# Statistical Natural Language Processing [COMP0087]

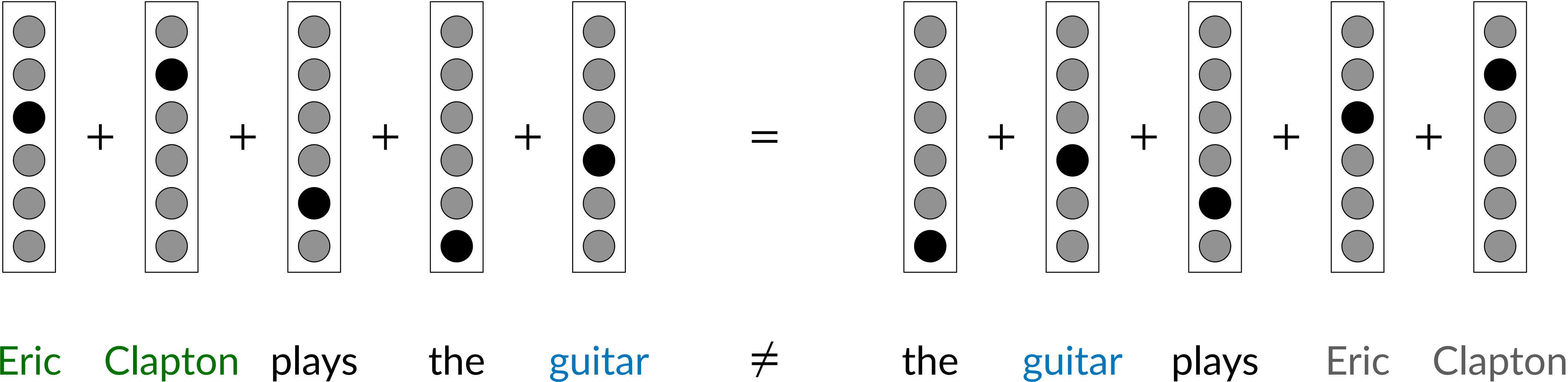## *Recurrent Neural Networks*

### Vasileios Lampos

Computer Science, UCL

lampos.net

▶ In this lecture:

— brief overview on language models (*more on this during the lecture by Dr. Oana-Maria Camburu*)

— Recurrent Neural Networks

— The Long Short-Term Memory (LSTM) architecture

— Applications and extensions

— slides: lampos.net/teaching

▶ **Reading / Lecture based on**: Chapters 3 (*less so*), 7 (*less so*), and 9 (*more so*) of "*Speech and Language Processing*" (SLP) by Jurafsky and Martin (2023) — web.stanford.edu/~jurafsky/slp3/

▶ Additional material

✳ Difficulties in training RNNs — proceedings.mlr.press/v28/pascanu13.pdf

✳ LSTMs — colah.github.io/posts/2015-08-Understanding-LSTMs/

Language is a **sequence** of "events" over time



Eric   Clapton   plays   the   guitar   ≠   the   guitar   plays   Eric   Clapton

A language model predicts the next word of a word sequence:

*... and all of a sudden Eric Clapton started to play the _____*

guitar ?

piano

game

drums

A language model predicts the next word of a word sequence:

*… and all of a sudden Eric Clapton started to play the _____*
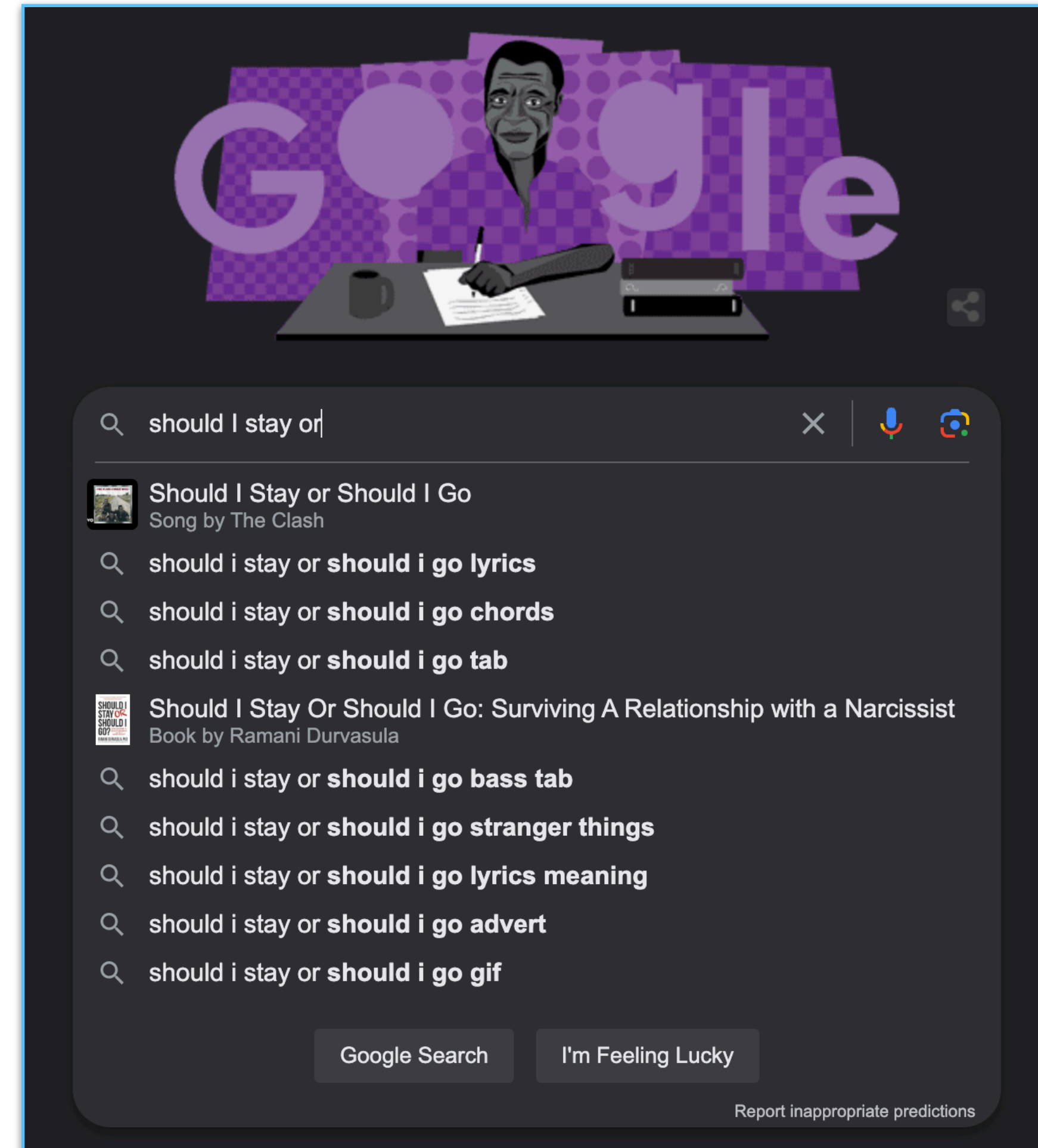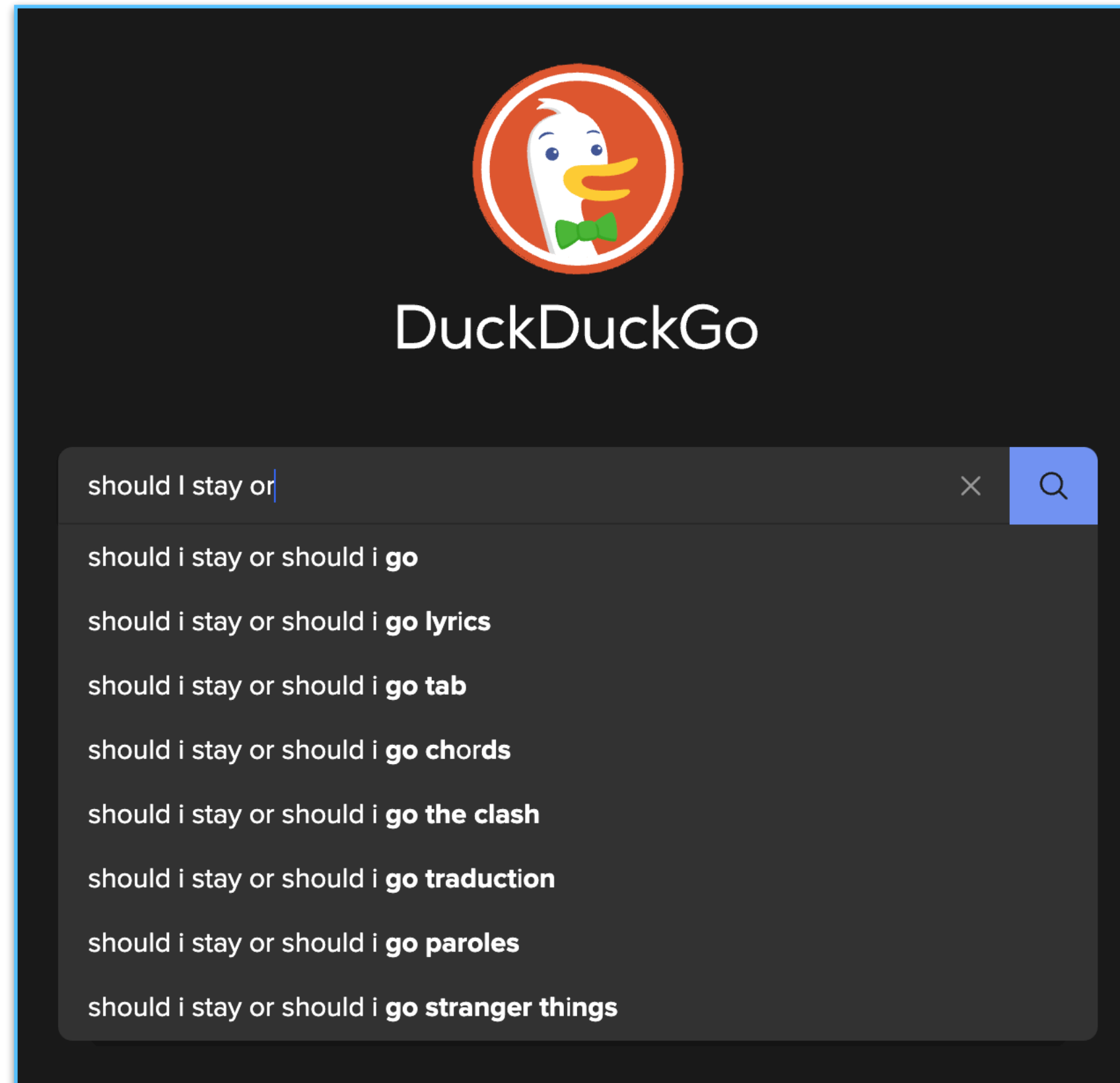
guitar ?

piano

game

drums

## Language model

Given a sequence of words $x_1, x_2, \ldots, x_t$

compute the probability of the next word $p\left(x_{t+1} \mid x_t, x_{t-1}, \ldots, x_1\right)$

where $x_i \in \mathscr{V}$ (a word from our vocabulary)

$$\text{PPL} = \prod_{t=1}^{N} \left( \frac{1}{p_\ell\left(x_{t+1} \mid x_t, \ldots, x_1\right)} \right)^{\frac{1}{N}}$$

number of tokens
in our corpus

lower is
better

inverse probability of the corpus, according to the language model $\ell$

$$PPL = \prod_{t=1}^{N} \left( \frac{1}{p_\ell(x_{t+1} \mid x_t, \ldots, x_1)} \right)^{\frac{1}{N}}$$

number of tokens
in our corpus

lower is
better

inverse probability of the corpus, according to the language model $\ell$

Intuition: if $PPL = \delta$, then our uncertainty about the next word is ~ equivalent to the uncertainty of tossing a $\delta$-sided dice and getting a $\delta$

$$PPL = \prod_{t=1}^{N} \left( \frac{1}{p_\ell\left(x_{t+1} \mid x_t, \ldots, x_1\right)} \right)^{\frac{1}{N}}$$

number of tokens
in our corpus

lower is
better

inverse probability of the corpus, according to the language model $\ell$

Intuition: if $PPL = \delta$, then our uncertainty about the next word is ~ equivalent to the uncertainty of tossing a $\delta$-sided dice and getting a $\delta$

$$PPL = \prod_{t=1}^{N} \left( \frac{1}{\hat{\mathbf{y}}^{[t]}_{x_{t+1}}} \right)^{\frac{1}{N}} = \cdots = \exp\left(L(\theta)\right)$$

the estimated prob. at word $t$ that the next word is $x_{t+1}$ based on the language model

see 3.8 in SLP

cross entropy loss of a language model parametrised by $\theta$

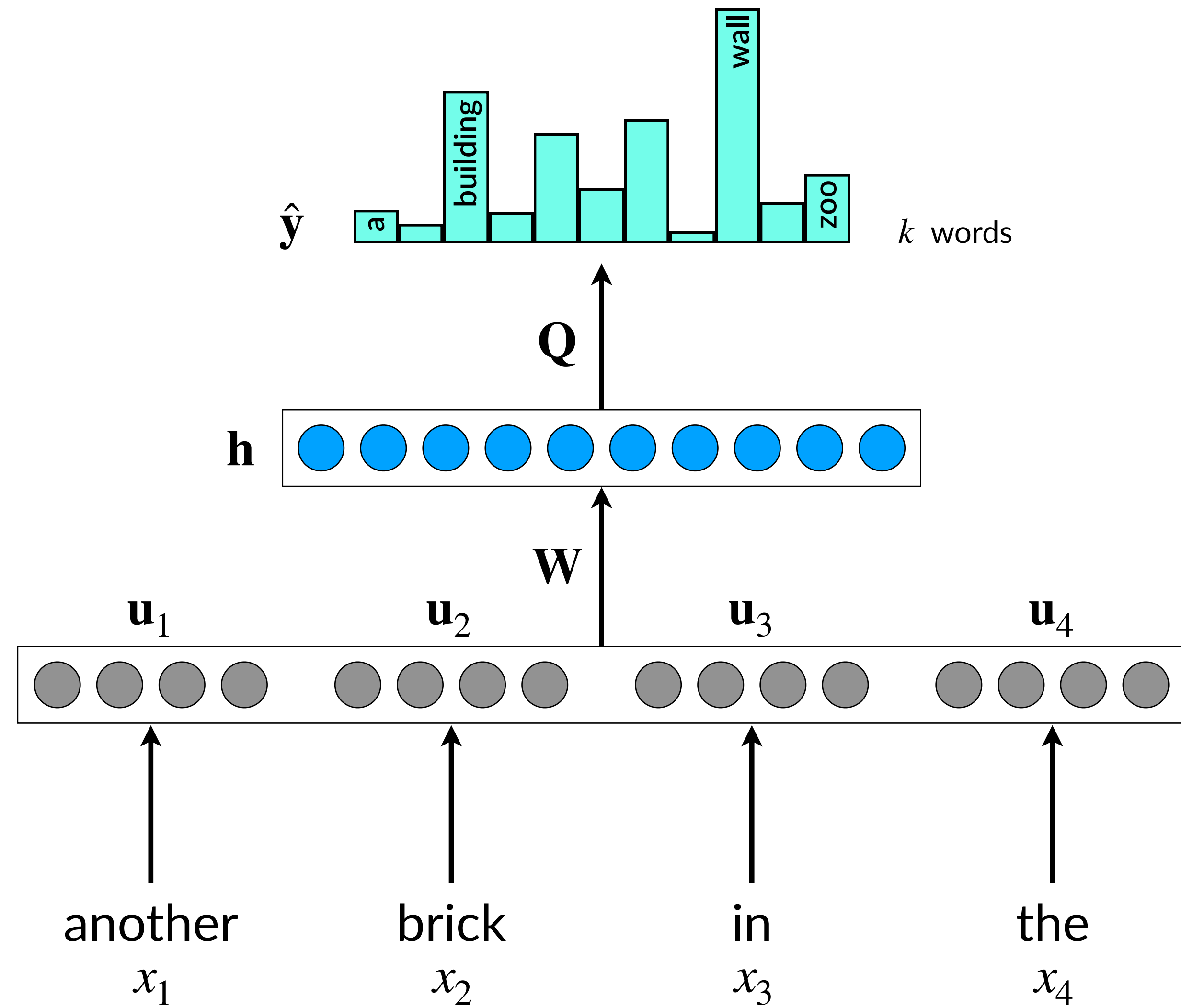| Model | PPL |
|---|---|
| Interpolated Kneser-Ney 5-gram (*2013*) | 67.6 |
| RNN-1024 + MaxEnt 9-gram (*2013*) | 51.3 |
| LSTM-2048 (*2016*) | 43.7 |
| 2-layer LSTM-8192 (*2016*) | 30 |
| Adaptive input Transformer (*2019*) | 23.02 |
| GPT-2 (*2019*) | 16.45 |

*But of course, there is a limit on how low perplexity can realistically be!*

Source 1: engineering.fb.com/2016/10/25/ml-applications/building-an-efficient-neural-language-model-over-a-billion-words/
Source 2: openreview.net/pdf?id=ByxZX20qFQ
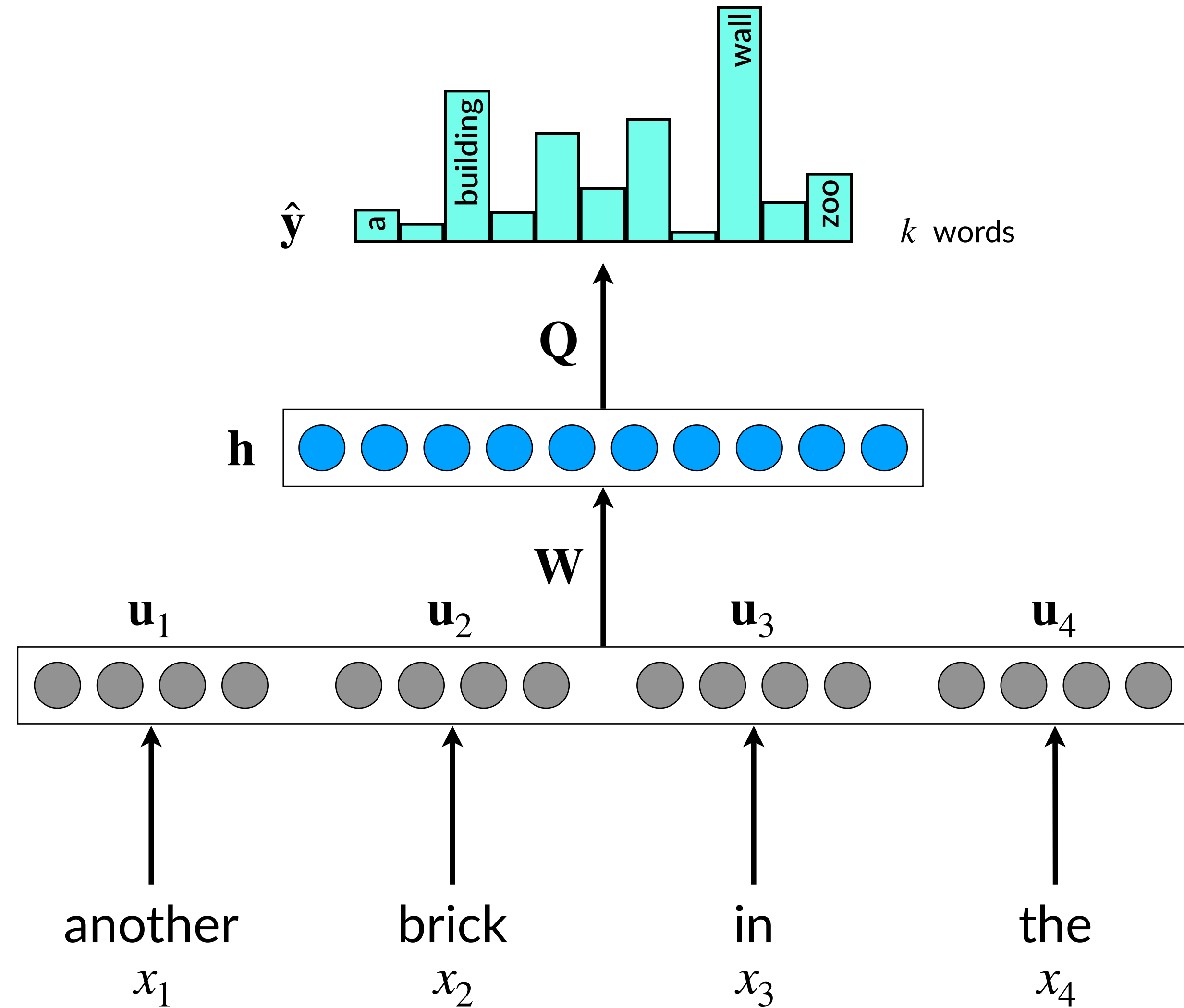Source 3: huggingface.co/docs/transformers/perplexity

