

Bayesian Sparsification of Neural Networks

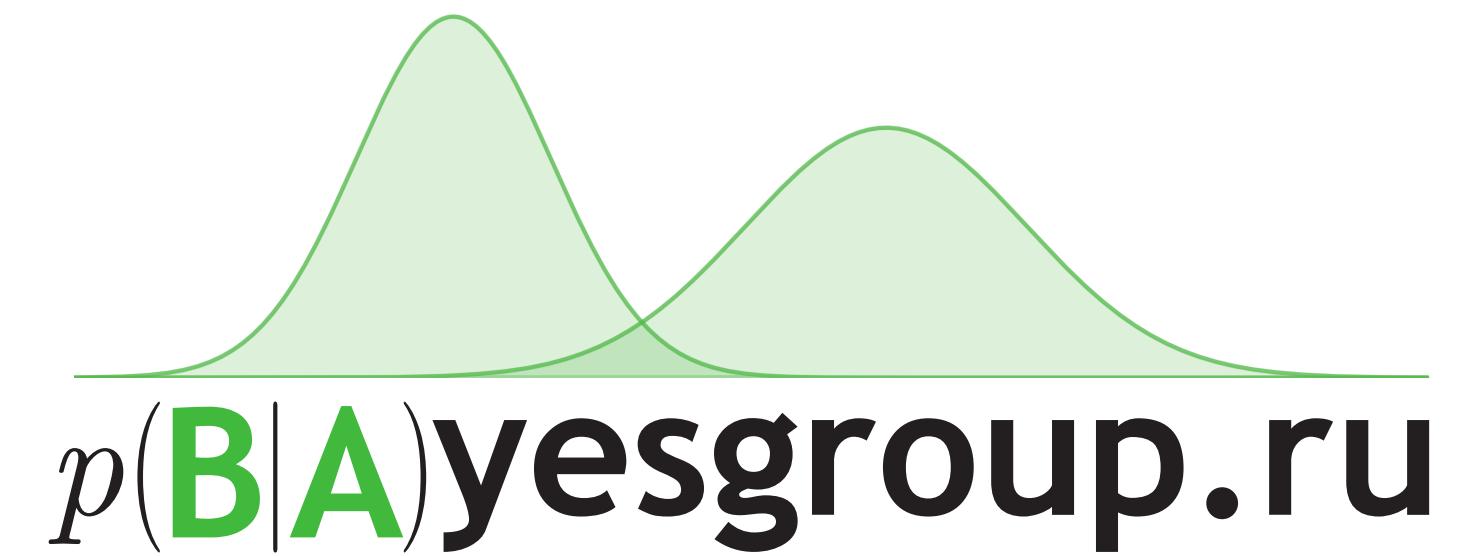


Nadezhda Chirkova

Higher School of Economics, Samsung AI Laboratory
Moscow, Russia

these slides
are here

Nordic Probabilistic AI school
June 3–7, 2019. Trondheim, Norway



Agenda

- Sparsification: what and why
- Bayesian neural networks
- Sparse variational dropout
- Practical assignment: implementation of SparseVD
- Model enhancements

Agenda

- Sparsification: what and why
- Bayesian neural networks
- Sparse variational dropout
- Practical assignment: implementation of SparseVD
- Model enhancements

Compression of neural networks

- Deep neural networks achieve state-of-the-art performance in a variety of domains
- Model quality scales with model and dataset size
- State-of-the-art models usually incorporate **tens of millions of parameters**
- But **resources** (memory, processing time) **may be limited**



Compression of neural networks

- Deep neural networks achieve state-of-the-art performance in a variety of domains
- Model quality scales with model and dataset size
- State-of-the-art models usually incorporate **tens of millions of parameters**
- But **resources** (memory, processing time) **may be limited**



Urgent industrial problem! well solved using Bayesian deep learning!

Compression of neural networks

- Deep neural networks achieve state-of-the-art performance in a variety of domains
- Model quality scales with model and dataset size
- State-of-the-art models usually incorporate **tens of millions of parameters**
- But **resources** (memory, processing time) **may be limited**

Methods for neural network compression:

- Quantization
- Matrix factorizations
- Sparsification

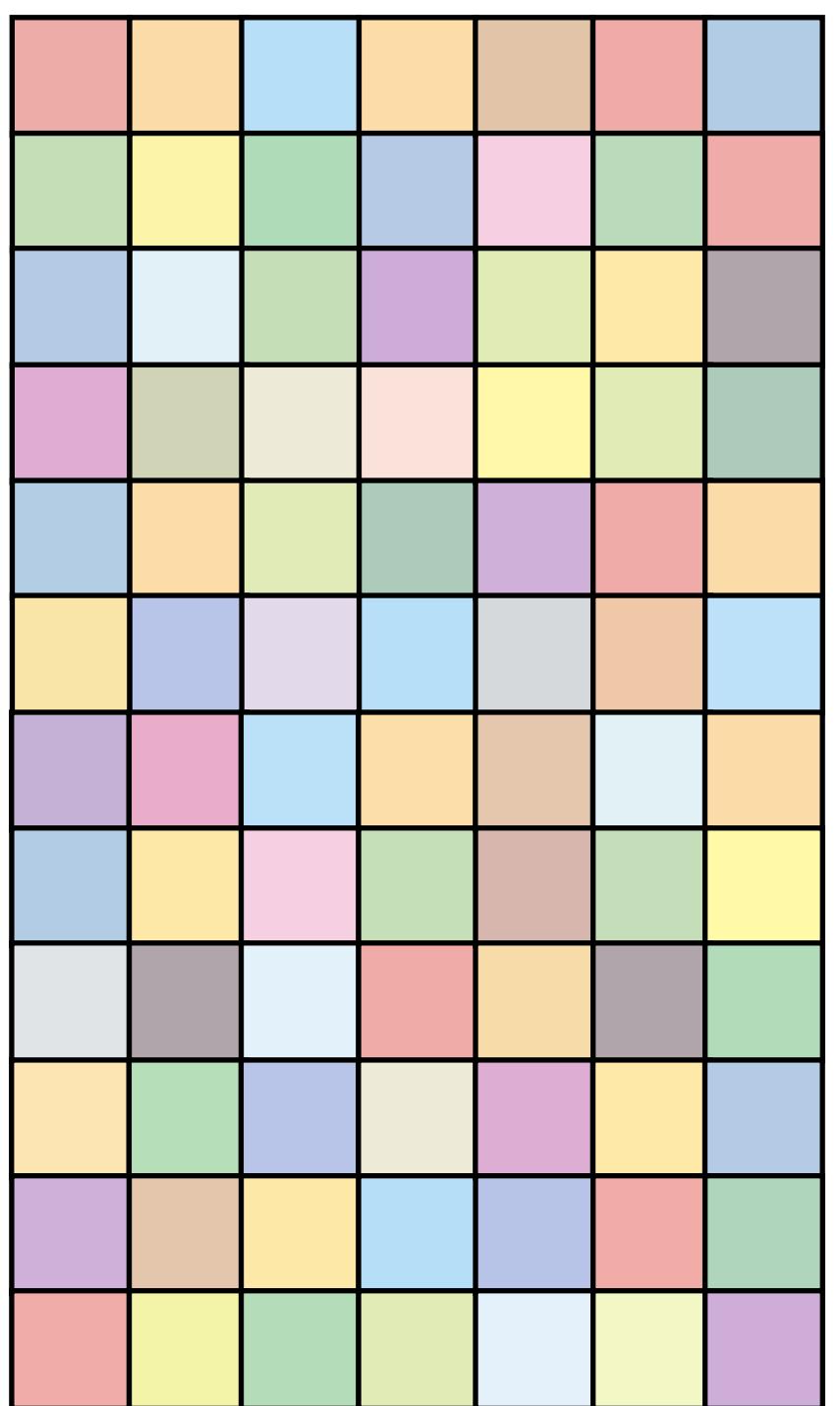
Compression of neural networks

- Deep neural networks achieve state-of-the-art performance in a variety of domains
- Model quality scales with model and dataset size
- State-of-the-art models usually incorporate **tens of millions of parameters**
- But **resources** (memory, processing time) **may be limited**

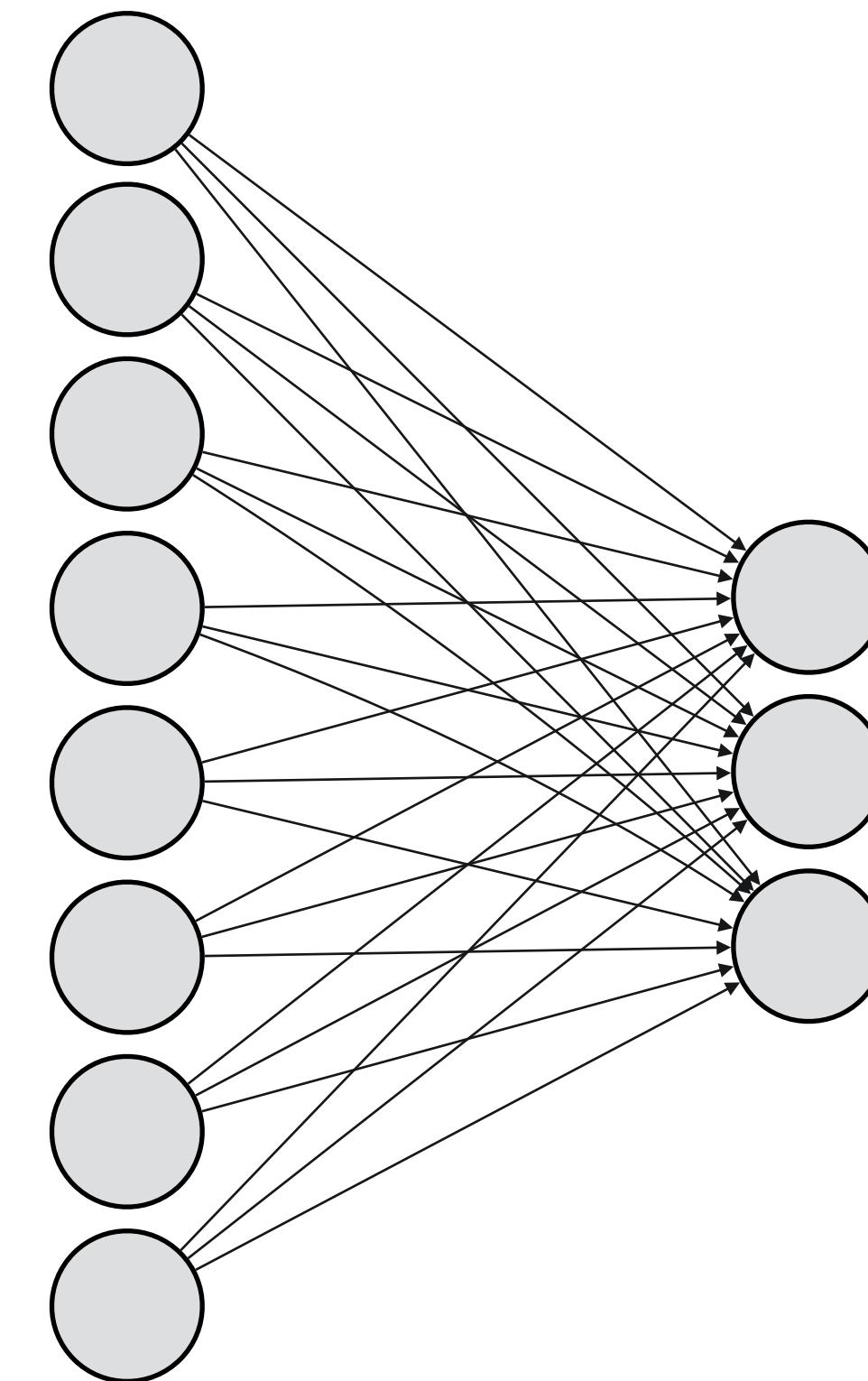
Methods for neural network compression:

- Quantization
- Matrix factorizations
- Sparsification

Neural network

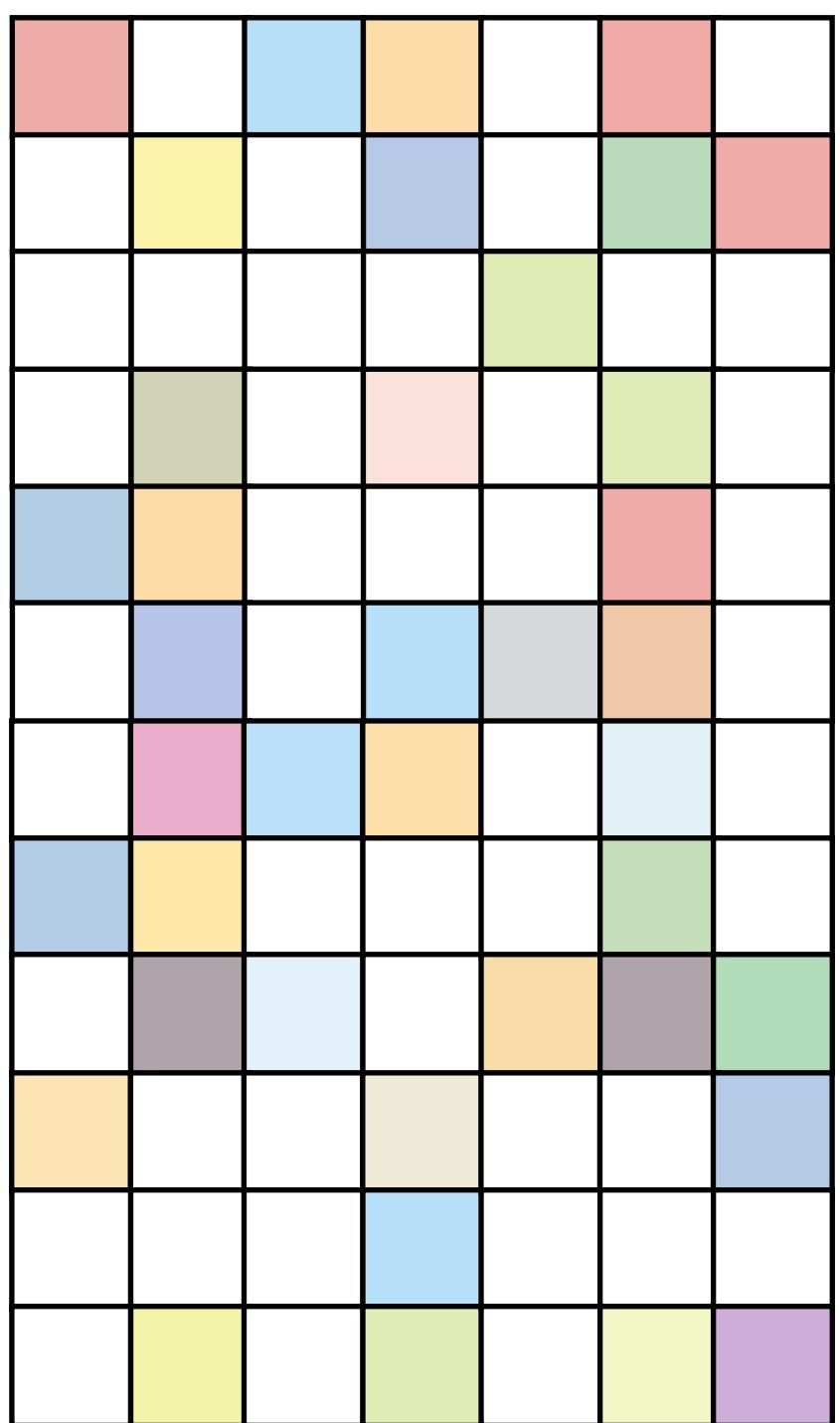


Weight matrix W



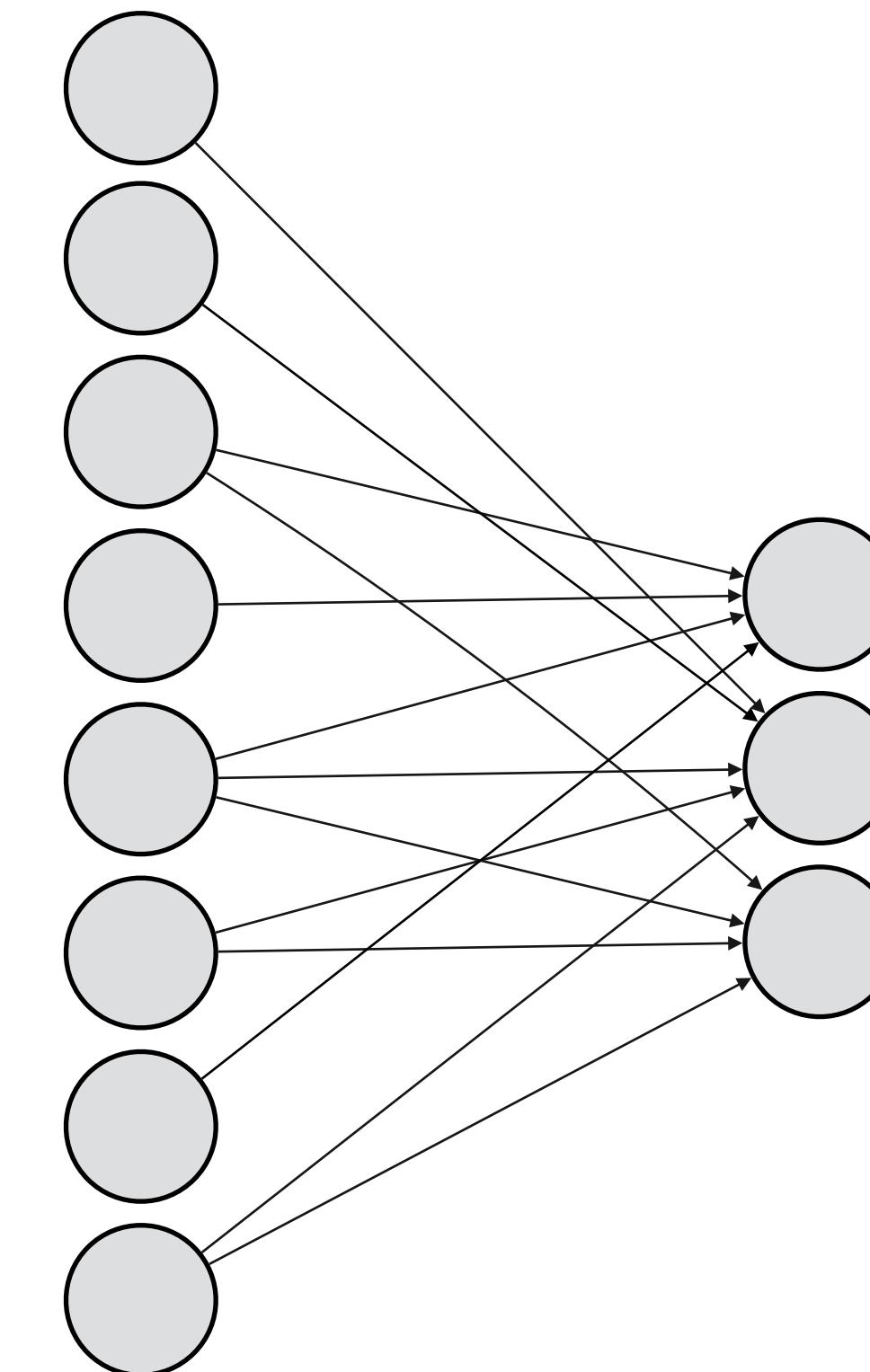
Computational graph

Sparse neural network



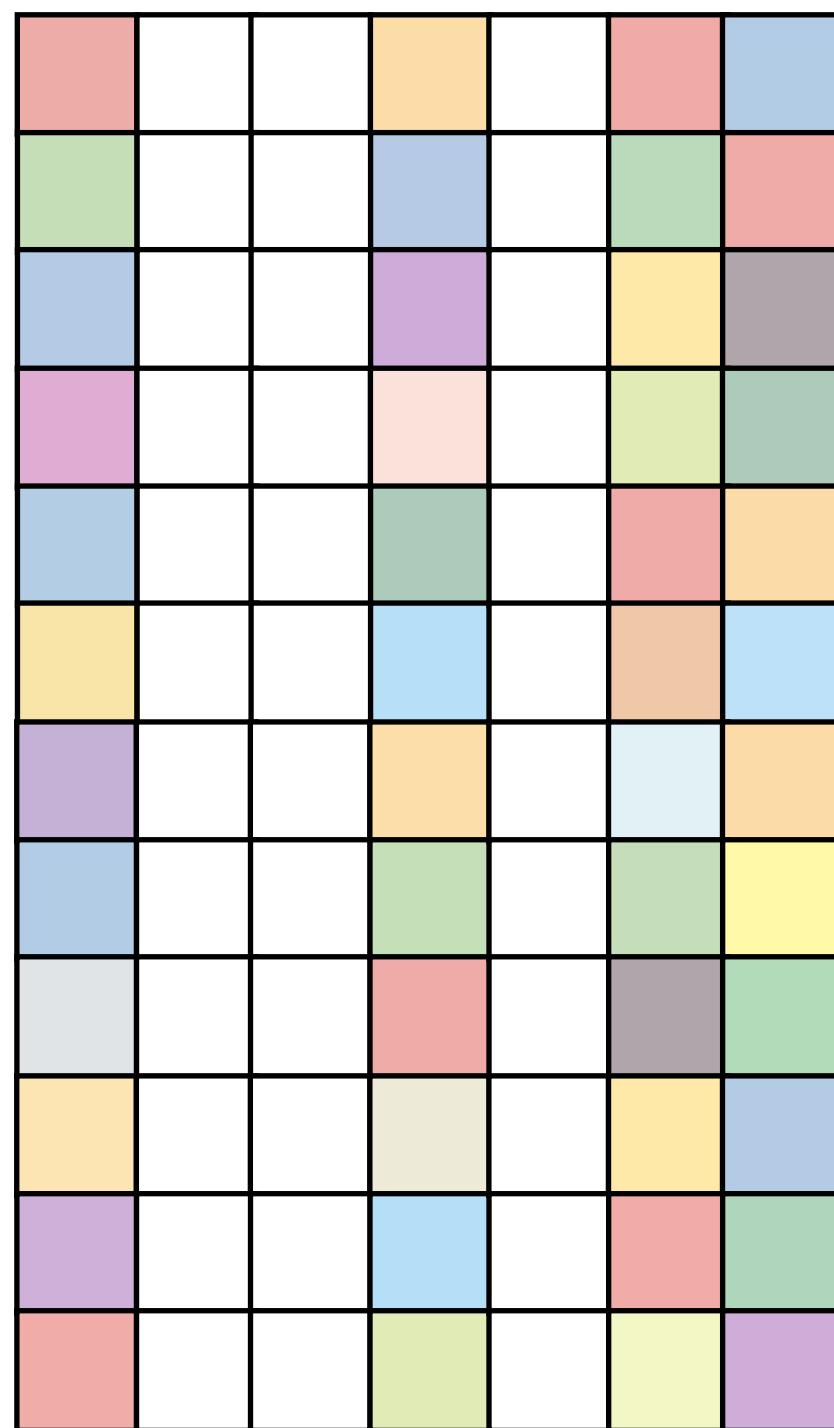
Weight matrix W

A lot of weights
set to zero



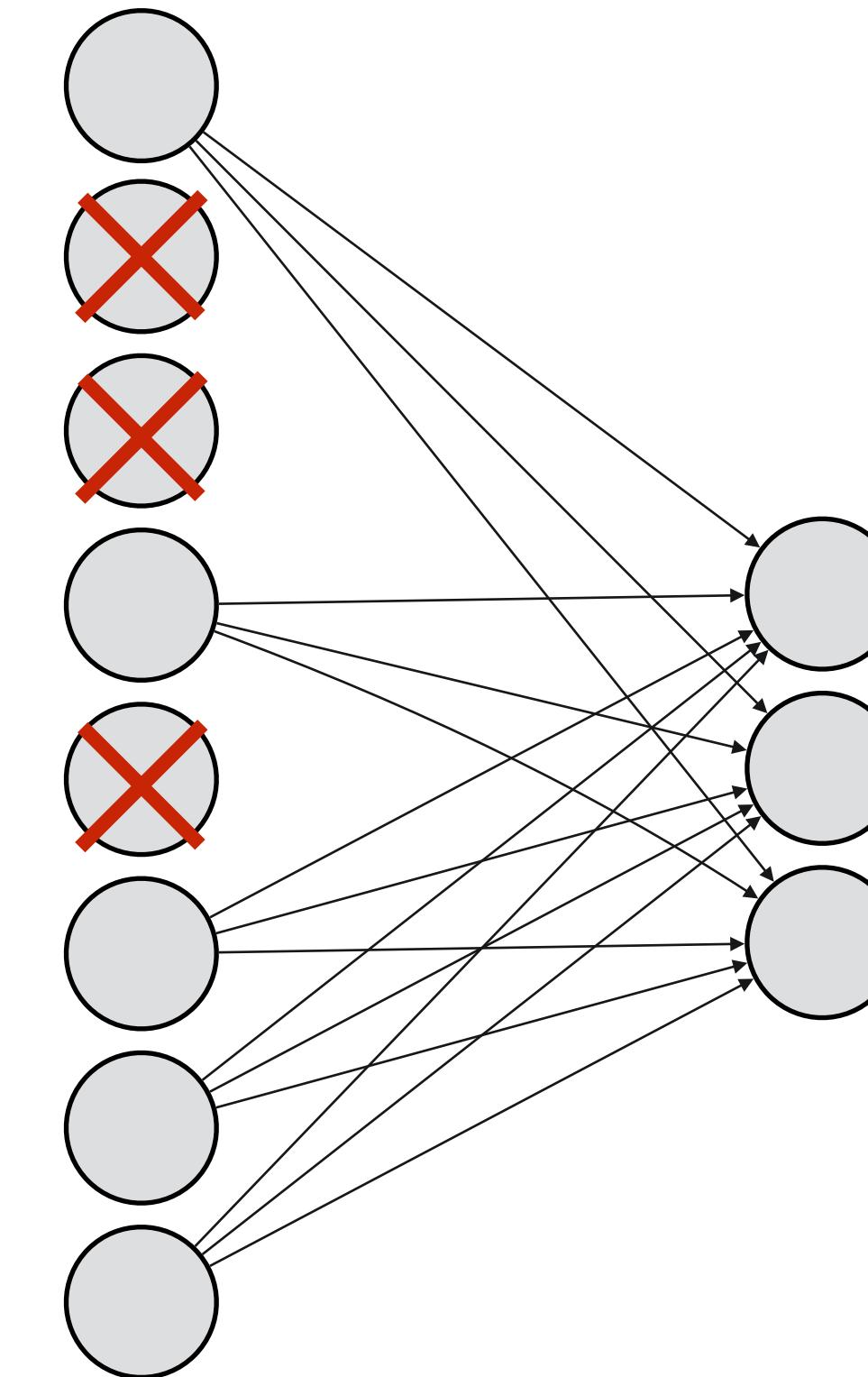
Computational graph

Structured sparsity



Weight matrix W

No outgoing edges
⇒ remove neuron



Computational graph

Sparsification of neural networks

Benefits:

- sparse matrices \Rightarrow less memory consumption
- structured sparsification \Rightarrow faster testing stage (prediction)
- Regularization
- A bit more interpretable model

Drawbacks:

- Sometimes leads to small quality drop

Applications:

- Mobile devices, smartphones
- Online services (where fast reply is needed)

Pruning

- Popular approach to NN sparsification
- Example training algorithm:

- Optimize L_1 -regularized loss (induces sparsity)
 - At the beginning of each epoch, set to 0 all weights satisfying $|w| < T$.

