

344.063 KV Special Topic:

Natural Language Processing with Deep Learning

Gated RNNs: LSTM & GRU



Navid Rekab-saz

navid.rekabsaz@jku.at

Agenda

- Backpropagation Through Time
- RNNs with Gates: LSTM, GRU

Element-wise Multiplication

- $a \odot b = c$

- dimensions: $1 \times d \odot 1 \times d = 1 \times d$

$$[1 \quad 2 \quad 3] \odot [3 \quad 0 \quad -2] = [3 \quad 0 \quad -6]$$

- $A \odot B = C$

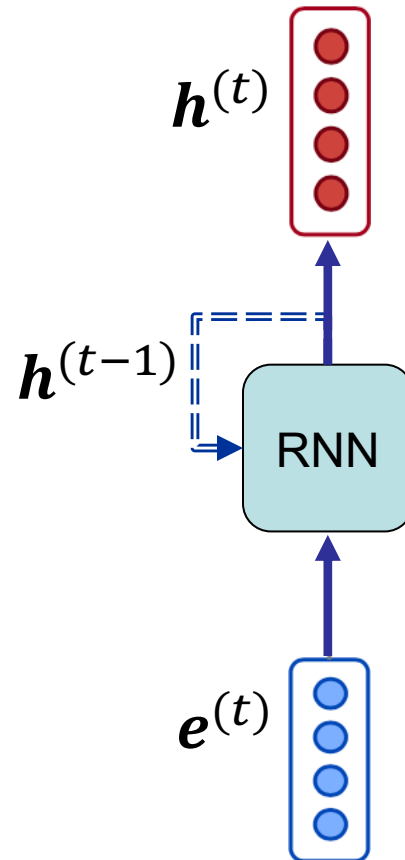
- dimensions: $l \times m \odot l \times m = l \times m$

$$\begin{bmatrix} 2 & 3 \\ 0 & 1 \\ 1 & -1 \end{bmatrix} \odot \begin{bmatrix} -1 & 0 \\ 0 & 2 \\ 0.5 & -1 \end{bmatrix} = \begin{bmatrix} -2 & 0 \\ 0 & 2 \\ 0.5 & 1 \end{bmatrix}$$

Agenda

- **Backpropagation Through Time**
- RNNs with Gates: LSTM, GRU

Recurrent Neural Networks – recap



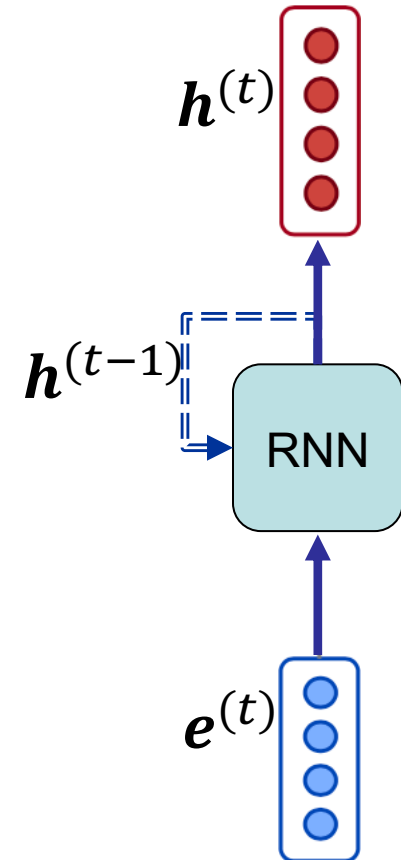
Vanilla (Elman) RNN – recap

- General form of an RNN function

$$\mathbf{h}^{(t)} = \text{RNN}(\mathbf{h}^{(t-1)}, \mathbf{e}^{(t)})$$

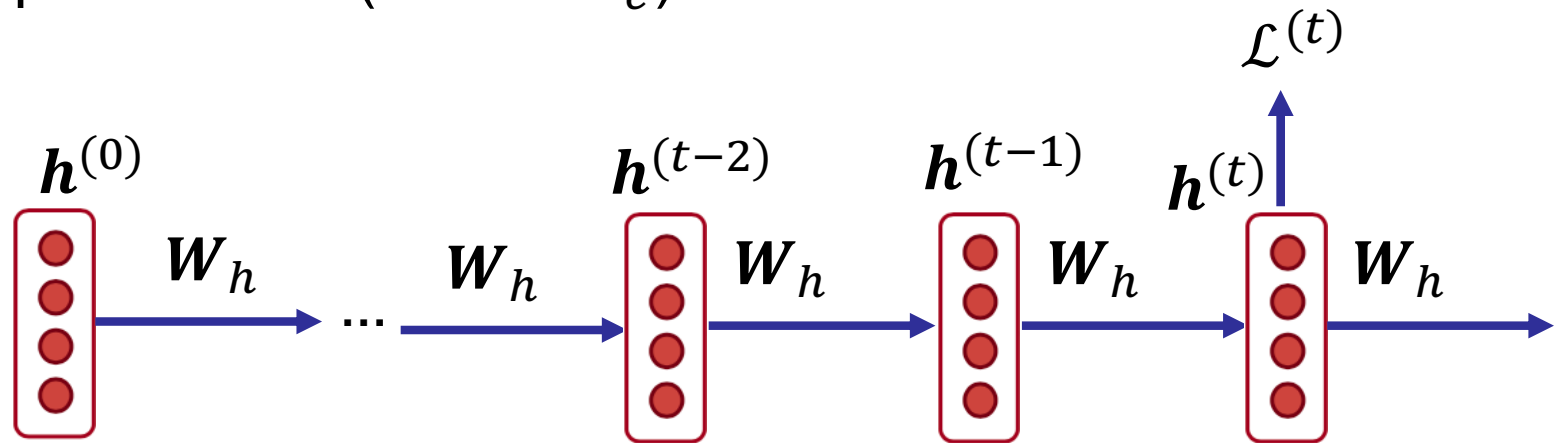
- Vanilla RNN:

$$\mathbf{h}^{(t)} = \sigma(\mathbf{h}^{(t-1)} \mathbf{W}_h + \mathbf{e}^{(t)} \mathbf{W}_e + \mathbf{b})$$