$\hat{\mathbf{y}}$    $k$ words

**Q**

**h**

**W**

$\mathbf{u}_1$     $\mathbf{u}_2$     $\mathbf{u}_3$     $\mathbf{u}_4$

another    brick    in    the
$x_1$     $x_2$     $x_3$     $x_4$

$$\mathbf{u} = [\mathbf{u}_1; \mathbf{u}_2; \mathbf{u}_3; \mathbf{u}_4] \quad \in \mathbb{R}^{4d}$$

concatenate
word representations

$\hat{\mathbf{y}}$

$k$ words

$\mathbf{Q}$

$\mathbf{h}$

$\mathbf{W}$

$\mathbf{u}_1 \qquad \mathbf{u}_2 \qquad \mathbf{u}_3 \qquad \mathbf{u}_4$

another     brick     in     the

$x_1 \qquad x_2 \qquad x_3 \qquad x_4$

$$\mathbf{h} = \sigma\big(\mathbf{W} \cdot \mathbf{u} + \mathbf{b}_W\big) \quad \in \mathbb{R}^m$$

$$\mathbf{W} \in \mathbb{R}^{m \times 4d}$$

$$\mathbf{u} = [\mathbf{u}_1; \mathbf{u}_2; \mathbf{u}_3; \mathbf{u}_4] \quad \in \mathbb{R}^{4d}$$

concatenate
word representations

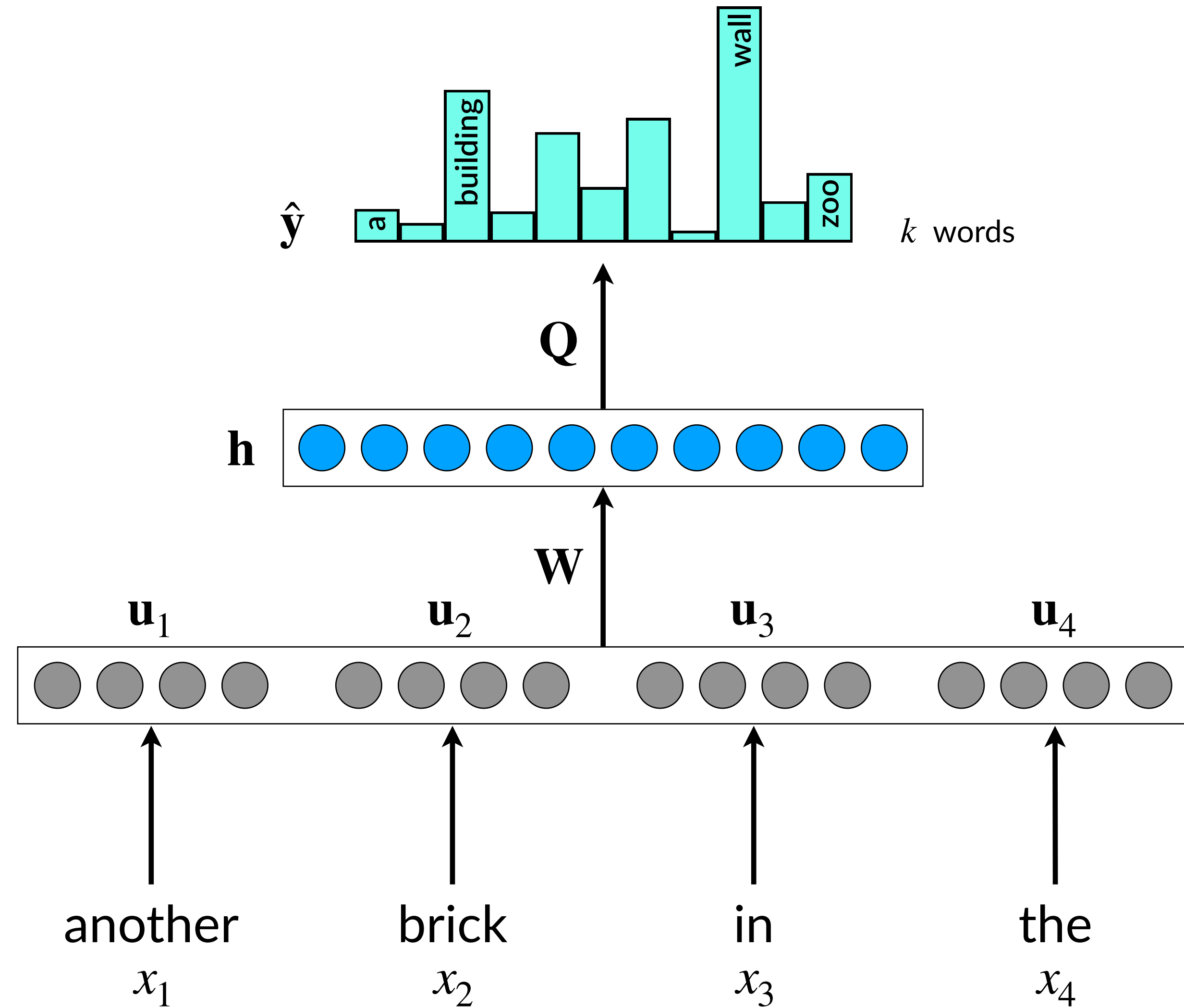$\hat{\mathbf{y}}$    $k$ words

$\mathbf{Q}$

$\mathbf{h}$

$\mathbf{W}$

$\mathbf{u}_1$    $\mathbf{u}_2$    $\mathbf{u}_3$    $\mathbf{u}_4$

another    brick    in    the

$x_1$    $x_2$    $x_3$    $x_4$

$$\hat{\mathbf{y}} = \mathrm{softmax}\left( \mathbf{Q} \cdot \mathbf{h} + \mathbf{b}_Q \right) \quad \in [0,1]^k$$
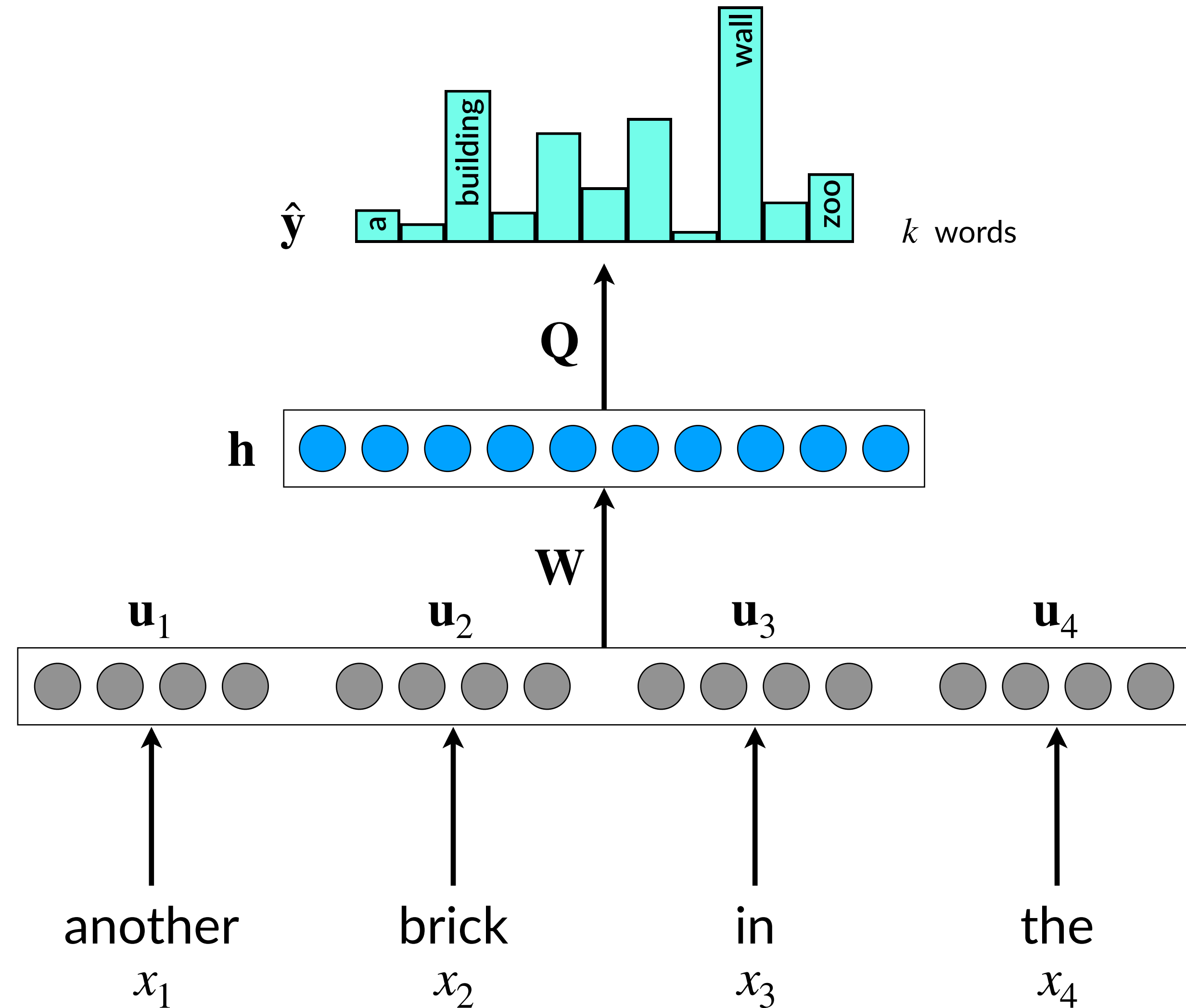
$$\mathbf{Q} \in \mathbb{R}^{k \times m}$$

$$\mathbf{h} = \sigma\left( \mathbf{W} \cdot \mathbf{u} + \mathbf{b}_W \right) \quad \in \mathbb{R}^m$$

$$\mathbf{W} \in \mathbb{R}^{m \times 4d}$$

$$\mathbf{u} = [\mathbf{u}_1; \mathbf{u}_2; \mathbf{u}_3; \mathbf{u}_4] \quad \in \mathbb{R}^{4d}$$
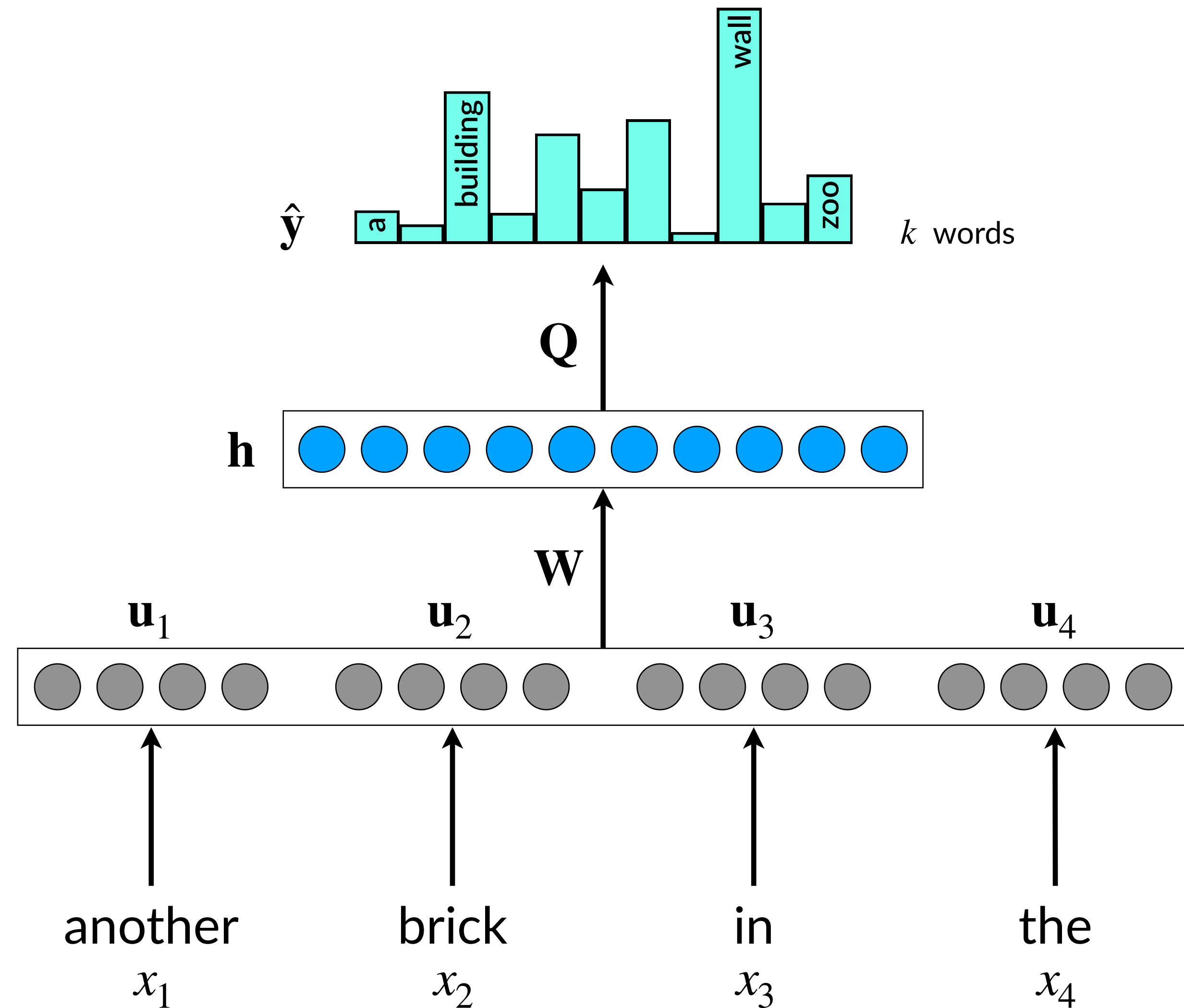
concatenate
word representations

$\hat{\mathbf{y}}$   $k$ words

$\mathbf{Q}$

$\mathbf{h}$

$\mathbf{W}$

$\mathbf{u}_1$   $\mathbf{u}_2$   $\mathbf{u}_3$   $\mathbf{u}_4$

another   brick   in   the
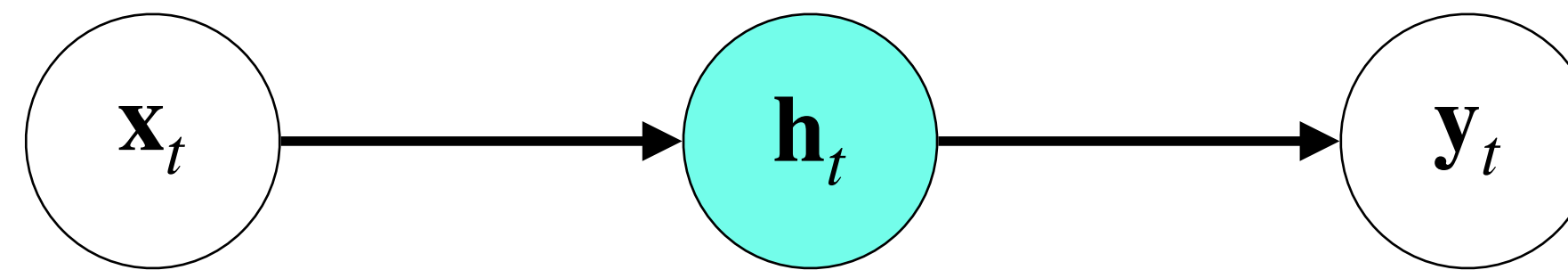$x_1$   $x_2$   $x_3$   $x_4$

$$\hat{\mathbf{y}} = \text{softmax}\left(\mathbf{Q} \cdot \mathbf{h} + \mathbf{b}_Q\right) \quad \in [0,1]^k$$

$$\mathbf{Q} \in \mathbb{R}^{k \times m}$$

$$\mathbf{h} = \sigma\left(\mathbf{W} \cdot \mathbf{u} + \mathbf{b}_W\right) \quad \in \mathbb{R}^m$$

$$\mathbf{W} \in \mathbb{R}^{m \times 4d}$$

$$\mathbf{u} = [\mathbf{u}_1; \mathbf{u}_2; \mathbf{u}_3; \mathbf{u}_4] \quad \in \mathbb{R}^{4d}$$
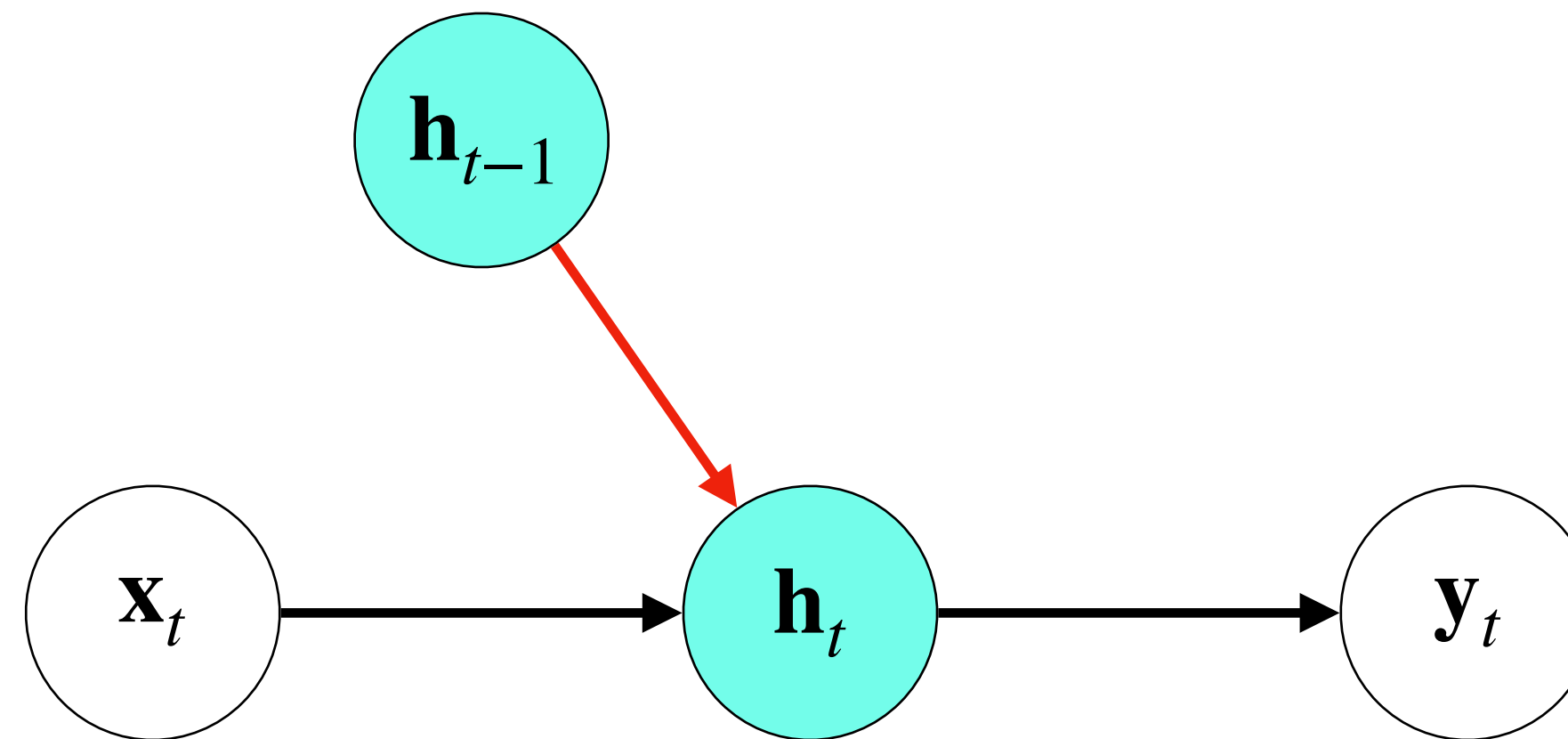
**Issues!**

‣ context / window size is fixed

‣ $\mathbf{W}$ grows if we increase the window

‣ word position is modelled explicitly and independently, i.e. there is no weight sharing between words

**Recurrency**
*The current hidden state $\mathbf{h}_t$ depends on the previous hidden state $\mathbf{h}_{t-1}$ and influences the next hidden state $\mathbf{h}_{t+1}$*
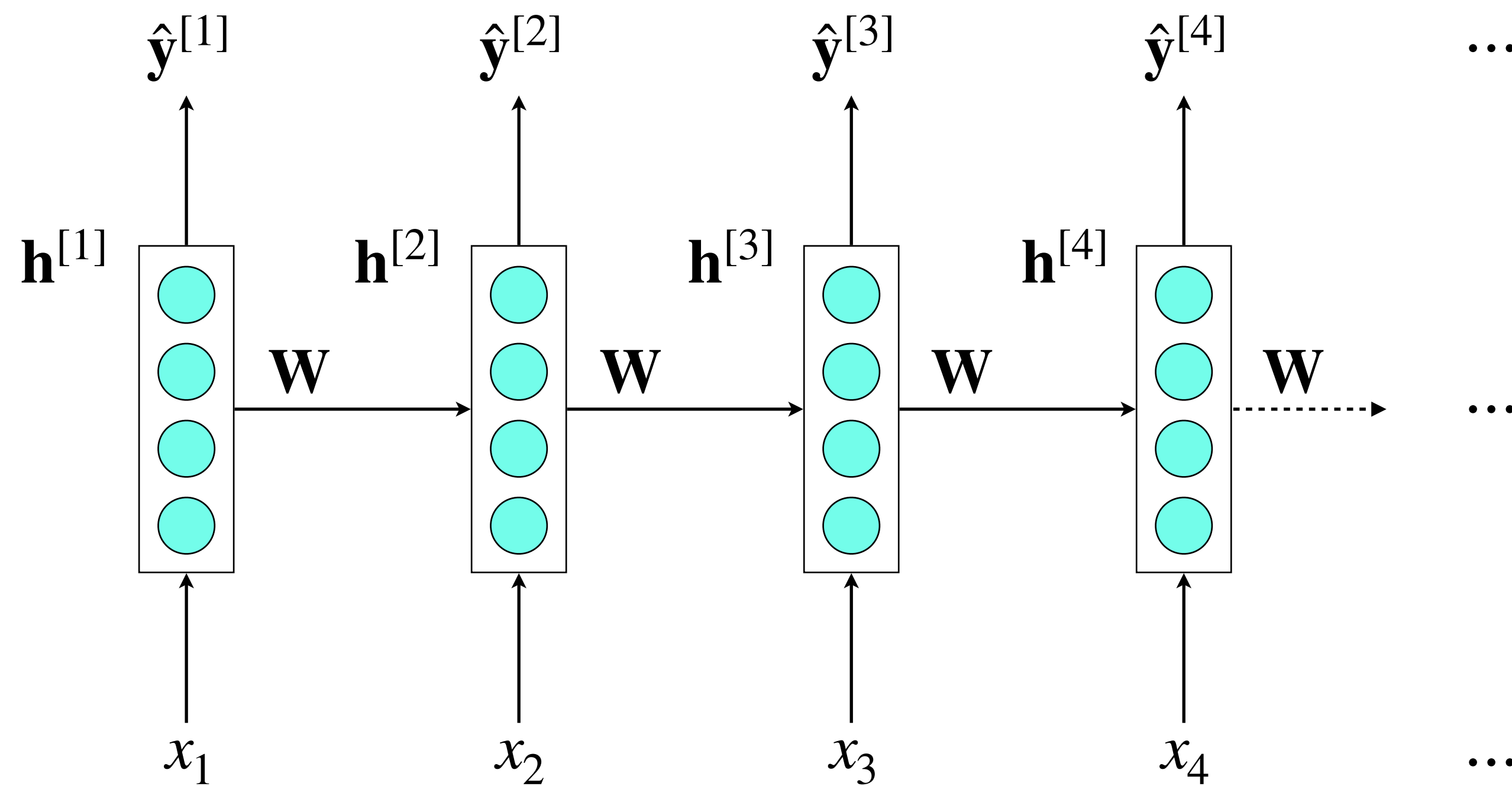
**Recurrency**

*The current hidden state $\mathbf{h}_t$ depends on the previous hidden state $\mathbf{h}_{t-1}$ and influences the next hidden state $\mathbf{h}_{t+1}$*

*The RNN unrolls to a theoretically unlimited number of time steps*

output at each
time step

$\hat{\mathbf{y}}^{[1]}$ $\hat{\mathbf{y}}^{[2]}$ $\hat{\mathbf{y}}^{[3]}$ $\hat{\mathbf{y}}^{[4]}$ ...

*optional*

hidden states
of the RNN
i.e. the *RNN layer*

$\mathbf{h}^{[1]}$ $\mathbf{h}^{[2]}$ $\mathbf{h}^{[3]}$ $\mathbf{h}^{[4]}$

$\mathbf{W}$ $\mathbf{W}$ $\mathbf{W}$ $\mathbf{W}$ ...

core property!

$$\mathbf{W}$$
*remains
the same*

sequence
of inputs

$x_1$ $x_2$ $x_3$ $x_4$ ...

