf{d}_i^\pm | i \in [1, n_p]\}$  at perturbed fiducial parameters  $\Delta\theta^\pm = \theta \pm \delta$

# Calculate the covariance



# Calculate the derivative of the mean of the summaries with respect to the parameters



# Calculate the Fisher information

$$\mathbf{F}_{\alpha\beta} = \frac{\partial \boldsymbol{\mu}_\ell^T}{\partial \theta_\alpha} \mathbf{C}_\ell^{-1} \frac{\partial \boldsymbol{\mu}_\ell}{\partial \theta_\beta}$$

This Gaussian form forces the summaries to be Gaussianised

**Optimise the tunable parameters of the network such that  $\ln |F_{\alpha\beta}|$  is maximised!**

Once converged the network compresses  
*and* Gaussianises the data *without* losing  
information\*

\*in the optimal case...

And we get free maximum likelihood estimates...

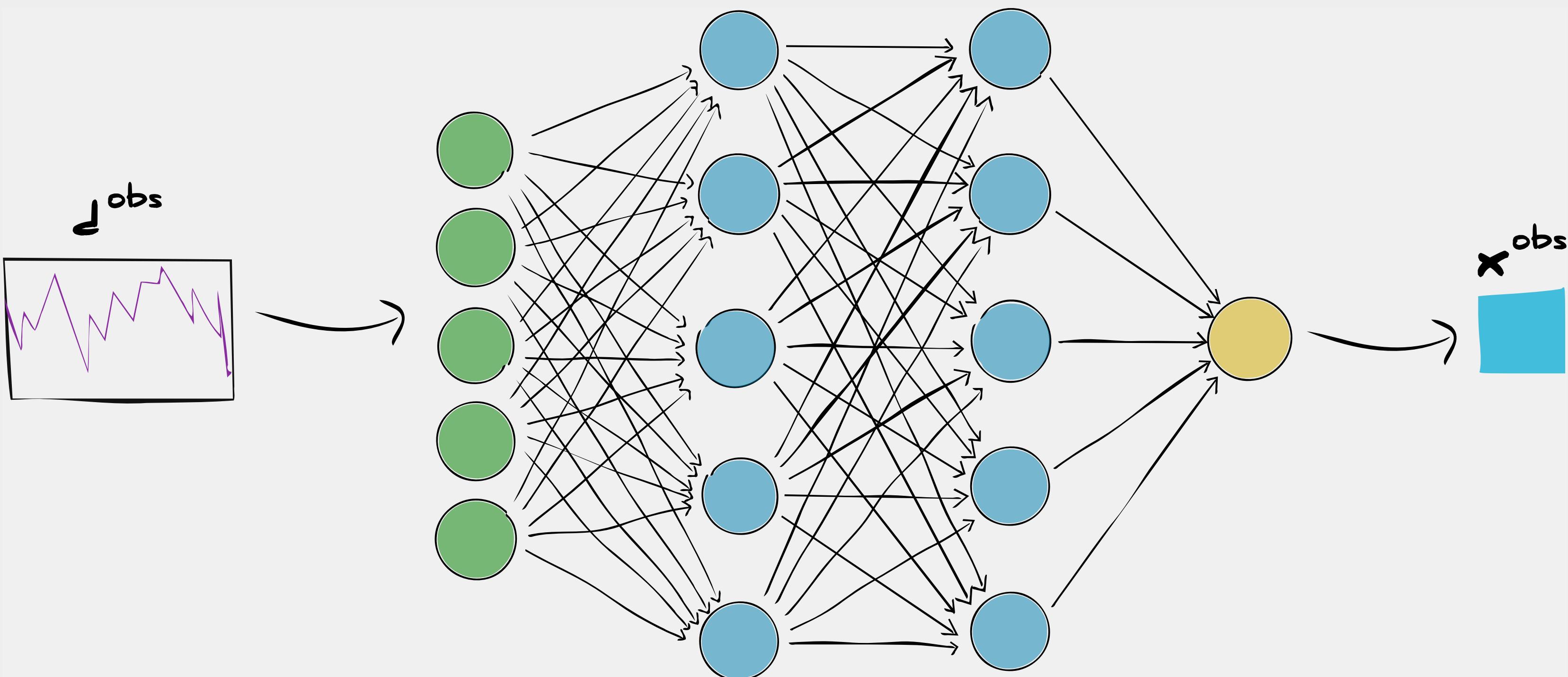
$$\theta_{\alpha}^{\text{MLE}} = \theta_{\alpha}^* + \mathbf{F}_{\alpha\beta}^{-1} \mathbf{C}_{\ell}^{-1} \frac{\partial \mu_{\ell}}{\partial \theta_{\beta}} (\mathbf{x} - \mu_{\ell})$$