Pruning

- Popular approach to NN sparsification
- Example training algorithm:
 - Optimize L_1 -regularized loss
 - At the beginning of each epoch, set to 0 all weights satisfying $|w| < T$.
- Variants:
 - Different pruning schedule (increase T gradually / constant T)
 - Prune based on threshold T / prune percent of weights
 - Propagate gradients through zero weights or not
 - Prune from scratch / train - prune - retrain
 - ...

huge number
of
hyperparameters



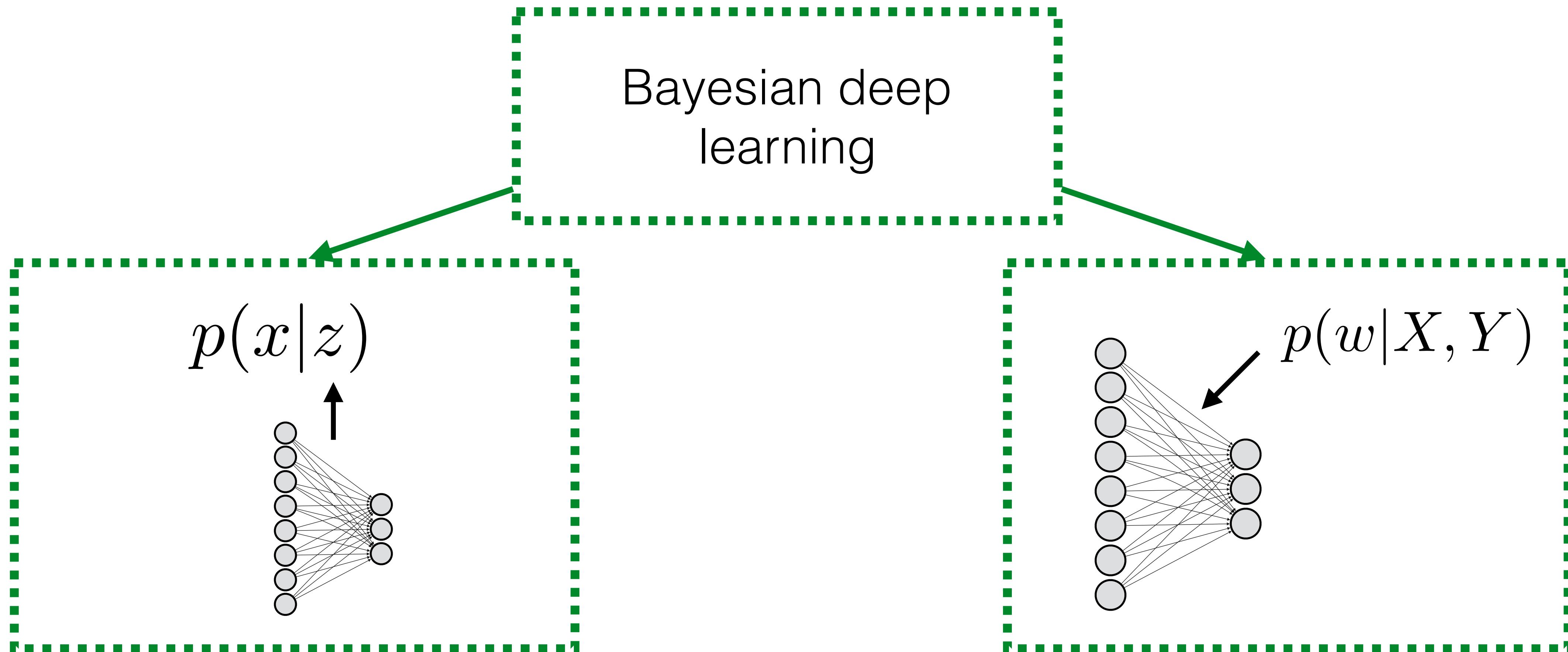
Two popular frameworks for sparsification

Pruning		Bayesian sparsification
A lot of method hyperparameters		(Almost) no method hyperparameters
Need to choose training schedule		Need to choose training schedule
non-Bayesian		Bayesian!

Agenda

- Sparsification: what and why
- Bayesian neural networks
- Sparse variational dropout
- Practical assignment: implementation of SparseVD
- Model enhancements

Bayesian deep learning



Deep generative models:
VAE, NF ...

Bayesian neural networks
(discriminative)

Regularization by noise

- Traditional (1943+) regularization: add some penalty for model complexity
 - L_2 , L_1 - regularization, max norm constraint

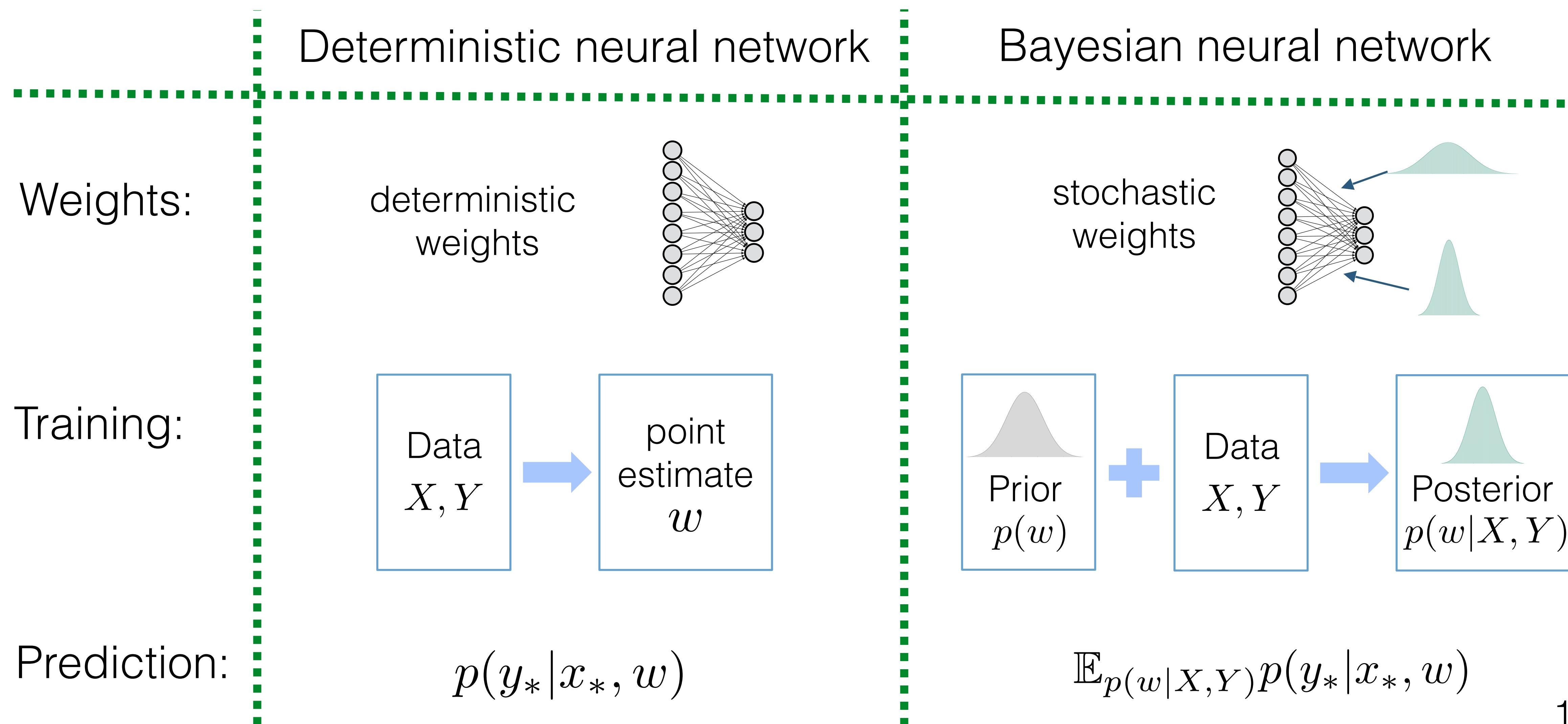
$$\text{Objective} = \text{DataLoss}(X, Y, w) + \text{Regularizer}(w)$$

- More recent (1990+) approaches: regularization by noise
 - Data augmentation, dropout, gradient noise

$$\text{Objective} = \mathbb{E}_{p(\Omega)} \text{DataLoss}(X, Y, w, \Omega)$$

Bayesian framework provides a principled approach to training with noise!

Bayesian neural networks



Why go Bayesian?

A principled framework with many useful applications

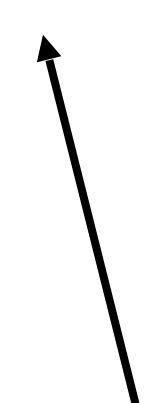
- Regularization
- Ensembling
- Uncertainty estimation
- On-line / continual learning
- Automatic hyperparameter choice
- Different prior \Rightarrow different properties of the network

Ensembling

A Bayesian neural network is **an infinite ensemble** of neural networks

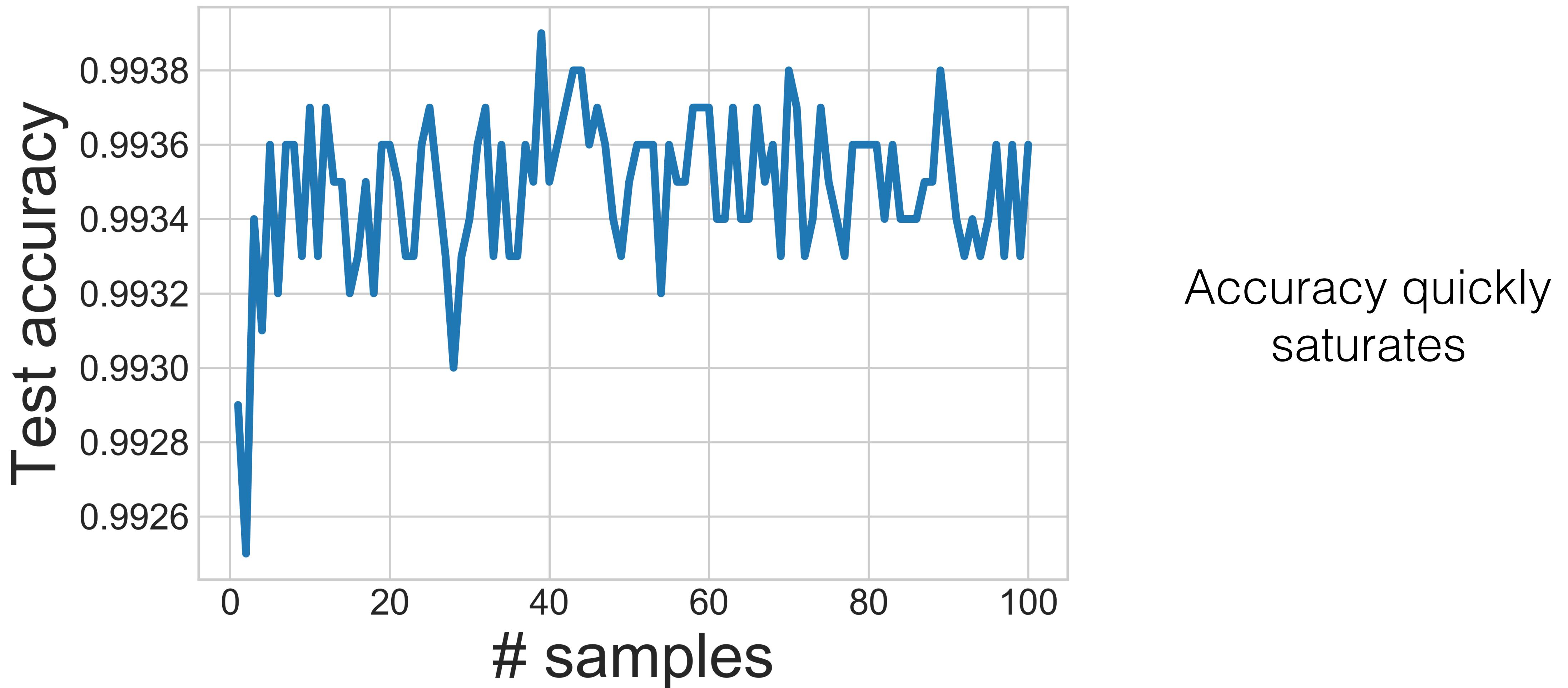
$$p(y_*|x_*, X, Y) = \int p(y_*|x_*, w)p(w|X, Y)dw \approx \frac{1}{K} \sum_{i=1}^K p(y_*|x_*, w^k)$$

$w^k \sim p(w|X, Y)$

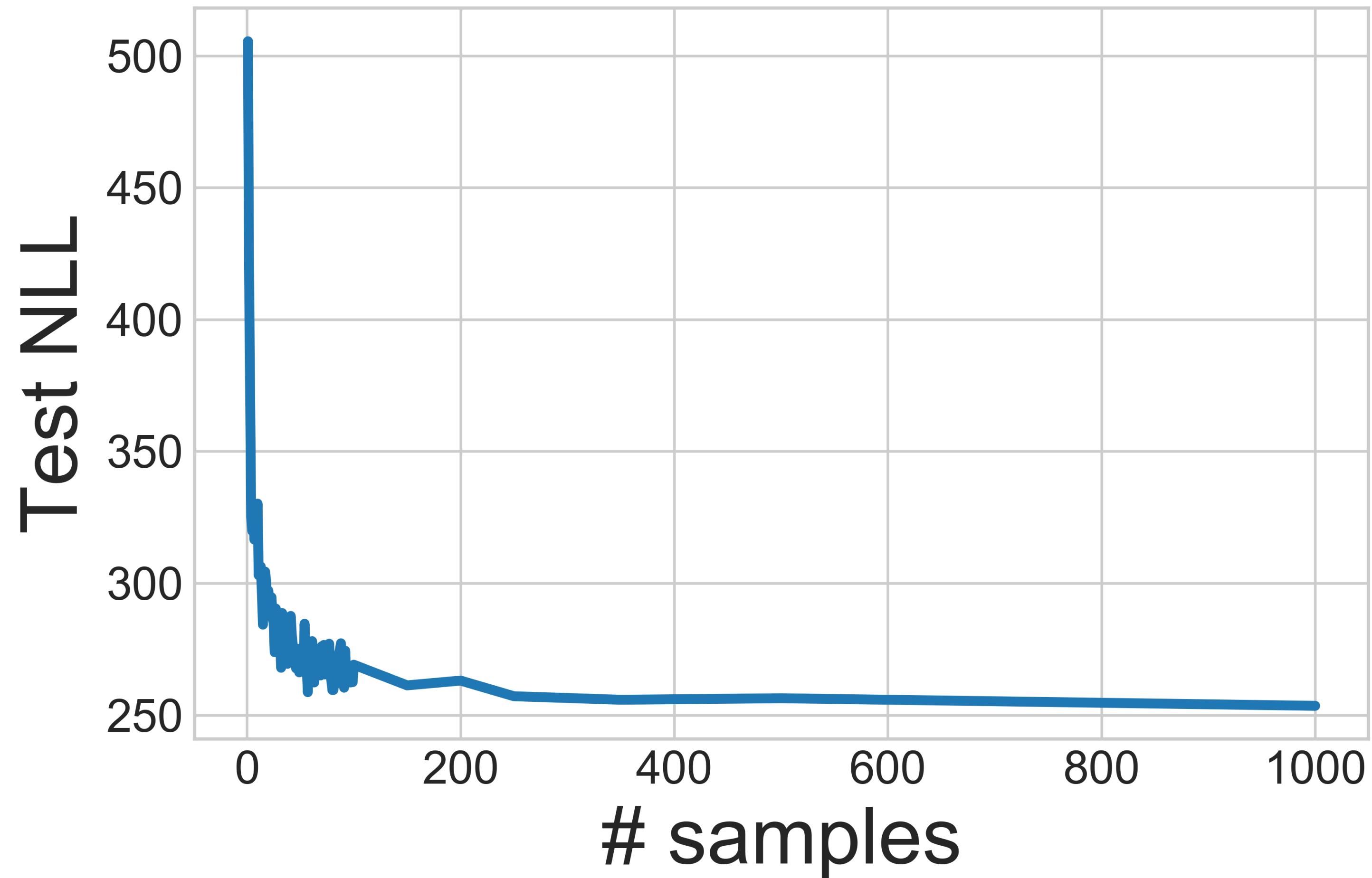


Average SoftMax outputs
across several samples

Ensembling



Ensembling

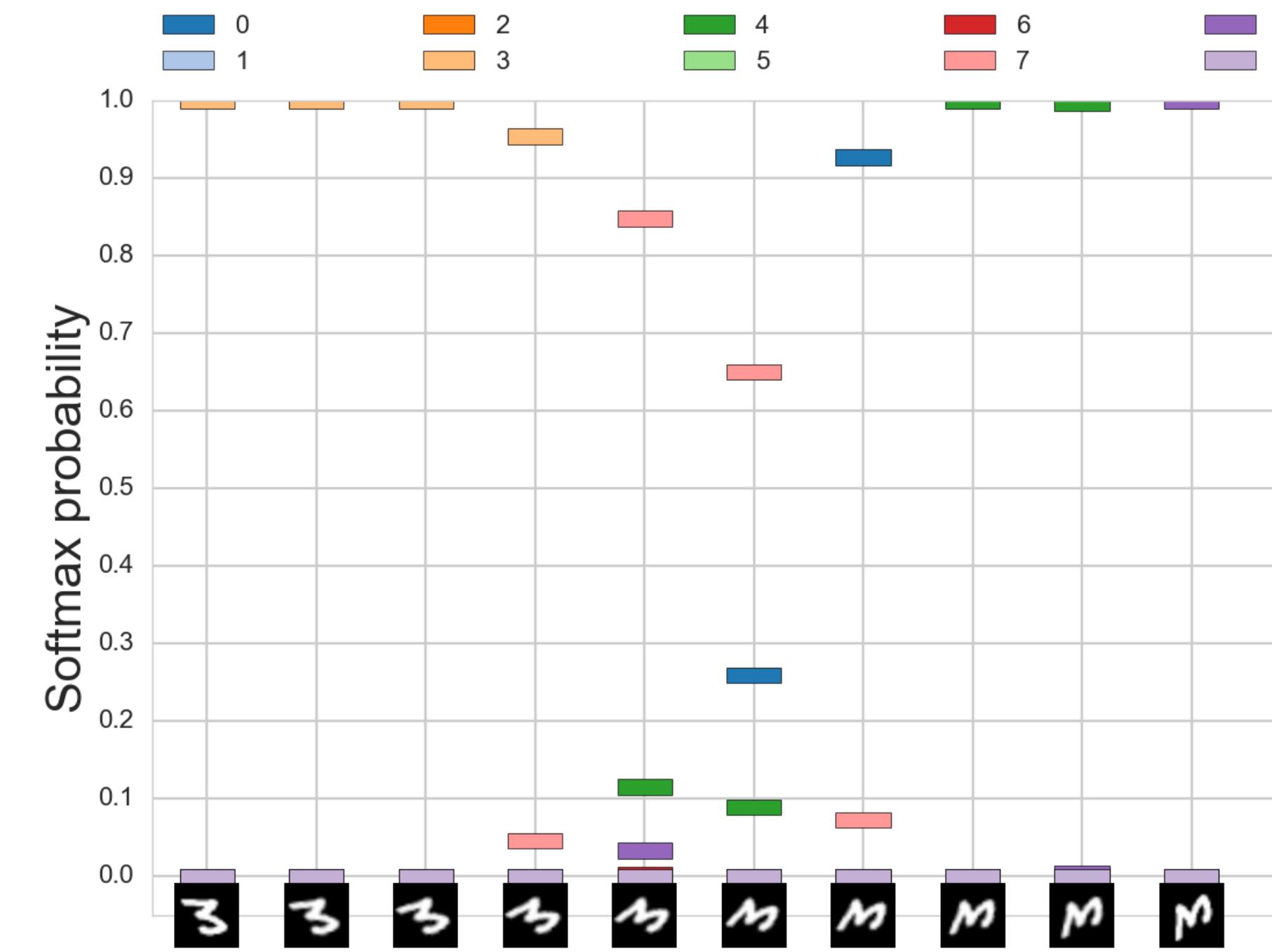


But the negative log—likelihood keeps improving!
This is a measure of “uncertainty”

