

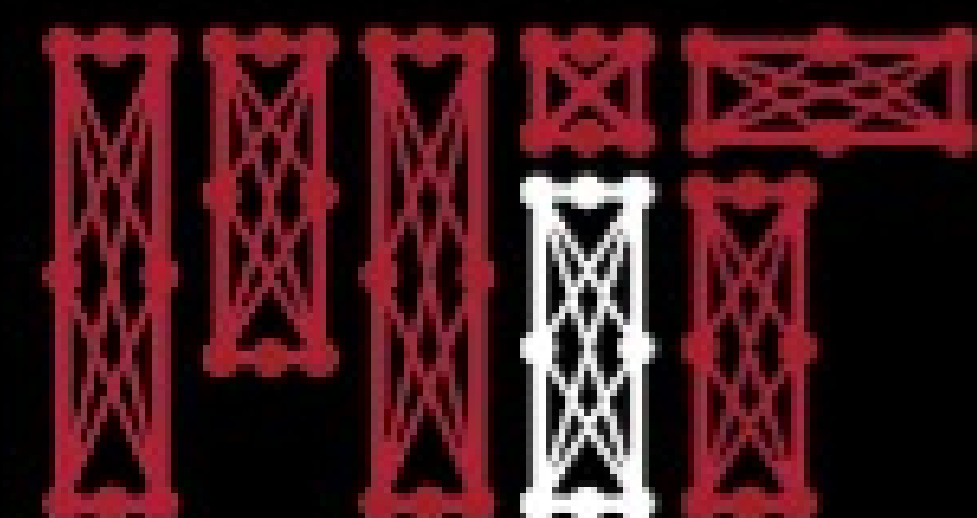


Evidential Deep Learning

Alexander Amini

MIT 6.S191

January 26, 2021

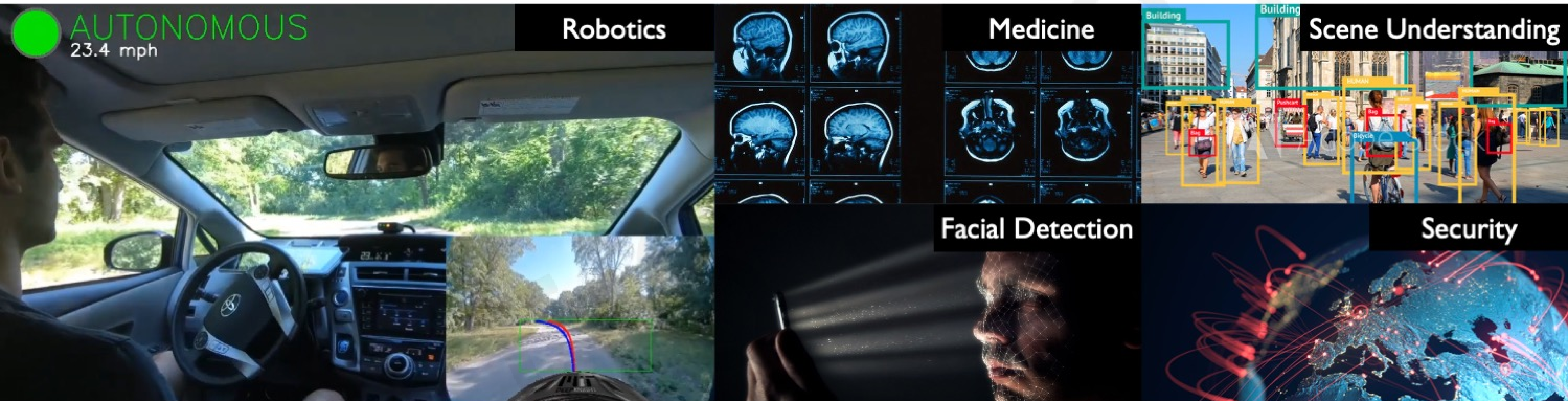


6.S191 Introduction to Deep Learning

introtodeeplearning.com [@MITDeepLearning](https://twitter.com/MITDeepLearning)



Motivation: uncertainty in learning



Safety critical domains require **fast, scalable, and calibrated** uncertainty estimation

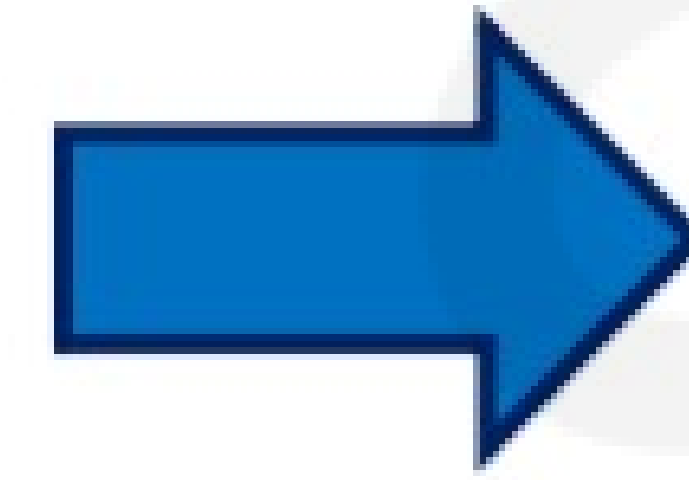
Evidential deep learning → predict answer and amount of evidence (confidence)

Neural networks: expectation vs reality

Expectation:
Training on a your dataset



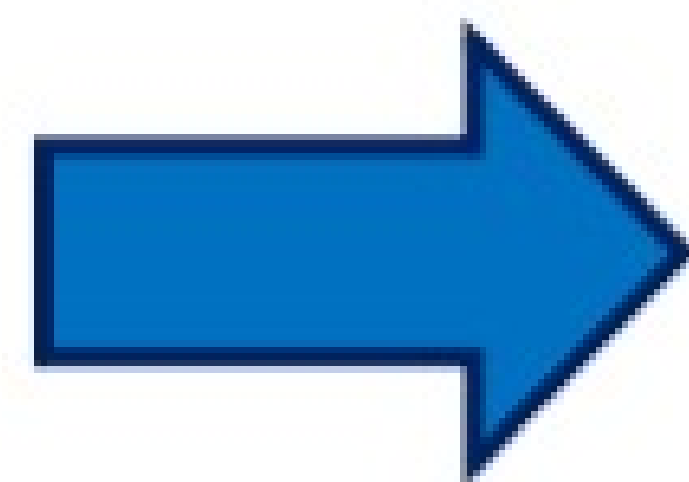
Dogs



Reality:
Testing in reality



Driving



Neural networks: expectation vs reality

Dogs

Expectation:
Training on a your dataset



Reality:
Testing in reality



"All models are wrong, but some — *that know when they can be trusted* — are useful!"

Driving



- George E.P. Box (Adapted)

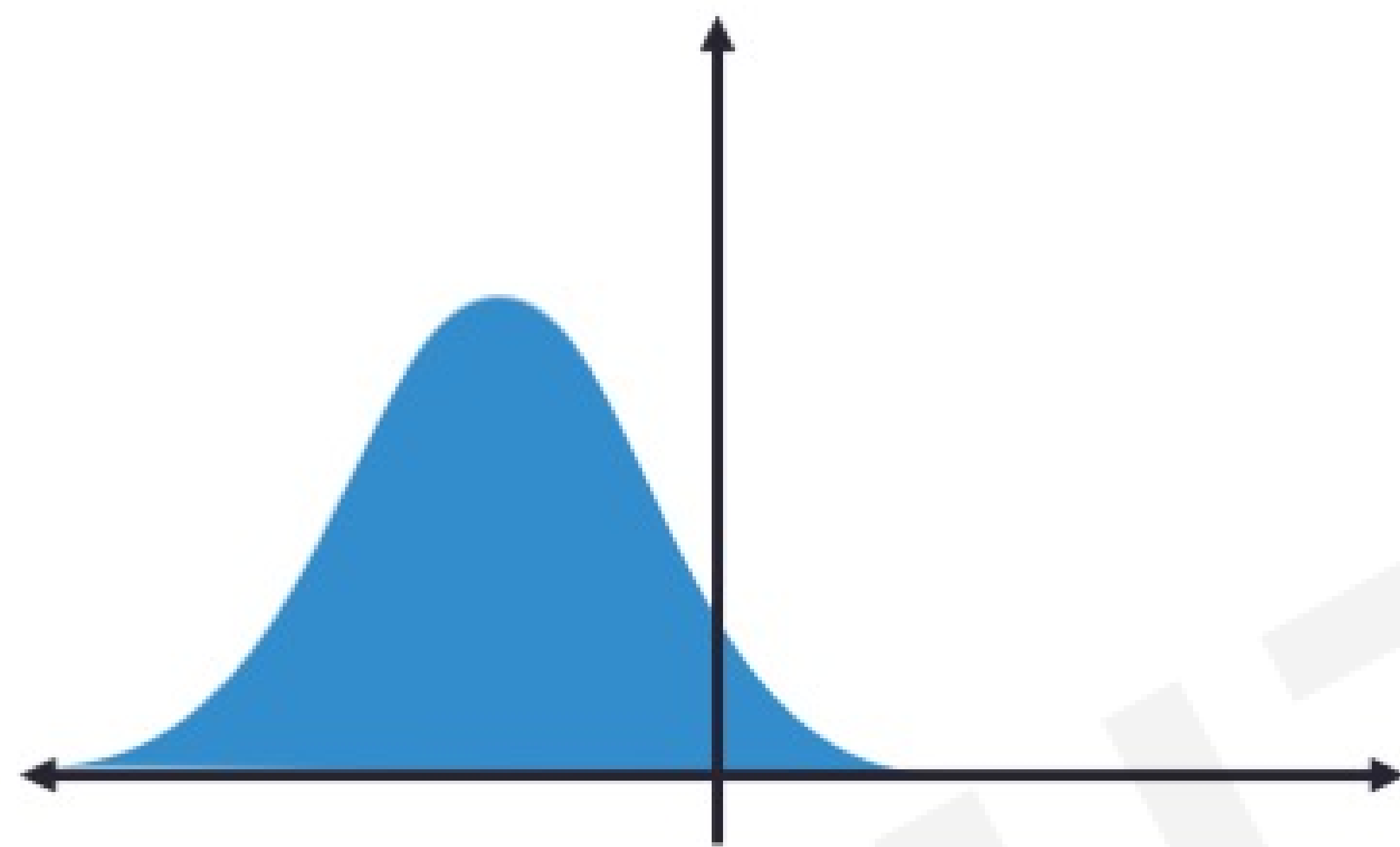
Knowing when we don't know is hard

...even for humans!



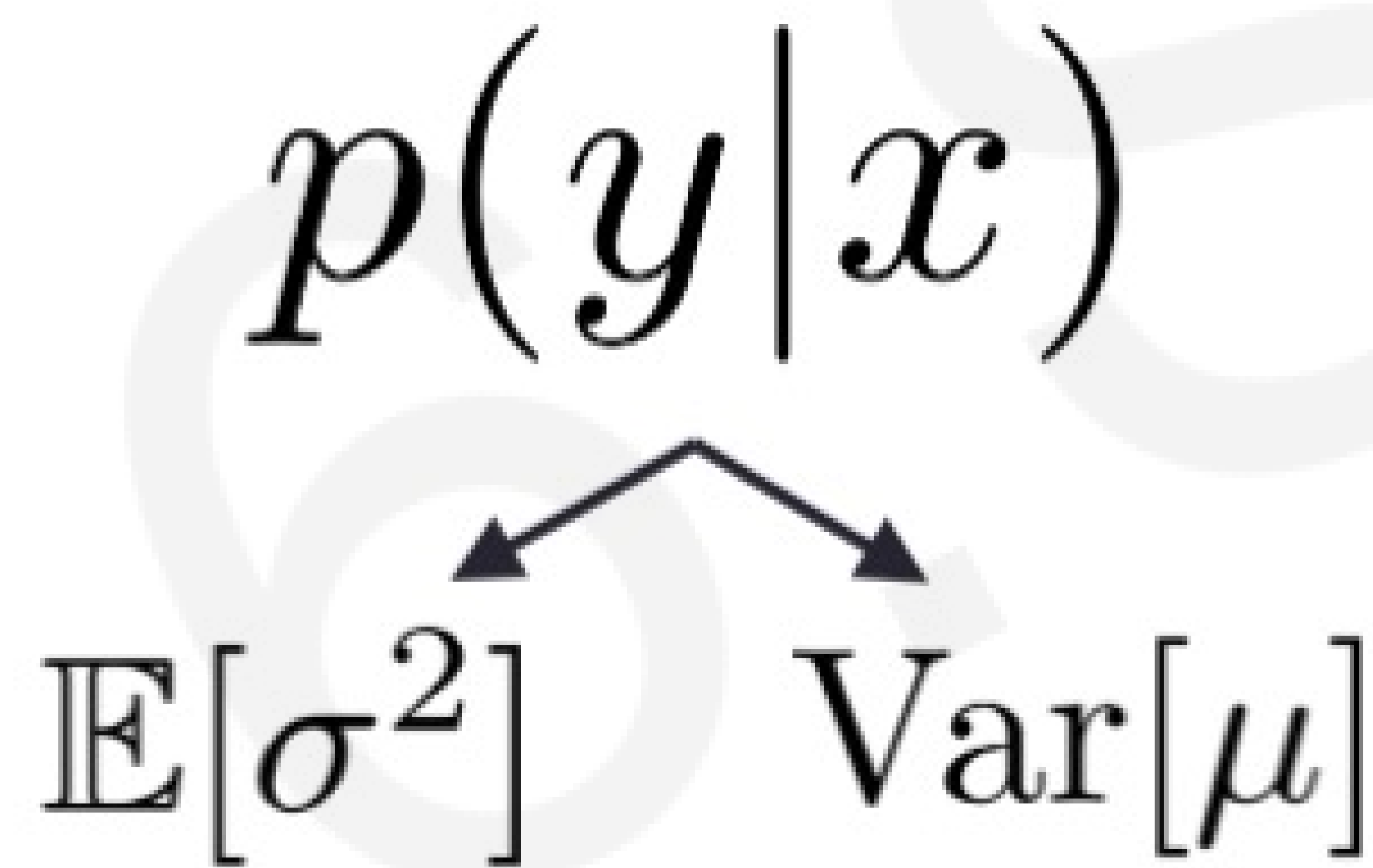
Topics for today

Learning probability distributions



Model distributions over labels:
Softmax (discrete) & Gaussian
(continuous)

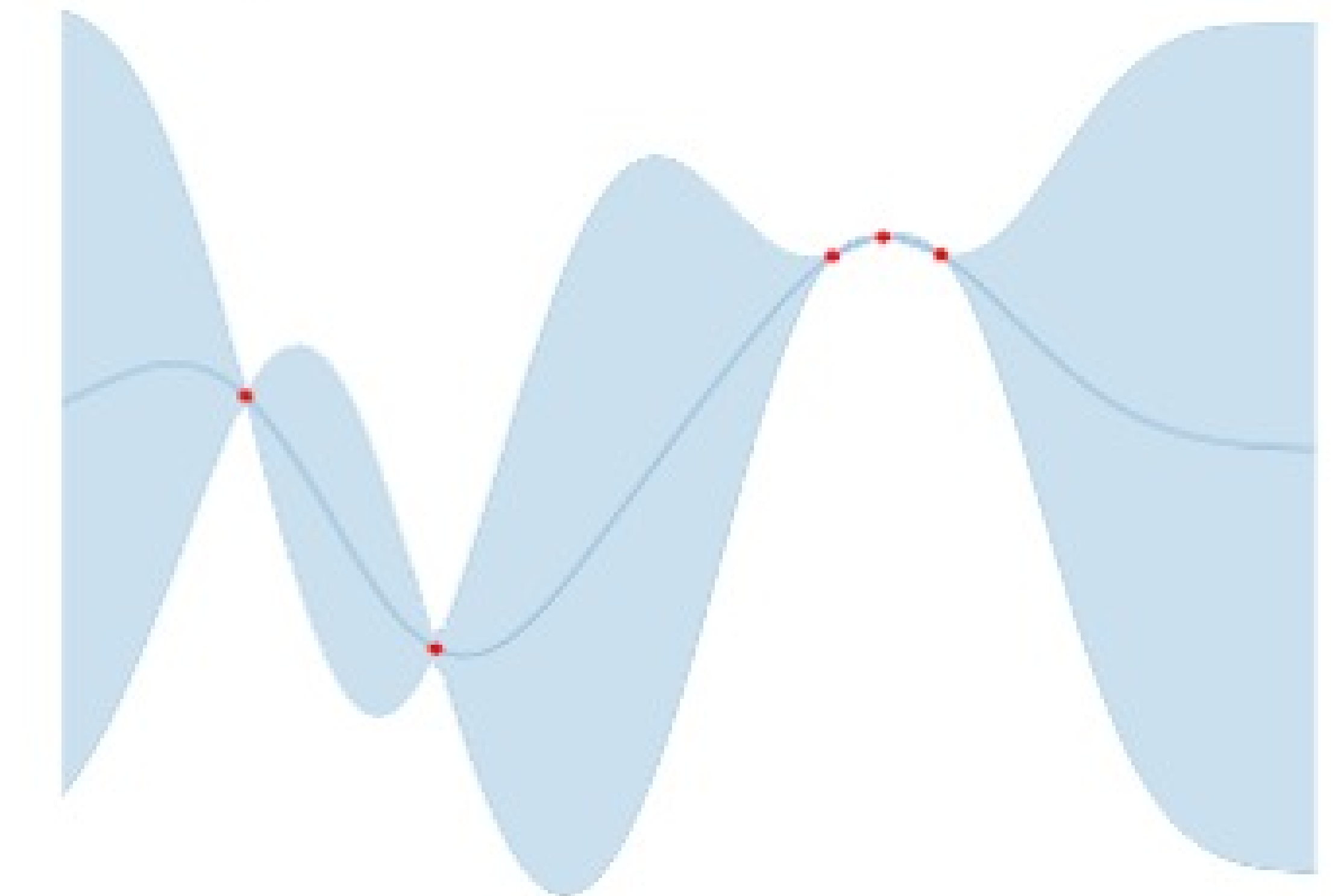
Different sources of uncertainty

$$p(y|x)$$


The diagram shows the conditional probability $p(y|x)$ at the top. Two arrows point downwards from it to the terms $\mathbb{E}[\sigma^2]$ and $\text{Var}[\mu]$.