*in theory* $\infty$

another   brick   in   the

$x_1$   $x_2$   $x_3$   $x_4$

*in theory sequence length can be unlimited*

$\mathbf{u}^{[1]}$  $\mathbf{u}^{[2]}$  $\mathbf{u}^{[3]}$  $\mathbf{u}^{[4]}$

$n$  $n$  $n$  $n$

another   brick   in   the

$x_1$   $x_2$   $x_3$   $x_4$

*in theory sequence length can be unlimited*

initial hidden state

$\mathbf{h}^{[0]}$     $\mathbf{h}^{[1]}$

$\mathbf{W}_h$

$m$     $m$

$\mathbf{W}_u$

$\mathbf{u}^{[1]}$     $\mathbf{u}^{[2]}$     $\mathbf{u}^{[3]}$     $\mathbf{u}^{[4]}$

$n$     $n$     $n$     $n$

another     brick     in     the

$x_1$     $x_2$     $x_3$     $x_4$

in theory sequence length can be unlimited

initial hidden state

$\mathbf{h}^{[0]}$     $\mathbf{h}^{[1]}$     $\mathbf{h}^{[2]}$

$\mathbf{W}_h$     $\mathbf{W}_h$

$m$     $m$     $m$

$\mathbf{W}_u$     $\mathbf{W}_u$

$\mathbf{u}^{[1]}$     $\mathbf{u}^{[2]}$     $\mathbf{u}^{[3]}$     $\mathbf{u}^{[4]}$

$n$     $n$     $n$     $n$

another     brick     in     the

$x_1$     $x_2$     $x_3$     $x_4$

in theory sequence length can be unlimited

$$\hat{\mathbf{y}} = p\left(x_5 \mid x_1, \ldots, x_4\right)$$

*initial hidden state*

zoo
wall
building
a

$k$ words

$\mathbf{h}^{[0]}$ $\mathbf{h}^{[1]}$ $\mathbf{h}^{[2]}$ $\mathbf{h}^{[3]}$ $\mathbf{h}^{[4]}$

$\mathbf{W}_y$

$m$ $\quad \mathbf{W}_h \quad$ $m$ $\quad \mathbf{W}_h \quad$ $m$ $\quad \mathbf{W}_h \quad$ $m$ $\quad \mathbf{W}_h \quad$ $m$

$\mathbf{W}_u$ $\mathbf{W}_u$ $\mathbf{W}_u$ $\mathbf{W}_u$

$\mathbf{u}^{[1]}$ $\mathbf{u}^{[2]}$ $\mathbf{u}^{[3]}$ $\mathbf{u}^{[4]}$

$n$ $n$ $n$ $n$

another brick in the

$x_1$ $x_2$ $x_3$ $x_4$

*in theory sequence length can be unlimited*

$$\hat{\mathbf{y}} = p\big(x_5 \,|\, x_1, \ldots, x_4\big)$$

*initial hidden state*
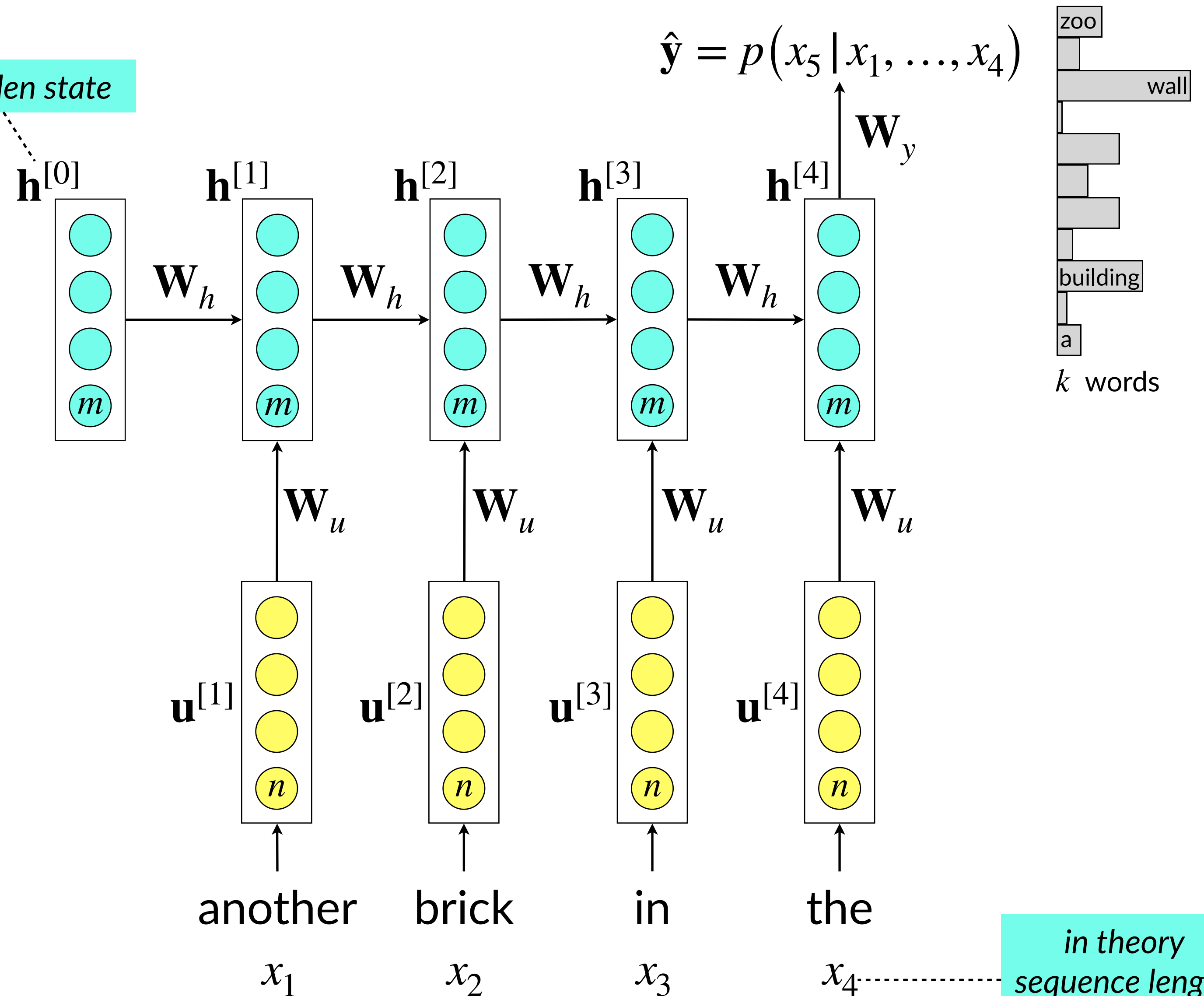
**Hidden states**

$$\mathbf{h}^{[t]} = \sigma\big(\mathbf{W}_u \cdot \mathbf{u}^{[t]} + \mathbf{W}_h \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_h\big)$$

or use $\tanh(\cdot)$

$\mathbf{h}^{[0]}$  $\mathbf{h}^{[1]}$  $\mathbf{h}^{[2]}$  $\mathbf{h}^{[3]}$  $\mathbf{h}^{[4]}$

$\mathbf{W}_h$  $\mathbf{W}_h$  $\mathbf{W}_h$  $\mathbf{W}_h$

$\mathbf{W}_y$

$m$  $m$  $m$  $m$  $m$

zoo
wall
building
a

$k$ words

$\mathbf{W}_u$  $\mathbf{W}_u$  $\mathbf{W}_u$  $\mathbf{W}_u$

$\mathbf{u}^{[1]}$  $\mathbf{u}^{[2]}$  $\mathbf{u}^{[3]}$  $\mathbf{u}^{[4]}$

$n$  $n$  $n$  $n$

another    brick    in    the

$x_1$    $x_2$    $x_3$    $x_4$

*in theory sequence length can be unlimited*

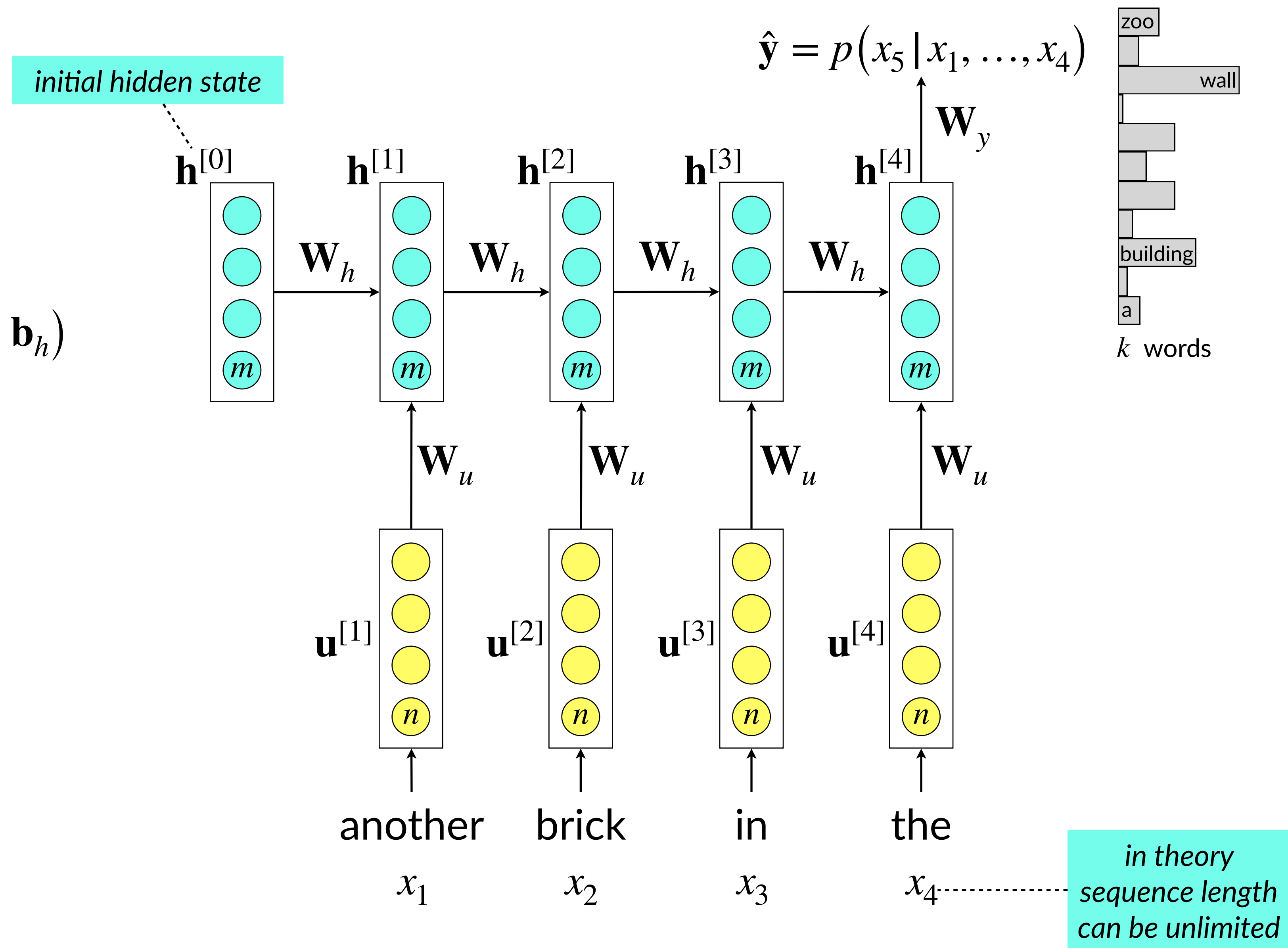**Output**

$$\hat{\mathbf{y}} = \text{softmax}\left(\mathbf{W}_y \cdot \mathbf{h}^{[4]} + \mathbf{b}_y\right)$$
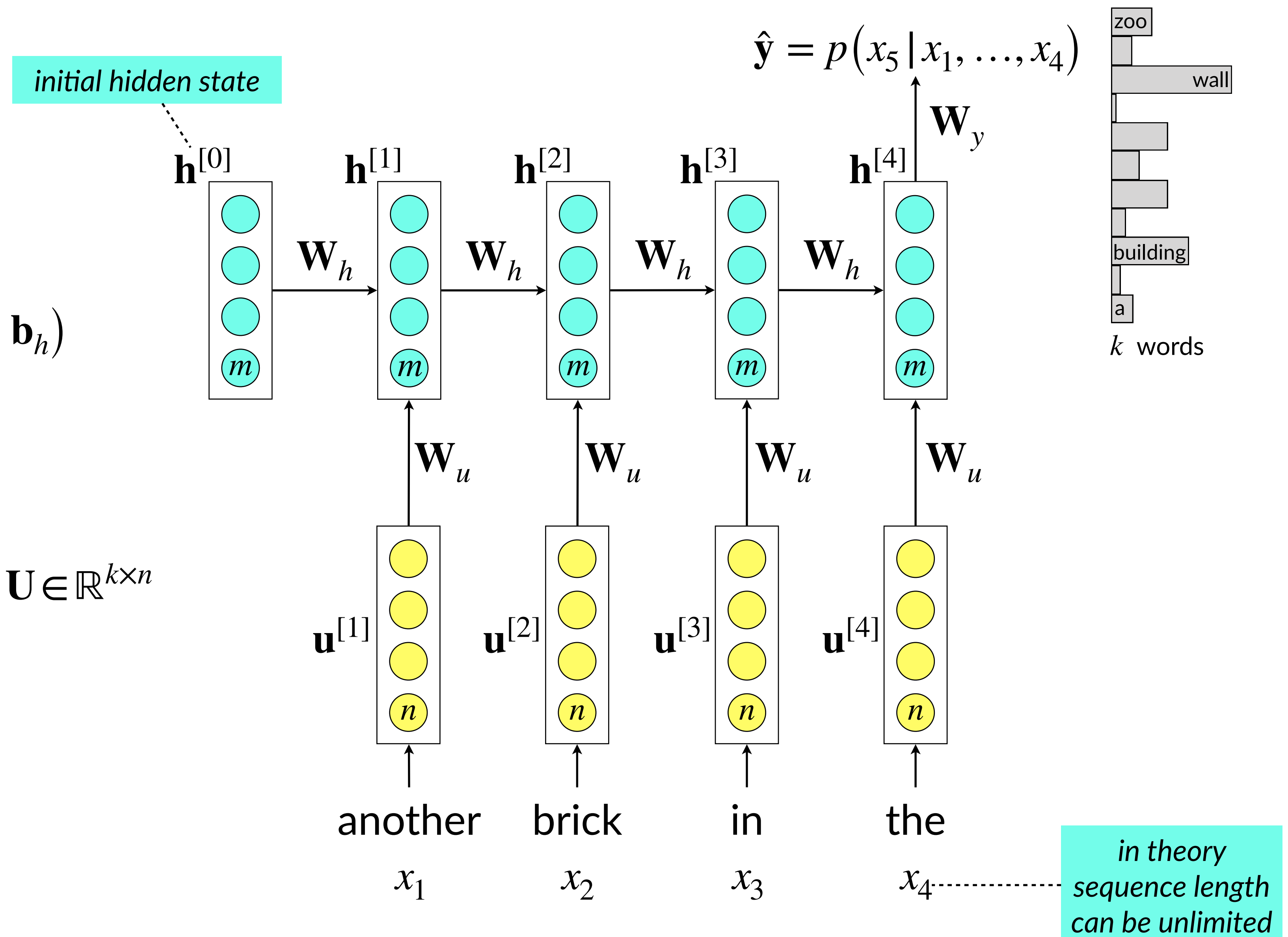
**Hidden states**

$$\mathbf{h}^{[t]} = \sigma\left(\mathbf{W}_u \cdot \mathbf{u}^{[t]} + \mathbf{W}_h \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_h\right)$$

or use $\tanh(\cdot)$

$$\hat{\mathbf{y}} = p\left(x_5 \,|\, x_1, \ldots, x_4\right)$$

*initial hidden state*

$\mathbf{h}^{[0]}$  $\mathbf{h}^{[1]}$  $\mathbf{h}^{[2]}$  $\mathbf{h}^{[3]}$  $\mathbf{h}^{[4]}$  $\mathbf{W}_y$

$\mathbf{W}_h$  $\mathbf{W}_h$  $\mathbf{W}_h$  $\mathbf{W}_h$

$m$  $m$  $m$  $m$  $m$

zoo

wall

building

a

$k$ words

$\mathbf{W}_u$  $\mathbf{W}_u$  $\mathbf{W}_u$  $\mathbf{W}_u$

$\mathbf{u}^{[1]}$  $\mathbf{u}^{[2]}$  $\mathbf{u}^{[3]}$  $\mathbf{u}^{[4]}$

$n$  $n$  $n$  $n$

another    brick    in    the

$x_1$    $x_2$    $x_3$    $x_4$

*in theory sequence length can be unlimited*

**Output**

$$\hat{\mathbf{y}} = \mathrm{softmax}\left(\mathbf{W}_y \cdot \mathbf{h}^{[4]} + \mathbf{b}_y\right)$$

**Hidden states**

$$\mathbf{h}^{[t]} = \sigma\left(\mathbf{W}_u \cdot \mathbf{u}^{[t]} + \mathbf{W}_h \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_h\right)$$
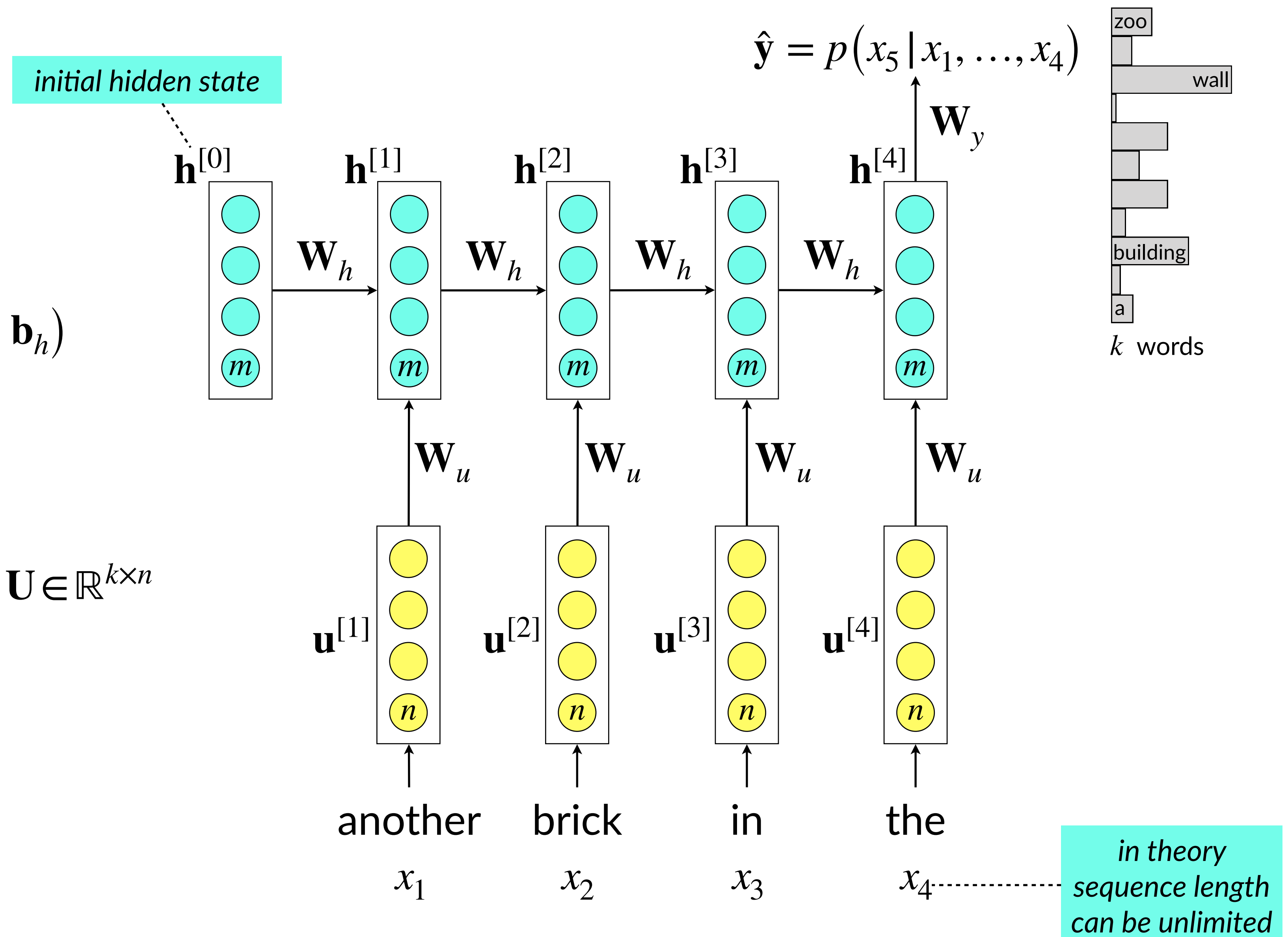
or use $\tanh(\cdot)$

**Dimensionalities?**

$\mathbf{u}^{[t]} \in \mathbb{R}^n$ embedding of $x_t$ from $\mathbf{U} \in \mathbb{R}^{k \times n}$

$\mathbf{h}^{[t]}, \mathbf{b}_h \in \mathbb{R}^m$

*initial hidden state*

$\hat{\mathbf{y}} = p\left(x_5 \mid x_1, \ldots, x_4\right)$

zoo

wall

building

a

$k$ words

$\mathbf{h}^{[0]}$  $\mathbf{h}^{[1]}$  $\mathbf{h}^{[2]}$  $\mathbf{h}^{[3]}$  $\mathbf{h}^{[4]}$

$\mathbf{W}_y$

$\mathbf{W}_h$  $\mathbf{W}_h$  $\mathbf{W}_h$  $\mathbf{W}_h$

$m$  $m$  $m$  $m$  $m$

$\mathbf{W}_u$  $\mathbf{W}_u$  $\mathbf{W}_u$  $\mathbf{W}_u$

$\mathbf{u}^{[1]}$  $\mathbf{u}^{[2]}$  $\mathbf{u}^{[3]}$  $\mathbf{u}^{[4]}$

$n$  $n$  $n$  $n$

another  brick  in  the

$x_1$  $x_2$  $x_3$  $x_4$

*in theory sequence length can be unlimited*

## Output

$$\hat{\mathbf{y}} = \mathrm{softmax}\left(\mathbf{W}_y \cdot \mathbf{h}^{[4]} + \mathbf{b}_y\right)$$