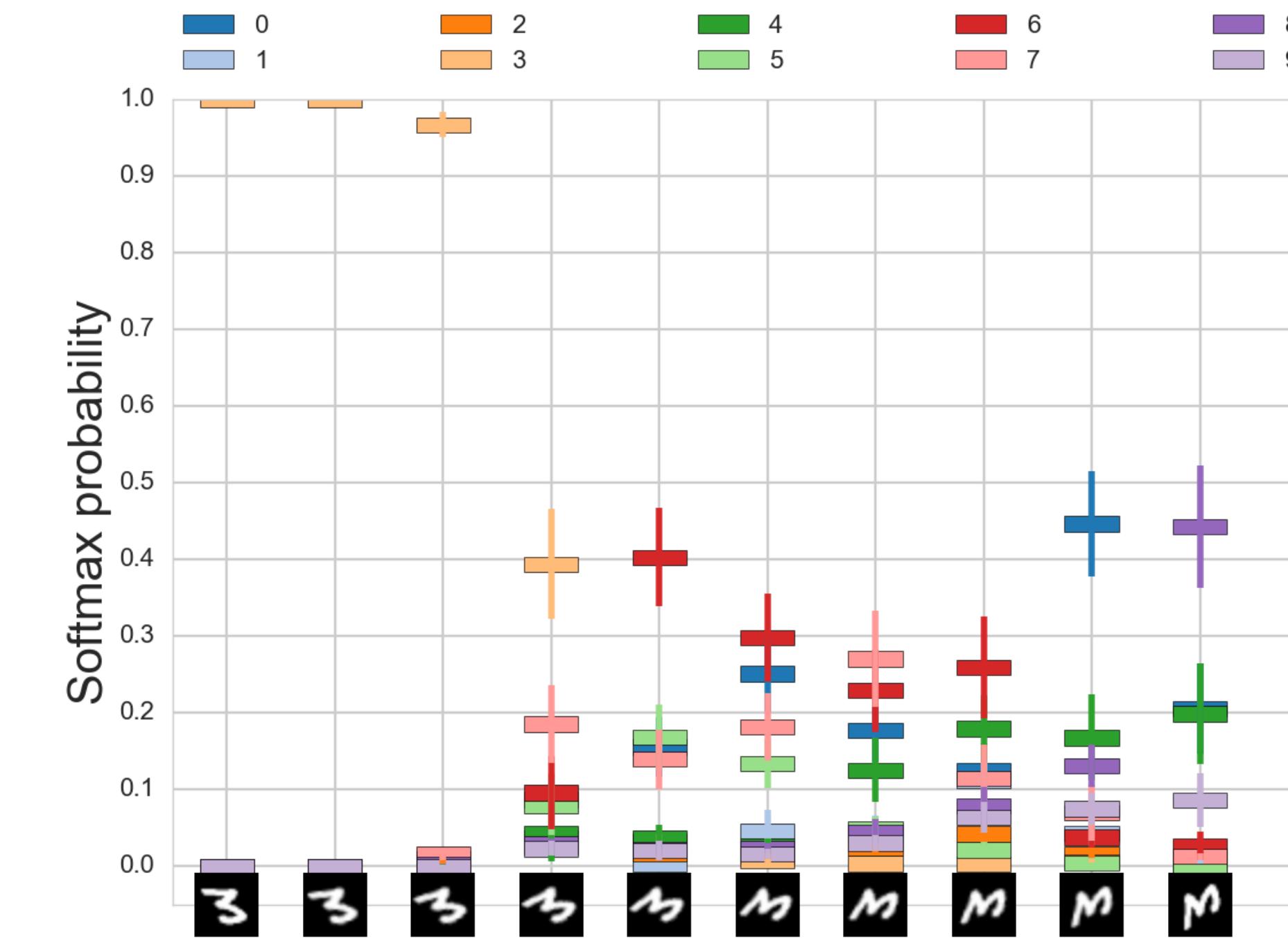
Uncertainty estimation

Deterministic NNs: a point estimate of the output, overconfident

Bayesian framework: a distribution over the outputs



(a) LeNet with weight decay



(b) LeNet with multiplicative formalizing flows

Model selection

- Empirical Bayes (maximum evidence)
 - Optimize w. r. t. prior parameters ⇒ Choose hyperparameters
- Discuss later

Online / incremental learning

Assume that dataset arrives gradually in independent parts:

$$\mathcal{D} = \mathcal{D}_1 \cup \mathcal{D}_2 \cup \mathcal{D}_3 \cup \dots \cup \mathcal{D}_M$$

We can first train on the first dataset as usual:

$$p(w|\mathcal{D}_1) = \frac{p(\mathcal{D}_1|w)p(w)}{\int p(\mathcal{D}_1|\tilde{w})p(\tilde{w})d\tilde{w}}$$

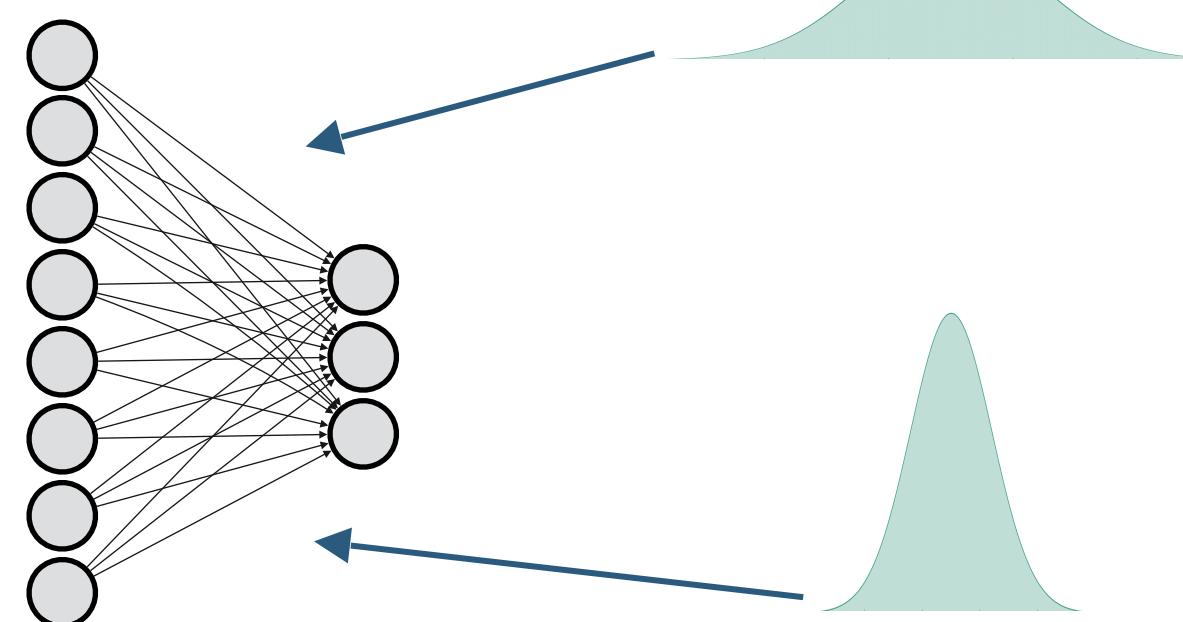
And then use the obtained posterior as the prior for the next step!

$$p(w|\mathcal{D}_2, \mathcal{D}_1) = \frac{p(\mathcal{D}_2|w)\cancel{p(\mathcal{D}_1|w)p(w)}}{\int p(\mathcal{D}_2|\tilde{w})\cancel{p(\mathcal{D}_1|\tilde{w})p(\tilde{w})}d\tilde{w}}$$

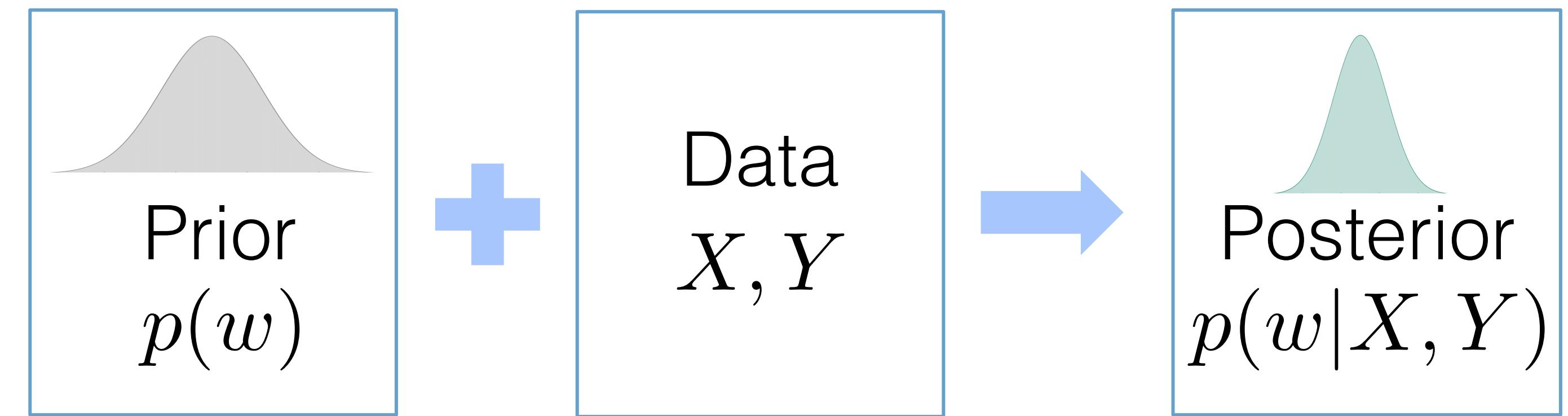
Using these sequential updates, we can find $p(w|\mathcal{D})$.

Training Bayesian neural networks

Stochastic weights:



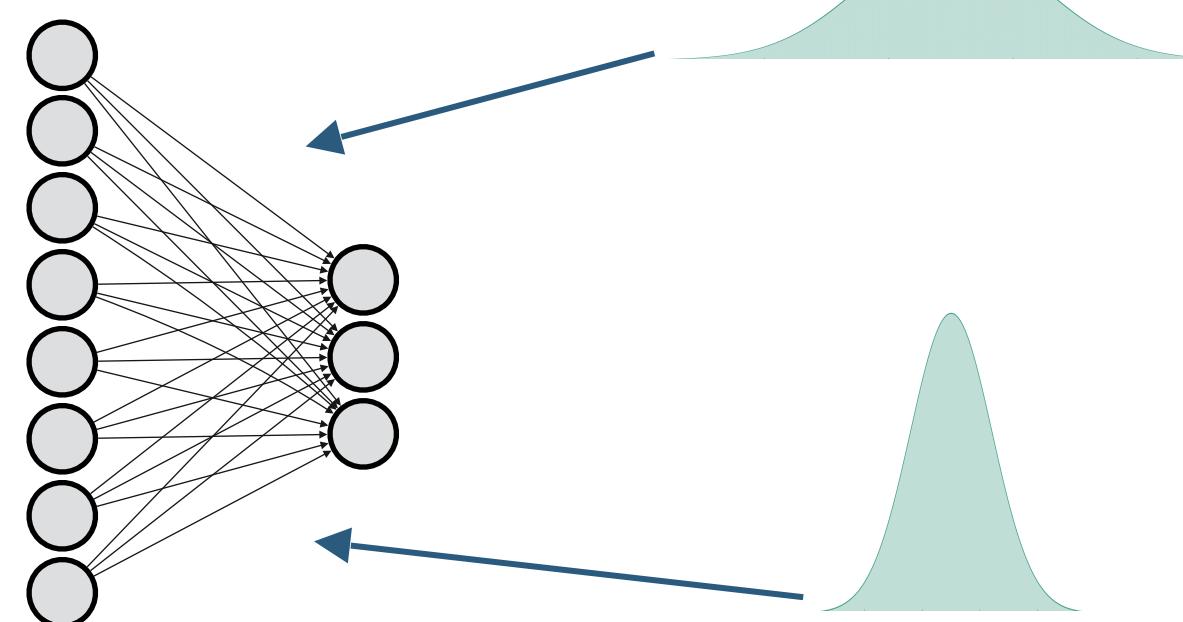
Bayesian Inference:



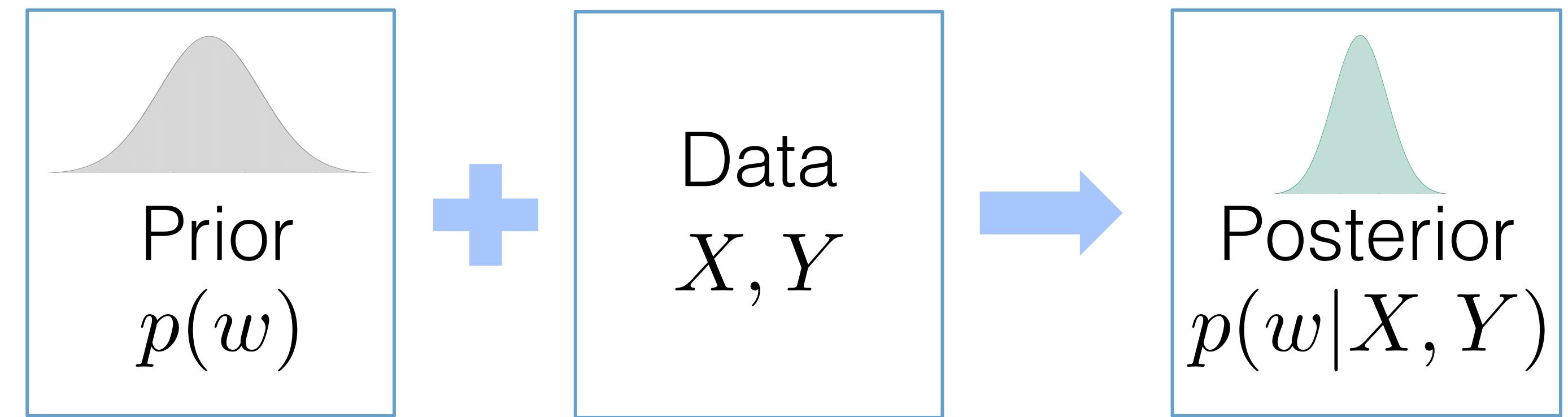
$$p(w|X, Y) = \frac{p(Y|X, w)p(w)}{\int p(Y|X, \tilde{w})p(\tilde{w})d\tilde{w}}$$

Training Bayesian neural networks

Stochastic weights:



Bayesian Inference:

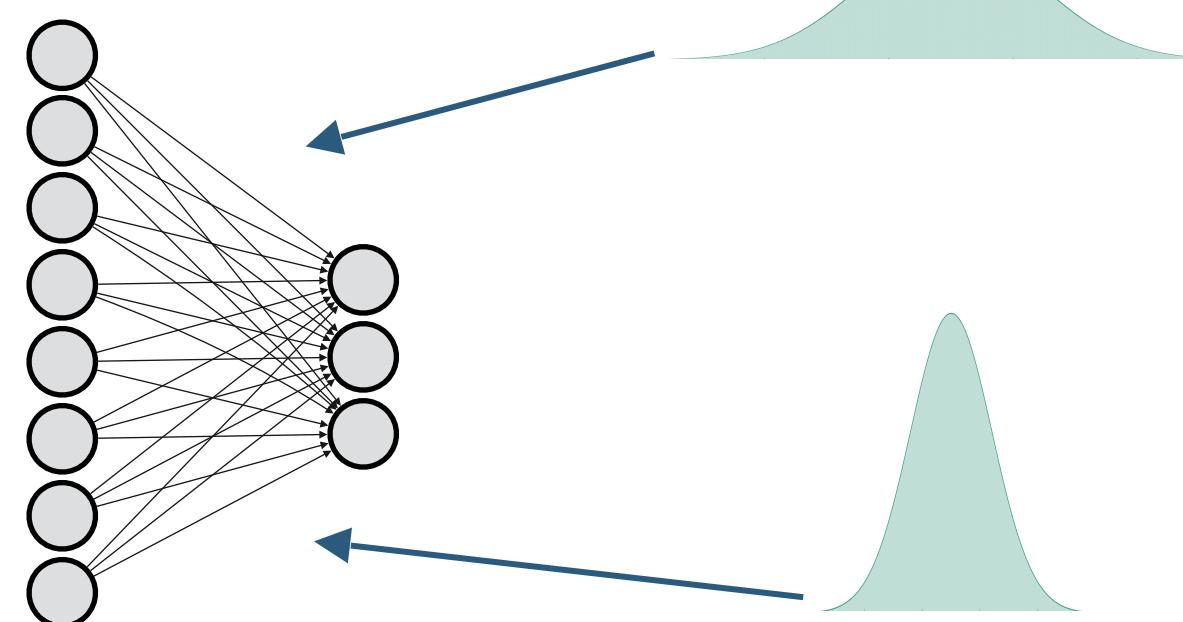


Posterior is intractable in neural networks → approximate it with $q(w|\lambda)$:

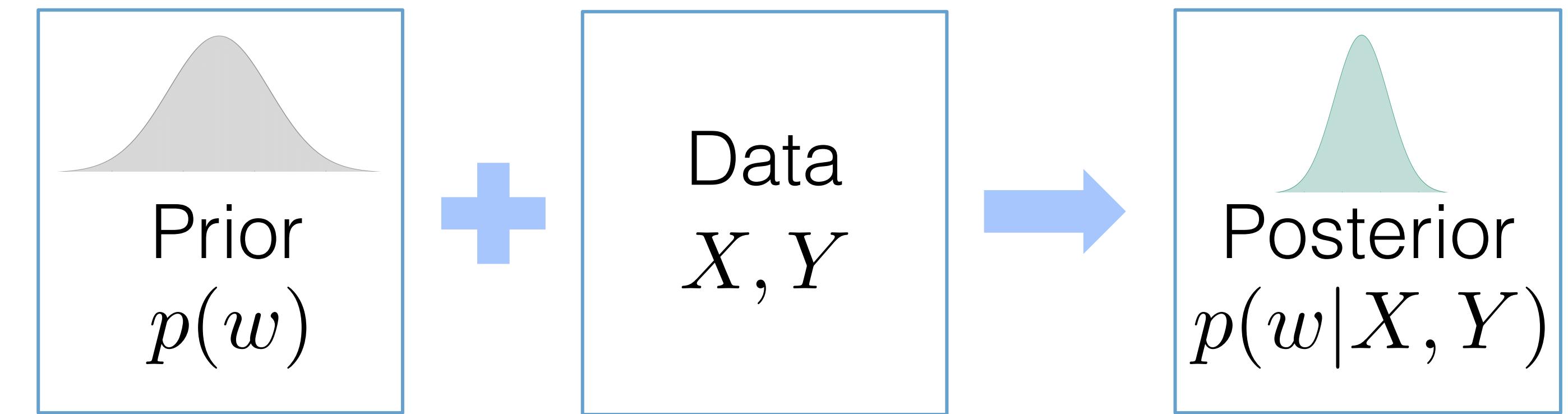
$$KL(q(w|\lambda)||p(w|X, y)) \rightarrow \min_{\lambda}$$

Training Bayesian neural networks

Stochastic weights:



Bayesian Inference:



Equivalently, we can optimize ELBO to find approximate posterior $q(w|\lambda)$:

$$\sum_{i=1}^N \underbrace{\mathbb{E}_{q(w|\lambda)} \log p(y^i|x^i, w)}_{\text{Data term}} - \underbrace{KL(q(w|\lambda)||p(w))}_{\text{Regularizer}} \rightarrow \max_{\lambda}$$

Training Bayesian neural networks

$$\sum_{i=1}^N \mathbb{E}_{q(w|\lambda)} \log p(y^i|x^i, w) - KL(q(w|\lambda)||p(w)) \rightarrow \max_{\lambda}$$

Variational autoencoders are also trained using variational inference.

Compared to them:

- KL-term is global
- Extremely high-dimensional approximate posterior

From general framework to particular method

$$\sum_{i=1}^N \mathbb{E}_{q(w|\lambda)} \log p(y^i|x^i, w) - KL(q(w|\lambda)||p(w)) \rightarrow \max_{\lambda}$$

Model specification:

- Choose particular prior

Training:

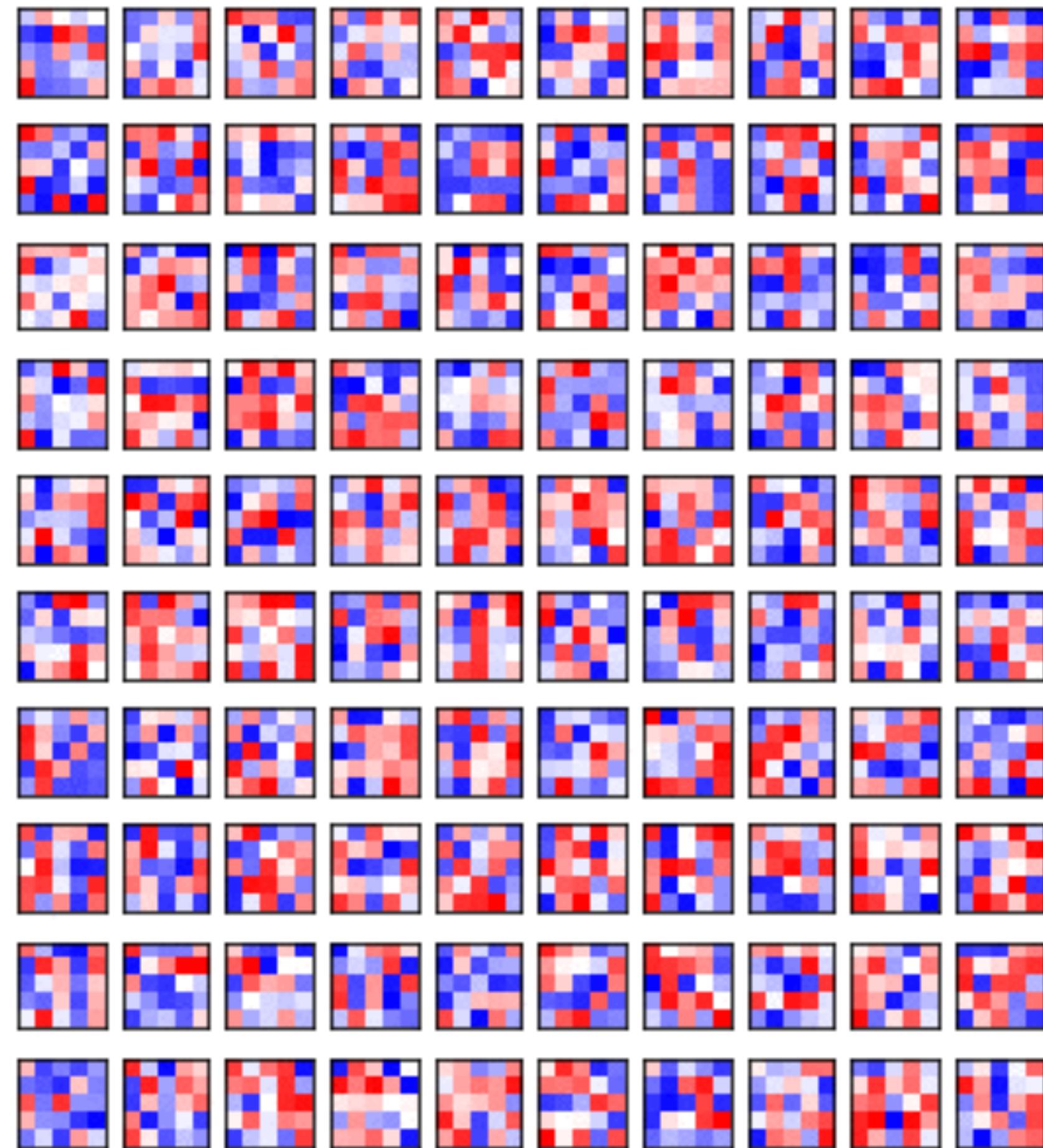
- Choose particular family for approximate posterior
- How to compute the KL-divergence?
- How to estimate the expectation?