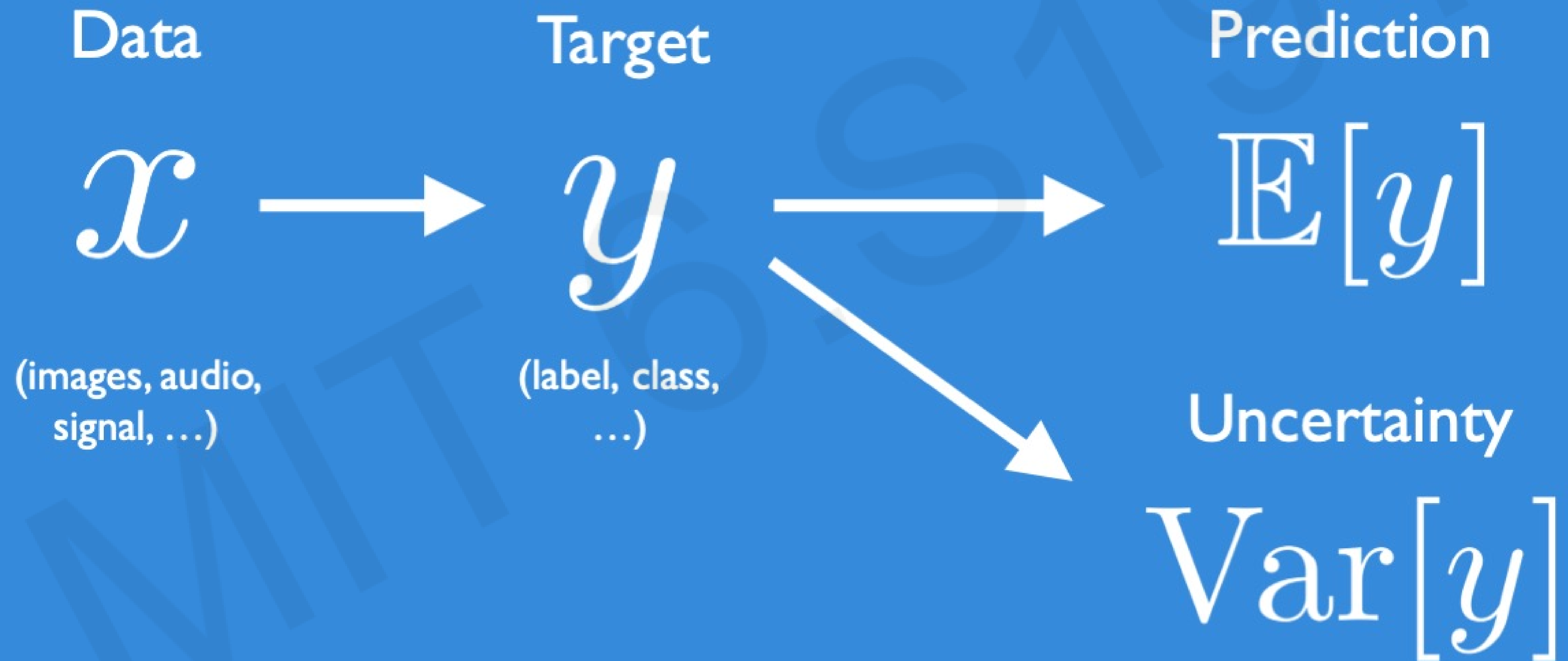
Data (aleatoric) uncertainty vs.
Model (epistemic) uncertainty.

Fast and scalable uncertainty estimation



Evidential deep learning
Uncertainty modelling for quickly
estimating confidence

Probabilistic learning



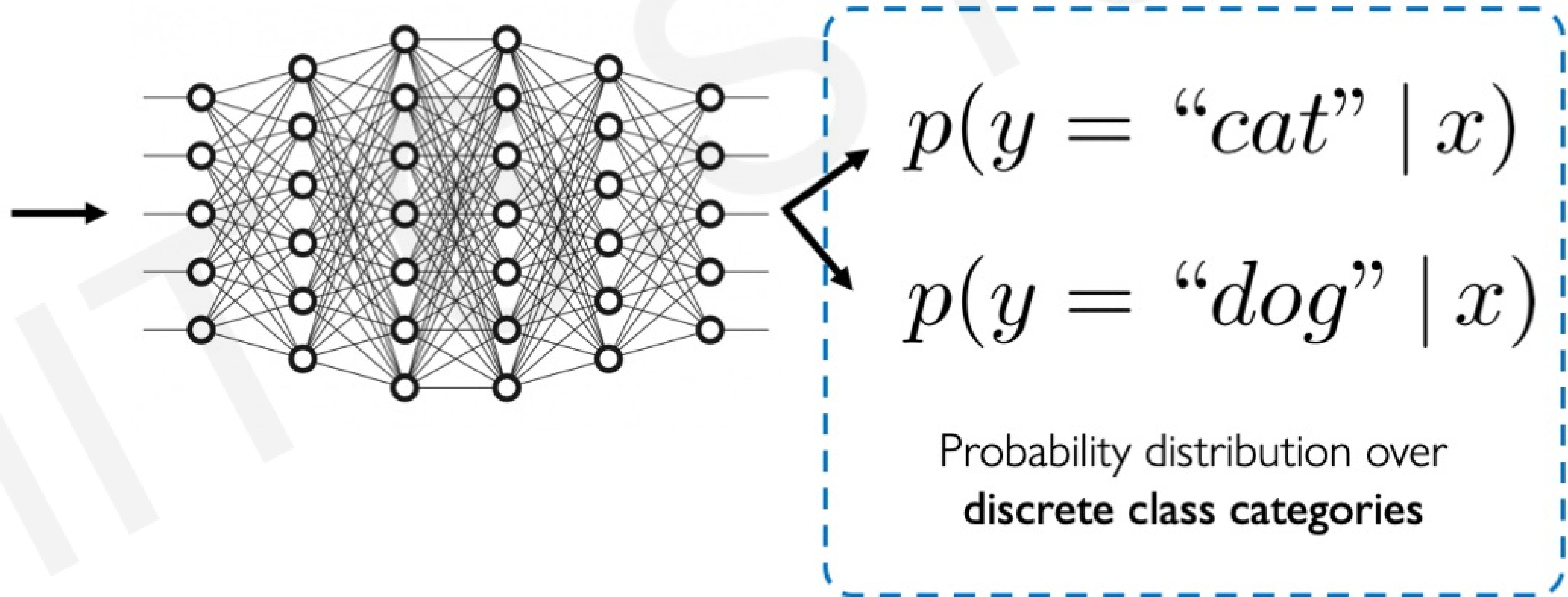
Learning probabilistic outputs



Wait, haven't we already learned this?!



x

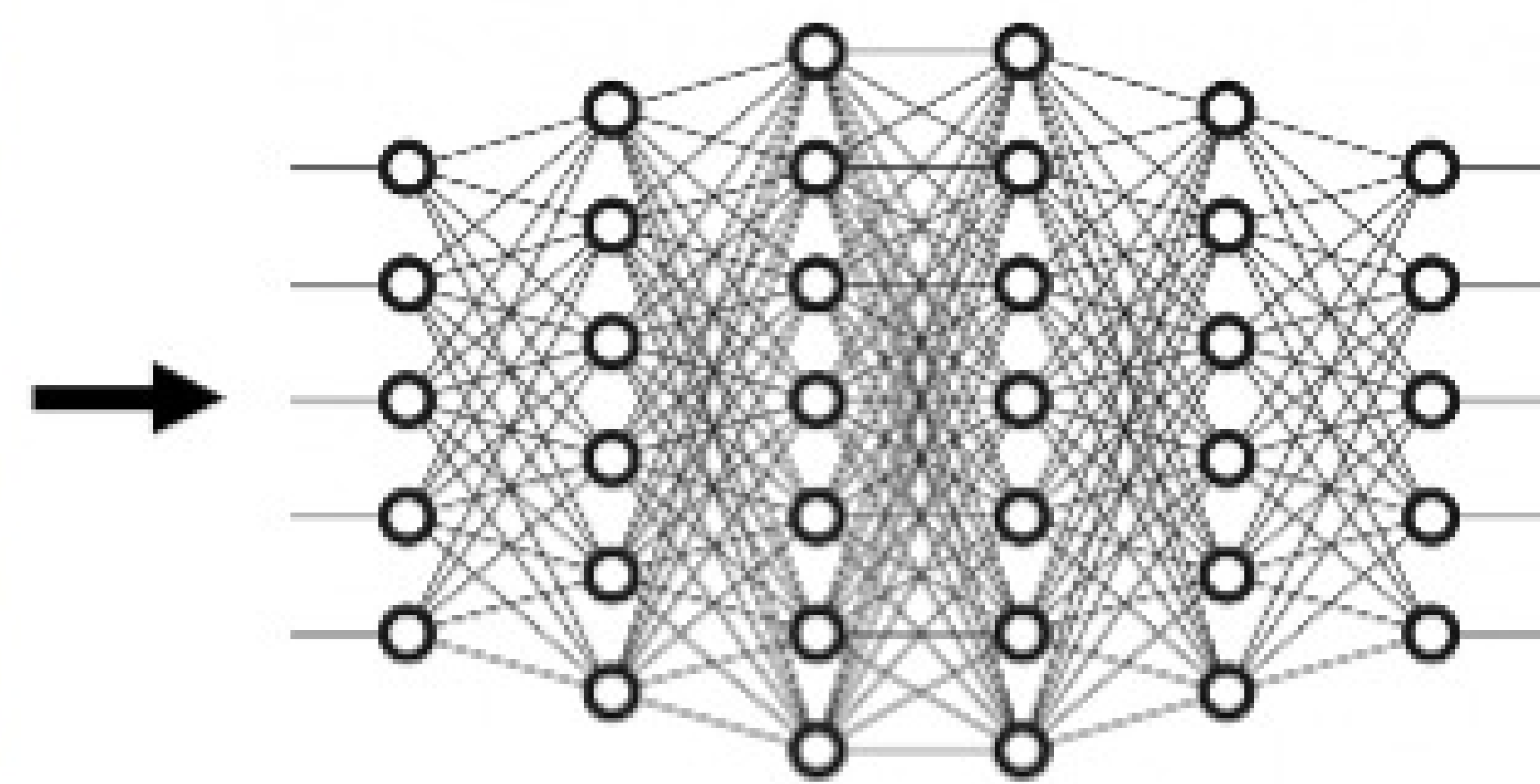


Learning discrete class targets

Classification



x



$$p(y = \text{"cat"} \mid x)$$

$$p(y = \text{"dog"} \mid x)$$

Why?

$$\underline{y} \sim \underline{\text{Categorical}}(\underline{p})$$

Class Labels

Likelihood function

Distribution parameters (probabilities)

$$f(y = y_i \mid \mathbf{p}) = p_i$$

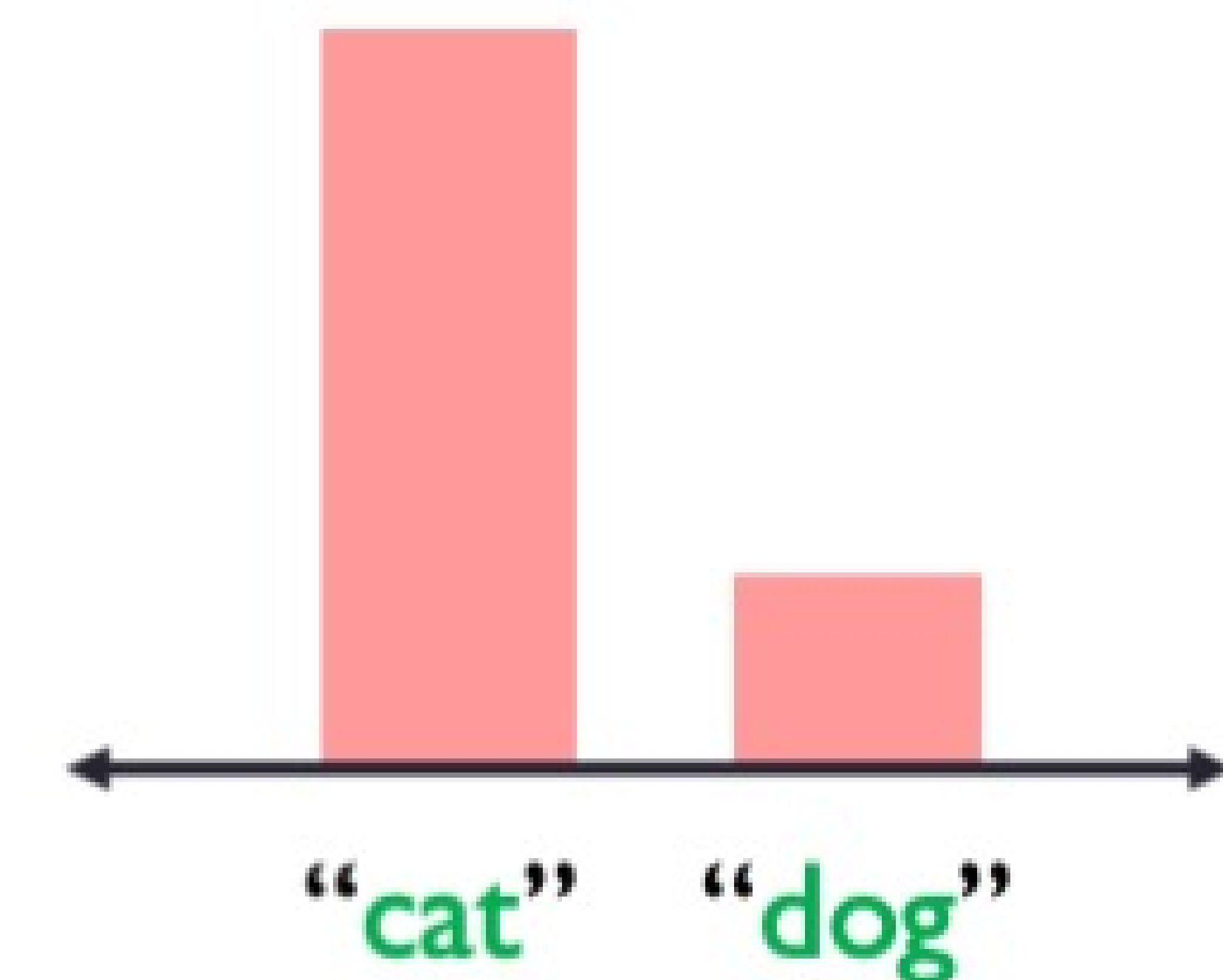
Activation: softmax(z)

$$\sigma(\vec{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

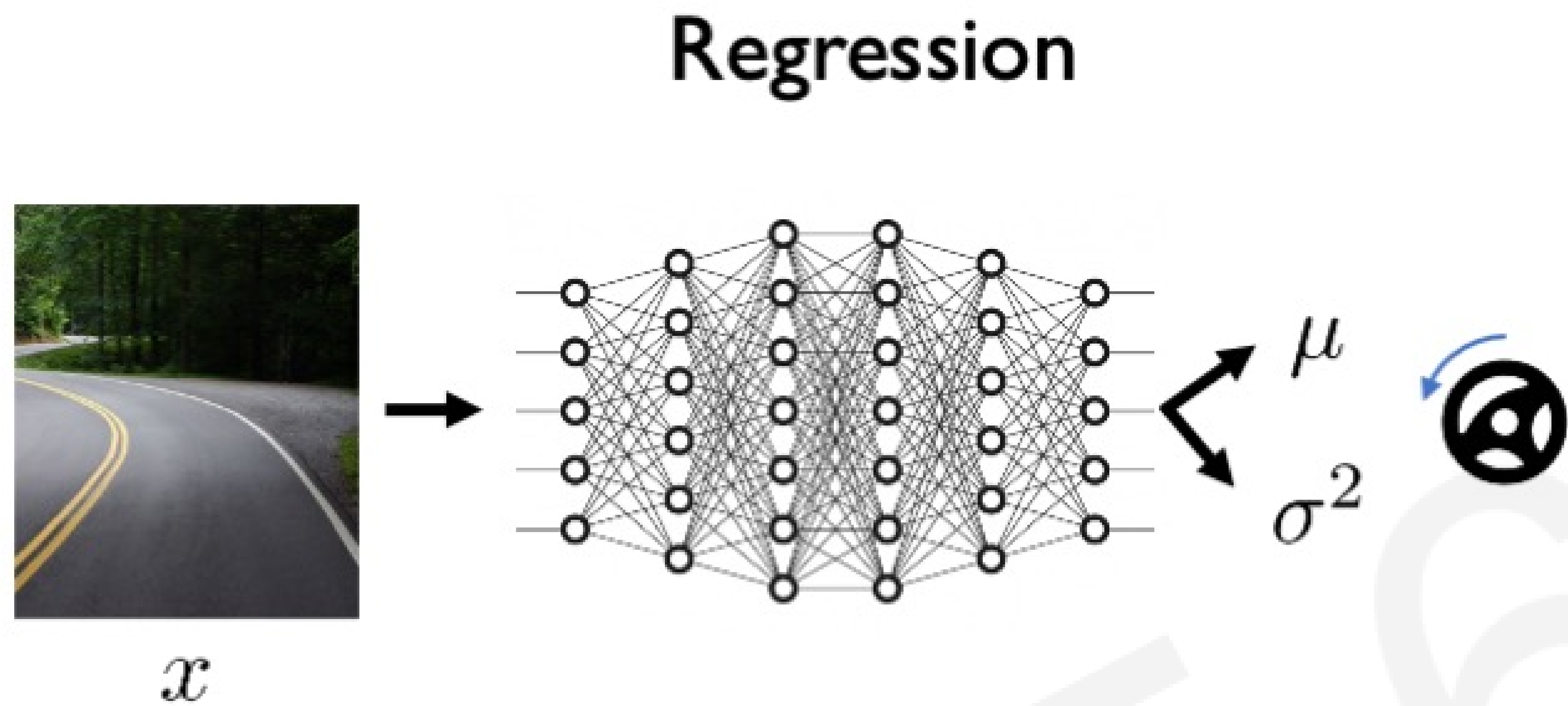
Loss:

Neg. Log Likelihood (Cross Entropy)

$$-\sum_{i=1}^K y_i \log p_i$$



Learning continuous class targets



Why?