Backpropagation for RNNs

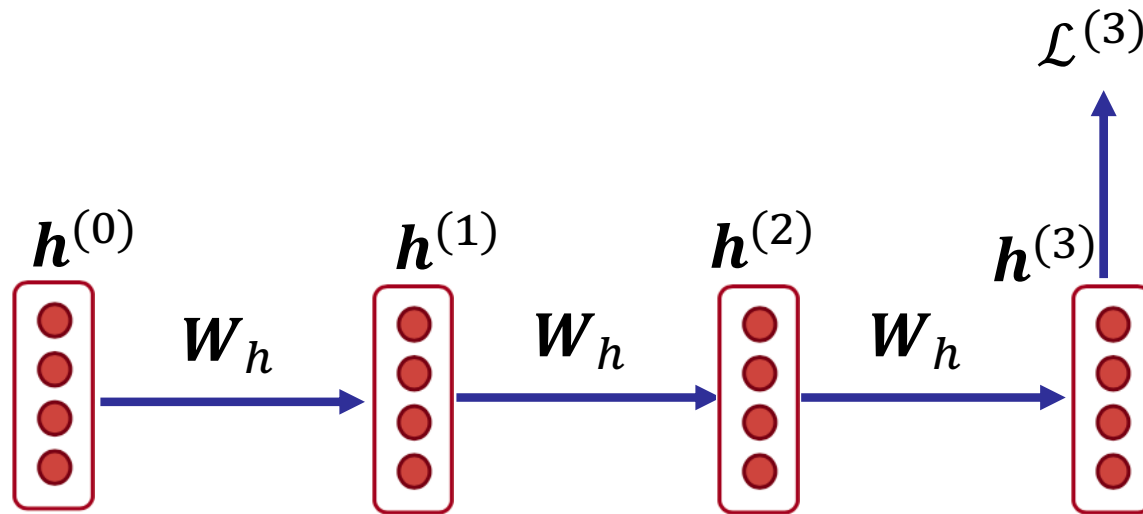
- Unrolling the computation graph of RNN
- Simplified: the interactions with U and also input parameters (E and W_e) are removed



- What is ...

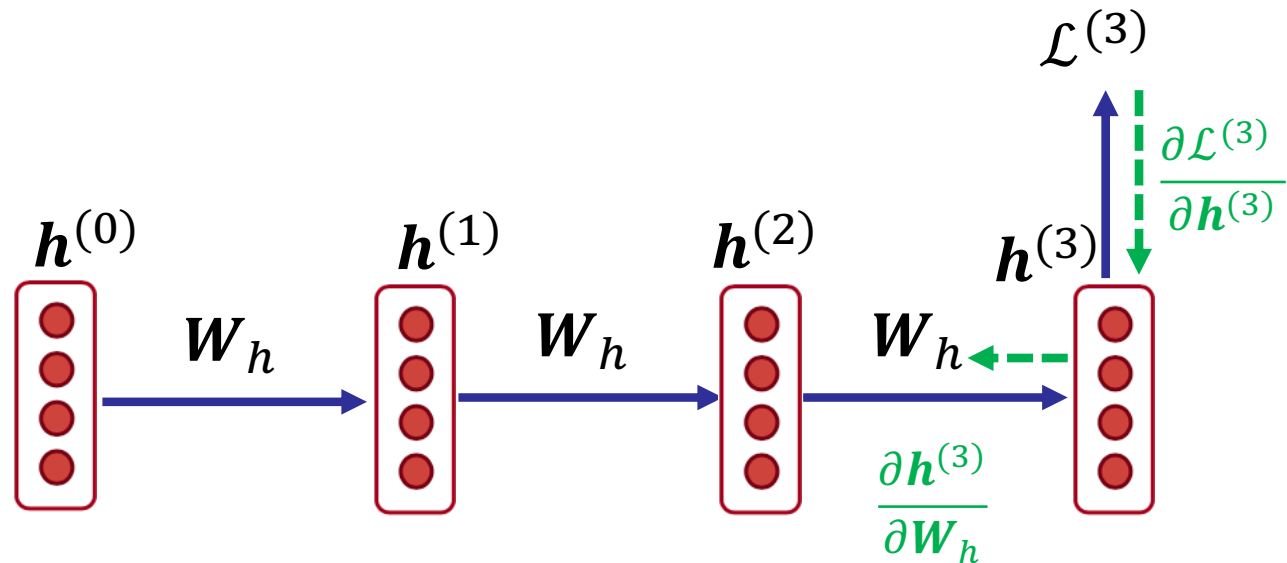
$$\frac{\partial \mathcal{L}^{(t)}}{\partial W_h} = ?$$