Agenda

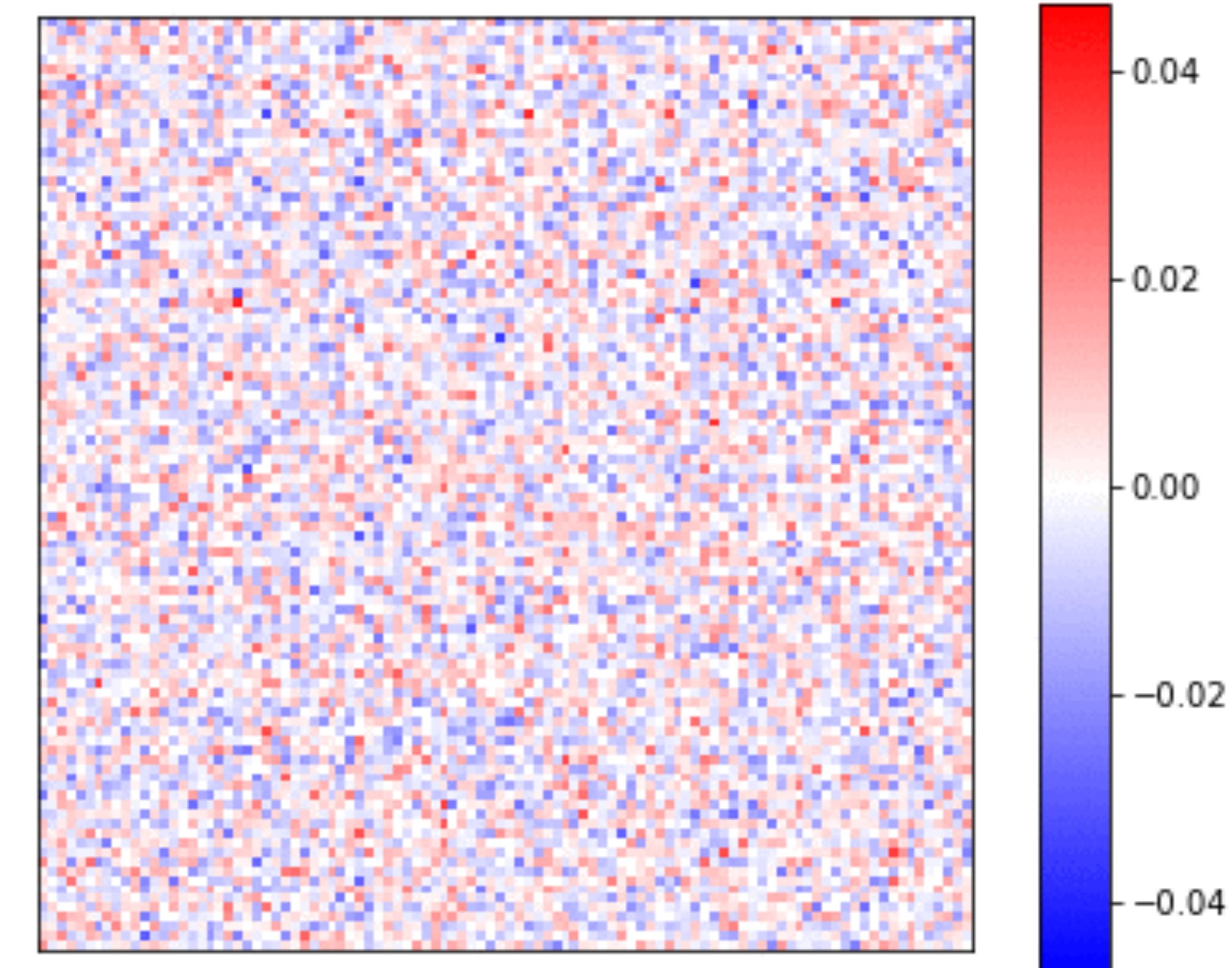
- Sparsification: what and why
- Bayesian neural networks
- Sparse variational dropout
- Practical assignment: implementation of SparseVD
- Model enhancements

Sparse variational dropout: visualization

Epoch: 0 Compression ratio: 1x Accuracy: 8.4



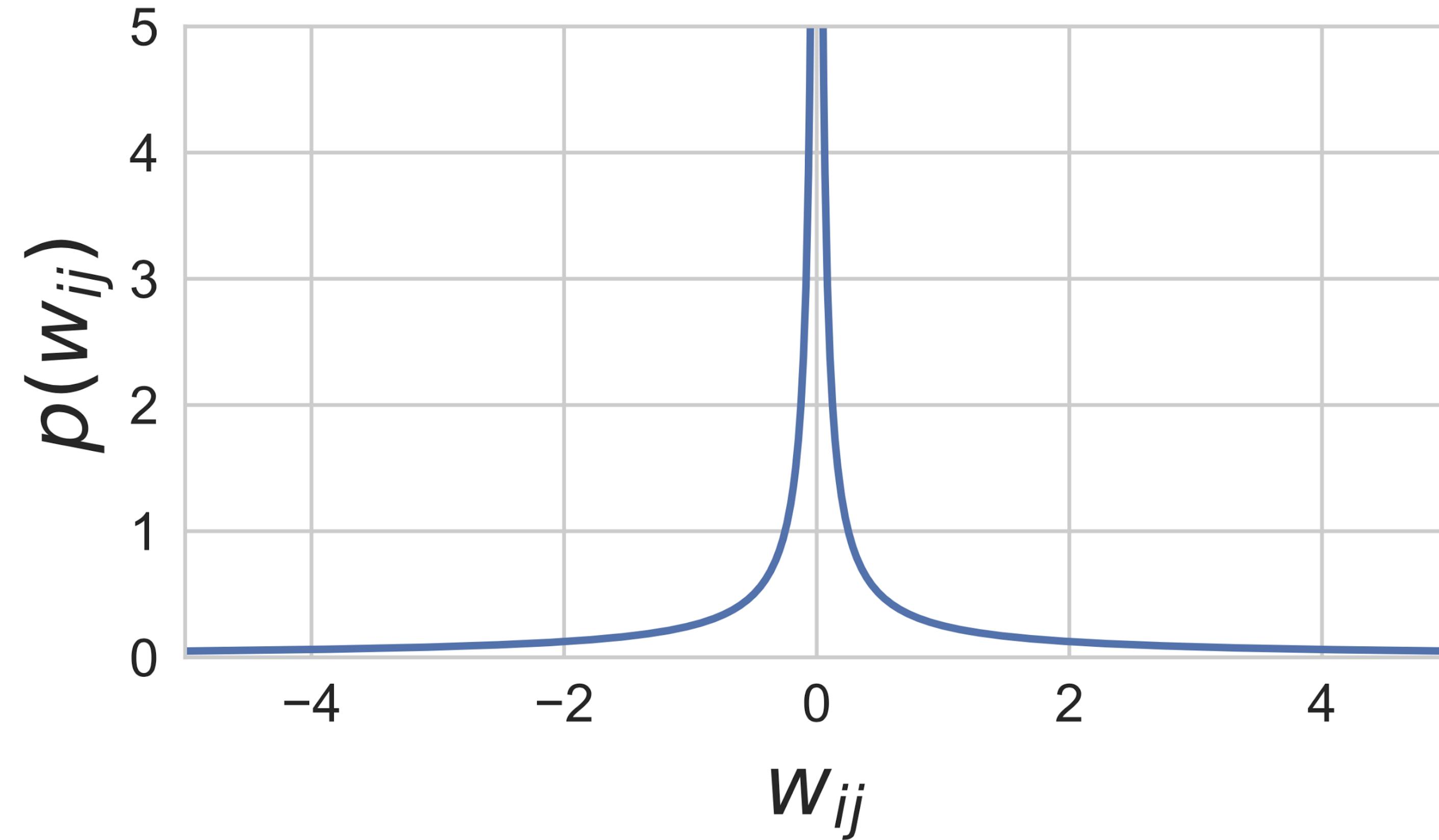
Epoch: 0 Compression ratio: 1x Accuracy: 8.4



LeNet-5: convolutional layer

LeNet-5: fully-connected layer
(100 x 100 patch)

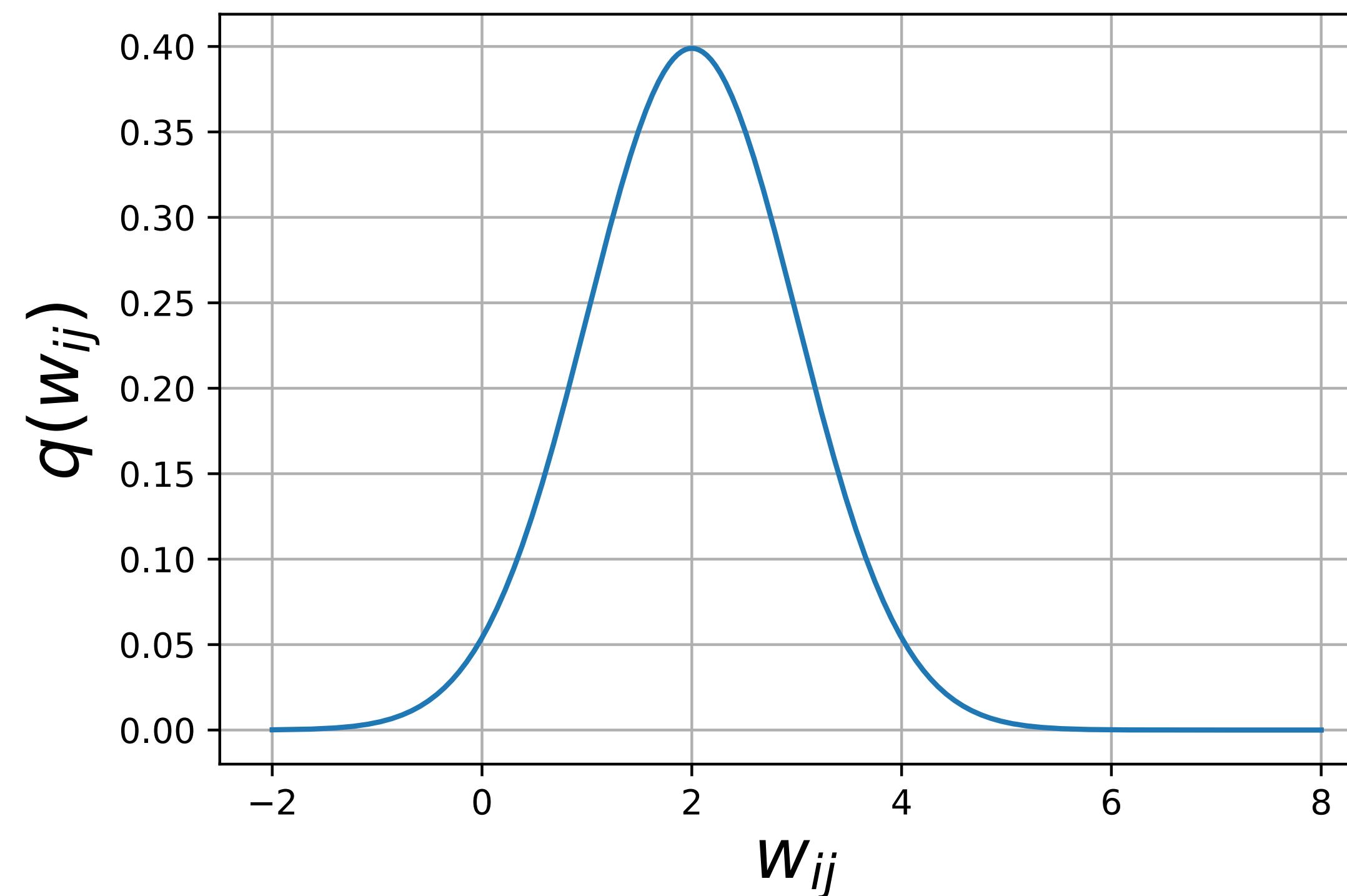
Log-uniform prior distribution



$$p(w_{ij}) \propto \frac{1}{|w_{ij}|}$$

Favors removing noisy weights
(few slides later)

Normal approximate posterior distribution (fully factorized)



$$q(w_{ij} | \mu_{ij}, \sigma_{ij}) = \mathcal{N}(\mu_{ij}, \sigma_{ij}^2)$$

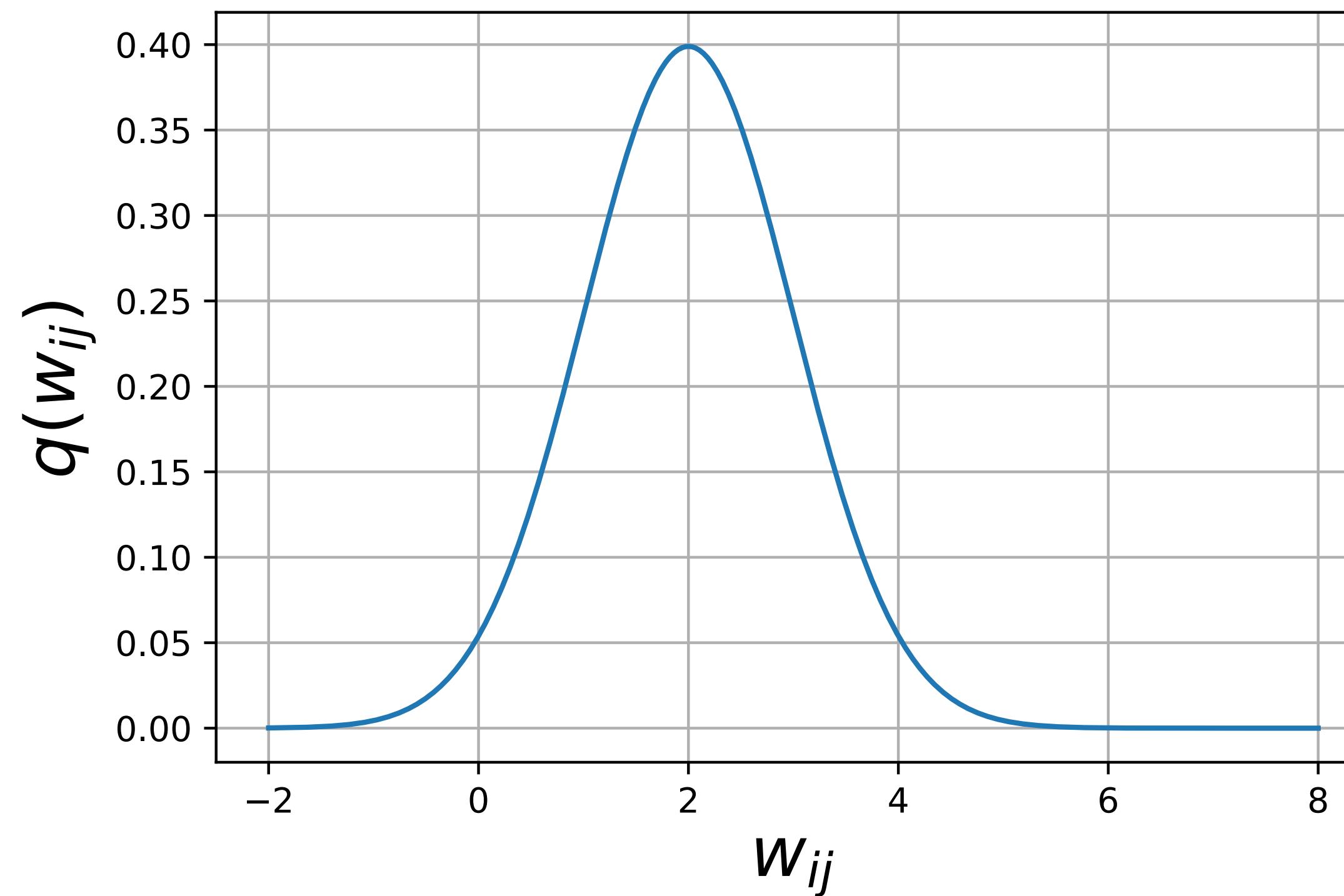
Reparametrization trick:

$$\hat{w}_{ij} = \mu_{ij} + \hat{\epsilon}_{ij} \sigma_{ij}, \quad \hat{\epsilon}_{ij} \sim \mathcal{N}(0, 1)$$

weight
mean

zero-centered noise
with learnable variance

Another possible parametrization (actually bad for sparsification)



Same approx. posterior distribution
but different parametrization:

$$q(w_{ij} | \mu_{ij}, \alpha_{ij}) = \mathcal{N}(\mu_{ij}, \alpha_{ij}\mu_{ij}^2)$$

Reparametrization trick:

$$\hat{w}_{ij} = \mu_{ij}(1 + \sqrt{\alpha_{ij}}\hat{\epsilon}_{ij}), \quad \hat{\epsilon}_{ij} \sim \mathcal{N}(0, 1)$$

one-centered
multiplicative noise
with learnable variance

Compare two parametrizations

$$\frac{\partial \mathcal{L}}{\partial \mu_{ij}} = \frac{\partial \mathcal{L}}{\partial \hat{w}_{ij}} \boxed{\frac{\partial \hat{w}_{ij}}{\partial \mu_{ij}}}$$

With multiplicative noise:

$$\hat{w}_{ij} = \mu_{ij}(1 + \sqrt{\alpha_{ij}}\hat{\epsilon}_{ij})$$

$$\frac{\partial \mathcal{L}}{\partial \mu_{ij}} = \frac{\partial \mathcal{L}}{\partial \hat{w}_{ij}} \boxed{(1 + \sqrt{\alpha_{ij}}\hat{\epsilon}_{ij})}$$

Very noisy!

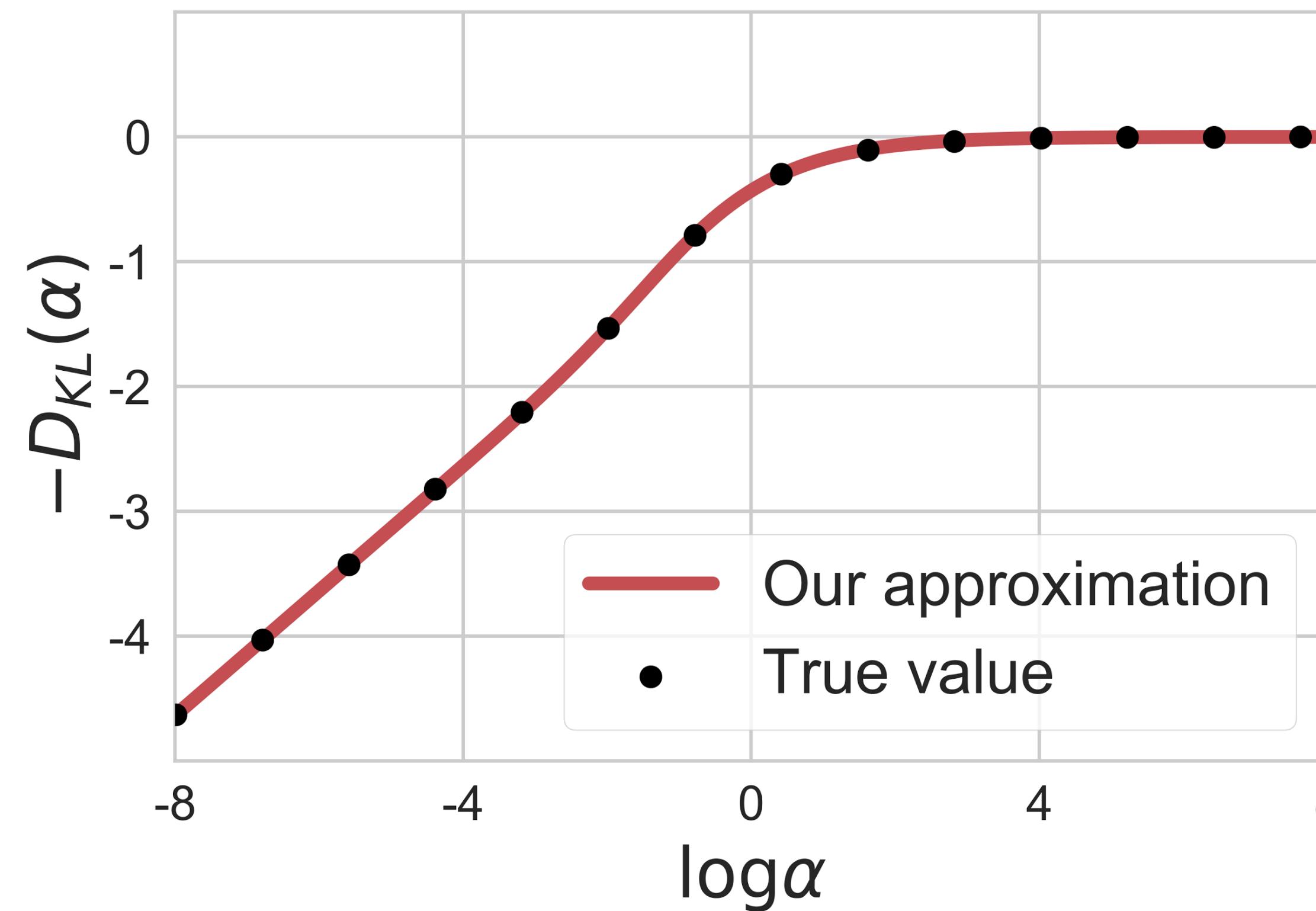
With additive noise:

$$\hat{w}_{ij} = \mu_{ij} + \hat{\epsilon}_{ij}\sigma_{ij}$$

$$\frac{\partial \mathcal{L}}{\partial \mu_{ij}} = \frac{\partial \mathcal{L}}{\partial \hat{w}_{ij}} \cdot \boxed{1}$$

No noise!

Approximating KL-divergence (fully factorized)



$$\begin{aligned}-KL(q(w_{ij}|\mu_{ij}, \sigma_{ij}) \| p(w_{ij})) &\approx \\ &\approx k_1 \sigma(k_2 + k_3 \log \alpha_{ij}) - 0.5 \log(1 + \alpha_{ij}^{-1}) + C\end{aligned}$$
$$k_1 = 0.63576 \quad k_2 = 1.87320 \quad k_3 = 1.48695$$

$$\alpha_{ij} = \frac{\sigma_{ij}^2}{\mu_{ij}^2}$$

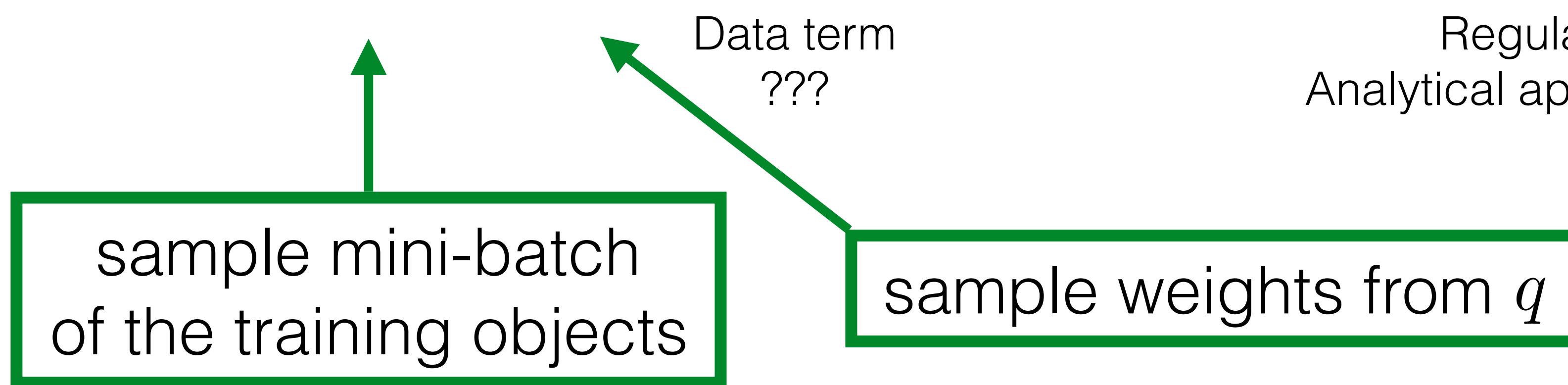
- KL depends only on α_{ij}
- Favors large $\alpha_{ij} \Rightarrow$ removing noisy weights

Doubly stochastic variational inference

$$\sum_{i=1}^N \underbrace{\mathbb{E}_{q(w|\mu,\sigma)} \log p(y^i|x^i, w)}_{\text{Data term} \atop ???} - \underbrace{KL(q(w|\mu,\sigma)||p(w))}_{\text{Regularizer} \atop \text{Analytical approximation}} \rightarrow \max_{\mu, \log \sigma}$$

Doubly stochastic variational inference

$$\sum_{i=1}^N \underbrace{\mathbb{E}_{q(w|\mu,\sigma)} \log p(y^i|x^i, w)}_{\text{Data term } ???} - \underbrace{KL(q(w|\mu,\sigma)||p(w))}_{\text{Regularizer}} \rightarrow \max_{\mu, \log \sigma}$$



Estimating expectation in ELBO

$$\sum_{i=1}^N \underbrace{\mathbb{E}_{q(w|\mu,\sigma)} \log p(y^i|x^i, w)}_{\substack{\text{Data term} \\ \text{Sample weights from } q}} - \underbrace{KL(q(w|\mu,\sigma)||p(w))}_{\substack{\text{Regularizer} \\ \text{Analytical approximation}}} \rightarrow \max_{\mu, \log \sigma}$$

Estimating expectation in ELBO

$$\sum_{i=1}^N \underbrace{\mathbb{E}_{q(w|\mu,\sigma)} \log p(y^i|x^i, w)}_{\substack{\text{Data term} \\ \text{Sample weights from } q}} - \underbrace{KL(q(w|\mu,\sigma)||p(w))}_{\substack{\text{Regularizer} \\ \text{Analytical approximation}}} \rightarrow \max_{\mu, \log \sigma}$$

