

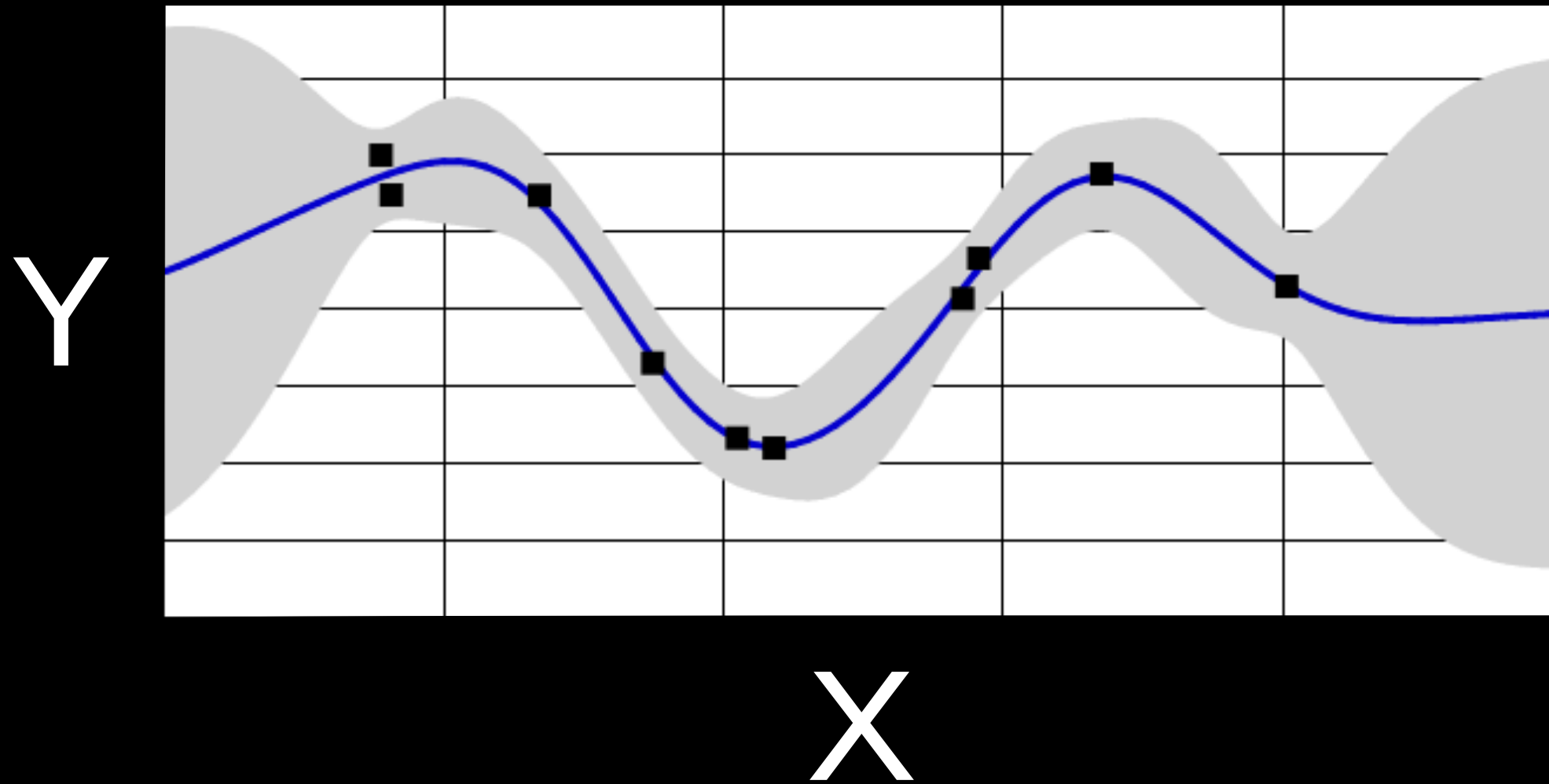
The Boons of Being Less Bayesian

a study of partially stochastic neural networks

Eric Nalisnick

Johns Hopkins University





Are stochastic parameters
always useful?

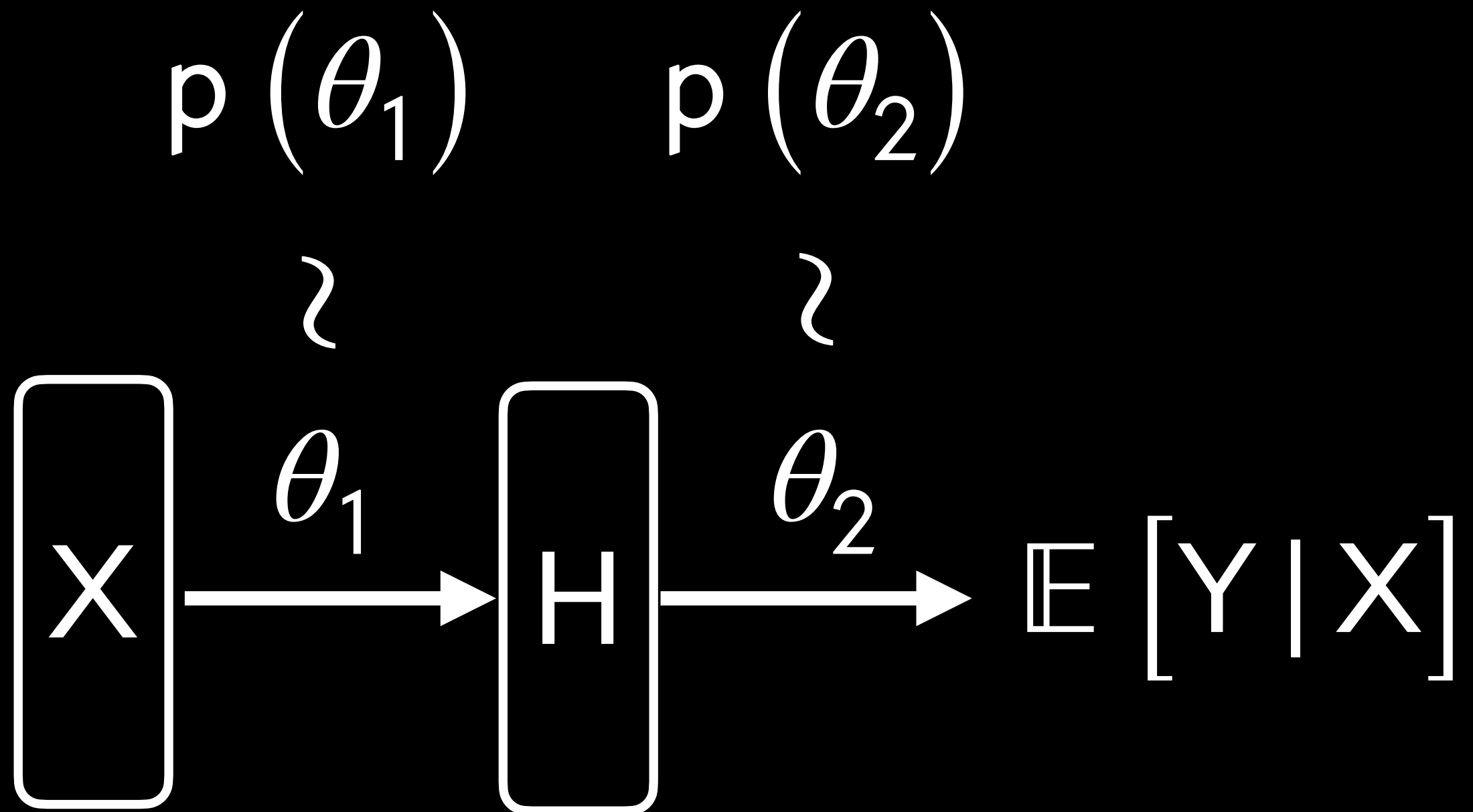
Outline

- ⊗ background
- ⊗ subnetwork inference algorithm
- ⊗ subnetworks are all you need

Outline

- ⊗ background
- ⊗ subnetwork inference algorithm
- ⊗ subnetworks are all you need

Neural Network



Bayes Rule

$$p(\theta_1, \dots, \theta_L | \mathfrak{D}) =$$

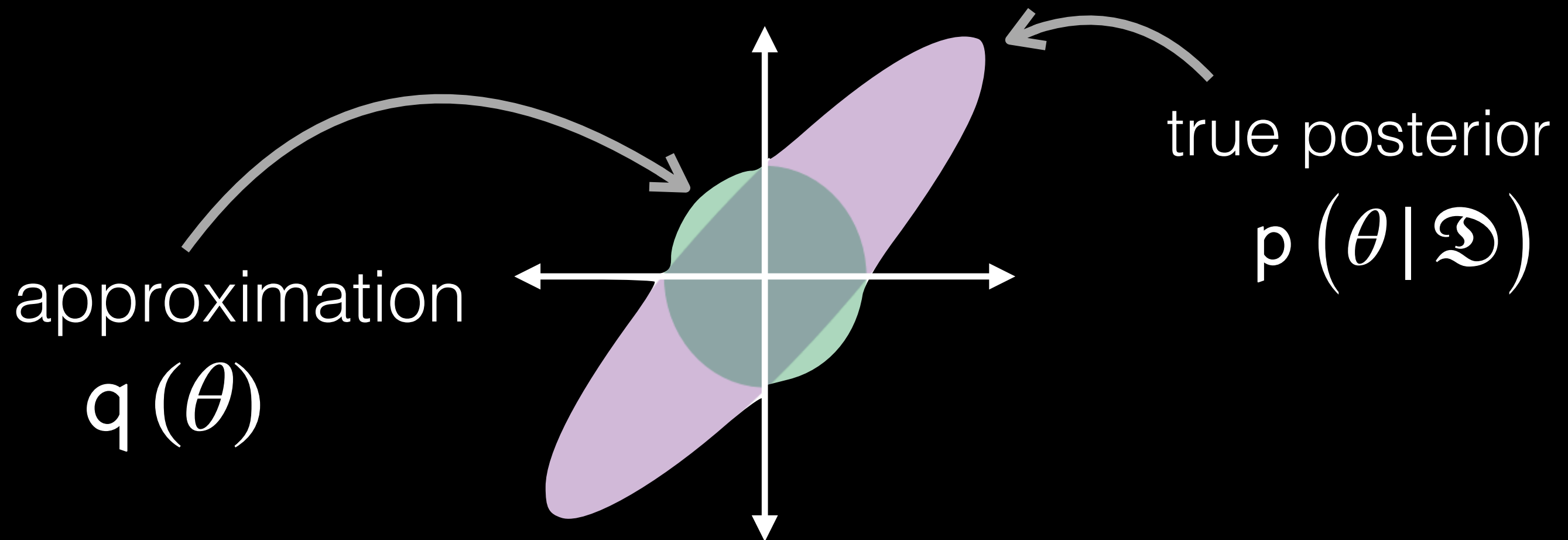
$$\frac{\prod_i p(\theta_i) \prod_{n=1}^N p(y_n | x_n, \theta_1, \dots, \theta_L)}{p(\mathfrak{D})}$$

Bayes Rule

$$\cancel{p(\theta_1, \dots, \theta_L | \mathfrak{D})} =$$

$$\frac{\prod_i p(\theta_i) \prod_{n=1}^N p(y_n | x_n, \theta_1, \dots, \theta_L)}{\cancel{p(\mathfrak{D})}}$$

Variational Inference



Mean Field Assumption

$$p(\theta_1, \dots, \theta_L | \mathfrak{D}) \approx q(\theta_1, \dots, \theta_L)$$

Mean Field Assumption

$$p(\theta_1, \dots, \theta_L | \mathfrak{D}) \approx q(\theta_1, \dots, \theta_L)$$