Backpropagation for RNNs



$$\frac{\partial \mathcal{L}^{(3)}}{\partial W_h} = ?$$

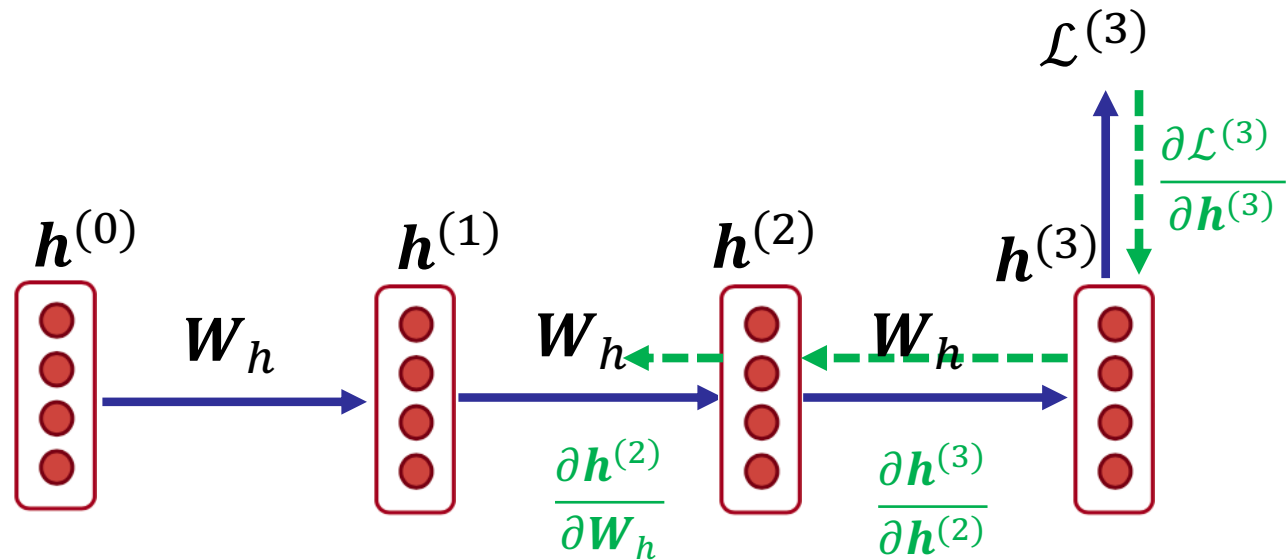
Backpropagation for RNNs



$$\left. \frac{\partial \mathcal{L}^{(3)}}{\partial W_h} \right|_{(3)} = \frac{\partial \mathcal{L}^{(3)}}{\partial h^{(3)}} \frac{\partial h^{(3)}}{\partial W_h}$$

- Gradient regarding W_h at time step 3

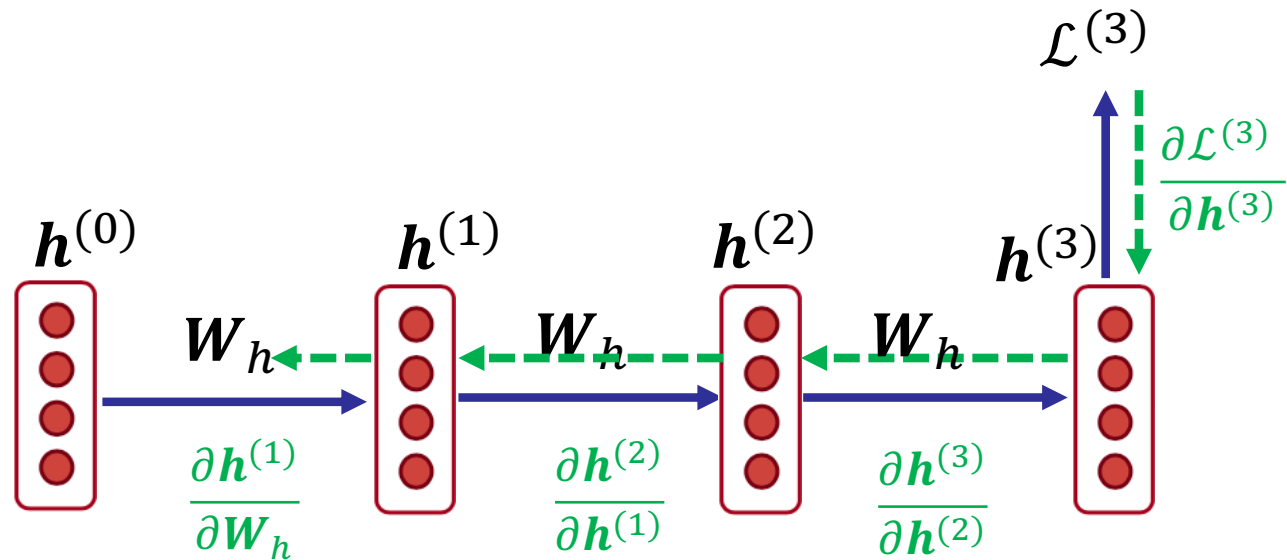
Backpropagation for RNNs



$$\left. \frac{\partial \mathcal{L}^{(3)}}{\partial W_h} \right|_{(2)} = \frac{\partial \mathcal{L}^{(3)}}{\partial h^{(3)}} \frac{\partial h^{(3)}}{\partial h^{(2)}} \frac{\partial h^{(2)}}{\partial W_h}$$

- Gradient regarding W_h at time step 2

Backpropagation for RNNs



$$\left. \frac{\partial \mathcal{L}^{(3)}}{\partial W_h} \right|_{(1)} = \frac{\partial \mathcal{L}^{(3)}}{\partial h^{(3)}} \frac{\partial h^{(3)}}{\partial h^{(2)}} \frac{\partial h^{(2)}}{\partial h^{(1)}} \frac{\partial h^{(1)}}{\partial W_h}$$

- Gradient regarding W_h at time step 1

Backpropagation Through Time (BPTT)

- Final gradient is the sum of the gradients regarding the model parameters (such as \mathbf{W}_h) from the current time step back to the beginning of corpus (or batch)