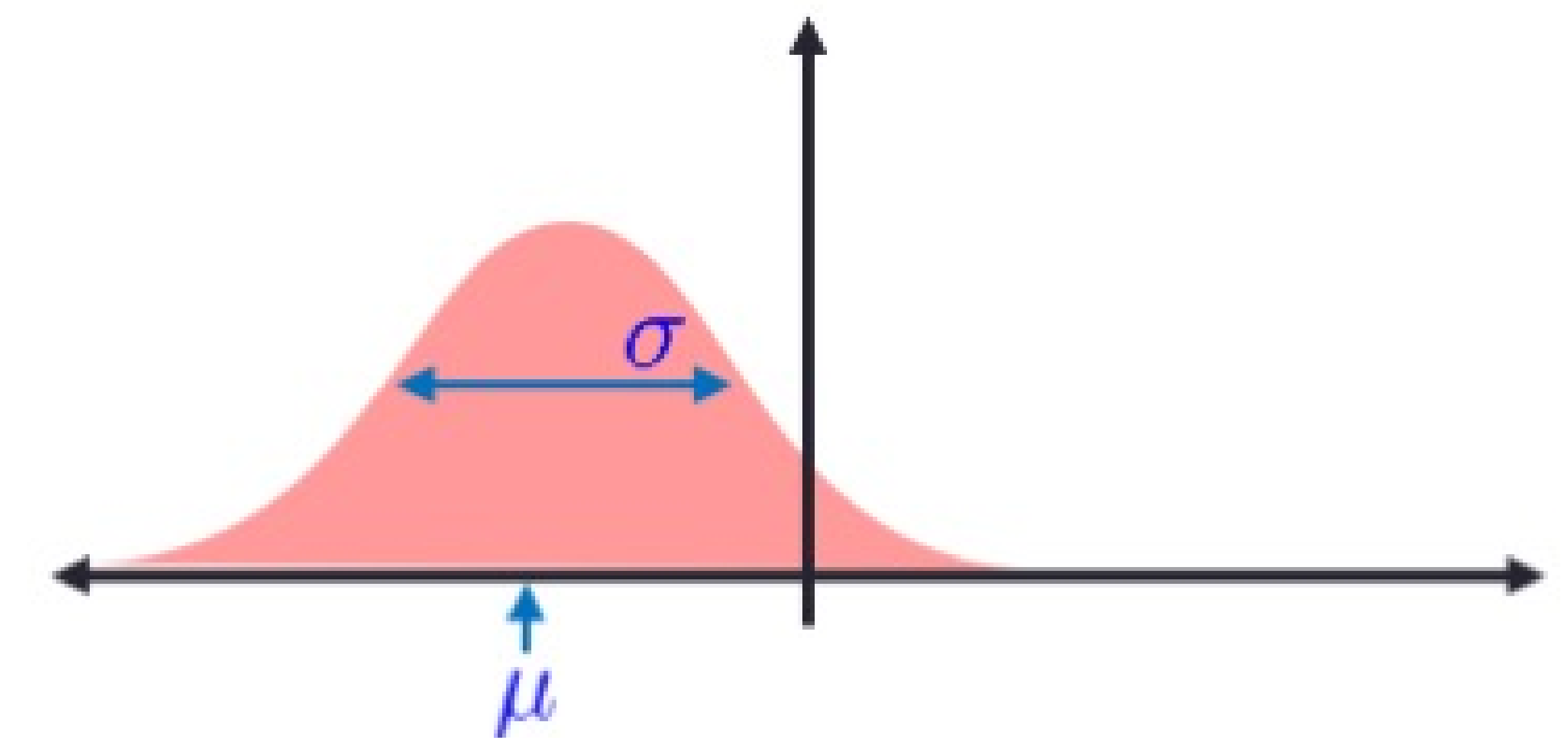
$y \sim \text{Normal}(\mu, \sigma^2)$

Target Labels Likelihood function Distribution parameters





$$f(y | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y-\mu)^2}{2\sigma^2}\right)$$

Activation: $\mu \in \mathbb{R}$
 $\sigma > 0$ $\mu = z_\mu$
 $\sigma = \exp(z_\sigma)$

Loss: Neg. Log Likelihood $-\log(\mathcal{N}(y|\mu, \sigma^2))$



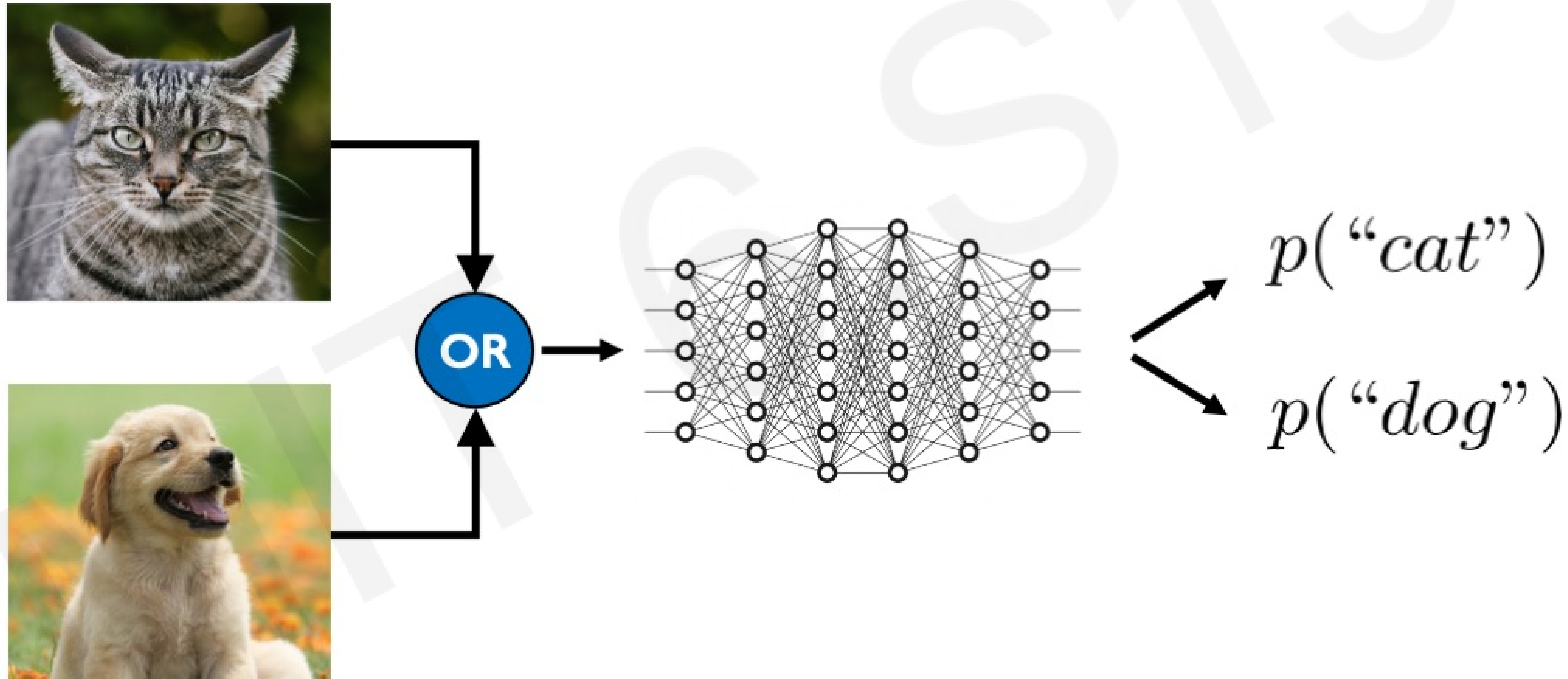
Likelihood estimation in deep learning

	Classification (discrete)	Regression (continuous)
Targets	$y \in \{1, \dots, K\}$	$y \in \mathbb{R}$
Likelihood	$y \sim \text{Categorical}(\mathbf{p})$  <code>tfp.distributions.Categorical(probs=p)</code>	$y \sim \text{Normal}(\mu, \sigma^2)$  <code>tfp.distributions.Normal(mu, sigma)</code>
Parameters	$\mathbf{p} = \{p_1, \dots, p_K\}$	(μ, σ^2)
Constraints	$\sum_i p_i = 1; \quad p_i > 0$	$\mu \in \mathbb{R}; \quad \sigma > 0$
Loss function	Cross Entropy $-\sum_{i=1}^K y_i \log p_i$  <code>dist.cross_entropy(y)</code>	Negative Log-Likelihood $-\log(\mathcal{N}(y \mu, \sigma^2))$  <code>-1 * dist.log_prob(y)</code>

Likelihood vs Confidence

⚠ WARNING: ⚠

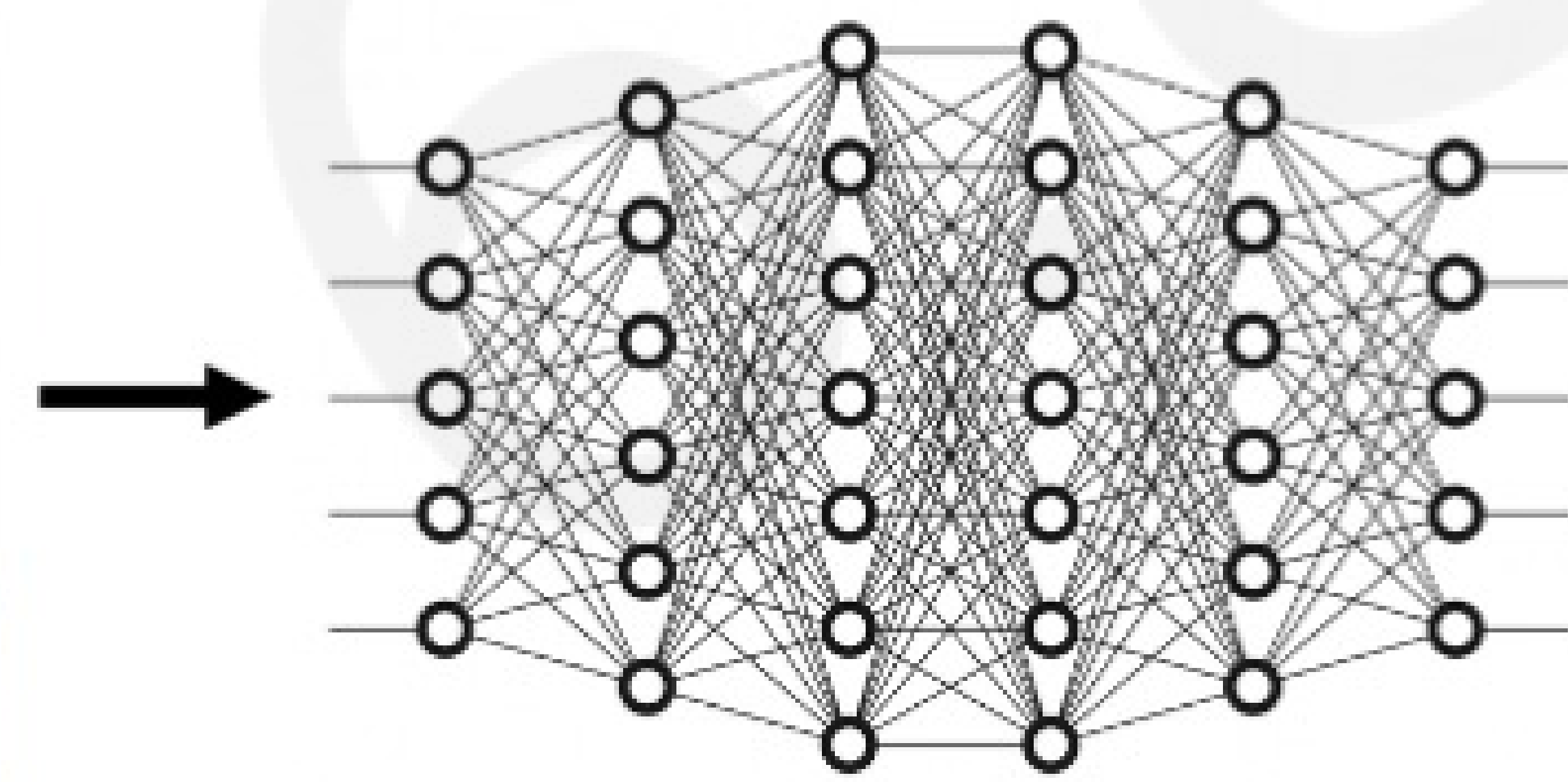
Do not mistake likelihood (probability) for model confidence



Likelihood vs Confidence

⚠ WARNING: ⚠

Do not mistake likelihood (probability) for model confidence



$$p(\text{"cat"}) = 0.5$$

$$p(\text{"dog"}) = 0.5$$

Likelihood vs Confidence

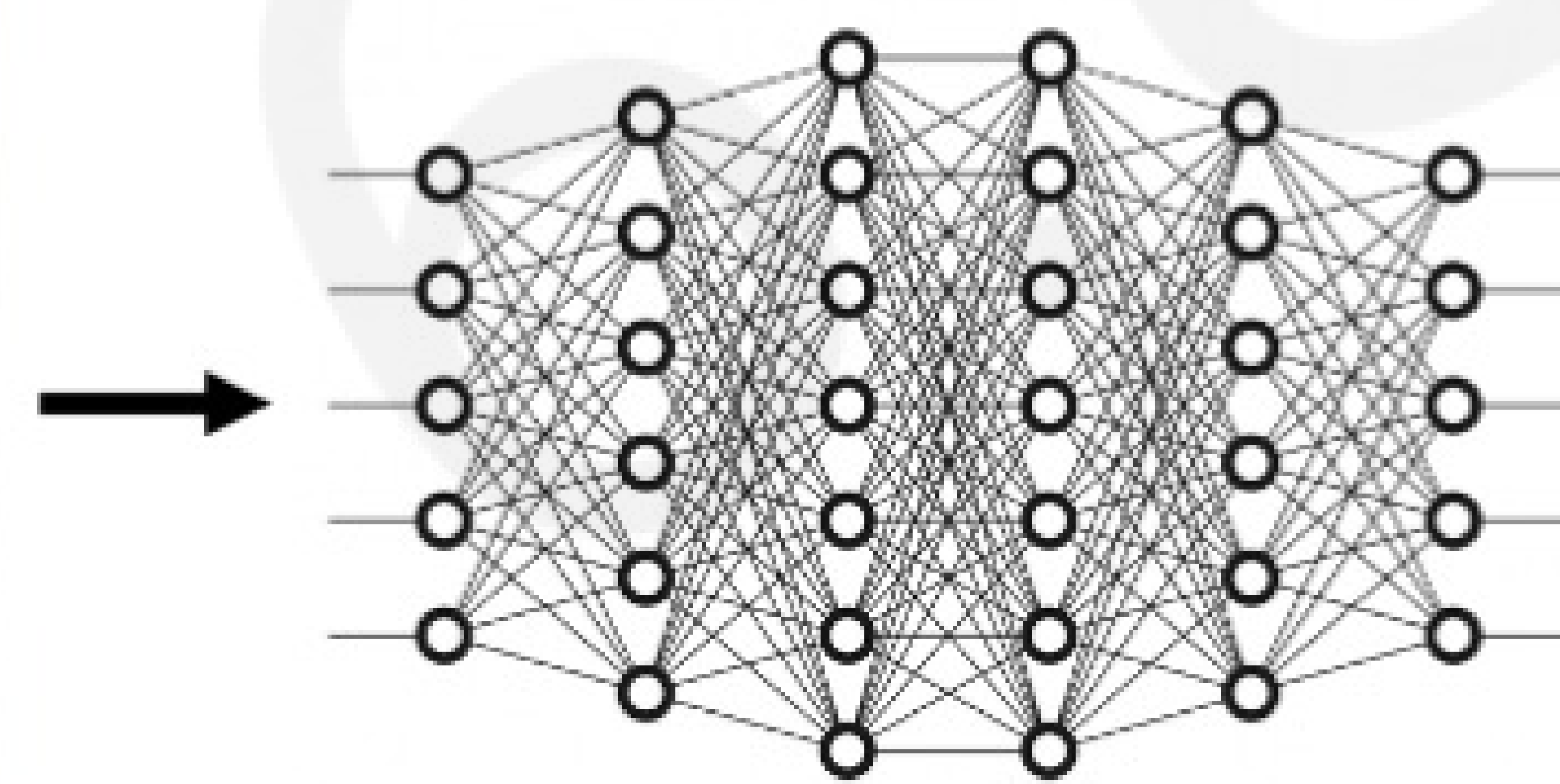


WARNING:



Do not mistake likelihood (probability) for model confidence

The output likelihoods will be unreliable if the input is unlike anything during training



$p(\text{"cat"})$

$p(\text{"dog"})$

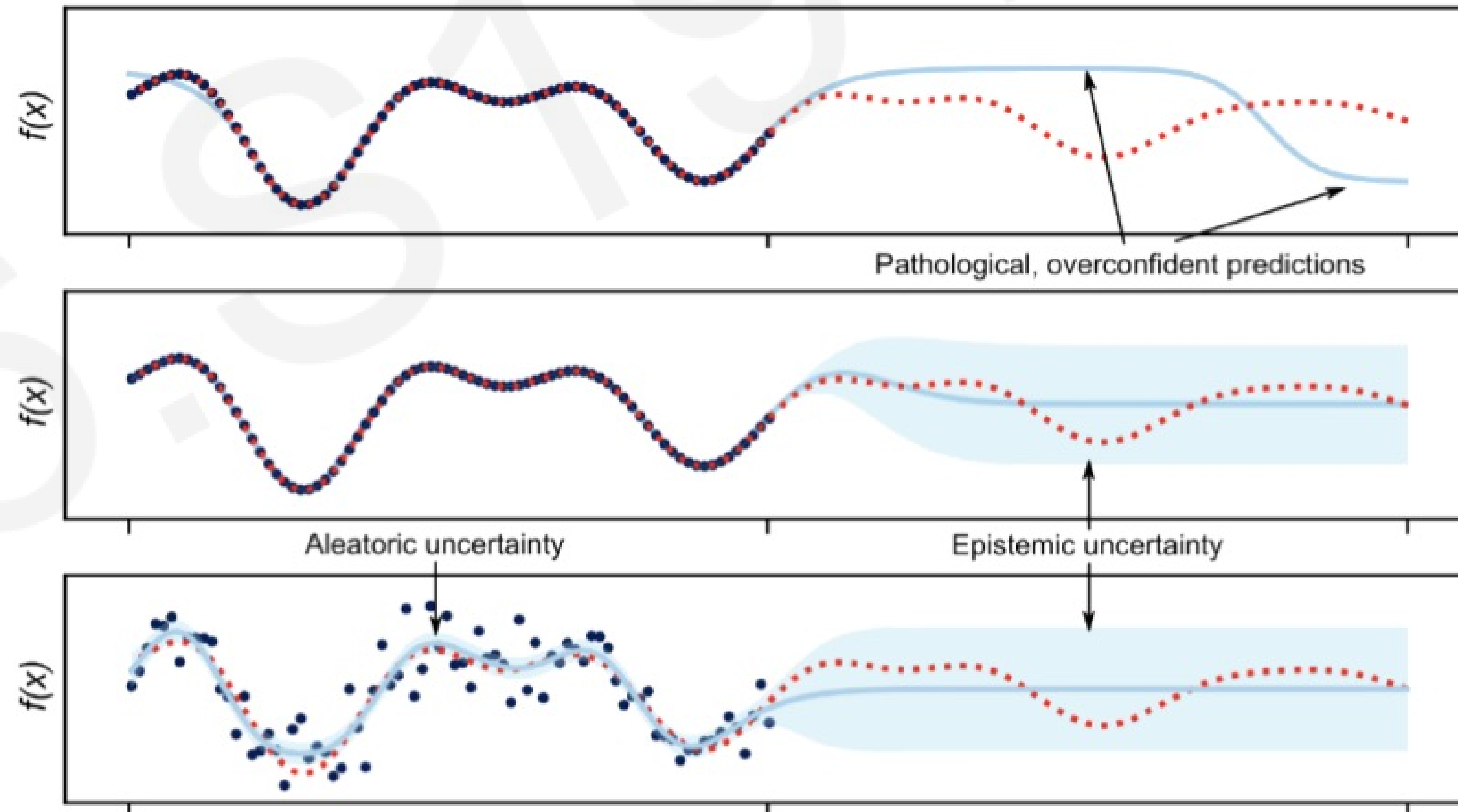
★ $p(\text{"cat"}) + p(\text{"dog"}) = 1$ ★

