

Machine Learning for HEP

Lecture II – Bayesian Neural Networks



Midjourney AI



UNIVERSITÀ
DEGLI STUDI
DI MILANO

BND Graduate School – Blankenberge 2024

Ramon Winterhalder

Lecture II

Bayesian Neural Networks

Reminder — Likelihood-Loss



Training should **optimize** probability that **network parameters** ω describe the **data** x

This is given by: $p(\omega | x)$ → Usually intractable 😞



Bayes' Theorem **Likelihood** → Known from simulation 😊

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

Labels for the equation:

- Likelihood** points to $p(x | \omega)$
- Prior** points to $p(\omega)$
- Normalization** points to $p(x)$
- Posterior** points to $p(\omega | x)$

Reminder — Likelihood-Loss



Bayes' Theorem

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

Likelihood-loss

$$\begin{aligned} \mathcal{L} &= -\log p(\omega | x) \\ &= -\log p(x | \omega) - \log p(\omega) + \text{const} \end{aligned}$$

Weight regularization

Heteroscedastic loss

$$\mathcal{L}_{\text{heteroscedastic}} = \sum_i \frac{|f(x_i) - f_\omega(x_i)|^2}{2\sigma(x_i)^2} + \log \sigma(x_i) + \dots$$

Normalization term

Regression with
Gaussian assumption

with point-wise
uncertainty

Homoscedastic loss = MSE

$$\mathcal{L}_{\text{MSE}} = \frac{1}{2\sigma} \sum_i |f(x_i) - f_\omega(x_i)|^2 + \dots$$

Same uncertainty

$\sigma(x) \rightarrow \sigma = \text{const}$

Is there more we can do?

Yes!

Bayesian Neural Networks

→ **There is nothing really Bayesian about them!**

What we have vs what we want



So far: NNs are deterministic \longrightarrow **fixed** input x \longrightarrow **fixed** output y

However, in **particle physics**:

\rightarrow we are not only interested in **results**, but also in **errorbars**

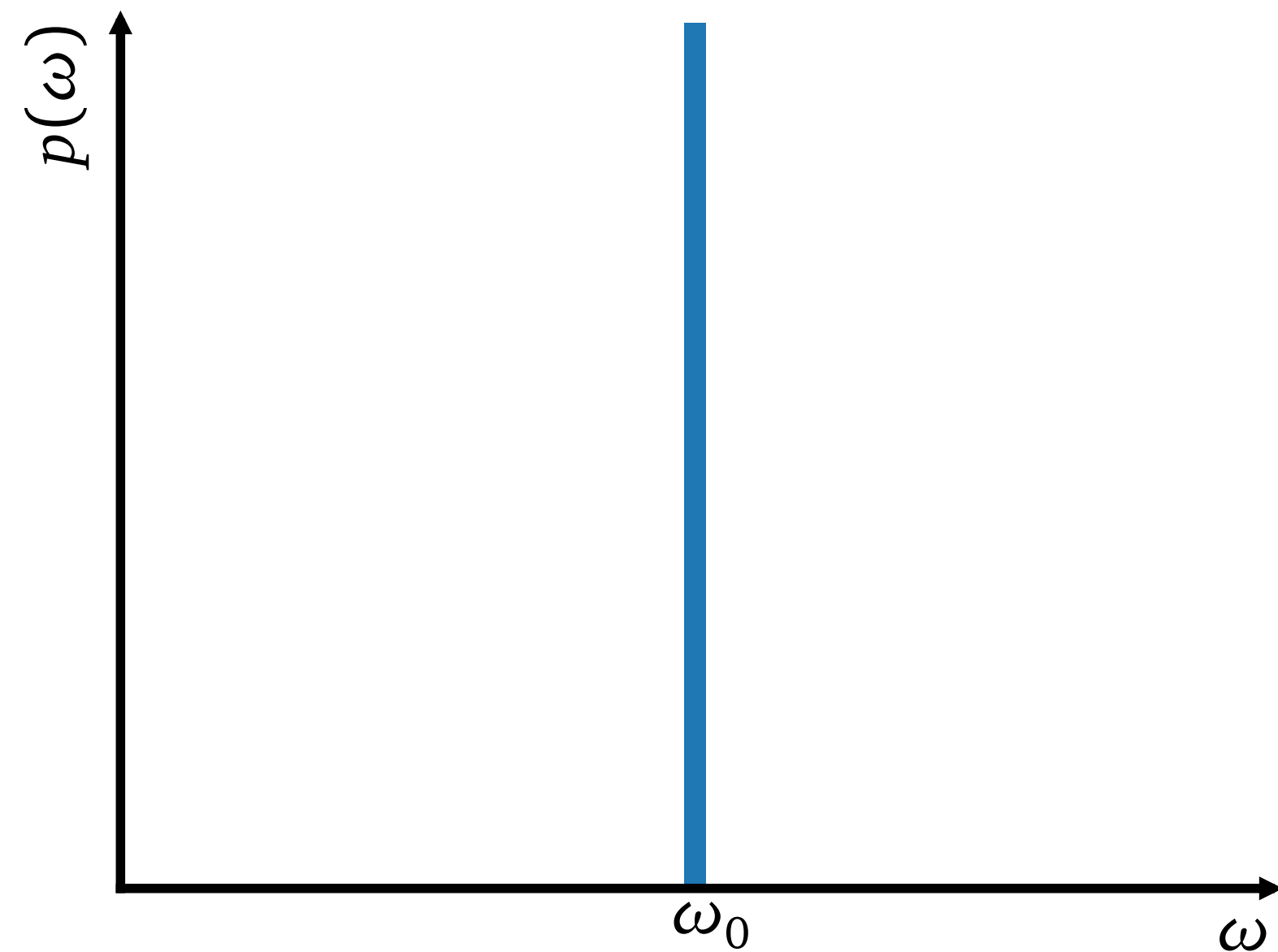
We need: **fixed** input x \longrightarrow **probabilistic** output $\langle y \rangle \pm \sigma_y$ \longrightarrow $y \sim p(y|x)$

How do we achieve this?

