

# Bayesian Sparsification of Gated Recurrent Neural Networks

Ekaterina Lobacheva<sup>\*1</sup> Nadezhda Chirkova<sup>\*1</sup> Dmitry Vetrov<sup>1,2</sup>

<sup>1</sup>Higher School of Economics, Samsung-HSE Laboratory

<sup>2</sup>Samsung AI Center, Moscow, Russia



SAMSUNG  
Research

$p(\mathbf{B}|\mathbf{A})$ yesgroup.ru

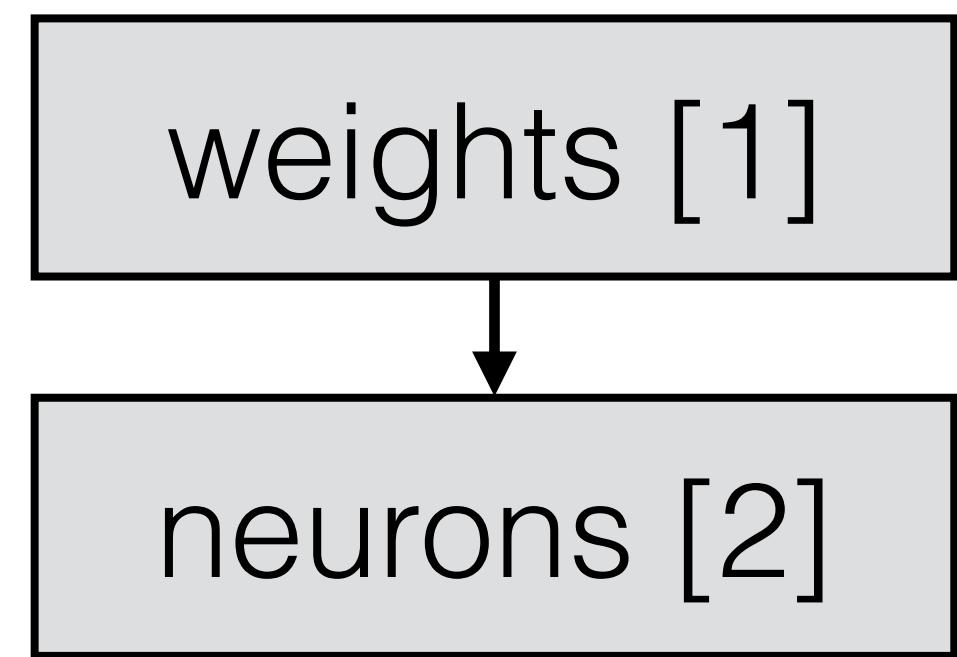
# Problem setting

Goal: learn structurally sparse gated RNNs (LSTM, GRU...)

# Problem setting

Goal: learn structurally sparse gated RNNs (LSTM, GRU...)

Which structure units to sparsify  
in vanilla RNN?



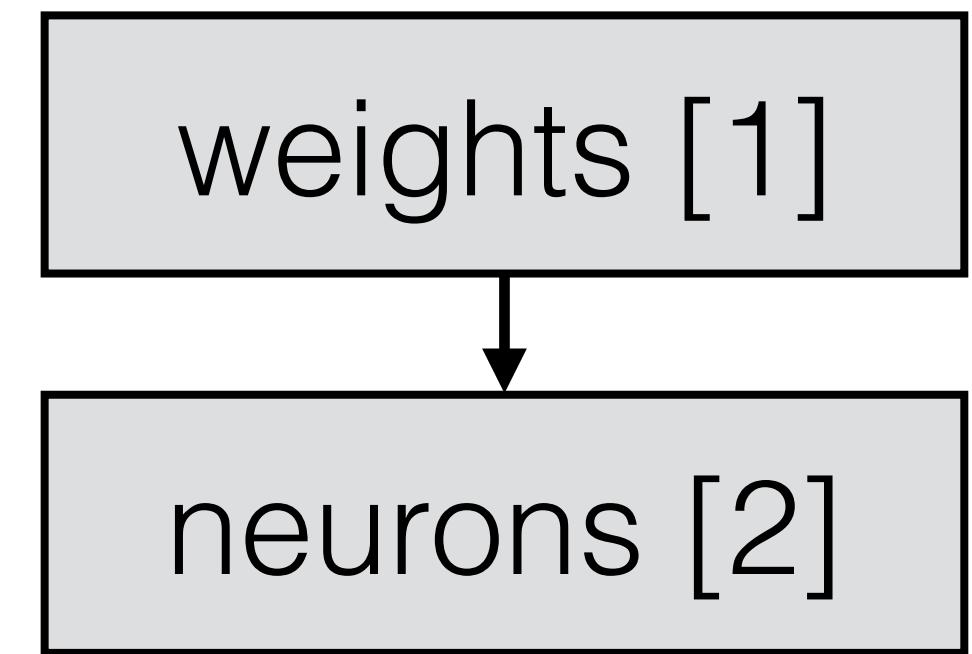
[1] Narang et al. Exploring sparsity in recurrent neural networks. ICLR 2017

[2] Wen, et al. Learning intrinsic sparse structures within long short-term memory. ICLR 2018

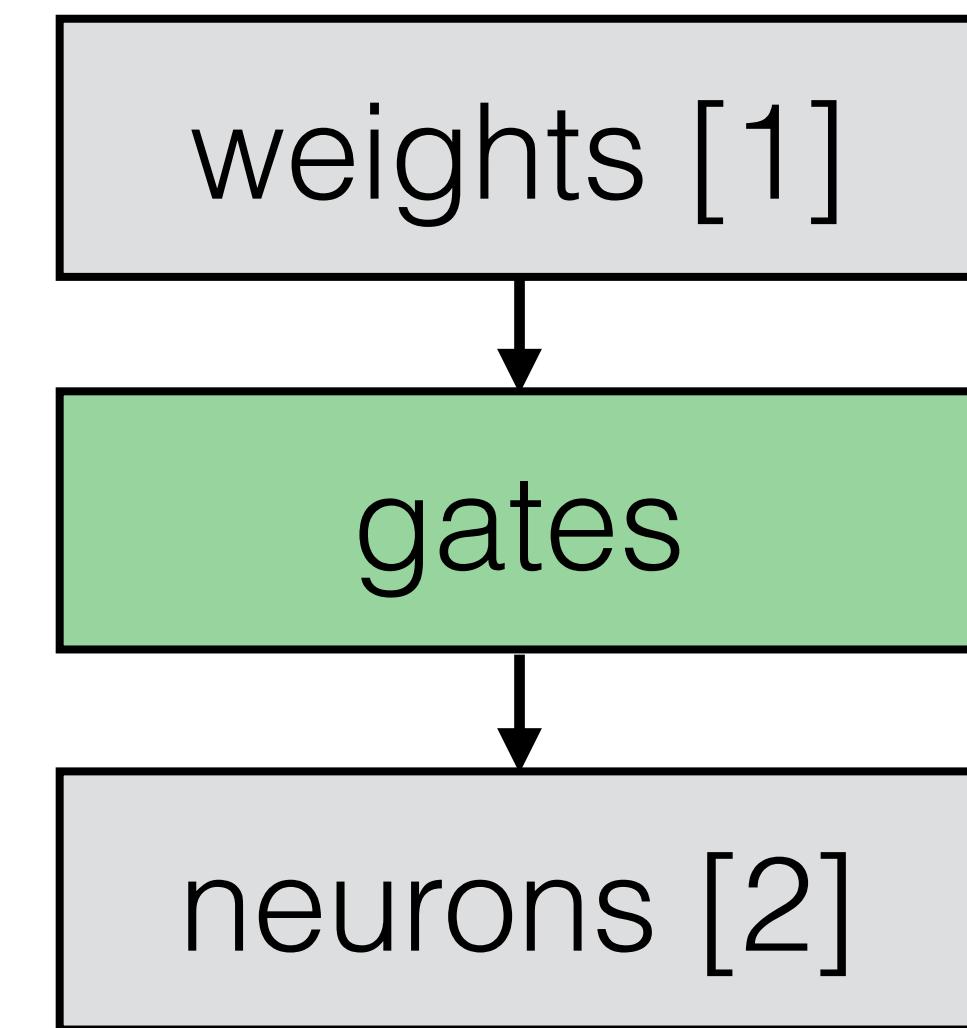
# Problem setting

Goal: learn structurally sparse gated RNNs (LSTM, GRU...)

Which structure units to sparsify  
in vanilla RNN?



in gated RNN?



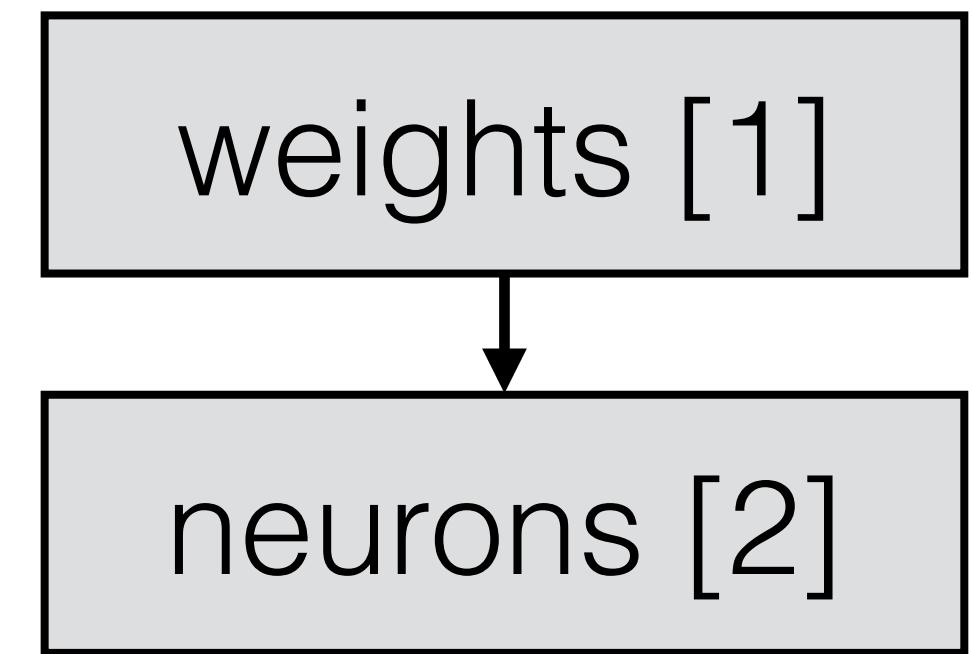
[1] Narang et al. Exploring sparsity in recurrent neural networks. ICLR 2017

[2] Wen, et al. Learning intrinsic sparse structures within long short-term memory. ICLR 2018

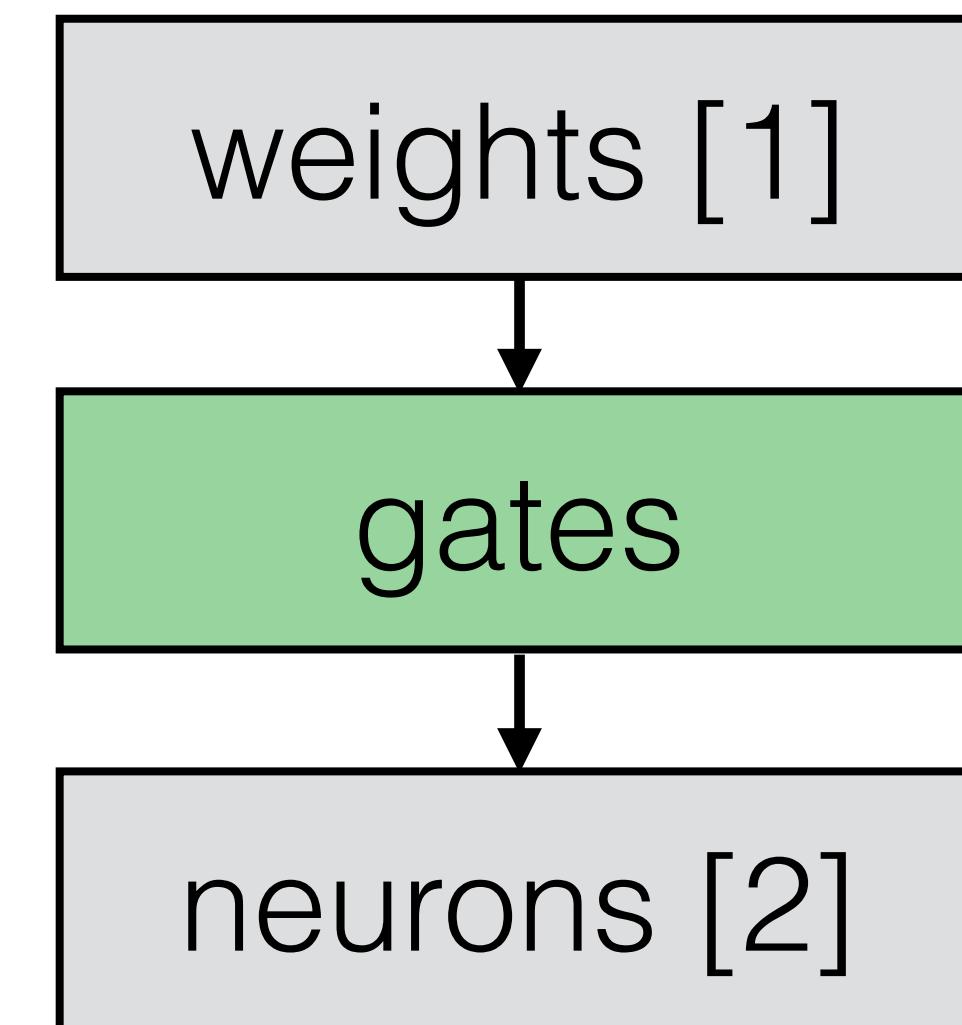
# Problem setting

Goal: learn structurally sparse gated RNNs (LSTM, GRU...)

Which structure units to sparsify  
in vanilla RNN?



in gated RNN?



Main idea: introduce intermediate level of sparsification — gates

[1] Narang et al. Exploring sparsity in recurrent neural networks. ICLR 2017

[2] Wen, et al. Learning intrinsic sparse structures within long short-term memory. ICLR 2018

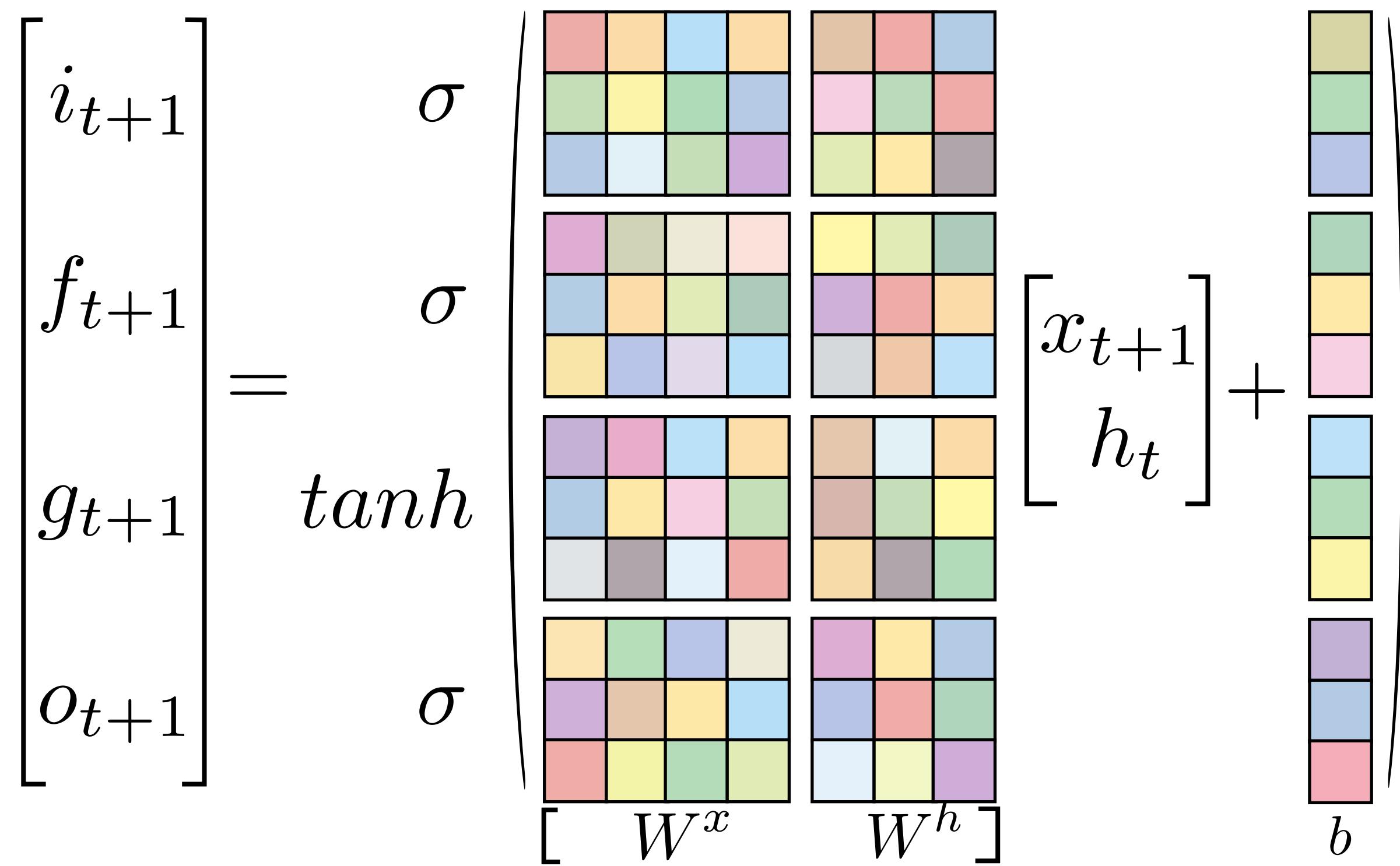
# Plan of the talk

1. **Idea:** introduce intermediate level of sparsification in RNN — gates
2. **Implementation** of the idea in Bayesian sparsification framework
3. **Results:** better compression and interpretable resulting gates structure