## Hidden states

$$\mathbf{h}^{[t]} = \sigma\left(\mathbf{W}_u \cdot \mathbf{u}^{[t]} + \mathbf{W}_h \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_h\right)$$

or use $\tanh(\cdot)$

## Dimensionalities?

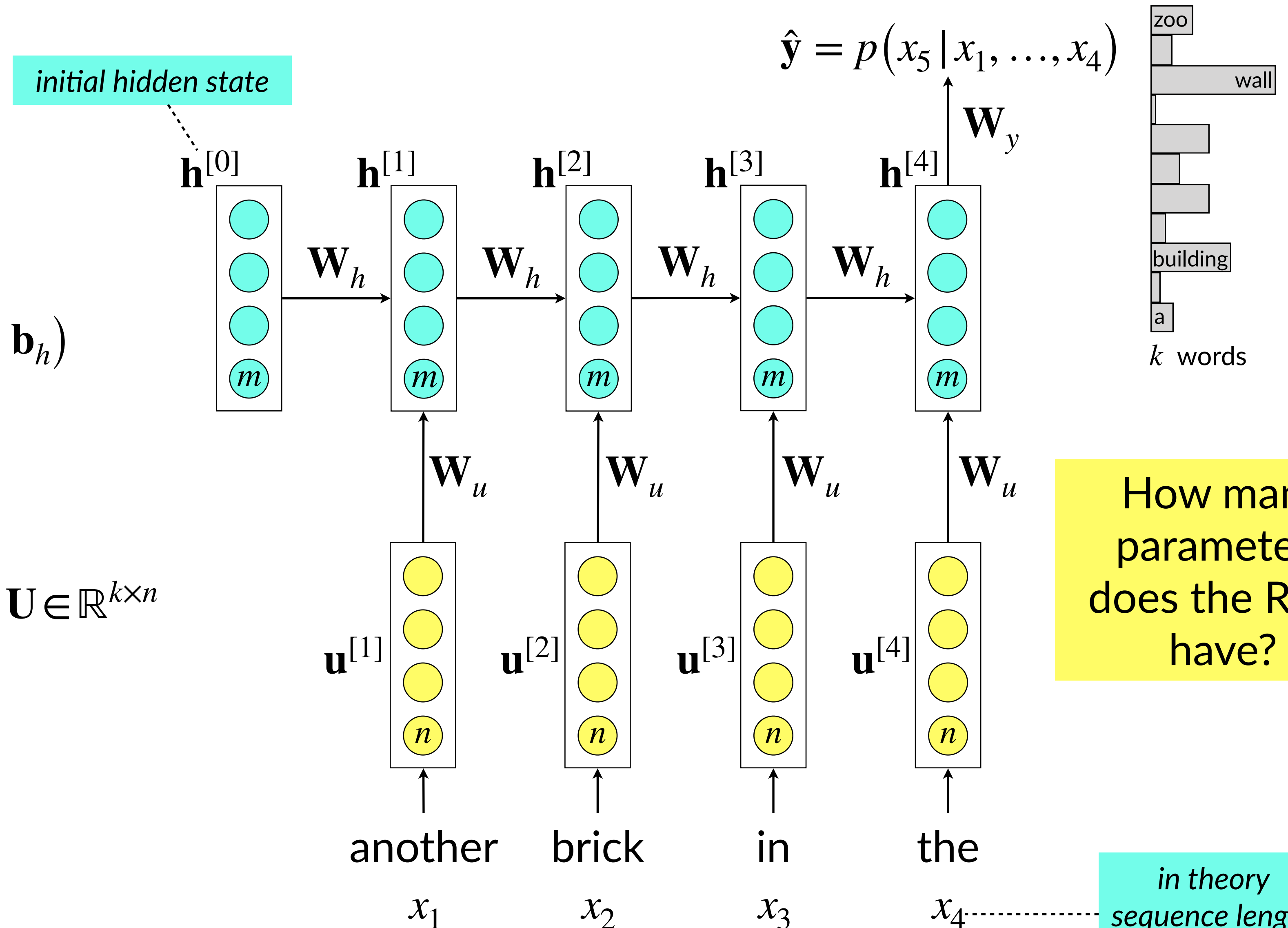$\mathbf{u}^{[t]} \in \mathbb{R}^n$ embedding of $x_t$ from $\mathbf{U} \in \mathbb{R}^{k \times n}$
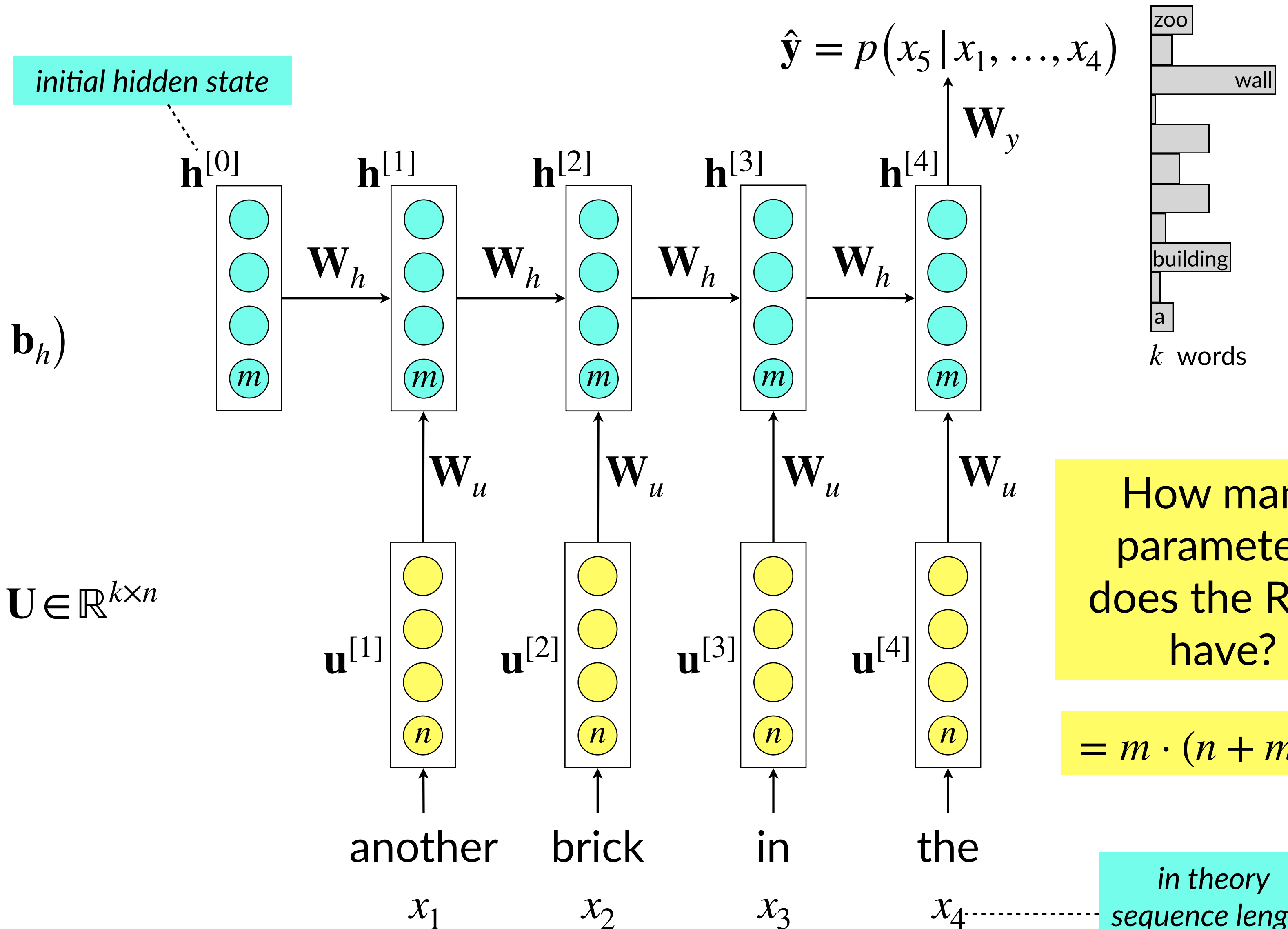
$\mathbf{h}^{[t]}, \mathbf{b}_h \in \mathbb{R}^m$

$\hat{\mathbf{y}}, \mathbf{b}_y \in \mathbb{R}^k$

$\mathbf{W}_u \in \mathbb{R}^{m \times n}, \mathbf{W}_h \in \mathbb{R}^{m \times m}$

$\mathbf{W}_y \in \mathbb{R}^{k \times m}$

initial hidden state

$$\hat{\mathbf{y}} = p\left(x_5 \mid x_1, \dots, x_4\right)$$

$\mathbf{h}^{[0]}$ $\mathbf{h}^{[1]}$ $\mathbf{h}^{[2]}$ $\mathbf{h}^{[3]}$ $\mathbf{h}^{[4]}$

$\mathbf{W}_h$ $\mathbf{W}_h$ $\mathbf{W}_h$ $\mathbf{W}_h$

$\mathbf{W}_y$

$m$ $m$ $m$ $m$ $m$

$\mathbf{W}_u$ $\mathbf{W}_u$ $\mathbf{W}_u$ $\mathbf{W}_u$

$\mathbf{u}^{[1]}$ $\mathbf{u}^{[2]}$ $\mathbf{u}^{[3]}$ $\mathbf{u}^{[4]}$

$n$ $n$ $n$ $n$

another   brick   in   the

$x_1$   $x_2$   $x_3$   $x_4$

zoo
wall
building
a

$k$ words

in theory
sequence length
can be unlimited

**Output**

$$\hat{\mathbf{y}} = \mathrm{softmax}\left(\mathbf{W}_y \cdot \mathbf{h}^{[4]} + \mathbf{b}_y\right)$$

**Hidden states**

$$\mathbf{h}^{[t]} = \sigma\left(\mathbf{W}_u \cdot \mathbf{u}^{[t]} + \mathbf{W}_h \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_h\right)$$

or use $\tanh(\cdot)$

**Dimensionalities?**

$\mathbf{u}^{[t]} \in \mathbb{R}^n$ embedding of $x_t$ from $\mathbf{U} \in \mathbb{R}^{k \times n}$

$\mathbf{h}^{[t]}, \mathbf{b}_h \in \mathbb{R}^m$

$\hat{\mathbf{y}}, \mathbf{b}_y \in \mathbb{R}^k$

$\mathbf{W}_u \in \mathbb{R}^{m \times n}, \mathbf{W}_h \in \mathbb{R}^{m \times m}$

$\mathbf{W}_y \in \mathbb{R}^{k \times m}$



$\hat{\mathbf{y}} = p\left(x_5 \mid x_1, \ldots, x_4\right)$

initial hidden state

How many parameters does the RNN have?

in theory sequence length can be unlimited

**Output**

$$\hat{\mathbf{y}} = \text{softmax}\left(\mathbf{W}_y \cdot \mathbf{h}^{[4]} + \mathbf{b}_y\right)$$

**Hidden states**

$$\mathbf{h}^{[t]} = \sigma\left(\mathbf{W}_u \cdot \mathbf{u}^{[t]} + \mathbf{W}_h \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_h\right)$$

or use $\tanh(\,\cdot\,)$

**Dimensionalities?**

$\mathbf{u}^{[t]} \in \mathbb{R}^n$ embedding of $x_t$ from $\mathbf{U} \in \mathbb{R}^{k \times n}$

$\mathbf{h}^{[t]}, \mathbf{b}_h \in \mathbb{R}^m$

$\hat{\mathbf{y}}, \mathbf{b}_y \in \mathbb{R}^k$

$\mathbf{W}_u \in \mathbb{R}^{m \times n}, \mathbf{W}_h \in \mathbb{R}^{m \times m}$

$\mathbf{W}_y \in \mathbb{R}^{k \times m}$

*initial hidden state*

$$\hat{\mathbf{y}} = p\left(x_5 \mid x_1, \ldots, x_4\right)$$

zoo

wall

building

a

$k$ words

$\mathbf{h}^{[0]}$ $\xrightarrow{\mathbf{W}_h}$ $\mathbf{h}^{[1]}$ $\xrightarrow{\mathbf{W}_h}$ $\mathbf{h}^{[2]}$ $\xrightarrow{\mathbf{W}_h}$ $\mathbf{h}^{[3]}$ $\xrightarrow{\mathbf{W}_h}$ $\mathbf{h}^{[4]}$ $\xrightarrow{\mathbf{W}_y}$

$m$ $m$ $m$ $m$ $m$

$\mathbf{W}_u$ $\mathbf{W}_u$ $\mathbf{W}_u$ $\mathbf{W}_u$

$\mathbf{u}^{[1]}$ $\mathbf{u}^{[2]}$ $\mathbf{u}^{[3]}$ $\mathbf{u}^{[4]}$

$n$ $n$ $n$ $n$

another    brick    in    the

$x_1$    $x_2$    $x_3$    $x_4$

**How many parameters does the RNN have?**

$$= m \cdot (n + m + 1)$$

*in theory sequence length can be unlimited*

$\theta = [\mathbf{W}_u, \mathbf{W}_h, \mathbf{W}_y]$

*Loss at each time step*

*at each time step we predict a prob. dist. over the entire corpus and determine the most probable next word*

$L^{[1]}(\theta)$

$L^{[2]}(\theta)$

$L^{[3]}(\theta)$

$L^{[4]}(\theta)$

$\hat{\mathbf{y}}^{[1]}$

$\hat{\mathbf{y}}^{[2]}$

$\hat{\mathbf{y}}^{[3]}$

$\hat{\mathbf{y}}^{[4]}$

$\mathbf{W}_y$

$\mathbf{W}_y$

$\mathbf{W}_y$

$\mathbf{W}_y$

$\mathbf{h}^{[0]}$

$\mathbf{h}^{[1]}$

$\mathbf{h}^{[2]}$

$\mathbf{h}^{[3]}$

$\mathbf{h}^{[4]}$

$\mathbf{W}_h$

$\mathbf{W}_h$

$\mathbf{W}_h$

$\mathbf{W}_h$

$m$

$\mathbf{W}_u$

$\mathbf{W}_u$

$\mathbf{W}_u$

$\mathbf{W}_u$

$\mathbf{u}^{[1]}$

$\mathbf{u}^{[2]}$

$\mathbf{u}^{[3]}$

$\mathbf{u}^{[4]}$

$n$

*text sequence / corpus*

| another | brick | in | the | Wall | ... |

$x_1$

$x_2$

$x_3$

$x_4$

$\theta = [\mathbf{W}_u, \mathbf{W}_h, \mathbf{W}_y]$

*Loss at each time step*

$-\log \hat{y}^{[1]}(\text{"brick"})$

$\|$

$L^{[1]}(\theta)$      $L^{[2]}(\theta)$      $L^{[3]}(\theta)$      $L^{[4]}(\theta)$

*at each time step we predict a prob. dist. over the entire corpus and determine the most probable next word*

$\hat{\mathbf{y}}^{[1]}$      $\hat{\mathbf{y}}^{[2]}$      $\hat{\mathbf{y}}^{[3]}$      $\hat{\mathbf{y}}^{[4]}$

$\mathbf{W}_y$      $\mathbf{W}_y$      $\mathbf{W}_y$      $\mathbf{W}_y$

$\mathbf{h}^{[0]}$   $\mathbf{h}^{[1]}$   $\mathbf{h}^{[2]}$   $\mathbf{h}^{[3]}$   $\mathbf{h}^{[4]}$

$m$

$\mathbf{W}_h$      $\mathbf{W}_h$      $\mathbf{W}_h$      $\mathbf{W}_h$

$m$      $m$      $m$      $m$

$\mathbf{W}_u$      $\mathbf{W}_u$      $\mathbf{W}_u$      $\mathbf{W}_u$

$\mathbf{u}^{[1]}$      $\mathbf{u}^{[2]}$      $\mathbf{u}^{[3]}$      $\mathbf{u}^{[4]}$

$n$      $n$      $n$      $n$

*text sequence / corpus*

another      brick      in      the      Wall      ...

$x_1$      $x_2$      $x_3$      $x_4$

$\theta = \left[ \mathbf{W}_u, \mathbf{W}_h, \mathbf{W}_y \right]$

*Loss at each time step*

$-\log \hat{y}^{[1]}(\text{"brick"})$

$-\log \hat{y}^{[2]}(\text{"in"})$

$\|$

$\|$

$L^{[1]}(\theta)$

$L^{[2]}(\theta)$

$L^{[3]}(\theta)$

$L^{[4]}(\theta)$

*at each time step we predict a prob. dist. over the entire corpus and determine the most probable next word*

$\hat{\mathbf{y}}^{[1]}$

$\hat{\mathbf{y}}^{[2]}$

$\hat{\mathbf{y}}^{[3]}$

$\hat{\mathbf{y}}^{[4]}$

$\mathbf{W}_y$

$\mathbf{W}_y$

$\mathbf{W}_y$

$\mathbf{W}_y$

$\mathbf{h}^{[0]}$

$\mathbf{h}^{[1]}$

$\mathbf{h}^{[2]}$

$\mathbf{h}^{[3]}$

$\mathbf{h}^{[4]}$

$m$

$\mathbf{W}_h$

$\mathbf{W}_h$

$\mathbf{W}_h$

$\mathbf{W}_h$

$m$

$m$

$m$

$m$

$\mathbf{W}_u$

$\mathbf{W}_u$

$\mathbf{W}_u$

$\mathbf{W}_u$

$\mathbf{u}^{[1]}$

$\mathbf{u}^{[2]}$

$\mathbf{u}^{[3]}$

$\mathbf{u}^{[4]}$

$n$

$n$

$n$

$n$

*text sequence / corpus*

another

brick

in

the

Wall

...

$x_1$

$x_2$

$x_3$

$x_4$

$\theta = [\mathbf{W}_u, \mathbf{W}_h, \mathbf{W}_y]$

*Loss at each time step*

$-\log \hat{y}^{[1]}(\text{"brick"})$

$-\log \hat{y}^{[2]}(\text{"in"})$

$-\log \hat{y}^{[3]}(\text{"the"})$

‖

‖

‖

$L^{[1]}(\theta)$

$L^{[2]}(\theta)$

$L^{[3]}(\theta)$

$L^{[4]}(\theta)$

*at each time step we predict a prob. dist. over the entire corpus and determine the most probable next word*

$\hat{\mathbf{y}}^{[1]}$

$\hat{\mathbf{y}}^{[2]}$

$\hat{\mathbf{y}}^{[3]}$

$\hat{\mathbf{y}}^{[4]}$

$\mathbf{W}_y$

$\mathbf{W}_y$

$\mathbf{W}_y$

$\mathbf{W}_y$

$\mathbf{h}^{[0]}$

$\mathbf{h}^{[1]}$

$\mathbf{h}^{[2]}$

$\mathbf{h}^{[3]}$

$\mathbf{h}^{[4]}$

$\mathbf{W}_h$

$\mathbf{W}_h$

$\mathbf{W}_h$

$\mathbf{W}_h$

$m$

$m$

$m$

$m$

$m$

$\mathbf{W}_u$

$\mathbf{W}_u$

$\mathbf{W}_u$

$\mathbf{W}_u$

$\mathbf{u}^{[1]}$

$\mathbf{u}^{[2]}$

$\mathbf{u}^{[3]}$

$\mathbf{u}^{[4]}$

$n$

$n$

$n$

$n$

*text sequence / corpus*

another

brick

in

the

Wall

...

$x_1$

$x_2$

$x_3$

$x_4$

$\theta = [\mathbf{W}_u, \mathbf{W}_h, \mathbf{W}_y]$

*Loss at each time step*

$-\log \hat{y}^{[1]}(\text{"brick"})$

$-\log \hat{y}^{[2]}(\text{"in"})$

$-\log \hat{y}^{[3]}(\text{"the"})$

$-\log \hat{y}^{[4]}(\text{"Wall"})$

$\|$

$\|$

$\|$

$\|$

$L^{[1]}(\theta)$

$L^{[2]}(\theta)$

$L^{[3]}(\theta)$

$L^{[4]}(\theta)$

*at each time step we predict a prob. dist. over the entire corpus and determine the most probable next word*

$\hat{\mathbf{y}}^{[1]}$

$\hat{\mathbf{y}}^{[2]}$

$\hat{\mathbf{y}}^{[3]}$

$\hat{\mathbf{y}}^{[4]}$

$\mathbf{W}_y$

$\mathbf{W}_y$

$\mathbf{W}_y$

$\mathbf{W}_y$

$\mathbf{h}^{[0]}$

$\mathbf{h}^{[1]}$

$\mathbf{h}^{[2]}$

$\mathbf{h}^{[3]}$

$\mathbf{h}^{[4]}$

$m$

$m$

$m$

$m$

$m$

$\mathbf{W}_h$

$\mathbf{W}_h$

$\mathbf{W}_h$

$\mathbf{W}_h$

$\mathbf{W}_u$

$\mathbf{W}_u$

$\mathbf{W}_u$

$\mathbf{W}_u$

$\mathbf{u}^{[1]}$

$\mathbf{u}^{[2]}$

$\mathbf{u}^{[3]}$

$\mathbf{u}^{[4]}$

$n$

$n$

$n$

$n$

*text sequence / corpus*

another

brick

in

the

Wall

...

$x_1$

$x_2$

$x_3$

$x_4$

$\theta = [\mathbf{W}_u, \mathbf{W}_h, \mathbf{W}_y]$

*Loss at each time step*

$-\log \hat{y}^{[1]}(\text{"brick"})$ $+$ $-\log \hat{y}^{[2]}(\text{"in"})$ $+$ $-\log \hat{y}^{[3]}(\text{"the"})$ $+$ $-\log \hat{y}^{[4]}(\text{"Wall"})$ $+$ ... $= L(\theta) = \dfrac{1}{T} \sum_{t=1}^{T} L^{[t]}(\theta)$

$\parallel$ $\qquad\qquad\qquad$ $\parallel$ $\qquad\qquad\qquad$ $\parallel$ $\qquad\qquad\qquad$ $\parallel$

$L^{[1]}(\theta)$ $\qquad\qquad$ $L^{[2]}(\theta)$ $\qquad\qquad$ $L^{[3]}(\theta)$ $\qquad\qquad$ $L^{[4]}(\theta)$

*at each time step we predict a prob. dist. over the entire corpus and determine the most probable next word*

*cumulative cross-entropy loss, i.e. the mean loss across all time steps*

$\hat{\mathbf{y}}^{[1]}$ $\qquad\qquad$ $\hat{\mathbf{y}}^{[2]}$ $\qquad\qquad$ $\hat{\mathbf{y}}^{[3]}$ $\qquad\qquad$ $\hat{\mathbf{y}}^{[4]}$

$\mathbf{W}_y$ $\qquad\qquad$ $\mathbf{W}_y$ $\qquad\qquad$ $\mathbf{W}_y$ $\qquad\qquad$ $\mathbf{W}_y$

$\mathbf{h}^{[0]}$ $\quad$ $\mathbf{h}^{[1]}$ $\qquad\quad$ $\mathbf{h}^{[2]}$ $\qquad\quad$ $\mathbf{h}^{[3]}$ $\qquad\quad$ $\mathbf{h}^{[4]}$

$\mathbf{W}_h$ $\qquad\quad$ $\mathbf{W}_h$ $\qquad\quad$ $\mathbf{W}_h$ $\qquad\quad$ $\mathbf{W}_h$

$m$ $\qquad m \qquad\quad m \qquad\quad m \qquad\quad m$

$\mathbf{W}_u$ $\qquad\qquad$ $\mathbf{W}_u$ $\qquad\qquad$ $\mathbf{W}_u$ $\qquad\qquad$ $\mathbf{W}_u$

$\mathbf{u}^{[1]}$ $\qquad\qquad$ $\mathbf{u}^{[2]}$ $\qquad\qquad$ $\mathbf{u}^{[3]}$ $\qquad\qquad$ $\mathbf{u}^{[4]}$

$n$ $\qquad\qquad n \qquad\qquad n \qquad\qquad n$

*text sequence / corpus*

another $\qquad\qquad$ brick $\qquad\qquad$ in $\qquad\qquad$ the $\qquad\qquad$ Wall $\qquad\qquad$ ...

$x_1$ $\qquad\qquad\quad$ $x_2$ $\qquad\qquad\quad$ $x_3$ $\qquad\qquad\quad$ $x_4$

▸ The **number of tokens**, $T$, across a large corpus is obviously quite large!

$$L(\theta) = \frac{1}{T} \sum_{t=1}^{T} L^{[t]}(\theta)$$

▸ Computing $L(\theta)$ becomes too **computationally expensive**…

▸ Instead we (*once again*) work with a specified **window of text**, say a sentence

▸ We compute $L(\theta)$ for a batch of sentences, then compute the gradient of the loss with respect to the parameters of the network, and then update the parameters.

▸ We repeat this on a new batch until we eventually pass across the entire corpus.

▸ And then we go back to the beginning and repeat the entire process (*a new training epoch*), if necessary.

▸ The **number of tokens**, $T$, across a large corpus is obviously quite large!

$$L(\theta) = \frac{1}{T} \sum_{t=1}^{T} L^{[t]}(\theta)$$

how?

▸ Computing $L(\theta)$ becomes too **computationally expensive**...

▸ Instead we (*once again*) work with a specified **window of text**, say a sentence

▸ We compute $L(\theta)$ for a batch of sentences, **then compute the gradient of the loss with respect to the parameters of the network**, and then update the parameters.

▸ We repeat this on a new batch until we eventually pass across the entire corpus.

▸ And then we go back to the beginning and repeat the entire process (*a new training epoch*), if necessary.