$$= \prod_i^L q(\theta_i)$$

factorize over layers

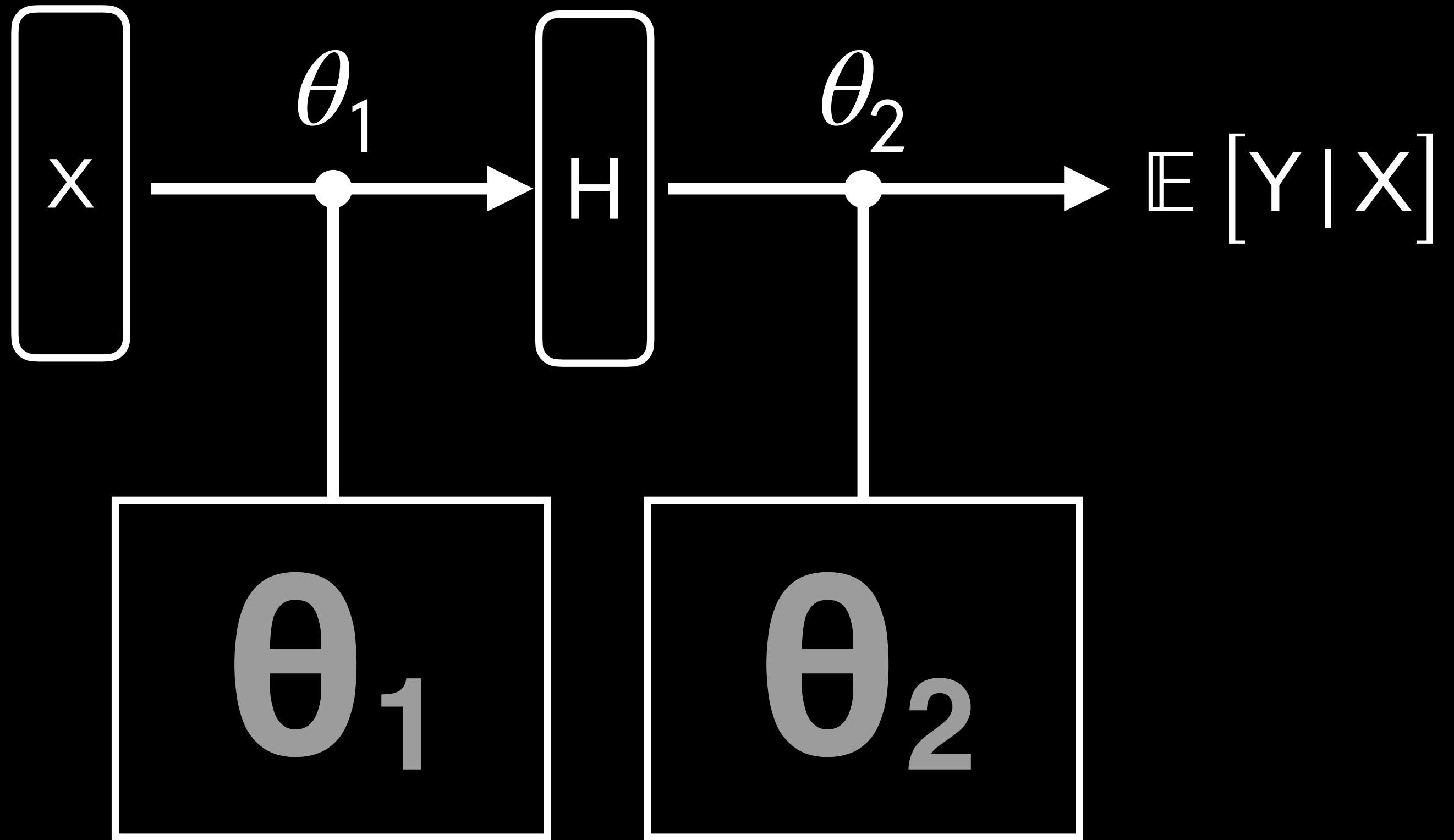
Mean Field Assumption

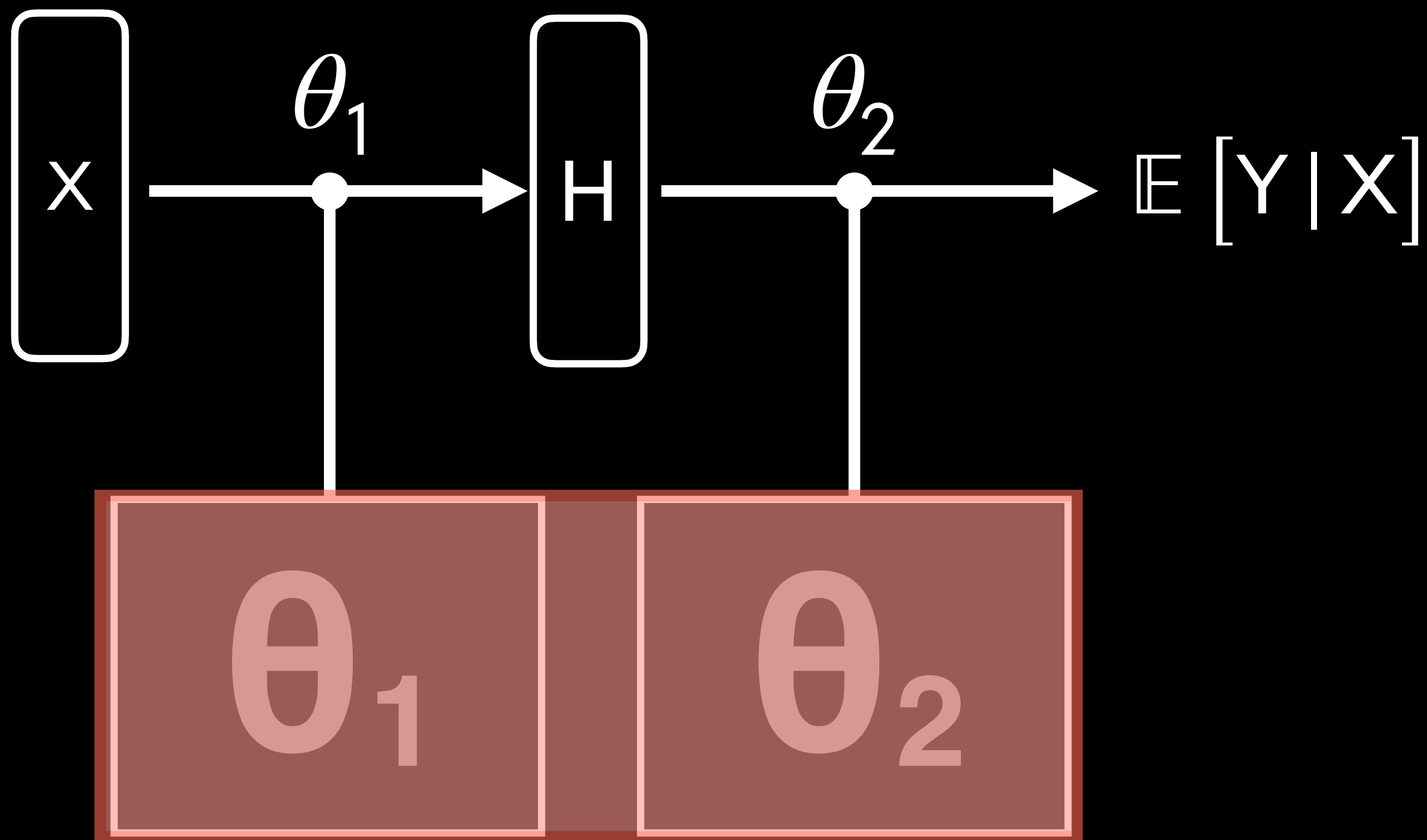
$$p(\theta_1, \dots, \theta_L | \mathfrak{D}) \approx q(\theta_1, \dots, \theta_L)$$

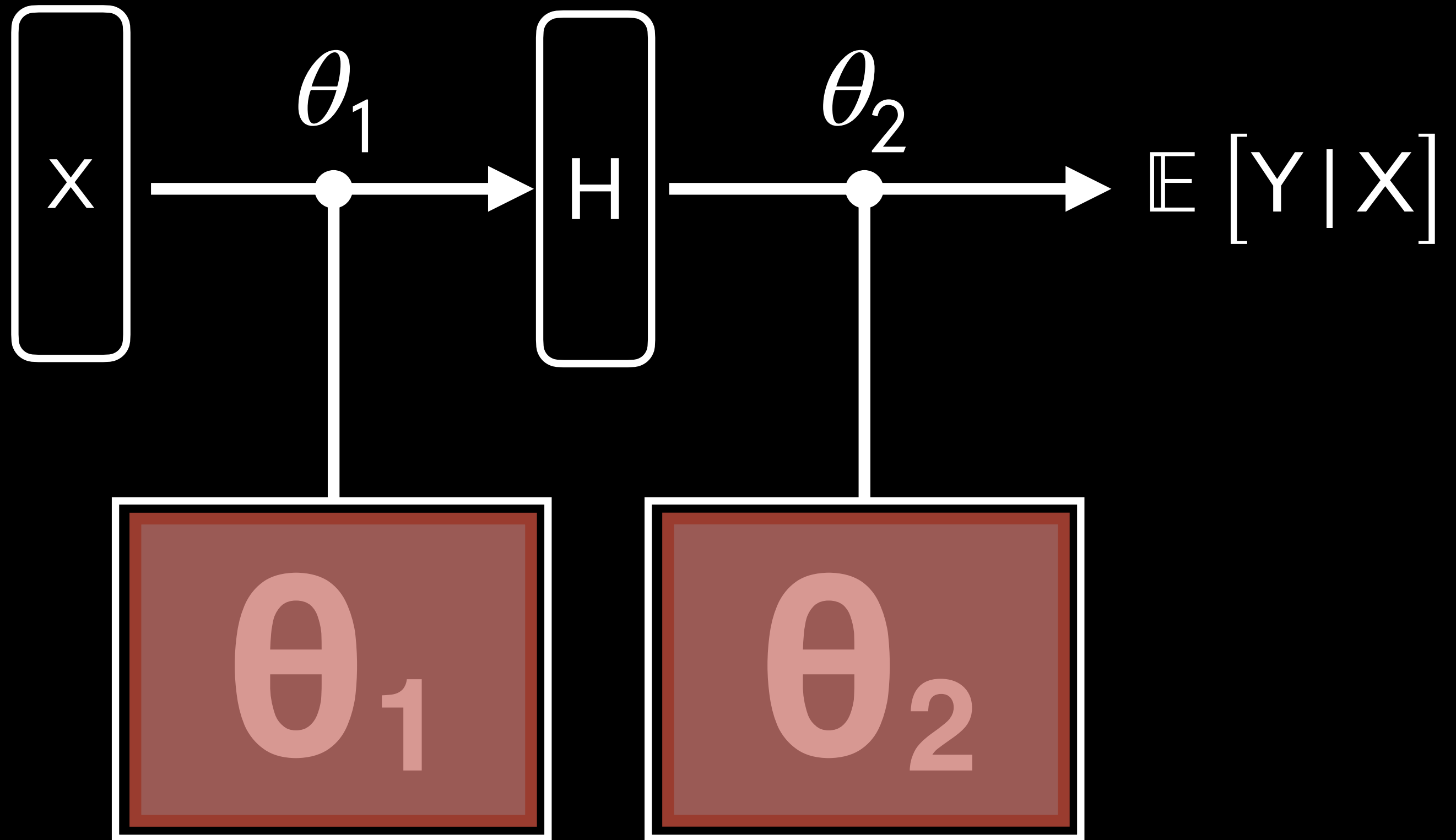
$$= \prod_l^L q(\theta_l) = \prod_l^L \prod_d^{D_l} q(\theta_{l,d})$$

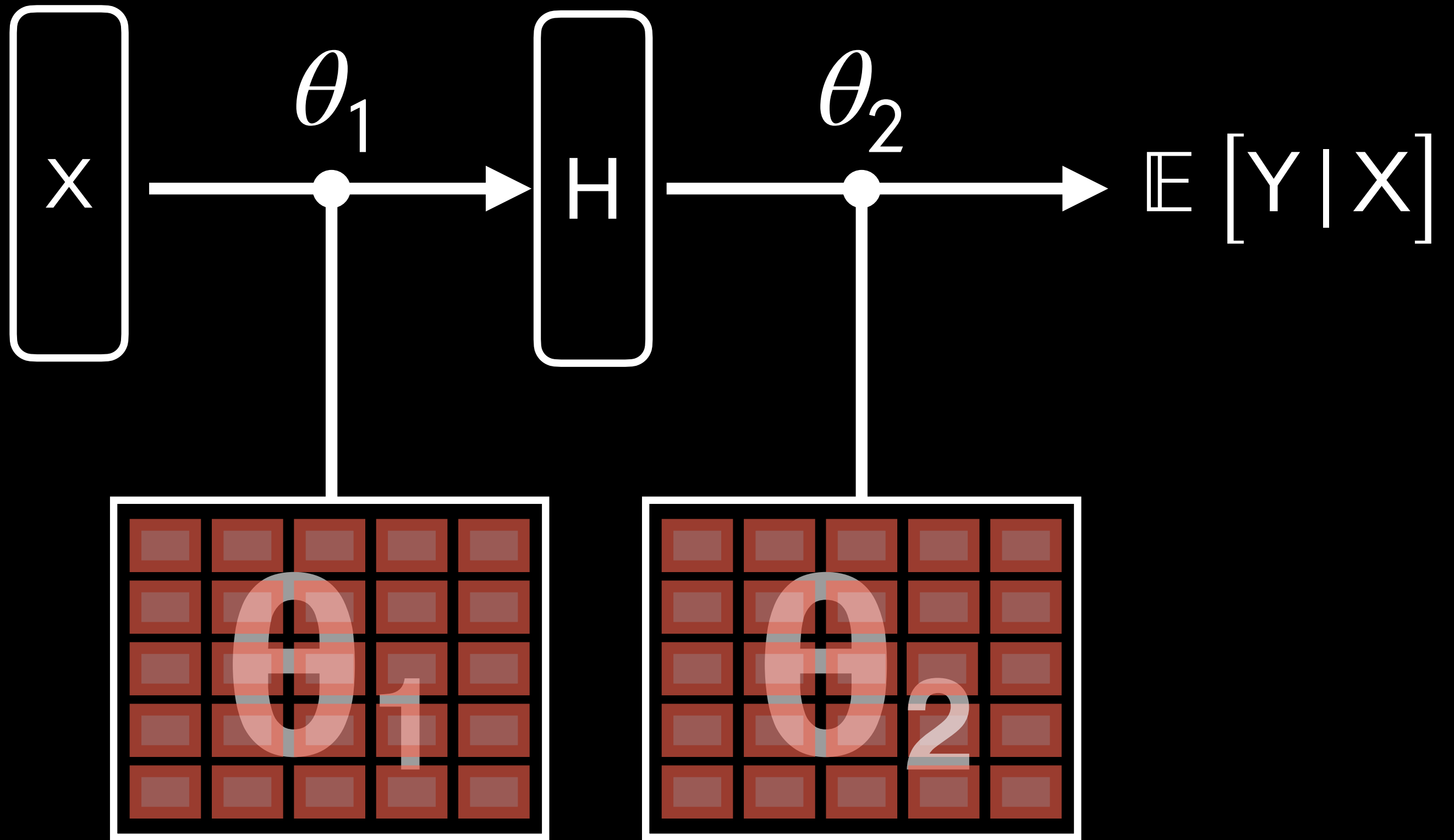
factorize over layers

factorize within layers









Outline

- ⊗ background
- ⊗ subnetwork inference algorithm
- ⊗ subnetworks are all you need