Types of uncertainty

..... Ground truth
 — Prediction
 Confidence interval
 ● Observations

Observed domain

Unobserved domain



Known Knowns

Things we are certain of

Known Unknowns

We know there are things we can't predict

Unknown Known

Others know but you don't know

Unknown Unknowns

Completely unexpected or unforeseeable events

Aleatoric vs Epistemic Uncertainty

Aleatoric Uncertainty



Data Uncertainty

Describes the confidence in the input data

High when input data is noisy

Cannot be reduced by adding more data

Epistemic Uncertainty



Model Uncertainty

Describes the confidence of the prediction

High when missing training data

Can be reduced by adding more data

Aleatoric vs Epistemic Uncertainty

Aleatoric Uncertainty



Data Uncertainty

Describes the confidence in the input data

High when input data is noisy

Cannot be reduced by adding more data

Epistemic Uncertainty



Model Uncertainty

Describes the confidence of the prediction

High when missing training data

Can be reduced by adding more data

Estimating epistemic uncertainty



Aleatoric uncertainty can be learned directly using neural networks



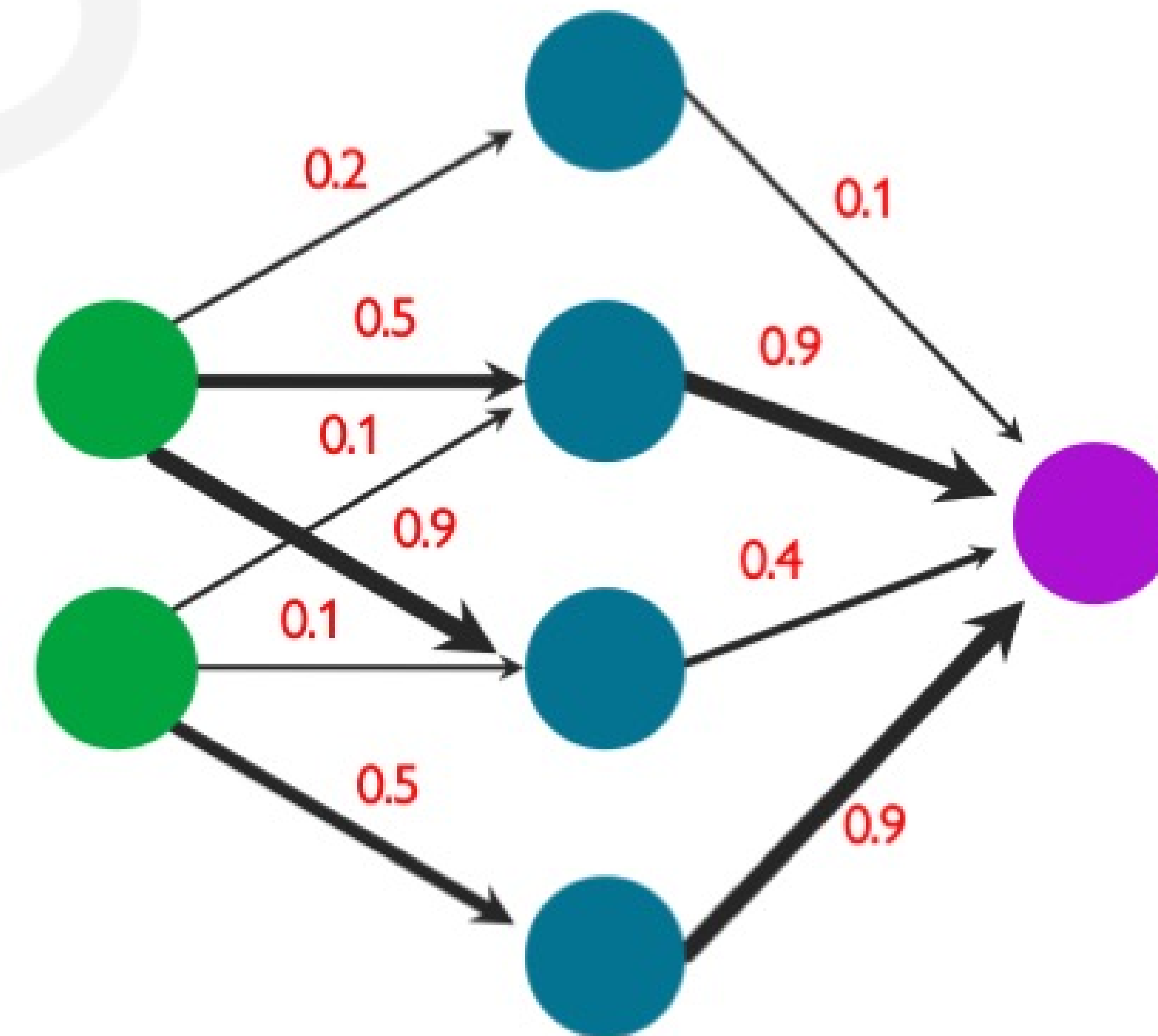
Epistemic uncertainty is **much more challenging** to estimate



How can a model understand when it does not know the answer?

One solution:

Don't train **deterministic NN**, but instead train a Bayesian NN!



Estimating epistemic uncertainty



Aleatoric uncertainty can be learned directly using neural networks



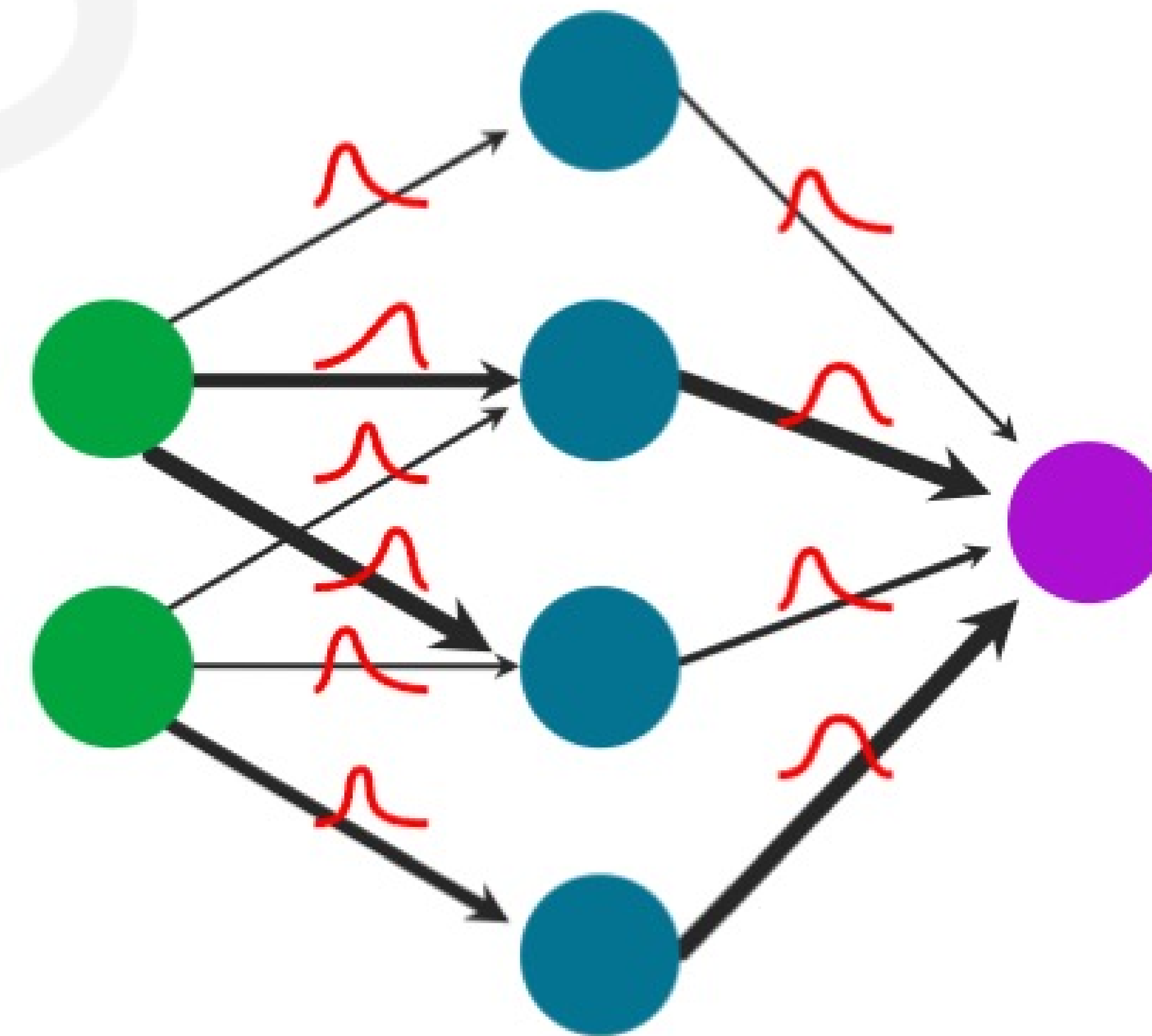
Epistemic uncertainty is **much more challenging** to estimate



How can a model understand when it does not know the answer?

One solution:

Don't train deterministic NN, but instead train a **Bayesian NN!**



Bayesian deep learning for uncertainty

Deterministic neural networks (NNs) learn a fixed set of weights,

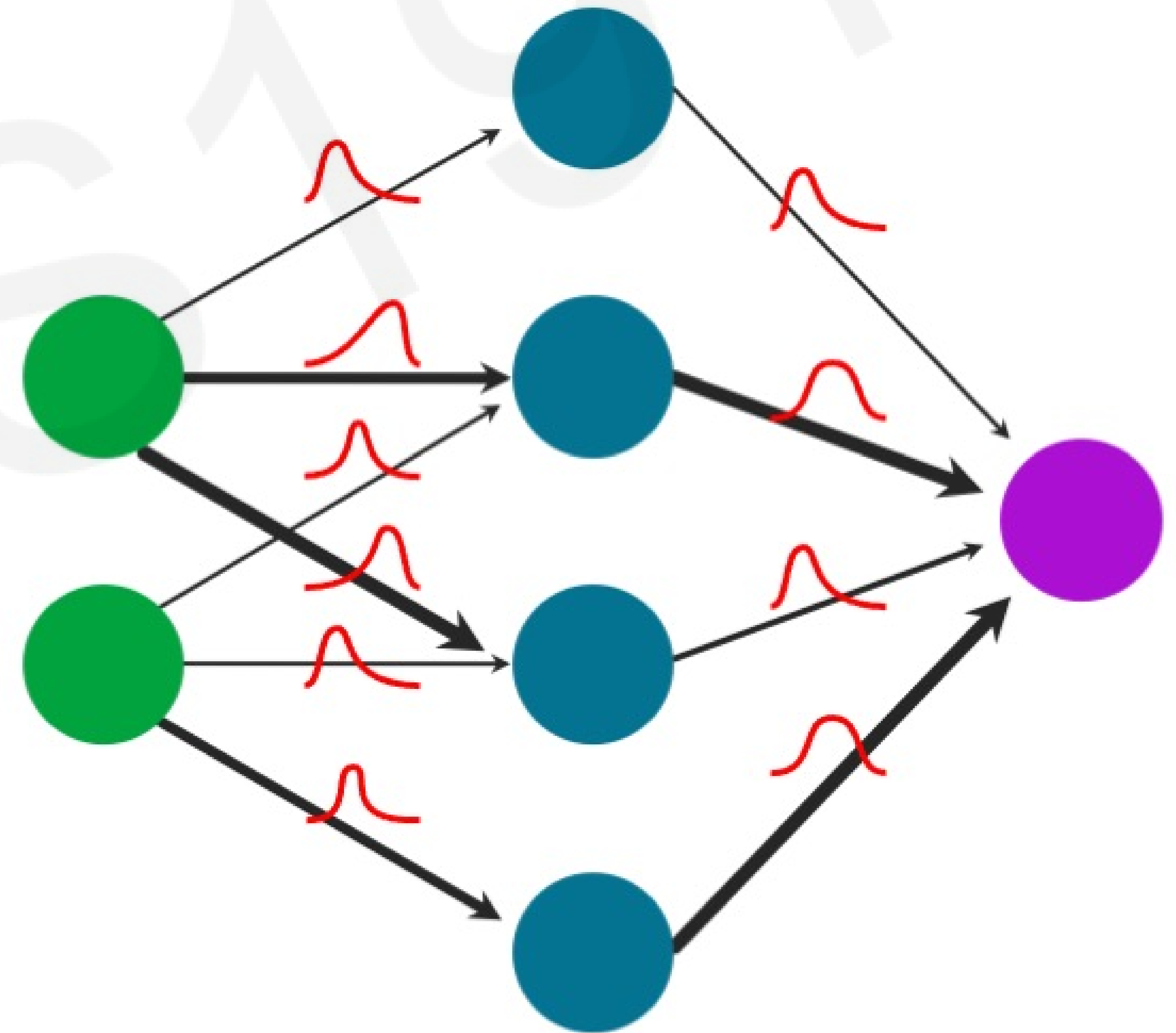
$$\mathbf{W}$$

Bayesian neural networks aim to learn a posterior over weights,

$$P(\mathbf{W}|\mathbf{X}, \mathbf{Y})$$

$$P(\mathbf{W}|\mathbf{X}, \mathbf{Y}) = \frac{P(\mathbf{Y}|\mathbf{X}, \mathbf{W}) P(\mathbf{W})}{P(\mathbf{Y}|\mathbf{X})}$$

Intractable!



Approximations through sampling

Evaluate T stochastic forward passes using different samples of weights $\{\mathbf{W}_t\}_{t=1}^T$

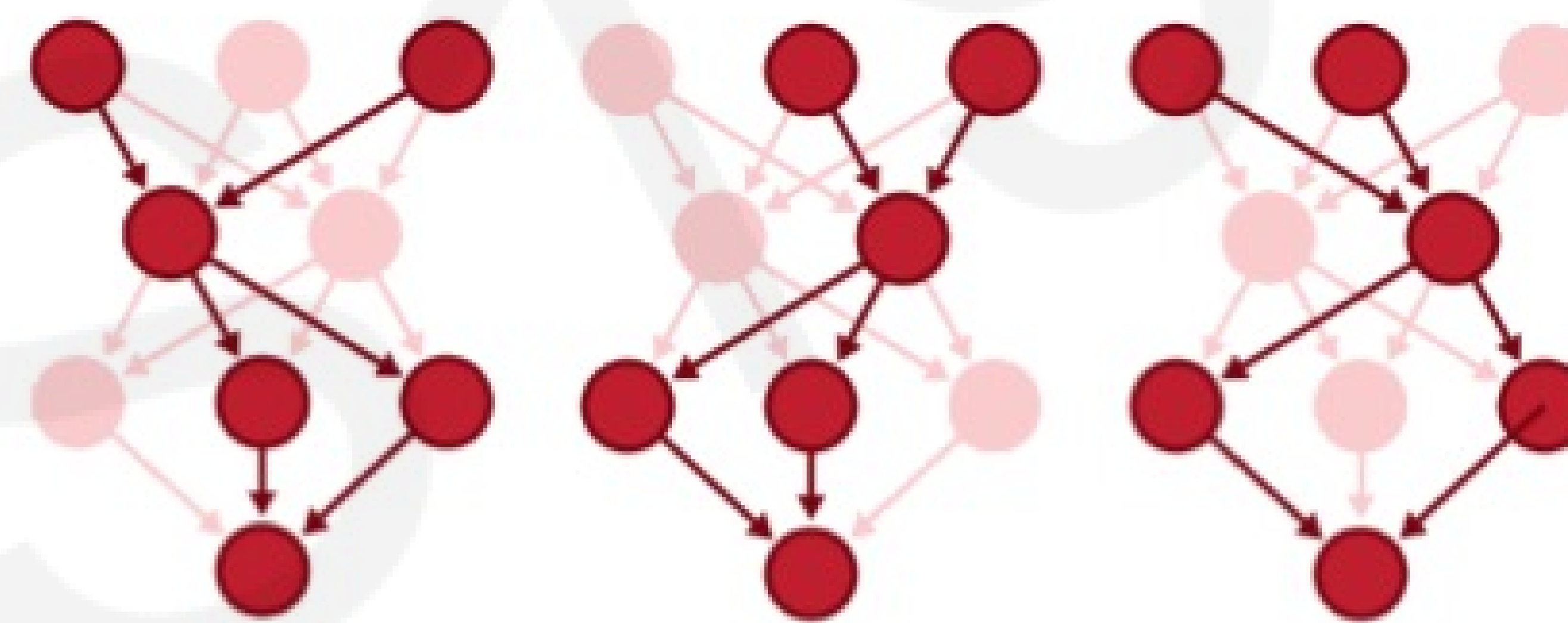
Dropout as a form of stochastic sampling

$$z_{w,t} \sim \text{Bernoulli}(p) \quad \forall w \in \mathbf{W}$$

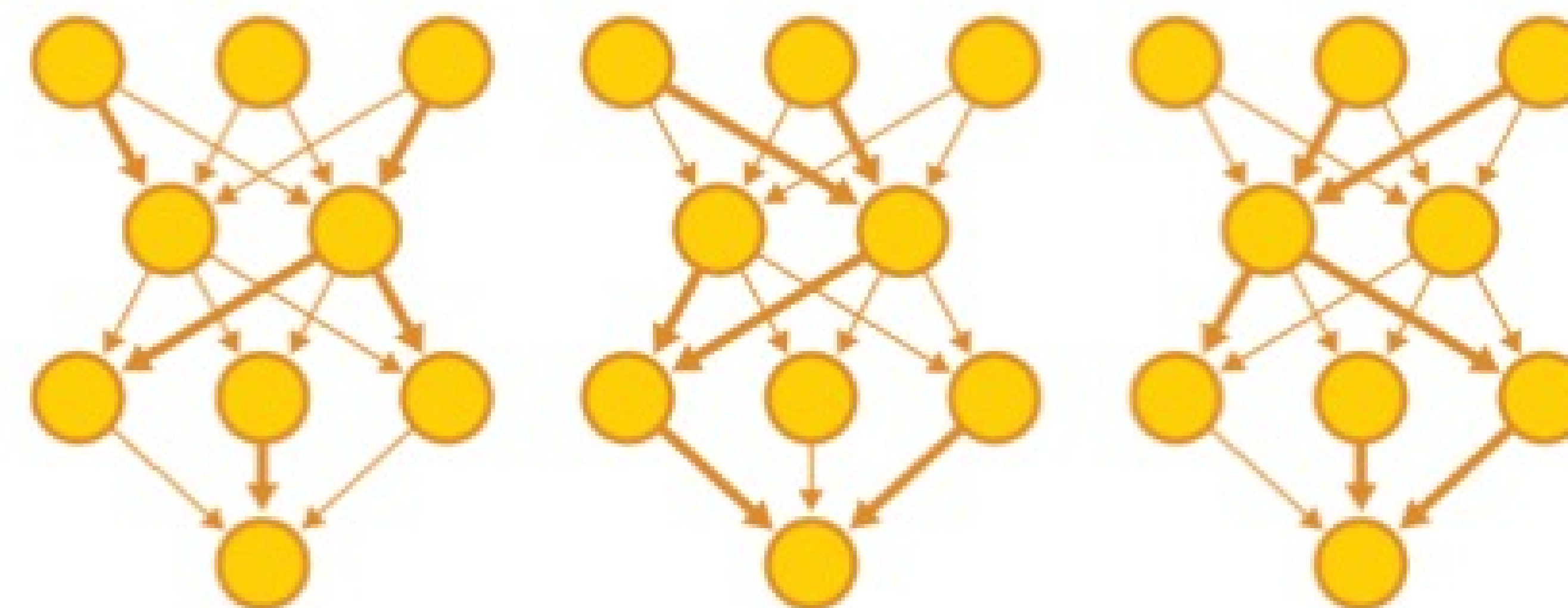
Ensemble of T independently trained models, each learning a unique \mathbf{W}_t

$$\mathbf{W}_t = \text{train}(f; \mathbf{X}, \mathbf{Y})$$

Monte Carlo Dropout



Model Ensembles



$$\mathbb{E}(\hat{\mathbf{Y}}|\mathbf{X}) = \frac{1}{T} \sum_{t=1}^T f(\mathbf{X}|\mathbf{W}_t)$$

$$\text{Var}(\hat{\mathbf{Y}}|\mathbf{X}) = \frac{1}{T} \sum_{t=1}^T f(\mathbf{X})^2 - \mathbb{E}(\hat{\mathbf{Y}}|\mathbf{X})^2$$