During training, one of the derivatives we need to estimate is: $\dfrac{\partial L^{[t]}}{\partial \mathbf{W}_h}$

During training, one of the derivatives we need to estimate is: $\dfrac{\partial L^{[t]}}{\partial \mathbf{W}_h}$

This is given by: $\dfrac{\partial L^{[t]}}{\partial \mathbf{W}_h} = \sum_{i=1}^{t} \dfrac{\partial L^{[t]}}{\partial \mathbf{W}_h}\Big|_{(i)}$ *we are summing up the gradients at each time step*

During training, one of the derivatives we need to estimate is: $\dfrac{\partial L^{[t]}}{\partial \mathbf{W}_h}$

This is given by: $\dfrac{\partial L^{[t]}}{\partial \mathbf{W}_h} = \sum_{i=1}^{t} \dfrac{\partial L^{[t]}}{\partial \mathbf{W}_h}\bigg|_{(i)}$    *we are summing up the gradients at each time step*

Total derivative of a multivariable function $f\big(x(t), y(t)\big)$ that depends on two single variable functions $x(t)$ and $y(t)$

$$\frac{d}{dt} f\big(x(t), y(t)\big) = \frac{\partial f}{\partial x} \cdot \frac{dx}{dt} + \frac{\partial f}{\partial y} \cdot \frac{dy}{dt}$$

Total derivative of a multivariable function $f\big(x(t), y(t)\big)$ that depends on two single variable functions $x(t)$ and $y(t)$

$$\frac{d}{dt} f\big(x(t), y(t)\big) = \frac{\partial f}{\partial x} \cdot \frac{dx}{dt} + \frac{\partial f}{\partial y} \cdot \frac{dy}{dt}$$

Example

$$f(x, y) = 3x + y^2$$

$$x(t) = t^2$$

$$y(t) = t - 1$$

Total derivative of a multivariable function $f\big(x(t), y(t)\big)$ that depends on two single variable functions $x(t)$ and $y(t)$

$$\frac{d}{dt} f\big(x(t), y(t)\big) = \frac{\partial f}{\partial x} \cdot \frac{dx}{dt} + \frac{\partial f}{\partial y} \cdot \frac{dy}{dt}$$

Example

trivial solution
(*not always possible*)

$$f(x, y) = 3x + y^2$$

$$x(t) = t^2$$

$$y(t) = t - 1$$

$$f(x, y) = 3x(t) + y(t)^2$$

Total derivative of a multivariable function $f\big(x(t), y(t)\big)$ that depends on two single variable functions $x(t)$ and $y(t)$

$$\frac{d}{dt} f\big(x(t), y(t)\big) = \frac{\partial f}{\partial x} \cdot \frac{dx}{dt} + \frac{\partial f}{\partial y} \cdot \frac{dy}{dt}$$

Example

$$f(x, y) = 3x + y^2$$

$$x(t) = t^2$$

$$y(t) = t - 1$$

trivial solution
(*not always possible*)

$$f(x, y) = 3x(t) + y(t)^2$$

$$= 3t^2 + (t - 1)^2$$

Total derivative of a multivariable function $f\big(x(t), y(t)\big)$ that depends on two single variable functions $x(t)$ and $y(t)$

$$\frac{d}{dt}f\big(x(t), y(t)\big) = \frac{\partial f}{\partial x} \cdot \frac{dx}{dt} + \frac{\partial f}{\partial y} \cdot \frac{dy}{dt}$$

Example

$$f(x, y) = 3x + y^2$$

$$x(t) = t^2$$

$$y(t) = t - 1$$

trivial solution
(*not always possible*)

$$f(x, y) = 3x(t) + y(t)^2$$
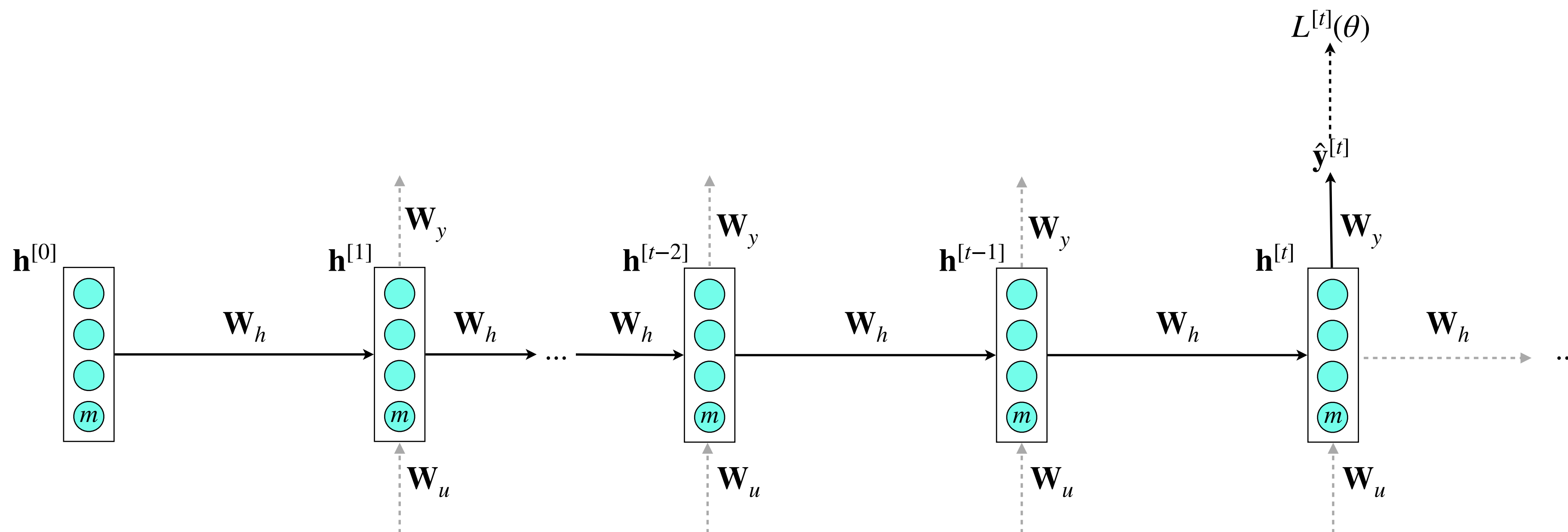
$$= 3t^2 + (t - 1)^2$$

$$= 4t^2 - 2t + 1$$

Total derivative of a multivariable function $f\big(x(t), y(t)\big)$ that depends on two single variable functions $x(t)$ and $y(t)$

$$\frac{d}{dt}f\big(x(t), y(t)\big) = \frac{\partial f}{\partial x} \cdot \frac{dx}{dt} + \frac{\partial f}{\partial y} \cdot \frac{dy}{dt}$$

Example

$$f(x, y) = 3x + y^2$$

$$x(t) = t^2$$

$$y(t) = t - 1$$

trivial solution
(*not always possible*)

$$f(x, y) = 3x(t) + y(t)^2$$

$$= 3t^2 + (t - 1)^2$$

$$= 4t^2 - 2t + 1$$

$$\frac{df}{dt} = 8t - 2$$

Total derivative of a multivariable function $f\big(x(t), y(t)\big)$ that depends on two single variable functions $x(t)$ and $y(t)$

$$\frac{d}{dt}f\big(x(t), y(t)\big) = \frac{\partial f}{\partial x} \cdot \frac{dx}{dt} + \frac{\partial f}{\partial y} \cdot \frac{dy}{dt}$$

**Example**

$$f(x, y) = 3x + y^2$$

$$x(t) = t^2$$

$$y(t) = t - 1$$

**trivial solution**
(*not always possible*)

$$f(x, y) = 3x(t) + y(t)^2$$

$$= 3t^2 + (t - 1)^2$$

$$= 4t^2 - 2t + 1$$

$$\frac{df}{dt} = 8t - 2$$

**multivariate chain rule**

Total derivative of a multivariable function $f\big(x(t), y(t)\big)$ that depends on two single variable functions $x(t)$ and $y(t)$

$$\frac{d}{dt}f\big(x(t), y(t)\big) = \frac{\partial f}{\partial x}\cdot\frac{dx}{dt} + \frac{\partial f}{\partial y}\cdot\frac{dy}{dt}$$

### Example

$$f(x, y) = 3x + y^2$$

$$x(t) = t^2$$

$$y(t) = t - 1$$

### trivial solution
### (*not always possible*)

$$f(x, y) = 3x(t) + y(t)^2$$

$$= 3t^2 + (t - 1)^2$$

$$= 4t^2 - 2t + 1$$

$$\frac{df}{dt} = 8t - 2$$

### multivariate chain rule

$$\frac{df}{dt} =$$

Total derivative of a multivariable function $f\big(x(t), y(t)\big)$ that depends on two single variable functions $x(t)$ and $y(t)$

$$\frac{d}{dt} f\big(x(t), y(t)\big) = \frac{\partial f}{\partial x} \cdot \frac{dx}{dt} + \frac{\partial f}{\partial y} \cdot \frac{dy}{dt}$$

Example

$$f(x, y) = 3x + y^2$$

$$x(t) = t^2$$

$$y(t) = t - 1$$

trivial solution
(*not always possible*)

$$f(x, y) = 3x(t) + y(t)^2$$

$$= 3t^2 + (t - 1)^2$$

$$= 4t^2 - 2t + 1$$

$$\frac{df}{dt} = 8t - 2$$

multivariate chain rule

$$\frac{df}{dt} = 3$$

Total derivative of a multivariable function $f\big(x(t), y(t)\big)$ that depends on two single variable functions $x(t)$ and $y(t)$

$$\frac{d}{dt}f\big(x(t), y(t)\big) = \frac{\partial f}{\partial x} \cdot \frac{dx}{dt} + \frac{\partial f}{\partial y} \cdot \frac{dy}{dt}$$

Example

$$f(x, y) = 3x + y^2$$

$$x(t) = t^2$$

$$y(t) = t - 1$$

trivial solution
(*not always possible*)

$$f(x, y) = 3x(t) + y(t)^2$$

$$= 3t^2 + (t - 1)^2$$

$$= 4t^2 - 2t + 1$$

$$\frac{df}{dt} = 8t - 2$$

multivariate
chain rule

$$\frac{df}{dt} = 3 \cdot 2t \; +$$

Total derivative of a multivariable function
$f\big(x(t), y(t)\big)$ that depends on two single
variable functions $x(t)$ and $y(t)$

$$\frac{d}{dt} f\big(x(t), y(t)\big) = \frac{\partial f}{\partial x} \cdot \frac{dx}{dt} + \frac{\partial f}{\partial y} \cdot \frac{dy}{dt}$$

Example

$$f(x, y) = 3x + y^2$$

$$x(t) = t^2$$

$$y(t) = t - 1$$

trivial solution
(*not always possible*)

$$f(x, y) = 3x(t) + y(t)^2$$

$$= 3t^2 + (t - 1)^2$$

$$= 4t^2 - 2t + 1$$

$$\frac{df}{dt} = 8t - 2$$

multivariate
chain rule

$$\frac{df}{dt} = 3 \cdot 2t \; + \; 2y$$

Total derivative of a multivariable function $f\big(x(t), y(t)\big)$ that depends on two single variable functions $x(t)$ and $y(t)$

$$\frac{d}{dt} f\big(x(t), y(t)\big) = \frac{\partial f}{\partial x} \cdot \frac{dx}{dt} + \frac{\partial f}{\partial y} \cdot \frac{dy}{dt}$$

Example

$$f(x, y) = 3x + y^2$$

$$x(t) = t^2$$

$$y(t) = t - 1$$

trivial solution
(*not always possible*)

$$f(x, y) = 3x(t) + y(t)^2$$

$$= 3t^2 + (t - 1)^2$$

$$= 4t^2 - 2t + 1$$

$$\frac{df}{dt} = 8t - 2$$

multivariate
chain rule

$$\frac{df}{dt} = 3 \cdot 2t + 2y \cdot 1$$

Total derivative of a multivariable function $f\big(x(t), y(t)\big)$ that depends on two single variable functions $x(t)$ and $y(t)$

$$\frac{d}{dt} f\big(x(t), y(t)\big) = \frac{\partial f}{\partial x} \cdot \frac{dx}{dt} + \frac{\partial f}{\partial y} \cdot \frac{dy}{dt}$$

Example

$$f(x, y) = 3x + y^2$$

$$x(t) = t^2$$

$$y(t) = t - 1$$

trivial solution
(*not always possible*)

$$f(x, y) = 3x(t) + y(t)^2$$

$$= 3t^2 + (t - 1)^2$$

$$= 4t^2 - 2t + 1$$

$$\frac{df}{dt} = 8t - 2$$

multivariate
chain rule

$$\frac{df}{dt} = 3 \cdot 2t + 2y \cdot 1$$

$$= 6t + 2(t - 1)$$

Total derivative of a multivariable function $f\big(x(t), y(t)\big)$ that depends on two single variable functions $x(t)$ and $y(t)$

$$\frac{d}{dt}f\big(x(t), y(t)\big) = \frac{\partial f}{\partial x} \cdot \frac{dx}{dt} + \frac{\partial f}{\partial y} \cdot \frac{dy}{dt}$$

### Example

$$f(x, y) = 3x + y^2$$

$$x(t) = t^2$$

$$y(t) = t - 1$$

### trivial solution
(*not always possible*)

$$f(x, y) = 3x(t) + y(t)^2$$

$$= 3t^2 + (t - 1)^2$$

$$= 4t^2 - 2t + 1$$

$$\frac{df}{dt} = 8t - 2$$

### multivariate chain rule

$$\frac{df}{dt} = 3 \cdot 2t + 2y \cdot 1$$

$$= 6t + 2(t - 1)$$

$$= 8t - 2$$

Total derivative of a multivariable function $f\big(x(t), y(t)\big)$ that depends on two single variable functions $x(t)$ and $y(t)$

$$\frac{d}{dt}f\big(x(t), y(t)\big) = \frac{\partial f}{\partial x} \cdot \frac{dx}{dt} + \frac{\partial f}{\partial y} \cdot \frac{dy}{dt}$$

### Example

$$f(x, y) = 3x + y^2$$

$$x(t) = t^2$$

$$y(t) = t - 1$$

*helpful when a function is unknown!*

### trivial solution
(*not always possible*)

$$f(x, y) = 3x(t) + y(t)^2$$

$$= 3t^2 + (t - 1)^2$$

$$= 4t^2 - 2t + 1$$

$$\frac{df}{dt} = 8t - 2$$

### multivariate chain rule

$$\frac{df}{dt} = 3 \cdot 2t + 2y \cdot 1$$

$$= 6t + 2(t - 1)$$

$$= 8t - 2$$

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

backpropagation over time steps $t, t-1, \ldots, 0$, summing gradients, a.k.a. **backpropagation through time (BPTT)**

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

backpropagation over time steps $t, t-1, \ldots, 0$, summing gradients, a.k.a. **backpropagation through time (BPTT)**

$$L(\theta) = \frac{1}{4} \sum_{t=1}^{4} L^{[t]}(\theta)$$

$$L(\theta) = \frac{1}{4} \sum_{t=1}^{4} L^{[t]}(\theta)$$

$$\frac{\partial L}{\partial \mathbf{W}_h} = \sum_{t=1}^{4} \frac{\partial L^{[t]}}{\partial \mathbf{W}_h}$$

$$L(\theta) = \frac{1}{4} \sum_{t=1}^{4} L^{[t]}(\theta)$$

$$\frac{\partial L}{\partial \mathbf{W}_h} = \sum_{t=1}^{4} \boxed{\frac{\partial L^{[t]}}{\partial \mathbf{W}_h}}$$

$$\frac{\partial L^{[t]}}{\partial \mathbf{W}_h} = \sum_{k=1}^{t} \frac{\partial L^{[t]}}{\partial \hat{\mathbf{y}}^{[t]}} \cdot \frac{\partial \hat{\mathbf{y}}^{[t]}}{\partial \mathbf{h}^{[t]}} \cdot \frac{\partial \mathbf{h}^{[t]}}{\partial \mathbf{h}^{[k]}} \cdot \frac{\partial \mathbf{h}^{[k]}}{\partial \mathbf{W}_h}$$

$$L(\theta) = \frac{1}{4} \sum_{t=1}^{4} L^{[t]}(\theta)$$

$$\frac{\partial L}{\partial \mathbf{W}_h} = \sum_{t=1}^{4} \boxed{\frac{\partial L^{[t]}}{\partial \mathbf{W}_h}}$$

$$\frac{\partial L^{[t]}}{\partial \mathbf{W}_h} = \sum_{k=1}^{t} \frac{\partial L^{[t]}}{\partial \hat{\mathbf{y}}^{[t]}} \cdot \frac{\partia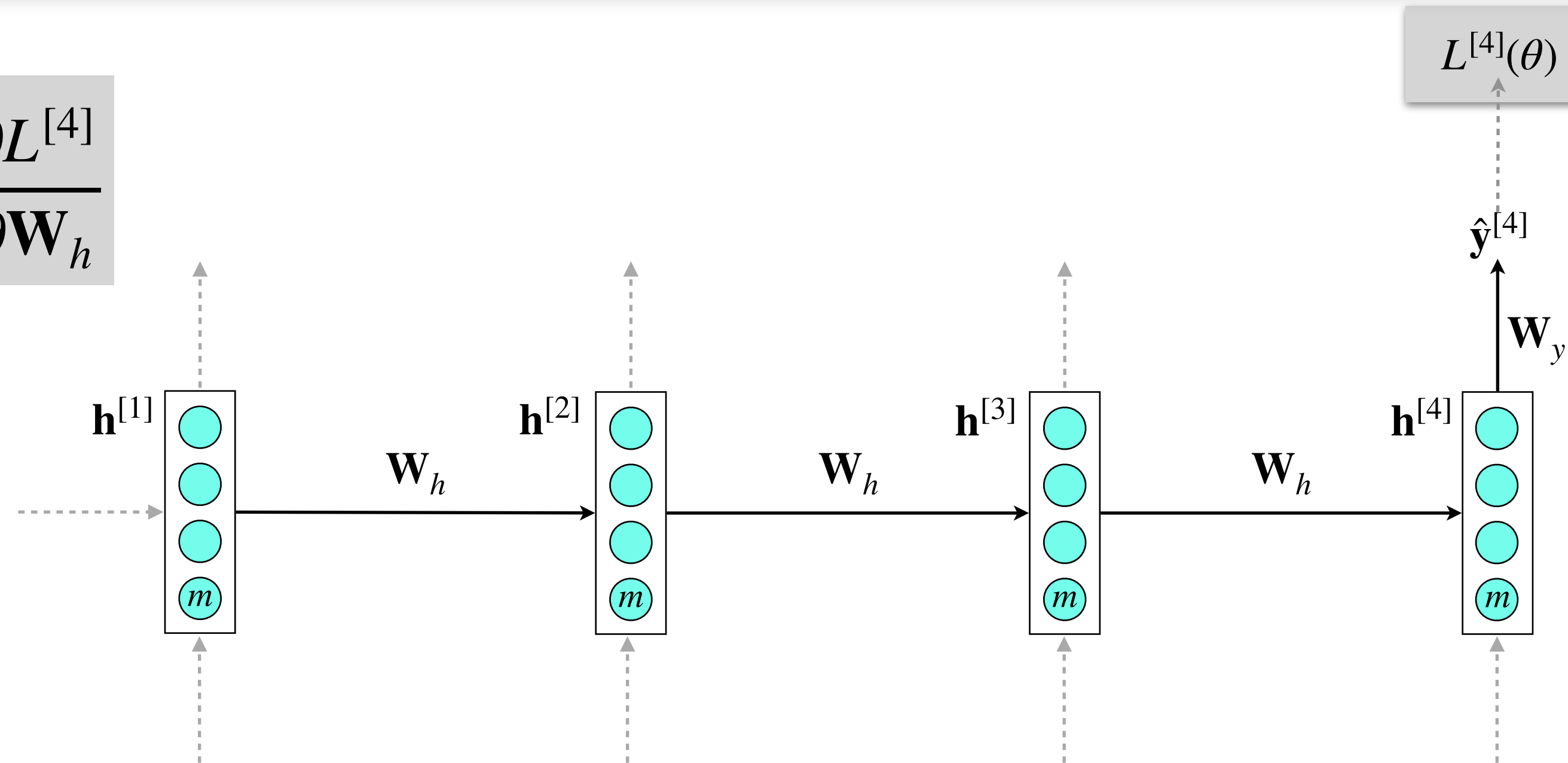l \hat{\mathbf{y}}^{[t]}}{\partial \mathbf{h}^{[t]}} \cdot \boxed{\frac{\partial \mathbf{h}^{[t]}}{\partial \mathbf{h}^{[k]}}} \cdot \frac{\partial \mathbf{h}^{[k]}}{\partial \mathbf{W}_h}$$

$$\frac{\partial \mathbf{h}^{[t]}}{\partial \mathbf{h}^{[k]}} = \prod_{j=k+1}^{t} \frac{\partial \mathbf{h}^{[j]}}{\partial \mathbf{h}^{[j-1]}}$$

$$L(\theta) = \frac{1}{4} \sum_{t=1}^{4} L^{[t]}(\theta)$$

$$\frac{\partial L}{\partial \mathbf{W}_h} = \sum_{t=1}^{4} \frac{\partial L^{[t]}}{\partial \mathbf{W}_h}$$

$$\frac{\partial L^{[t]}}{\partial \mathbf{W}_h} = \sum_{k=1}^{t} \frac{\partial L^{[t]}}{\partial \hat{\mathbf{y}}^{[t]}} \cdot \frac{\partial \hat{\mathbf{y}}^{[t]}}{\partial \mathbf{h}^{[t]}} \cdot \frac{\partial \mathbf{h}^{[t]}}{\partial \mathbf{h}^{[k]}} \cdot \frac{\partial \mathbf{h}^{[k]}}{\partial \mathbf{W}_h}$$

$$\frac{\partial \mathbf{h}^{[t]}}{\partial \mathbf{h}^{[k]}} = \prod_{j=k+1}^{t} \frac{\partial \mathbf{h}^{[j]}}{\partial \mathbf{h}^{[j-1]}}$$