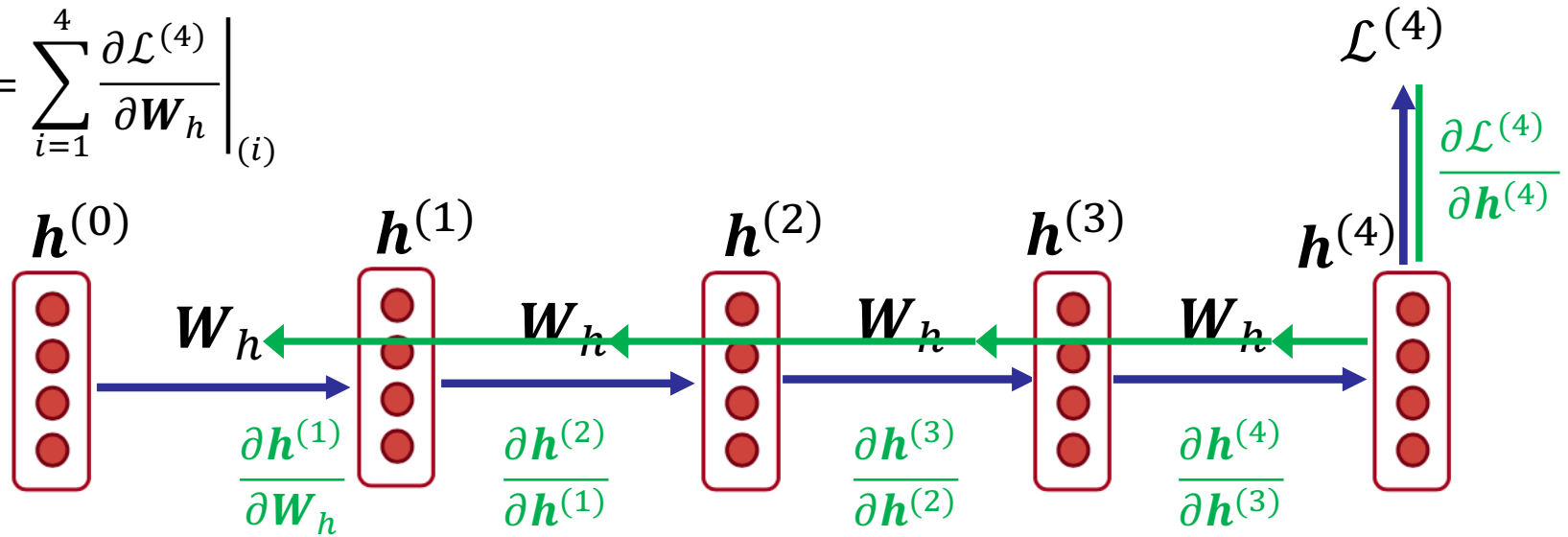
$$\frac{\partial \mathcal{L}^{(t)}}{\partial \mathbf{W}_h} = \sum_{i=1}^t \frac{\partial \mathcal{L}^{(t)}}{\partial \mathbf{W}_h} \Big|_{(i)}$$

- In this simplified case, this can be written as:

$$\frac{\partial \mathcal{L}^{(t)}}{\partial \mathbf{W}_h} = \sum_{i=1}^t \frac{\partial \mathcal{L}^{(t)}}{\partial \mathbf{h}^{(t)}} \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}} \cdots \frac{\partial \mathbf{h}^{(i)}}{\partial \mathbf{W}_h}$$

Backpropagation Through Time (BPTT) – all in one!

$$\frac{\partial \mathcal{L}^{(4)}}{\partial \mathbf{W}_h} = \sum_{i=1}^4 \frac{\partial \mathcal{L}^{(4)}}{\partial \mathbf{W}_h} \Big|_{(i)}$$



$$\frac{\partial \mathcal{L}^{(4)}}{\partial \mathbf{W}_h} \Big|_{(4)} = \frac{\partial \mathcal{L}^{(4)}}{\partial \mathbf{h}^{(4)}} \frac{\partial \mathbf{h}^{(4)}}{\partial \mathbf{W}_h}$$

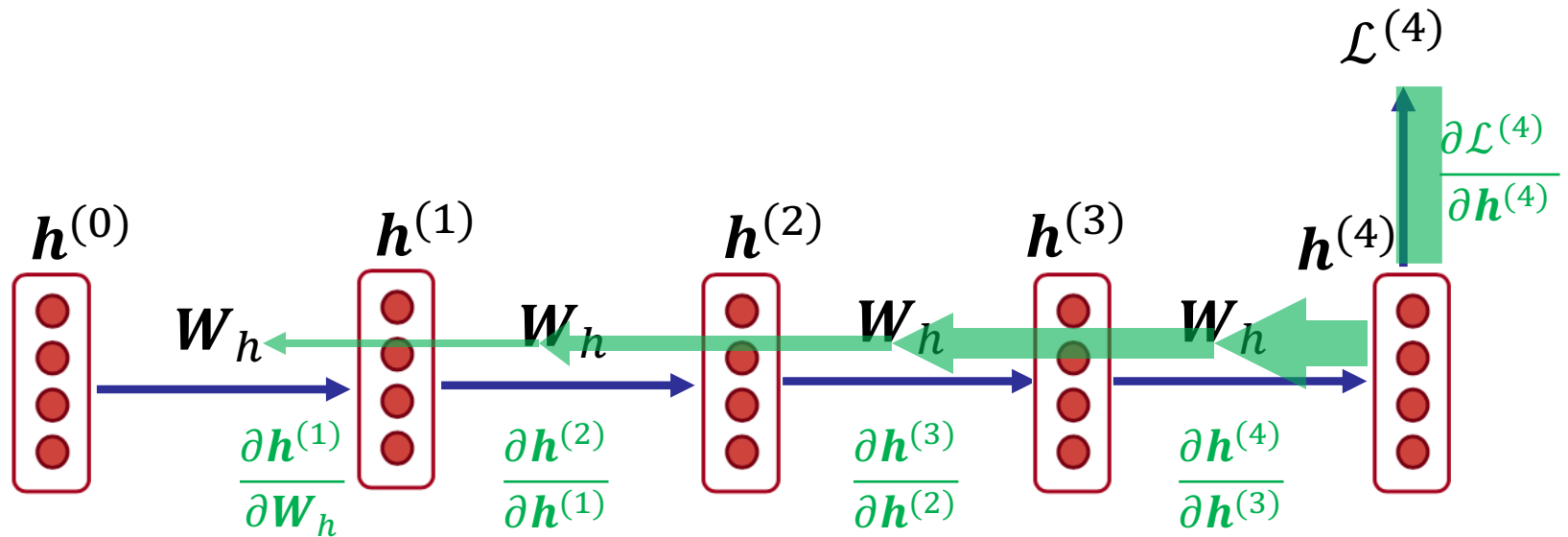
$$\frac{\partial \mathcal{L}^{(4)}}{\partial \mathbf{W}_h} \Big|_{(3)} = \frac{\partial \mathcal{L}^{(4)}}{\partial \mathbf{h}^{(4)}} \frac{\partial \mathbf{h}^{(4)}}{\partial \mathbf{h}^{(3)}} \frac{\partial \mathbf{h}^{(3)}}{\partial \mathbf{W}_h}$$