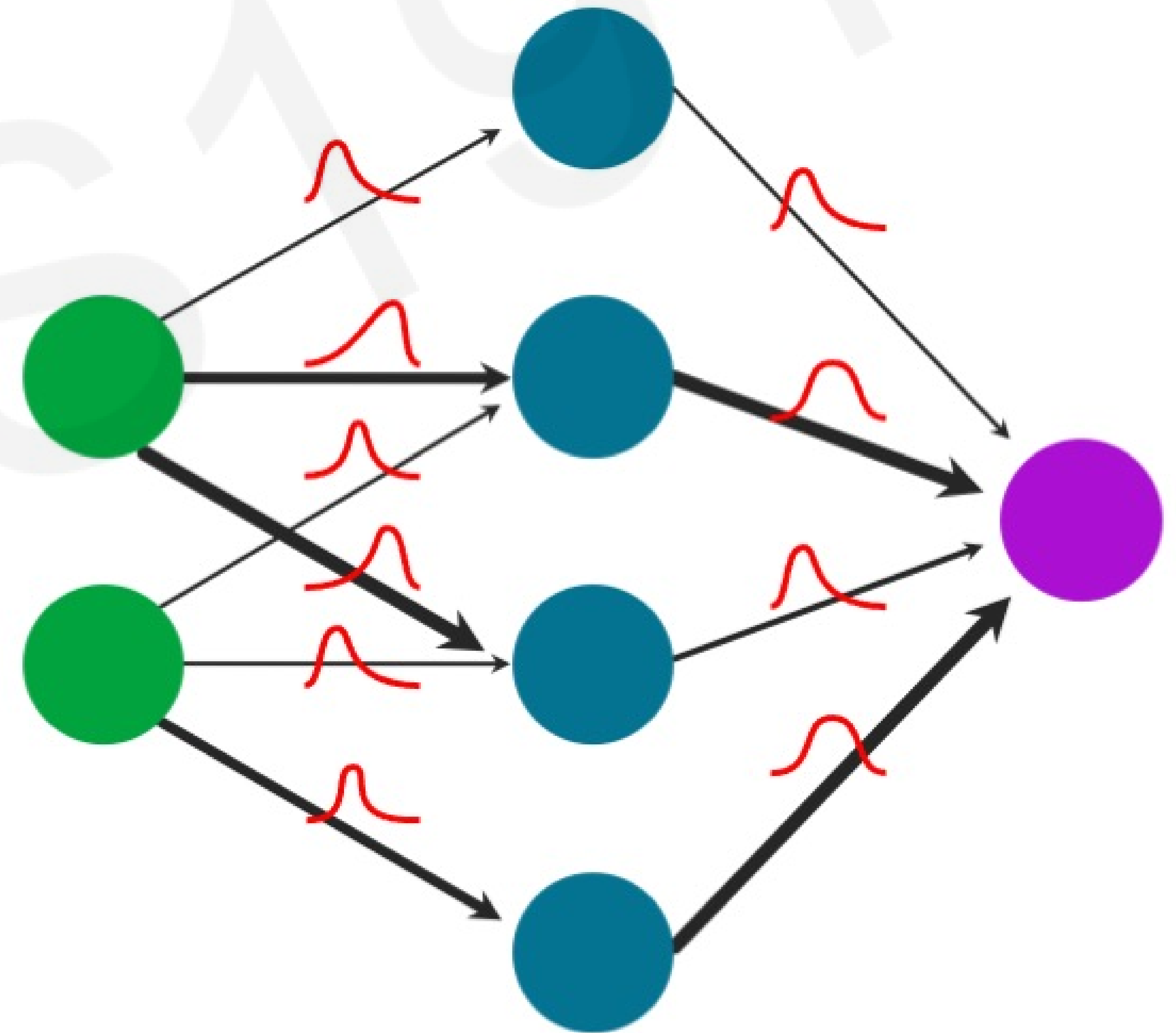
Downsides of Bayesian deep learning

Slow: Requires running the network T times for every input

Memory: Store T copies of the network in parallel

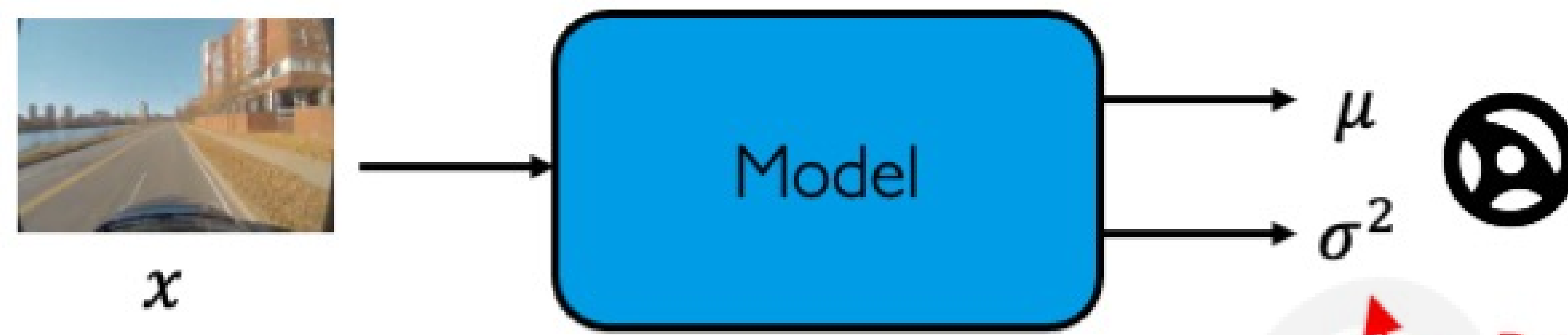
Efficiency: Sampling hinders real-time ability on edge devices (robotics)

Calibration: Sensitive to choice of prior and is often over-confident



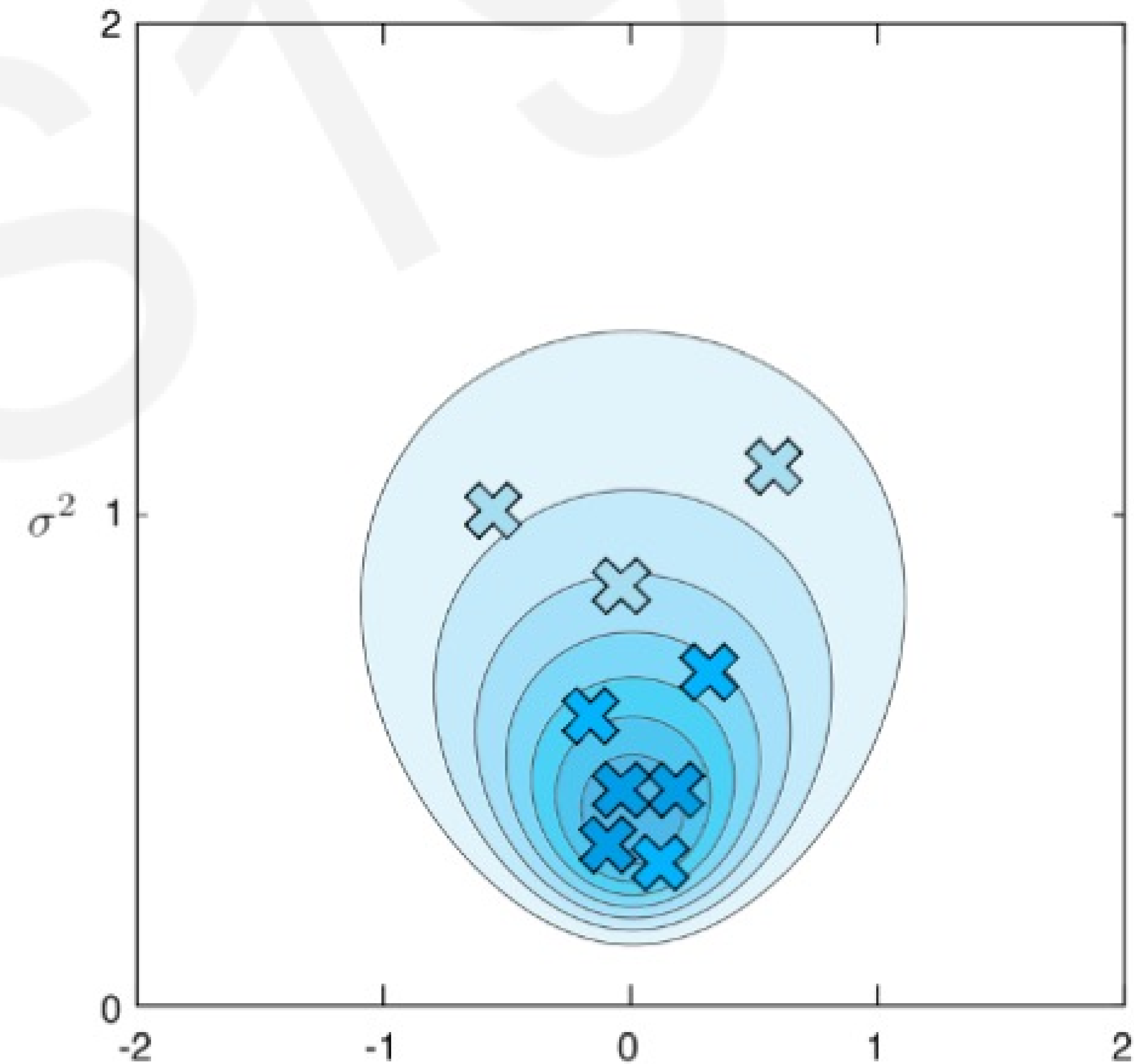
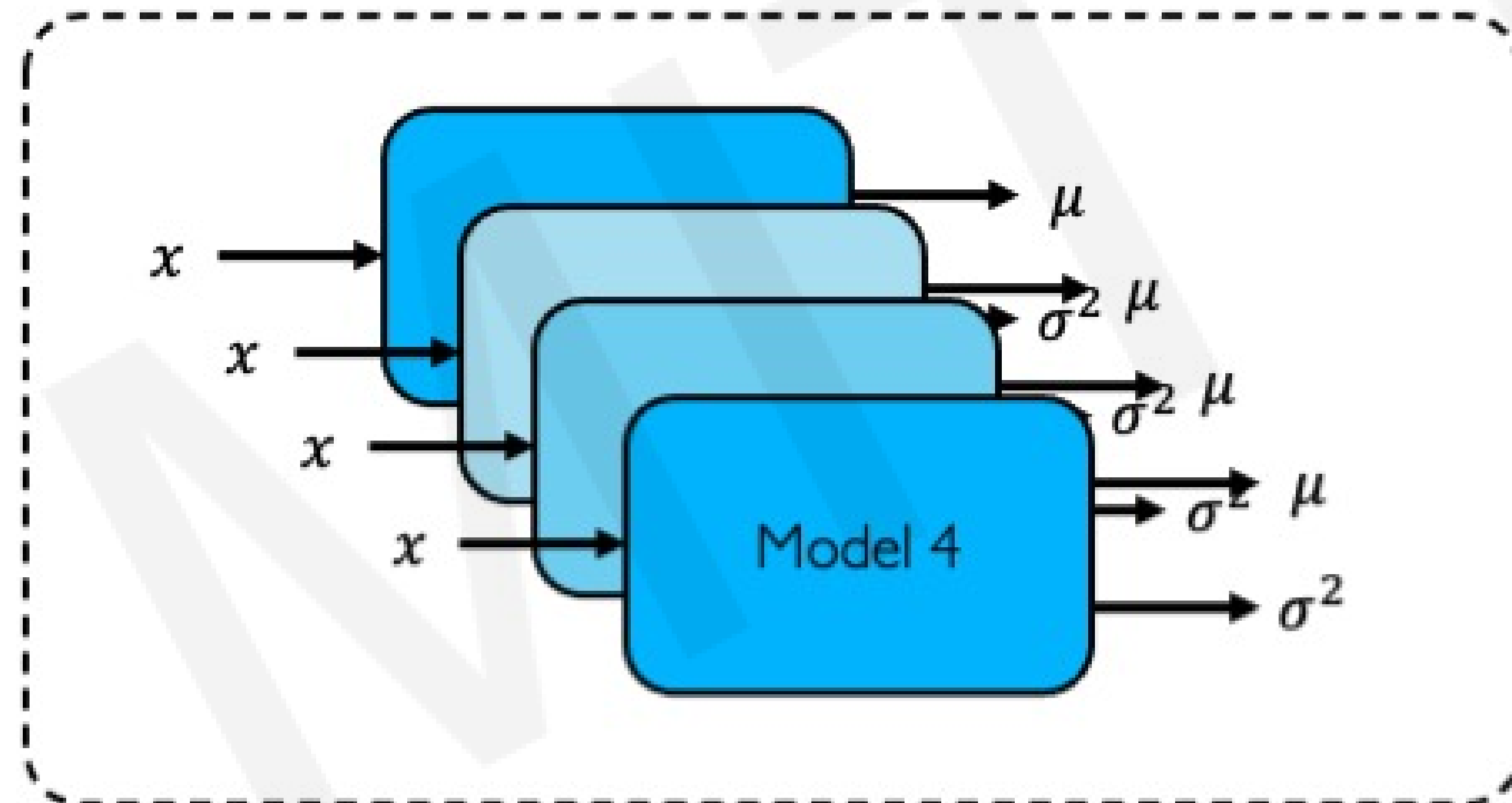
Beyond sampling for approximating uncertainty