e.g. if $t = 4$ and $k = 1$
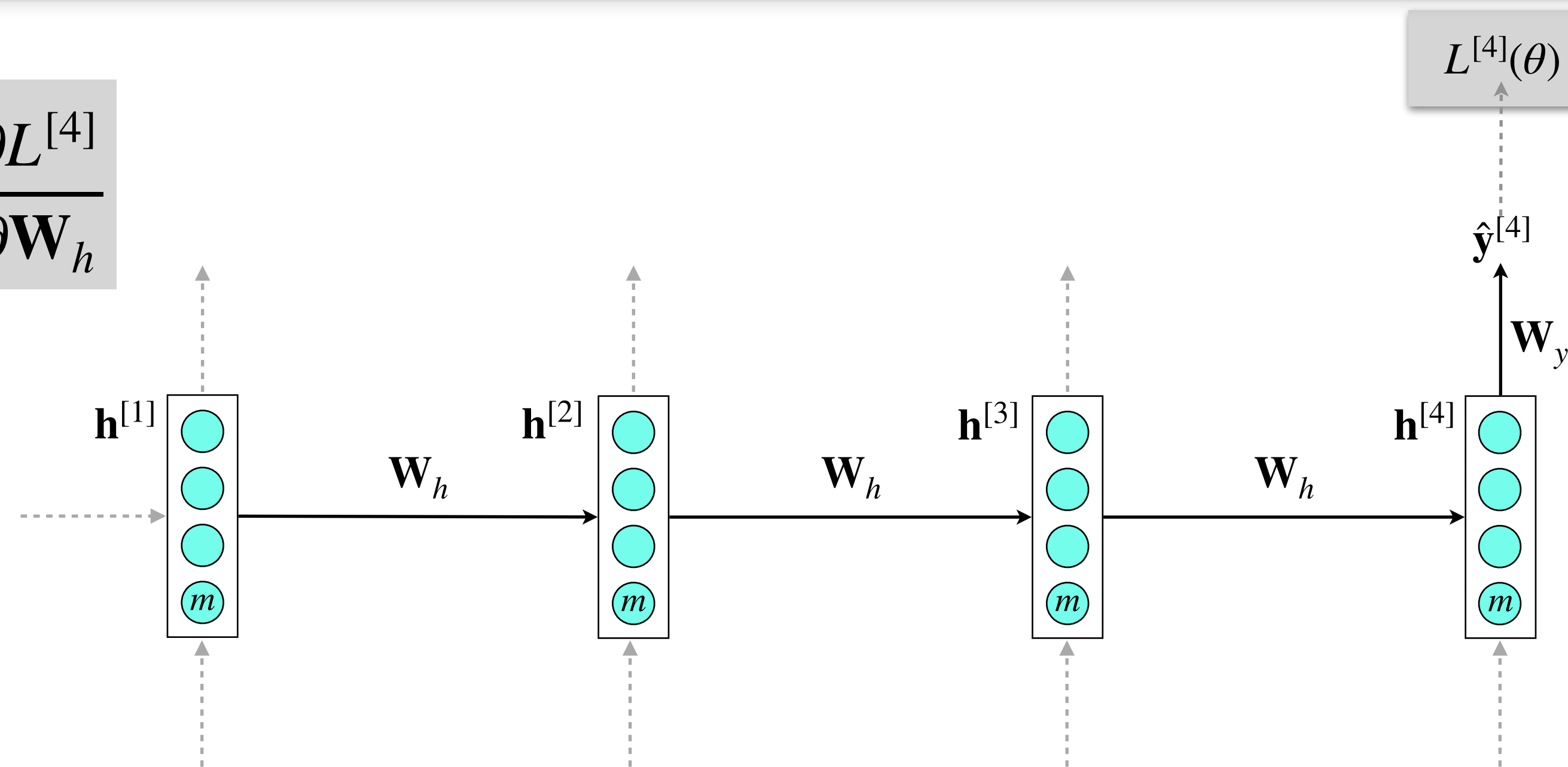
$$\frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[1]}} = \prod_{j=2}^{4} \frac{\partial \mathbf{h}^{[j]}}{\partial \mathbf{h}^{[j-1]}} = \frac{\partial \mathbf{h}^{[2]}}{\partial \mathbf{h}^{[1]}} \cdot \frac{\partial \mathbf{h}^{[3]}}{\partial \mathbf{h}^{[2]}} \cdot \frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[3]}}$$

$$\frac{\partial L}{\partial \mathbf{W}_h} = \sum_{t=1}^{4} \frac{\partial L^{[t]}}{\partial \mathbf{W}_h} = \frac{\partial L^{[1]}}{\partial \mathbf{W}_h} + \frac{\partial L^{[2]}}{\partial \mathbf{W}_h} + \frac{\partial L^{[3]}}{\partial \mathbf{W}_h} + \boxed{\frac{\partial L^{[4]}}{\partial \mathbf{W}_h}}$$

$$\frac{\partial L^{[4]}}{\partial \mathbf{W}_h} = \sum_{k=1}^{4} \frac{\partial L^{[4]}}{\partial \hat{\mathbf{y}}^{[4]}} \cdot \frac{\partial \hat{\mathbf{y}}^{[4]}}{\partial \mathbf{h}^{[4]}} \cdot \boxed{\frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[k]}}} \cdot \frac{\partial \mathbf{h}^{[k]}}{\partial \mathbf{W}_h}$$

let's focus on this component of the sum

$$\propto \frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[1]}} + \frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[2]}} + \frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[3]}} + \frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[4]}}$$

$$\frac{\partial L}{\partial \mathbf{W}_h} = \sum_{t=1}^{4} \frac{\partial L^{[t]}}{\partial \mathbf{W}_h} = \frac{\partial L^{[1]}}{\partial \mathbf{W}_h} + \frac{\partial L^{[2]}}{\partial \mathbf{W}_h} + \frac{\partial L^{[3]}}{\partial \mathbf{W}_h} + \boxed{\frac{\partial L^{[4]}}{\partial \mathbf{W}_h}}$$

$$\frac{\partial L^{[4]}}{\partial \mathbf{W}_h} = \sum_{k=1}^{4} \frac{\partial L^{[4]}}{\partial \hat{\mathbf{y}}^{[4]}} \cdot \frac{\partial \hat{\mathbf{y}}^{[4]}}{\partial \mathbf{h}^{[4]}} \cdot \boxed{\frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[k]}}} \cdot \frac{\partial \mathbf{h}^{[k]}}{\partial \mathbf{W}_h}$$
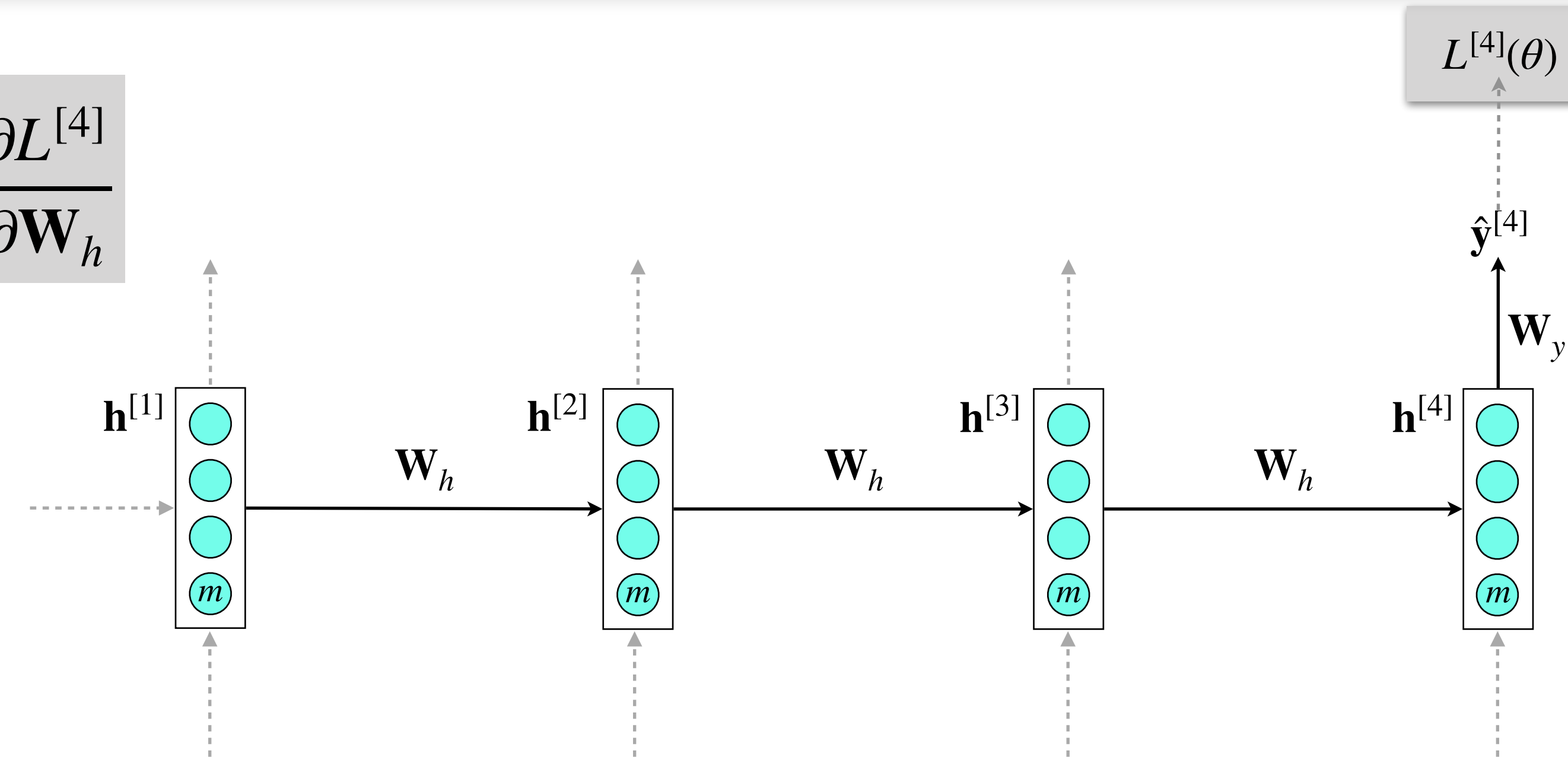
let's focus on this
component of the sum

$$\propto \frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[1]}} + \frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[2]}} + \frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[3]}} + \frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[4]}}$$

$$\propto \sum_{k=1}^{4} \prod_{j=k+1}^{4} \frac{\partial \mathbf{h}^{[j]}}{\partial \mathbf{h}^{[j-1]}}$$

recall
$$\frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[k]}} = \prod_{j=k+1}^{4} \frac{\partial \mathbf{h}^{[j]}}{\partial \mathbf{h}^{[j-1]}}$$

$L^{[4]}(\theta)$

$\hat{\mathbf{y}}^{[4]}$

$\mathbf{W}_y$

$\mathbf{h}^{[1]}$ $\xrightarrow{\mathbf{W}_h}$ $\mathbf{h}^{[2]}$ $\xrightarrow{\mathbf{W}_h}$ $\mathbf{h}^{[3]}$ $\xrightarrow{\mathbf{W}_h}$ $\mathbf{h}^{[4]}$

$m$ $m$ $m$ $m$

$$\frac{\partial L}{\partial \mathbf{W}_h} = \sum_{t=1}^{4} \frac{\partial L^{[t]}}{\partial \mathbf{W}_h} = \frac{\partial L^{[1]}}{\partial \mathbf{W}_h} + \frac{\partial L^{[2]}}{\partial \mathbf{W}_h} + \frac{\partial L^{[3]}}{\partial \mathbf{W}_h} + \boxed{\frac{\partial L^{[4]}}{\partial \mathbf{W}_h}}$$

$L^{[4]}(\theta)$

$\hat{\mathbf{y}}^{[4]}$

$\mathbf{W}_y$

$\mathbf{h}^{[1]}$  $\mathbf{h}^{[2]}$  $\mathbf{h}^{[3]}$  $\mathbf{h}^{[4]}$

$\mathbf{W}_h$   $\mathbf{W}_h$   $\mathbf{W}_h$

$$\frac{\partial L^{[4]}}{\partial \mathbf{W}_h} = \sum_{k=1}^{4} \frac{\partial L^{[4]}}{\partial \hat{\mathbf{y}}^{[4]}} \cdot \frac{\partial \hat{\mathbf{y}}^{[4]}}{\partial \mathbf{h}^{[4]}} \cdot \boxed{\frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[k]}}} \cdot \frac{\partial \mathbf{h}^{[k]}}{\partial \mathbf{W}_h}$$

let's focus on this
component of the sum

$$\propto \frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[1]}} + \frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[2]}} + \frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[3]}} + \frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[4]}}$$

recall

$$\frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[k]}} = \prod_{j=k+1}^{4} \frac{\partial \mathbf{h}^{[j]}}{\partial \mathbf{h}^{[j-1]}}$$

$$\propto \sum_{k=1}^{4} \prod_{j=k+1}^{4} \frac{\partial \mathbf{h}^{[j]}}{\partial \mathbf{h}^{[j-1]}}$$

$$\implies \frac{\partial L^{[4]}}{\partial \mathbf{W}_h} \propto \frac{\partial \mathbf{h}^{[2]}}{\partial \mathbf{h}^{[1]}} \cdot \frac{\partial \mathbf{h}^{[3]}}{\partial \mathbf{h}^{[2]}} \cdot \frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[3]}} + \frac{\partial \mathbf{h}^{[3]}}{\partial \mathbf{h}^{[2]}} \cdot \frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[3]}} + \frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[3]}}$$

$$\frac{\partial L}{\partial \mathbf{W}_h} = \sum_{t=1}^{4} \frac{\partial L^{[t]}}{\partial \mathbf{W}_h} = \frac{\partial L^{[1]}}{\partial \mathbf{W}_h} + \frac{\partial L^{[2]}}{\partial \mathbf{W}_h} + \frac{\partial L^{[3]}}{\partial \mathbf{W}_h} + \boxed{\frac{\partial L^{[4]}}{\partial \mathbf{W}_h}}$$

$L^{[4]}(\theta)$

$\hat{\mathbf{y}}^{[4]}$

$\mathbf{W}_y$

$\mathbf{h}^{[1]}$  $\mathbf{W}_h$  $\mathbf{h}^{[2]}$  $\mathbf{W}_h$  $\mathbf{h}^{[3]}$  $\mathbf{W}_h$  $\mathbf{h}^{[4]}$

$m$  $m$  $m$  $m$

$$\frac{\partial L^{[4]}}{\partial \mathbf{W}_h} = \sum_{k=1}^{4} \frac{\partial L^{[4]}}{\partial \hat{\mathbf{y}}^{[4]}} \cdot \frac{\partial \hat{\mathbf{y}}^{[4]}}{\partial \mathbf{h}^{[4]}} \cdot \boxed{\frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[k]}}} \cdot \frac{\partial \mathbf{h}^{[k]}}{\partial \mathbf{W}_h}$$

let's focus on this component of the sum

$$\propto \frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[1]}} + \frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[2]}} + \frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[3]}} + \frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[4]}}$$

*what if these are small (or large)?*

$$\propto \sum_{k=1}^{4} \prod_{j=k+1}^{4} \frac{\partial \mathbf{h}^{[j]}}{\partial \mathbf{h}^{[j-1]}}$$

recall

$$\frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[k]}} = \prod_{j=k+1}^{4} \frac{\partial \mathbf{h}^{[j]}}{\partial \mathbf{h}^{[j-1]}}$$

$$\Longrightarrow \quad \frac{\partial L^{[4]}}{\partial \mathbf{W}_h} \propto \frac{\partial \mathbf{h}^{[2]}}{\partial \mathbf{h}^{[1]}} \cdot \frac{\partial \mathbf{h}^{[3]}}{\partial \mathbf{h}^{[2]}} \cdot \frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[3]}} + \frac{\partial \mathbf{h}^{[3]}}{\partial \mathbf{h}^{[2]}} \cdot \frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[3]}} + \frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[3]}}$$

$$\frac{\partial L}{\partial \mathbf{W}_h} = \sum_{t=1}^{4} \frac{\partial L^{[t]}}{\partial \mathbf{W}_h} = \frac{\partial L^{[1]}}{\partial \mathbf{W}_h} + \frac{\partial L^{[2]}}{\partial \mathbf{W}_h} + \frac{\partial L^{[3]}}{\partial \mathbf{W}_h} + \boxed{\frac{\partial L^{[4]}}{\partial \mathbf{W}_h}}$$

$L^{[4]}(\theta)$

$\hat{\mathbf{y}}^{[4]}$

$\mathbf{W}_y$

$$\frac{\partial L^{[4]}}{\partial \mathbf{W}_h} = \sum_{k=1}^{4} \frac{\partial L^{[4]}}{\partial \hat{\mathbf{y}}^{[4]}} \cdot \frac{\partial \hat{\mathbf{y}}^{[4]}}{\partial \mathbf{h}^{[4]}} \cdot \boxed{\frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[k]}}} \cdot \frac{\partial \mathbf{h}^{[k]}}{\partial \mathbf{W}_h}$$

$\mathbf{h}^{[1]}$    $\mathbf{W}_h$    $\mathbf{h}^{[2]}$    $\mathbf{W}_h$    $\mathbf{h}^{[3]}$    $\mathbf{W}_h$    $\mathbf{h}^{[4]}$

$m$     $m$     $m$     $m$

let's focus on this
component of the sum

$$\propto \frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[1]}} + \frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[2]}} + \frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[3]}} + \frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[4]}}$$

*what if these are
small (or large)?*

*vanishing (or exploding)
gradient as we
backpropagate!*

$$\propto \sum_{k=1}^{4} \prod_{j=k+1}^{4} \frac{\partial \mathbf{h}^{[j]}}{\partial \mathbf{h}^{[j-1]}}$$

recall
$$\frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[k]}} = \prod_{j=k+1}^{4} \frac{\partial \mathbf{h}^{[j]}}{\partial \mathbf{h}^{[j-1]}}$$

$$\implies \frac{\partial L^{[4]}}{\partial \mathbf{W}_h} \propto \frac{\partial \mathbf{h}^{[2]}}{\partial \mathbf{h}^{[1]}} \cdot \frac{\partial \mathbf{h}^{[3]}}{\partial \mathbf{h}^{[2]}} \cdot \frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[3]}} + \frac{\partial \mathbf{h}^{[3]}}{\partial \mathbf{h}^{[2]}} \cdot \frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[3]}} + \frac{\partial \mathbf{h}^{[4]}}{\partial \mathbf{h}^{[3]}}$$

$$\mathbf{h}^{[t]} = \sigma\big(\mathbf{W}_u \cdot \mathbf{u}^{[t]} + \mathbf{W}_h \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_h\big)$$

$$\mathbf{h}^{[t]} = \sigma\left(\mathbf{W}_u \cdot \mathbf{u}^{[t]} + \mathbf{W}_h \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_h\right)$$
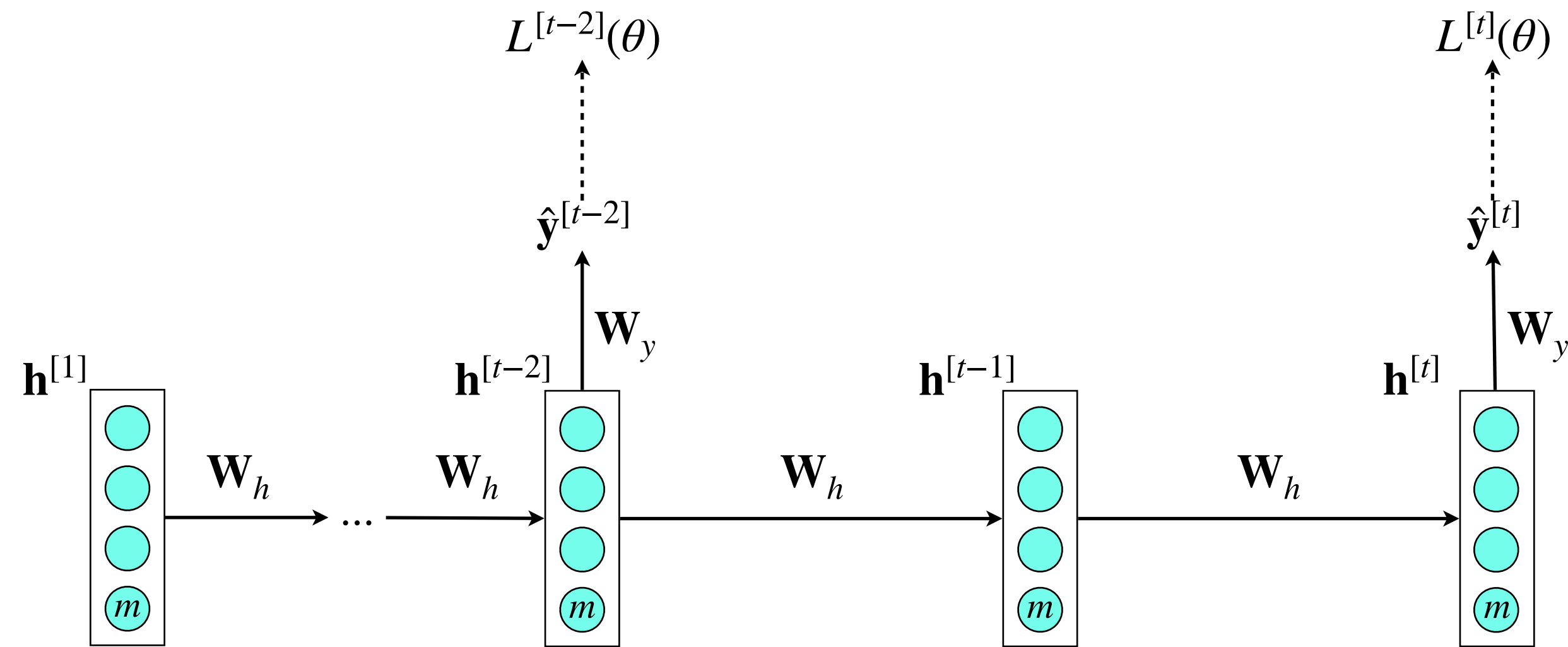
let's ignore the activation function $\sigma$

$$\frac{\partial \mathbf{h}^{[t]}}{\partial \mathbf{h}^{[t-1]}} =$$

$$\mathbf{h}^{[t]} = \sigma\big(\mathbf{W}_u \cdot \mathbf{u}^{[t]} + \mathbf{W}_h \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_h\big)$$

let's ignore the activation function $\sigma$

$$\frac{\partial \mathbf{h}^{[t]}}{\partial \mathbf{h}^{[t-1]}} = \mathbf{W}_h$$

$$\mathbf{h}^{[t]} = \sigma\big(\mathbf{W}_u \cdot \mathbf{u}^{[t]} + \mathbf{W}_h \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_h\big)$$

$$\frac{\partial \mathbf{h}^{[t]}}{\partial \mathbf{h}^{[t-1]}} = \mathbf{W}_h$$

$$\frac{\partial \mathbf{h}^{[t]}}{\partial \mathbf{h}^{[t-\xi]}} = \frac{\partial \mathbf{h}^{[t]}}{\partial \mathbf{h}^{[t-1]}} \cdot \ldots \cdot \frac{\partial \mathbf{h}^{[t-\xi-1]}}{\partial \mathbf{h}^{[t-\xi]}}$$

let's ignore the activation function $\sigma$

let's now see what happens when we compute the partial derivative of hidden state $\mathbf{h}^{[t]}$ w.r.t. the hidden state $\xi$ time steps before it, i.e. $\mathbf{h}^{[t-\xi]}$

$$\mathbf{h}^{[t]} = \sigma\left(\mathbf{W}_u \cdot \mathbf{u}^{[t]} + \mathbf{W}_h \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_h\right)$$

let's ignore the activation function $\sigma$

$$\frac{\partial \mathbf{h}^{[t]}}{\partial \mathbf{h}^{[t-1]}} = \mathbf{W}_h$$

let's now see what happens when we compute the partial derivative of hidden state $\mathbf{h}^{[t]}$ w.r.t. the hidden state $\xi$ time steps before it, i.e. $\mathbf{h}^{[t-\xi]}$

$$\frac{\partial \mathbf{h}^{[t]}}{\partial \mathbf{h}^{[t-\xi]}} = \underbrace{\frac{\partial \mathbf{h}^{[t]}}{\partial \mathbf{h}^{[t-1]}} \cdot \dots \cdot \frac{\partial \mathbf{h}^{[t-\xi-1]}}{\partial \mathbf{h}^{[t-\xi]}}}_{\xi \text{ components}} = \mathbf{W}_h^{\xi}$$

$$\mathbf{h}^{[t]} = \sigma\big(\mathbf{W}_u \cdot \mathbf{u}^{[t]} + \mathbf{W}_h \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_h\big)$$

let's ignore the activation function $\sigma$

$$\frac{\partial \mathbf{h}^{[t]}}{\partial \mathbf{h}^{[t-1]}} = \mathbf{W}_h$$

let's now see what happens when we compute
the partial derivative of hidden state $\mathbf{h}^{[t]}$ w.r.t.
the hidden state $\xi$ time steps before it, i.e. $\mathbf{h}^{[t-\xi]}$

$$\frac{\partial \mathbf{h}^{[t]}}{\partial \mathbf{h}^{[t-\xi]}} = \underbrace{\frac{\partial \mathbf{h}^{[t]}}{\partial \mathbf{h}^{[t-1]}} \cdot \dots \cdot \frac{\partial \mathbf{h}^{[t-\xi-1]}}{\partial \mathbf{h}^{[t-\xi]}}}_{\xi \text{ components}} = \mathbf{W}_h^{\xi}$$

▸ If $\mathbf{W}_h$ has eigenvalues $< 1$, gradients become exponentially smaller
  as time steps $\xi$ increase $\implies$ gradients will become $0$, i.e. vanish

$$\mathbf{h}^{[t]} = \sigma\big(\mathbf{W}_u \cdot \mathbf{u}^{[t]} + \mathbf{W}_h \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_h\big)$$

let's ignore the activation function $\sigma$

$$\frac{\partial \mathbf{h}^{[t]}}{\partial \mathbf{h}^{[t-1]}} = \mathbf{W}_h$$

let's now see what happens when we compute
the partial derivative of hidden state $\mathbf{h}^{[t]}$ w.r.t.
the hidden state $\xi$ time steps before it, i.e. $\mathbf{h}^{[t-\xi]}$

$$\frac{\partial \mathbf{h}^{[t]}}{\partial \mathbf{h}^{[t-\xi]}} = \underbrace{\frac{\partial \mathbf{h}^{[t]}}{\partial \mathbf{h}^{[t-1]}} \cdot \ldots \cdot \frac{\partial \mathbf{h}^{[t-\xi-1]}}{\partial \mathbf{h}^{[t-\xi]}}}_{\xi \text{ components}} = \mathbf{W}_h^{\xi}$$

- ▸ If $\mathbf{W}_h$ has eigenvalues $< 1$, gradients become exponentially smaller
  as time steps $\xi$ increase $\implies$ gradients will become $0$, i.e. vanish

- ▸ If $\mathbf{W}_h$ has eigenvalues $> 1$ $\implies$ gradients will explode

$$\mathbf{h}^{[t]} = \sigma\big(\mathbf{W}_u \cdot \mathbf{u}^{[t]} + \mathbf{W}_h \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_h\big)$$

let's ignore the activation function $\sigma$

$$\frac{\partial \mathbf{h}^{[t]}}{\partial \mathbf{h}^{[t-1]}} = \mathbf{W}_h$$

let's now see what happens when we compute the partial derivative of hidden state $\mathbf{h}^{[t]}$ w.r.t. the hidden state $\xi$ time steps before it, i.e. $\mathbf{h}^{[t-\xi]}$

$$\frac{\partial \mathbf{h}^{[t]}}{\partial \mathbf{h}^{[t-\xi]}} = \underbrace{\frac{\partial \mathbf{h}^{[t]}}{\partial \mathbf{h}^{[t-1]}} \cdot \ldots \cdot \frac{\partial \mathbf{h}^{[t-\xi-1]}}{\partial \mathbf{h}^{[t-\xi]}}}_{\xi \text{ components}} = \mathbf{W}_h^{\xi}$$

- ▸ If $\mathbf{W}_h$ has eigenvalues $< 1$, gradients become exponentially smaller as time steps $\xi$ increase $\implies$ gradients will become $0$, i.e. vanish

- ▸ If $\mathbf{W}_h$ has eigenvalues $> 1$ $\implies$ gradients will explode

- ▸ Similar outcome when we re-introduce an activation function

▸ Signal (gradient) from early states that are distant to the current state is lost $\Longrightarrow$ long-terms effects are not captured

▸ Signal (gradient) from early states that are distant to the current state is lost $\Longrightarrow$ long-terms effects are not captured

▸ NB: Parameters will still be updated, but based on shorter-term gradients that have not vanished.