# Plan of the talk

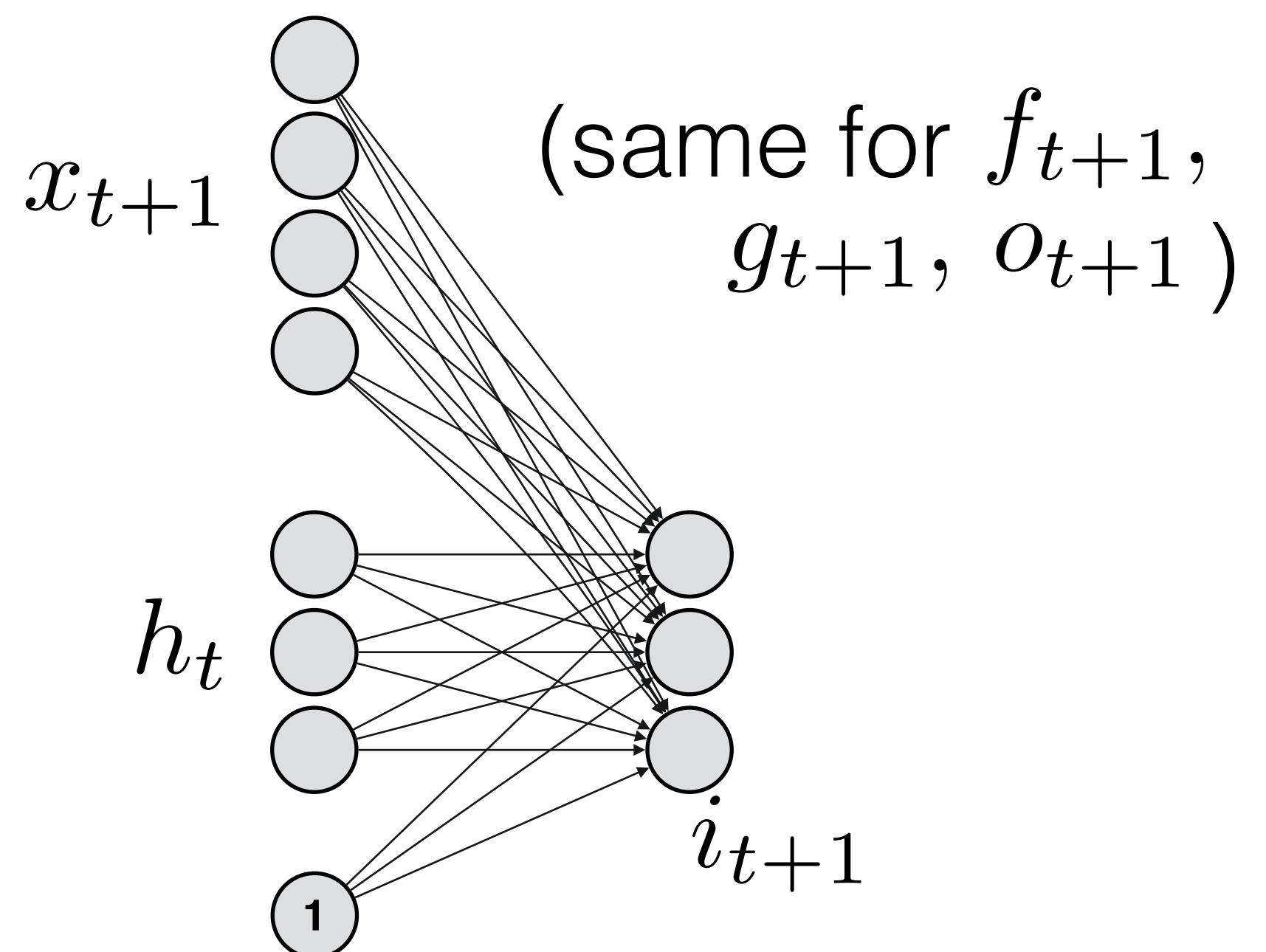
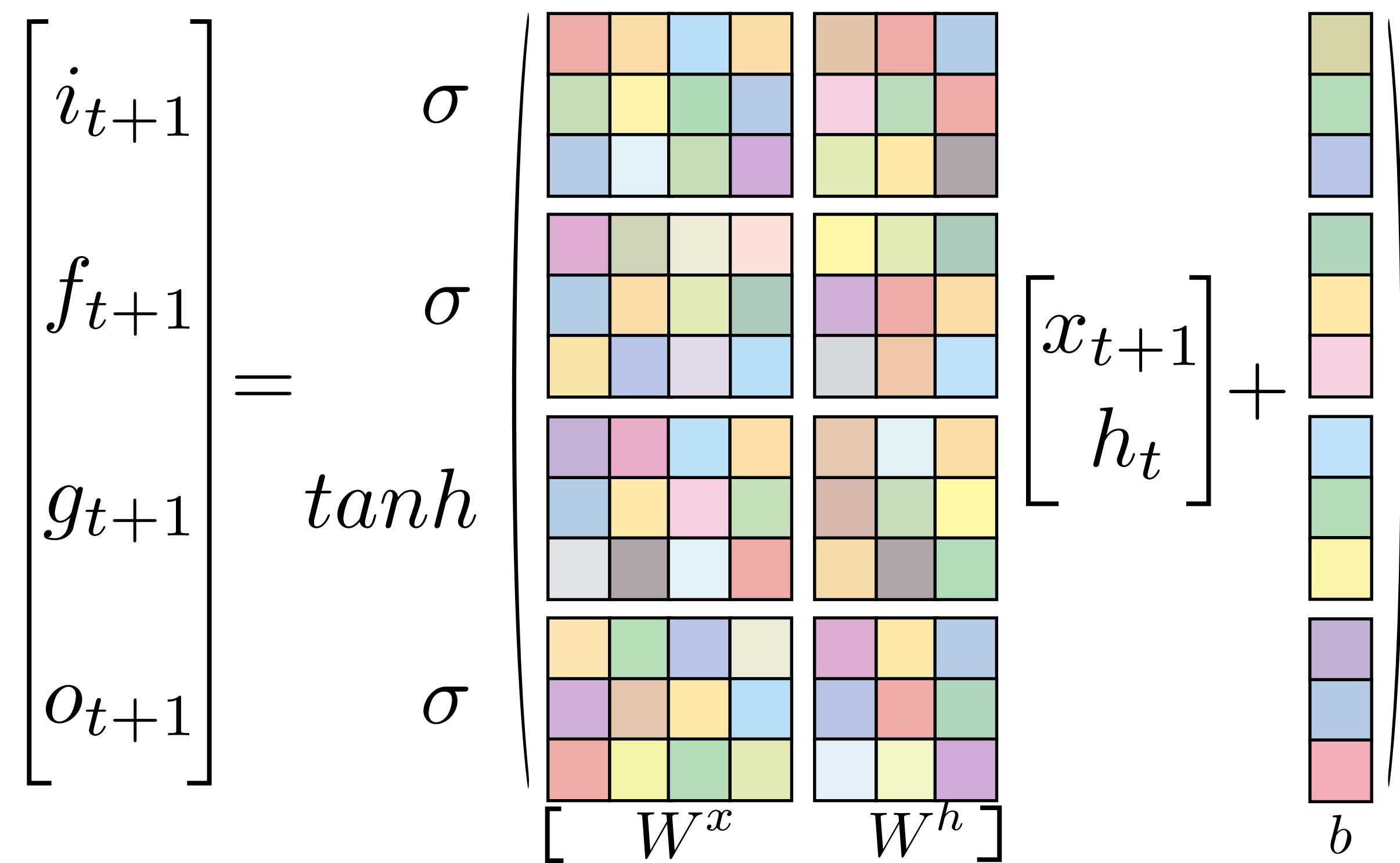
1. **Idea:** introduce intermediate level of sparsification in RNN — gates
2. **Implementation** of the idea in Bayesian sparsification framework
3. **Results:** better compression and interpretable resulting gates structure

# LSTM



$$c_{t+1} = f_{t+1} \odot c_t + i_{t+1} \odot g_{t+1}$$
$$h_{t+1} = o_{t+1} \odot g_{t+1}$$

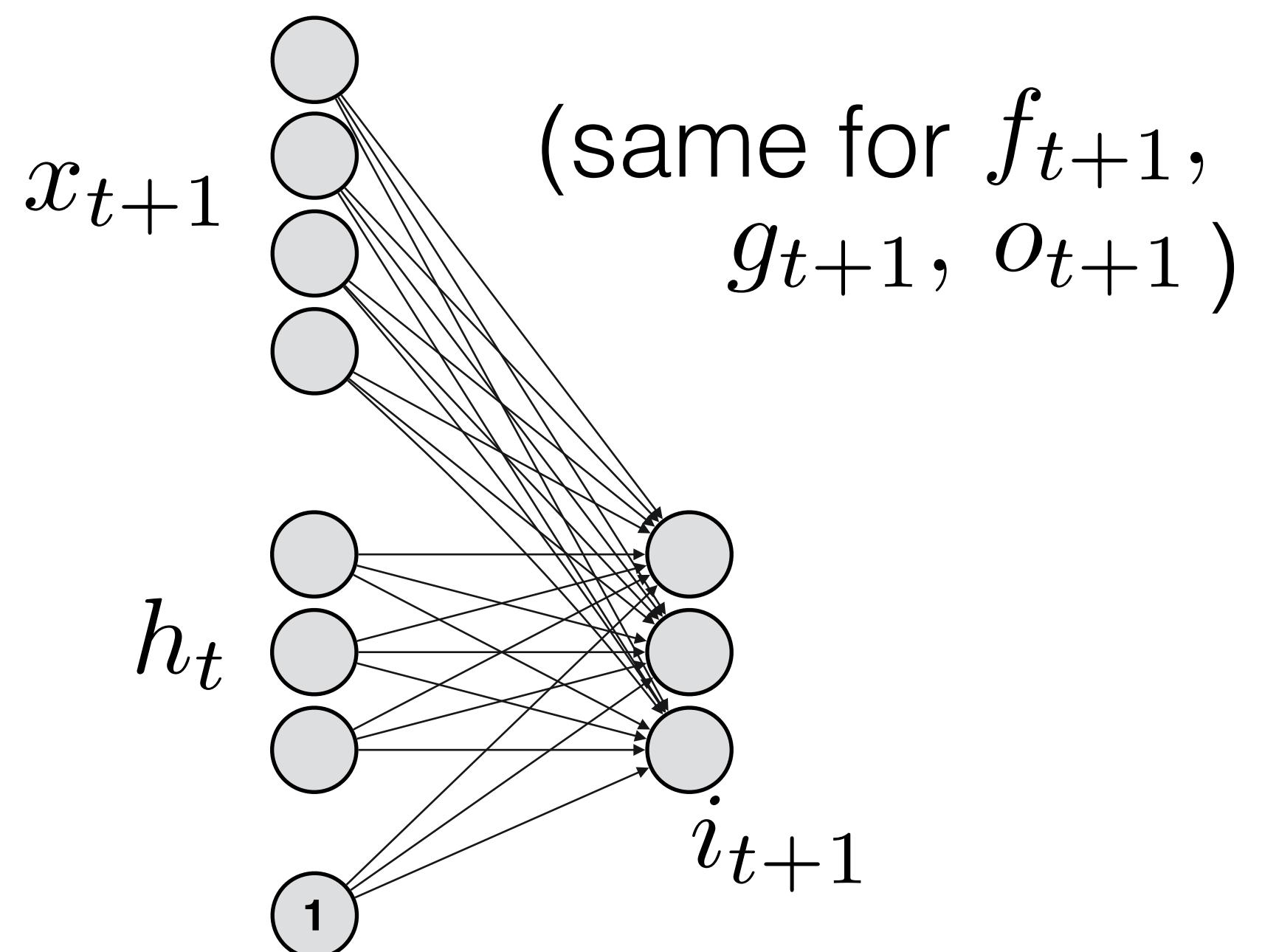
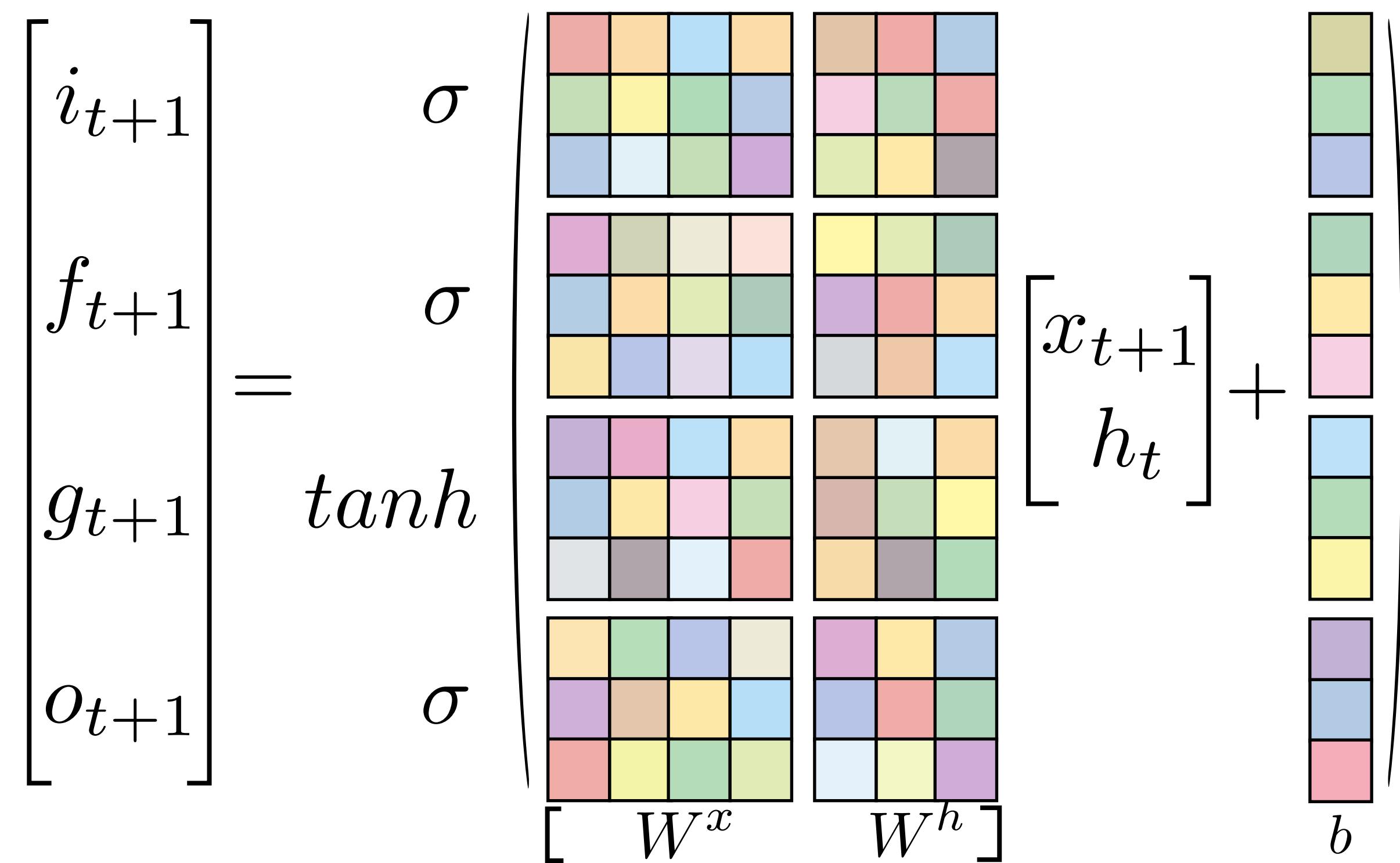
# LSTM