Sampling an ensemble of models to approximate the uncertainty



Data Uncertainty

Ensemble

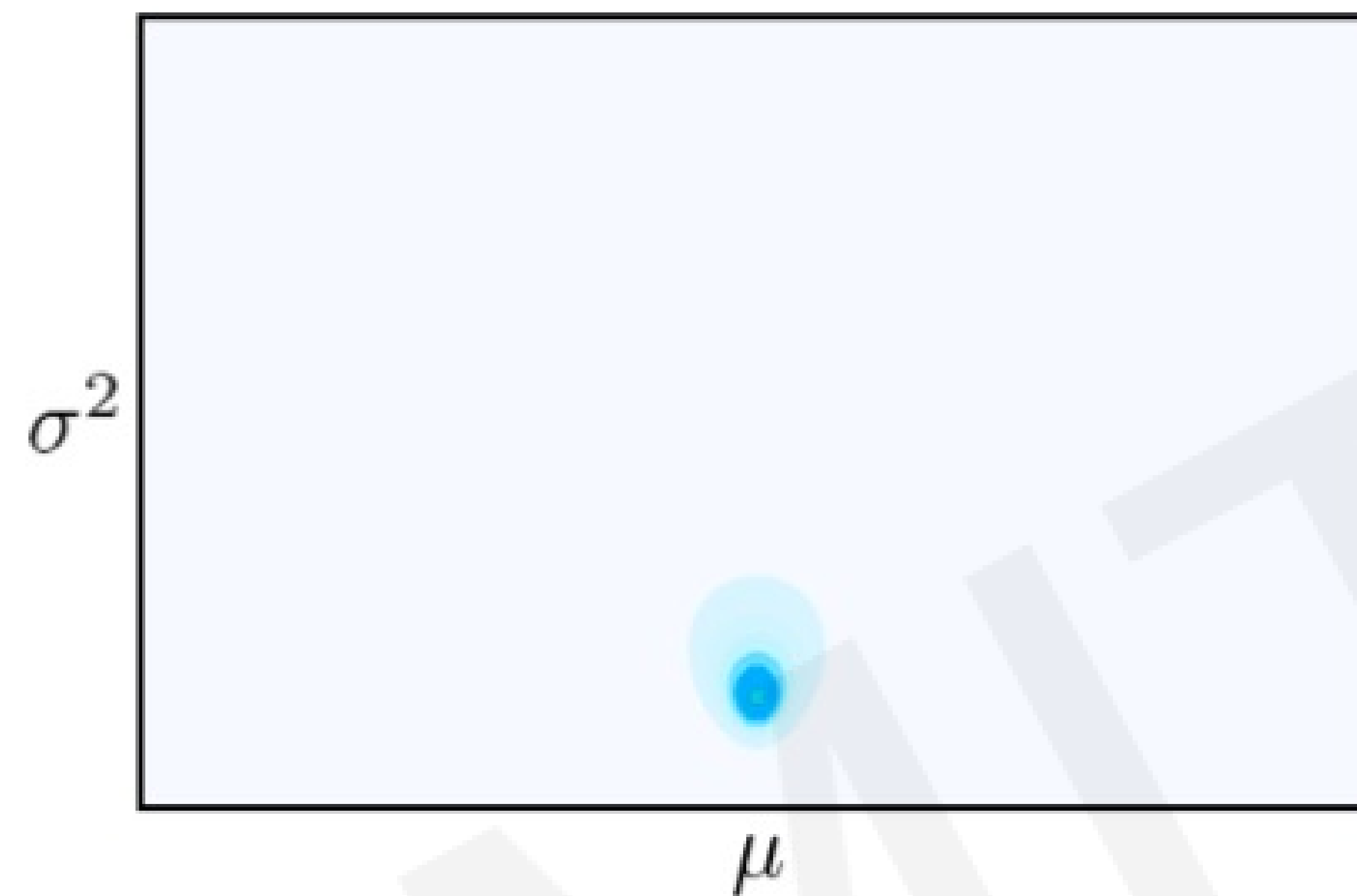


$\text{Var}[\mu]$ Model Uncertainty

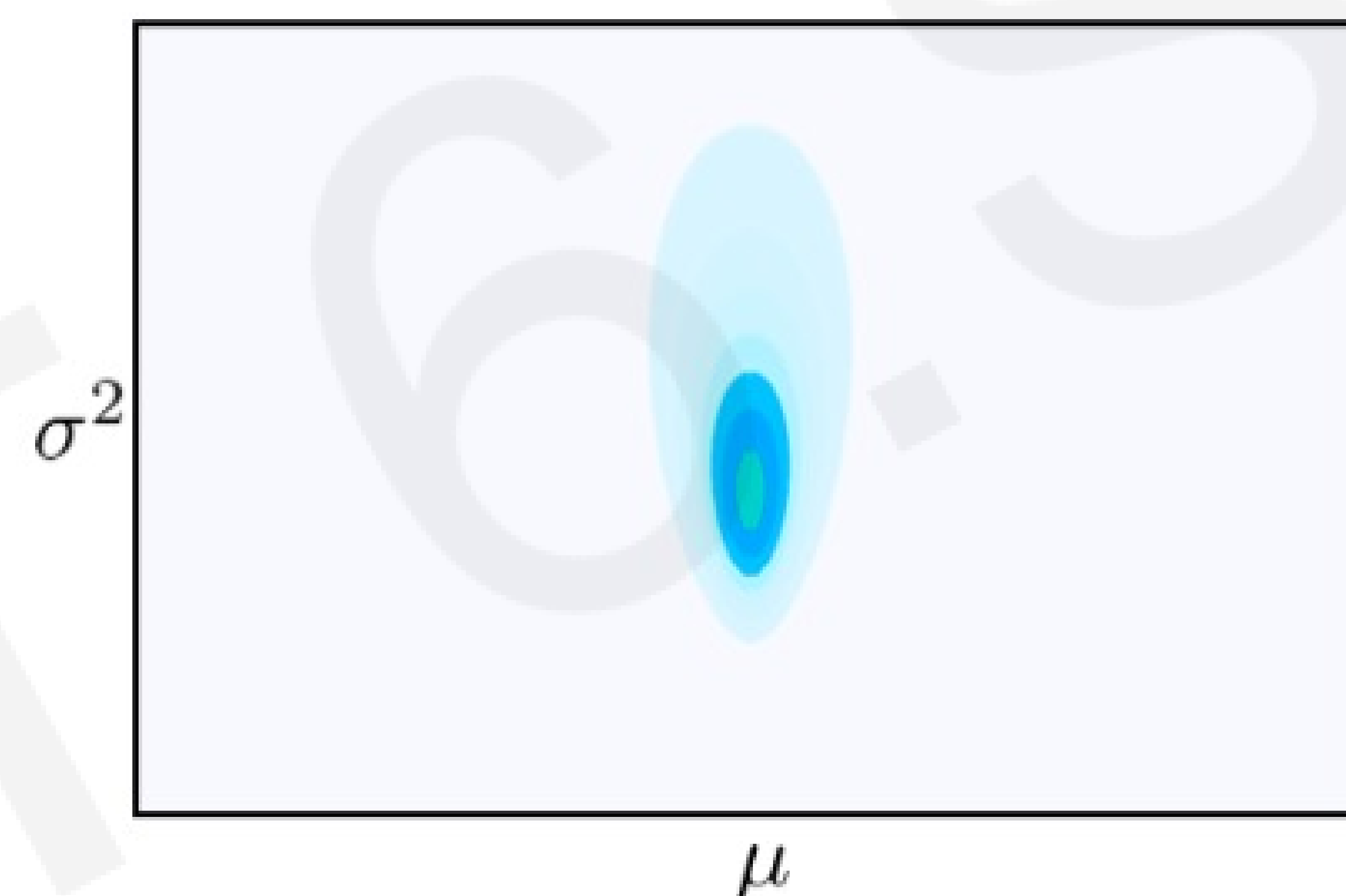
Evidential deep learning

Treat learning as an **evidence acquisition** process

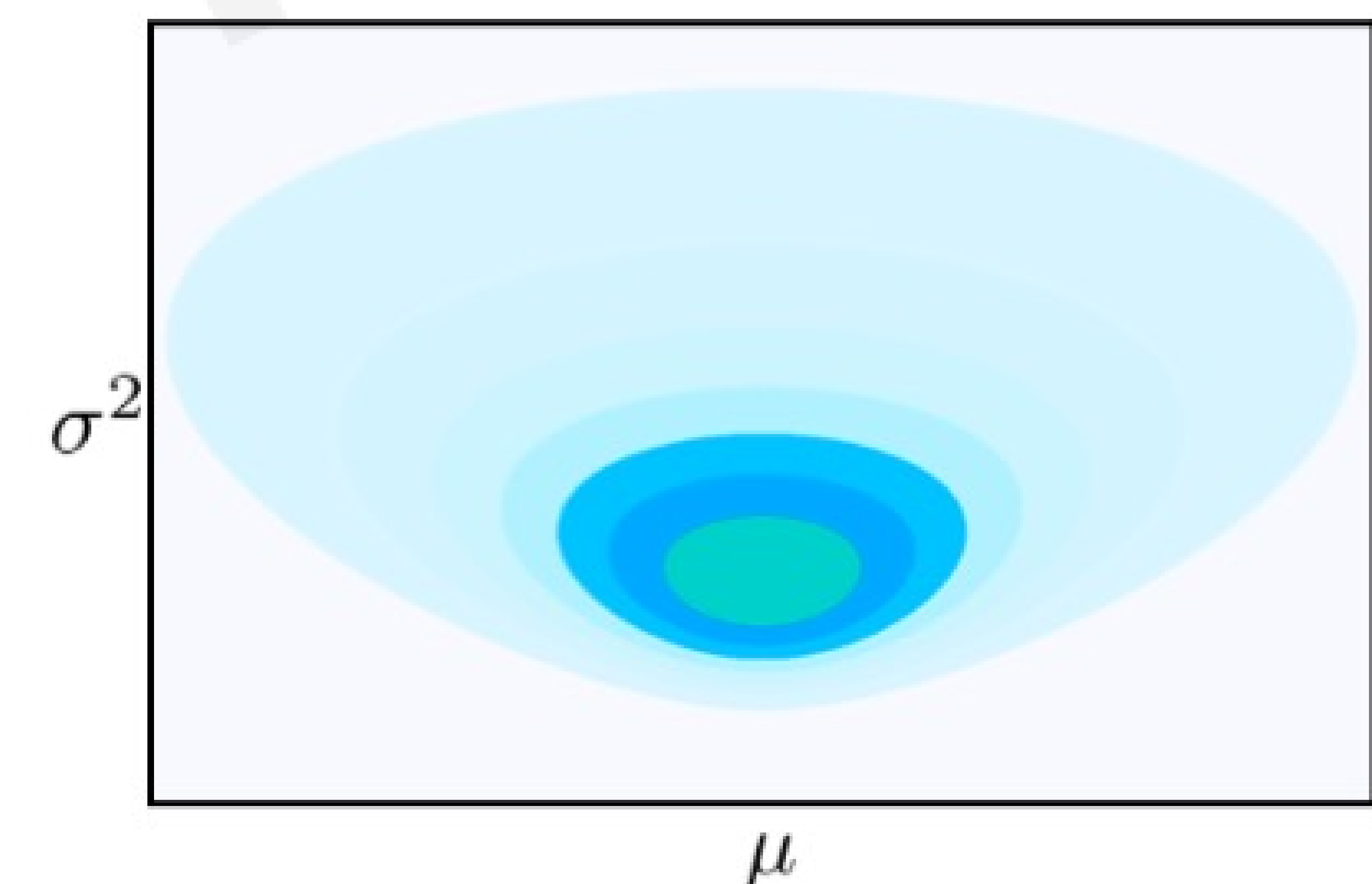
More evidence \rightarrow increased predictive confidence



Low uncertainty
High confidence



High aleatoric (data)
uncertainty



High epistemic (model)
uncertainty

Evidential learning for regression

Sampling from an evidential distribution yields individual new distributions over the data

$$y \sim \text{Normal}(\mu, \sigma^2)$$

Target
Labels

Likelihood
function

Distribution
parameters

Assume the **distribution parameters** are not known, place priors over each and probabilistically estimate!

$$\mu \sim \text{Normal}(\gamma, \sigma^2 v^{-1})$$

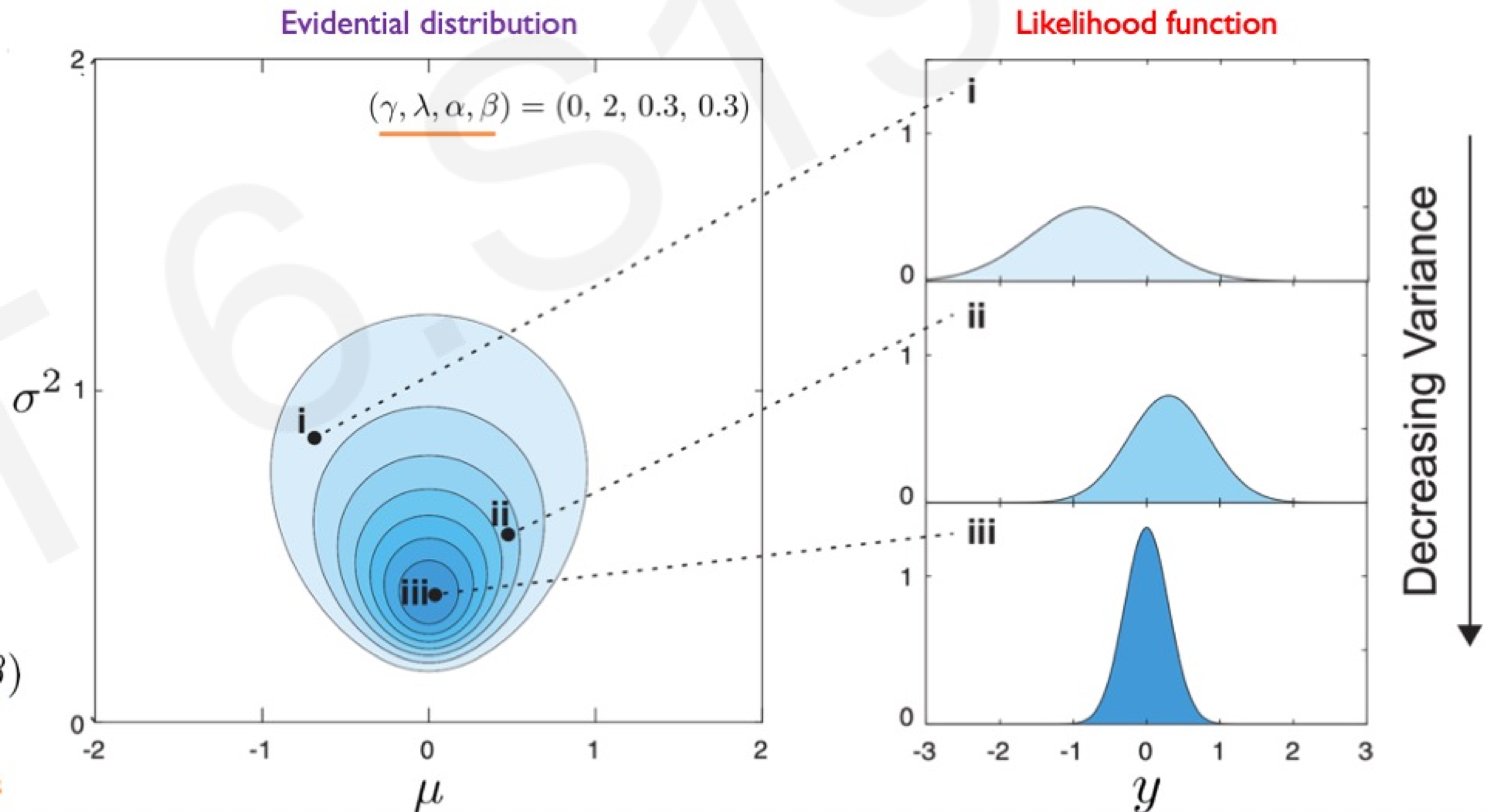
$$\sigma^2 \sim \Gamma^{-1}(\alpha, \beta)$$

$$\mu, \sigma^2 \sim \text{NormalInvGamma}(\gamma, v, \alpha, \beta)$$

Distribution
parameters

Evidential Prior

Model
parameters



Evidential learning for classification

Sampling from an evidential distribution yields individual new distributions over the data

$$y \in \{1, \dots, K\}$$

$$y \sim \text{Categorical}(p)$$

Class Labels

Likelihood function

Distribution parameters (probabilities)

$$p \sim \text{Dirichlet}(\alpha)$$

Distribution parameters

Evidential Prior

Model parameters

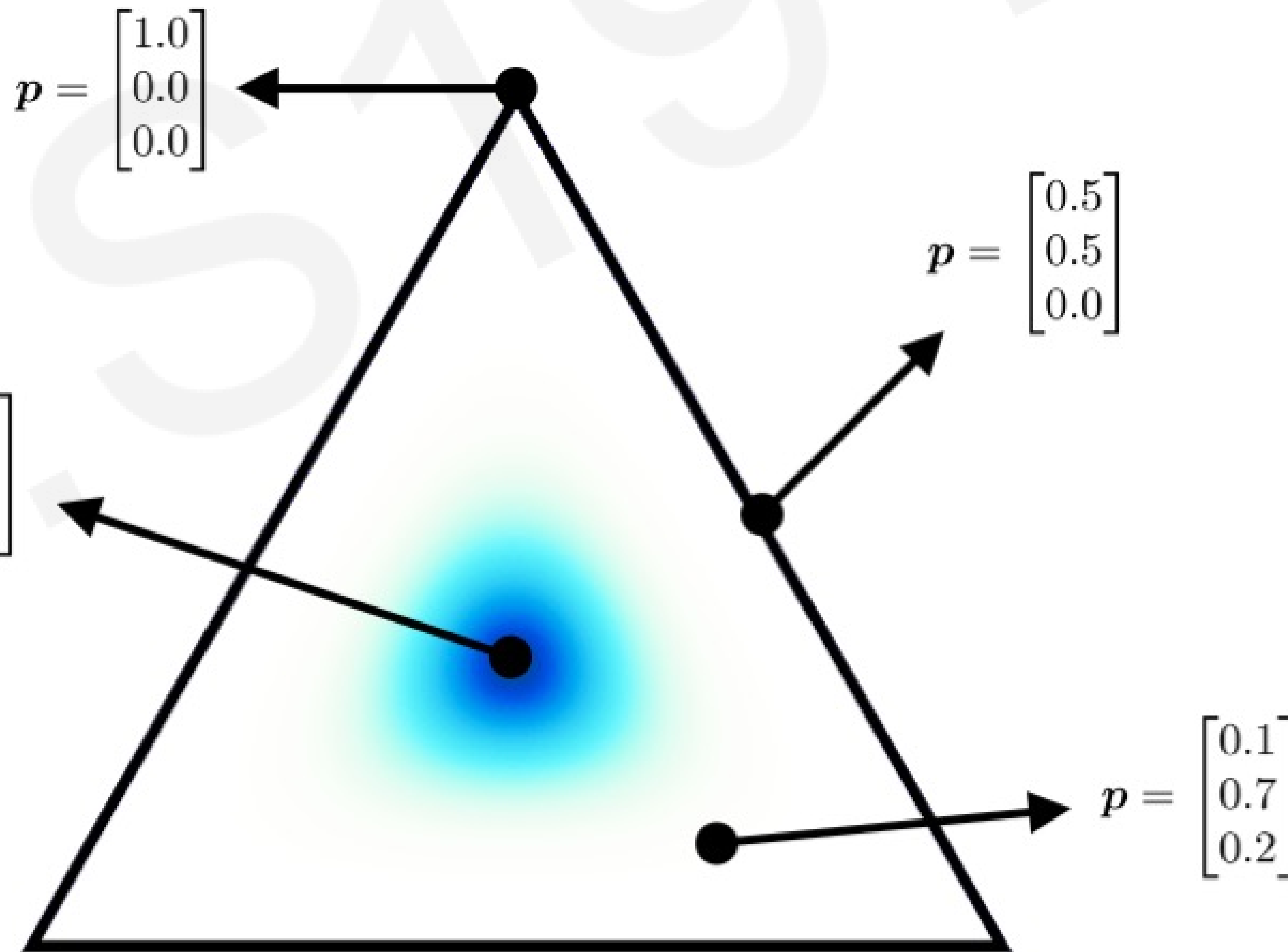
$$p = \begin{bmatrix} 0.33 \\ 0.33 \\ 0.33 \end{bmatrix}$$

$$p = \begin{bmatrix} 1.0 \\ 0.0 \\ 0.0 \end{bmatrix}$$

$$p = \begin{bmatrix} 0.5 \\ 0.5 \\ 0.0 \end{bmatrix}$$

$$p = \begin{bmatrix} 0.1 \\ 0.7 \\ 0.2 \end{bmatrix}$$

$$K = 3; \quad \alpha = (5, 5, 5)$$



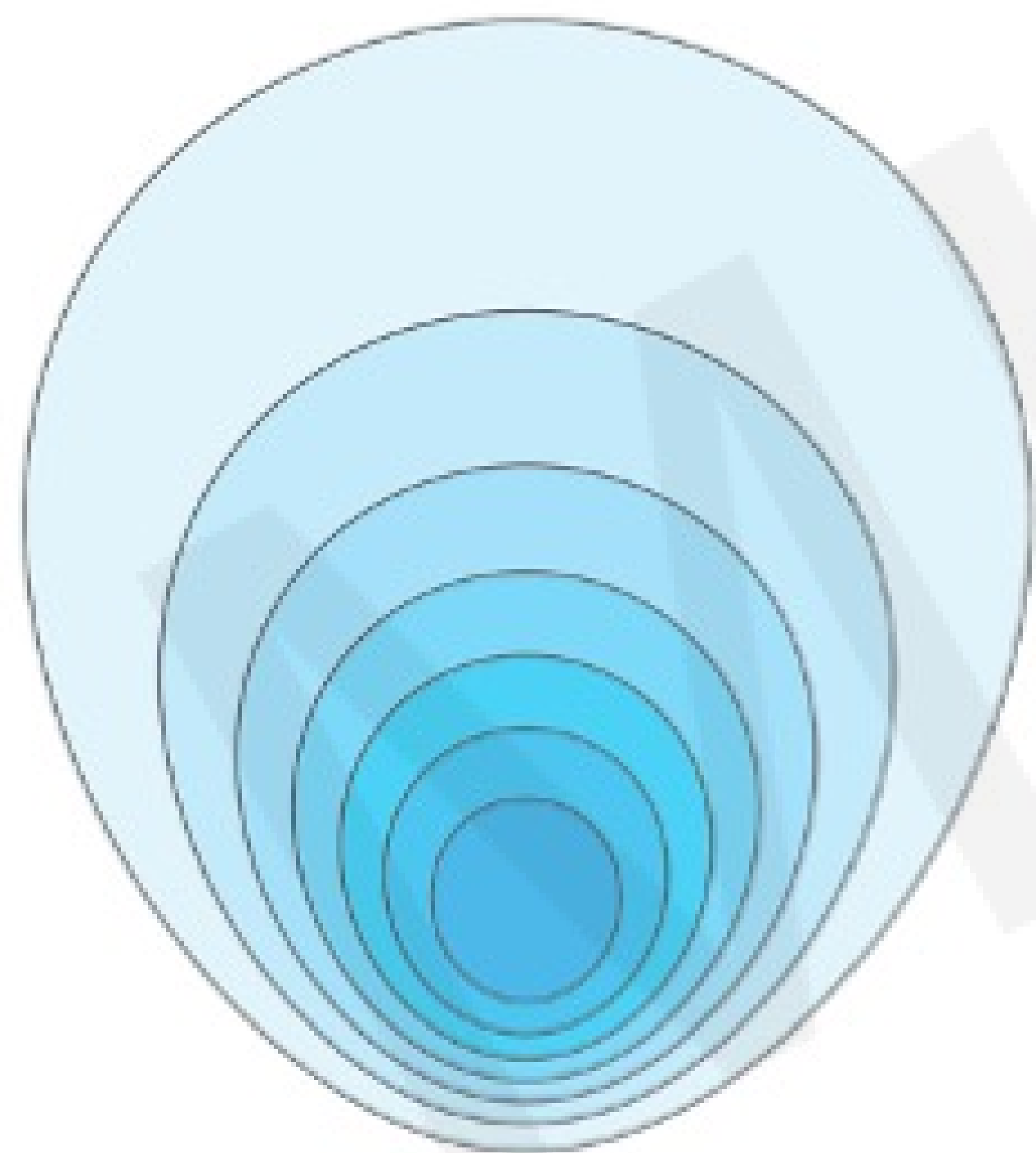
Evidential distributions for regression and classification