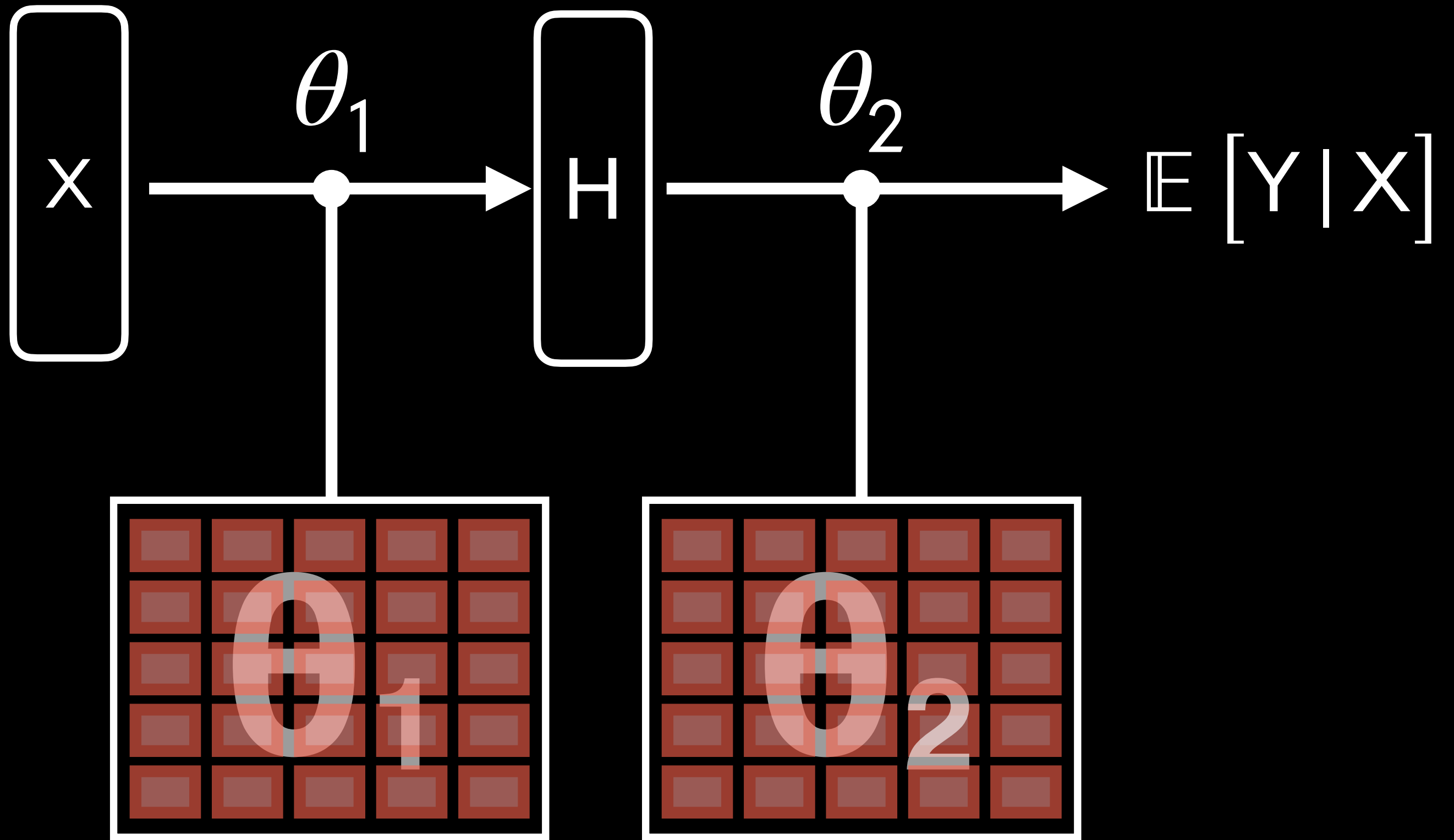
Outline

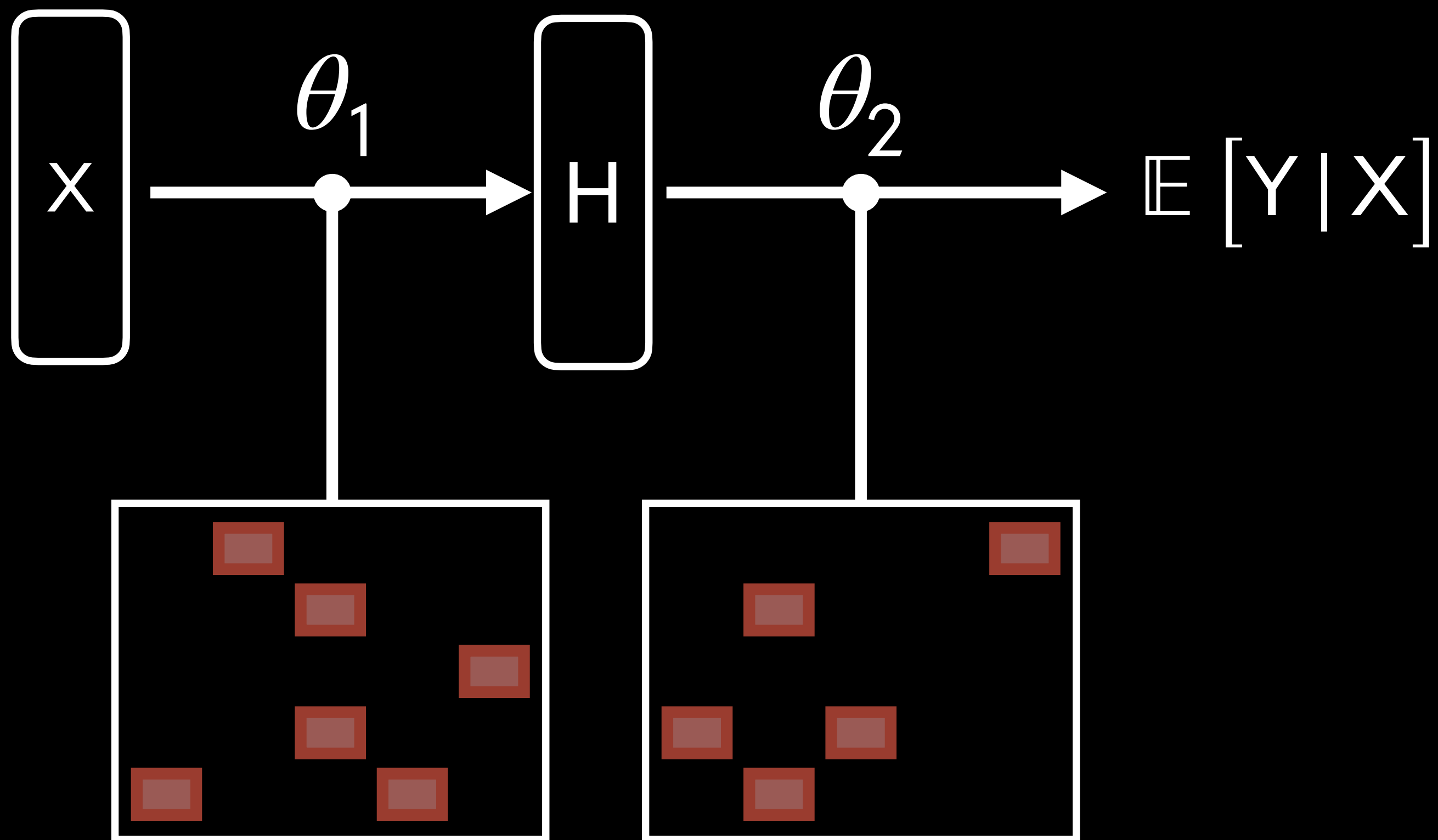
- ⊗ background
- ⊗ subnetwork inference algorithm
- ⊗ subnetworks are all you need

Lottery Ticket Hypothesis

Lottery Ticket Hypothesis: *feed-forward networks contain subnetworks ("winning tickets") that—when trained in isolation—reach test accuracy comparable to the original network.*

[Frankle & Carbin, ICLR 2019 Best Paper]





Lottery Ticket Hypothesis for BNNs

Can the posterior distribution over all weights be represented as a posterior over a subnetwork? (in terms of inducing equivalent predictive distributions)

Lottery Ticket Hypothesis for BNNs

$$p(\theta_1, \dots, \theta_L | \mathfrak{D}) \stackrel{?}{=}$$

$$p(\{\theta_s | s \in \mathcal{S}\} | \mathfrak{D}) \cdot \prod_{r \in \mathbb{R}} \delta[\theta_r - \bar{\theta}_r]$$

Problem

We can't simply find the true, complete posterior and prune it, analogously to how Frankle & Carbin [2019] find their subnetworks.

So we gave up on investigating a LTH for BNNs...but I'll return to this topic later.

Subnetwork Variational Inference

$$p(\theta_1, \dots, \theta_L | \mathfrak{D}) \stackrel{?}{=}$$

$$p\left(\{\theta_s | s \in \mathbb{S}\} | \mathfrak{D}\right) \cdot \prod_{r \in \mathbb{R}} \delta[\theta_r - \bar{\theta}_r]$$

$$\approx q\left(\{\theta_s | s \in \tilde{\mathbb{S}}\}\right) \cdot \prod_{r \in \tilde{\mathbb{R}}} \delta[\theta_r - \bar{\theta}_r]$$

Subnetwork Variational Inference

Can we have a posterior approximation whose structure is data- / learning- driven?

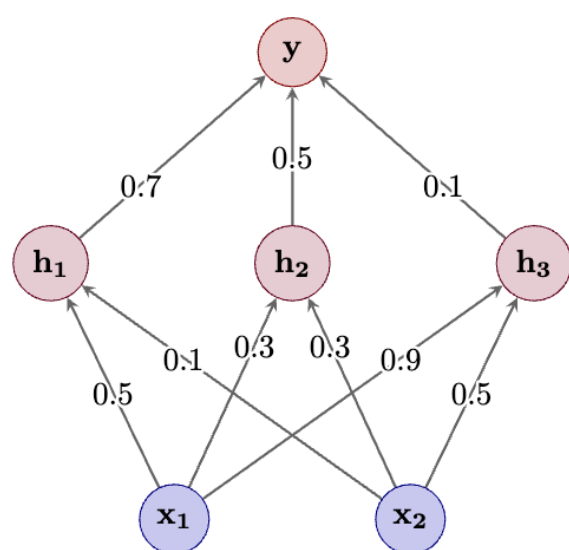
Is it better to perform high-quality, expensive inference over a few parameters than poor, cheap inference over many parameters?

Bayesian Deep Learning via Subnetwork Inference

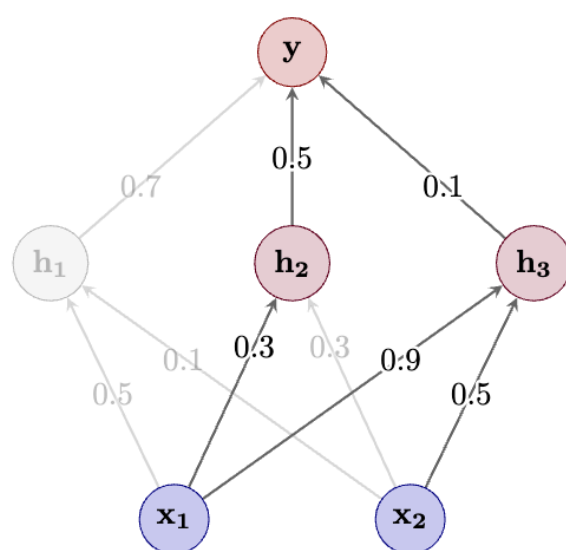
Erik Daxberger^{1,2} Eric Nalisnick^{*3} James Urquhart Allingham^{*1} Javier Antorán^{*1}
José Miguel Hernández-Lobato^{1,4,5}