$$c_{t+1} = f_{t+1} \odot c_t + i_{t+1} \odot g_{t+1}$$

$$h_{t+1} = o_{t+1} \odot g_{t+1}$$

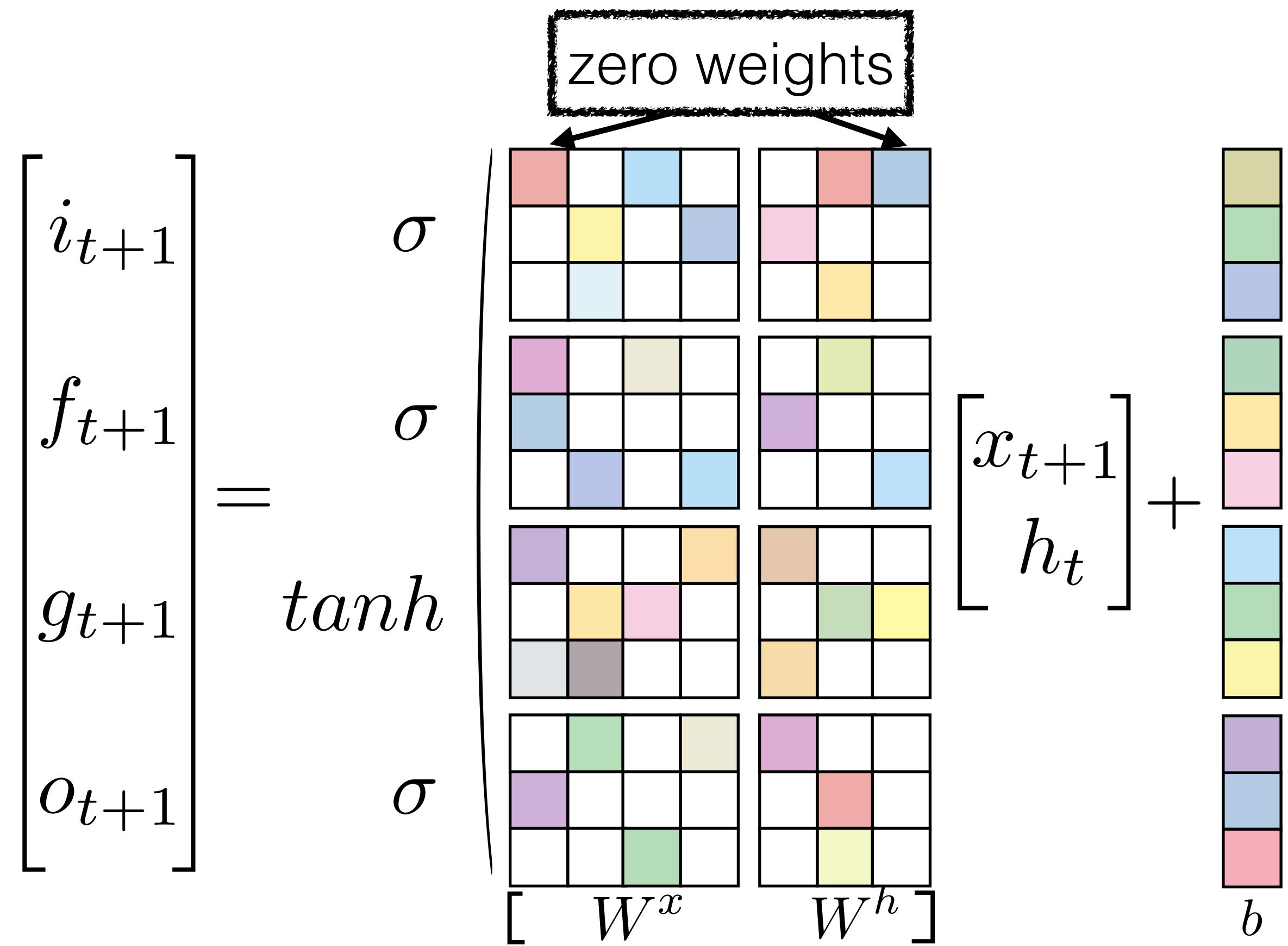
# LSTM



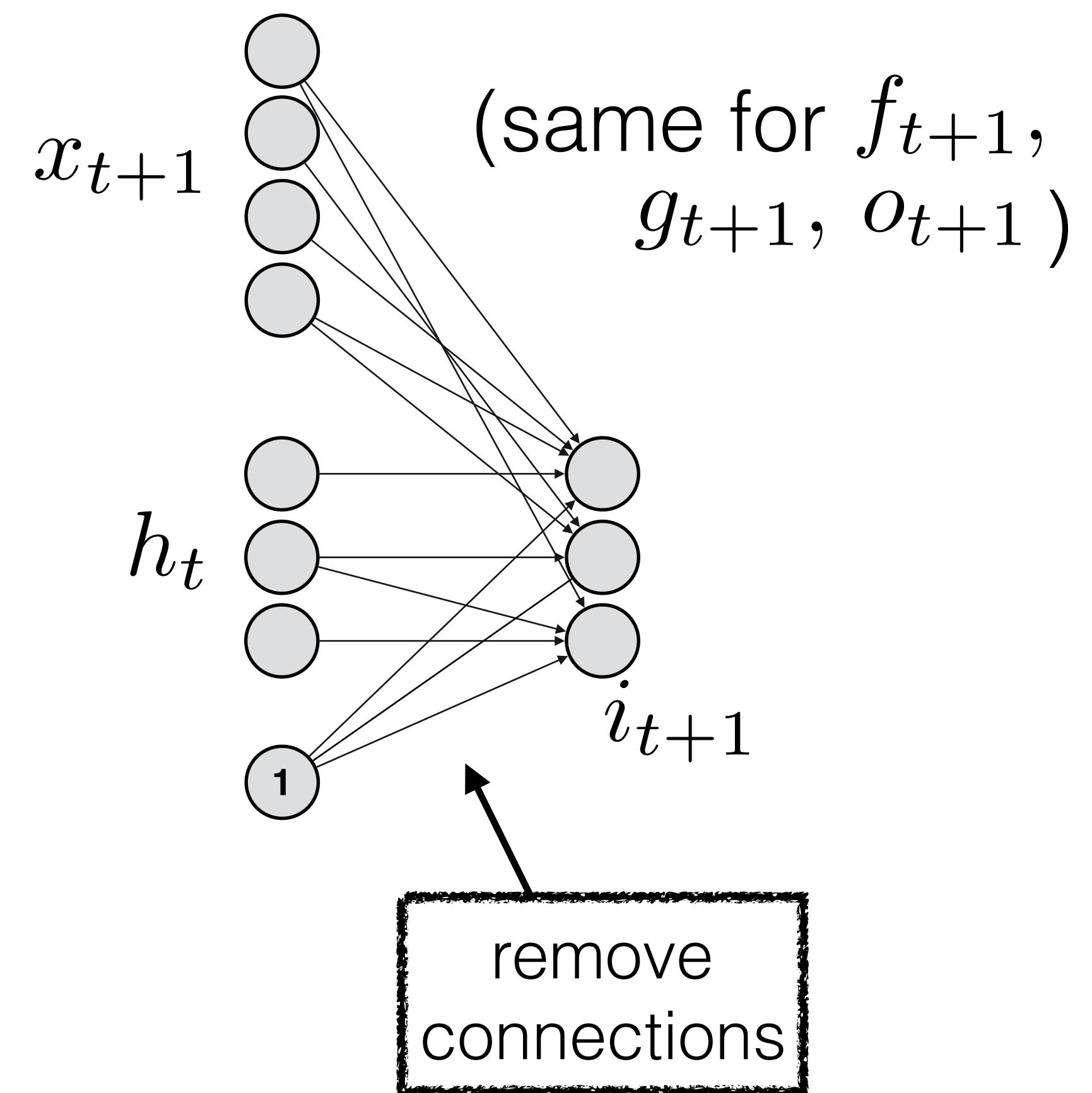
$$c_{t+1} = f_{t+1} \odot c_t + i_{t+1} \odot g_{t+1}$$

$$h_{t+1} = o_{t+1} \odot g_{t+1}$$

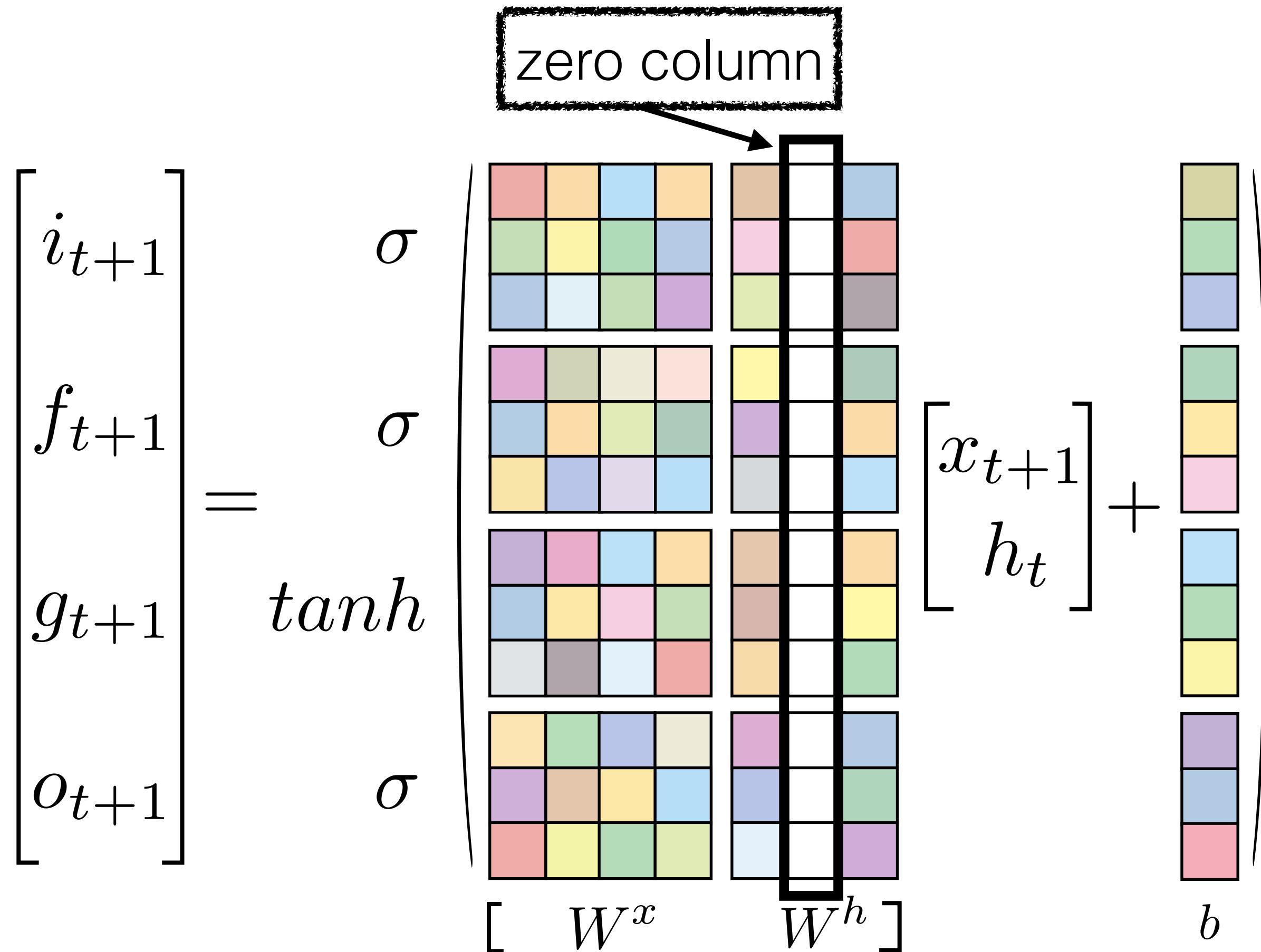
# LSTM: remove individual connections



$$c_{t+1} = f_{t+1} \odot c_t + i_{t+1} \odot g_{t+1}$$
$$h_{t+1} = o_{t+1} \odot g_{t+1}$$