Regression (continuous)

$$y \in \mathbb{R}$$

$$y \sim \text{Normal}(\mu, \sigma^2)$$

$$\mu, \sigma^2 \sim \text{NormalInvGamma}(\gamma, \nu, \alpha, \beta)$$

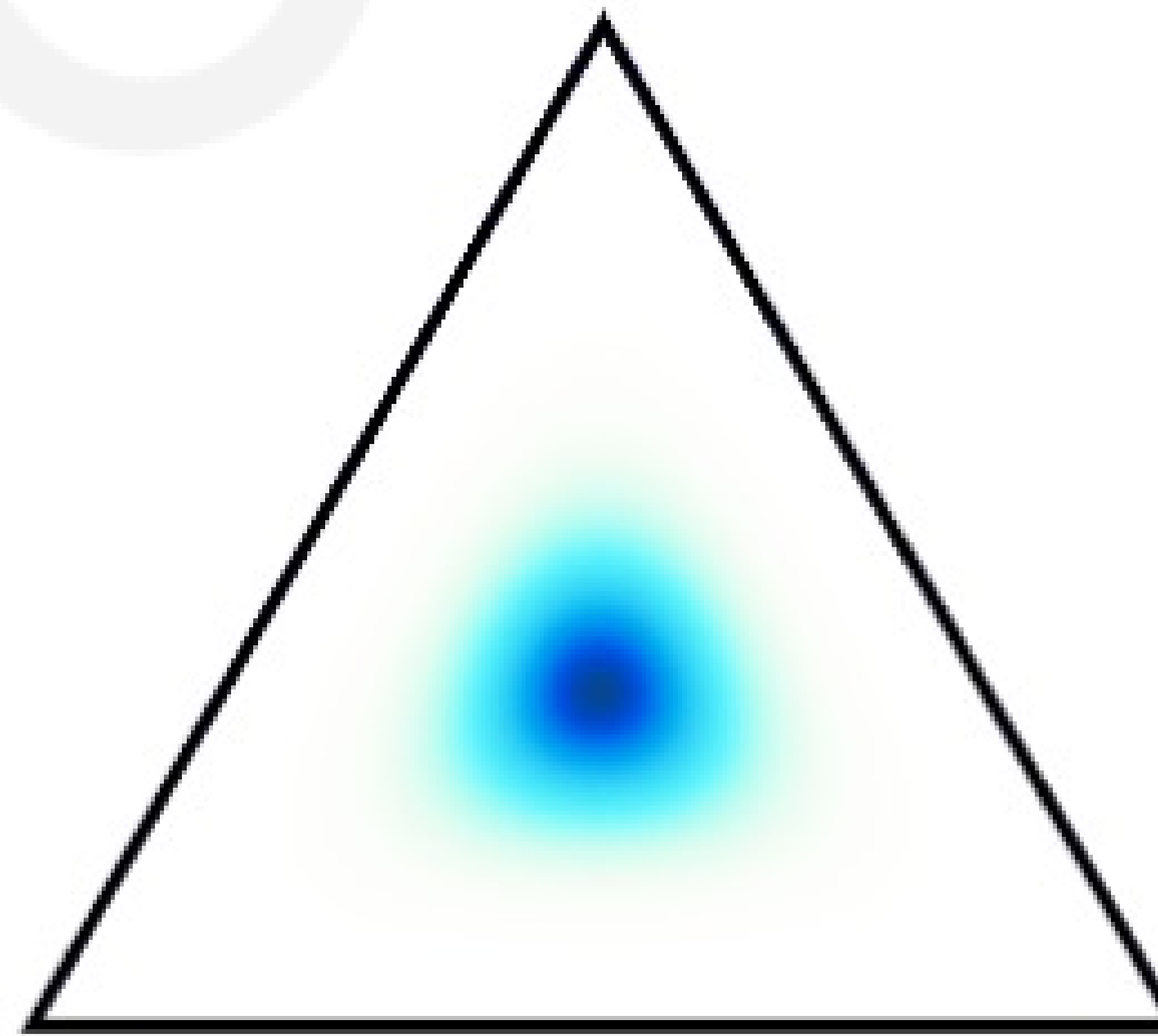


Classification (discrete)

$$y \in \{1, \dots, K\}$$

$$y \sim \text{Categorical}(\mathbf{p})$$

$$\mathbf{p} \sim \text{Dirichlet}(\boldsymbol{\alpha})$$



Side note:

Choice of evidential distribution is closely related to conjugate priors in Bayesian inference.

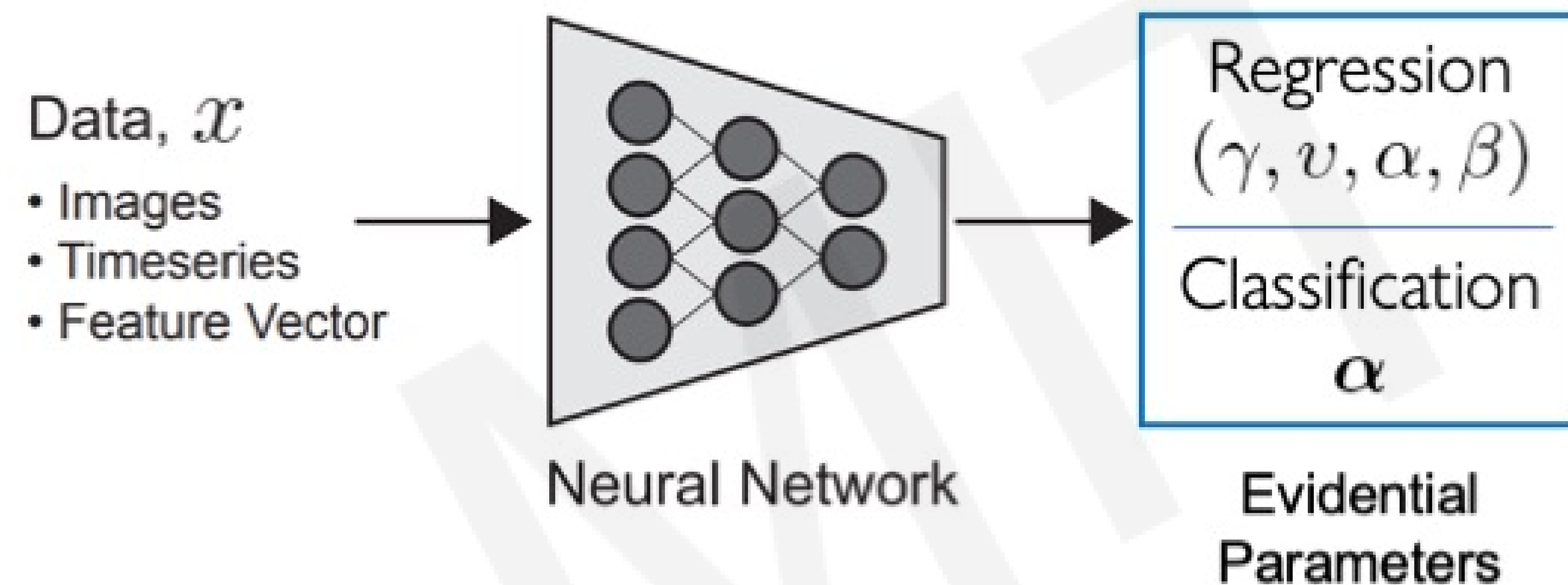
It is often easiest to pick your evidential distribution to be a conjugate prior of your likelihood

$$p(\theta|y) = \frac{p(y|\theta) p(\theta)}{\int_{\theta'} p(y|\theta') p(\theta') d\theta'}$$

Model and training

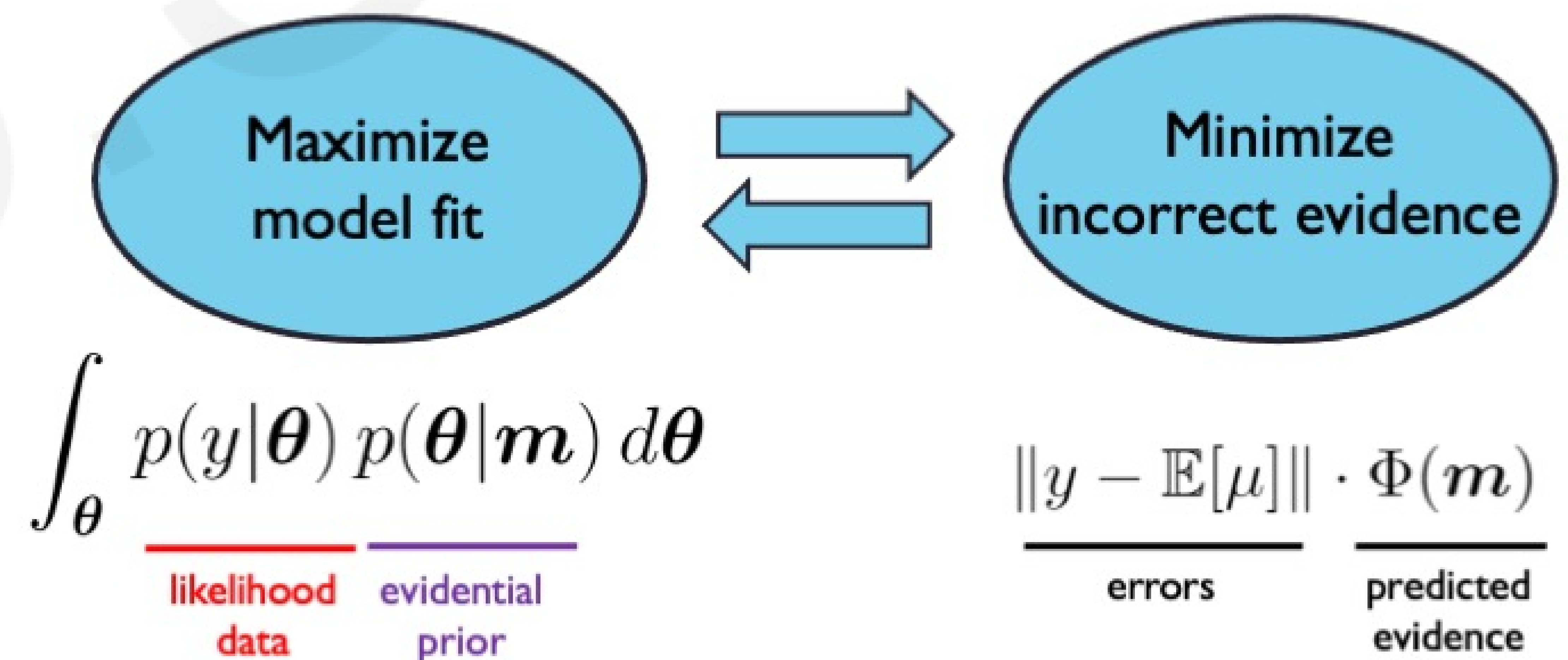
Model

Train the network to output the parameters of an evidential distribution



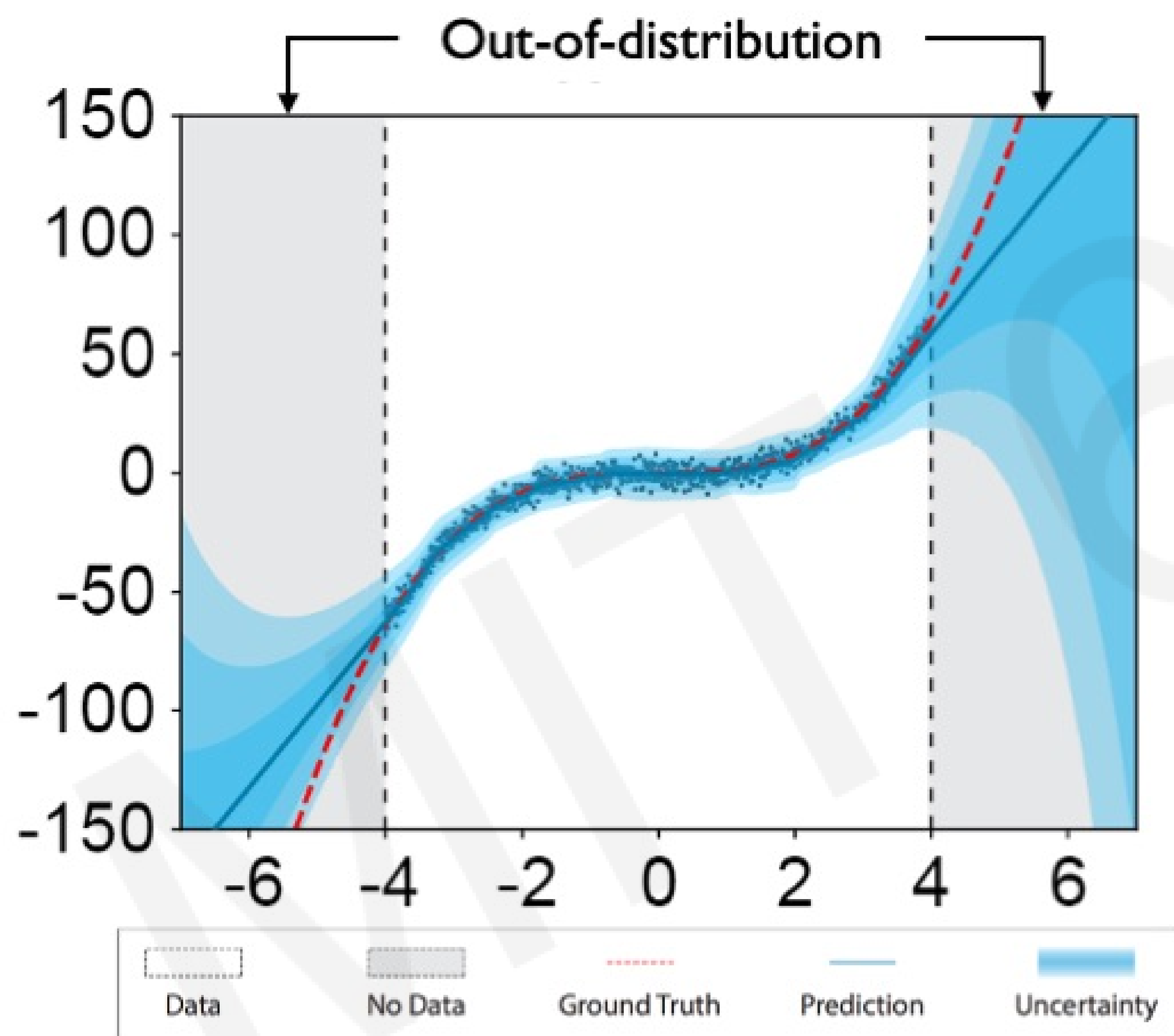
Optimization

Multi-objective training:

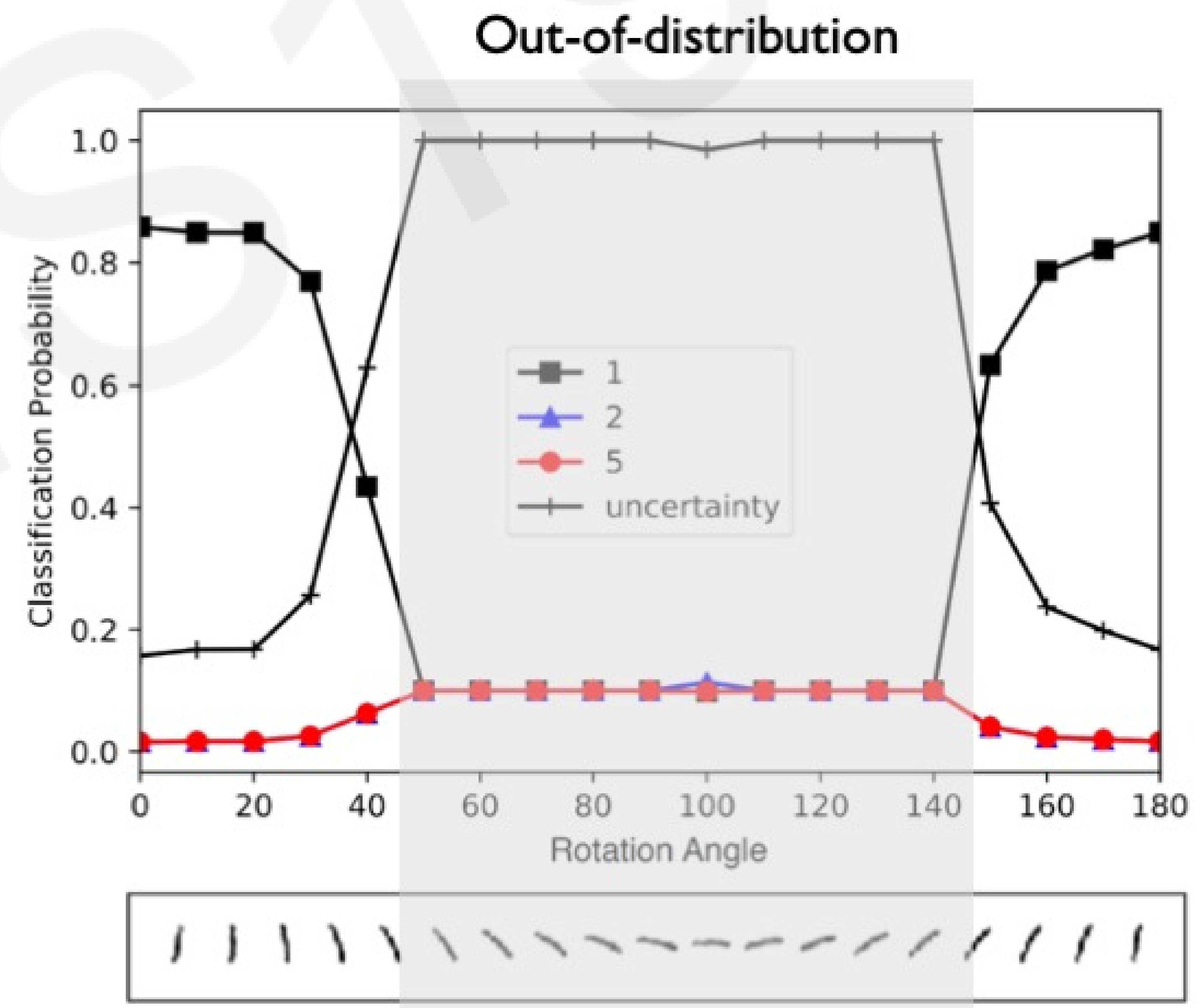


Toy learning problems

Regression (continuous)



Classification (discrete)