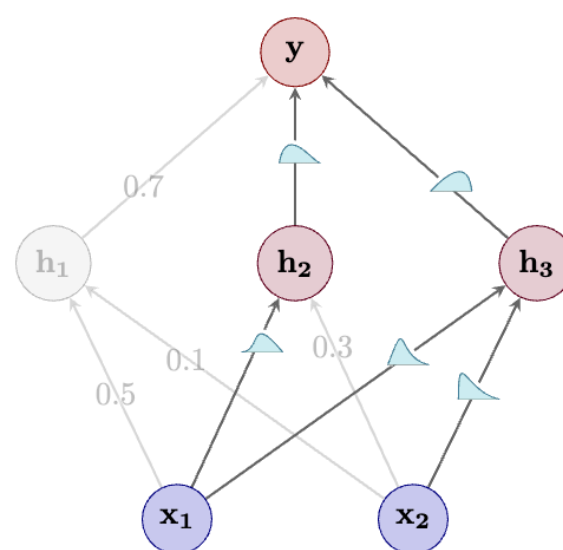
Erik Daxberger



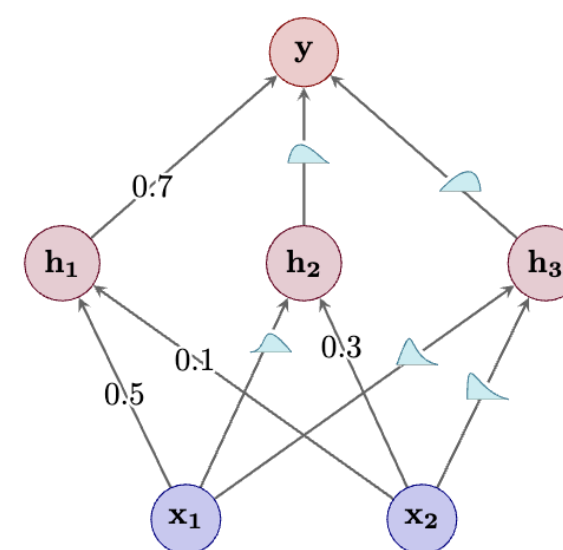
(a) Point Estimation



(b) Subnetwork Selection

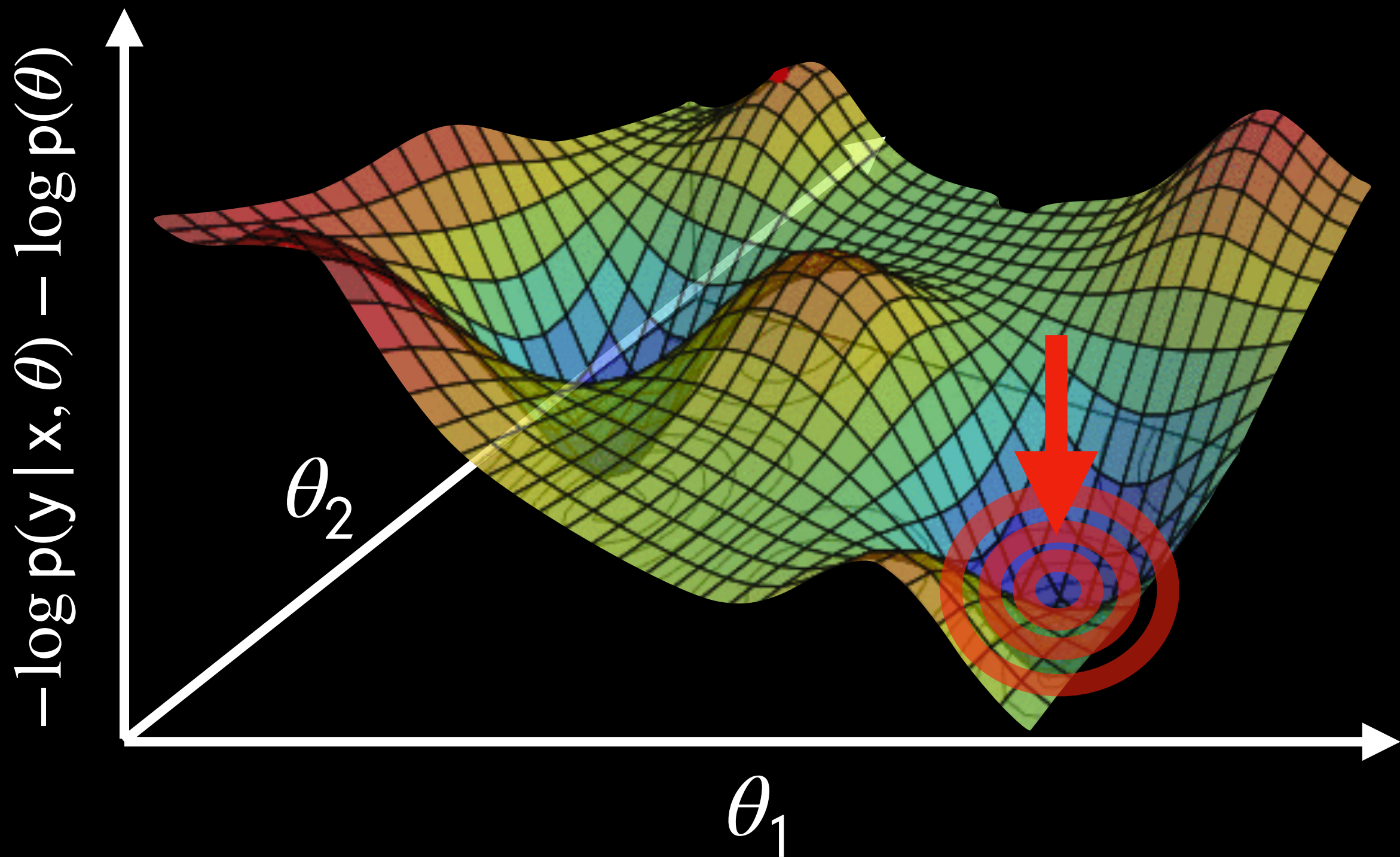


(c) Bayesian Inference

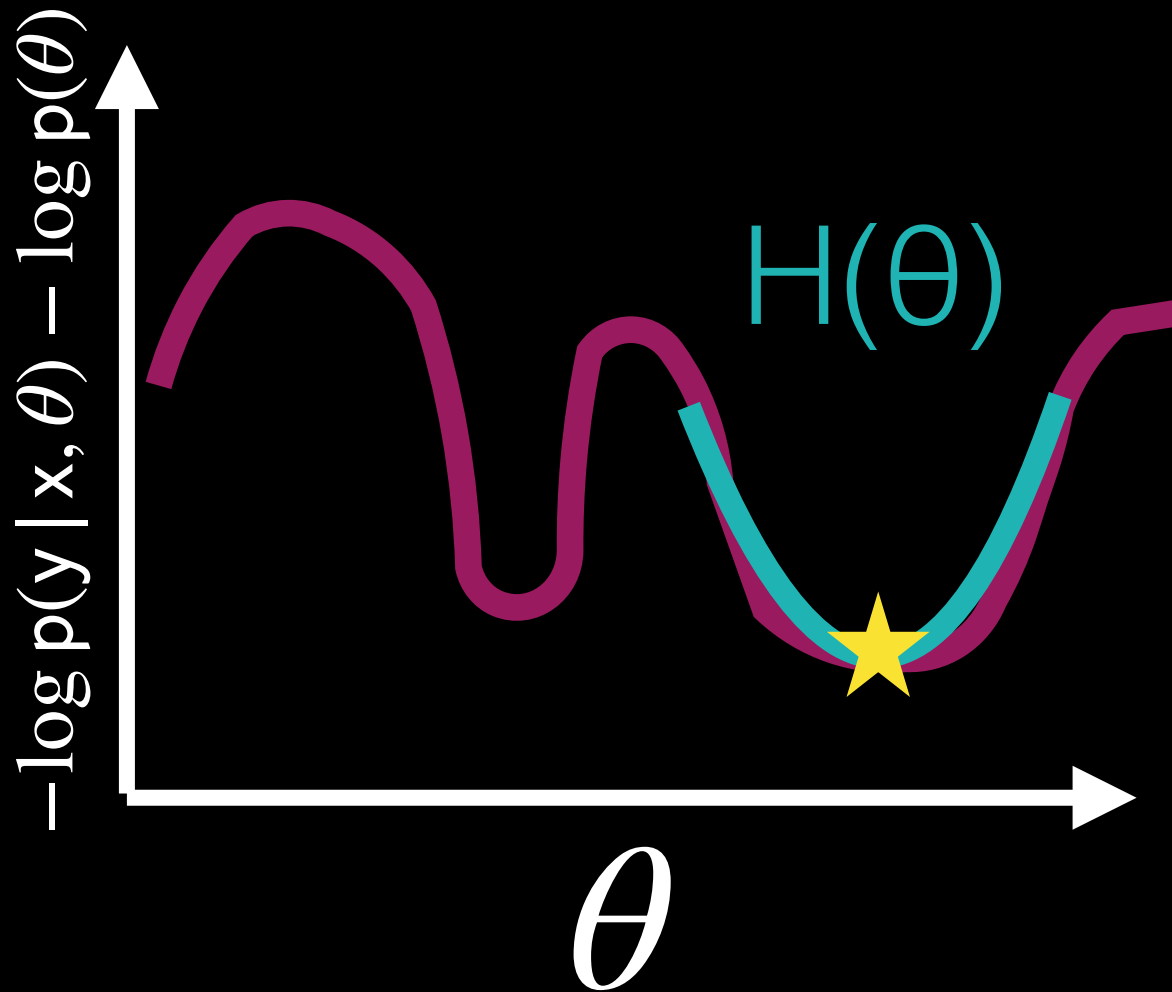


(d) Prediction

Laplace Approximation



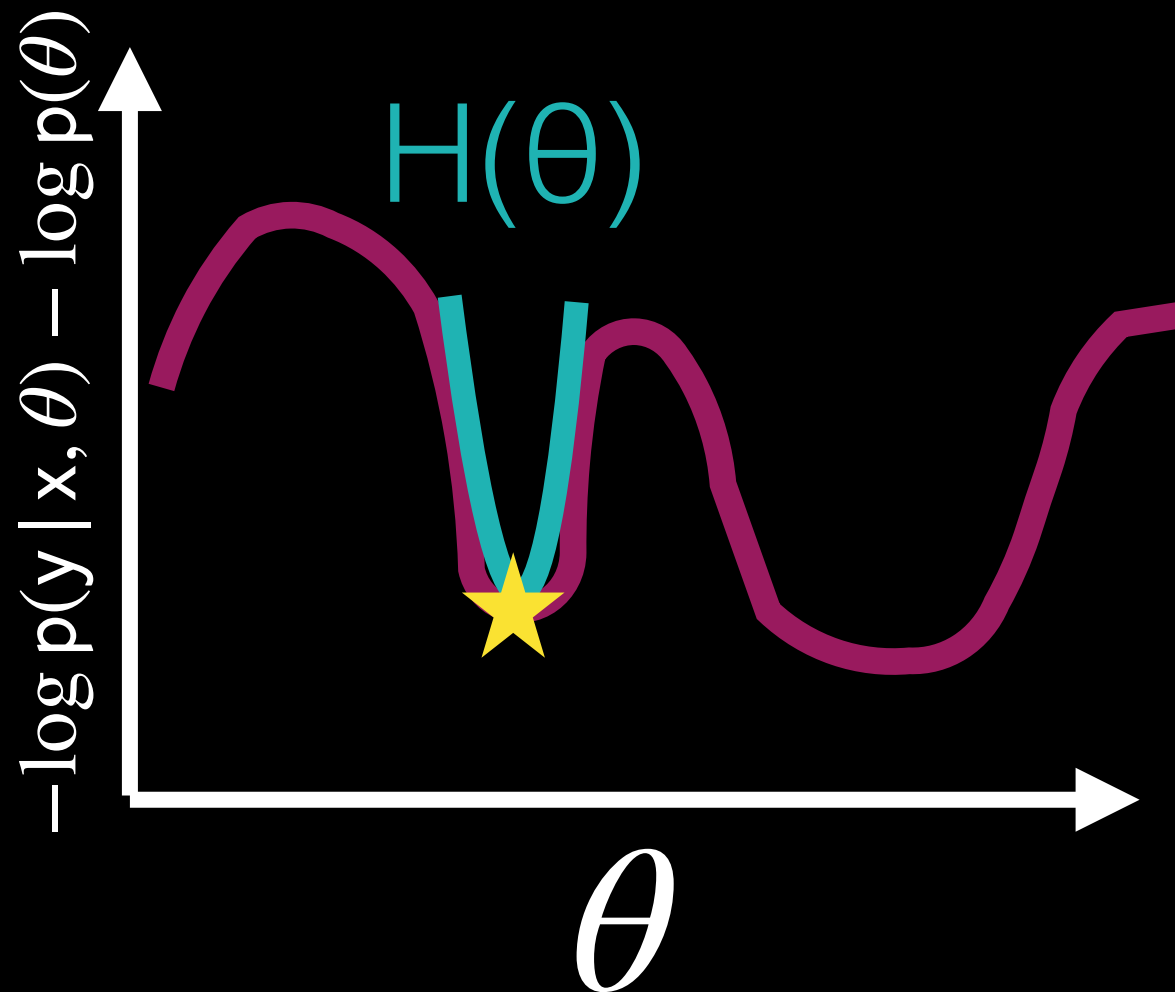
Laplace Approximation



$$N \left(\hat{\theta}_{\text{MAP}}, H^{-1}(\hat{\theta}_{\text{MAP}}) \right)$$

small curvature,
large posterior variance

Laplace Approximation



$$N \left(\hat{\theta}_{\text{MAP}}, H^{-1}(\hat{\theta}_{\text{MAP}}) \right)$$

large curvature,
small posterior variance

Subnetwork Laplace Approximation

1. Find MAP estimate for all weights
2. Select subnetwork via heuristic
3. Construct the Laplace approximation, with a full covariance matrix, over the subnetwork
4. Compute predictive distribution as usual with the posterior from step #3.

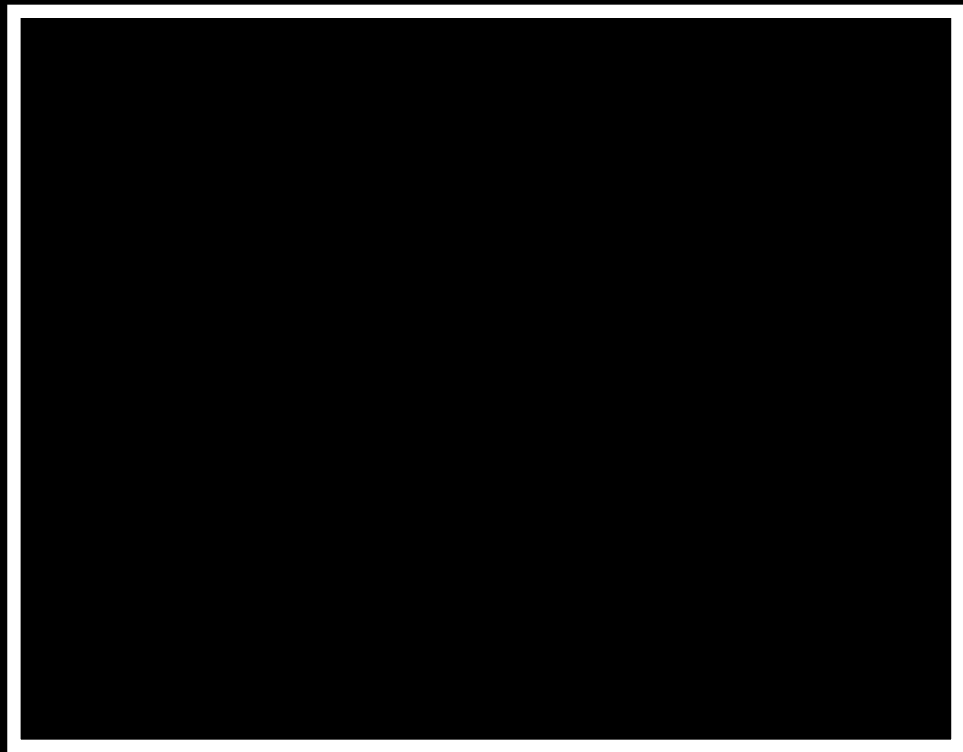
Subnetwork Laplace Approximation

1. Find MAP estimate for all weights
2. Select subnetwork via heuristic
3. Construct the Laplace approximation, with a full covariance matrix, over the subnetwork
4. Compute predictive distribution as usual with the posterior from step #3.