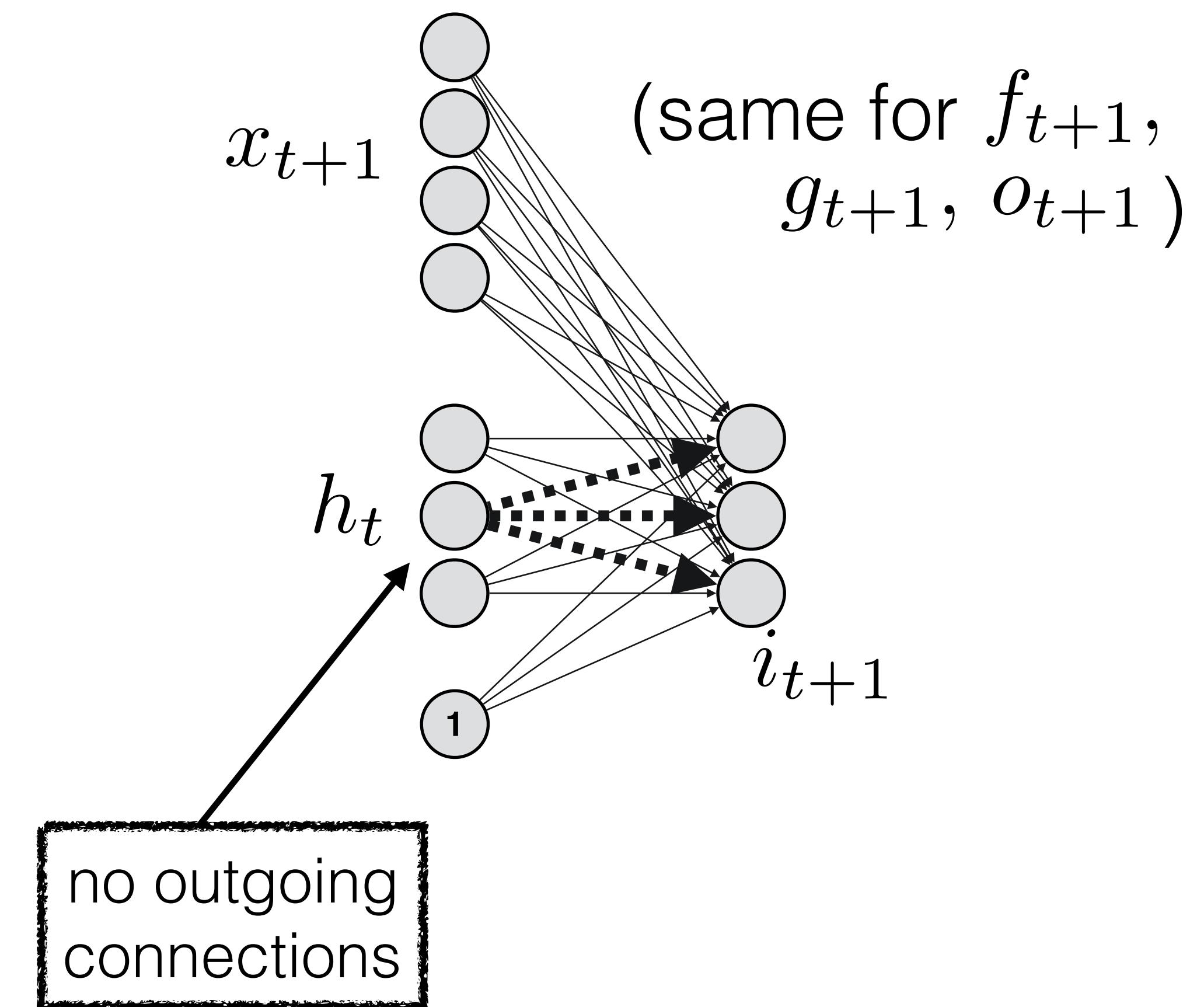


# LSTM: remove neurons

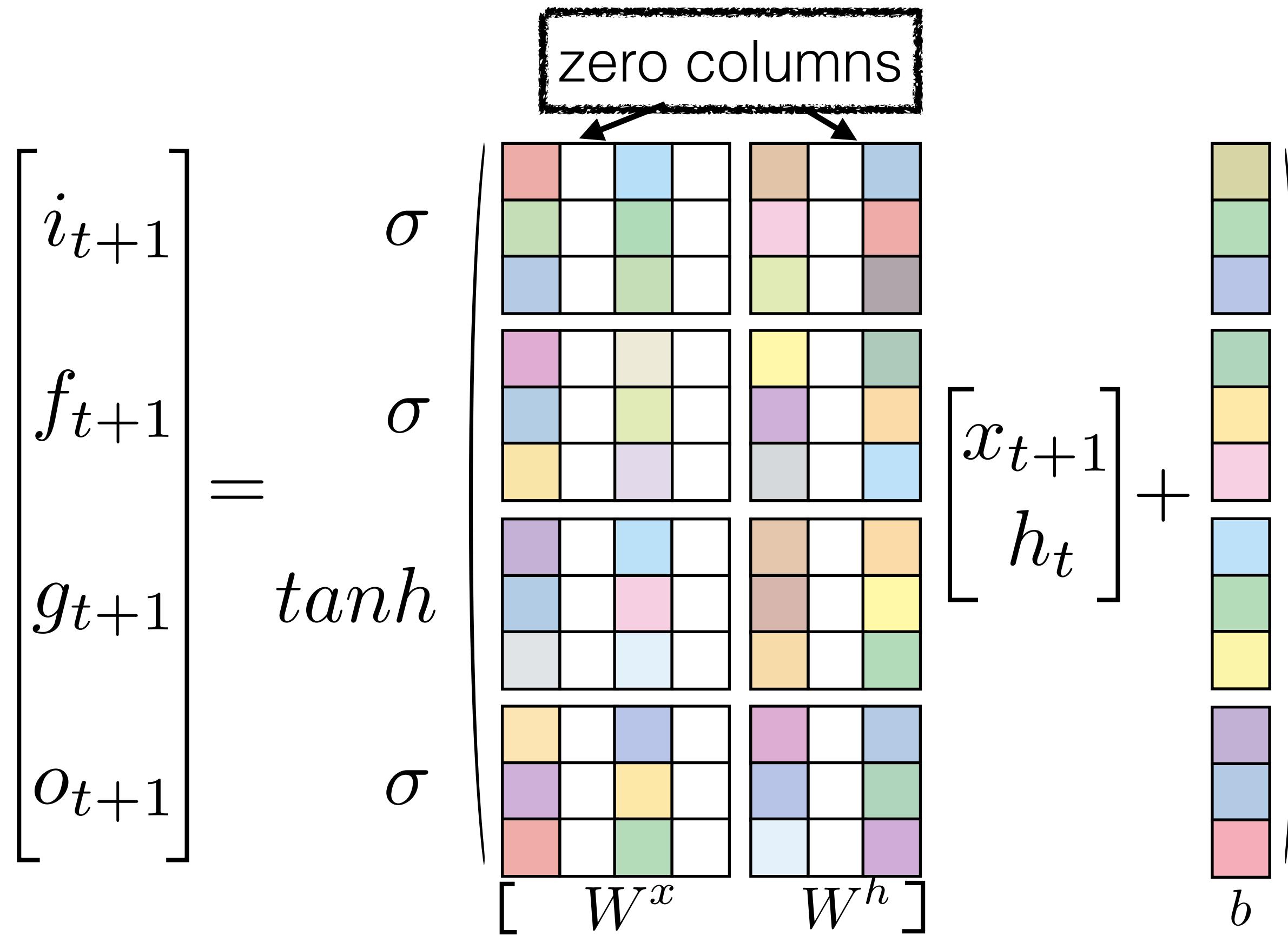


$$c_{t+1} = f_{t+1} \odot c_t + i_{t+1} \odot g_{t+1}$$

$$h_{t+1} = o_{t+1} \odot g_{t+1}$$

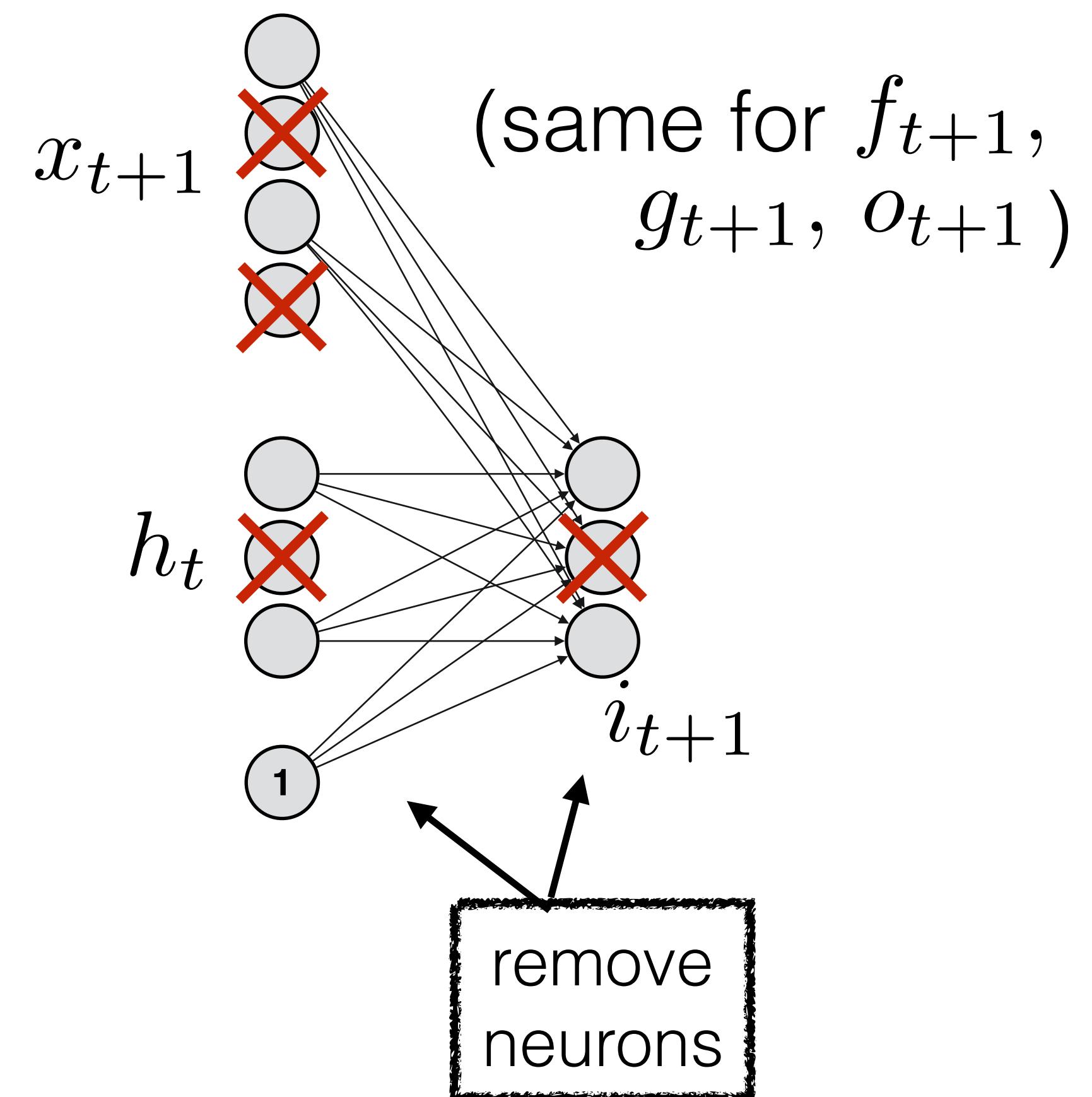


# LSTM: remove neurons

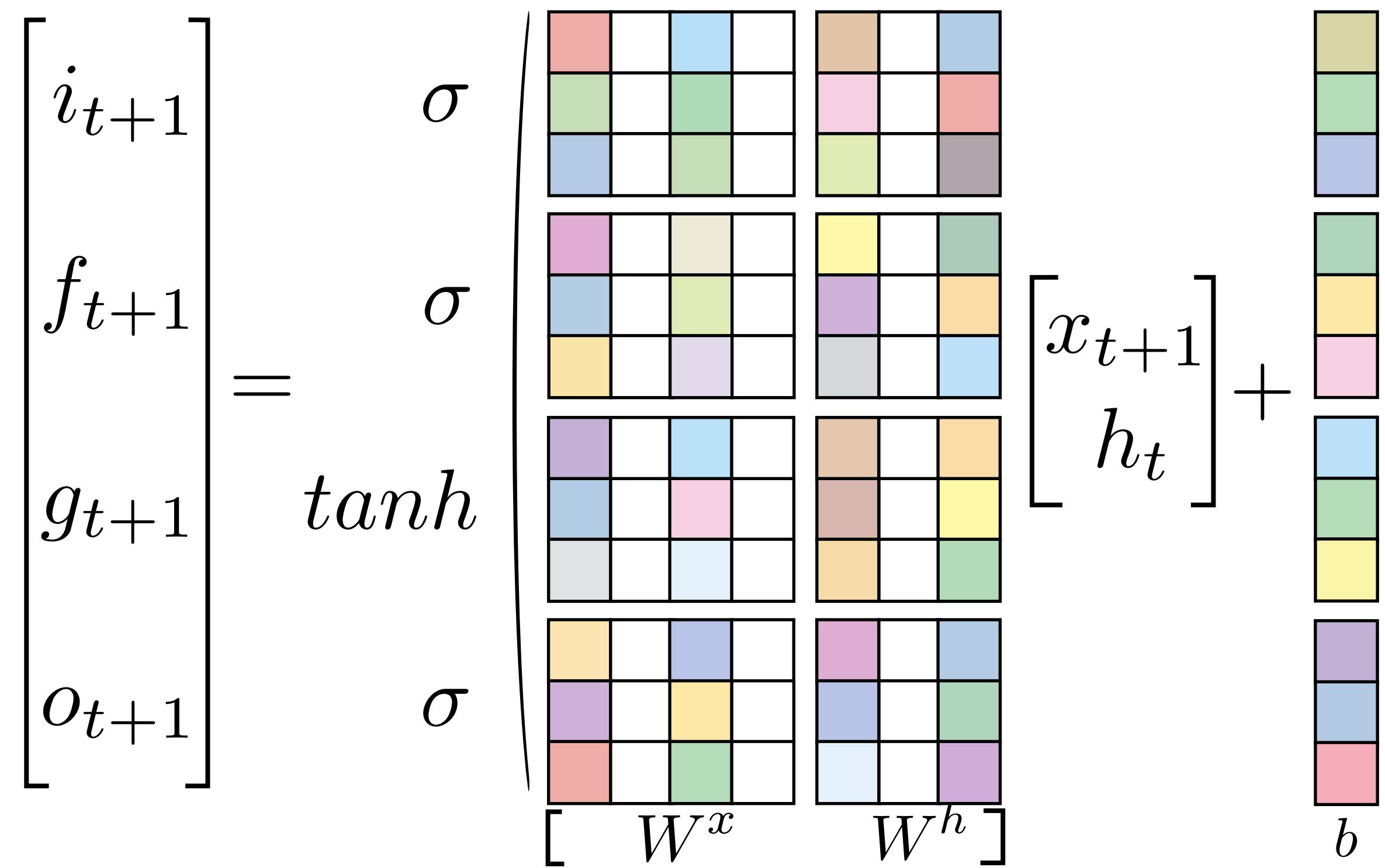


$$c_{t+1} = f_{t+1} \odot c_t + i_{t+1} \odot g_{t+1}$$

$$h_{t+1} = o_{t+1} \odot g_{t+1}$$

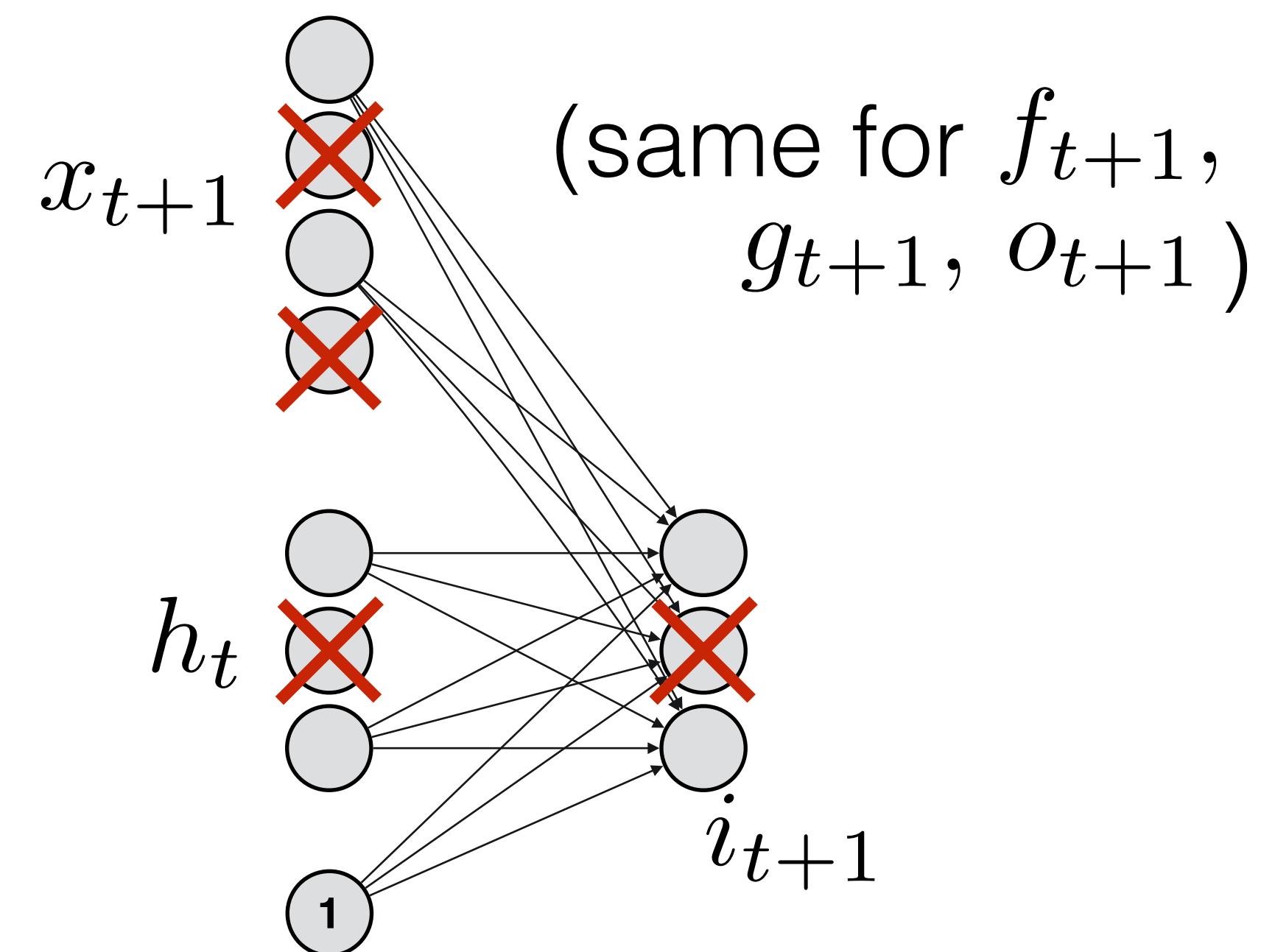


# LSTM: remove neurons



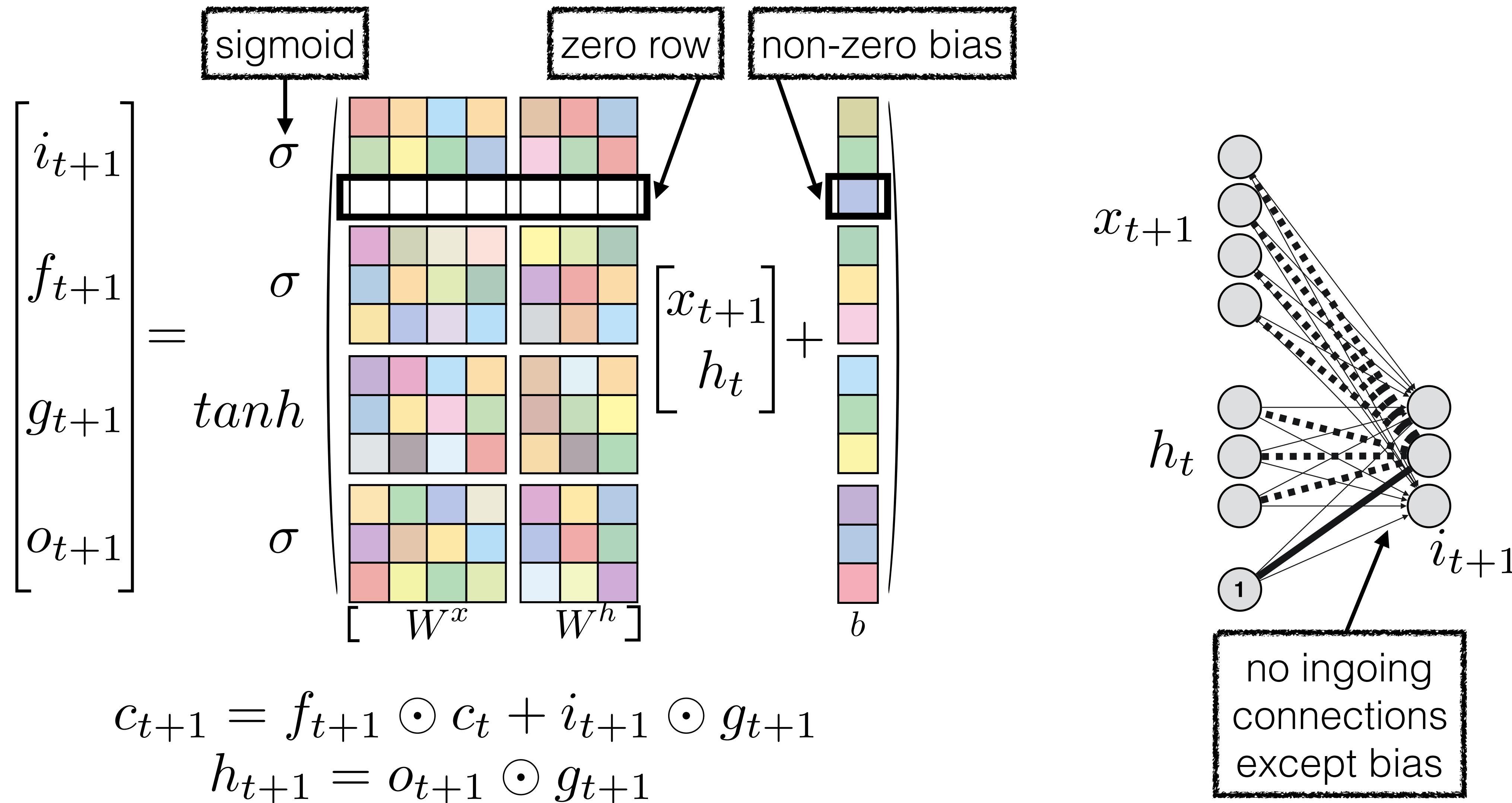
$$c_{t+1} = f_{t+1} \odot c_t + i_{t+1} \odot g_{t+1}$$