▸ Large gradients, $\dfrac{\partial L}{\partial \theta_j}$, mean large learning steps

during optimisation

$$\theta_{j+1} = \theta_j - \eta \frac{\partial L}{\partial \theta_j}$$

▸ Large gradients, $\dfrac{\partial L}{\partial \theta_j}$, mean large learning steps during optimisation

$$\theta_{j+1} = \theta_j - \eta \frac{\partial L}{\partial \theta_j}$$

▸ This would possibly result in a poor parameter setting from which we might not be able to recover, especially while using large learning steps

▸ Large gradients, $\dfrac{\partial L}{\partial \theta_j}$, mean large learning steps during optimisation

$$\theta_{j+1} = \theta_j - \eta \frac{\partial L}{\partial \theta_j}$$

▸ This would possibly result in a poor parameter setting from which we might not be able to recover, especially while using large learning steps

▸ The worst penalty to pay would be NaN / Inf errors in the NN parameters; training will have to be restarted



*COMP0087 - Recurrent Neural Networks*

▸ If the L2 norm of the gradient is greater than a threshold $\gamma$, simply scale the gradient down, i.e. clip it!

$$\mathbf{q} = \frac{\partial L}{\partial \theta}$$

**if** $\|\mathbf{q}\| \geq \gamma$ **then**

$$\mathbf{q} = \frac{\gamma}{\|\mathbf{q}\|} \cdot \mathbf{q}$$

**endif**

▸ We are still taking a step in the same direction, albeit a smaller one

▸ We need to learn / set the threshold $\gamma$; a good heuristic $0.5$ to $10$ times the average norm of the gradient over a sufficient number of updates

- ‣ Simple RNNs **fail to maintain information over many time steps** as their architecture does not have explicit components to do so

- ‣ Long Short-Term Memory (LSTM) is an update to the RNN architecture with the aim of solving the problem of vanishing gradients

- ‣ The LSTM has a hidden state like the simple RNN, but also a "cell" state, both being $n$-dimensional vectors

- ‣ The cell is designed to store more long-term information and acts like a memory module — the LSTM can read, delete, and write information to the cell

- ‣ 3 new $n$-dimensional vectors control what is read, deleted, and written; however their decisions are "probabilistic" $\in [0,1]$ for each of the $n$ dimensions (not $0$ or $1$) and are learned during optimisation

$\mathbf{h}^{[t]}$



LSTM

$\mathbf{c}^{[t]}$

$\mathbf{u}^{[1]}$  $\mathbf{u}^{[2]}$  $\mathbf{u}^{[3]}$  ...  $\mathbf{u}^{[t]}$

$n$  $n$  $n$  $n$

another  brick  in  the

**New content**: similarly to simple RNN, there is an input sequence $\mathbf{u}^{[1]}, \ldots, \mathbf{u}^{[t]}$ and there is a dependency to the previous hidden state $\mathbf{h}^{[t-1]}$

$$\tilde{\mathbf{c}}^{[t]} = \tanh\left(\mathbf{U}_c \cdot \mathbf{u}^{[t]} + \mathbf{W}_c \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_c\right)$$

$h^{[t]}$

$\begin{array}{|c|}\hline \text{LSTM} \\ \hline\end{array}$ $c^{[t]}$

$u^{[1]}$ $u^{[2]}$ $u^{[3]}$ ... $u^{[t]}$

$n$ $n$ $n$ $n$

another brick in the

**New content:** similarly to simple RNN, there is an input sequence $\mathbf{u}^{[1]}, \ldots, \mathbf{u}^{[t]}$ and there is a dependency to the previous hidden state $\mathbf{h}^{[t-1]}$

$$\tilde{\mathbf{c}}^{[t]} = \tanh\left(\mathbf{U}_c \cdot \mathbf{u}^{[t]} + \mathbf{W}_c \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_c\right)$$

**Forget gate:** what should be forgotten from the previous cell state; $0 \rightarrow 1 \sim$ forget $\rightarrow$ keep.
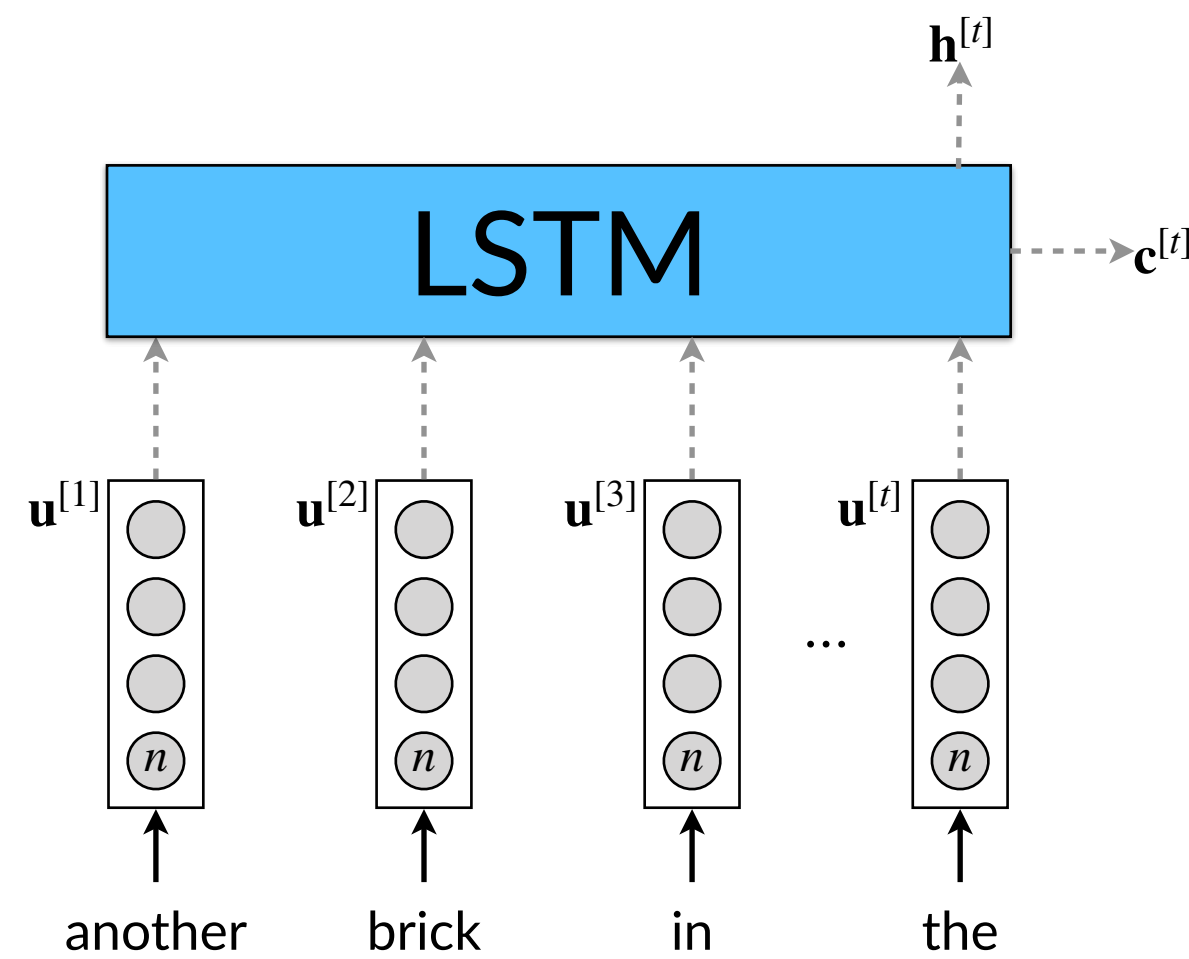
$$\mathbf{f}^{[t]} = \sigma\left(\mathbf{U}_f \cdot \mathbf{u}^{[t]} + \mathbf{W}_f \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_f\right)$$

$$\mathbf{h}^{[t]}$$



LSTM $\dashrightarrow \mathbf{c}^{[t]}$

$\mathbf{u}^{[1]}$  $\mathbf{u}^{[2]}$  $\mathbf{u}^{[3]}$  ...  $\mathbf{u}^{[t]}$

another    brick    in    the

**New content**: similarly to simple RNN, there is an input sequence $\mathbf{u}^{[1]}, \ldots, \mathbf{u}^{[t]}$ and there is a dependency to the previous hidden state $\mathbf{h}^{[t-1]}$

$$\tilde{\mathbf{c}}^{[t]} = \tanh\left(\mathbf{U}_c \cdot \mathbf{u}^{[t]} + \mathbf{W}_c \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_c\right)$$

**Forget gate**: what should be forgotten from the previous cell state; $0 \rightarrow 1 \sim$ forget $\rightarrow$ keep.

$$\mathbf{f}^{[t]} = \sigma\left(\mathbf{U}_f \cdot \mathbf{u}^{[t]} + \mathbf{W}_f \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_f\right)$$

**Input gate**: what should be kept from the new content?
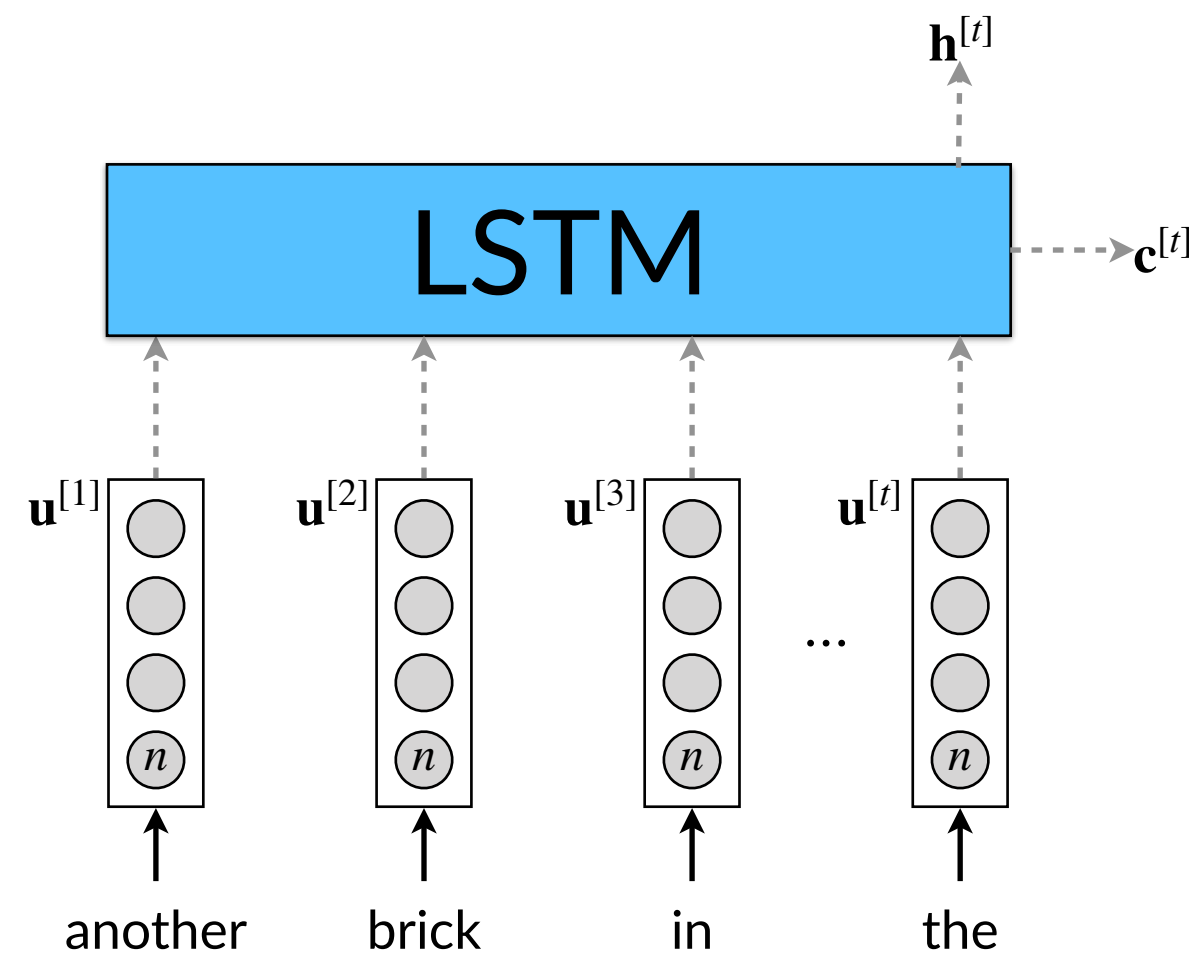
$$\mathbf{i}^{[t]} = \sigma\left(\mathbf{U}_i \cdot \mathbf{u}^{[t]} + \mathbf{W}_i \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_i\right)$$

**New content**: similarly to simple RNN, there is an input sequence $\mathbf{u}^{[1]}, \ldots, \mathbf{u}^{[t]}$ and there is a dependency to the previous hidden state $\mathbf{h}^{[t-1]}$

$$\tilde{\mathbf{c}}^{[t]} = \tanh\left(\mathbf{U}_c \cdot \mathbf{u}^{[t]} + \mathbf{W}_c \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_c\right)$$

**Forget gate**: what should be forgotten from the previous cell state; $0 \to 1 \sim$ forget $\to$ keep.

$$\mathbf{f}^{[t]} = \sigma\left(\mathbf{U}_f \cdot \mathbf{u}^{[t]} + \mathbf{W}_f \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_f\right)$$

**Input gate**: what should be kept from the new content?

$$\mathbf{i}^{[t]} = \sigma\left(\mathbf{U}_i \cdot \mathbf{u}^{[t]} + \mathbf{W}_i \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_i\right)$$

**Cell state**: forget some past, keep some present

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

$\mathbf{h}^{[t]}$

$$\text{LSTM}$$

$\mathbf{c}^{[t]}$

$\mathbf{u}^{[1]}$    $\mathbf{u}^{[2]}$    $\mathbf{u}^{[3]}$    $\mathbf{u}^{[t]}$

$n$    $n$    $n$    $n$

another    brick    in    the

**New content**: similarly to simple RNN, there is an input sequence $\mathbf{u}^{[1]}, \ldots, \mathbf{u}^{[t]}$ and there is a dependency to the previous hidden state $\mathbf{h}^{[t-1]}$

$$\tilde{\mathbf{c}}^{[t]} = \tanh\big(\mathbf{U}_c \cdot \mathbf{u}^{[t]} + \mathbf{W}_c \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_c\big)$$

**Forget gate**: what should be forgotten from the previous cell state; $0 \to 1$ ~ forget $\to$ keep.

$$\mathbf{f}^{[t]} = \sigma\big(\mathbf{U}_f \cdot \mathbf{u}^{[t]} + \mathbf{W}_f \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_f\big)$$
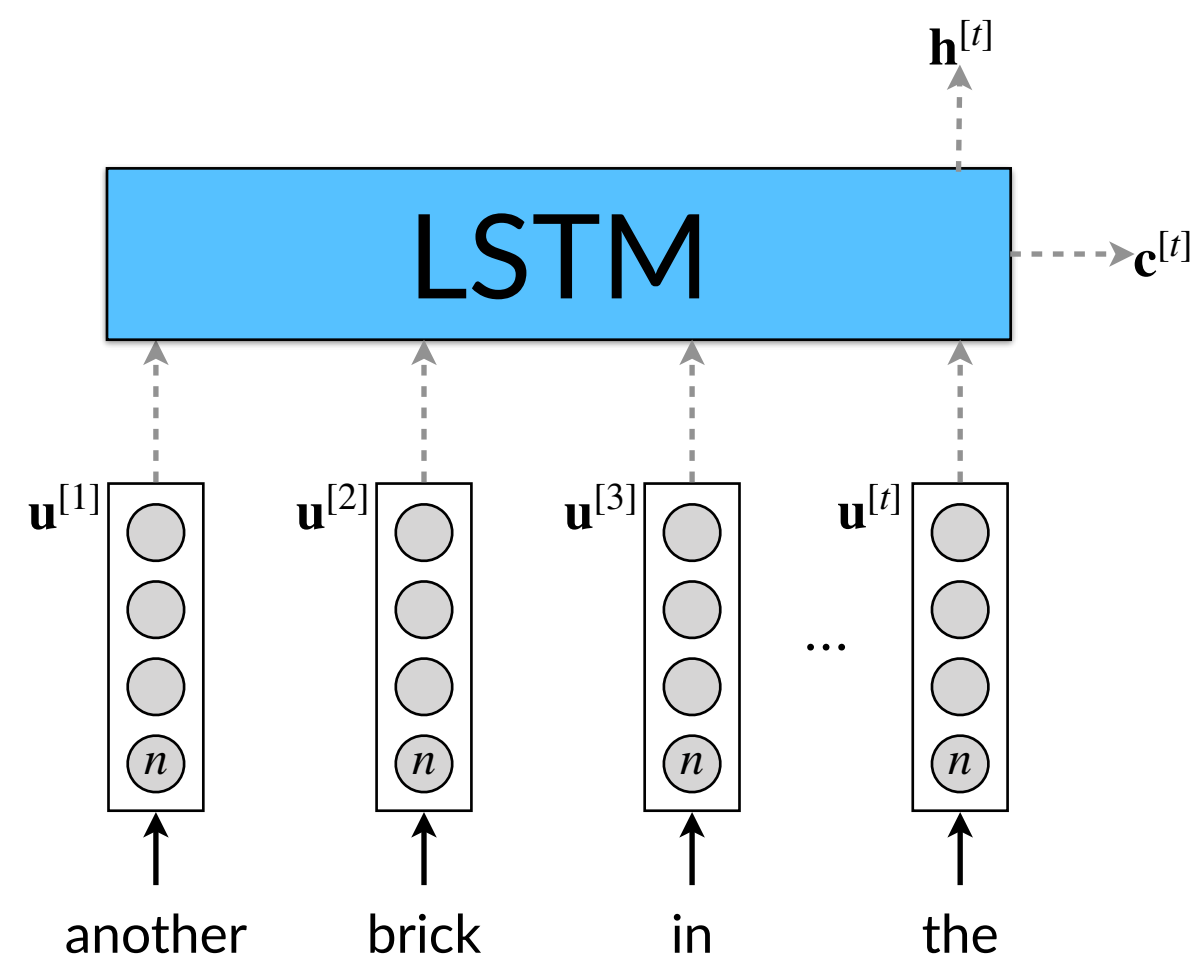
**Input gate**: what should be kept from the new content?

$$\mathbf{i}^{[t]} = \sigma\big(\mathbf{U}_i \cdot \mathbf{u}^{[t]} + \mathbf{W}_i \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_i\big)$$

**Cell state**: forget some past, keep some present

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

**Output gate**: what parts of the cell state will be passed on to the hidden state

$$\mathbf{o}^{[t]} = \sigma\big(\mathbf{U}_o \cdot \mathbf{u}^{[t]} + \mathbf{W}_o \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_o\big)$$

**New content:** similarly to simple RNN, there is an input sequence $\mathbf{u}^{[1]}, \ldots, \mathbf{u}^{[t]}$ and there is a dependency to the previous hidden state $\mathbf{h}^{[t-1]}$

$$\tilde{\mathbf{c}}^{[t]} = \tanh\left(\mathbf{U}_c \cdot \mathbf{u}^{[t]} + \mathbf{W}_c \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_c\right)$$

**Forget gate:** what should be forgotten from the previous cell state; $0 \to 1$ ~ forget $\to$ keep.

$$\mathbf{f}^{[t]} = \sigma\left(\mathbf{U}_f \cdot \mathbf{u}^{[t]} + \mathbf{W}_f \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_f\right)$$

**Input gate:** what should be kept from the new content?

$$\mathbf{i}^{[t]} = \sigma\left(\mathbf{U}_i \cdot \mathbf{u}^{[t]} + \mathbf{W}_i \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_i\right)$$

**Cell state:** forget some past, keep some present

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

**Output gate:** what parts of the cell state will be passed on to the hidden state

$$\mathbf{o}^{[t]} = \sigma\left(\mathbf{U}_o \cdot \mathbf{u}^{[t]} + \mathbf{W}_o \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_o\right)$$