Bayesian neural networks



Neural network

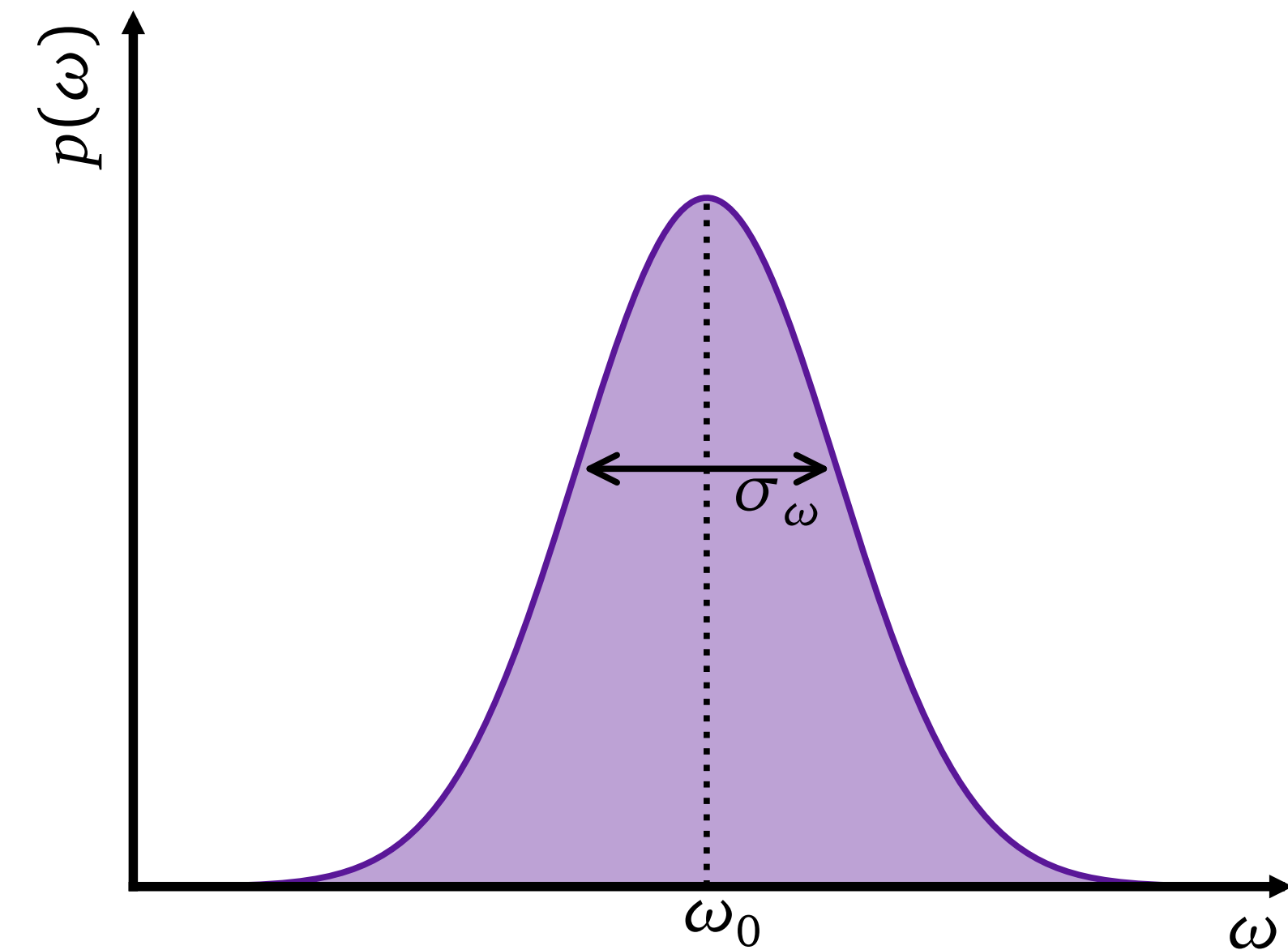


Network weights are deterministic

$$\omega = \omega_0$$



Bayesian Neural network



Network weights are drawn from distribution

$$\omega \sim p(\omega)$$

Basics of BNNs



Predictive distribution in Classification

$$p(y|x) = \int d\omega p(y|\omega, x) p(\omega|x)$$

Classifier output

→ average over $p(\omega|x)$ → Intractable 😞

Bayes' theorem

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

Need to Calculate normalization

$$p(x) = \int d\omega p(x|\omega) p(\omega)$$

Also Intractable 😞

Prediction in Regression

$$\begin{aligned} y(x) \equiv \langle y \rangle &= \int dy p(y|x) y \\ &= \int d\omega p(\omega|x) \int dy p(y|x, \omega) y \\ &= \int d\omega p(\omega|x) \bar{y}(x, \omega) \end{aligned}$$

Regression output

So what do we do instead?

Variational approximation

Approximate posterior

$$p(\omega | x) \simeq q_\alpha(\omega | x) \approx q_\alpha(\omega)$$

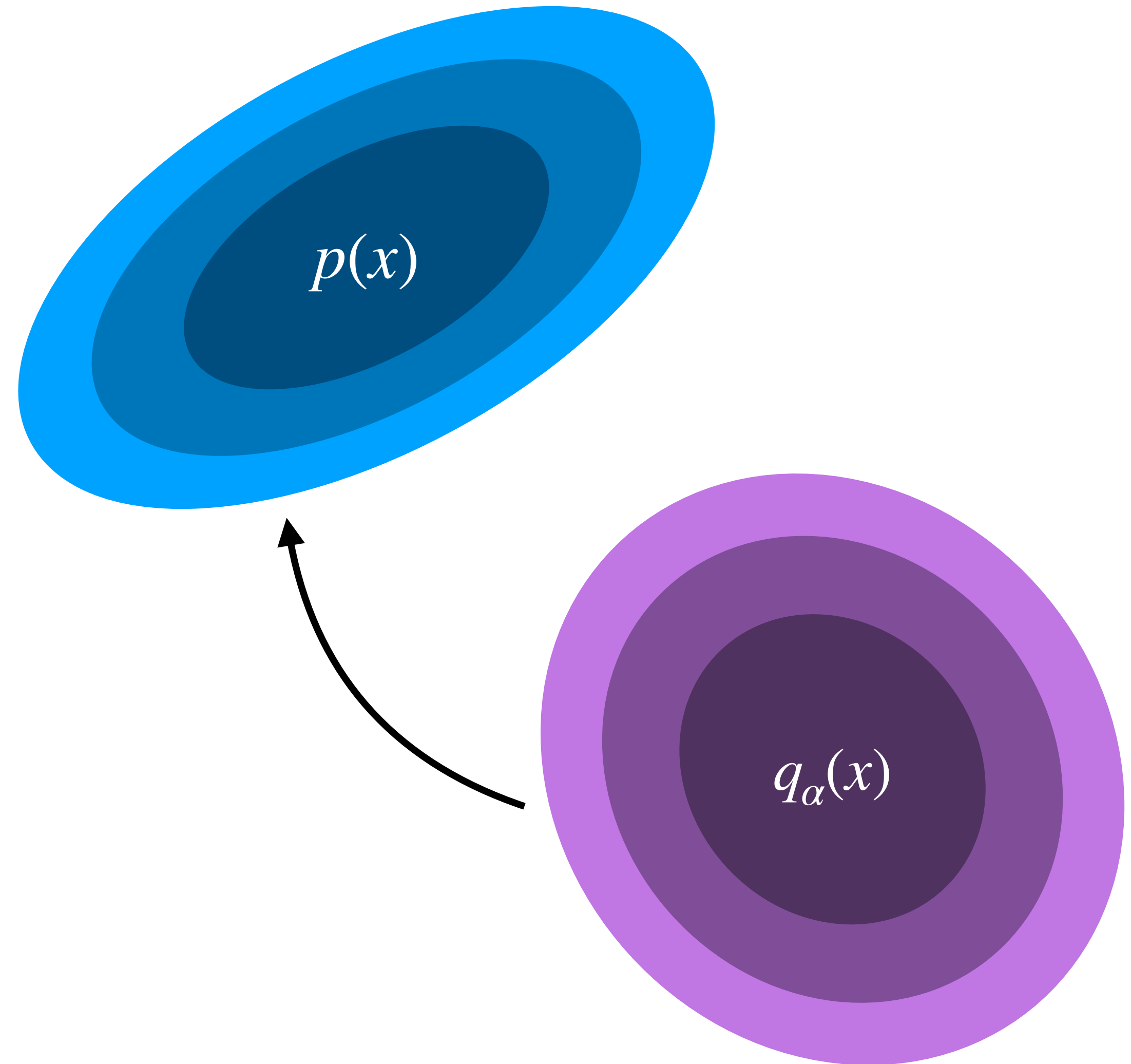
Simplify (see later)