$$h_{t+1} = o_{t+1} \odot g_{t+1}$$

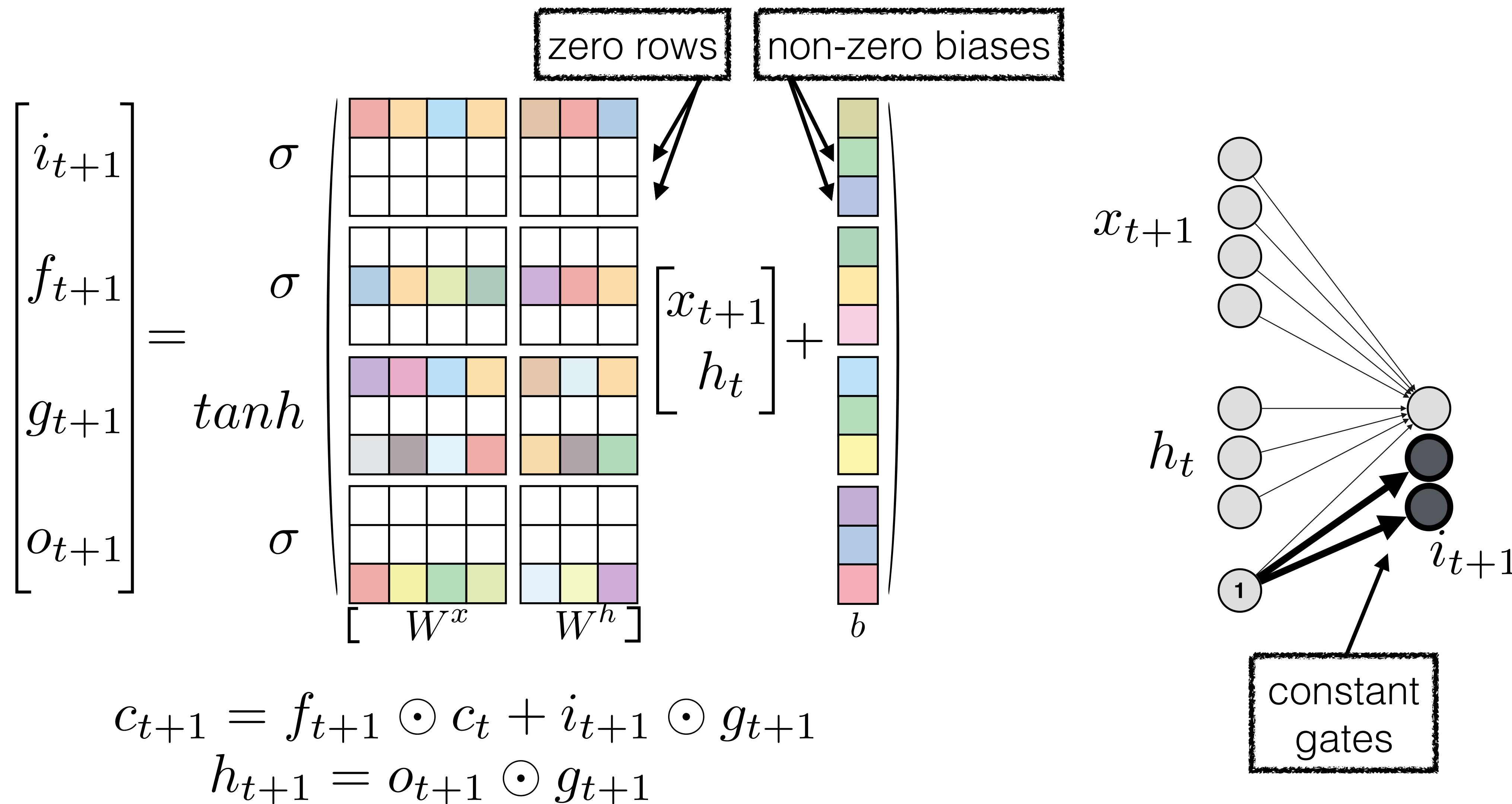


there are other ways to group weights to remove neurons

# LSTM: make gates constant

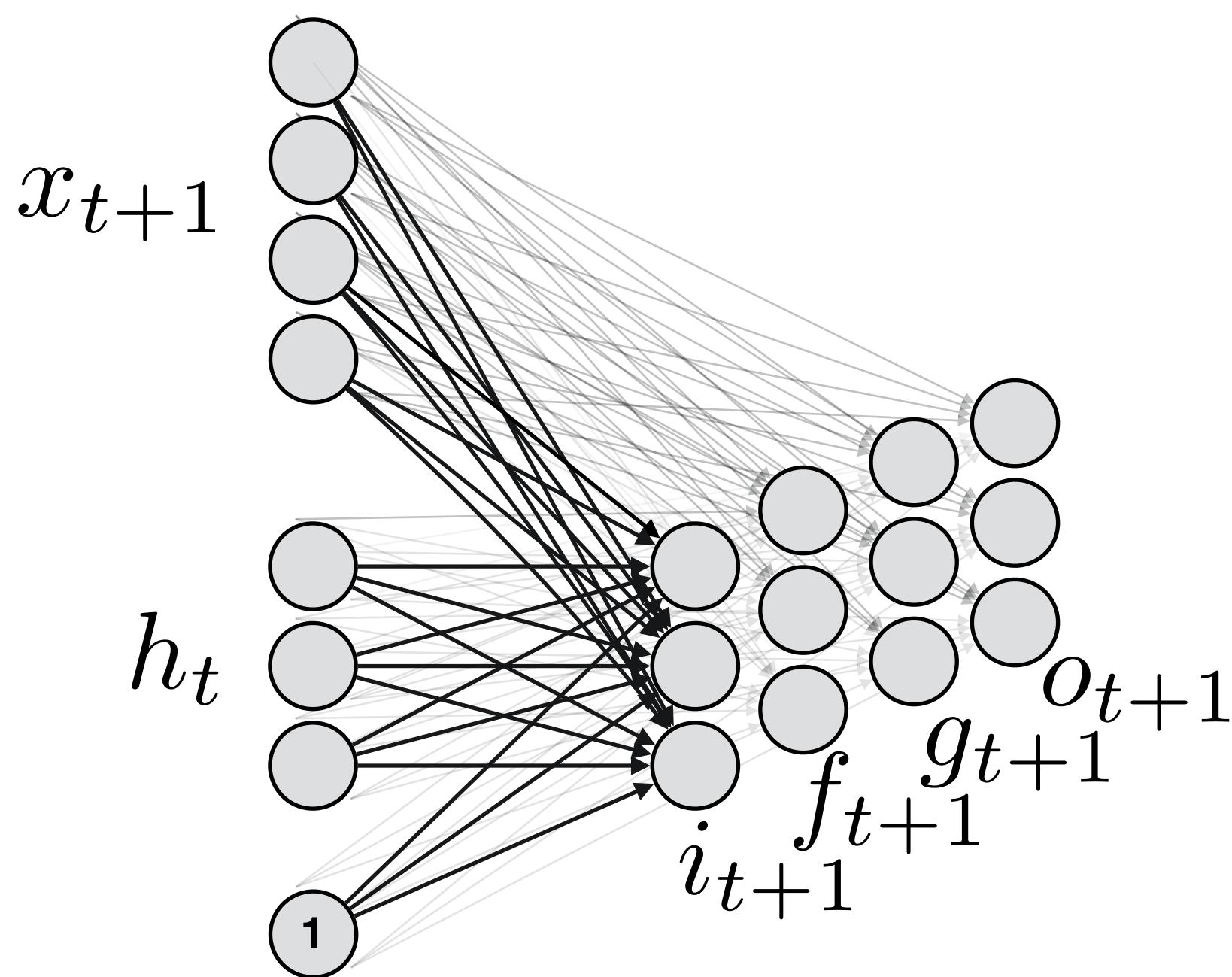


# LSTM: make gates constant



# Putting everything together

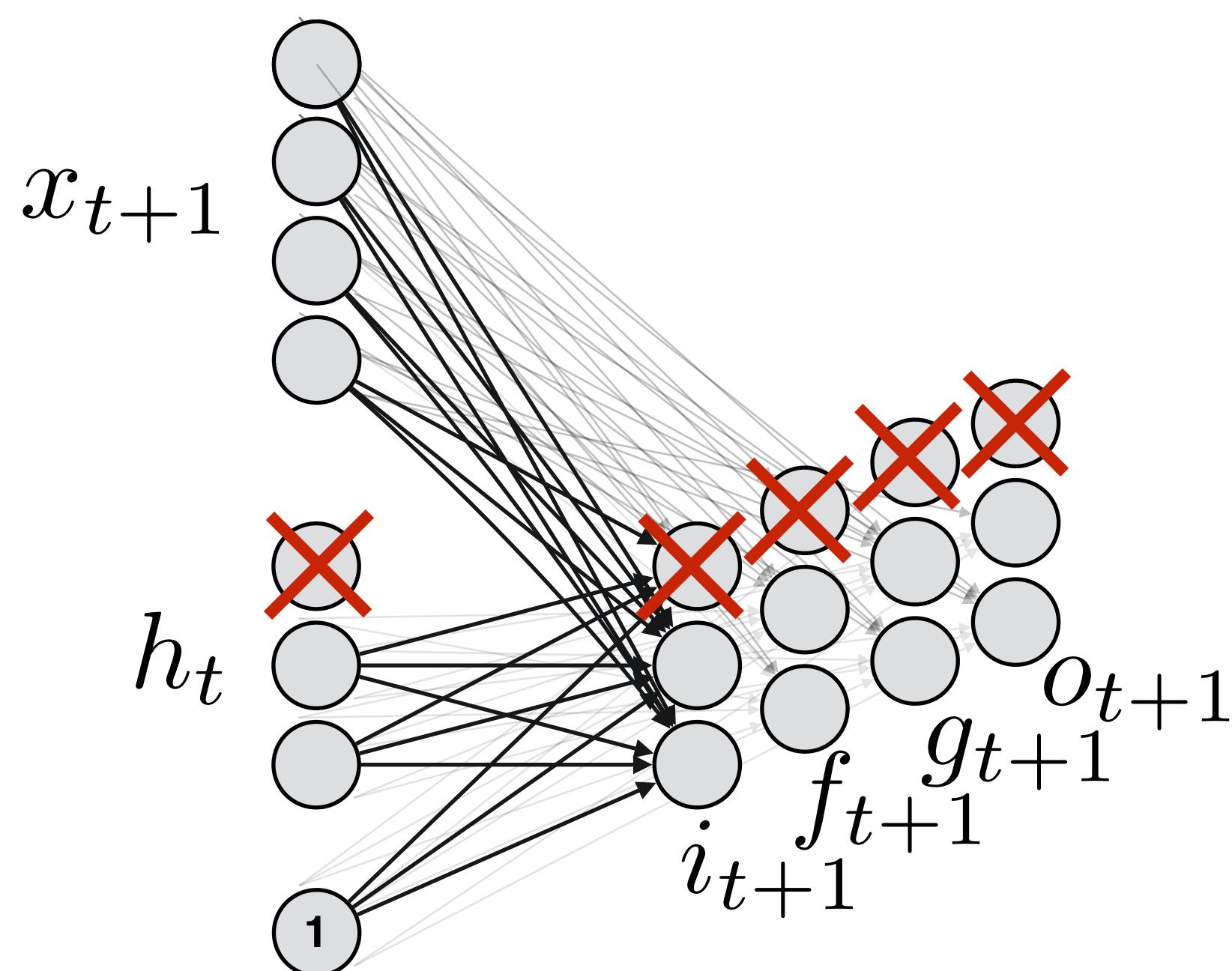
LSTM:



Levels of sparsification:

# Putting everything together

LSTM:

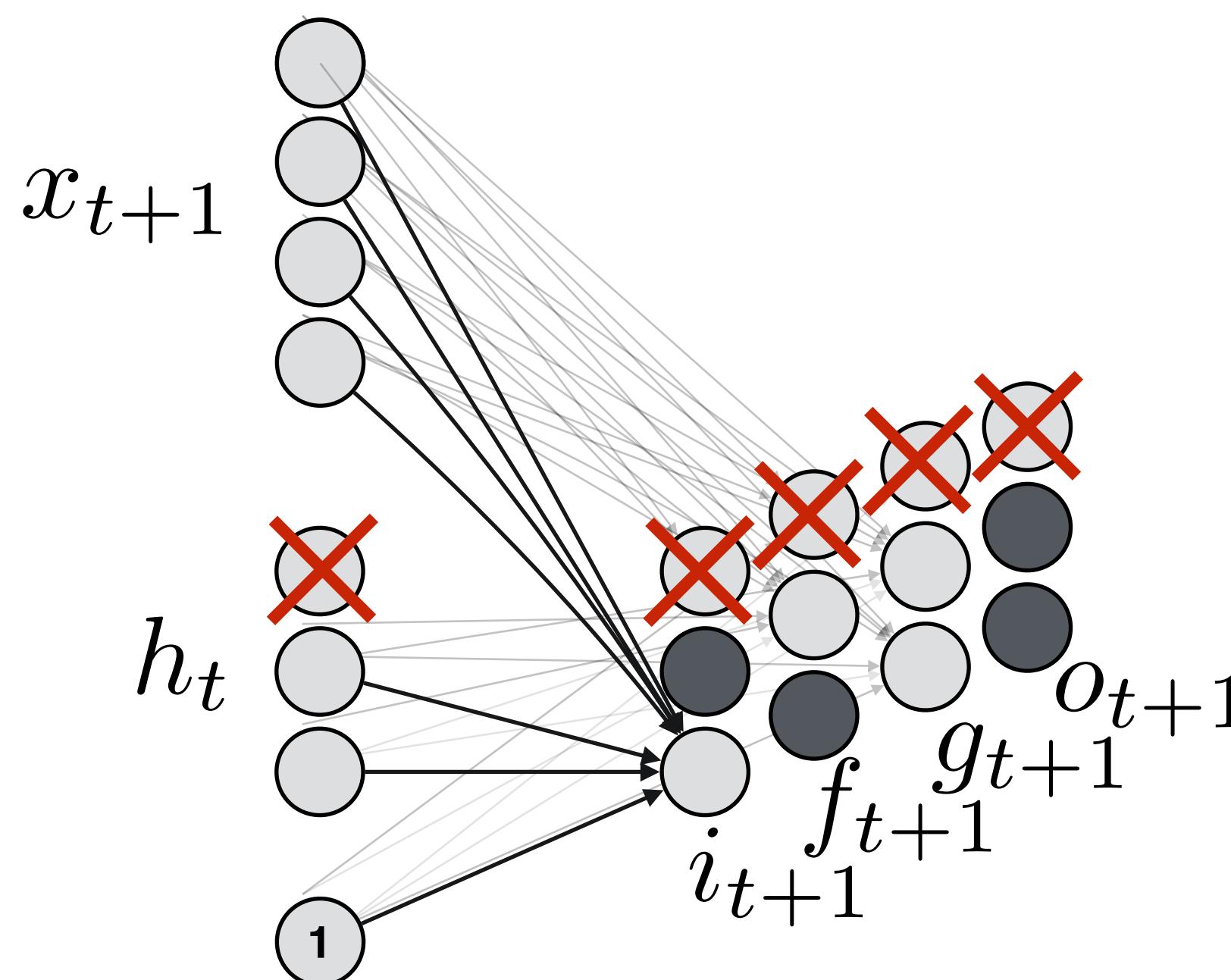


Levels of sparsification:

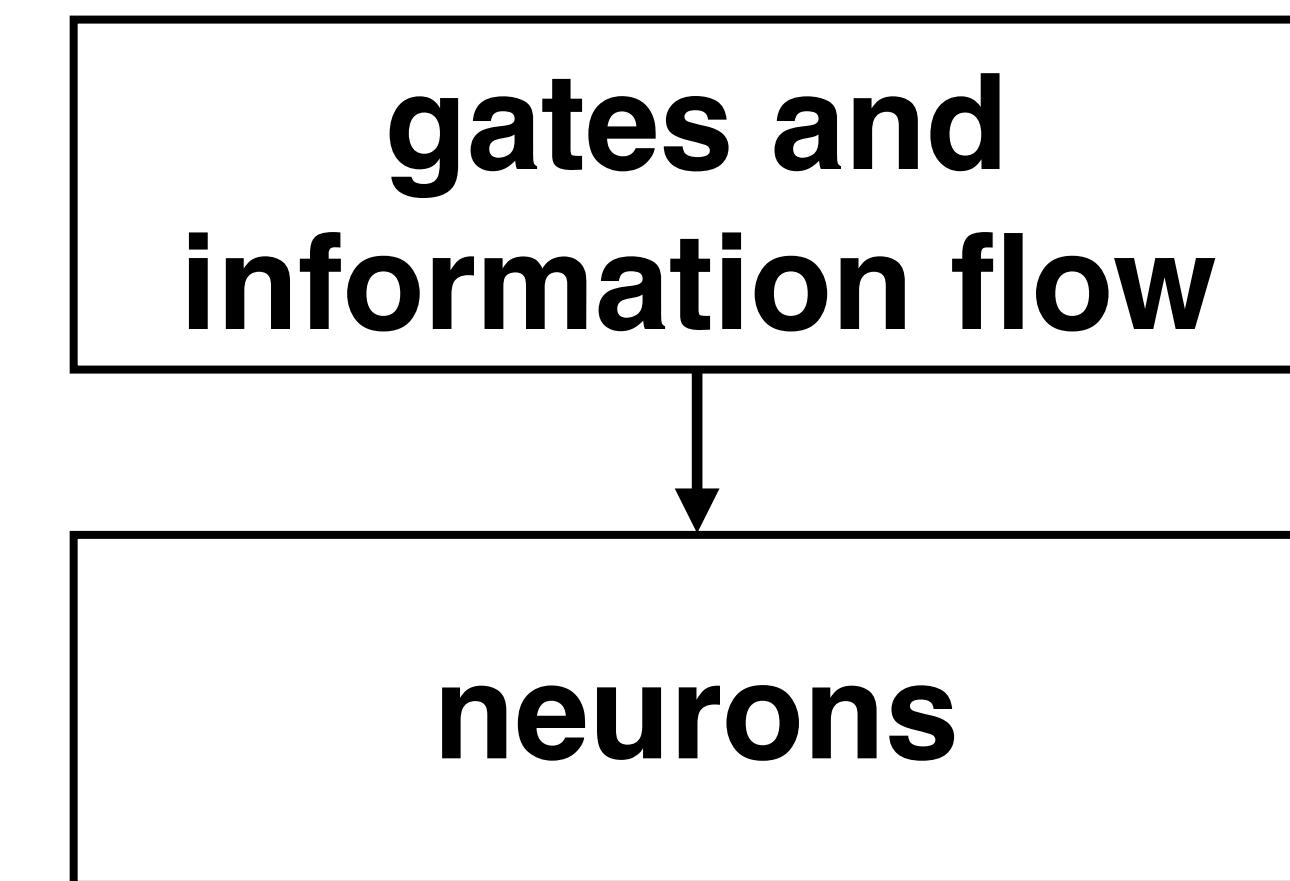


# Putting everything together

LSTM:

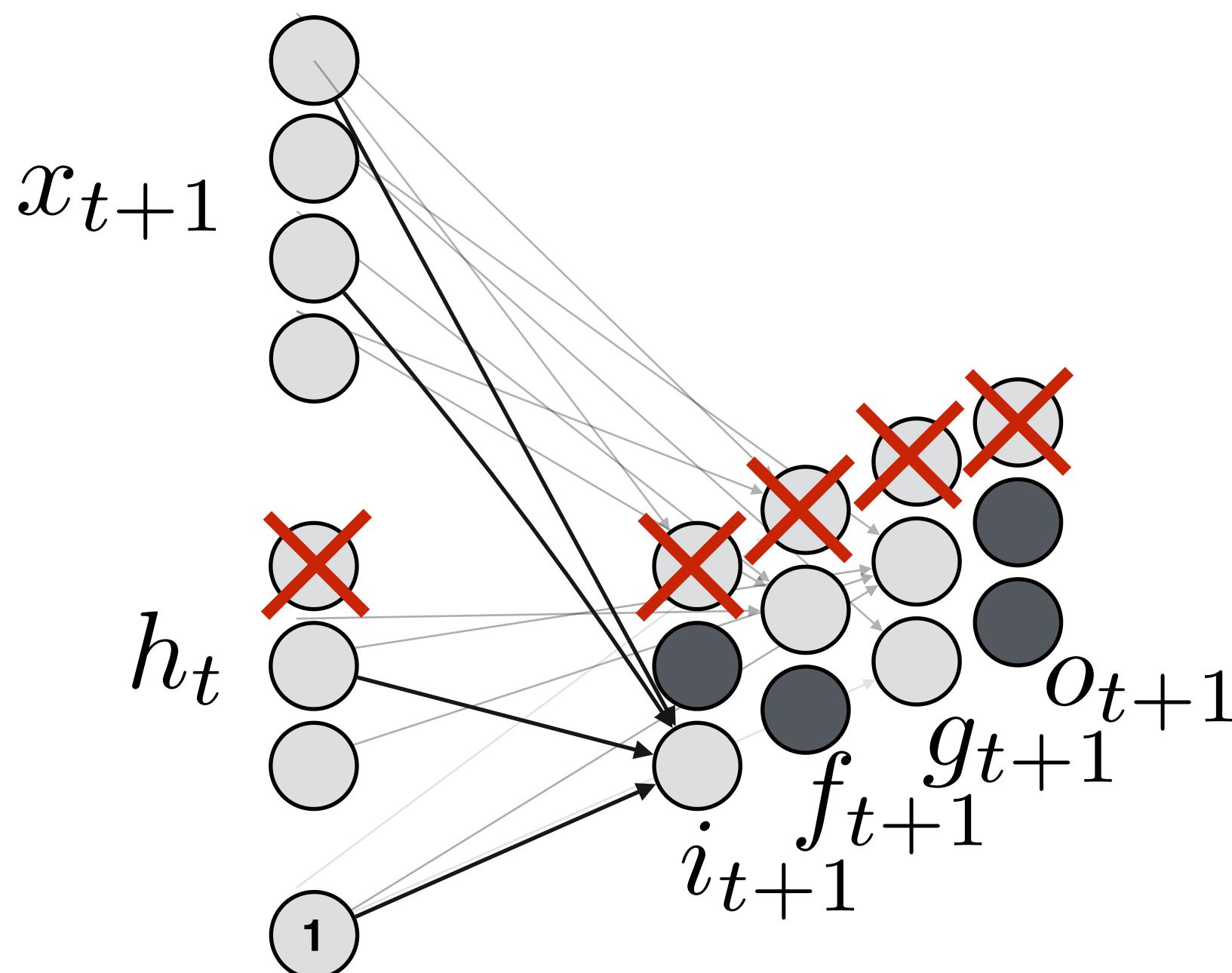


Levels of sparsification:

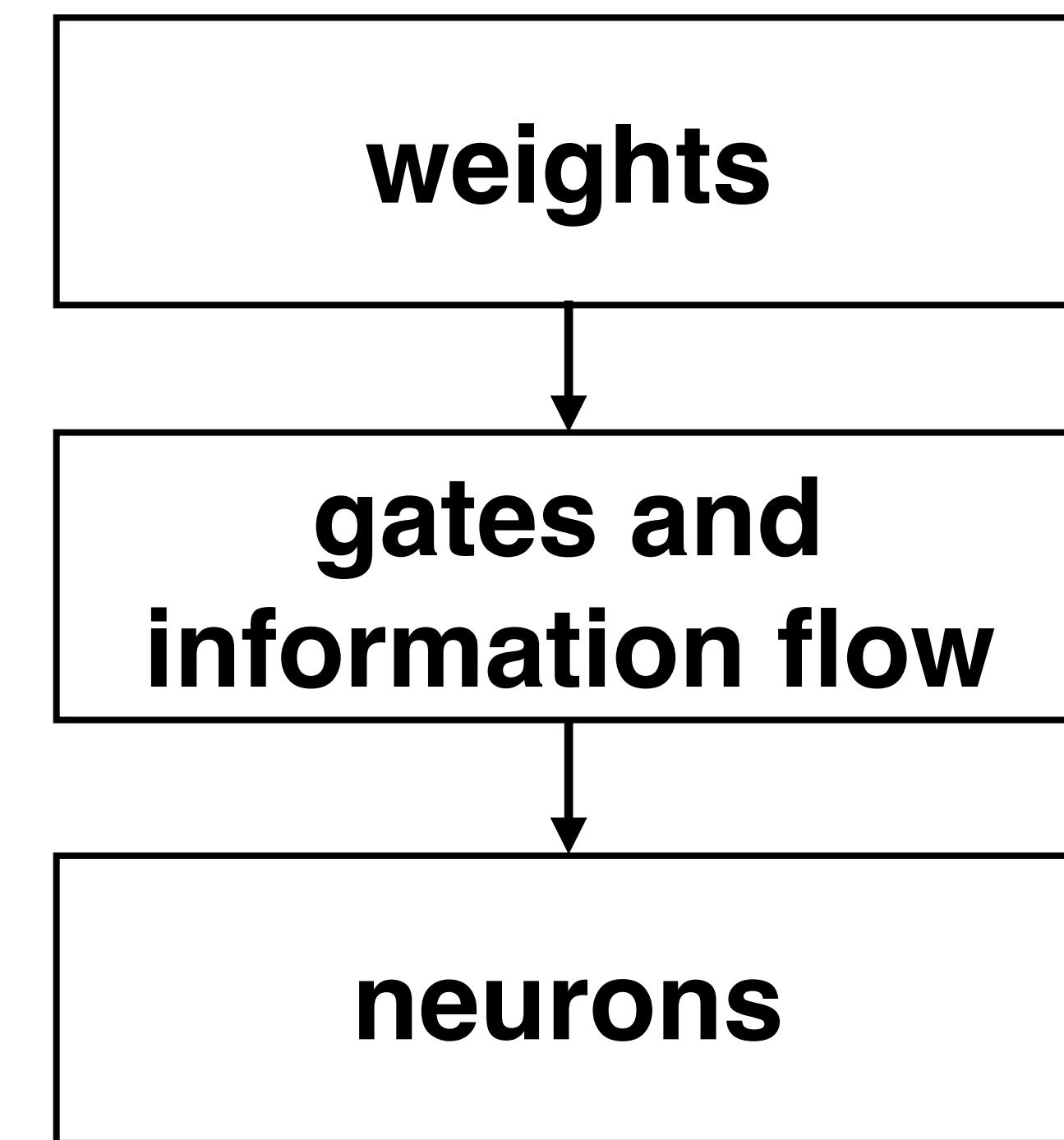


# Putting everything together

LSTM:



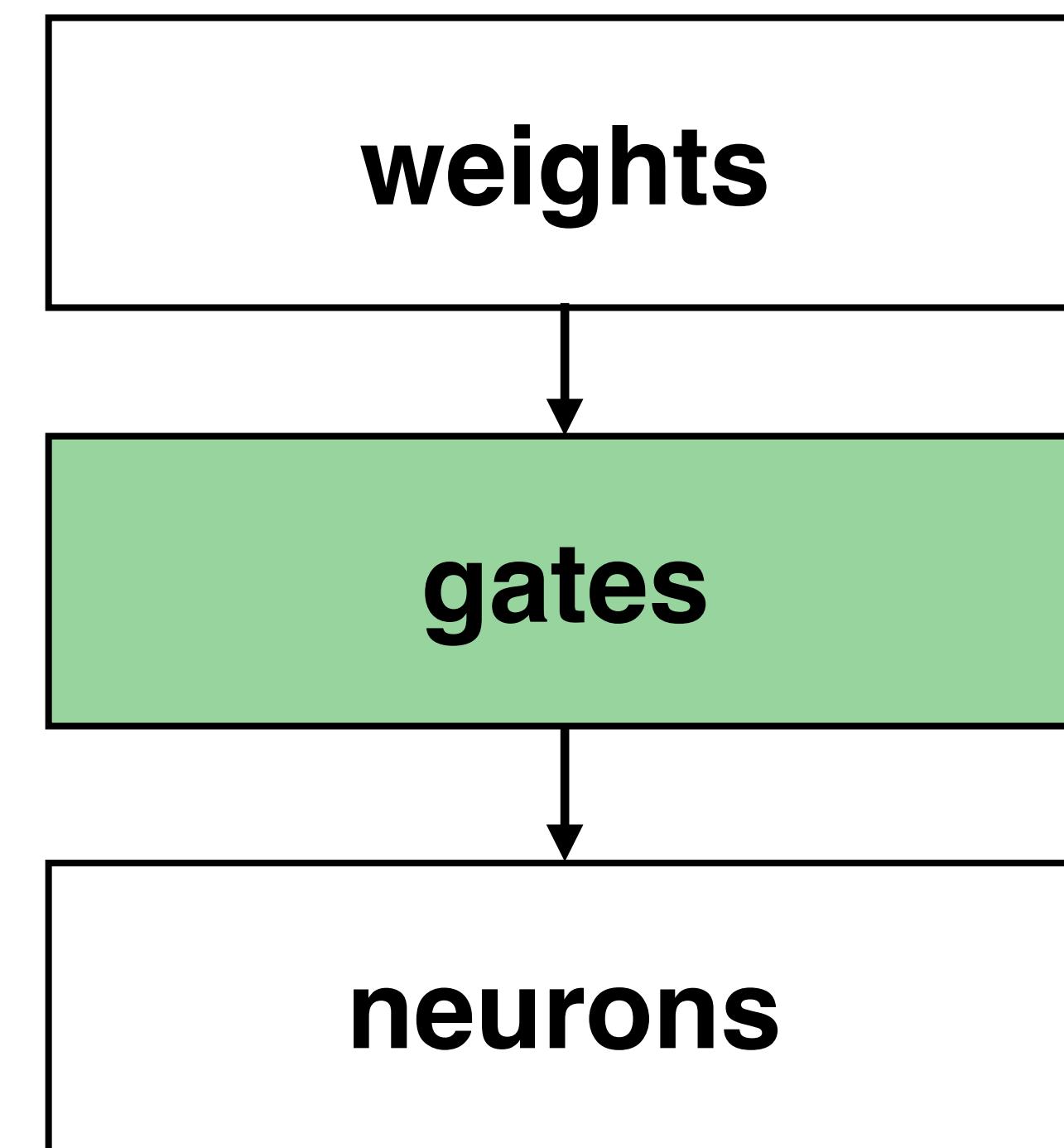
Levels of sparsification:



# Putting everything together

- Less computations on testing stage
- Better overall compression on testing stage
- Show that LSTM works in different manner on different tasks

Levels of sparsification:

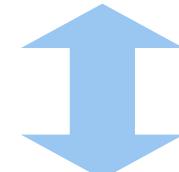


# Plan of the talk

1. **Idea:** introduce intermediate level of sparsification in RNN — gates
2. **Implementation** of the idea in Bayesian sparsification framework
3. **Results:** better compression and interpretable resulting gates structure

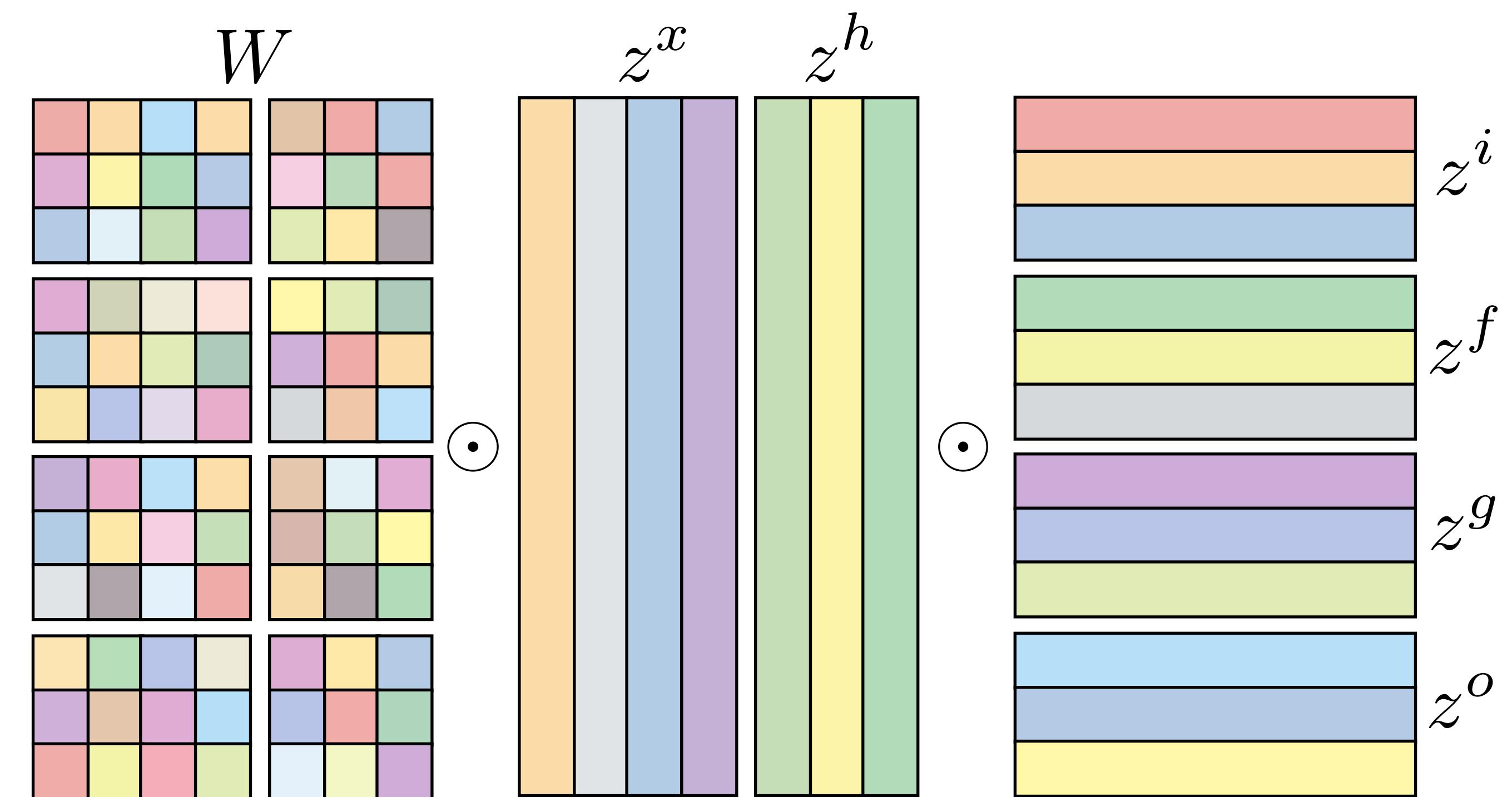
# Implementation: introducing group weights

$$f_{t+1} = \sigma \left( W_f^x (x_{t+1} \odot z^x) + (W_f^h (h_t \odot z^h)) \odot z^f + b_f \right)$$



$$\tilde{W}_{f,ij}^h = W_{f,ij}^h \cdot z_i^h \cdot z_j^f$$

$$\tilde{W}_{f,ij}^x = W_{f,ij}^x \cdot z_i^x \cdot z_j^x$$



# Implementation: introducing group weights

zero components → remove neuron      zero component → constant gate

similarly for  
 $i, g, o$

$$f_{t+1} = \sigma \left( W_f^x (x_{t+1} \odot z^x) + (W_f^h (h_t \odot z^h)) \odot z^f + b_f \right)$$

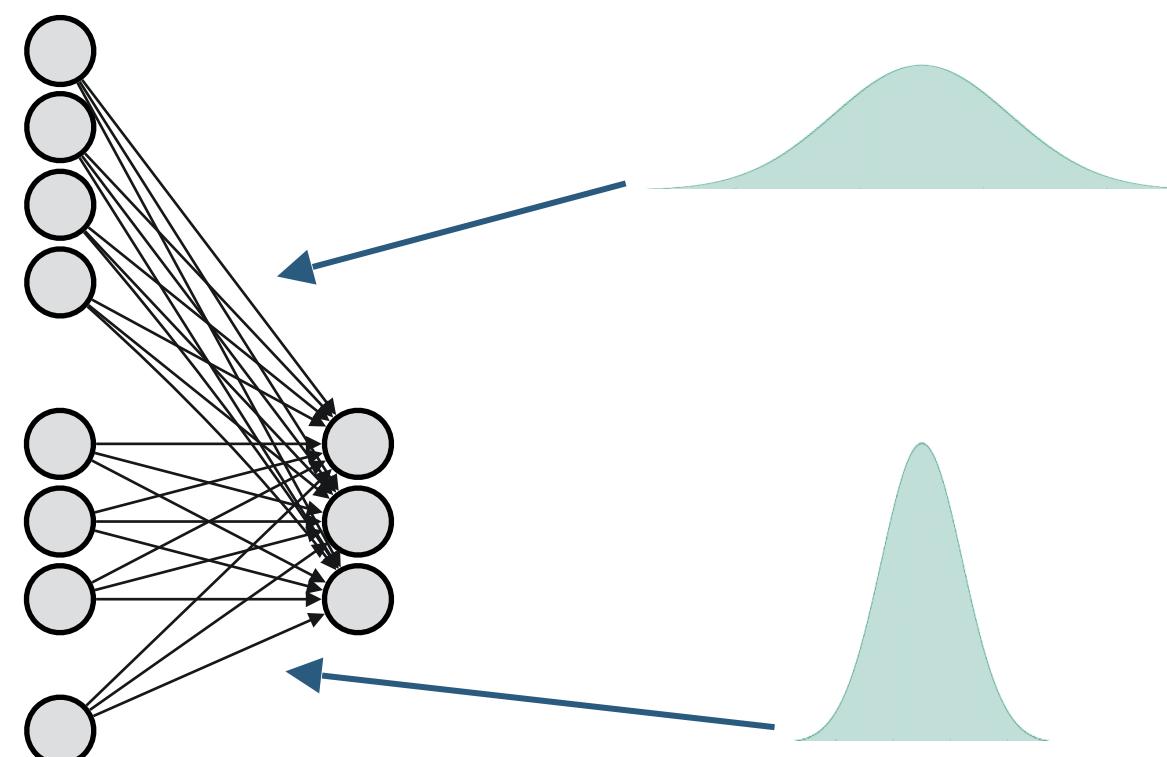
$\Updownarrow$

$$\tilde{W}_{f,ij}^h = W_{f,ij}^h \cdot z_i^h \cdot z_j^f$$

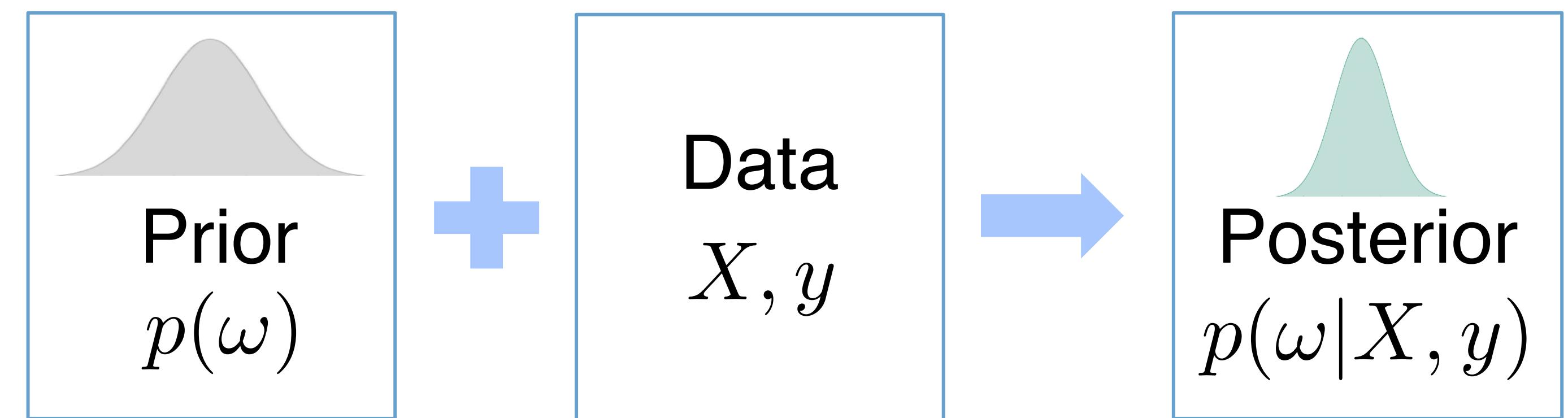
$$\tilde{W}_{f,ij}^x = W_{f,ij}^x \cdot z_i^x \cdot z_j^x$$