1. Reparametrization trick \Rightarrow unbiased estimate of the gradients

$$\hat{w}_{ij} = \mu_{ij} + \hat{\epsilon}_{ij}\sigma_{ij}, \quad \hat{\epsilon}_{ij} \sim \mathcal{N}(0, 1)$$

Estimating expectation in ELBO

$$\sum_{i=1}^N \underbrace{\mathbb{E}_{q(w|\mu,\sigma)} \log p(y^i|x^i, w)}_{\substack{\text{Data term} \\ \text{Sample weights from } q}} - \underbrace{KL(q(w|\mu,\sigma)||p(w))}_{\substack{\text{Regularizer} \\ \text{Analytical approximation}}} \rightarrow \max_{\mu, \log \sigma}$$

1. Reparametrization trick \Rightarrow unbiased estimate of the gradients

$$\hat{w}_{ij} = \mu_{ij} + \hat{\epsilon}_{ij}\sigma_{ij}, \quad \hat{\epsilon}_{ij} \sim \mathcal{N}(0, 1)$$

- Too expensive! ($|w| \times$ mini-batch size)

Estimating expectation in ELBO

$$\sum_{i=1}^N \underbrace{\mathbb{E}_{q(w|\mu,\sigma)} \log p(y^i|x^i, w)}_{\substack{\text{Data term} \\ \text{Sample weights from } q}} - \underbrace{KL(q(w|\mu,\sigma)||p(w))}_{\substack{\text{Regularizer} \\ \text{Analytical approximation}}} \rightarrow \max_{\mu, \log \sigma}$$

1. Reparametrization trick \Rightarrow unbiased estimate of the gradients

$$\hat{w}_{ij} = \mu_{ij} + \hat{\epsilon}_{ij}\sigma_{ij}, \quad \hat{\epsilon}_{ij} \sim \mathcal{N}(0, 1)$$

- Too expensive! ($|w| \times$ mini-batch size)
- One sample per mini-batch? High variance of stochastic gradients & correlated predictions

Estimating expectation in ELBO

$$\sum_{i=1}^N \underbrace{\mathbb{E}_{q(w|\mu,\sigma)} \log p(y^i|x^i, w)}_{\substack{\text{Data term} \\ \text{Sample weights from } q}} - \underbrace{KL(q(w|\mu,\sigma)||p(w))}_{\substack{\text{Regularizer} \\ \text{Analytical approximation}}} \rightarrow \max_{\mu, \log \sigma}$$

1. Reparametrization trick \Rightarrow unbiased estimate of the gradients

$$\hat{w}_{ij} = \mu_{ij} + \hat{\epsilon}_{ij}\sigma_{ij}, \quad \hat{\epsilon}_{ij} \sim \mathcal{N}(0, 1)$$

2. Local reparametrization trick: sample preactivations instead of weights
 \Rightarrow reduced variance of the gradients & uncorrelated predictions

$$w_{ij} \sim \mathcal{N}(\mu_{ij}, \sigma_{ij}^2)$$

$$B = AW$$

$$\mathbb{E}B = A\mu \quad \text{Var}B = A^2\sigma^2$$

$$B \sim \mathcal{N}(A\mu, A^2\sigma^2)$$

$$B = A\mu + \sqrt{A^2\sigma^2} \odot \epsilon$$

blue means
element-wise

The local reparametrization trick

Consider a linear layer with weight matrix W , input A and output B :

$$w_{ij} \sim \mathcal{N}(\mu_{ij}, \sigma_{ij}^2)$$

$$B = AW$$

Predictions have **high**
correlation because there is
one weight sample **per mini-batch**

The local reparametrization trick

Consider a linear layer with weight matrix W , input A and output B :

$$w_{ij} \sim \mathcal{N}(\mu_{ij}, \sigma_{ij}^2)$$

$$B = AW$$

Predictions have **high** correlation because there is one weight sample **per mini-batch**

Less correlation \Rightarrow lower variance of the stochastic gradients **for a mini-batch**

$$\mathbb{E}B = A\mu \quad \text{Var}B = A^2\sigma^2$$

$$B \sim \mathcal{N}(A\mu, A^2\sigma^2)$$

$$B = A\mu + \sqrt{A^2\sigma^2} \odot \epsilon$$

blue means element-wise

Predictions have **zero** correlation because there is one weight sample **per object**

The local reparametrization trick

The LRT also reduces the variance of the stochastic gradient w.r.t $\sigma_{i,j}$ for **an object**.

Sketch of the proof:

$$\text{Var}_{q_\phi, \mathcal{D}} \left[\frac{\partial L_{\mathcal{D}}^{\text{SGVB}}}{\partial \sigma_{i,j}^2} \right] = \text{Var}_{b_{m,j}} \left[\mathbb{E}_{q_\phi, \mathcal{D}} \left[\frac{\partial L_{\mathcal{D}}^{\text{SGVB}}}{\partial \sigma_{i,j}^2} | b_{m,j} \right] \right] + \mathbb{E}_{b_{m,j}} \left[\text{Var}_{q_\phi, \mathcal{D}} \left[\frac{\partial L_{\mathcal{D}}^{\text{SGVB}}}{\partial \sigma_{i,j}^2} | b_{m,j} \right] \right]$$

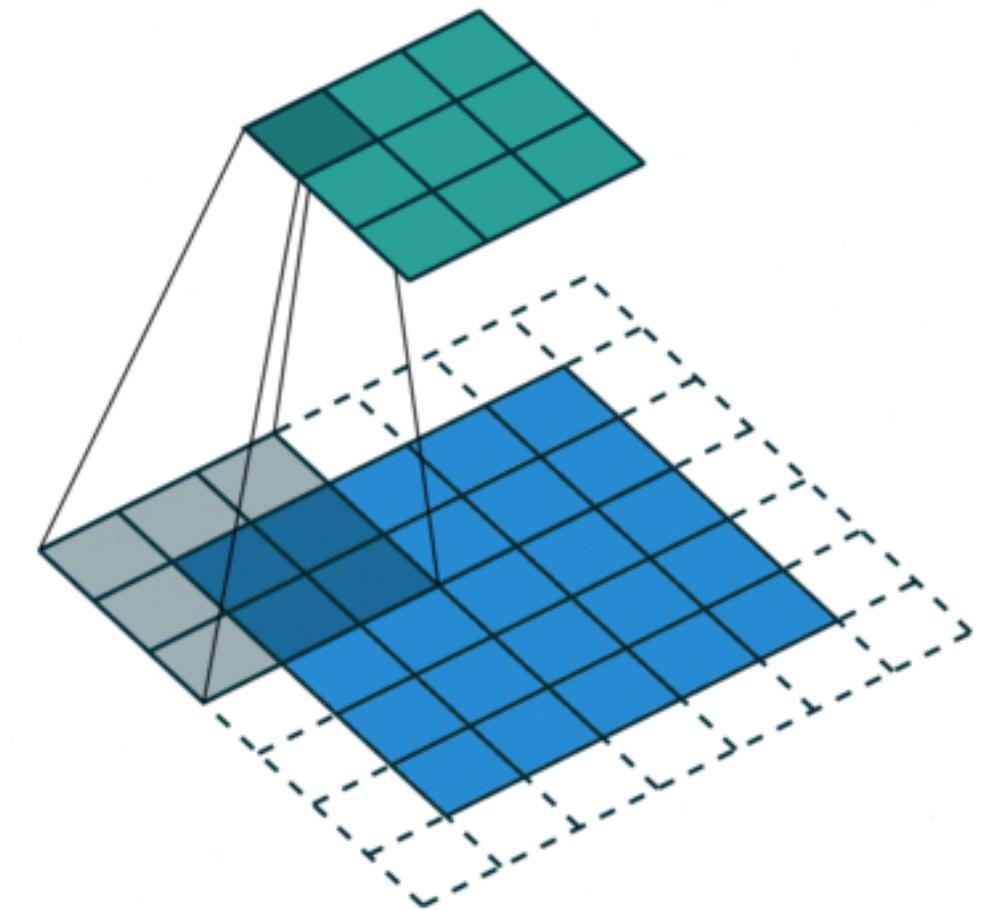
ELBO

same for both parametrizations

for LRT: vanishes because $b_{m,j} = \dots + \dots \cdot \zeta_{m,j}$ and $\zeta_{m,j}$ is uniquely determined given $b_{m,j}$
for RT: doesn't vanish (and positive)

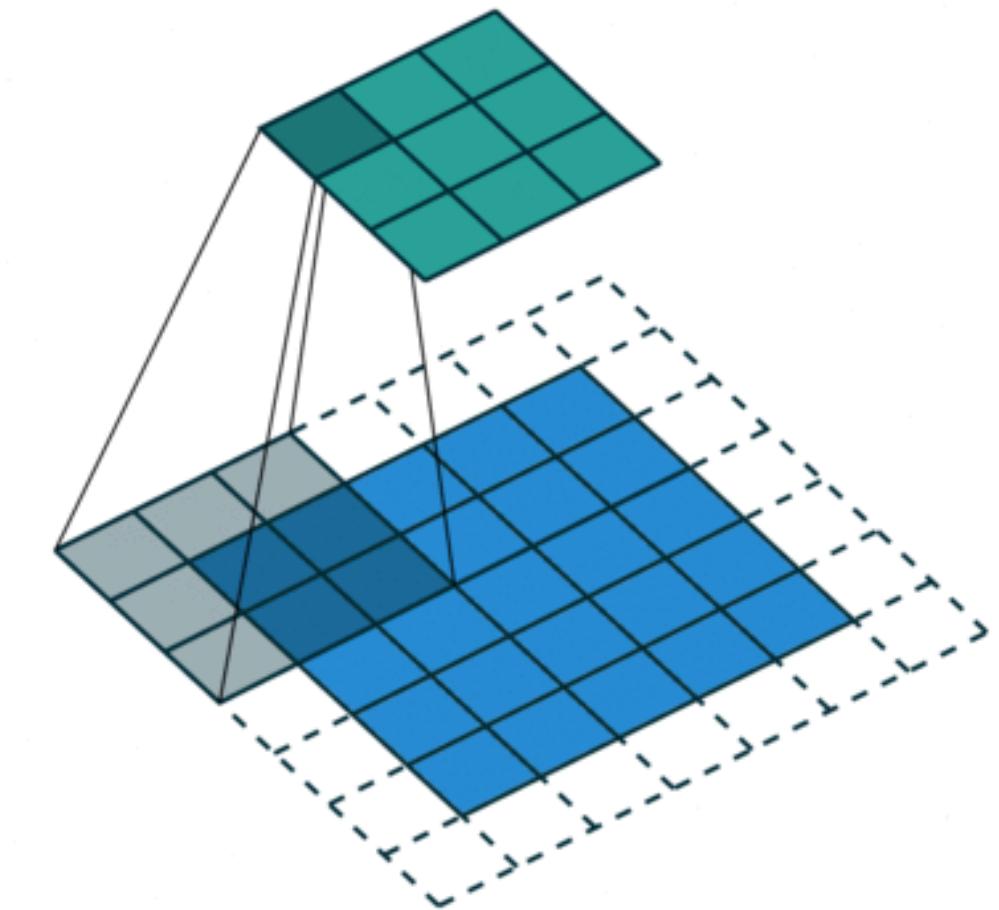
LRT for convolutions

- B no longer factorizes in convolutional layers
 - Same weights sample should be used for different spatial positions



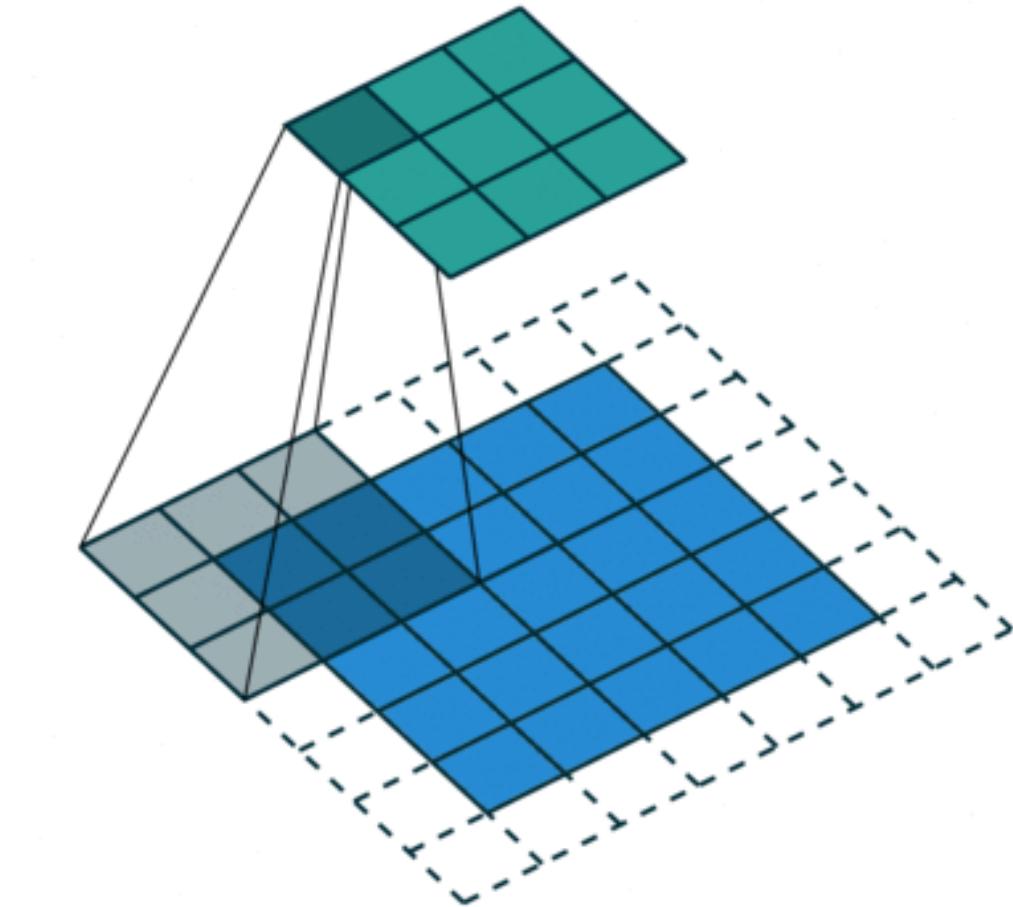
LRT for convolutions

- B no longer factorizes in convolutional layers
 - Same weights sample should be used for different spatial positions
- Exact local reparametrization is too complex
 - Preactivations are normal but correlated
⇒ we need to estimate the full covariance matrix for each preactivation



LRT for convolutions

- B no longer factorizes in convolutional layers
 - Same weights sample should be used for different spatial positions
- Exact local reparametrization is too complex
 - Preactivations are normal but correlated
⇒ we need to estimate the full covariance matrix for each preactivation
- Mean-field approximation performs much better than sampling weights:



$$B = A \star \mu + \sqrt{A^2 \star \sigma^2} \odot \epsilon$$

From general framework to particular method

$$\sum_{i=1}^N \mathbb{E}_{q(w|\mu,\sigma)} \log p(y^i|x^i, w) - KL(q(w|\mu,\sigma)||p(w)) \rightarrow \max_{\mu, \log \sigma}$$

Model specification:

- Choose particular prior  log-uniform distribution

Training:

- Choose particular family for approximate posterior  normal distribution
- How to compute KL-divergence?  analytical approximation
- How to estimate the expectation?  RT & LRT

From general framework to particular method

$$\sum_{i=1}^N \mathbb{E}_{q(w|\mu,\sigma)} \log p(y^i|x^i, w) - KL(q(w|\mu,\sigma)||p(w)) \rightarrow \max_{\mu, \log \sigma}$$

What about biases? batch-norm parameters?

From general framework to particular method

$$\sum_{i=1}^N \mathbb{E}_{q(w|\mu,\sigma)} \log p(y^i|x^i, w) - KL(q(w|\mu,\sigma)||p(w)) \rightarrow \max_{\mu, \log \sigma}$$

What about biases? batch-norm parameters?

Treat them as deterministic parameters and find point estimate:

$$\sum_{i=1}^N \mathbb{E}_{q(w|\mu,\sigma)} \log p(y^i|x^i, w, \mathbf{b}) - KL(q(w|\mu,\sigma)||p(w)) \rightarrow \max_{\mu, \log \sigma, \mathbf{b}}$$

Eventually, it assumes a flat prior and a delta-peak posterior.

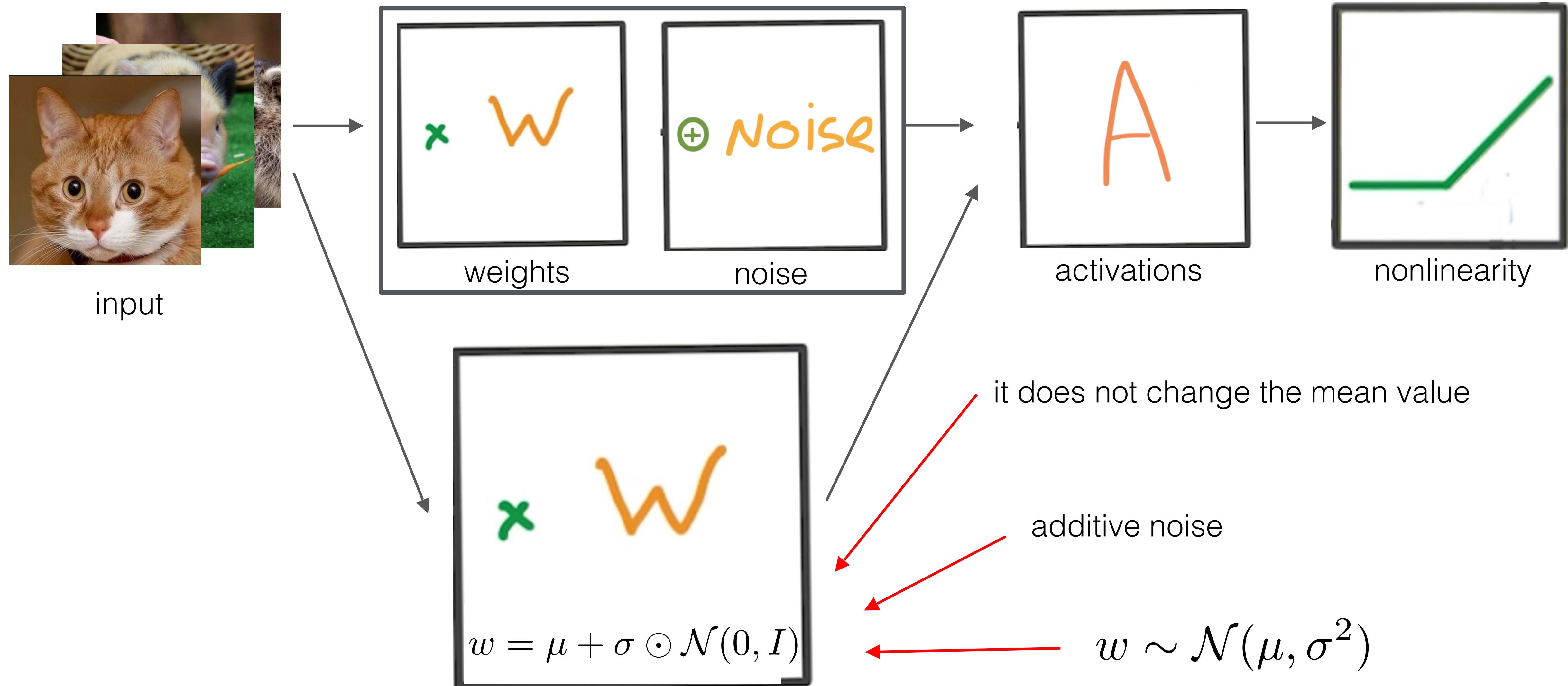
Final algorithm (without LRT)

Training on a mini-batch X with labels Y :

1. Sample weights: $\hat{w} = \mu + \sigma \odot \epsilon$, $\epsilon \sim \mathcal{N}(0, I)$
2. Forward pass: $Y_{\text{pred}} = NN(X, \hat{w}, b)$
3. Backward pass: compute stochastic gradients of ELBO:

$$\nabla_{\mu, \log \sigma, b} \left(N \cdot \text{Loss}(Y, Y_{\text{pred}}) + \text{SparseReg}(\sigma/\mu) \right)$$

Forward pass with stochastic weights



Final algorithm (without LRT)

Training on a mini-batch X with labels Y :

1. Sample weights: $\hat{w} = \mu + \sigma \odot \epsilon$, $\epsilon \sim \mathcal{N}(0, I)$
2. Forward pass: $Y_{\text{pred}} = NN(X, \hat{w}, b)$
3. Backward pass: compute stochastic gradients of ELBO:

$$\nabla_{\mu, \log \sigma, b} \left(N \cdot \text{Loss}(Y, Y_{\text{pred}}) + \text{SparseReg}(\sigma/\mu) \right)$$

Final algorithm (with LRT)

Training on a mini-batch X with labels Y :

1. Sample noise $\epsilon \sim \mathcal{N}(0, I)$
2. Forward pass with LRT: $Y_{\text{pred}} = NN(X, \mu, \sigma, \epsilon, b)$
3. Backward pass: compute stochastic gradients of ELBO:

$$\nabla_{\mu, \log \sigma, b} \left(N \cdot \text{Loss}(Y, Y_{\text{pred}}) + \text{SparseReg}(\sigma/\mu) \right)$$

Final algorithm (with LRT)

Training on a mini-batch X with labels Y :

1. Sample noise $\epsilon \sim \mathcal{N}(0, I)$
2. Forward pass with LRT: $Y_{\text{pred}} = NN(X, \mu, \sigma, \epsilon, b)$
3. Backward pass: compute stochastic gradients of ELBO:

$$\nabla_{\mu, \log \sigma, b} \left(N \cdot \text{Loss}(Y, Y_{\text{pred}}) + \text{SparseReg}(\sigma/\mu) \right)$$

Pruning after training:

If $\mu_{ij}^2/\sigma_{ij}^2 < \text{threshold}$:

$$\theta_{ij} = 0, \sigma_{ij} = 0$$

signal-to-noise ratio

Final algorithm (with LRT)

Training on a mini-batch X with labels Y :

1. Sample noise $\epsilon \sim \mathcal{N}(0, I)$
2. Forward pass with LRT: $Y_{\text{pred}} = NN(X, \mu, \sigma, \epsilon, b)$
3. Backward pass: compute stochastic gradients of ELBO:

$$\nabla_{\mu, \log \sigma, b} \left(N \cdot \text{Loss}(Y, Y_{\text{pred}}) + \text{SparseReg}(\sigma/\mu) \right)$$

Pruning after training:

If $\mu_{ij}^2/\sigma_{ij}^2 < \text{threshold}$:

$$\theta_{ij} = 0, \sigma_{ij} = 0$$

Prediction for a mini-batch X :

Return $Y_{\text{pred}} = NN(X, \mu, b)$

do not ensemble because we want
the most compact and fast network

Automatic relevance determination

(other prior but very similar method)

- Empirical Bayes
 - optimize w. r. t. prior parameters \Rightarrow automatically choose hyperparameters

Automatic relevance determination

(other prior but very similar method)

- Empirical Bayes
 - optimize w. r. t. prior parameters \Rightarrow automatically choose hyperparameters
- Normal prior with **learnable** variance for **each weight**:

$$p(w_{ij} | \tau_{ij}) = \mathcal{N}(0, \tau_{ij}^{-1})$$

Automatic relevance determination

(other prior but very similar method)

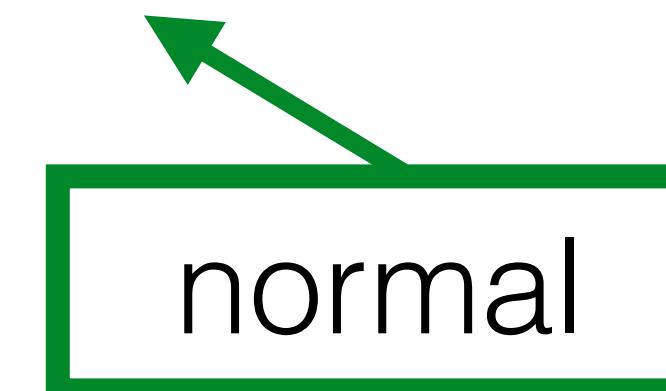
- Empirical Bayes
 - optimize w. r. t. prior parameters \Rightarrow automatically choose hyperparameters