Parameters α to adjust

Kullback Leibler Divergence

$$\text{KL}(q_\alpha(x), p(x)) = \int dx q_\alpha(x) \log \left(\frac{q_\alpha(x)}{p(x)} \right)$$

- Measures information overlap
- Always positive \rightarrow zero iff $p(x) = q_\alpha(x)$
- Non symmetric



BNN Loss function

Kullback Leibler Divergence

$$\begin{aligned}
 \text{KL}(q_\alpha(\omega), p(\omega | x)) &= \int dx q_\alpha(\omega) \log \left(\frac{q_\alpha(\omega)}{p(\omega | x)} \right) \\
 &= \int dx q_\alpha(\omega) \log \left(\frac{q_\alpha(\omega) p(x)}{p(\omega) p(x | \omega)} \right) \\
 &= \text{KL}(q_\alpha(\omega), p(\omega)) - \int d\omega q_\alpha(\omega) \log p(x | \omega) + \log p(x) \int d\omega q_\alpha(\omega)
 \end{aligned}$$

Intractable norm

= 1

BNN loss function

Problem specific

$$\mathcal{L}_{\text{BNN}} = \text{KL}(q_\alpha(\omega), p(\omega)) - \int d\omega q_\alpha(\omega) \log p(x | \omega)$$

- ① Neg log-likelihood averaged over q_α
- ② q_α should not deviate too much from prior!

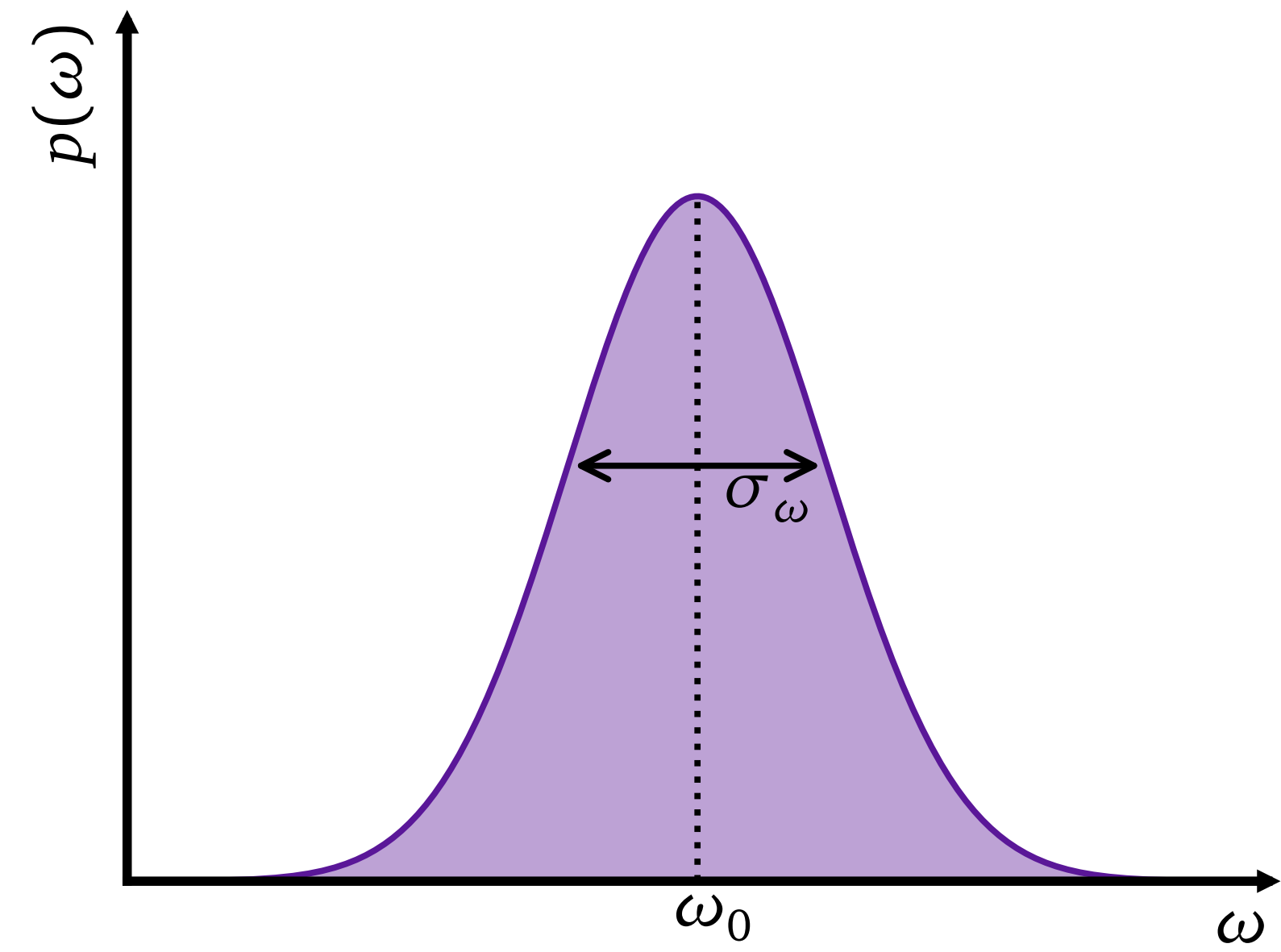
How to choose q_α and the prior?