$$\frac{\partial \mathcal{L}^{(4)}}{\partial \mathbf{W}_h} \Big|_{(2)} = \frac{\partial \mathcal{L}^{(4)}}{\partial \mathbf{h}^{(4)}} \frac{\partial \mathbf{h}^{(4)}}{\partial \mathbf{h}^{(3)}} \frac{\partial \mathbf{h}^{(3)}}{\partial \mathbf{h}^{(2)}} \frac{\partial \mathbf{h}^{(2)}}{\partial \mathbf{W}_h}$$

$$\frac{\partial \mathcal{L}^{(4)}}{\partial \mathbf{W}_h} \Big|_{(1)} = \frac{\partial \mathcal{L}^{(4)}}{\partial \mathbf{h}^{(4)}} \frac{\partial \mathbf{h}^{(4)}}{\partial \mathbf{h}^{(3)}} \frac{\partial \mathbf{h}^{(3)}}{\partial \mathbf{h}^{(2)}} \frac{\partial \mathbf{h}^{(2)}}{\partial \mathbf{h}^{(1)}} \frac{\partial \mathbf{h}^{(1)}}{\partial \mathbf{W}_h}$$

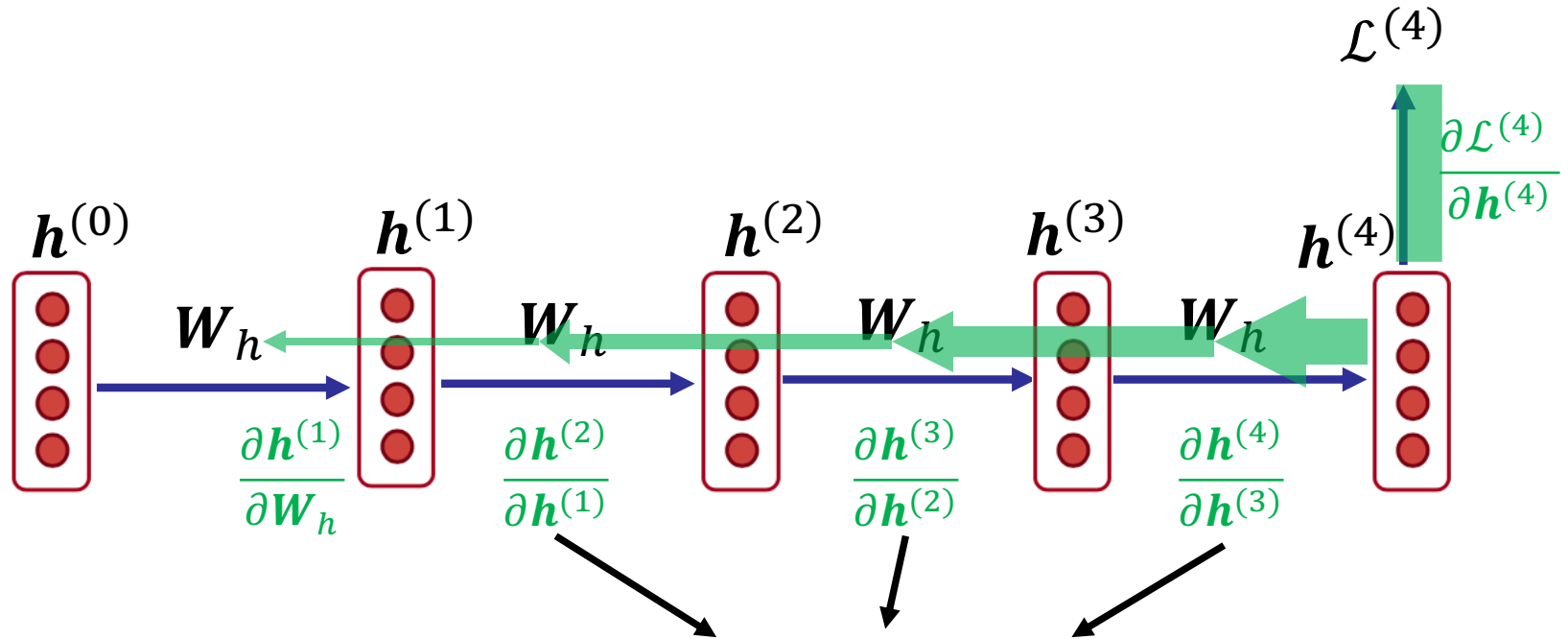
$$\frac{\partial \mathcal{L}^{(t)}}{\partial \mathbf{W}_h} = \sum_{i=1}^t \frac{\partial \mathcal{L}^{(t)}}{\partial \mathbf{W}_h} \Big|_{(i)} = \sum_{i=1}^t \frac{\partial \mathcal{L}^{(t)}}{\partial \mathbf{h}^{(t)}} \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}} \cdots \frac{\partial \mathbf{h}^{(i)}}{\partial \mathbf{W}_h}$$

Vanishing/Exploding gradient



- In practice, the gradient regarding each time step becomes smaller and smaller as it goes back in time → **Vanishing gradient**
- While less often, this may also happen other way around: the gradient regarding further time steps becomes larger and larger → **Exploding gradient**

Vanishing/Exploding gradient – why?



If these gradients are small, their multiplication gets smaller. As we go further back, the final gradient contains more of these!

$$\left. \frac{\partial \mathcal{L}^{(4)}}{\partial W_h} \right|_{(1)} = \frac{\partial \mathcal{L}^{(4)}}{\partial h^{(4)}} \frac{\partial h^{(4)}}{\partial h^{(3)}} \frac{\partial h^{(3)}}{\partial h^{(2)}} \frac{\partial h^{(2)}}{\partial h^{(1)}} \frac{\partial h^{(1)}}{\partial W_h}$$

Vanishing/Exploding gradient – why?

- What is $\frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}}$?!

- Recall the definition of RNN:

$$\mathbf{h}^{(t)} = \sigma(\mathbf{h}^{(t-1)}\mathbf{W}_h + \mathbf{e}^{(t)}\mathbf{W}_e + \mathbf{b})$$

- Let's replace sigmoid (σ) with a simple linear activation ($y = x$) function.

$$\mathbf{h}^{(t)} = \mathbf{h}^{(t-1)}\mathbf{W}_h + \mathbf{e}^{(t)}\mathbf{W}_e + \mathbf{b}$$

- In this case:

$$\frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}} = \mathbf{W}_h$$

Vanishing/Exploding gradient – why?