And we can get free approximations of the posterior

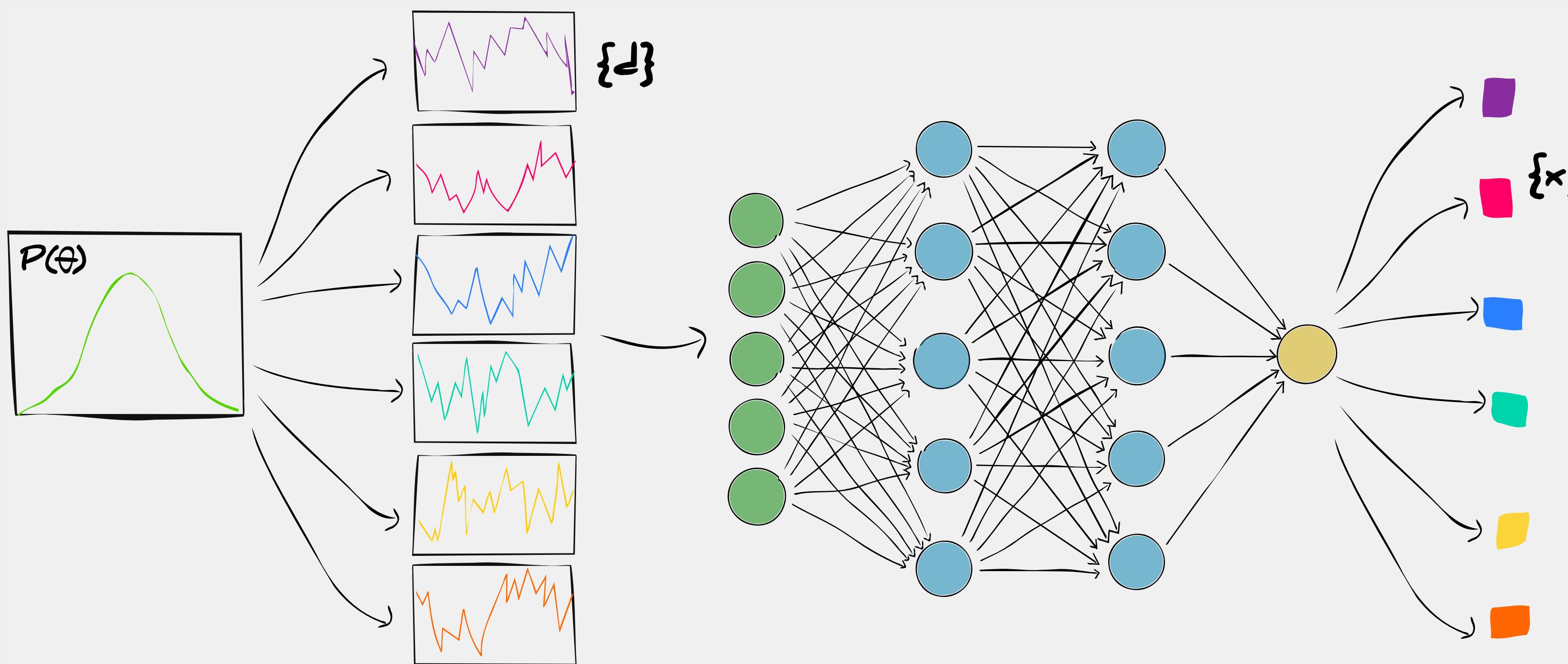
$$\text{Cov}[\theta_\alpha^{\text{MLE}}, \theta_\beta^{\text{MLE}}] \geq F_{\alpha\beta}^{-1}$$

**Or do likelihood-free inference**

We pass our observed data through the network  
(put it to the side)