Standard choice

In practice: p, q_α are chosen to be Gaussian

$$q_\alpha(\omega) \equiv q_{\mu, \sigma}(\omega) = \frac{1}{\sqrt{2\pi\sigma_q^2}} \exp\left[-\frac{(\omega - \mu_q)^2}{2\sigma_q^2}\right]$$
$$p(\omega) = \frac{1}{\sqrt{2\pi\sigma_p^2}} \exp\left[-\frac{(\omega - \mu_p)^2}{2\sigma_p^2}\right]$$

Usually $\mu_p = 0$



KL divergence

$$\text{KL}(q_\alpha(\omega), p(\omega)) = \frac{\sigma_q^2 - \sigma_p^2}{2\sigma_p^2} + \frac{\mu_q^2}{2\sigma_p^2} + \log \frac{\sigma_p}{\sigma_q}$$

2

2 serves as regularization!

Deterministic limit

Gaussian probability

$$q_{\alpha}(\omega) = \frac{1}{\sqrt{2\pi\sigma_q^2}} \exp\left[-\frac{(\omega - \mu_q)^2}{2\sigma_q^2}\right]$$

$$\begin{array}{l} \sigma_q \rightarrow 0 \\ \mu_q \rightarrow \omega_0 \end{array}$$

Delta distribution = fixed weight

$$q_{\alpha}(\omega) = \delta(\omega - \omega_0)$$

Loss function

$$\mathcal{L}_{\text{BNN}} \rightarrow \mathcal{L}_{\text{NN}} = -\log p(x | \omega_0) + \frac{\omega_0^2}{2\sigma_p} + \text{const}$$

L2 norm: $\lambda\omega_0^2$

→ Usually λ is a free parameter

Mean-field approximation

For K network parameters we do factorization ansatz:

$$q_{\alpha}(\omega) = \prod_{i=1}^K \mathcal{N}(\omega_i | \mu_i, \sigma_i)$$

→ BNN has $2K$ parameters → **Mean-field approximation** [1505.05424, 1601.00670]

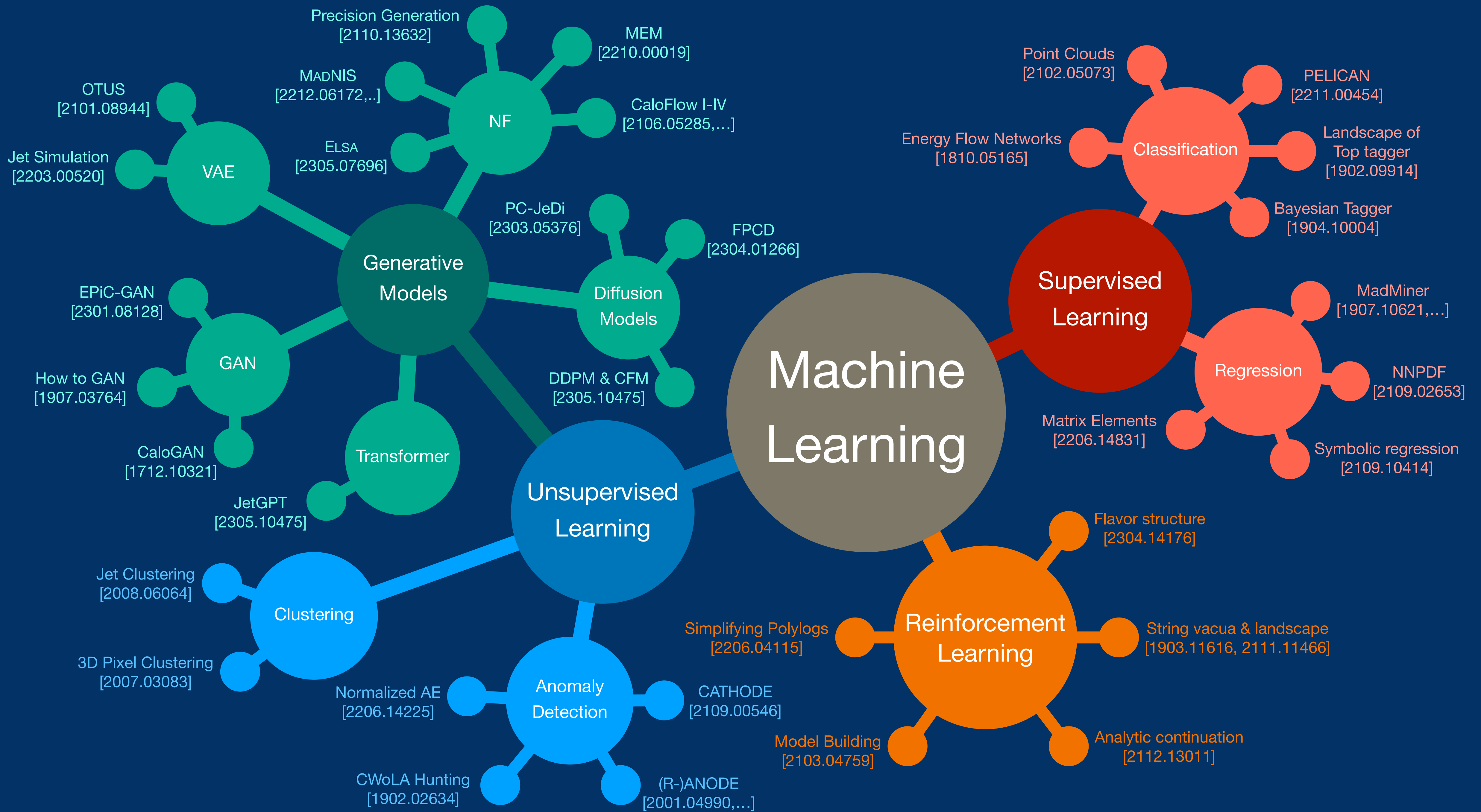
If correlated Gaussians → BNN $\sim K^2$ parameters → bad scaling 😞