Subnetwork Laplace Approximation

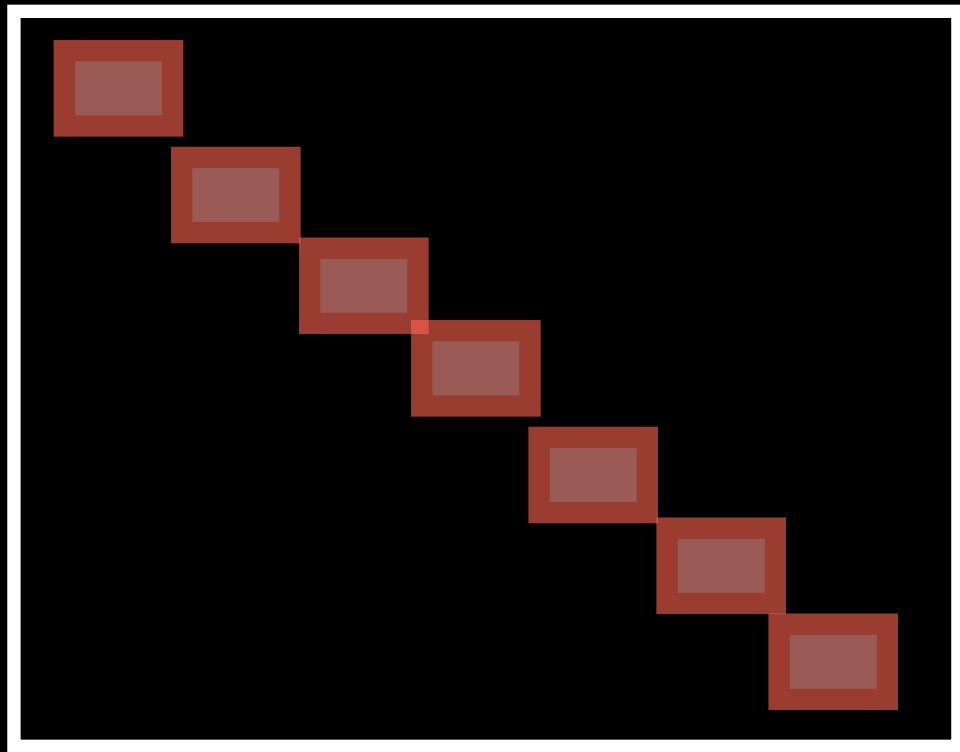
1. Find MAP estimate for all weights
2. Select subnetwork via heuristic
3. Construct the Laplace approximation, with a full covariance matrix, over the subnetwork
4. Compute predictive distribution as usual with the posterior from step #3.

Subnetwork Selection



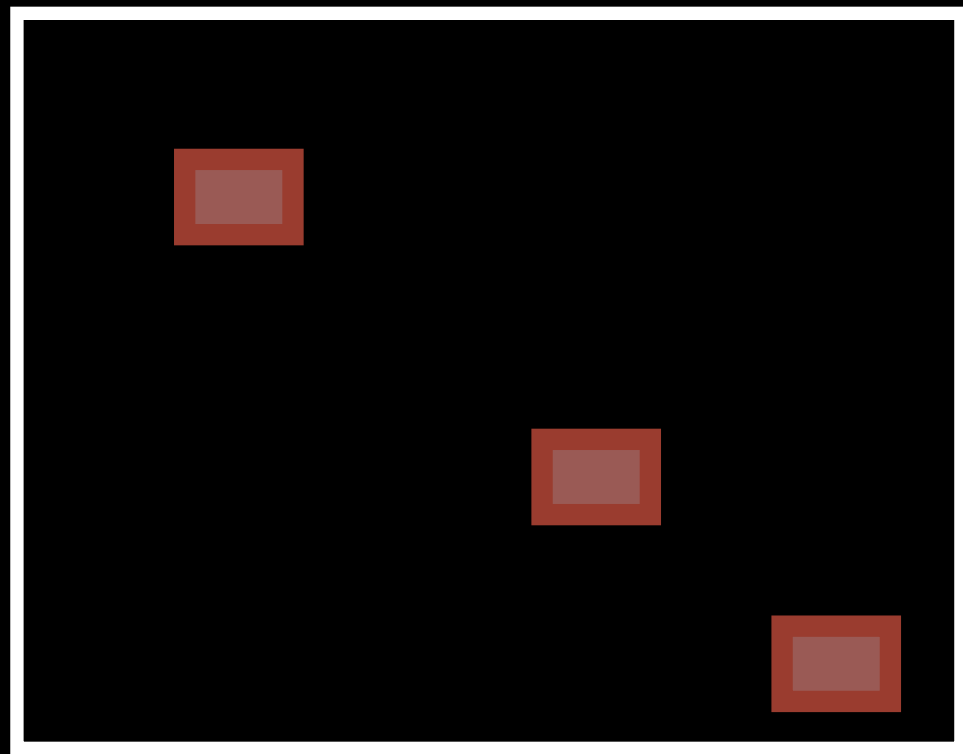
$$\mathbf{H}^{-1} \left(\hat{\theta}_{\text{MAP}} \right)$$

Subnetwork Selection



$$\mathbf{H}^{-1} \left(\hat{\theta}_{\text{MAP}} \right)$$

Subnetwork Selection



$$H^{-1} \left(\hat{\theta}_{\text{MAP}} \right)$$

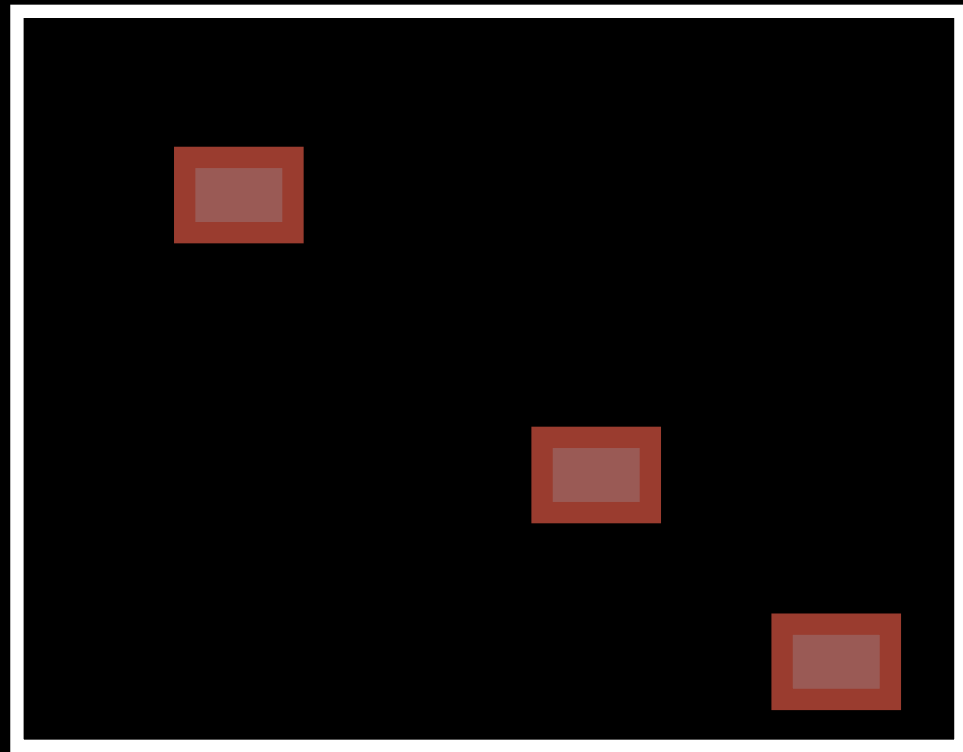
Subnetwork Laplace Approximation

1. Find MAP estimate for all weights
2. Select subnetwork via heuristic
3. Construct the Laplace approximation, with a full covariance matrix, over the subnetwork
4. Compute predictive distribution as usual with the posterior from step #3.

Subnetwork Laplace Approximation

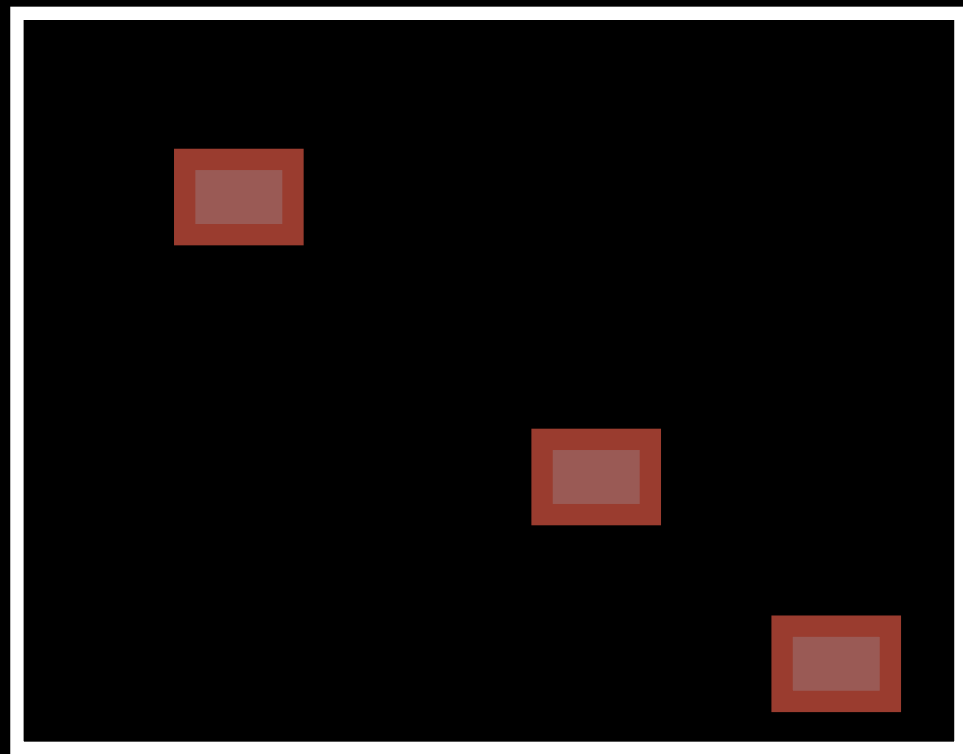
1. Find MAP estimate for all weights
2. Select subnetwork via heuristic
3. Construct the Laplace approximation, with a full covariance matrix, over the subnetwork
4. Compute predictive distribution as usual with the posterior from step #3.

Posterior Construction

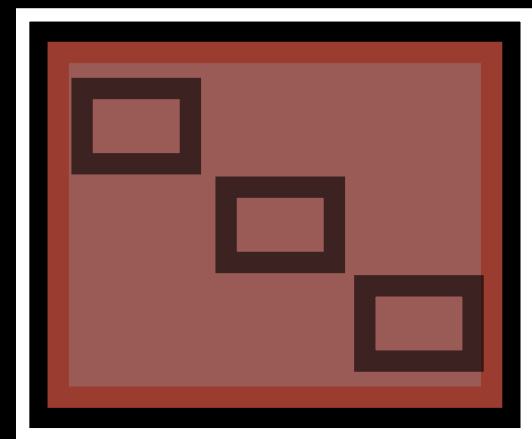


$$H^{-1} \left(\hat{\theta}_{\text{MAP}} \right)$$

Posterior Construction



$$H^{-1} \left(\hat{\theta}_{\text{MAP}} \right)$$



$$H^{-1} \left(\hat{\theta}_{\text{S}} \right)$$

Posterior Construction

$$p(\theta | \mathfrak{D}) \approx$$

$$N\left(\hat{\theta}_{\mathcal{S}}, H^{-1}(\hat{\theta}_{\mathcal{S}})\right) \cdot \prod_{r \in \mathbb{R}} \delta\left[\theta_r - \hat{\theta}_r\right]$$

Subnetwork Laplace Approximation

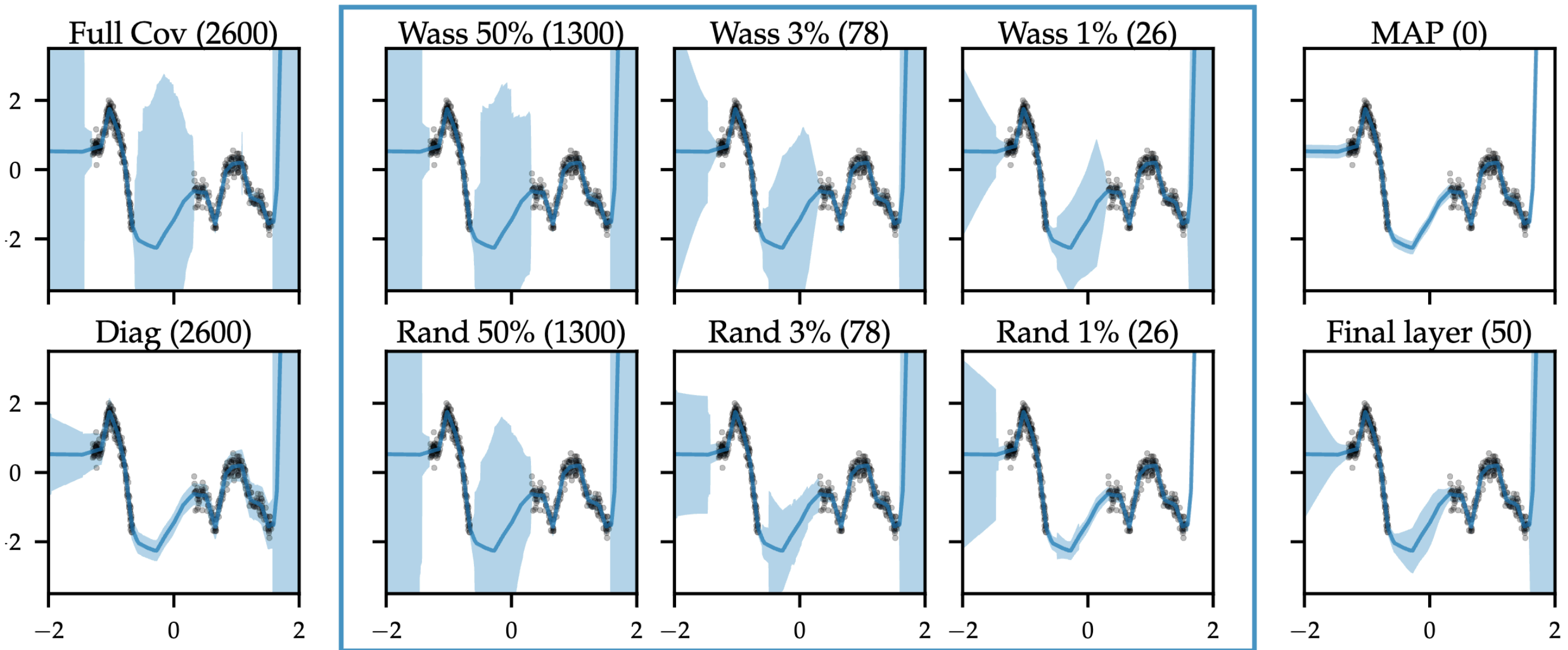
1. Find MAP estimate for all weights
2. Select subnetwork via heuristic
3. Construct the Laplace approximation, with a full covariance matrix, over the subnetwork
4. Compute predictive distribution as usual with the posterior from step #3.