- Normal prior with **learnable** variance for **each weight**:

$$p(w_{ij} | \tau_{ij}) = \mathcal{N}(0, \tau_{ij}^{-1})$$

- KL-divergence:

$$KL(q(w_{ij} | \mu_{ij}, \sigma_{ij}) || p(w_{ij} | \tau_{ij})) = -\log \sigma_{ij} - \frac{1}{2} \log \tau_{ij} + \frac{1}{2} \tau_{ij} (\mu_{ij}^2 + \sigma_{ij}^2)$$



Automatic relevance determination

(other prior but very similar method)

- Empirical Bayes
 - optimize w. r. t. prior parameters \Rightarrow automatically choose hyperparameters

- Normal prior with **learnable** variance for **each weight**:

$$p(w_{ij} | \tau_{ij}) = \mathcal{N}(0, \tau_{ij}^{-1})$$

- KL-divergence:

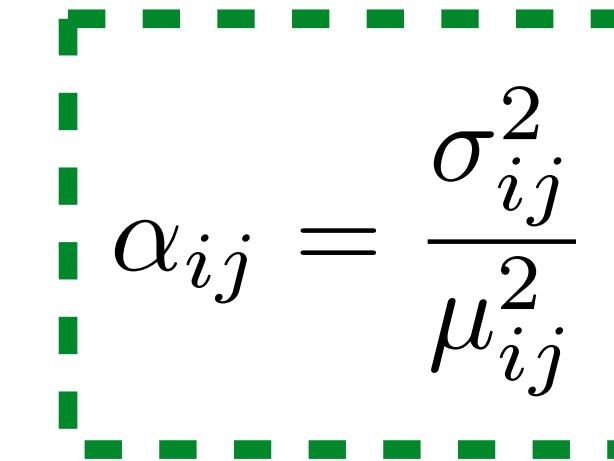
$$KL(q(w_{ij} | \mu_{ij}, \sigma_{ij}) || p(w_{ij} | \tau_{ij})) = -\log \sigma_{ij} - \frac{1}{2} \log \tau_{ij} + \frac{1}{2} \tau_{ij} (\mu_{ij}^2 + \sigma_{ij}^2) \equiv$$

- Analytical optimization w. r. t. prior parameters τ_{ij} :

$$\boxed{\tau_{ij}^* = (\mu_{ij}^2 + \sigma_{ij}^2)^{-1}}$$

$$\stackrel{\tau_{ij}^*}{\equiv} -\frac{1}{2} \log(\mu_{ij}^2 + \sigma_{ij}^2)^{-1} = -\frac{1}{2} \log\left(1 + \frac{\mu_{ij}^2}{\sigma_{ij}^2}\right)$$

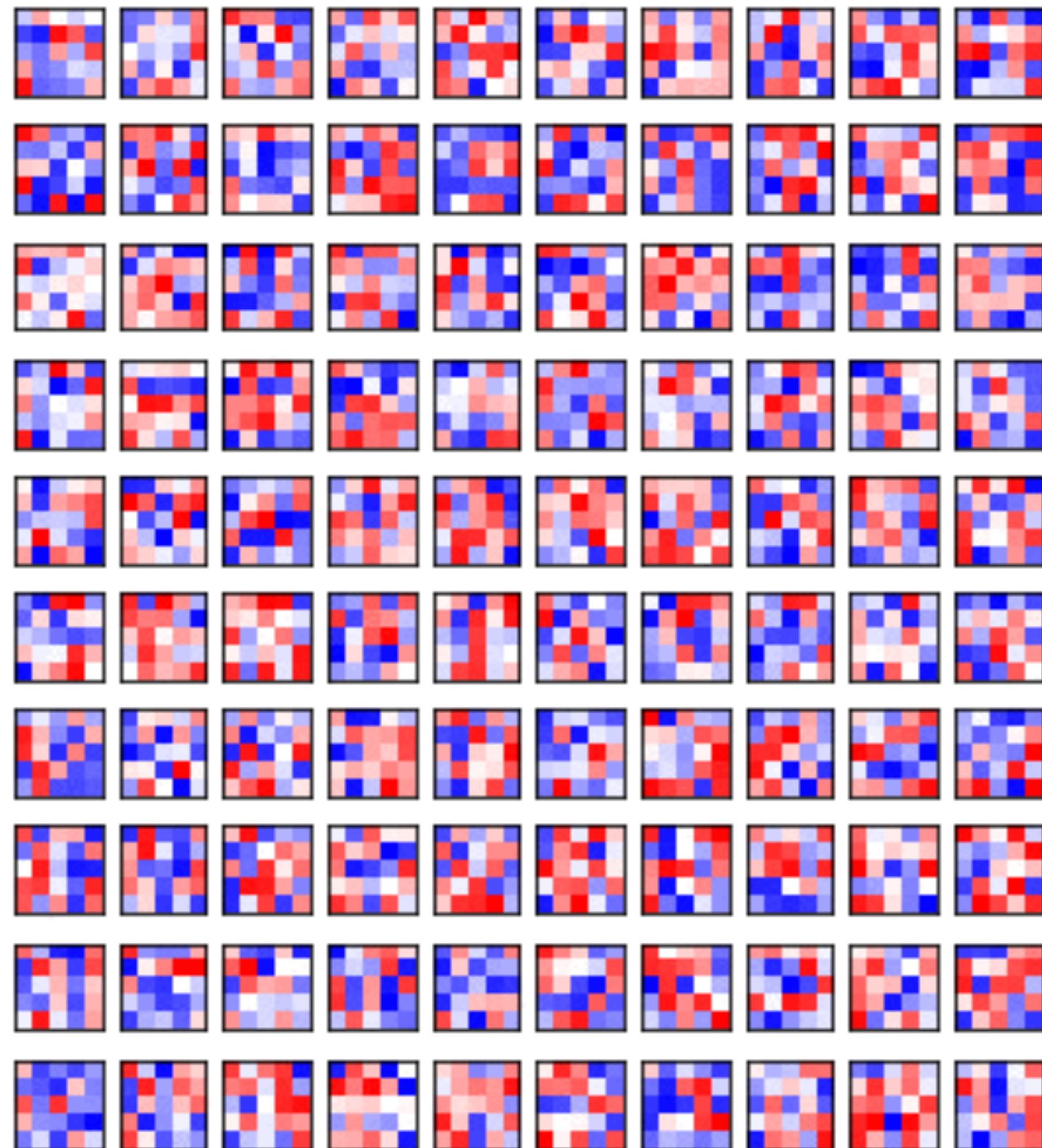
Comparing two models

	SparseVD	ARD
Prior	$p(w_{ij}) \propto \frac{1}{ w_{ij} }$	$p(w_{ij} \tau_{ij}) = \mathcal{N}(0, \tau_{ij}^{-1})$
Prior parameters	None	τ_{ij}
KL-divergence	$k_1 \sigma(k_2 + k_3 \log \alpha_{ij})) -$ $-0.5 \log(1 + \alpha_{ij}^{-1}) + C$ 	$-0.5 \log(1 + \alpha_{ij}^{-1})$ 

Training procedure is the same (approx. posterior, RT, LRT etc.)

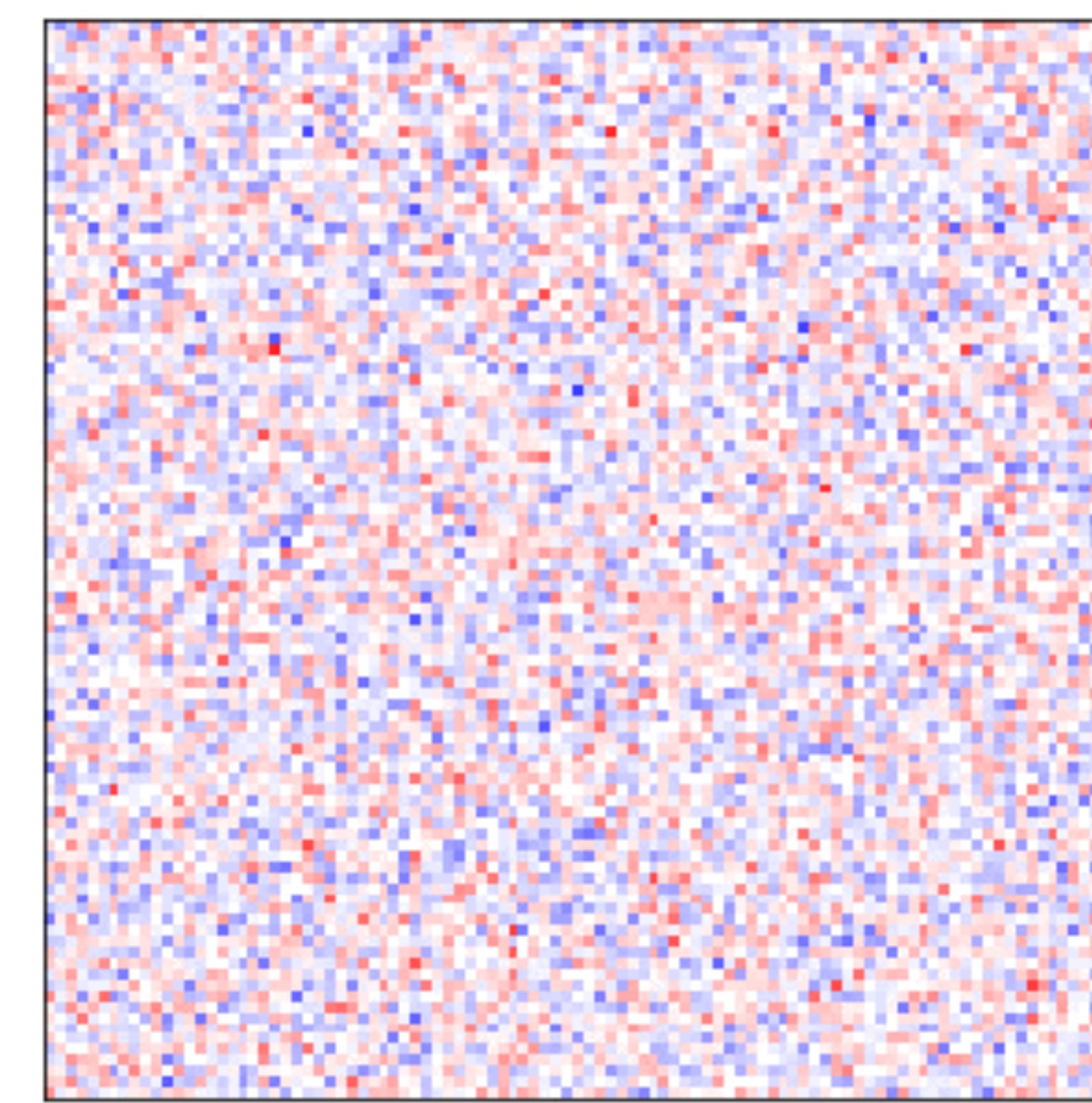
Sparse variational dropout: visualization

Epoch: 0 Compression ratio: 1x Accuracy: 8.4



LeNet-5: convolutional layer

Epoch: 0 Compression ratio: 1x Accuracy: 8.4

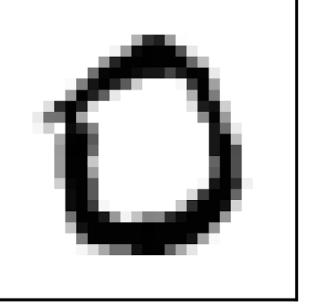
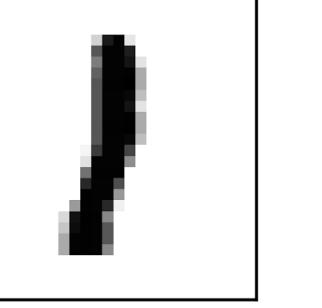
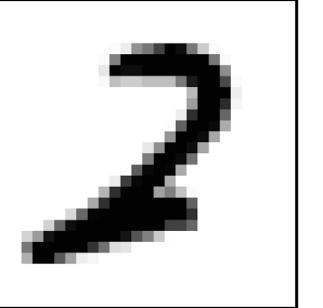
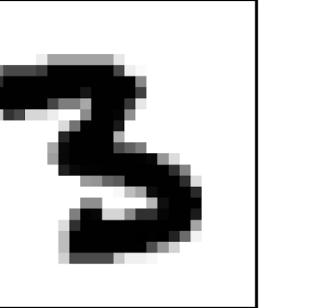
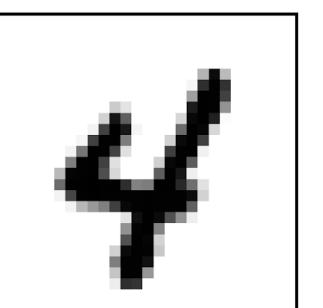
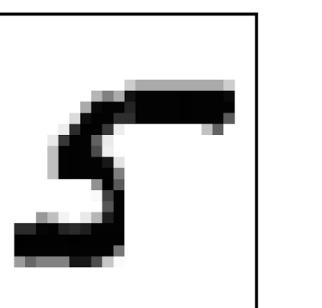
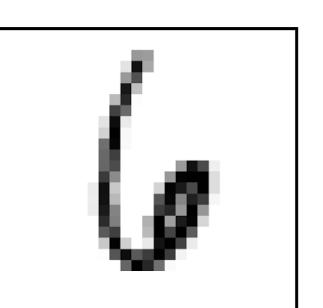
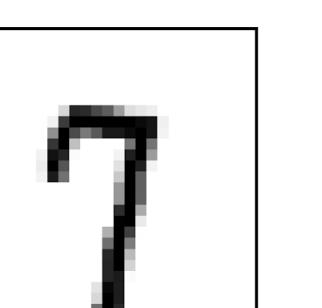
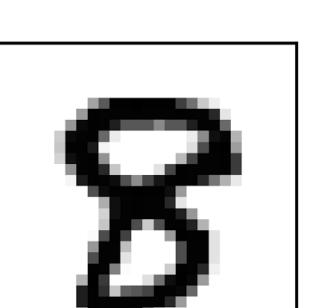
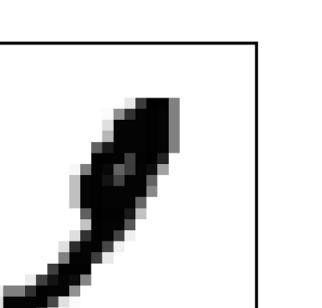


LeNet-5: fully-connected layer
(100 x 100 patch)

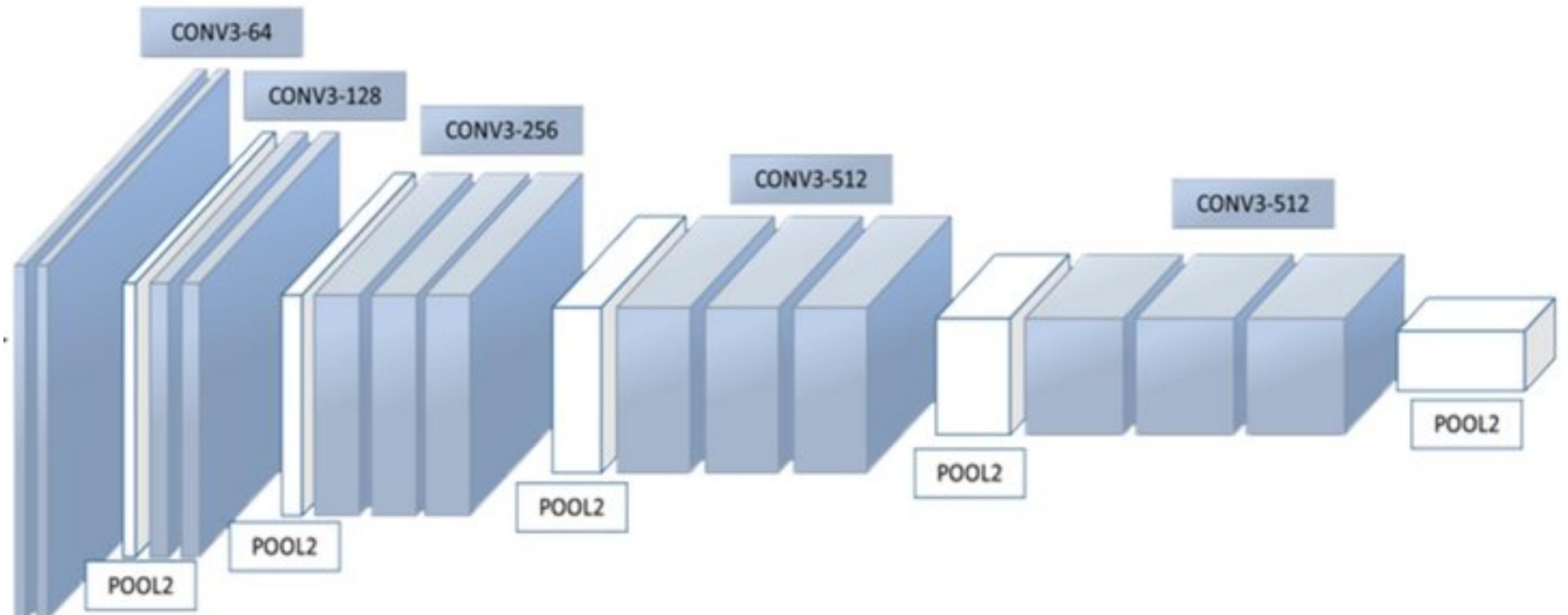
Lenet-5-Caffe and Lenet-300-100 on MNIST

Fully Connected network: LeNet-300-100

Convolutional network: Lenet-5-Caffe

Network	Method	Error %	Sparsity per Layer %	$\frac{ \mathbf{W} }{ \mathbf{W}_{\neq 0} }$		
LeNet-300-100	Original	1.64		1		
	Pruning	1.59	92.0 – 91.0 – 74.0	12		
	DNS	1.99	98.2 – 98.2 – 94.5	56		
	SWS	1.94		23		
	(ours) Sparse VD	1.92	98.9 – 97.2 – 62.0	68		
LeNet-5-Caffe	Original	0.80		1		
	Pruning	0.77	34 – 88 – 92.0 – 81	12		
	DNS	0.91	86 – 97 – 99.3 – 96	111		
	SWS	0.97		200		
	(ours) Sparse VD	0.75	67 – 98 – 99.8 – 95	280		

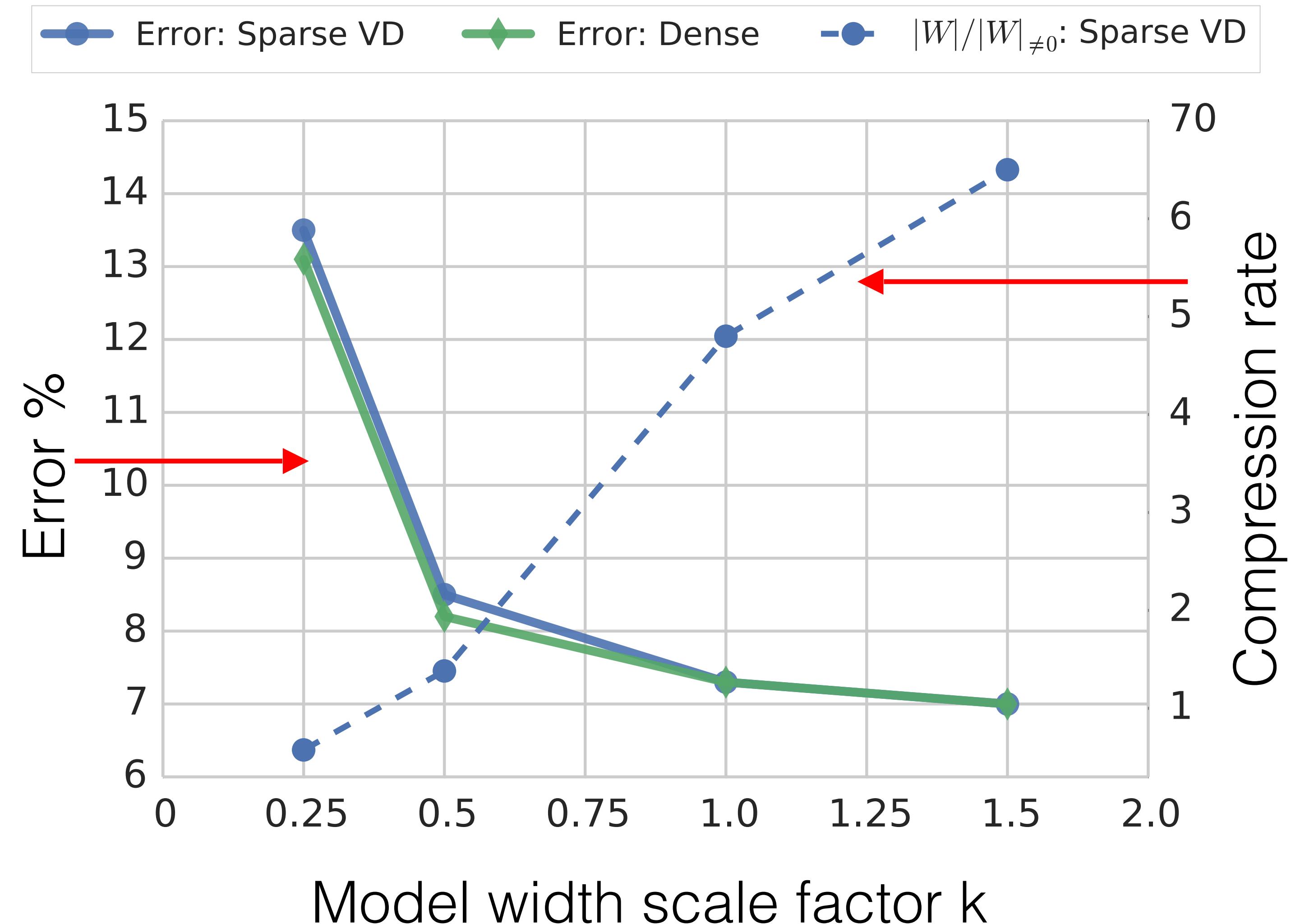
VGG-like on CIFAR-10



- 13 Convolutional layers and 2 Fully-Connected layers
- Pre-Activation Batch Norm and Binary Dropout after each layer

VGG-like on CIFAR-10

Number of filters / neurons is linearly scaled by k (the width of the network)



Random Labeling



Dataset	Architecture	Train Acc.	Test Acc.	Sparsity
MNIST	FC + BD	100%	10%	—
MNIST	FC + Sparse VD	10%	10%	100%
CIFAR-10	VGG + BD	100%	10%	—
CIFAR-10	VGG + Sparse VD	10%	10%	100%

No dependency between data and labels \Rightarrow Sparse VD yields an empty model
where conventional models easily overfit.

Agenda

- Sparsification: what and why
- Bayesian neural networks
- Sparse variational dropout
- Practical assignment: implementation of SparseVD
- Model enhancements

Practical assignment



1. Take a look at an example:
FC for a MNIST in PyTorch (or skip it :)
2. Implement two classes for SparseVD and incorporate them into model training
3. Train SparseVD and visualize weights

Agenda

- Sparsification: what and why
- Bayesian neural networks
- Sparse variational dropout
- Practical assignment: implementation of SparseVD
- Model enhancements

What's next?

- SparseVD for RNNs
- Structured Bayesian sparsification
- How does it all work on large datasets and networks

What's next?

- SparseVD for RNNs
- Structured Bayesian sparsification
- How does it all work on large datasets and networks

SparseVD for recurrent neural networks

$$\sum_{i=1}^N \mathbb{E}_{q(w|\mu,\sigma)} \log p(y^i|x^i, w) - KL(q(w|\mu,\sigma)||p(w)) \rightarrow \max_{\mu, \log \sigma}$$

Model specification:

- Choose particular prior  log-uniform distribution

Training:

- Choose particular family for approximate posterior  normal distribution
- How to compute KL-divergence?  analytical approximation
- How to estimate the expectation? 