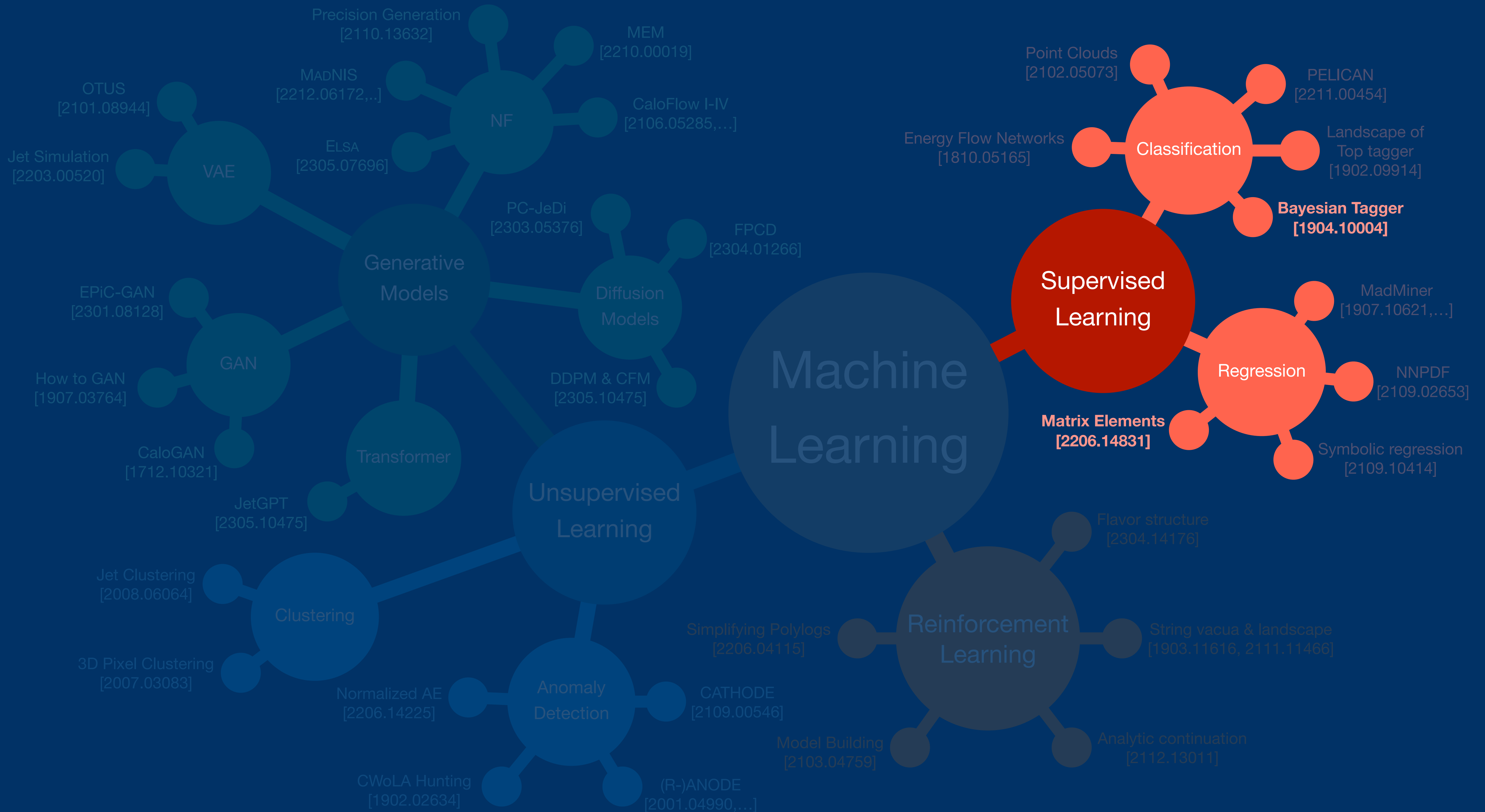
→ Deepness of NN compensates approximation?

→ **Open question in ML community!**



Questions?

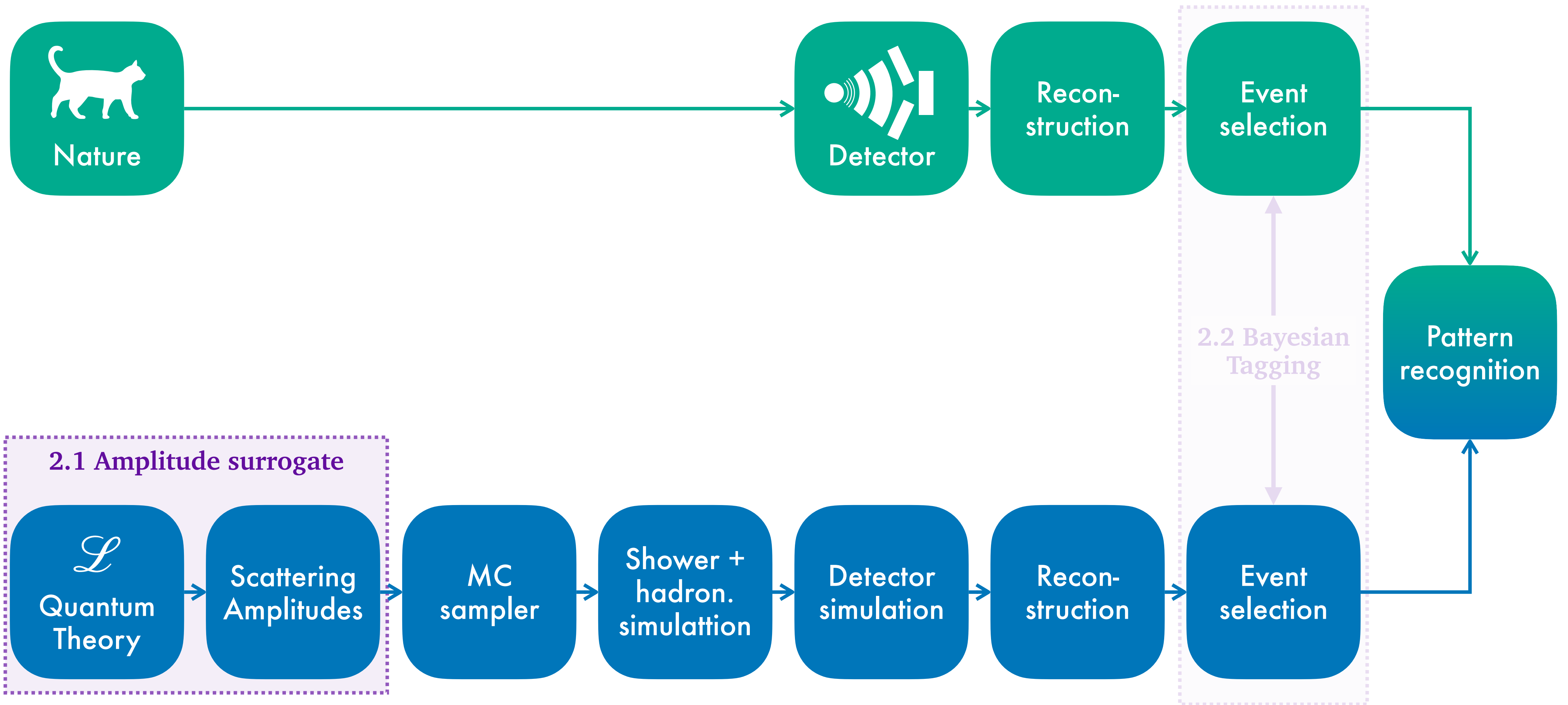




Example I

Bayesian Amplitude Regression

LHC analysis + ML



Bayesian Regression

We want to fit a set of amplitudes $A(x)$

Training data: $\{x, A(x)\}$

We can define: $p(A | x) \sim$ prob of A @ position x ($\rightarrow p(y | x)$ before)

Prediction in regression

$$A(x) \equiv \langle A \rangle = \int dA A p(A | x)$$

with $p(A | x) = \int d\omega p(A | \omega, x) p(\omega | x)$

Variational approximation

$$p(\omega | x) \simeq q_\alpha(\omega)$$

BNN



We want to describe **mean** and
stddev of prediction

BNN mean prediction

Mean prediction

$$\begin{aligned} A(x) \equiv \langle A \rangle &= \int dA A p(A | x) \quad (*) \\ &= \int dA \int d\omega A p(A | \omega, x) p(\omega | x) \\ &= \int dA \int d\omega A p(A | \omega, x) q_\alpha(\omega) \quad \text{Variational approximation} \\ &= \int d\omega q_\alpha(\omega) \int dA A p(A | \omega, x) \\ &= \int d\omega q_\alpha(\omega) \bar{A}(x, \omega) \quad \text{network output} \end{aligned}$$

What about the uncertainty?