Subnetwork Laplace Approximation

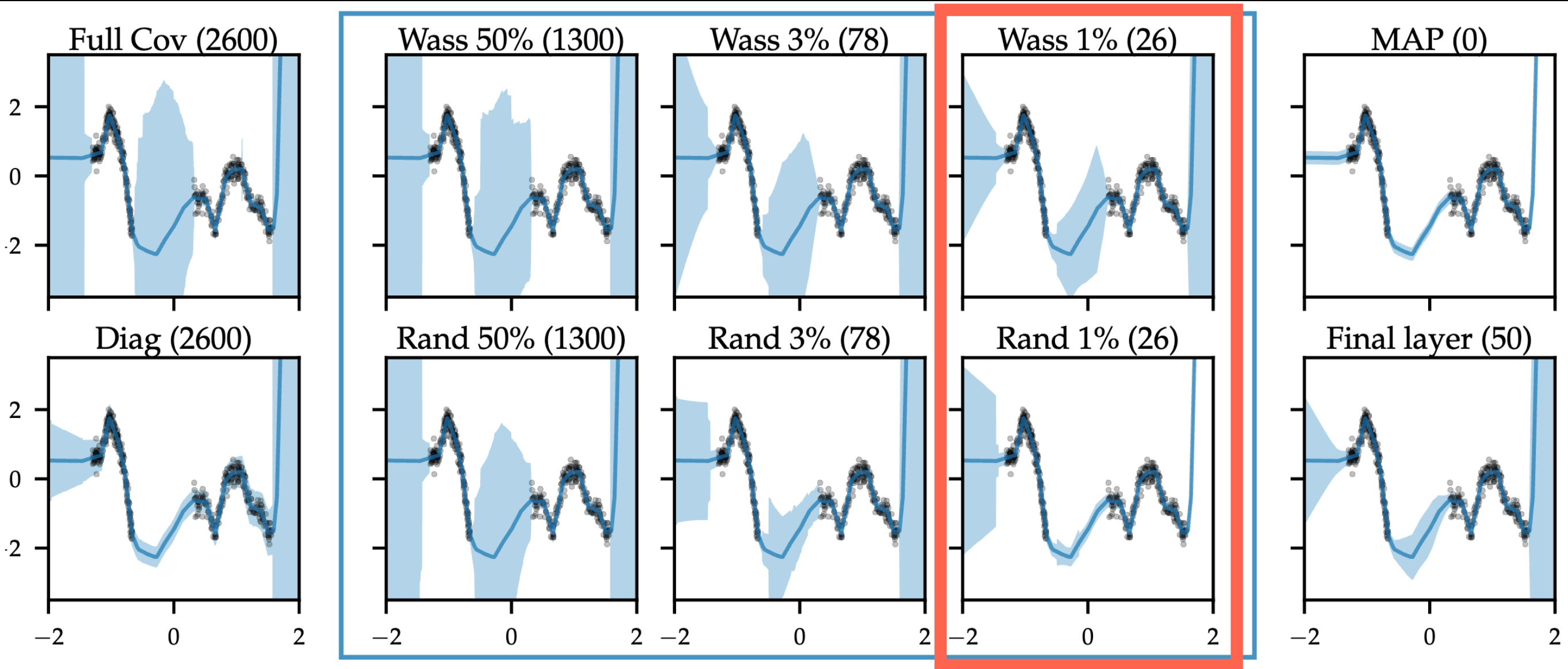
1. Find MAP estimate for all weights
2. Select subnetwork via heuristic
3. Construct the Laplace approximation, with a full covariance matrix, over the subnetwork
4. Compute predictive distribution as usual with the posterior from step #3.

Simulation Results

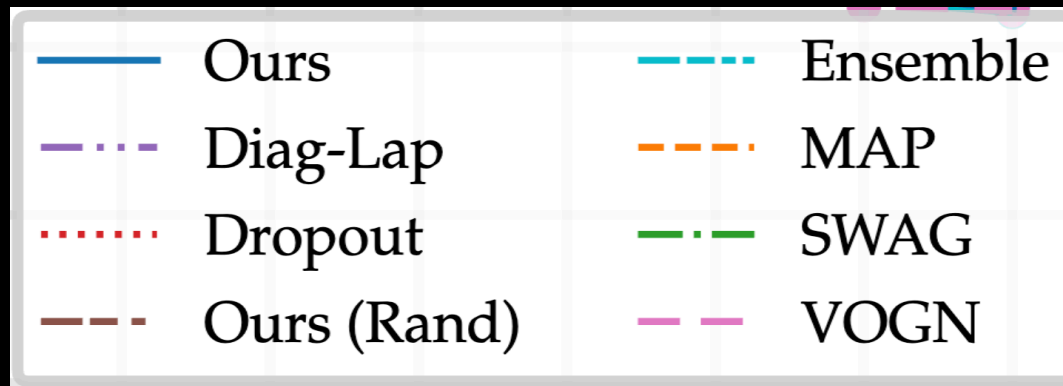
Simulation Results



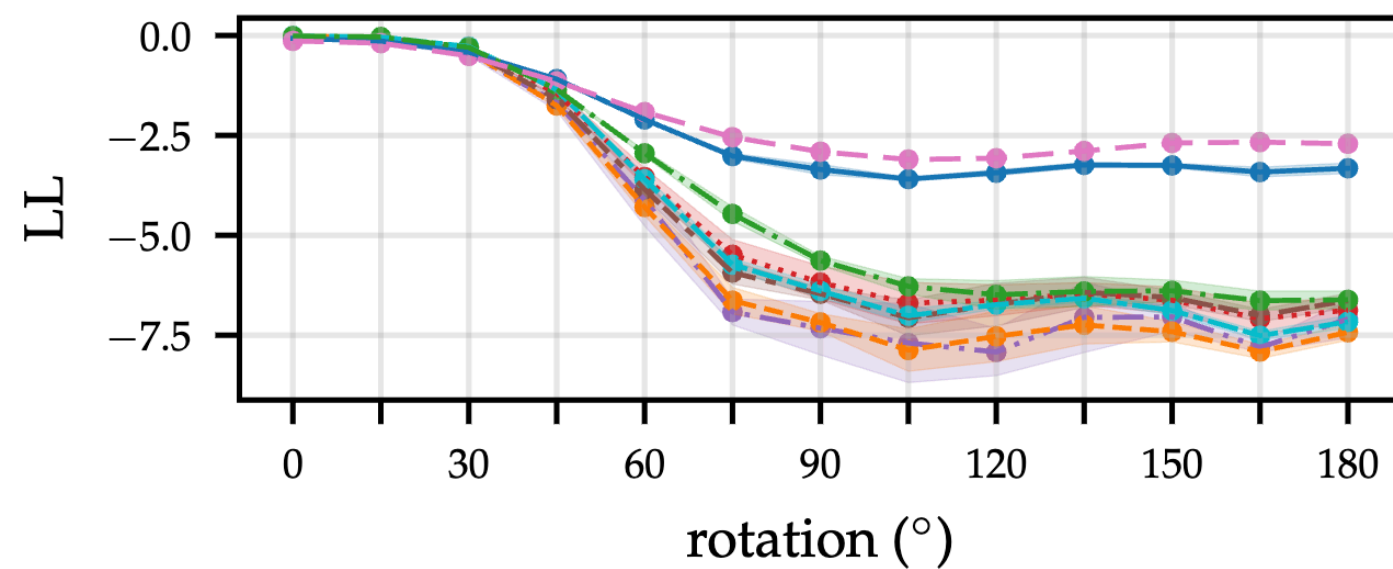
Simulation Results



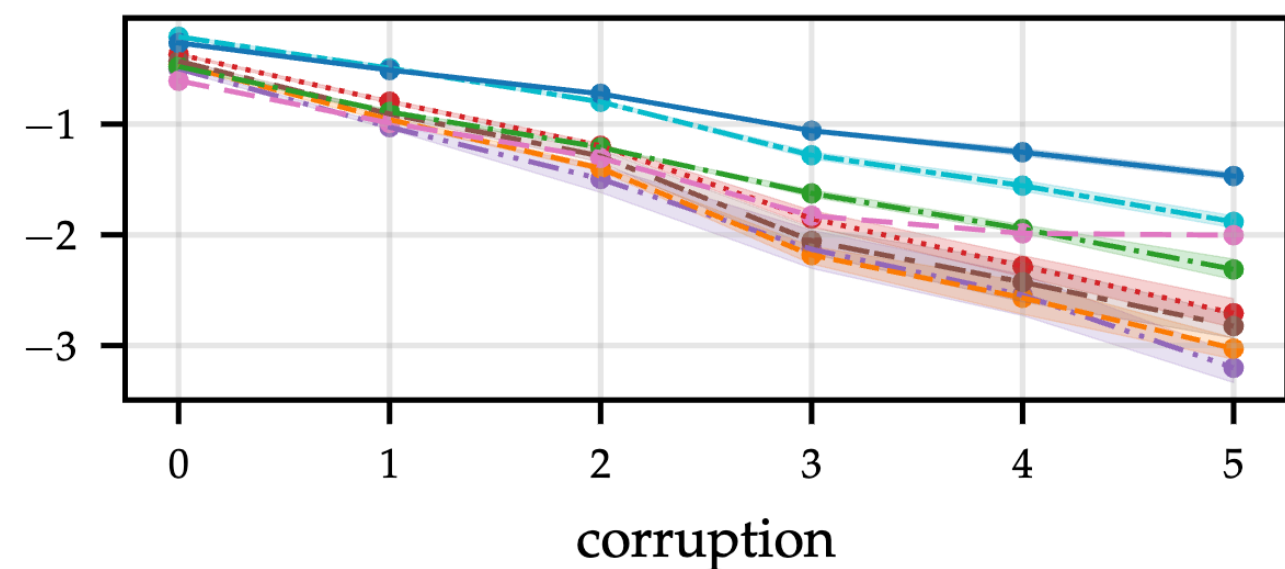
Robustness Results



Rotated MNIST



Corrupted CIFAR-10



Outline

- ⊗ background
- ⊗ subnetwork inference algorithm
- ⊗ subnetworks are all you need

Outline

- ⊗ background
- ⊗ subnetwork inference algorithm
- ⊗ subnetworks are all you need

Do Bayesian Neural Networks Need To Be Fully Stochastic?

Mrinank Sharma
University of Oxford

Sebastian Farquhar
University of Oxford

Eric Nalisnick
University of Amsterdam

Tom Rainforth
University of Oxford



Mrinank Sharma

$$p(\theta_1, \dots, \theta_L | \mathfrak{D}) \stackrel{?}{=}$$

$$p(\{\theta_s | s \in \mathcal{S}\} | \mathfrak{D}) \cdot \prod_{r \in \mathbb{R}} \delta[\theta_r - \bar{\theta}_r]$$

Do Bayesian Neural Networks Need To Be Fully Stochastic?

Mrinank Sharma
University of Oxford

Sebastian Farquhar
University of Oxford

Eric Nalisnick
University of Amsterdam

Tom Rainforth
University of Oxford



Mrinank Sharma

$$p(\theta_1, \dots, \theta_L | \mathfrak{D}) \stackrel{?}{=}$$

$$p(\{\theta_s | s \in \mathcal{S}\} | \mathfrak{D}) \cdot \prod_{r \in \mathbb{R}} \delta[\theta_r - \bar{\theta}_r]$$



Do Bayesian Neural Networks Need To Be Fully Stochastic?

Mrinank Sharma
University of Oxford

Sebastian Farquhar
University of Oxford

Eric Nalisnick
University of Amsterdam

Tom Rainforth
University of Oxford



Mrinank Sharma

$$p(\theta_1, \dots, \theta_L | \mathfrak{D}) \stackrel{?}{=}$$

$$p(\{\theta_s | s \in \mathcal{S}\} | \mathfrak{D}) \cdot \prod_{r \in \mathbb{R}} \delta[\theta_r - \bar{\theta}_r]$$