Important RNN specifics

Input is a sequence: $x = x_1 \dots x_T$:

$$\begin{aligned} & - \sum_{i=1}^N \int q(w|\theta, \sigma) \log p(y^i | \overbrace{x_0^i, \dots, x_T^i}^{\text{sequence}}, w, B) dw + \\ & + \sum_{w_{ij} \in \omega} KL(q(w_{ij}|\theta_{ij}, \sigma_{ij}) || p(w_{ij})) \rightarrow \min_{\theta, \sigma, B} \end{aligned}$$

- 1 Use one sample w for all t for 1 sequence

Important RNN specifics

Input is a sequence: $x = x_1 \dots x_T$:

$$\begin{aligned} & - \sum_{i=1}^N \int q(w|\theta, \sigma) \log p(y^i | \overbrace{x_0^i, \dots, x_T^i}^{\text{sequence}}, w, B) dw + \\ & + \sum_{w_{ij} \in \omega} KL(q(w_{ij}|\theta_{ij}, \sigma_{ij}) || p(w_{ij})) \rightarrow \min_{\theta, \sigma, B} \end{aligned}$$

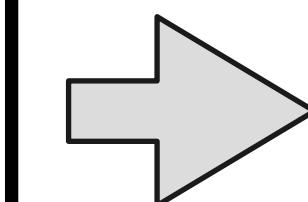
- 1 Use one sample w for all t for 1 sequence
- 2 LRT is not applicable \Rightarrow sample one w for a whole mini-batch

Why LRT is not applicable in RNNs?

$$h_{t+1} = g(W^x x_{t+1} + W^h h_t + b^h)$$

For W^h :

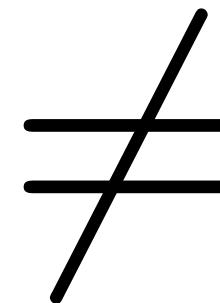
h_t is a random variable
depending on W^h



$W^h h_t$ is not
normal

For W^x :

Sampling same
noise on weights
for all t



Sampling same
noise on preactivation
for all t

Experiments: text classification

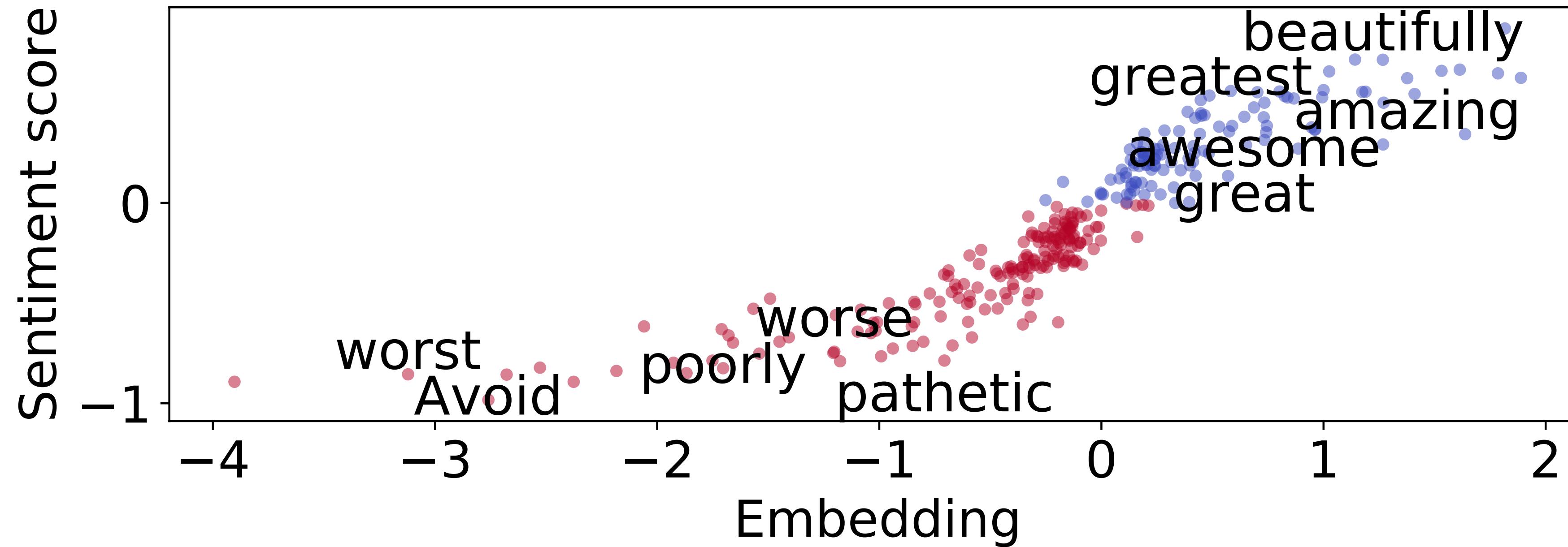
Datasets: IMDb (25K texts, 2 classes), AGNews (120K texts, 4 classes)

Architecture: Embedding → LSTM → FC on h_T

Task	Method	Accuracy %	Compression	Vocabulary
IMDb	Original	84.1	1x	20000
	SparseVD	85.1	1135x	4611
	SparseVD-Voc	83.6	12985x	292
AGNews	Original	90.6	1x	20000
	SparseVD	88.8	322x	5727
	SparseVD-Voc	89.2	469x	2444

Experiments: text classification

IMDb: the only remaining embedding component for remaining words



sentiment score = (#pos. — #neg.) / # all texts with word

Experiments: text generation

PTB Corpus: char-level (6M chars), word-level (0.9M words)

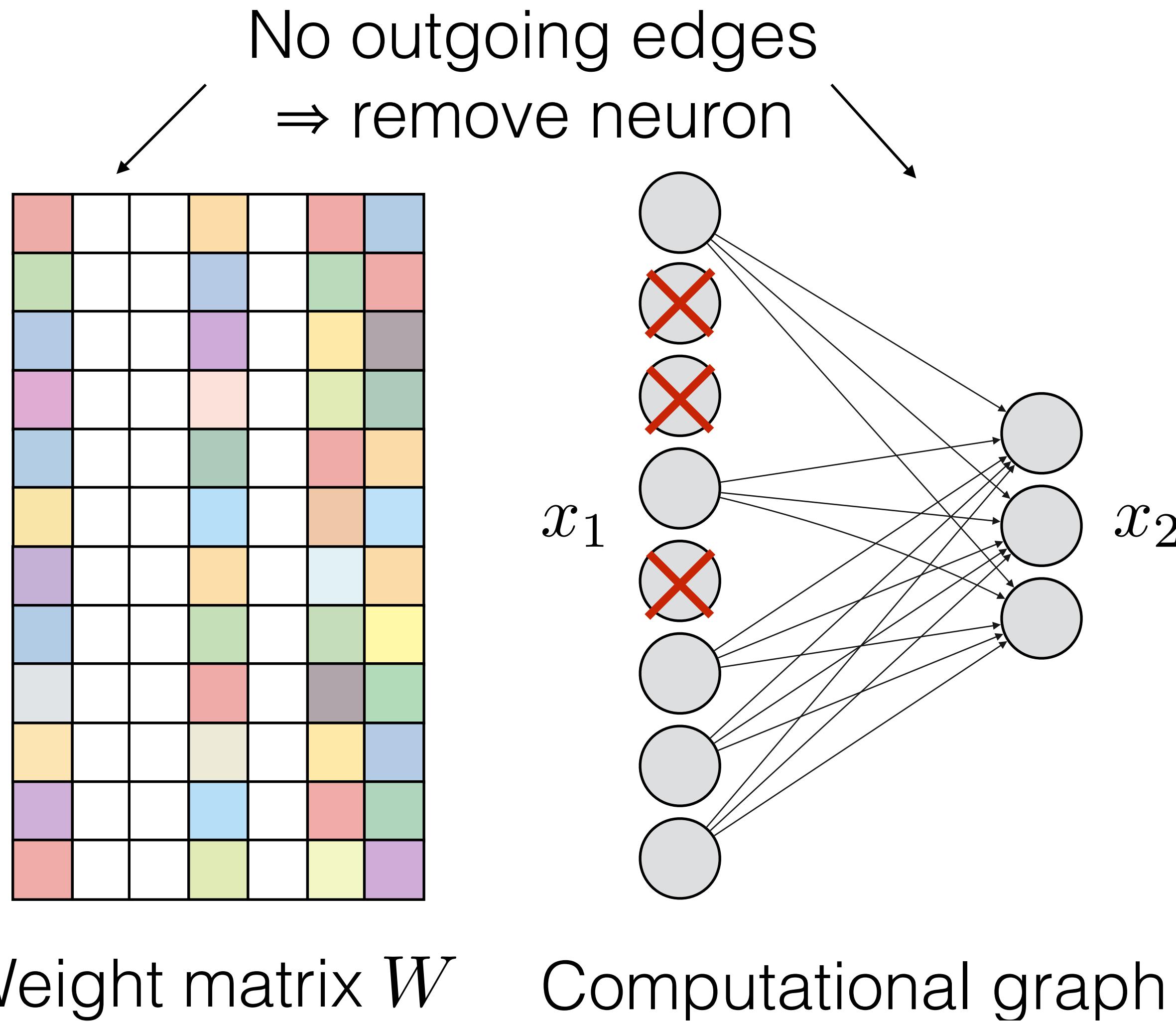
Architecture: LSTM → FC

Task	Method	Valid	Test	Compression	Vocabulary
Char PTB (Bits-per-char)	Original	1.498	1.454	1x	50
	SparseVD	1.472	1.429	7.6x	50
	SparseVD-Voc	1.4584	1.4165	5.8x	48
Word PTB (Perplexity)	Original	135.6	129.5	1x	10000
	SparseVD	115.0	109.0	14.0x	9985
	SparseVD-Voc	126.3	120.6	11.1x	4353

What's next?

- SparseVD for RNNs
- Structured Bayesian sparsification
- How does it all work on large datasets and networks

Structured Bayesian sparsification



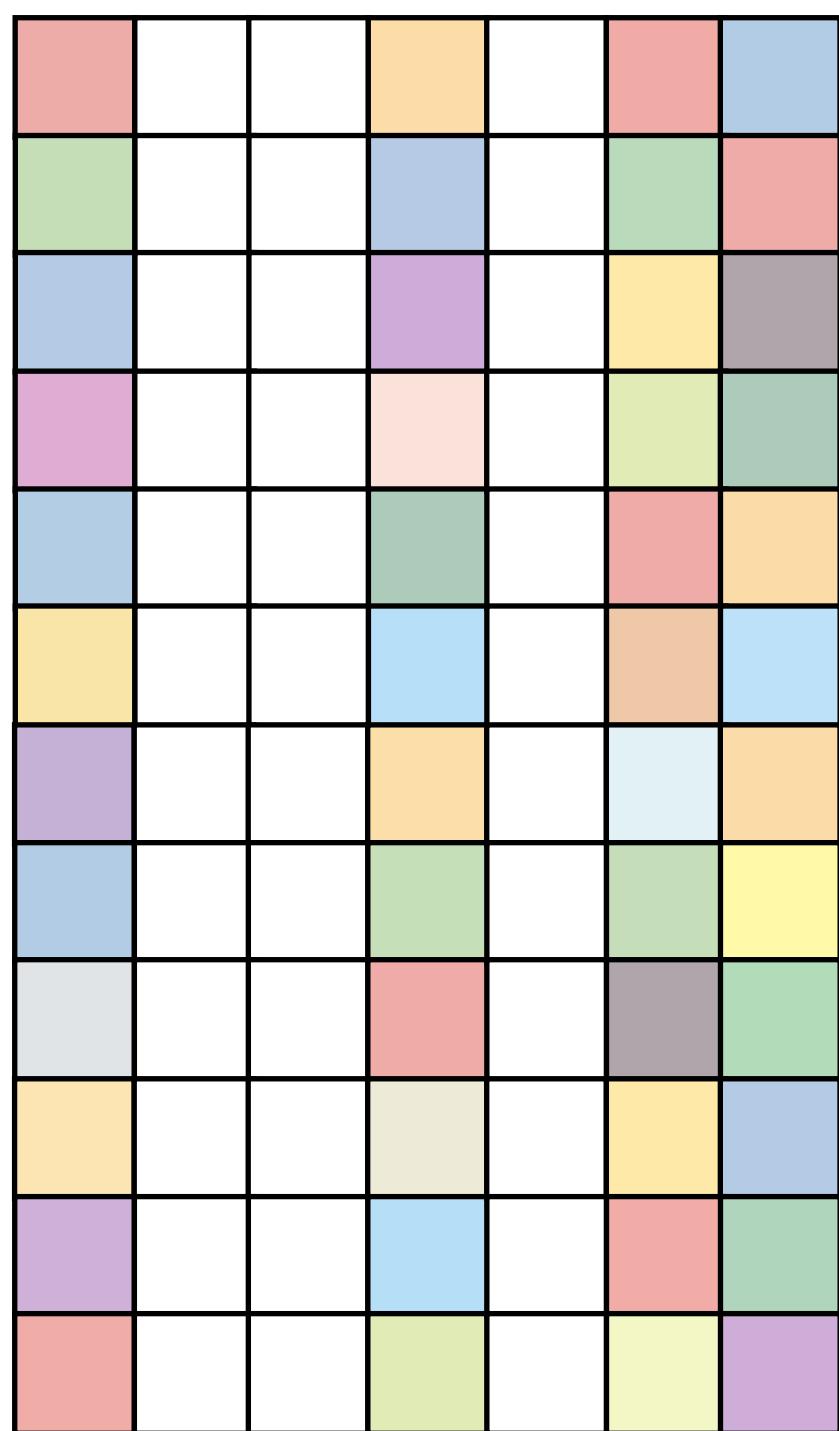
Benefits of structured sparsification:

- more efficient compression
- speed-up of forward pass (faster testing stage)