BNN variance of prediction

Variance of prediction

$$\begin{aligned}\sigma_{\text{tot}}^2(x) &\equiv \langle (A - \langle A \rangle)^2 \rangle = \int dA (A - \langle A \rangle)^2 p(A | x) \\ &= \int dA \int d\omega (A - \langle A \rangle)^2 p(A | \omega, x) q_\alpha(\omega) \\ &= \int d\omega q_\alpha(\omega) \left[\int dA A^2 p(A | \omega, x) - 2\langle A \rangle \int dA A p(A | \omega, x) + \langle A \rangle^2 \int dA p(A | \omega, x) \right] \\ &= \underbrace{\int d\omega q_\alpha(\omega) \left(\overline{A^2}(x, \omega) - \overline{A}(x, \omega)^2 \right)}_{\text{1} = \sigma_{\text{model}}^2(x)} + \underbrace{\int d\omega q_\alpha(\omega) \left(\overline{A}(x, \omega) - \langle A \rangle \right)^2}_{\text{2} = \sigma_{\text{pred}}^2(x)}\end{aligned}$$

Types of uncertainties

1 Predictive uncertainty

$$\sigma_{\text{pred}}^2(x) = \int d\omega q_{\alpha}(\omega) (\bar{A}(x, \omega) - \langle A \rangle)^2$$

→ following (*) this **vanishes** for $q(\omega) \rightarrow \delta(\omega - \omega_0)$

→ with **precise** training data

→ requires more and better training → **decreases** with more training data

→ represents **statistical uncertainty** (*epistemic uncertainty*)

for physicists

$$\sigma_{\text{pred}} \rightarrow \sigma_{\text{stat}}$$

Types of uncertainties

2 Likelihood related uncertainty

$$\sigma_{\text{like}}^2(x) = \int d\omega q_\alpha(\omega) \left(\overline{A^2}(x, \omega) - \overline{A}(x, \omega)^2 \right) = \int d\omega q_\alpha(\omega) \sigma_{\text{like}}^2(\omega, x) \equiv \langle \sigma_{\text{like}}^2(\omega) \rangle$$

→ already occurs without sampling → does **not vanish** for $q(\omega) \rightarrow \delta(\omega - \omega_0)$

→ However, **vanishes** if amplitude is **perfectly known** → $p(A | \omega, x) \rightarrow \delta(A - A_0)$

$$\overline{A}(x, \omega) = \int dA A p(A | \omega, x) \rightarrow A_0 \quad \overline{A^2}(x, \omega) = \int dA A^2 p(A | \omega, x) \rightarrow A_0^2$$

→ Corresponds to the **Gaussian uncertainty** in the **heteroscedastic loss**

→ reaches constant for perfect training

→ reflects stochastic training data, bad hyperparameters, model expressivity,...

→ represents **systematic uncertainty** (*aleatoric uncertainty*)