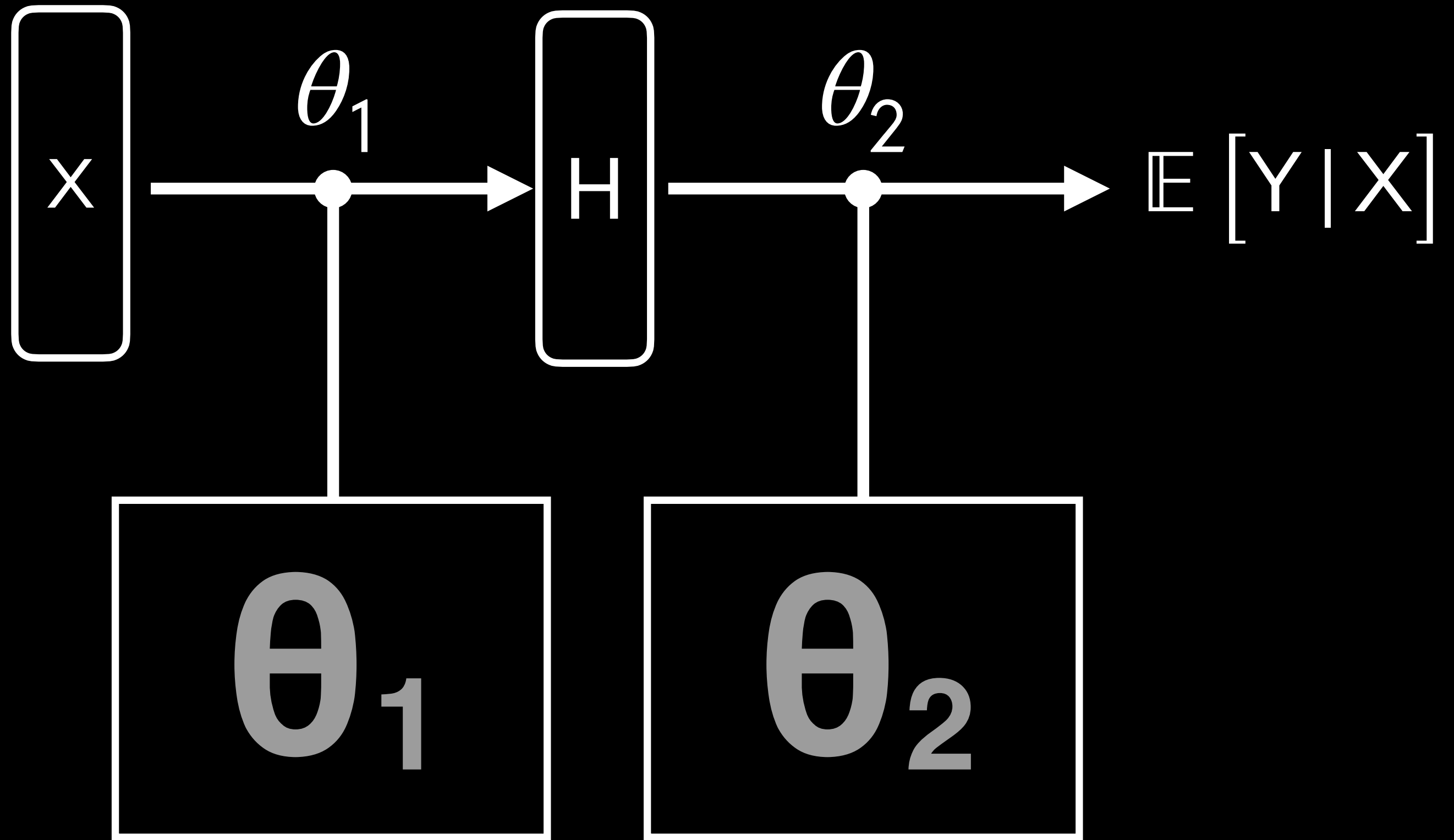
Theoretical Result

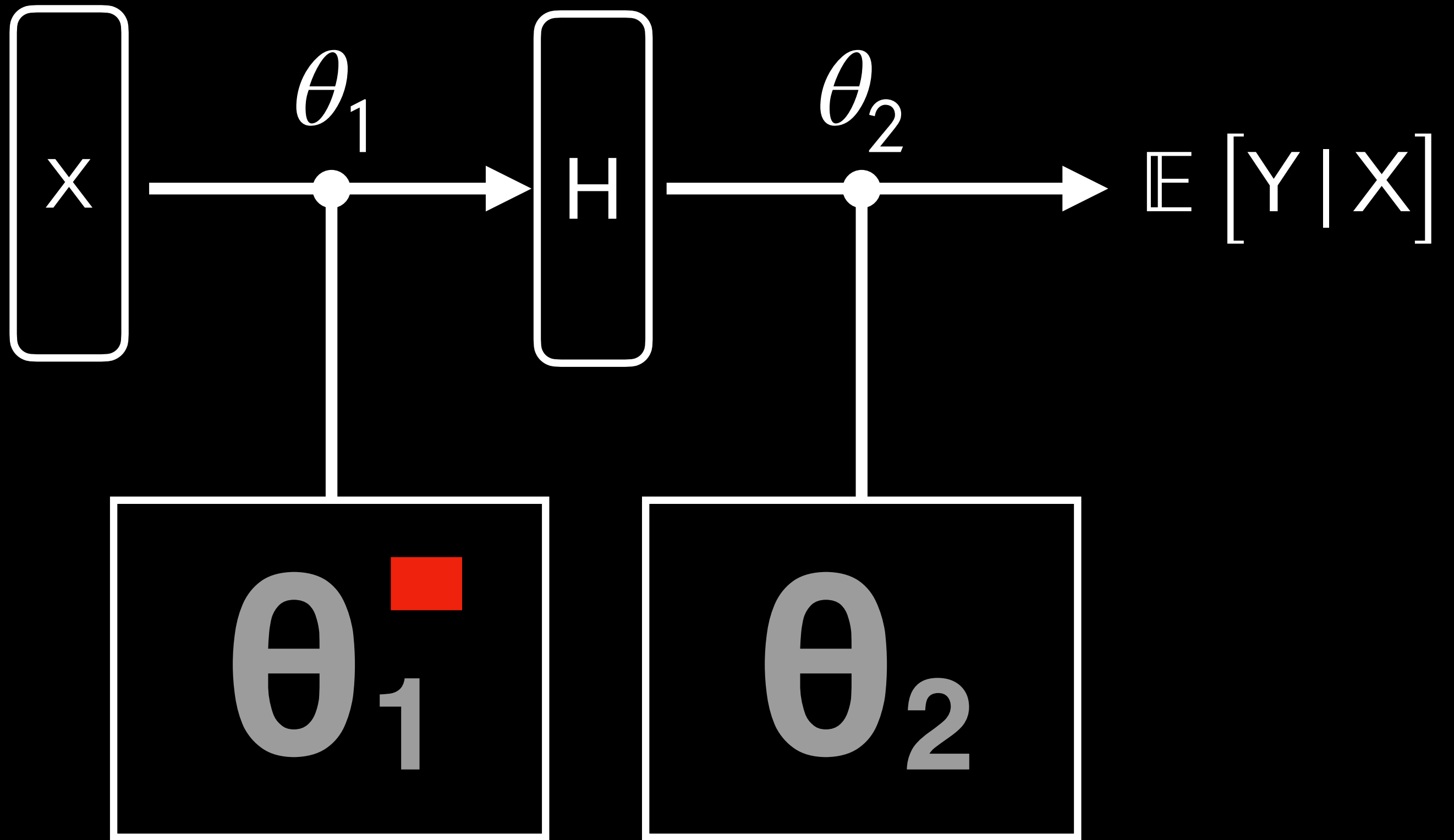
Informal: a multi-layer perceptron, with at least one hidden layer, can represent arbitrary predictive distributions, as long as there is at least one hidden layer between its stochastic variable(s) and output.

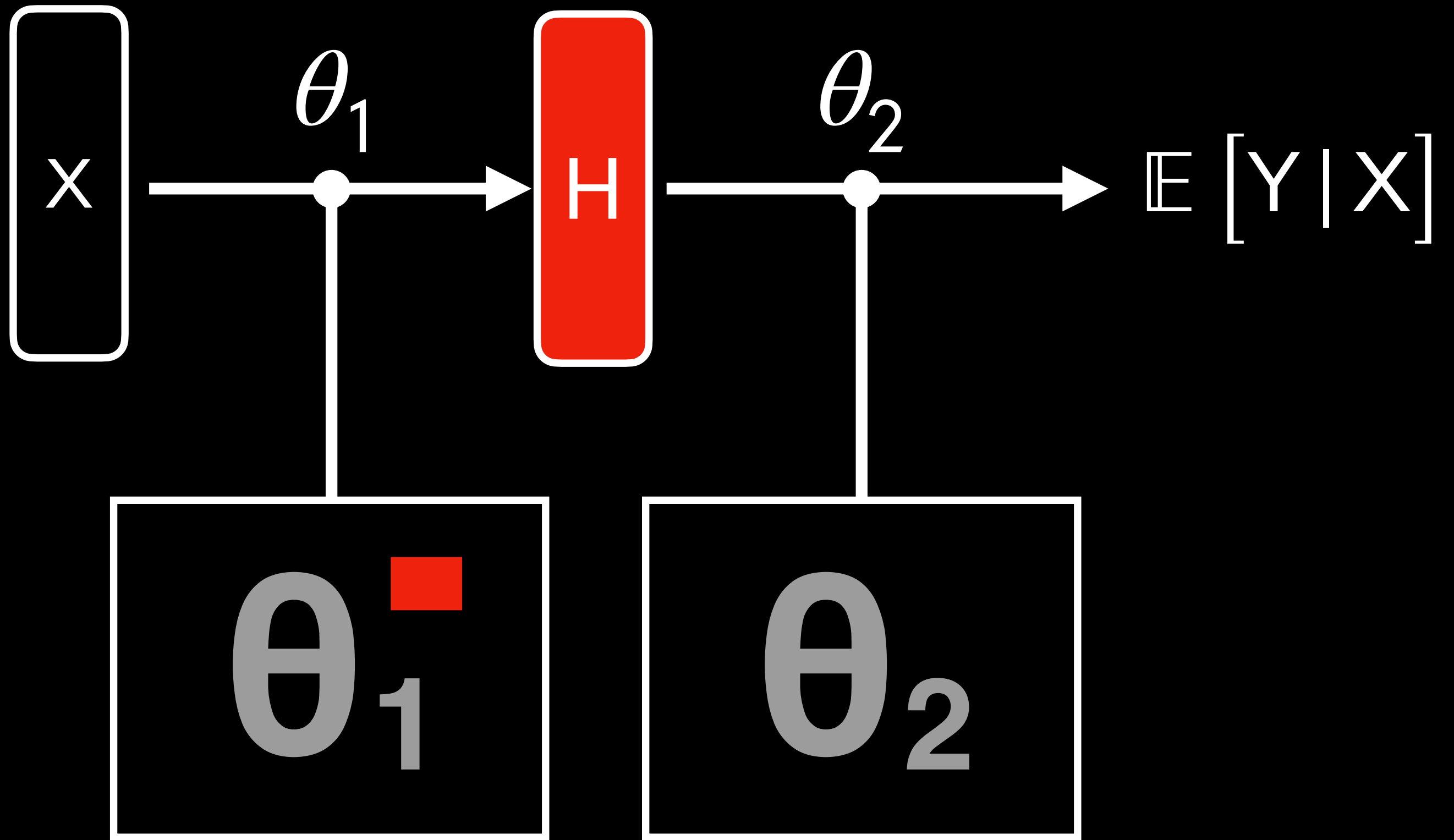
Theoretical Result

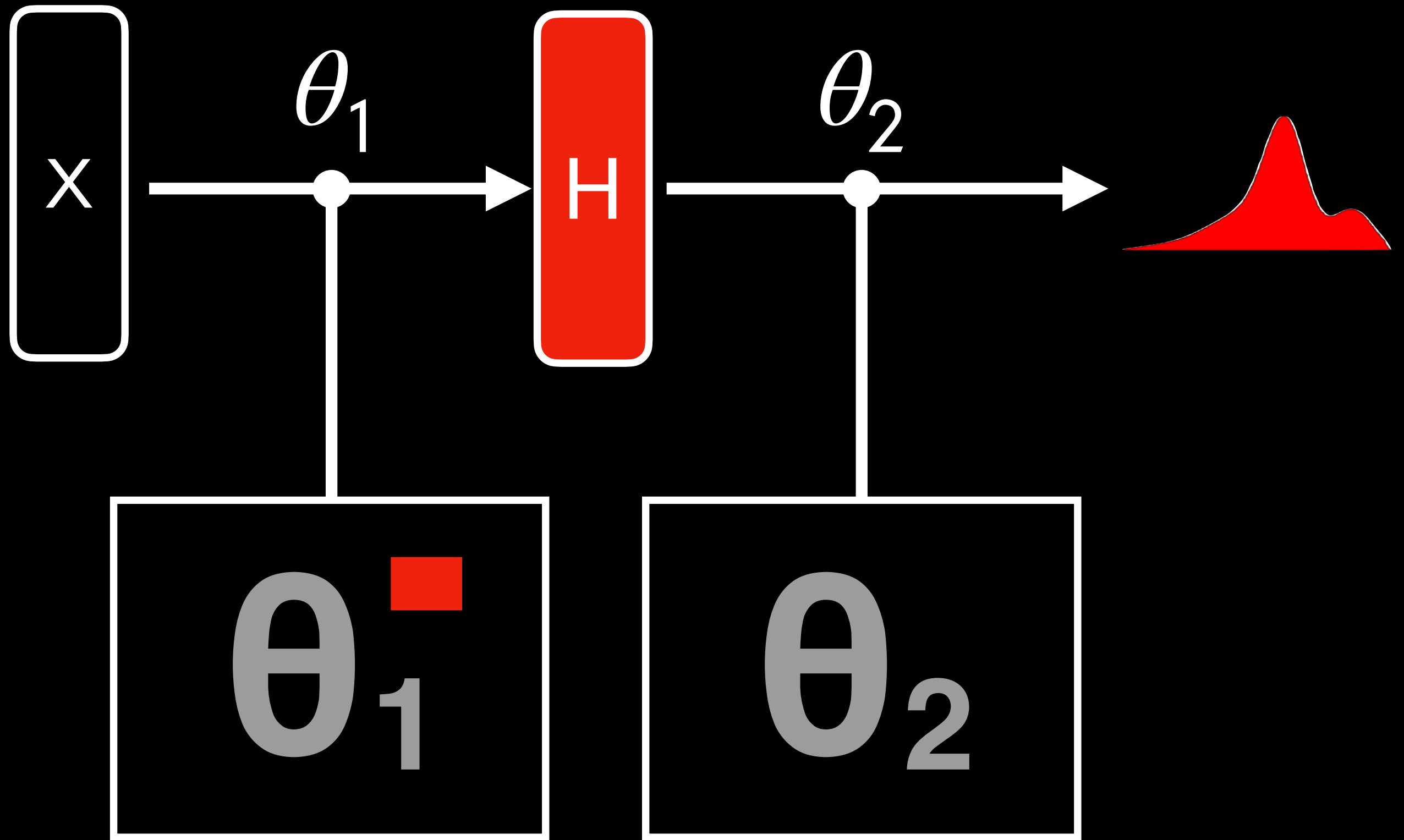
Informal: a multi-layer perceptron, with at least one hidden layer, can represent arbitrary predictive distributions, as long as there is at least one hidden layer between its stochastic variable(s) and output.

Proof sketch: Combine the *noise outsourcing lemma* [Austin, 2012] with the *universal approximation theorem* [Leshno, 1993]









Consequence

Doing posterior inference for many / all of a BNN's parameters is overkill!

Unfortunately, the theory is too blunt to give any more advice about how many stochastic variables to use and where to place them (except for not in the last layer).