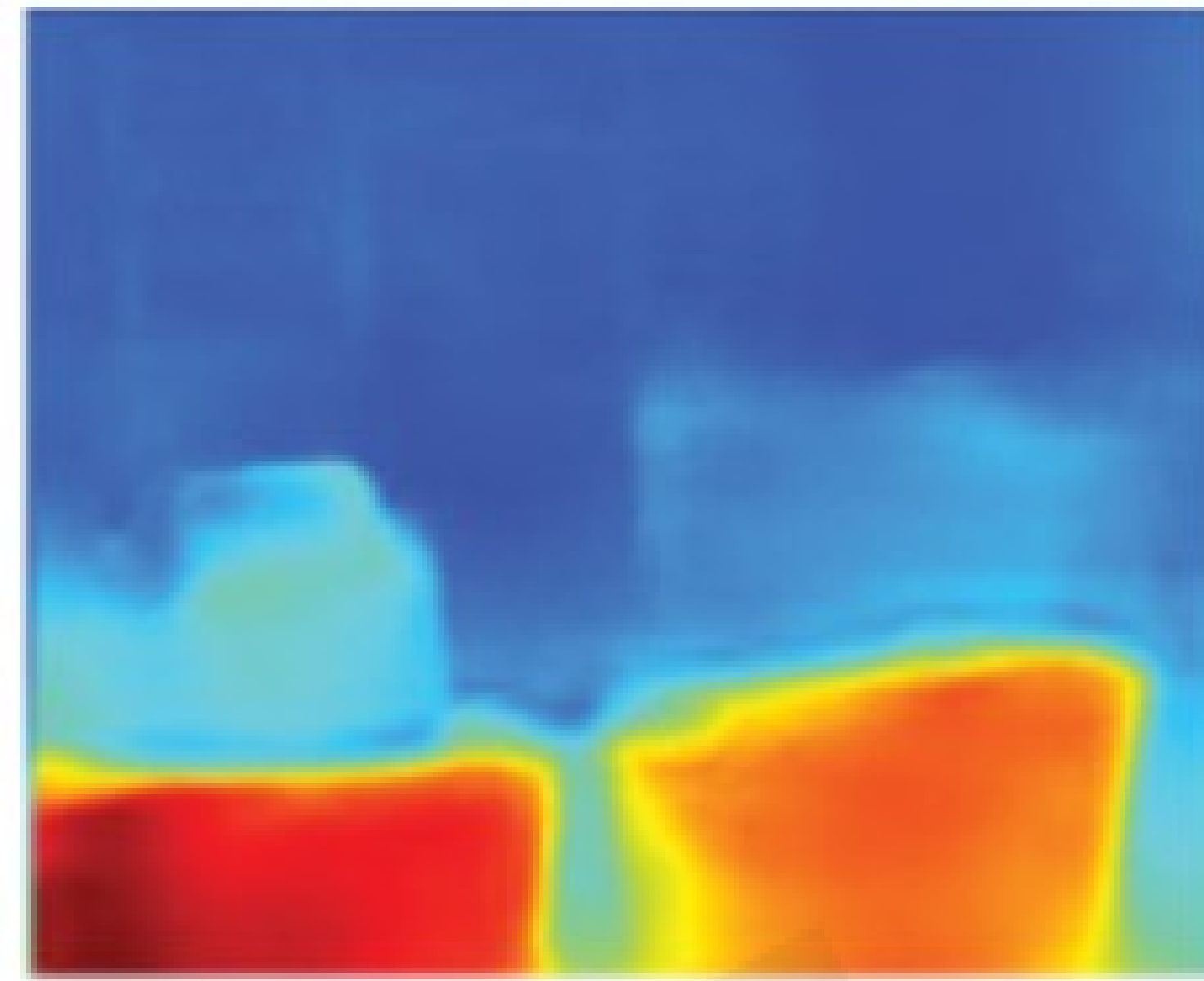


Applications of evidential learning

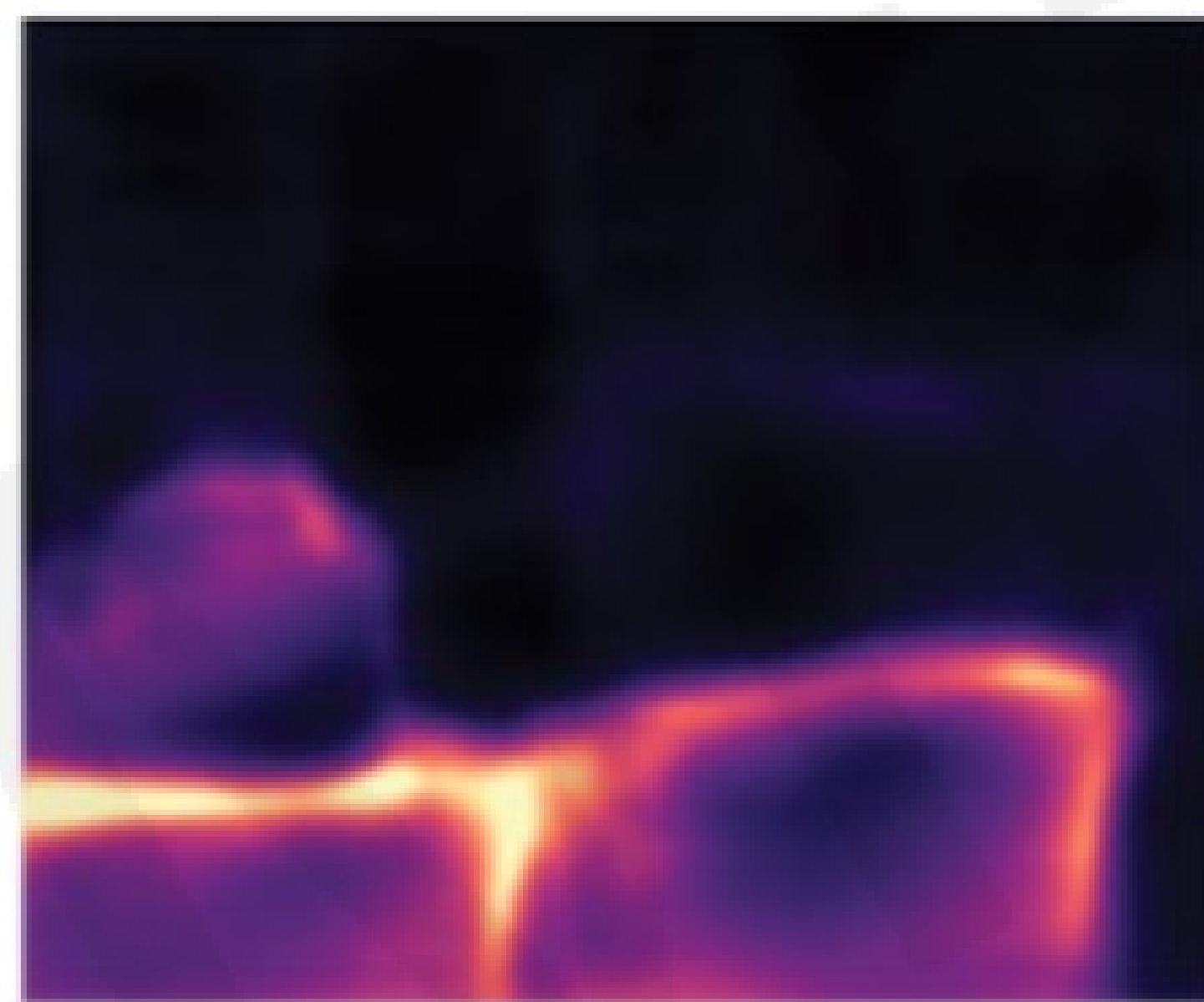
Monocular Depth Estimation



RGB input

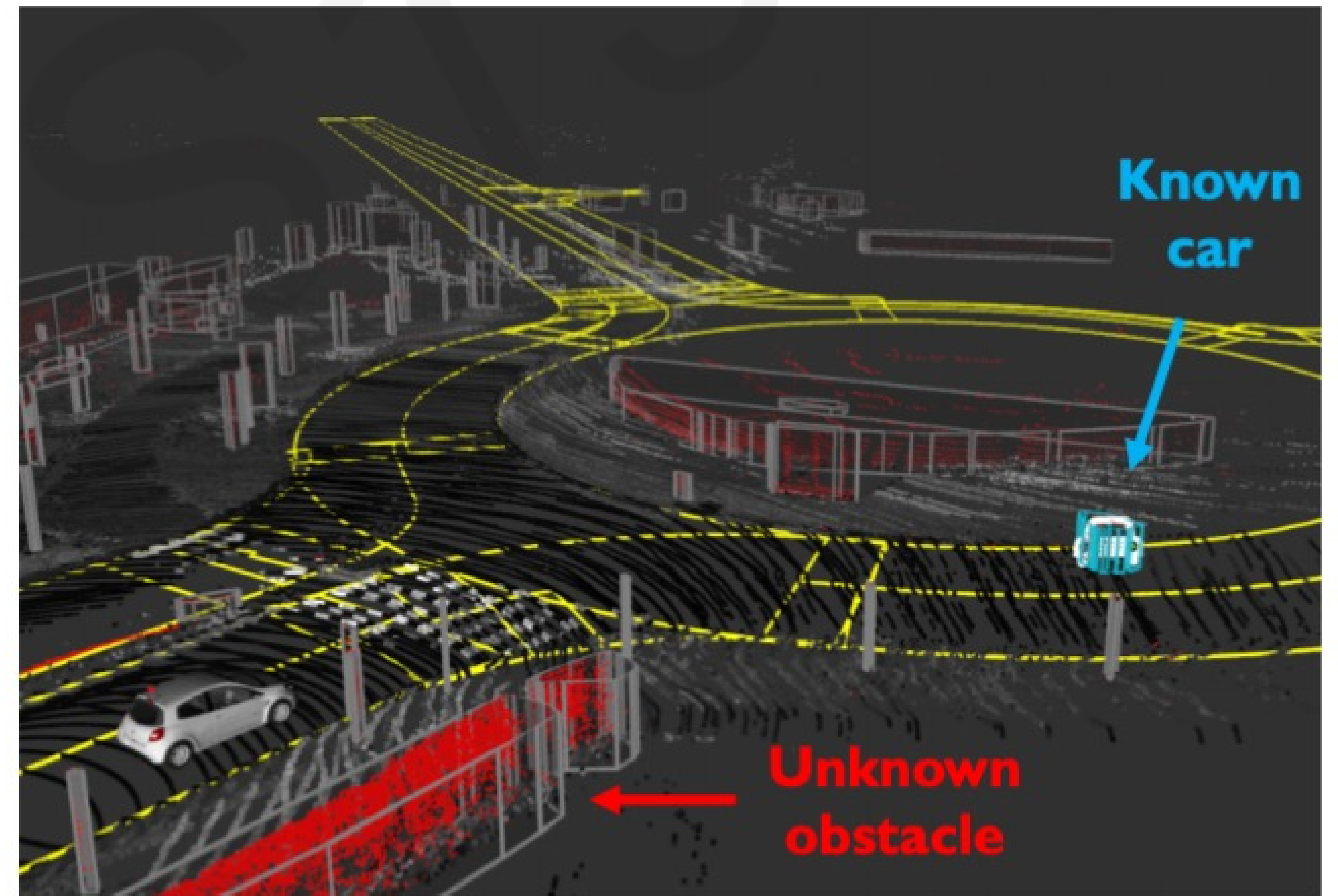


Predicted depth



Predicted uncertainty

LiDAR Object Classification



Known car

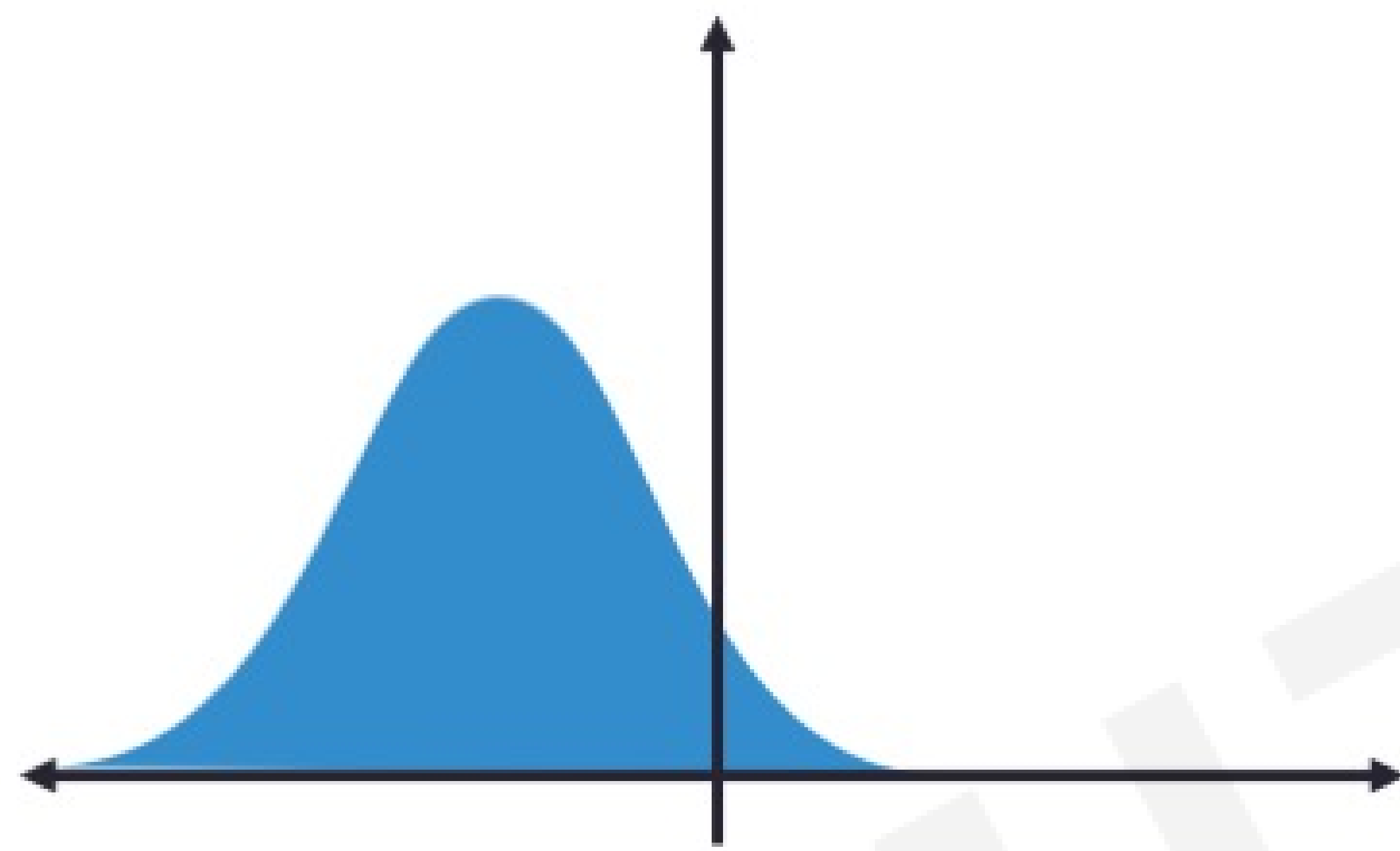
Unknown obstacle

Comparison of uncertainty estimation approaches

	Likelihood estimation	Bayesian NN	Evidential NN
Prior placed over:	Data	Weights	Likelihood
Weights are:	Deterministic	Stochastic	Deterministic
Fast (no sampling)	✓		✓
Captures epistemic uncertainty		✓	✓

Summary

Learning probability distributions



Model distributions over labels:
Softmax (discrete) & Gaussian
(continuous)

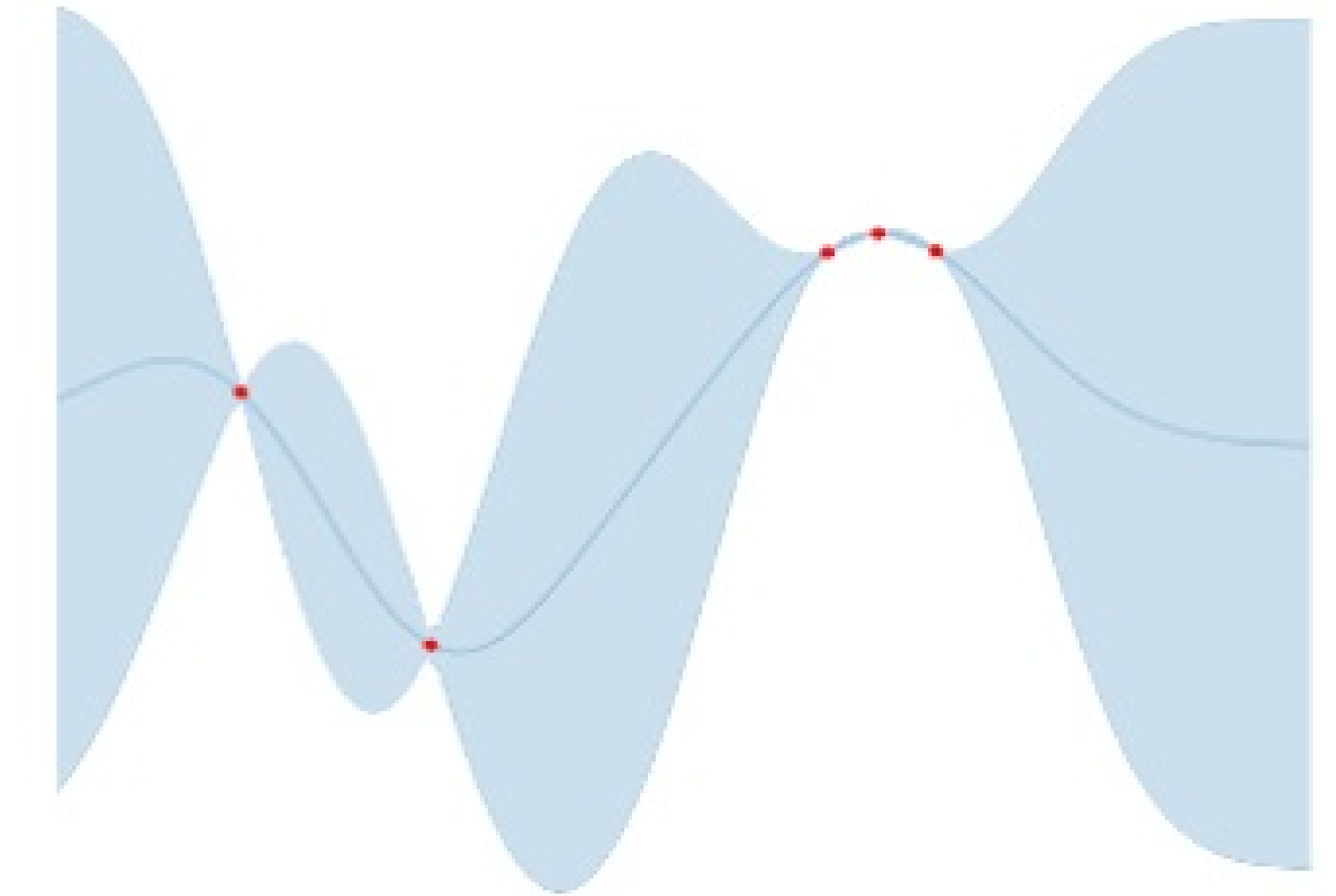
Different sources of uncertainty

$$p(y|x)$$

$\mathbb{E}[\sigma^2]$ $\text{Var}[\mu]$

Data (aleatoric) uncertainty vs.
Model (epistemic) uncertainty.

Fast and scalable uncertainty estimation



Evidential deep learning
Uncertainty modelling for quickly
estimating confidence