- Recall the BPTT formula (for the simplified case):

$$\left. \frac{\partial \mathcal{L}^{(t)}}{\partial \mathbf{W}_h} \right|_{(i)} = \frac{\partial \mathcal{L}^{(t)}}{\partial \mathbf{h}^{(t)}} \underbrace{\frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}} \cdots \frac{\partial \mathbf{h}^{(i+1)}}{\partial \mathbf{h}^{(i)}}}_{\text{product of Jacobians}} \frac{\partial \mathbf{h}^{(i)}}{\partial \mathbf{W}_h}$$

- Given $l = t - i$, the BPTT formula can be rewritten as:

$$\left. \frac{\partial \mathcal{L}^{(t)}}{\partial \mathbf{W}_h} \right|_{(i)} = \frac{\partial \mathcal{L}^{(t)}}{\partial \mathbf{h}^{(t)}} \boxed{(\mathbf{W}_h)^l} \frac{\partial \mathbf{h}^{(i)}}{\partial \mathbf{W}_h}$$

If weights in \mathbf{W}_h are small (i.e. eigenvalues of \mathbf{W}_h are smaller than 1), this term gets *exponentially* smaller

Why is vanishing/exploding gradient a problem?

- Vanishing gradient

- Gradient signal from faraway “fades away” and becomes insignificant in comparison with the gradient signal from close-by
 - Long-term dependencies are not captured, since model weights are updated only with respect to near effects
- one approach to address it: RNNs with gates – LSTM, GRU

- Exploding gradient

- Gradients become too big → SGD update steps become too large
 - This causes (large loss values and) large updates on parameters, and eventually unstable training
- main approach to address it: Gradient clipping

Gradient clipping

- Gradient clipping: if the **norm of the gradient** is greater than some **threshold**, scale the gradient down

Algorithm 1 Pseudo-code for norm clipping

$$\hat{\mathbf{g}} \leftarrow \frac{\partial \mathcal{E}}{\partial \theta}$$

if $\|\hat{\mathbf{g}}\| \geq \textit{threshold}$ **then**

$$\hat{\mathbf{g}} \leftarrow \frac{\textit{threshold}}{\|\hat{\mathbf{g}}\|} \hat{\mathbf{g}}$$

end if

- **Intuition**: take the step in the same direction, but with a smaller step

Problem with vanilla RNN – summary

- It is too difficult for the hidden state of vanilla RNN to learn and preserve information of several time steps
 - In particular as new contents are constantly added to the hidden state in every step

$$\mathbf{h}^{(t)} = \sigma(\mathbf{h}^{(t-1)}\mathbf{W}_h + \mathbf{e}^{(t)}\mathbf{W}_e + \mathbf{b})$$



In every step, **input vector** “adds” new content to hidden state

Agenda

- Backpropagation Through Time
- **RNNs with Gates: LSTM, GRU**

Gate vector



- Gate vector:
 - A vector with values between 0 and 1
 - Gate vector acts as “gate-keeper”, such that it controls the content flow of another vector
- Gate vectors are typically defined using sigmoid:

$$\mathbf{g} = \sigma(\text{some vector})$$

... and are applied to a vector \mathbf{v} with element-wise multiplication to control its contents:

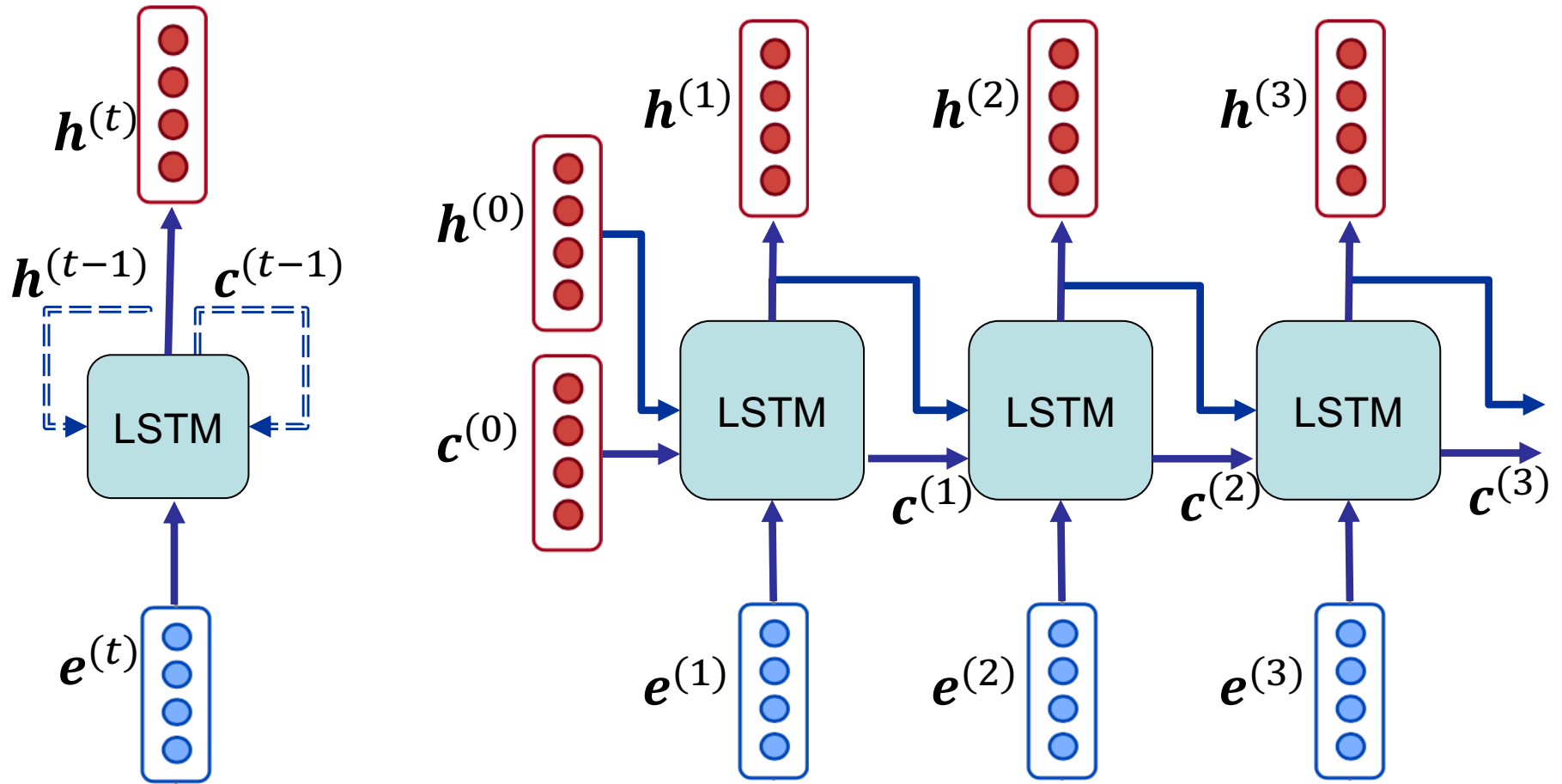
$$\mathbf{g} \odot \mathbf{v}$$

- For each element (feature) i of the vectors:
 - If g_i is 1 $\rightarrow v_i$ remains the same; everything passes; *open gate!*
 - If g_i is 0 $\rightarrow v_i$ becomes 0; nothing passes; *closed gate!*

Long Short-Term Memory (LSTM)

- Proposed by Hochreiter and Schmidhuber in 1997
- LSTM exploits a new vector **cell state** $c^{(t)}$ to carry the memory of previous states
 - The cell state stores **long-term information**
 - As in vanilla RNN, hidden states $h^{(t)}$ is used as **output vector**
- LSTM controls the process of **reading**, **writing**, and **erasing** information in/from memory states
 - These controls are done using **gate vectors**
 - Gates are **dynamic** and defined based on the input vector and hidden state

LSTM – unrolled



LSTM definition – gates

- Gates are functions of input vector $\mathbf{e}^{(t)}$ and previous hidden state $\mathbf{h}^{(t-1)}$

$$\mathbf{i}^{(t)} = \text{function}(\mathbf{h}^{(t-1)}, \mathbf{e}^{(t)})$$

$$\mathbf{i}^{(t)} = \sigma(\mathbf{h}^{(t-1)} \mathbf{W}_{hi} + \mathbf{e}^{(t)} \mathbf{W}_{xi} + \mathbf{b}_i)$$

input gate: controls what parts of the new cell content are written to cell

$$\mathbf{f}^{(t)} = \text{function}(\mathbf{h}^{(t-1)}, \mathbf{e}^{(t)})$$

$$\mathbf{f}^{(t)} = \sigma(\mathbf{h}^{(t-1)} \mathbf{W}_{hf} + \mathbf{e}^{(t)} \mathbf{W}_{xf} + \mathbf{b}_f)$$

forget gate: controls what is kept vs forgotten, from previous cell state

$$\mathbf{o}^{(t)} = \text{function}(\mathbf{h}^{(t-1)}, \mathbf{e}^{(t)})$$

$$\mathbf{o}^{(t)} = \sigma(\mathbf{h}^{(t-1)} \mathbf{W}_{ho} + \mathbf{e}^{(t)} \mathbf{W}_{xo} + \mathbf{b}_o)$$

output gate: controls what parts of cell are output to hidden state

Parameters are shown in red

LSTM definition – states

$$\tilde{\mathbf{c}}^{(t)} = \text{function}(\mathbf{h}^{(t-1)}, \mathbf{e}^{(t)})$$

$$\tilde{\mathbf{c}}^{(t)} = \tanh(\mathbf{h}^{(t-1)} \mathbf{W}_{hc} + \mathbf{e}^{(t)} \mathbf{W}_{xc} + \mathbf{b}_c)$$

new cell content: the new content to be used for cell and hidden (output) state

$$\mathbf{c}^{(t)} = \mathbf{f}^{(t)} \odot \mathbf{c}^{(t-1)} + \mathbf{i}^{(t)} \odot \tilde{\mathbf{c}}^{(t)}$$

cell state: erases (“forgets”) some content from last cell state, and writes (“inputs”) some new cell content

$$\mathbf{h}^{(t)} = \mathbf{o}^{(t)} \odot \tanh(\mathbf{c}^{(t)})$$

hidden state: reads (“outputs”) some content from the current cell state

LSTM definition – all together

$$\mathbf{i}^{(t)} = \sigma(\mathbf{h}^{(t-1)}\mathbf{W}_{hi} + \mathbf{e}^{(t)}\mathbf{W}_{xi} + \mathbf{b}_i)$$

input gate: controls what parts of the new cell content are written to cell

$$\mathbf{f}^{(t)} = \sigma(\mathbf{h}^{(t-1)}\mathbf{W}_{hf} + \mathbf{e}^{(t)}\mathbf{W}_{xf} + \mathbf{b}_f)$$

forget gate: controls what is kept vs forgotten, from previous cell state

$$\mathbf{o}^{(t)} = \sigma(\mathbf{h}^{(t-1)}\mathbf{W}_{ho} + \mathbf{e}^{(t)}\mathbf{W}_{xo} + \mathbf{b}_o)$$

output gate: controls what parts of cell are output to hidden state

$$\tilde{\mathbf{c}}^{(t)} = \tanh(\mathbf{h}^{(t-1)}\mathbf{W}_{hc} + \mathbf{e}^{(t)}\mathbf{W}_{xc} + \mathbf{b}_c)$$