Structured Bayesian sparsification

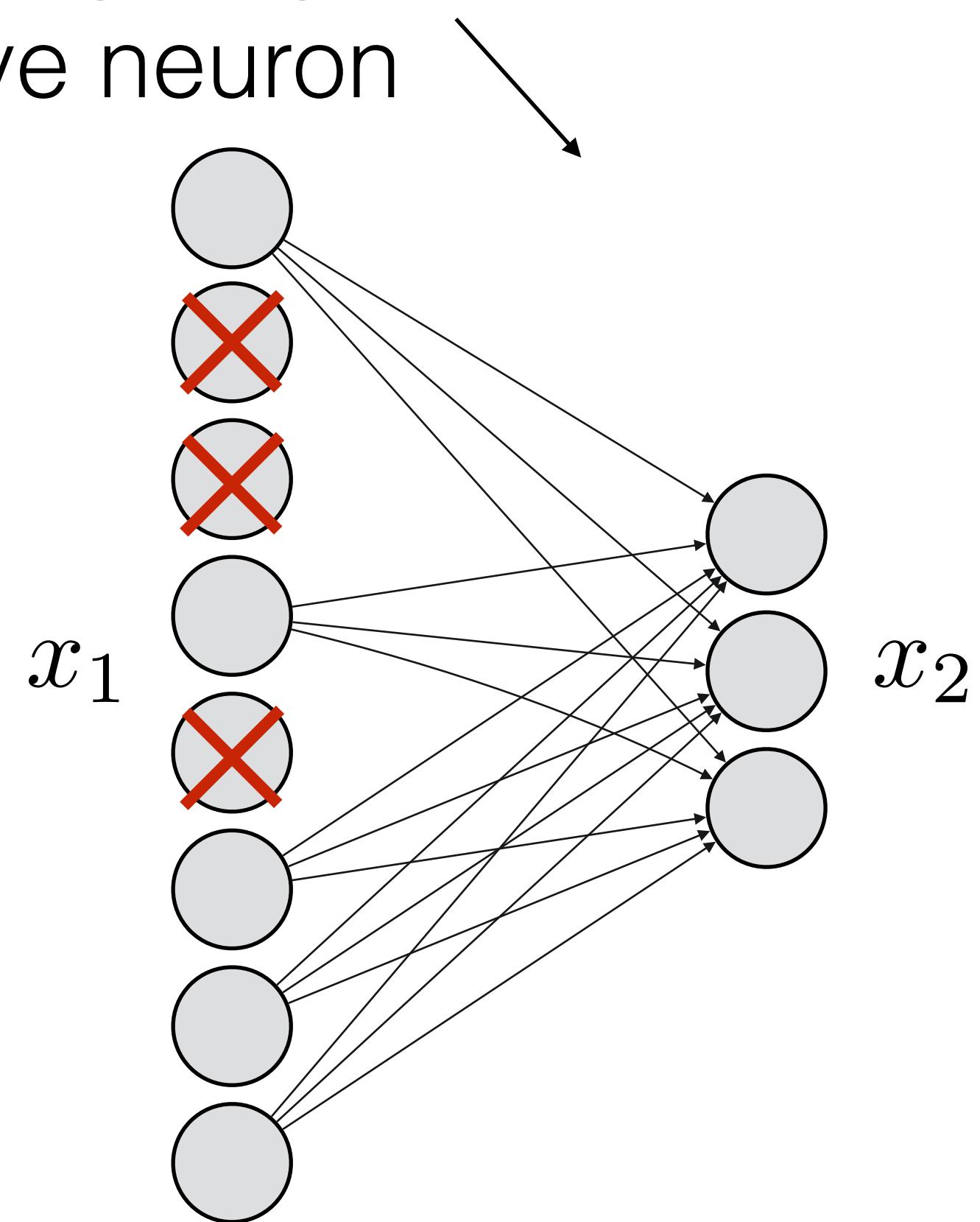
No outgoing edges

⇒ remove neuron



Weight matrix W

Computational graph

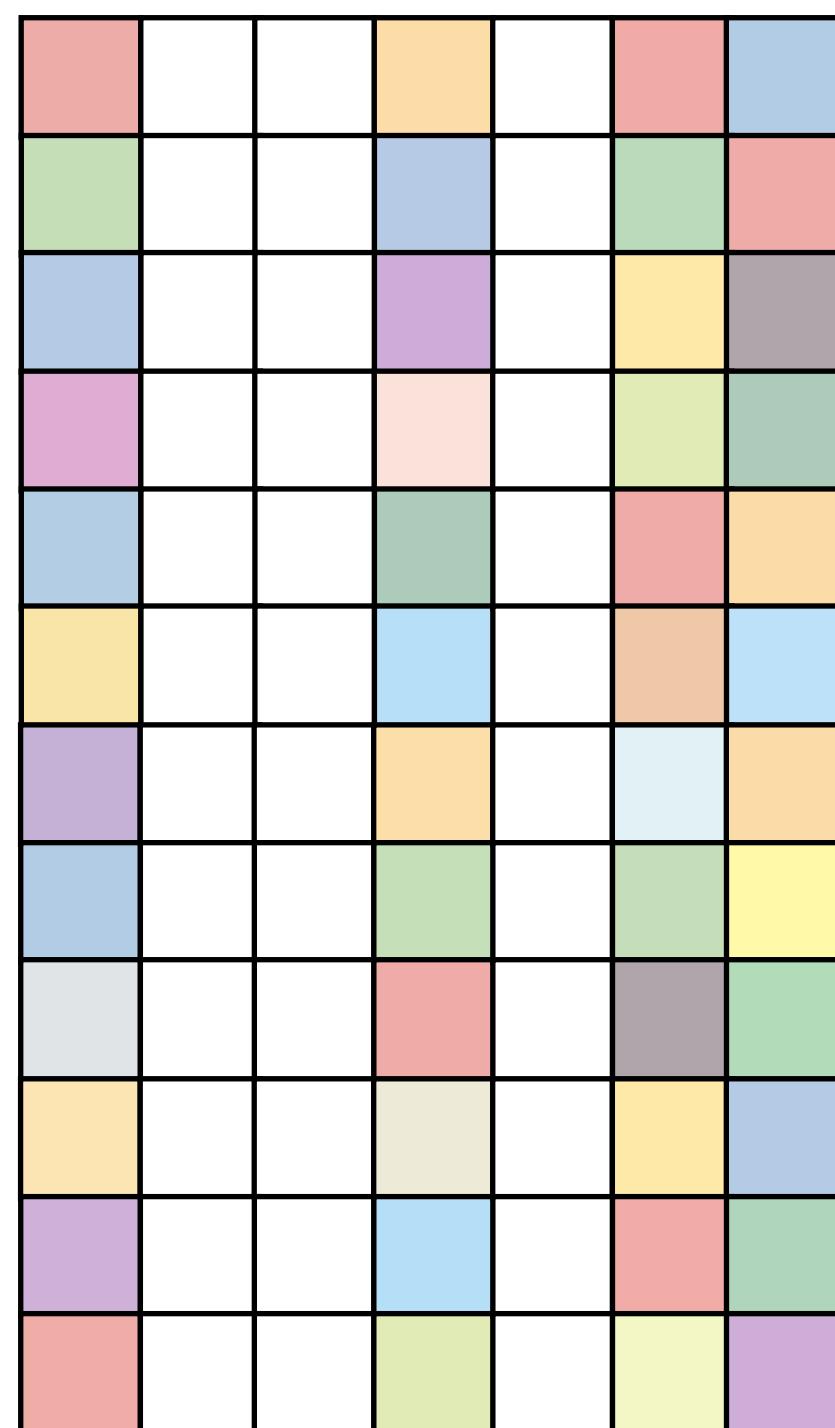


$$x_2 = \sigma(Wx_1 + b)$$

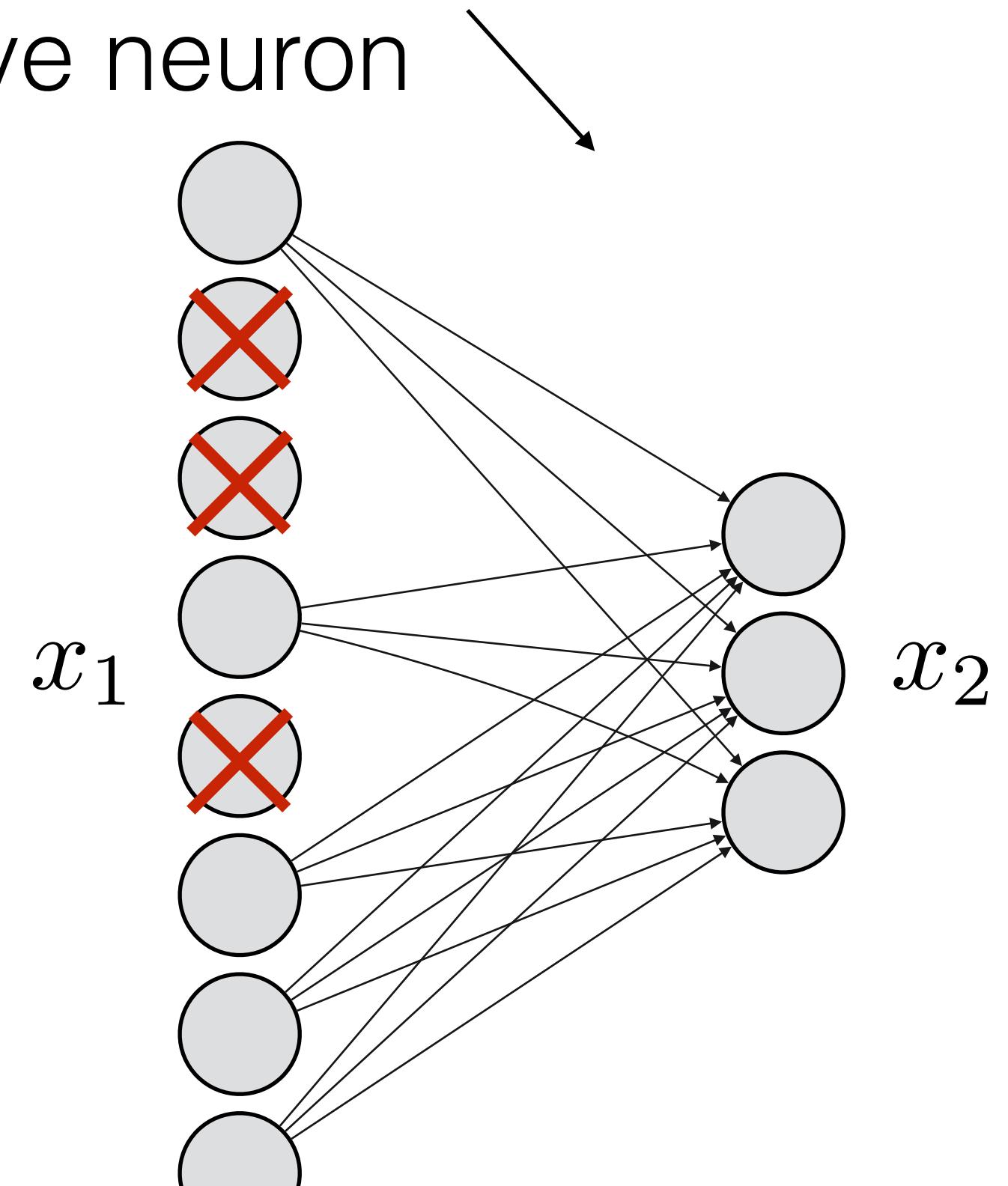
Structured Bayesian sparsification

No outgoing edges

⇒ remove neuron



Weight matrix W



Computational graph

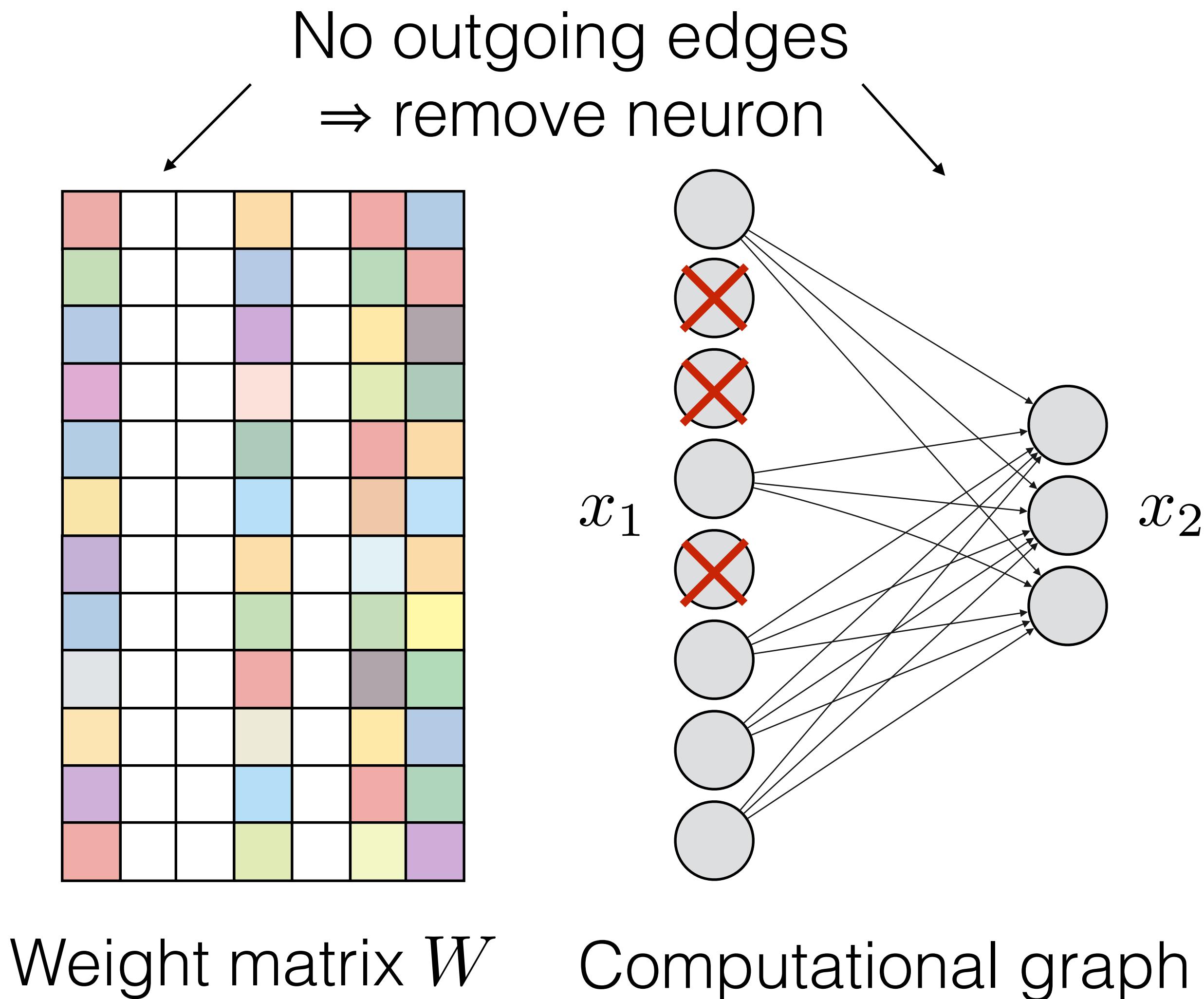
$$x_2 = \sigma(Wx_1 + b)$$

Multiply x_1 by group variable z :

$$x_2 = \sigma(W(x_1 \odot z) + b)$$

zero component in $z \Rightarrow$
no outgoing edges ⇒
remove neuron

Structured Bayesian sparsification: training

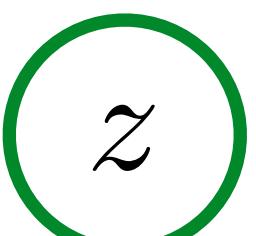


- Treat z in **the same way** as W :
- Log-uniform prior on z and on W
 - Normal approx. posterior for both
 - RT & LRT for W
 - RT for z (LRT is not needed)

A bit more mathematical motivation

Let's choose standard normal prior on W :

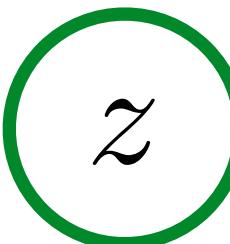
$$x_2 = \sigma(W(x_1 \odot z) + b) \quad p(z_i) \propto \frac{1}{|z_i|}; \quad p(w_{ij}) = \mathcal{N}(0, 1)$$

 
independent

A bit more mathematical motivation

Let's choose standard normal prior on W :

$$x_2 = \sigma(W(x_1 \odot z) + b) \quad p(z_i) \propto \frac{1}{|z_i|}; \quad p(w_{ij}) = \mathcal{N}(0, 1)$$


 
independent

$$p(\underbrace{w_{ij} z_i}_{\tilde{w}_{ij}} | z_i) = \mathcal{N}(0, z_i^2)$$

A bit more mathematical motivation

Let's choose standard normal prior on W :

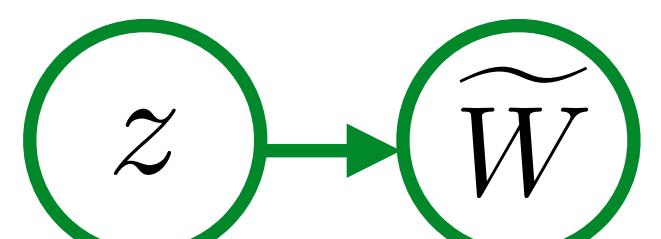
$$x_2 = \sigma(W(x_1 \odot z) + b) \quad p(z_i) \propto \frac{1}{|z_i|}; \quad p(w_{ij}) = \mathcal{N}(0, 1)$$



independent

Equivalent model:

$$x_2 = \sigma(\tilde{W}x_1 + b) \quad p(z_i) \propto \frac{1}{|z_i|}; \quad p(\tilde{w}_{ij}|z_i) = \mathcal{N}(0, z_i^2)$$



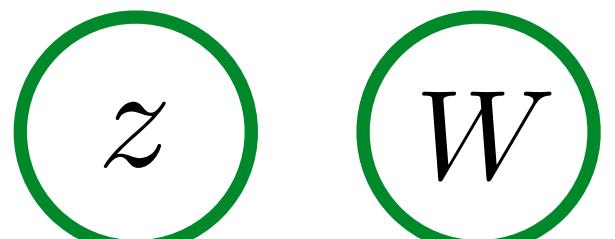
dependent!

A bit more mathematical motivation

Let's choose standard normal prior on W :

$$x_2 = \sigma(W(x_1 \odot z) + b)$$

$$p(z_i) \propto \frac{1}{|z_i|}; \quad p(w_{ij}) = \mathcal{N}(0, 1)$$

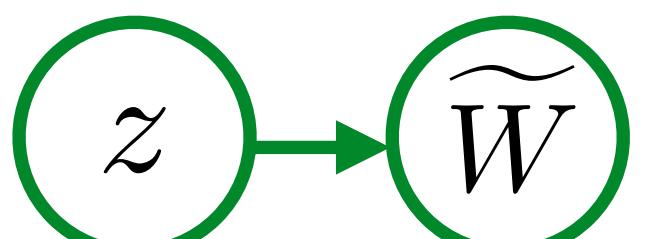


independent

Equivalent model:

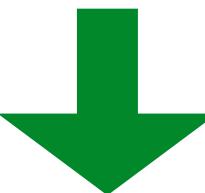
$$x_2 = \sigma(\tilde{W}x_1 + b)$$

$$p(z_i) \propto \frac{1}{|z_i|}; \quad p(\tilde{w}_{ij}|z_i) = \mathcal{N}(0, z_i^2)$$



dependent!

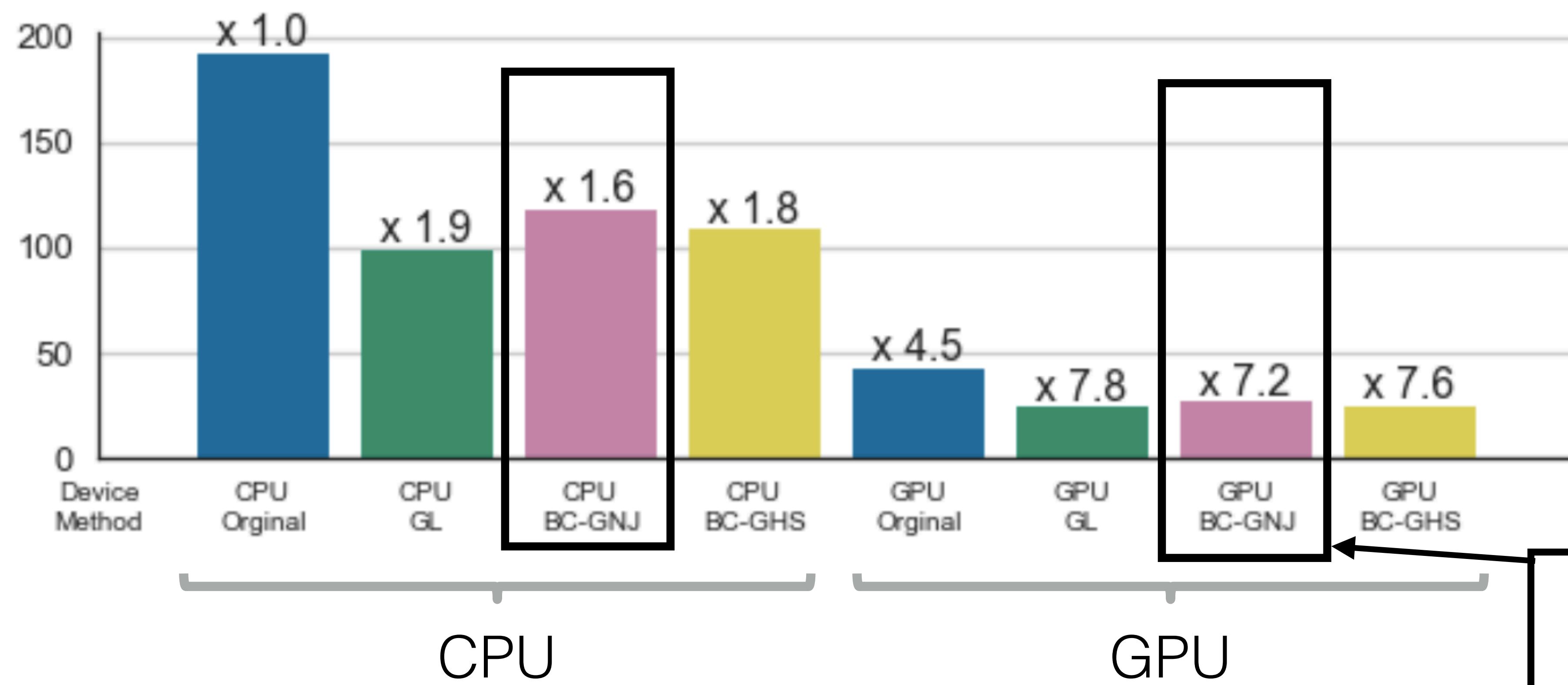
$$p(\tilde{w}_{ij}) \propto \int \frac{1}{|z_i|} \mathcal{N}(0, z_i^2) dz_i = \frac{1}{|\tilde{w}_{ij}|}$$



SparseVD model!
But with structured sparsity

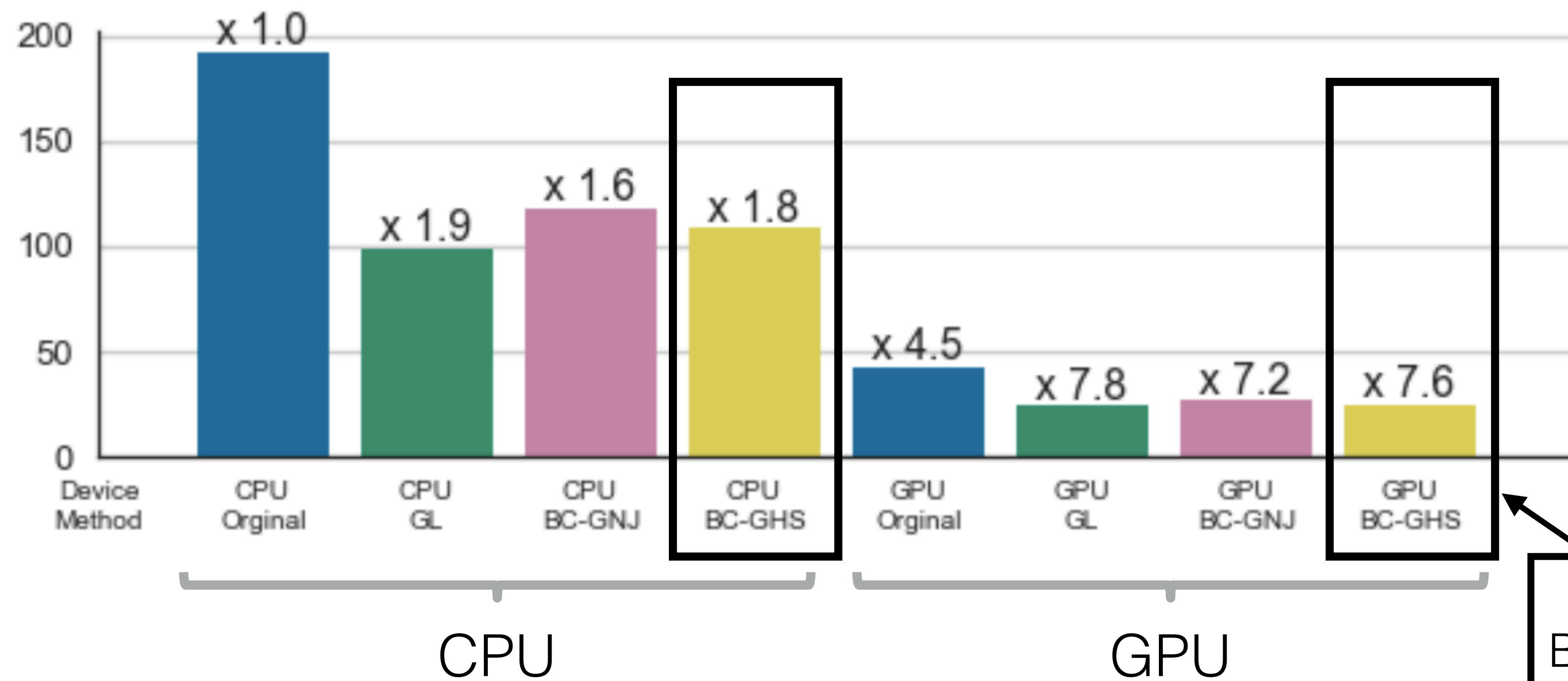
Structured Bayesian sparsification: results

Speed-up of forward pass (testing stage) for Lenet-5-Caffe



Structured Bayesian sparsification: results

Speed-up of forward pass (testing stage) for Lenet-5-Caffe



What's next?

- SparseVD for RNNs
- Structured Bayesian sparsification
- How does it all work on large datasets and networks

Two popular frameworks for sparsification

Magnitude pruning	Bayesian sparsification
A lot of method hyperparameters	(Almost) no method hyperparameters
Need to choose training schedule	Need to choose training schedule
non-Bayesian	Bayesian!

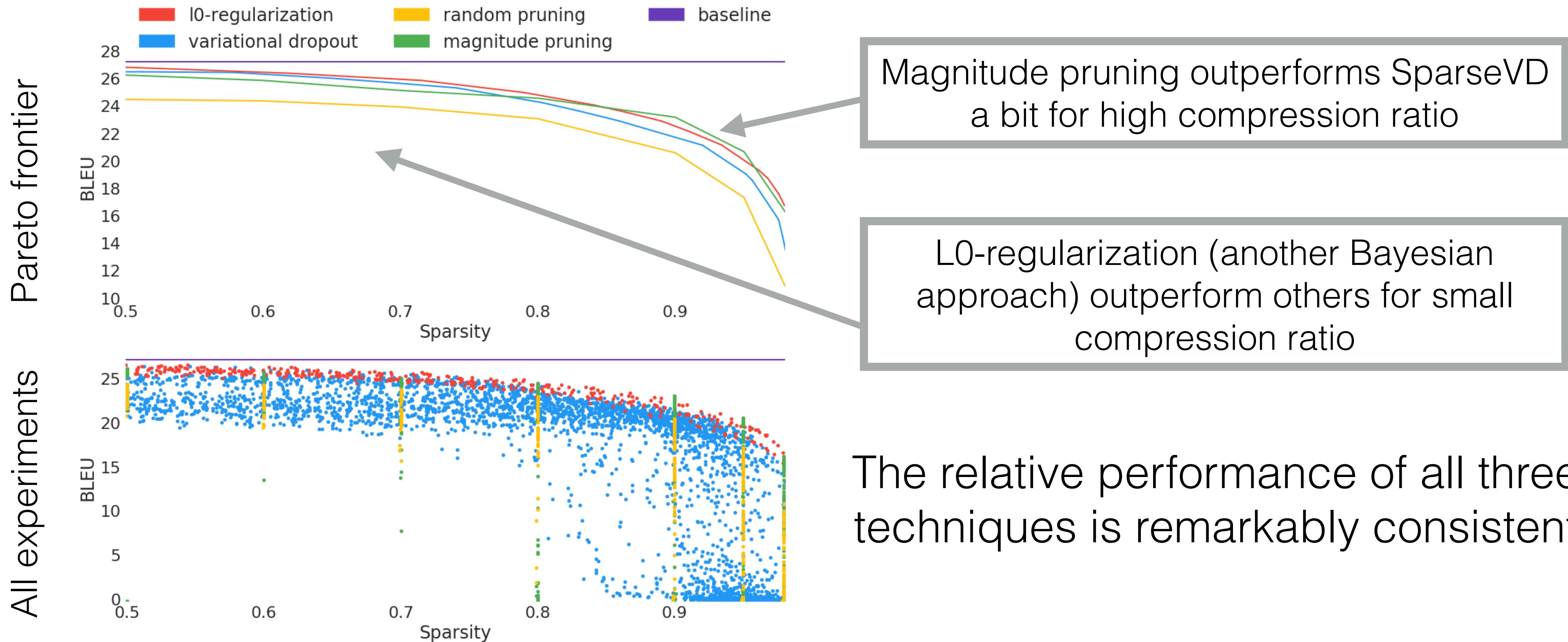
The state of sparsity in deep neural networks

- Recent (2019) work on comparing different (unstructured) sparsification techniques on **large** datasets and models:
 - Transformer for neural machine translation (WMT 2014 English-to-German)
 - ResNet-50 for ImageNet
- Open-source code and top performing model checkpoints for reproducibility
- **Thousands** of experiments!

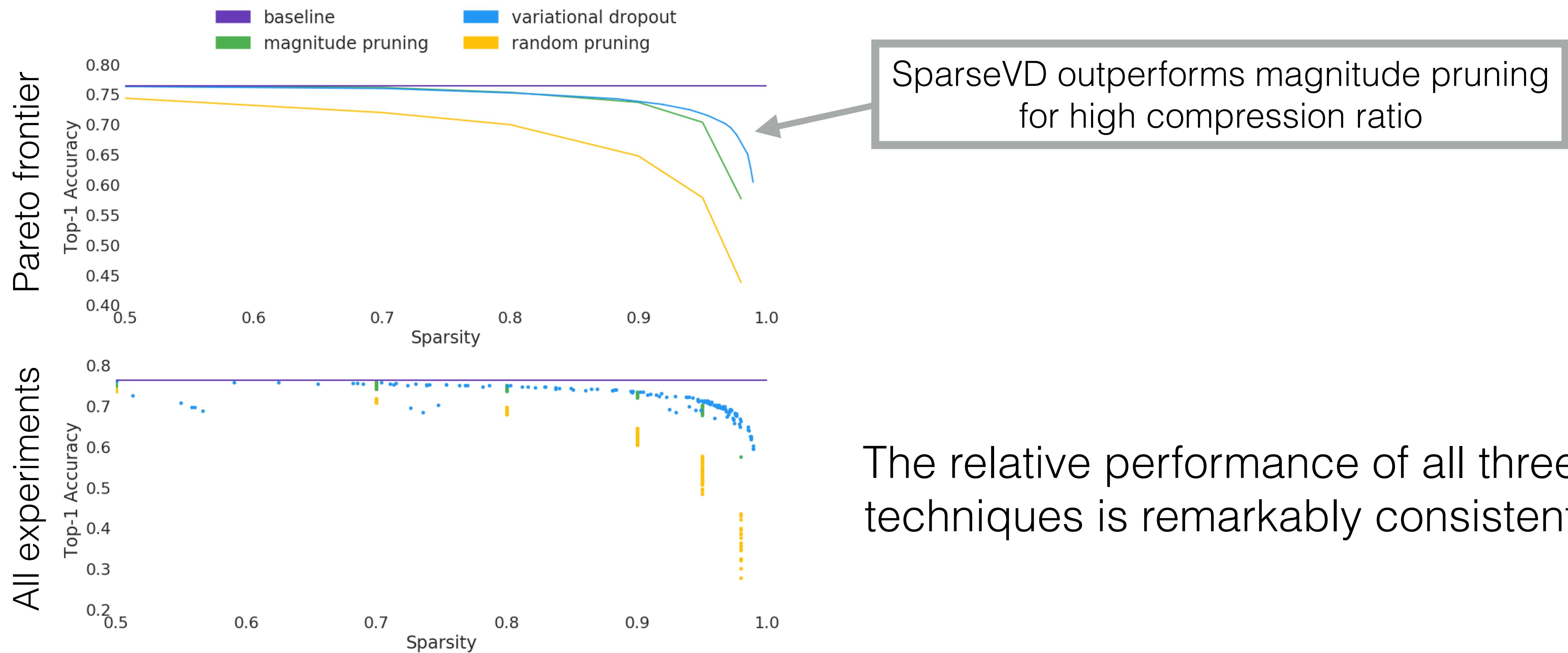
Hyperparameters for both tasks

- Hyperparameters tuned in Sparse variational dropout:
 - KL-divergence weight & annealing schedule (not grounded theoretically, but works in practice)
 - Sparsification threshold
- Hyperparameters tuned in Magnitude pruning:
 - target sparsity
 - starting iteration of the sparsification process
 - ending iteration of the sparsification process
 - frequency of pruning steps

Transformer for neural machine translation



ResNet-50 for ImageNet



Other Bayesian sparsification models

- Group horseshoe with half-Cauchy scale priors:
more levels is model hierarchy & usually better compression

$$s \sim \mathcal{C}^+(0, \tau_0); \quad \tilde{z}_i \sim \mathcal{C}^+(0, 1); \quad \tilde{w}_{ij} \sim \mathcal{N}(0, 1); \quad w_{ij} = \tilde{w}_{ij} \tilde{z}_i s$$

- L₀ regularization
relax non-differentiable L₀ regularizer & obtain exact zero weights

$$\mathcal{R}(\boldsymbol{\theta}) = \frac{1}{N} \left(\sum_{i=1}^N \mathcal{L}(h(\mathbf{x}_i; \boldsymbol{\theta}), \mathbf{y}_i) \right) + \lambda \|\boldsymbol{\theta}\|_0$$

relaxation is based on spike-and-slab prior (as well as approx. posterior)

$$p(z) = \text{Bernoulli}(\pi), \quad p(\theta|z=0) = \delta(\theta), \quad p(\theta|z=1) = \mathcal{N}(\theta|0, 1)$$

Summary

- Sparsification of neural networks is an urgent industrial problem that is well solved using Bayesian deep learning
- Bayesian sparsification relies on Bayesian neural networks that provide regularization, fast ensembling, uncertainty estimation etc.
- Bayesian sparsification is a theoretically grounded approach but several tricks are needed to train the model