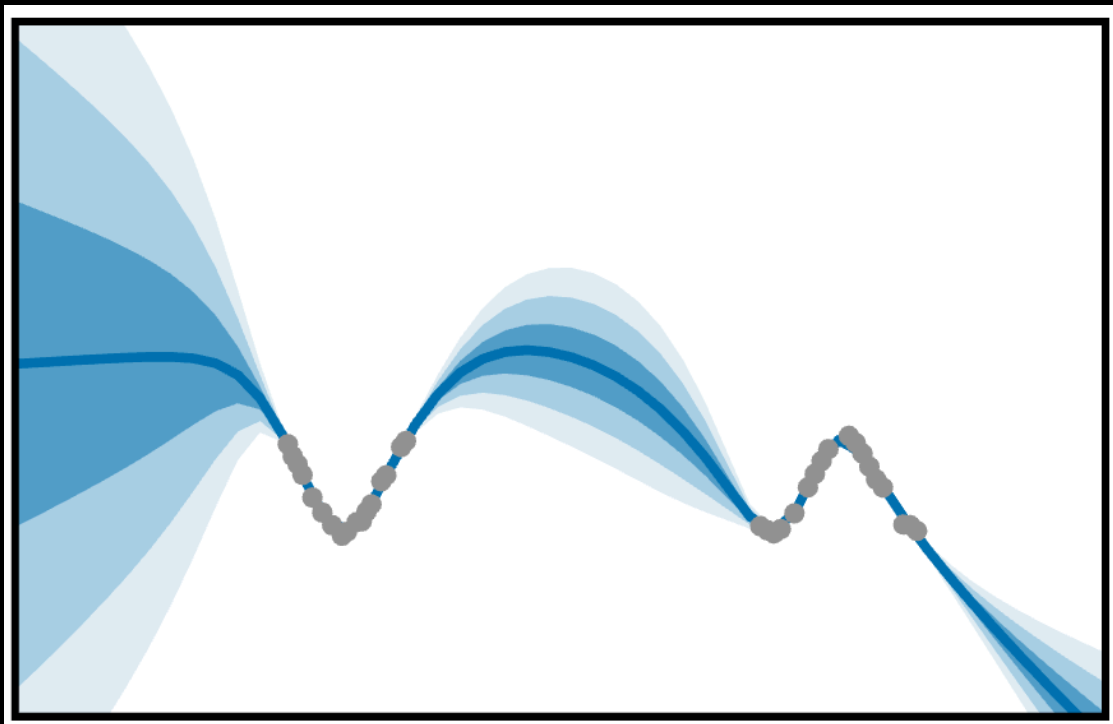
Experimental Results

Do we ever see a systematic benefit to having more stochastic variables?

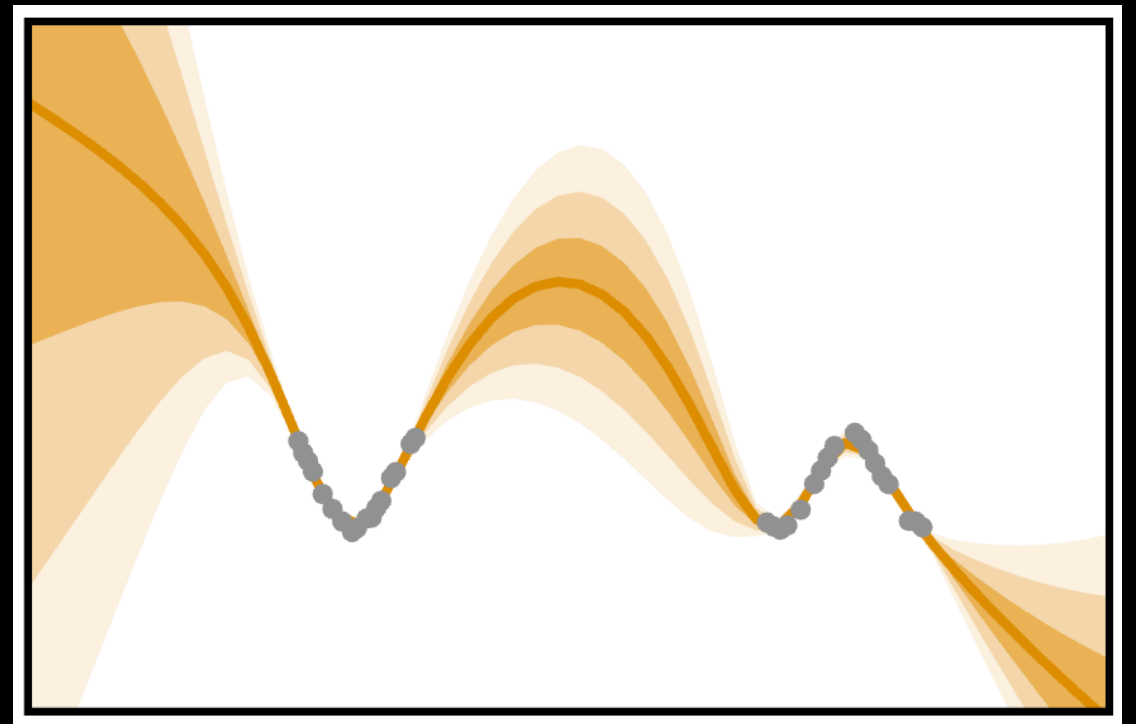
Experimental Results

Do we ever see a systematic benefit to having more stochastic variables?

Predictive Distributions



HMC on All Layers

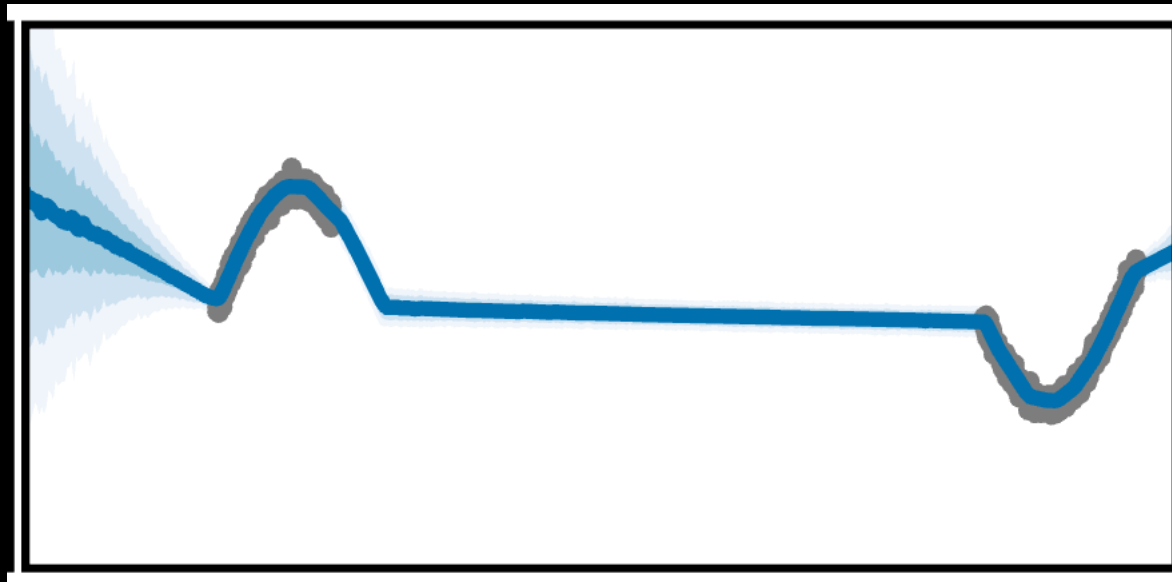


HMC on 1st Layer

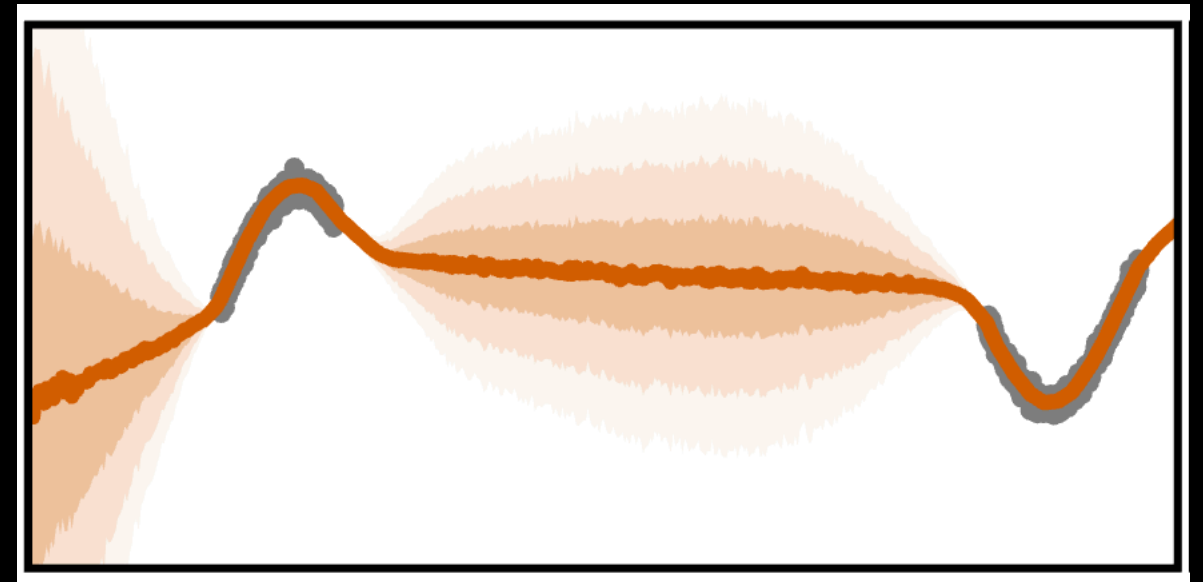
Experimental Results

Do we ever see a systematic benefit to having more stochastic variables?

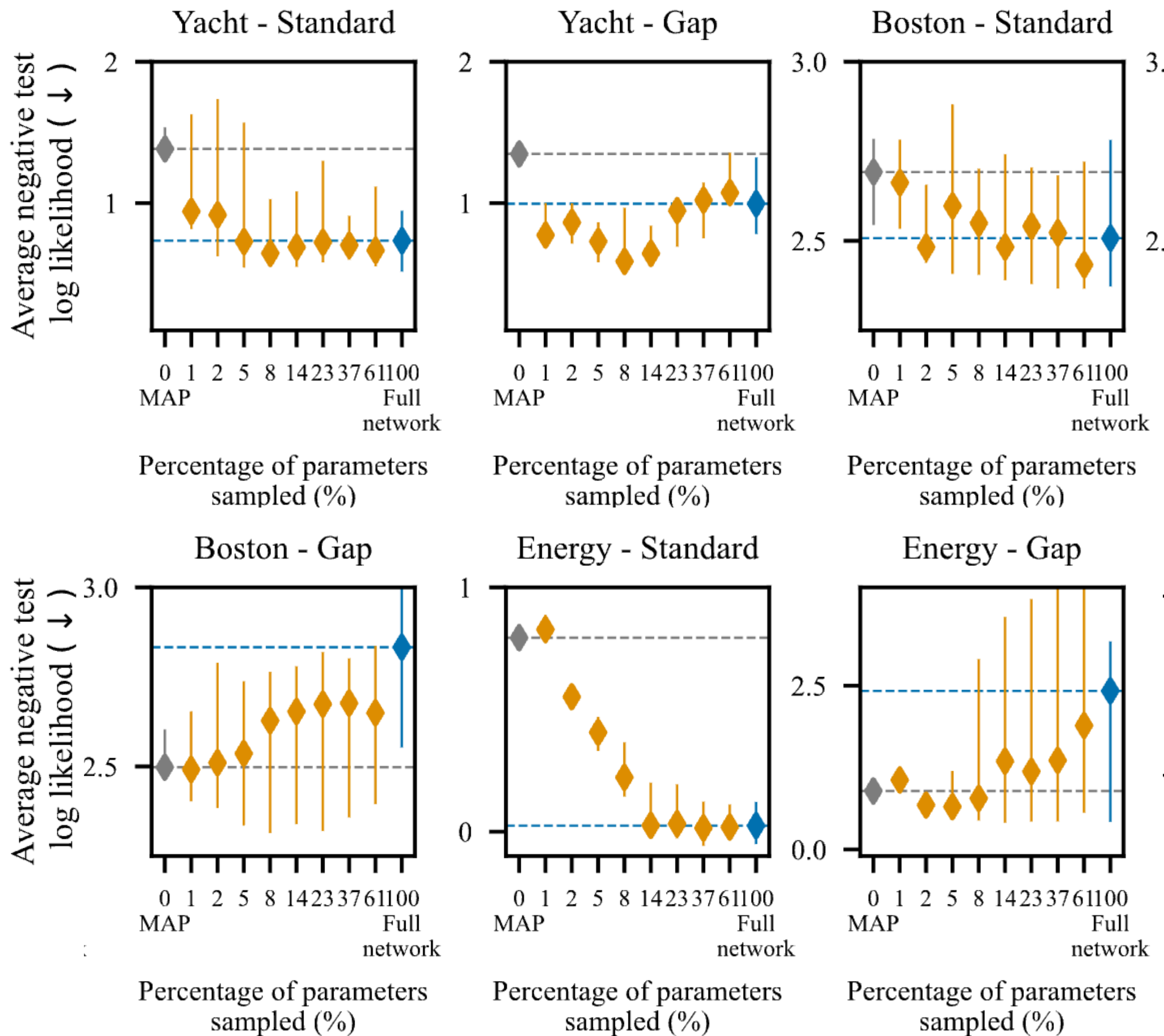
Predictive Distributions



MF-VI on All Layers



MF-VI on Last Layer



Model	CIFAR10		CIFAR100	
	Acc (%)	NLL	Acc (%)	NLL
Deterministic	95.61 ± 0.01	0.187 ± 0.001	79.33 ± 0.45	0.862 ± 0.014
Fully stochastic	94.69 ± 0.07	0.214 ± 0.002	77.68 ± 0.29	0.944 ± 0.002
Input layer stochastic	95.70 ± 0.08	0.187 ± 0.002	79.49 ± 0.15	0.861 ± 0.021
Output layer stochastic	95.60 ± 0.05	0.189 ± 0.001	78.92 ± 0.34	0.933 ± 0.010
Output layer and last block stochastic	95.59 ± 0.08	0.168 ± 0.0005	79.00 ± 0.091	0.834 ± 0.0007

Outline

- ⊗ background
- ⊗ subnetwork inference algorithm
- ⊗ subnetworks are all you need

Outline

- ⊗ background
- ⊗ subnetwork inference algorithm
- ⊗ subnetworks are all you need

Conclusions

1. Subnetwork inference is justified both experimentally and theoretically
2. Open problems:
 1. Better methods for choosing the subnetwork
 2. Principled, unified inference algorithms for models with stochastic and deterministic parameters