new cell content: the new content to be used for cell and hidden (output) state

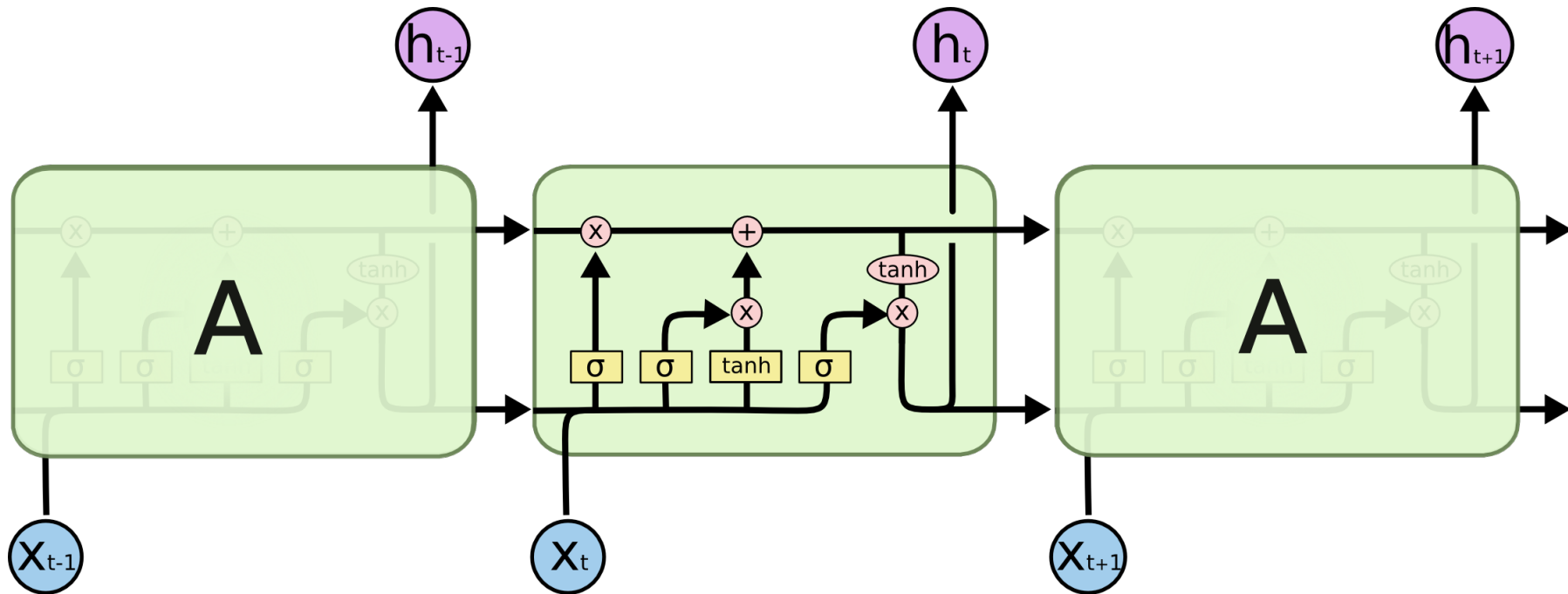
$$\mathbf{c}^{(t)} = \mathbf{f}^{(t)} \odot \mathbf{c}^{(t-1)} + \mathbf{i}^{(t)} \odot \tilde{\mathbf{c}}^{(t)}$$


cell state: erases (“forgets”) some content from last cell state, and writes (“inputs”) some new cell content

$$\mathbf{h}^{(t)} = \mathbf{o}^{(t)} \odot \tanh(\mathbf{c}^{(t)})$$


hidden state: reads (“outputs”) some content from the current cell state

LSTM definition – visually!




 Neural Network
 Layer


 Pointwise
 Operation


 Vector
 Transfer


 Concatenate


 Copy

Gated Recurrent Unit (GRU)

$$\mathbf{u}^{(t)} = \sigma(\mathbf{h}^{(t-1)} \mathbf{W}_{hu} + \mathbf{e}^{(t)} \mathbf{W}_{xu} + \mathbf{b}_u)$$

update gate: controls what parts of hidden state are updated vs preserved

$$\mathbf{r}^{(t)} = \sigma(\mathbf{h}^{(t-1)} \mathbf{W}_{hr} + \mathbf{e}^{(t)} \mathbf{W}_{xr} + \mathbf{b}_r)$$

reset gate: controls what parts of previous hidden state are used to compute new content

new hidden state content: (1) reset gate selects useful parts of previous hidden state. (2) Use this and current input to compute new hidden content.

$$\tilde{\mathbf{h}}^{(t)} = \tanh((\mathbf{r}^{(t)} \odot \mathbf{h}^{(t-1)}) \mathbf{W}_{hh} + \mathbf{e}^{(t)} \mathbf{W}_{xh} + \mathbf{b}_h)$$

$$\mathbf{h}^{(t)} = (1 - \mathbf{u}^{(t)}) \odot \mathbf{h}^{(t-1)} + \mathbf{u}^{(t)} \odot \tilde{\mathbf{h}}^{(t)}$$

hidden state: update gate simultaneously controls what is kept from previous hidden state, and what is updated to new hidden state content

Parameters are shown in red

RNNs with gates – counting parameters

- Parameters in LSTM (bias terms discarded)
 - $W_{hi}, W_{hf}, W_{ho}, W_{hc} \rightarrow h \times h * 4$
 - $W_{xi}, W_{xf}, W_{xo}, W_{xc} \rightarrow d \times h * 4$
- Parameters in GRU (bias terms discarded)
 - $W_{hu}, W_{hr}, W_{hh} \rightarrow h \times h * 3$
 - $W_{xu}, W_{xr}, W_{xh} \rightarrow d \times h * 3$
- If also considering encoder and decoder embeddings (e.g., in a Language Modeling network)
 - $E \rightarrow N \times d$
 - $U \rightarrow h \times N$

d : dimension of input embedding

h : dimension of hidden vectors and output embedding

RNNs with gates – summary

- LSTM (and GRU) with dynamic gate mechanisms makes it easier to **preserve necessary information** over many timesteps
- LSTM does not *guarantee* that there is no vanishing/exploding gradient, but its large success in practice has shown that it can **learn long-distance dependencies**
- LSTM vs. GRU: LSTM is usually the **default choice**. Especially, when enough training data is available and capturing longer distances is important. GRU is faster and more suited for settings with low computation resources