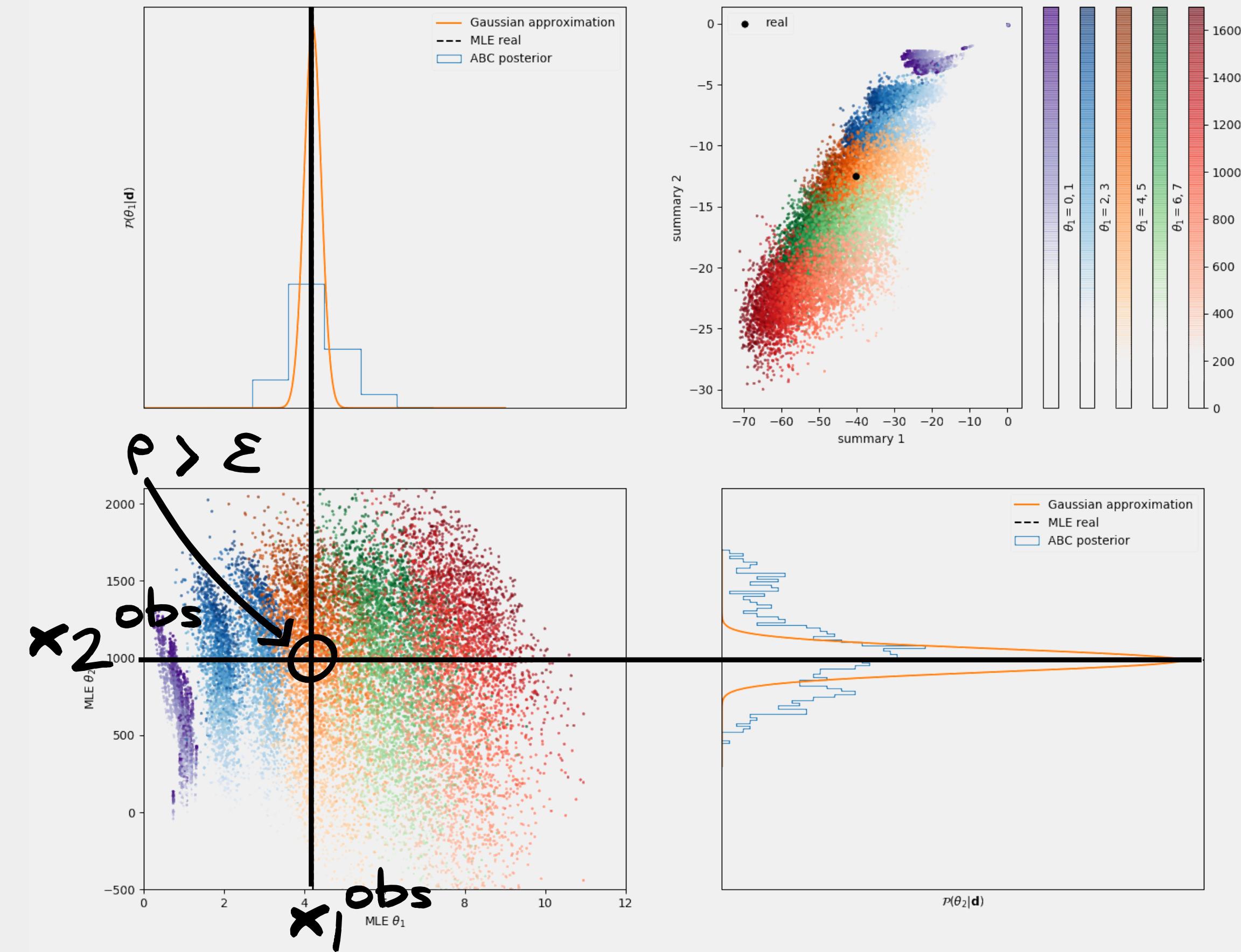
Draw simulations (cleverly?) from the prior model parameters,  $\mathcal{P}(\theta)$  and pass them through the network



Measure the difference between the summarised simulations and the observed summary

$$\sqrt{(\{ \cdot \textcolor{violet}{\cdot} \cdot \textcolor{red}{\cdot} \cdot \textcolor{blue}{\cdot} \cdot \textcolor{teal}{\cdot} \cdot \textcolor{yellow}{\cdot} \cdot \textcolor{orange}{\cdot} \} - \textcolor{cyan}{x}^{\text{obs}})^2} = R$$

# Compare distance between observed summaries and simulation summaries and select results within $\epsilon$



# Conclusions

# Conclusions

**Neural networks are not to be trusted**

**They can make trusty companions - when the correct framework is introduced**

**Using statistics we can build them into the forward model to give us lossless summaries**

# Take home message

# Take home message

Stop doing machine learning, think, then start doing machine learning again!