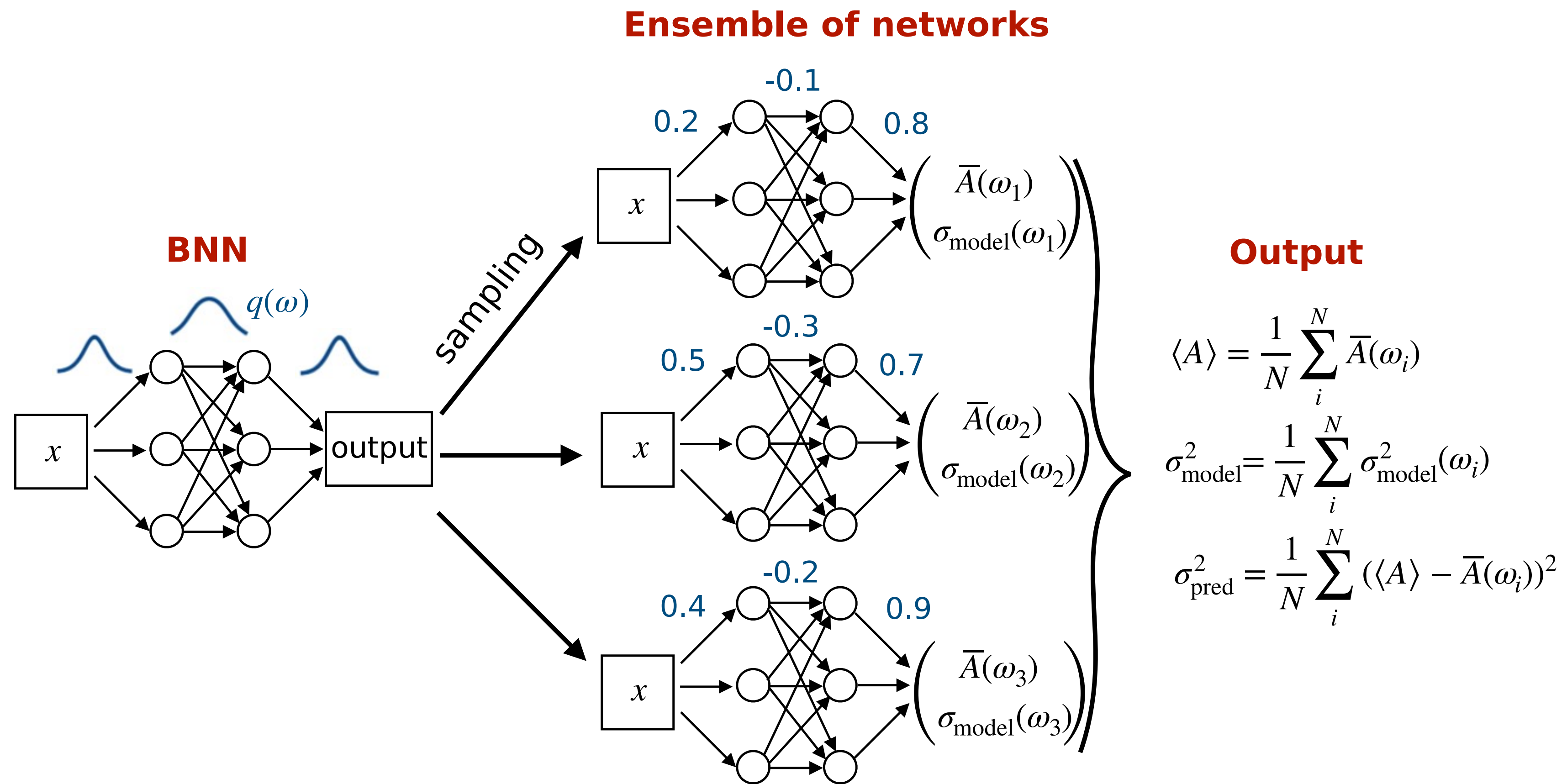
for physicists

$$\sigma_{\text{like}} \rightarrow \sigma_{\text{sys}}$$

network constructs
uncertainty

Also works if
training data
has no uncertainty

Bayesian Regression

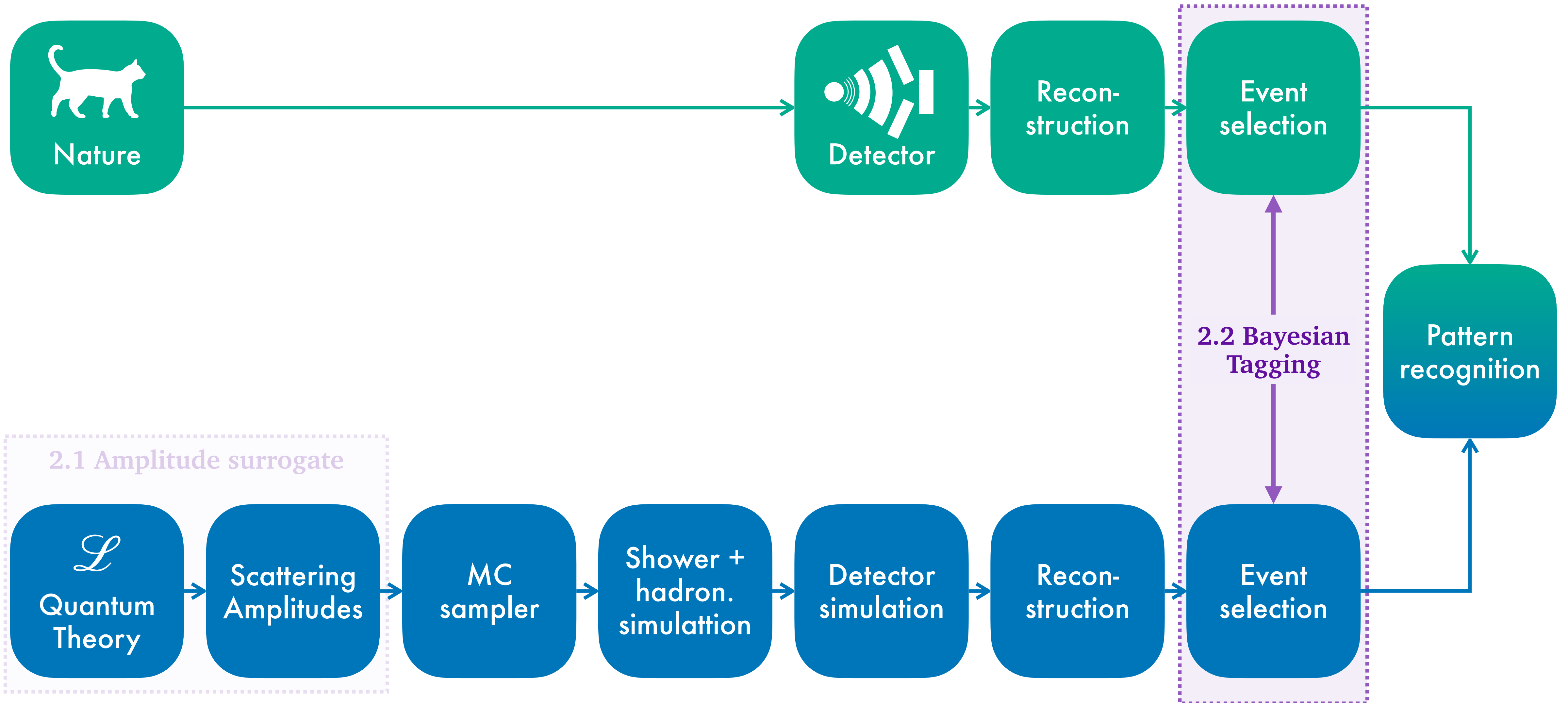


$$\mathcal{L}_{\text{BNN}} = \text{KL}(q_\alpha(\omega), p(\omega)) + \int d\omega q_\alpha(\omega) \sum_j \left[\frac{(\bar{A}(x_j, \omega) - A_j^{\text{truth}})^2}{2\sigma_{\text{syst}}(x_j, \omega)} - \log \sigma_{\text{syst}}(x_j, \omega) \right]$$

Example II

Bayesian Tagger

LHC analysis + ML



Bayesian Jet Tagger

Prediction of Bayesian classifier

$$\mu_{\text{pred}} \equiv p(c | x) = \int d\omega q_{\alpha}(\omega) p(c | \omega, x)$$

BNN classifier output

Class prediction
is only output

Statistical uncertainty

$$\sigma_{\text{pred}}^2 \equiv \sigma_{\text{stat}}^2 = \int d\omega q_{\alpha}(\omega) (p(c | x, \omega) - \mu_{\text{pred}})^2$$