# Bayesian neural networks

Stochastic weights:

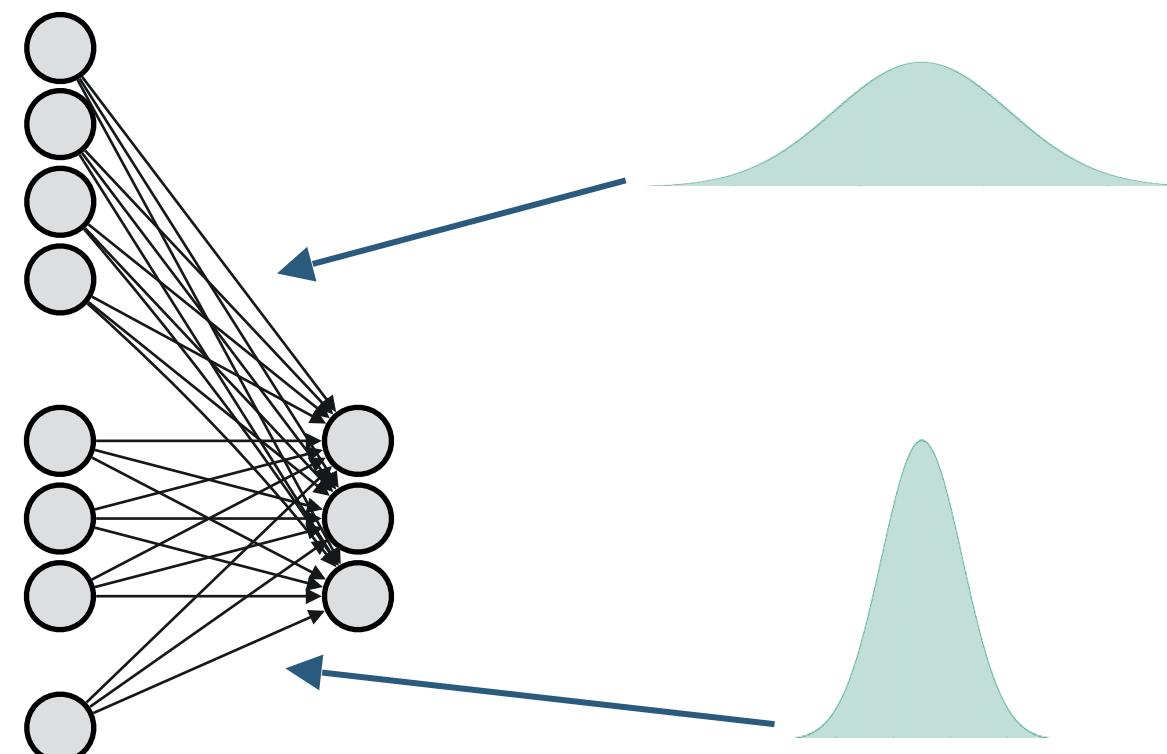


Bayesian Inference:

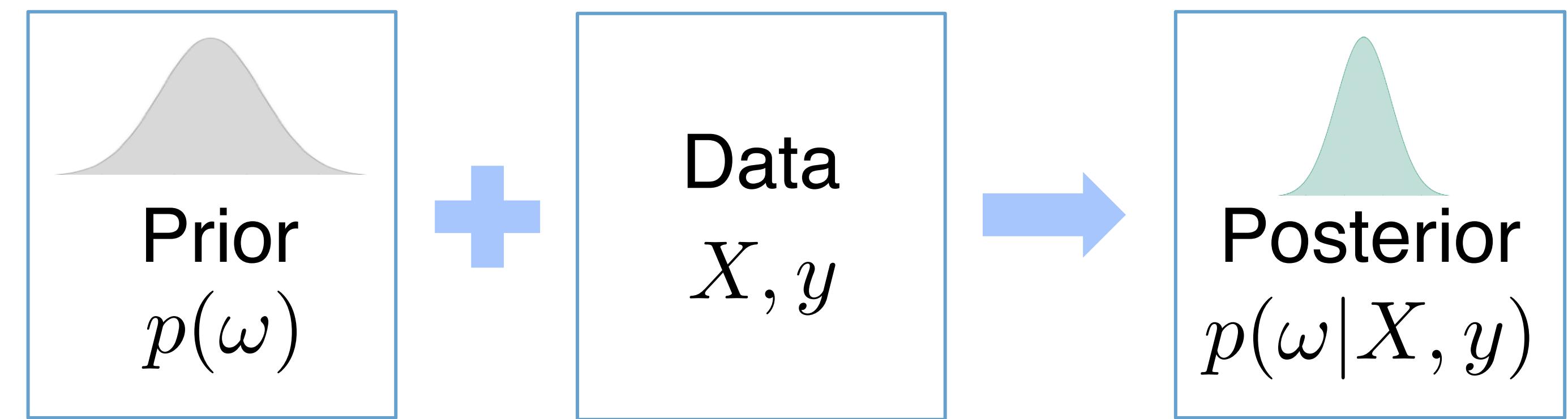


# Bayesian neural networks

Stochastic weights:



Bayesian Inference:

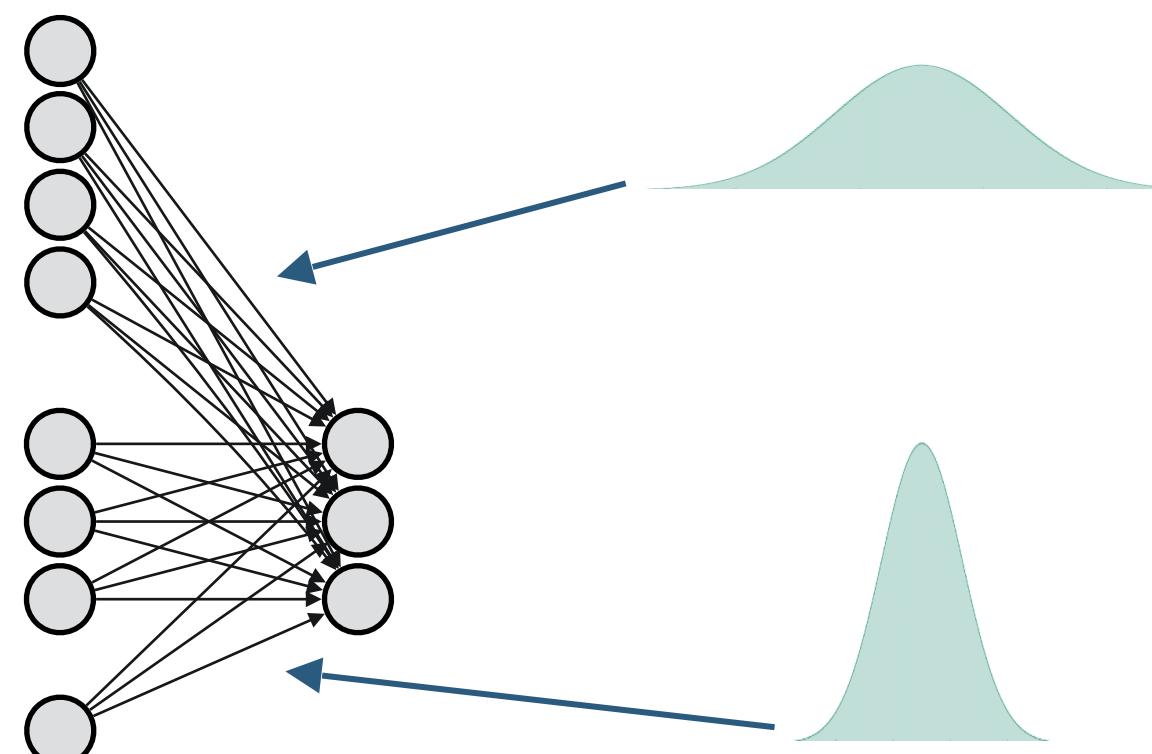


Posterior is intractable in neural networks → approximate it with  $q(\omega|\lambda)$ :

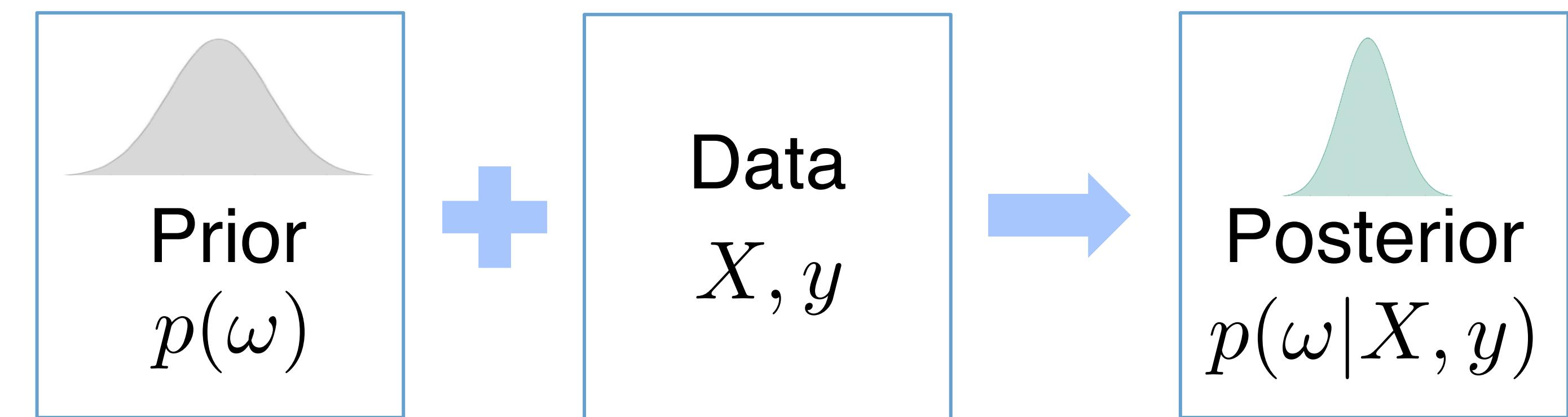
$$KL(q(\omega|\lambda)||p(\omega|X, y)) \rightarrow \min_{\lambda}$$

# Bayesian neural networks

# Stochastic weights:



# Bayesian Inference



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

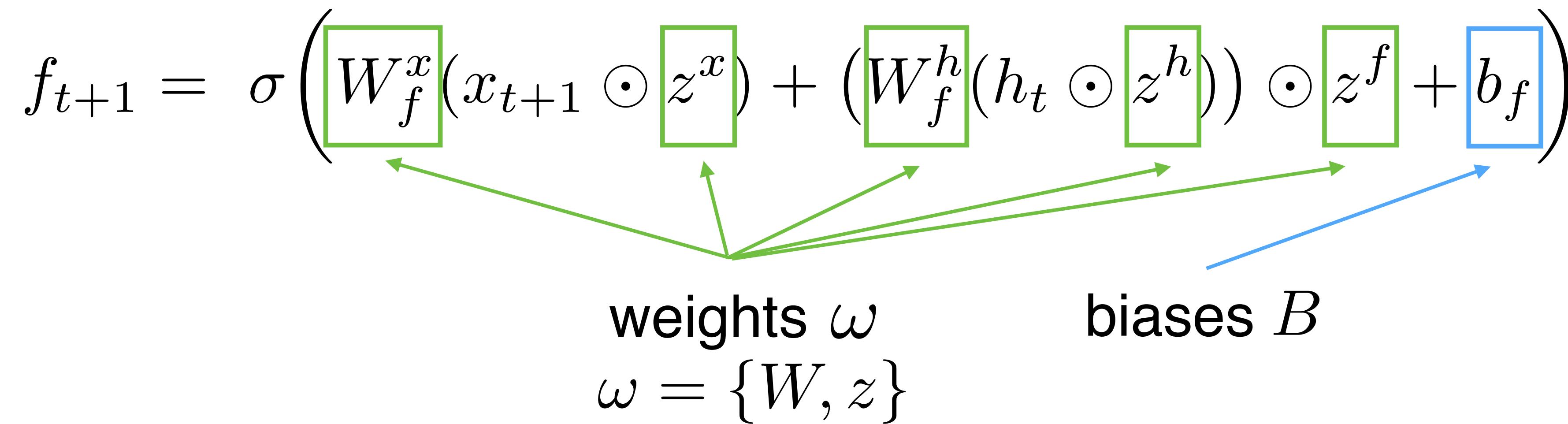
$$\sum_{i=1}^N \mathbb{E}_{q(\omega|\lambda)} \log p(y^i|x^i, \omega) - KL(q(\omega|\lambda)||p(\omega)) \rightarrow \max_{\lambda}$$

# LSTM with group weights

$$f_{t+1} = \sigma \left( W_f^x (x_{t+1} \odot z^x) + (W_f^h (h_t \odot z^h)) \odot z^f + b_f \right)$$

**weights  $\omega$**   
 $\omega = \{W, z\}$

**biases  $B$**



# Bayesian LSTM with group weights

$$f_{t+1} = \sigma \left( W_f^x (x_{t+1} \odot z^x) + (W_f^h (h_t \odot z^h)) \odot z^f + b_f \right)$$

The diagram illustrates the Bayesian LSTM update rule. It shows the calculation of the forget gate  $f_{t+1}$  as a sigmoid function of the weighted sum of inputs and hidden states, plus a bias. The weights  $W_f^x$  and  $W_f^h$  are highlighted with green boxes, and the bias  $b_f$  is highlighted with a blue box. Green arrows point from the text "stochastic weights  $\omega$ " to the weight terms  $W_f^x (x_{t+1} \odot z^x)$  and  $(W_f^h (h_t \odot z^h)) \odot z^f$ . A blue arrow points from the text "deterministic biases  $B$ " to the bias term  $b_f$ .

stochastic weights  $\omega$

deterministic biases  $B$

# Sparse Variational Dropout [1] for LSTM

$$f_{t+1} = \sigma \left( W_f^x (x_{t+1} \odot z^x) + (W_f^h (h_t \odot z^h)) \odot z^f + b_f \right)$$

stochastic weights  $\omega$

deterministic biases  $B$

Prior:  $p(\omega_k) \propto \frac{1}{|\omega_k|}$

$k = 1, \dots, |\omega|$

Approximate posterior:  $q(\omega_k | \theta_k, \sigma_k) = \mathcal{N}(\omega_k | \theta_k, \sigma_k^2)$

# Training Sparse Variational Dropout

Optimize ELBO to find  $\theta, \sigma, B$ :

$$\begin{aligned} & - \sum_{i=1}^N \overbrace{\int q(\omega|\theta, \sigma) \log p(y^i|x_0^i, \dots, x_T^i, \omega, B) d\omega}^{\text{data term}} + \\ & + \sum_{k=1}^{|\omega|} \underbrace{KL(q(\omega_k|\theta_k, \sigma_k)||p(\omega_k))}_{\text{sparsity inducing regularizer}} \rightarrow \min_{\theta, \sigma, B} \end{aligned}$$

# Training Sparse Variational Dropout

Optimize ELBO to find  $\theta, \sigma, B$ :

$$\begin{aligned} & - \sum_{i=1}^N \overbrace{\int q(\omega|\theta, \sigma) \log p(y^i|x_0^i, \dots, x_T^i, \omega, B) d\omega}^{\text{data term}} + \\ & + \underbrace{\sum_{k=1}^{|\omega|} KL(q(\omega_k|\theta_k, \sigma_k)||p(\omega_k))}_{\text{sparsity inducing regularizer}} \rightarrow \min_{\theta, \sigma, B} \end{aligned}$$

$\prod_{k=1}^{|\omega|} \mathcal{N}(\omega_k|\theta_k, \sigma_k^2)$

- First term: Monte-Carlo estimate with 1 sample, reparametrization trick:

$$\hat{\omega}_k = \theta_k + \hat{\epsilon}_k \sigma_k, \quad \hat{\epsilon}_k \sim \mathcal{N}(\epsilon_k|0, 1)$$

# Training Sparse Variational Dropout

Optimize ELBO to find  $\theta, \sigma, B$ :

$$\begin{aligned} & - \underbrace{\sum_{i=1}^N \log p(y^i | x_0^i, \dots, x_T^i, \hat{\omega}, B)}_{\text{data term}} + \\ & + \underbrace{\sum_{k=1}^{|\omega|} KL(q(\omega_k | \theta_k, \sigma_k) || p(\omega_k))}_{\text{sparsity inducing regularizer}} \rightarrow \min_{\theta, \sigma, B} \end{aligned}$$