only statistic
uncertainty

no additional
parameter for
syst. uncertainty

BCE loss

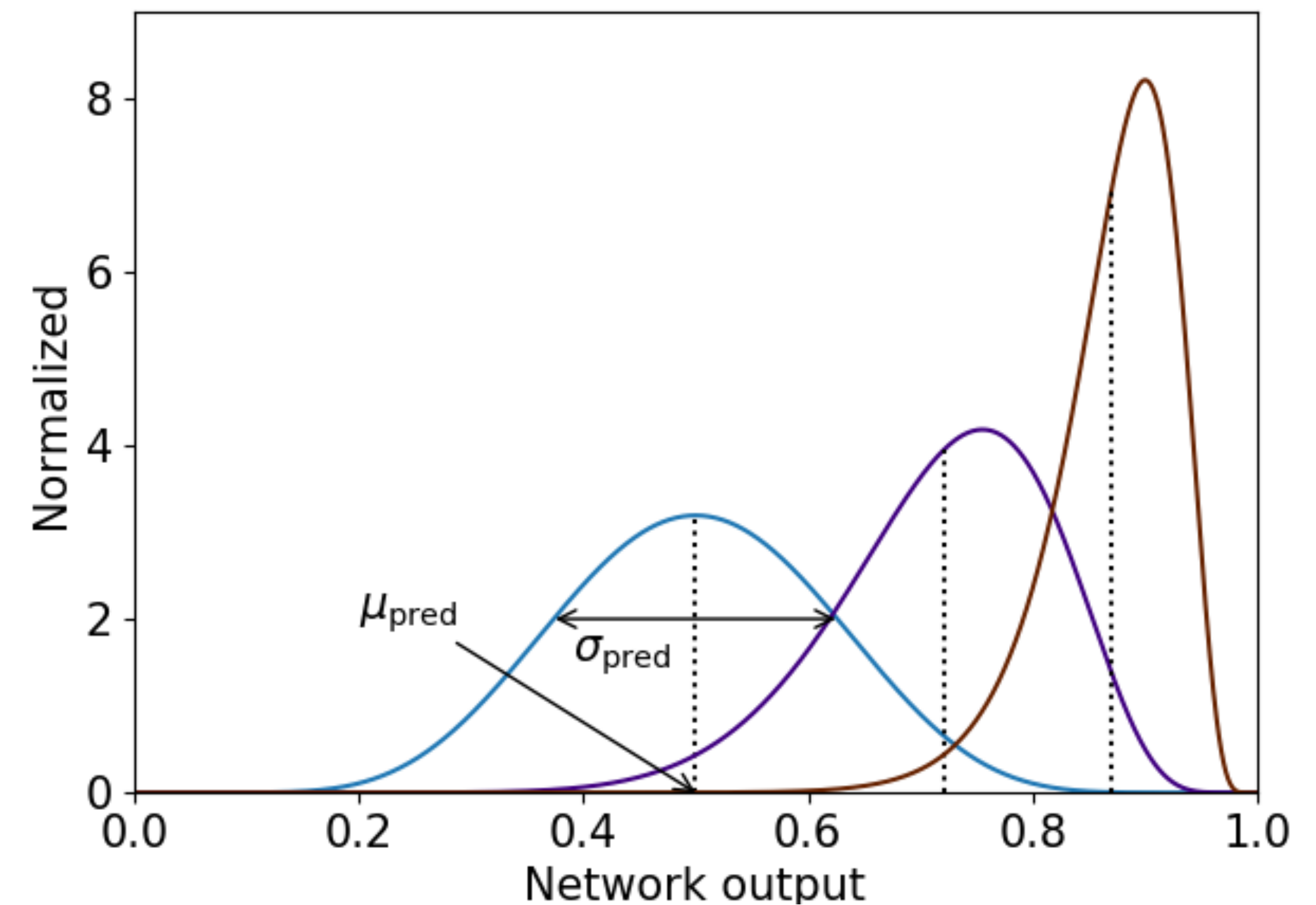
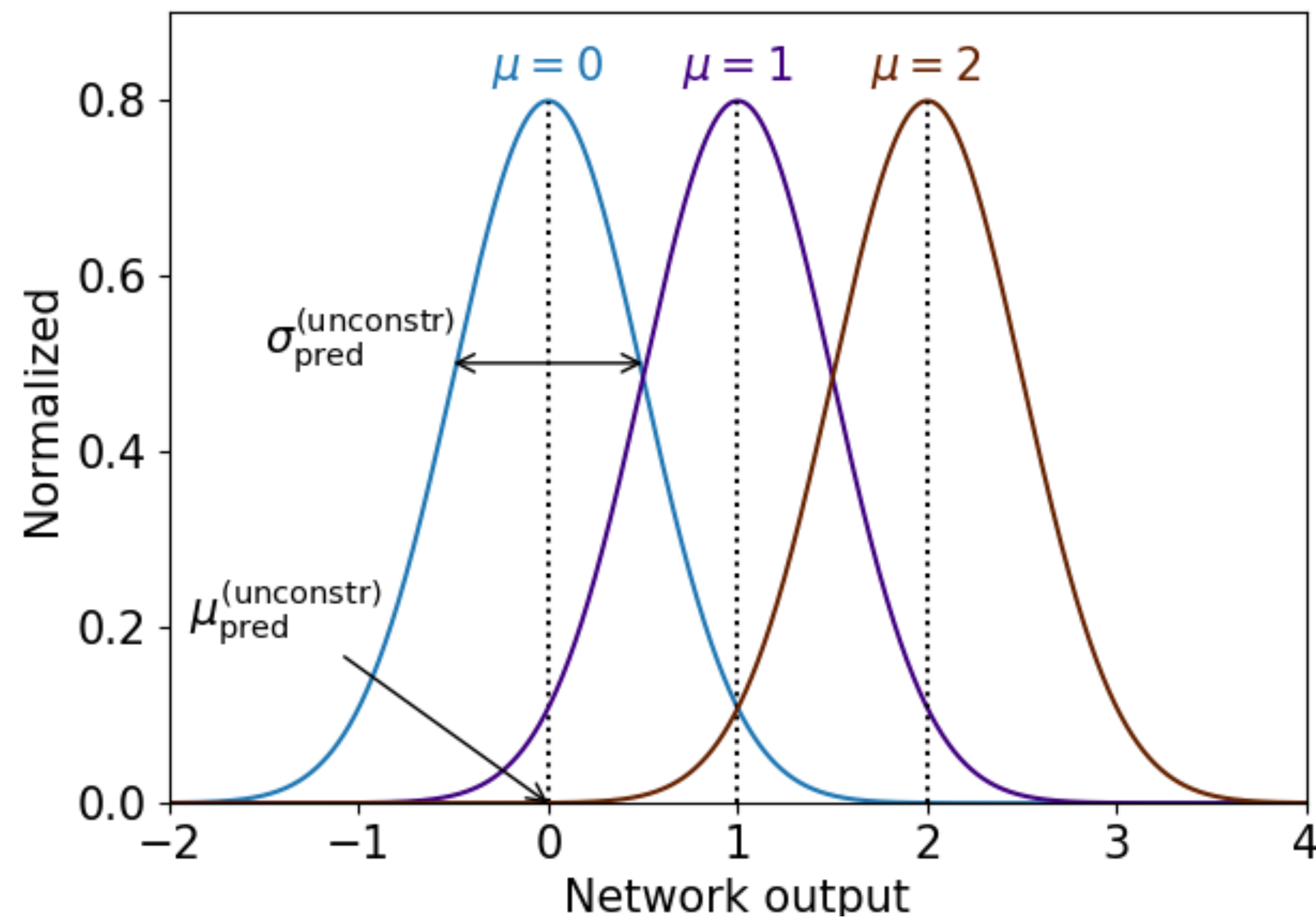
$$\mathcal{L}_{\text{BNN}} = \text{KL}(q_{\alpha}(\omega), p(\omega)) + \int d\omega q_{\alpha}(\omega) \left[\sum_j y_j^{\text{truth}} \log p(c | x_j, \omega) + (1 - y_j^{\text{truth}}) \log(1 - p(c | x_j, \omega)) \right]$$

We have to be careful though...

Bayesian Jet Tagger

Classifiers have a non-linearity (sigmoid) in last layer to get $p(c | x) \in [0,1]$

→ Leads to distortion of Gaussian shapes in general ⚠



Have to check if still Gaussian after sigmoid!!!

Bayesian Jet Tagger