- First term: 1 sample, reparametrization trick:

$$\hat{\omega}_k = \theta_k + \hat{\epsilon}_k \sigma_k, \quad \hat{\epsilon}_k \sim \mathcal{N}(\epsilon_k | 0, 1)$$

# Training Sparse Variational Dropout

Optimize ELBO to find  $\theta, \sigma, B$ :

$$\begin{aligned} & - \underbrace{\sum_{i=1}^N \log p(y^i | x_0^i, \dots, x_T^i, \hat{\omega}, B)}_{\text{data term}} + \\ & + \underbrace{\sum_{k=1}^{|\omega|} KL(q(\omega_k | \theta_k, \sigma_k) || p(\omega_k))}_{\text{sparsity inducing regularizer}} \rightarrow \min_{\theta, \sigma, B} \end{aligned}$$

- First term: 1 sample, reparametrization trick:

$$\hat{\omega}_k = \theta_k + \hat{\epsilon}_k \sigma_k, \quad \hat{\epsilon}_k \sim \mathcal{N}(\epsilon_k | 0, 1)$$

- Second term: analytical approximation with  $f_{KL}(\sigma_k^2 / \theta_k^2)$

# Training Sparse Variational Dropout

Optimize ELBO to find  $\theta, \sigma, B$ :

$$\begin{aligned} & - \underbrace{\sum_{i=1}^N \log p(y^i | x_0^i, \dots, x_T^i, \hat{\omega}, B)}_{\text{data term}} + \\ & + \underbrace{\sum_{k=1}^{|\omega|} f_{KL}(\sigma_k^2 / \theta_k^2)}_{\text{sparsity inducing regularizer}} \rightarrow \min_{\theta, \sigma, B} \end{aligned}$$

- First term: 1 sample, reparametrization trick:

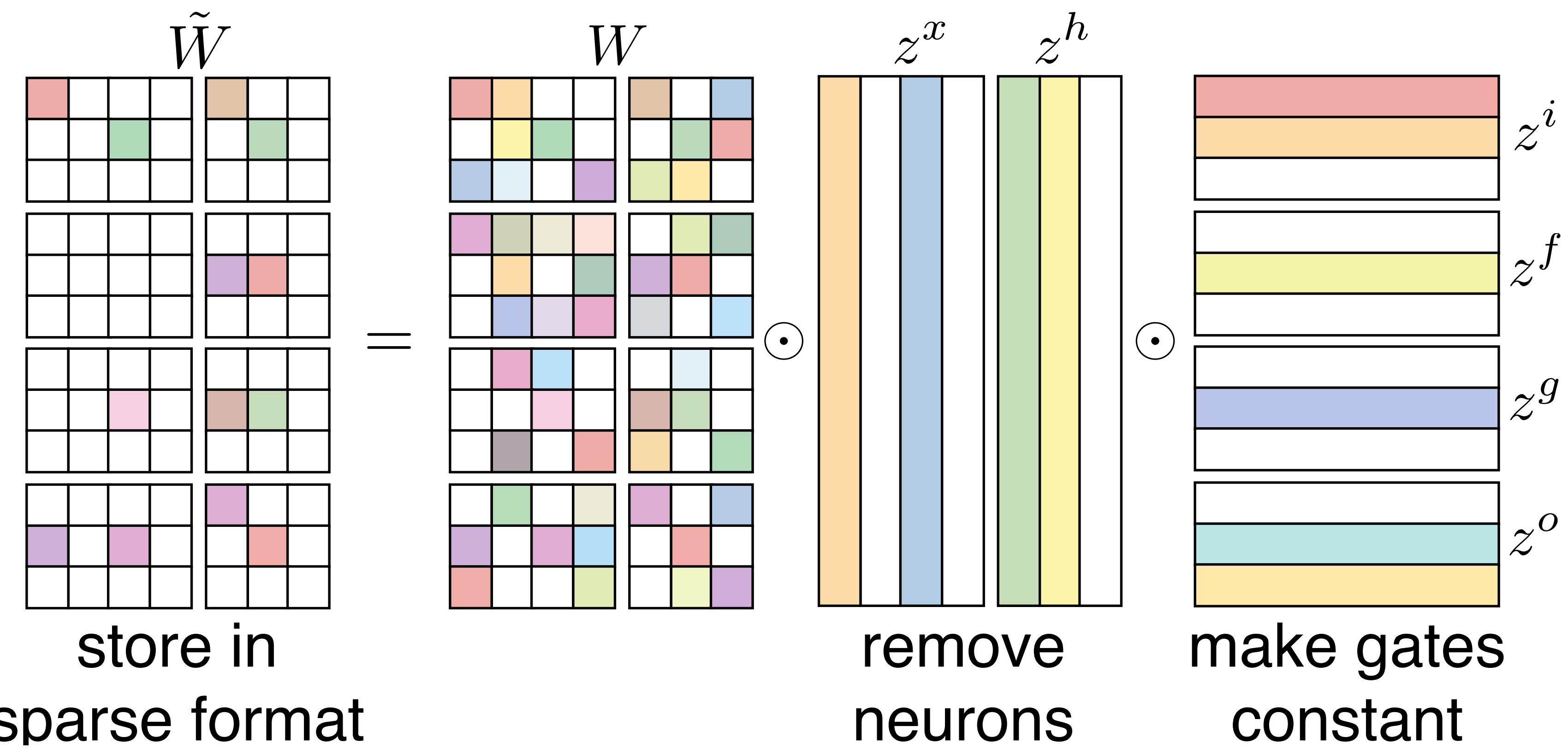
$$\hat{\omega}_k = \theta_k + \hat{\epsilon}_k \sigma_k, \quad \hat{\epsilon}_k \sim \mathcal{N}(\epsilon_k | 0, 1)$$

- Second term: analytical approximation with  $f_{KL}(\sigma_k^2 / \theta_k^2)$

# Compression after training

If  $\theta_k^2/\sigma_k^2 < \text{threshold}$ :  $\theta_k = 0, \sigma_k = 0$        $k = 1, \dots, |\omega|$        $\omega = \{W, z\}$

- On testing stage: use deterministic LSTM with  $\omega = \theta$



# Final algorithm

Training on a sequence  $x = x_1 \dots x_T$  with label  $y$  :

1. Sample weights:  $\hat{\omega} = \theta + \hat{\epsilon} \odot \sigma, \quad \hat{\epsilon} \sim \mathcal{N}(\epsilon|0, I)$
2. Forward pass:  $y_{pred} = \text{LSTM}(x, \hat{\omega}, B)$
3. Backward pass: compute stochastic gradients of ELBO:

$$\nabla_{\theta, \log \sigma, B} \left( N\text{Loss}(y, y_{pred}) + \text{SparseReg}(\theta, \sigma) \right)$$

Pruning after training:

If  $\theta_k^2 / \sigma_k^2 < \text{threshold}$ :

$$\theta_k = 0, \quad \sigma_k = 0$$

Prediction for a sequence  $x = x_1 \dots x_T$  :

Return  $\text{LSTM}(x, \theta, B)$

# Plan of the talk

1. **Idea:** introduce intermediate level of sparsification in RNN — gates
2. **Implementation** of the idea in Bayesian sparsification framework
3. **Results:** better compression and interpretable resulting gates structure

# Experiments: text classification

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

Architecture: Embedding → LSTM → FC on  $h_T$

Task	Method	Quality	Compression	Neurons $x - h$	Gates
IMDb	Original	84.1	1x	300 – 128	512
	SparseVD W	<b>84.47</b>	1135x	7 – 9	22
	SparseVD W+N	83.98	17874x	1 – 5	12
AGNews	SparseVD W+G+N	83.98	<b>19747x</b>	1 – 4	<b>6</b>
	Original	<b>90.6</b>	1x	300 – 512	2048
	SparseVD W	89.01	350x	193 – 65	72
Accuracy %	SparseVD W+N	88.55	645x	43 – 17	62
	SparseVD W+G+N	88.41	<b>647x</b>	<b>43 – 14</b>	<b>39</b>

W: weights    N: neurons    G: gates (ours)

# Experiments: text generation

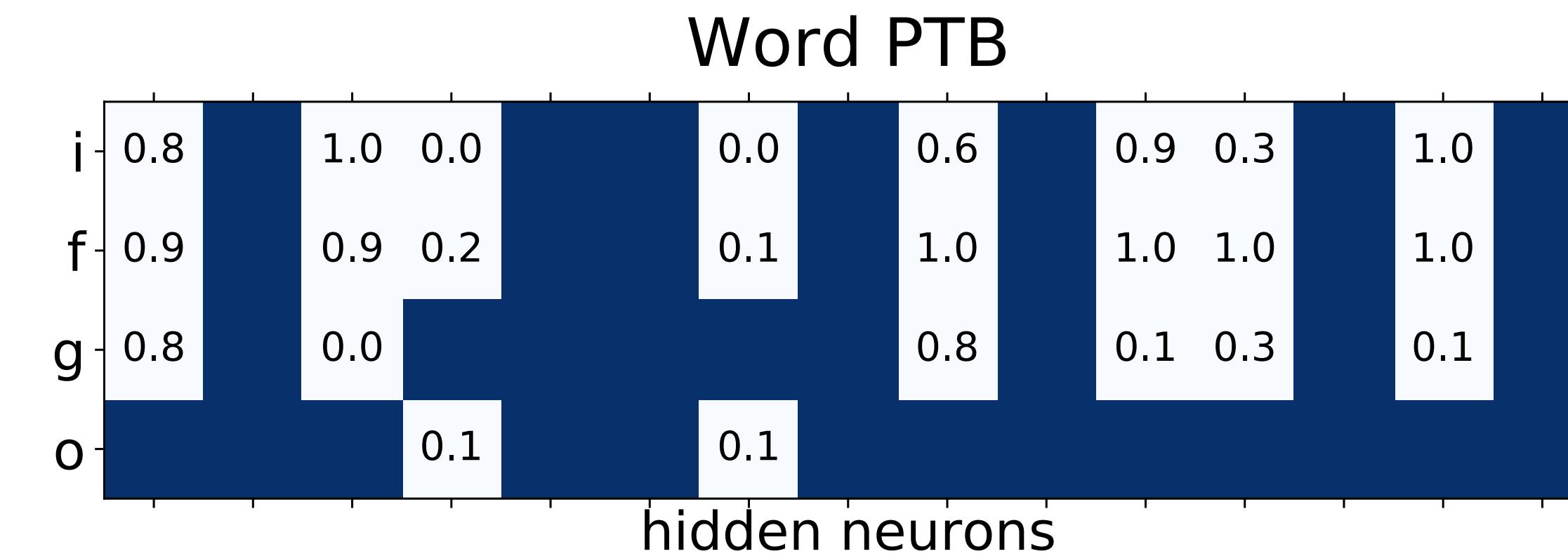
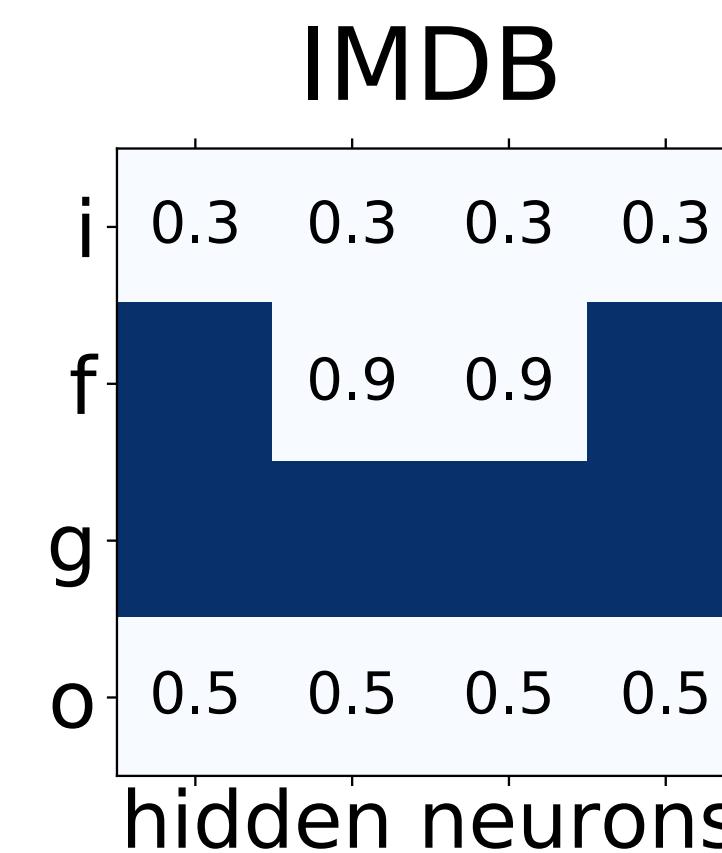
PTB Corpus: char-level (6M chars), word-level (0.9M words)

Architecture: LSTM → FC

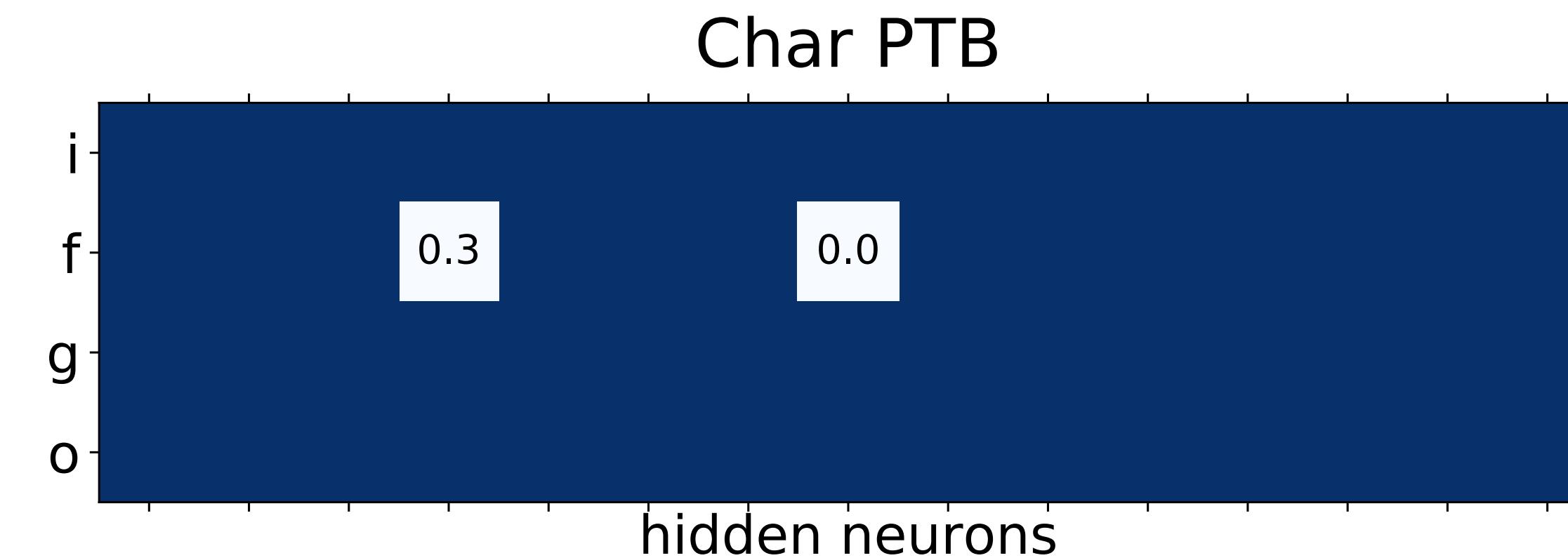
Task	Method	Quality	Compression	Neurons $x - h$	Gates
Char PTB	Original	1.498 – 1.455	1x	50 – 1000	4000
	SparseVD W	1.472 – 1.429	7.9x	50 – 431	1718
	SparseVD W+N	1.478 – 1.430	<b>10.0x</b>	<b>50 – 390</b>	<b>1560</b>
	SparseVD W+G+N	<b>1.467 – 1.425</b>	9.8x	50 – 404	1563
Word PTB	Original	135.6 – 129.3	1x	10000 – 256	1024
	SparseVD W	<b>115.0 – 109.0</b>	22.0x	9985 – 153	281
	SparseVD W+N	116.2 – 111.0	23.4x	9993 – 134	335
	SparseVD W+G+N	122.2 – 116.5	<b>24.4x</b>	<b>9973 – 114</b>	<b>220</b>

W: weights    N: neurons    G: gates (ours)

# Experiments: resulting gates structure



- C constant gate with value c
- non-constant gate

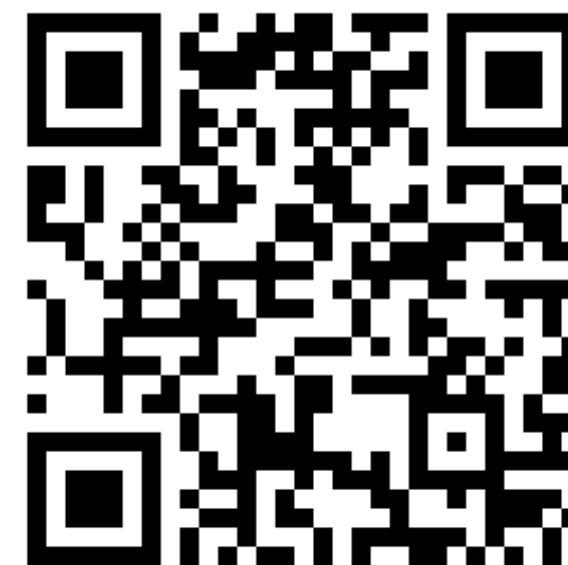


Different gates are important for different tasks

# Summary

- Three levels of sparsity for gated RNNs:  
(1) sparsify weights; (2) make gates constant; (3) remove neurons
- Can be easily implemented in any sparsification framework  
Our implementation: introduce group weights and use Bayesian sparsification
- Improve compression and reveal LSTM specifics on different tasks

**Extended abstract**



**Source code**