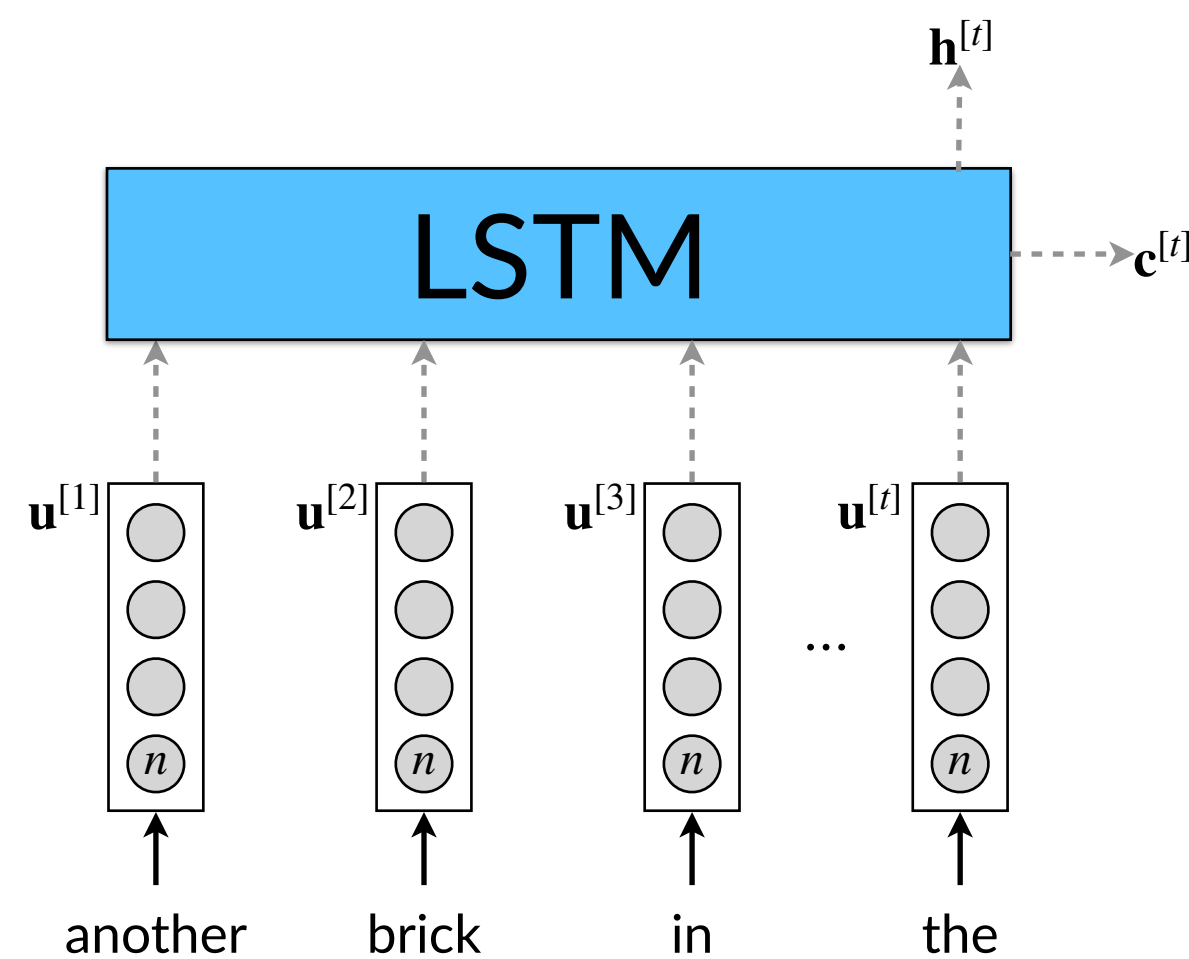**Hidden state:** maintains the part of the cell state that will pass from the output gate

$$\mathbf{h}^{[t]} = \mathbf{o}^{[t]} \odot \tanh\left(\mathbf{c}^{[t]}\right)$$

$$\mathbf{h}^{[t]}$$

LSTM

$$\mathbf{c}^{[t]}$$

$$\mathbf{u}^{[1]} \quad \mathbf{u}^{[2]} \quad \mathbf{u}^{[3]} \qquad \mathbf{u}^{[t]}$$

$$\dots$$

$$n \qquad n \qquad n \qquad n$$

another   brick   in   the

all gate values range from $0$ to $1$
given the sigmoid activation ($\sigma$)
$$\mathbf{f}^{[t]}, \mathbf{i}^{[t]}, \mathbf{o}^{[t]} \in (0,1)^n$$

$$\tilde{\mathbf{c}}^{[t]}, \mathbf{h}^{[t]} \in (-1,1)^n \text{ and } \mathbf{c}^{[t]} \in \mathbb{R}^n$$

$$\odot$$  Hadamard or element-wise product

**New content**: similarly to simple RNN, there is an input sequence $\mathbf{u}^{[1]}, \dots, \mathbf{u}^{[t]}$ and there is a dependency to the previous hidden state $\mathbf{h}^{[t-1]}$

$$\tilde{\mathbf{c}}^{[t]} = \tanh\left(\mathbf{U}_c \cdot \mathbf{u}^{[t]} + \mathbf{W}_c \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_c\right)$$

**Forget gate**: what should be forgotten from the previous cell state; $0 \to 1$ ~ forget $\to$ keep.

$$\mathbf{f}^{[t]} = \sigma\left(\mathbf{U}_f \cdot \mathbf{u}^{[t]} + \mathbf{W}_f \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_f\right)$$

**Input gate**: what should be kept from the new content?

$$\mathbf{i}^{[t]} = \sigma\left(\mathbf{U}_i \cdot \mathbf{u}^{[t]} + \mathbf{W}_i \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_i\right)$$
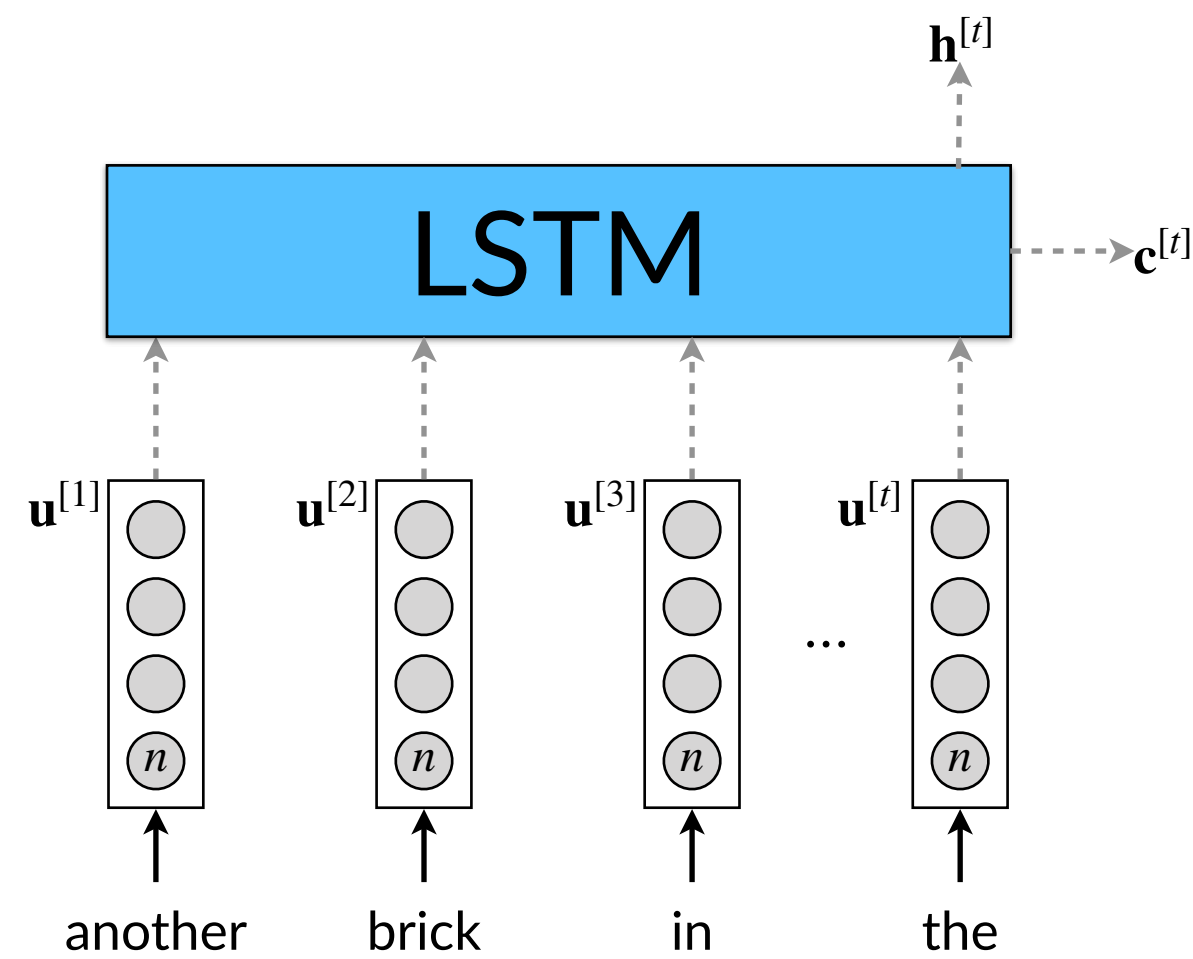
**Cell state**: forget some past, keep some present

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

**Output gate**: what parts of the cell state will be passed on to the hidden state

$$\mathbf{o}^{[t]} = \sigma\left(\mathbf{U}_o \cdot \mathbf{u}^{[t]} + \mathbf{W}_o \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_o\right)$$

**Hidden state**: maintains the part of the cell state that will pass from the output gate

$$\mathbf{h}^{[t]} = \mathbf{o}^{[t]} \odot \tanh\left(\mathbf{c}^{[t]}\right)$$

LSTM

$\mathbf{h}^{[t]}$

$\mathbf{c}^{[t]}$

$\mathbf{u}^{[1]}$  $\mathbf{u}^{[2]}$  $\mathbf{u}^{[3]}$  ...  $\mathbf{u}^{[t]}$

another    brick    in    the

all gate values range from $0$ to $1$
given the sigmoid activation ($\sigma$)
$\mathbf{f}^{[t]}, \mathbf{i}^{[t]}, \mathbf{o}^{[t]} \in (0,1)^n$

$\tilde{\mathbf{c}}^{[t]}, \mathbf{h}^{[t]} \in (-1,1)^n$ and $\mathbf{c}^{[t]} \in \mathbb{R}^n$

$\odot$  Hadamard or element-wise product

independent from each other
could be computed in parallel

**New content:** similarly to simple RNN, there is an input sequence $\mathbf{u}^{[1]}, \ldots, \mathbf{u}^{[t]}$ and there is a dependency to the previous hidden state $\mathbf{h}^{[t-1]}$

$$\tilde{\mathbf{c}}^{[t]} = \tanh\left(\mathbf{U}_c \cdot \mathbf{u}^{[t]} + \mathbf{W}_c \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_c\right)$$

**Forget gate:** what should be forgotten from the previous cell state; $0 \to 1$ ~ forget $\to$ keep.

$$\mathbf{f}^{[t]} = \sigma\left(\mathbf{U}_f \cdot \mathbf{u}^{[t]} + \mathbf{W}_f \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_f\right)$$

**Input gate:** what should be kept from the new content?

$$\mathbf{i}^{[t]} = \sigma\left(\mathbf{U}_i \cdot \mathbf{u}^{[t]} + \mathbf{W}_i \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_i\right)$$

**Cell state:** forget some past, keep some present

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

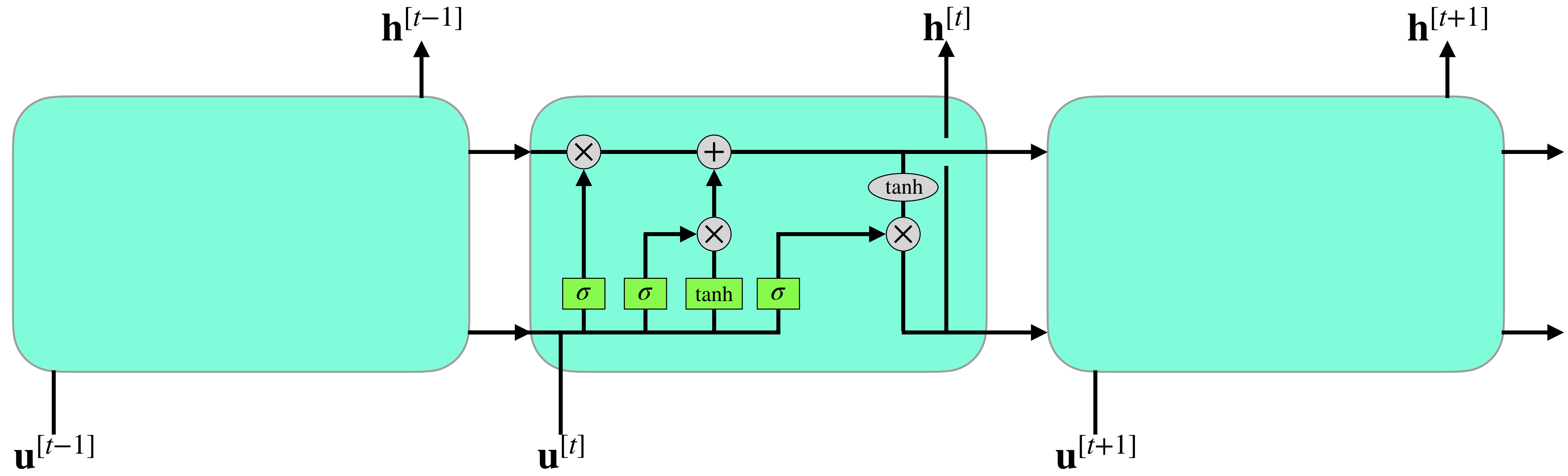**Output gate:** what parts of the cell state will be passed on to the hidden state

$$\mathbf{o}^{[t]} = \sigma\left(\mathbf{U}_o \cdot \mathbf{u}^{[t]} + \mathbf{W}_o \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_o\right)$$

**Hidden state:** maintains the part of the cell state that will pass from the output gate

$$\mathbf{h}^{[t]} = \mathbf{o}^{[t]} \odot \tanh\left(\mathbf{c}^{[t]}\right)$$

## LSTM

$\mathbf{h}^{[t]}$

$\mathbf{c}^{[t]}$

$\mathbf{u}^{[1]}$   $\mathbf{u}^{[2]}$   $\mathbf{u}^{[3]}$   ...   $\mathbf{u}^{[t]}$

another    brick    in    the

all gate values range from $0$ to $1$
given the sigmoid activation ($\sigma$)
$\mathbf{f}^{[t]}, \mathbf{i}^{[t]}, \mathbf{o}^{[t]} \in (0,1)^n$

------

$\tilde{\mathbf{c}}^{[t]}, \mathbf{h}^{[t]} \in (-1,1)^n$ and $\mathbf{c}^{[t]} \in \mathbb{R}^n$

$\odot$   Hadamard or element-wise product

**New content**: similarly to simple RNN, there is an input sequence $\mathbf{u}^{[1]}, \ldots, \mathbf{u}^{[t]}$ and there is a dependency to the previous hidden state $\mathbf{h}^{[t-1]}$

$$\tilde{\mathbf{c}}^{[t]} = \tanh\left(\mathbf{U}_c \cdot \mathbf{u}^{[t]} + \mathbf{W}_c \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_c\right)$$

independent from each other could be computed in parallel

**Forget gate**: what should be forgotten from the previous cell state; $0 \rightarrow 1 \sim$ forget $\rightarrow$ keep.

$$\mathbf{f}^{[t]} = \sigma\left(\mathbf{U}_f \cdot \mathbf{u}^{[t]} + \mathbf{W}_f \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_f\right)$$

**Input gate**: what should be kept from the new content?

$$\mathbf{i}^{[t]} = \sigma\left(\mathbf{U}_i \cdot \mathbf{u}^{[t]} + \mathbf{W}_i \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_i\right)$$

**Cell state**: forget some past, keep some present

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

**Output gate**: what parts of the cell state will be passed on to the hidden state

$$\mathbf{o}^{[t]} = \sigma\left(\mathbf{U}_o \cdot \mathbf{u}^{[t]} + \mathbf{W}_o \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_o\right)$$

**Hidden state**: maintains the part of the cell state that will pass from the output gate

$$\mathbf{h}^{[t]} = \mathbf{o}^{[t]} \odot \tanh\left(\mathbf{c}^{[t]}\right)$$

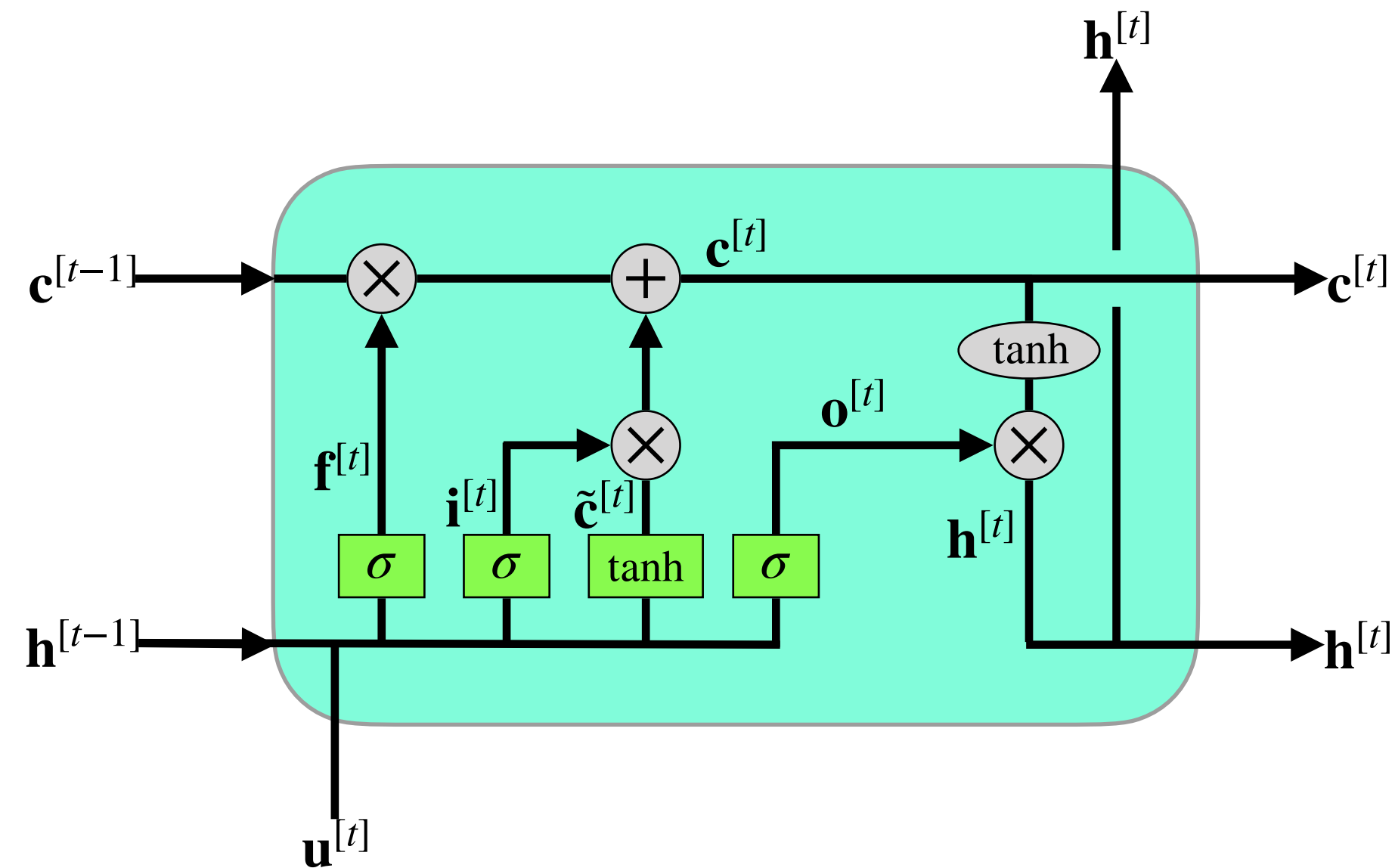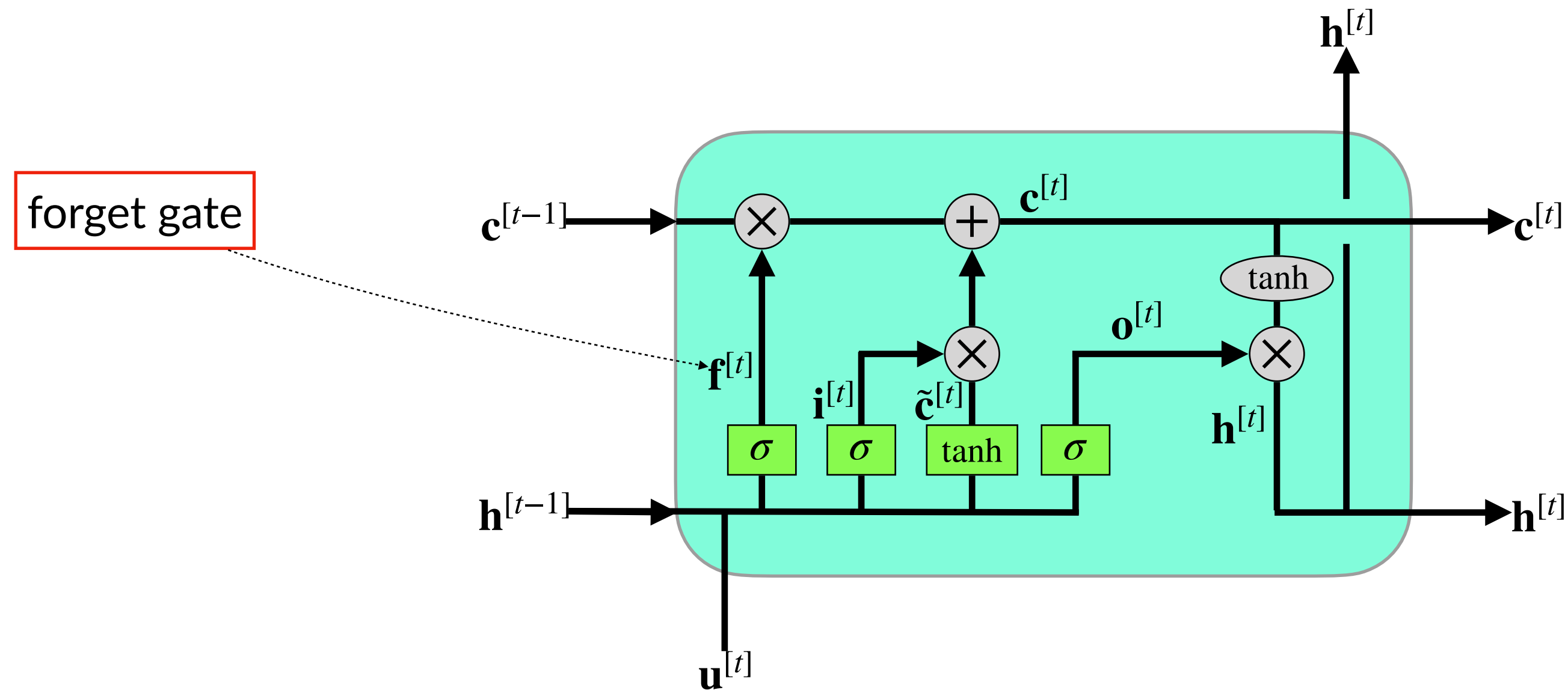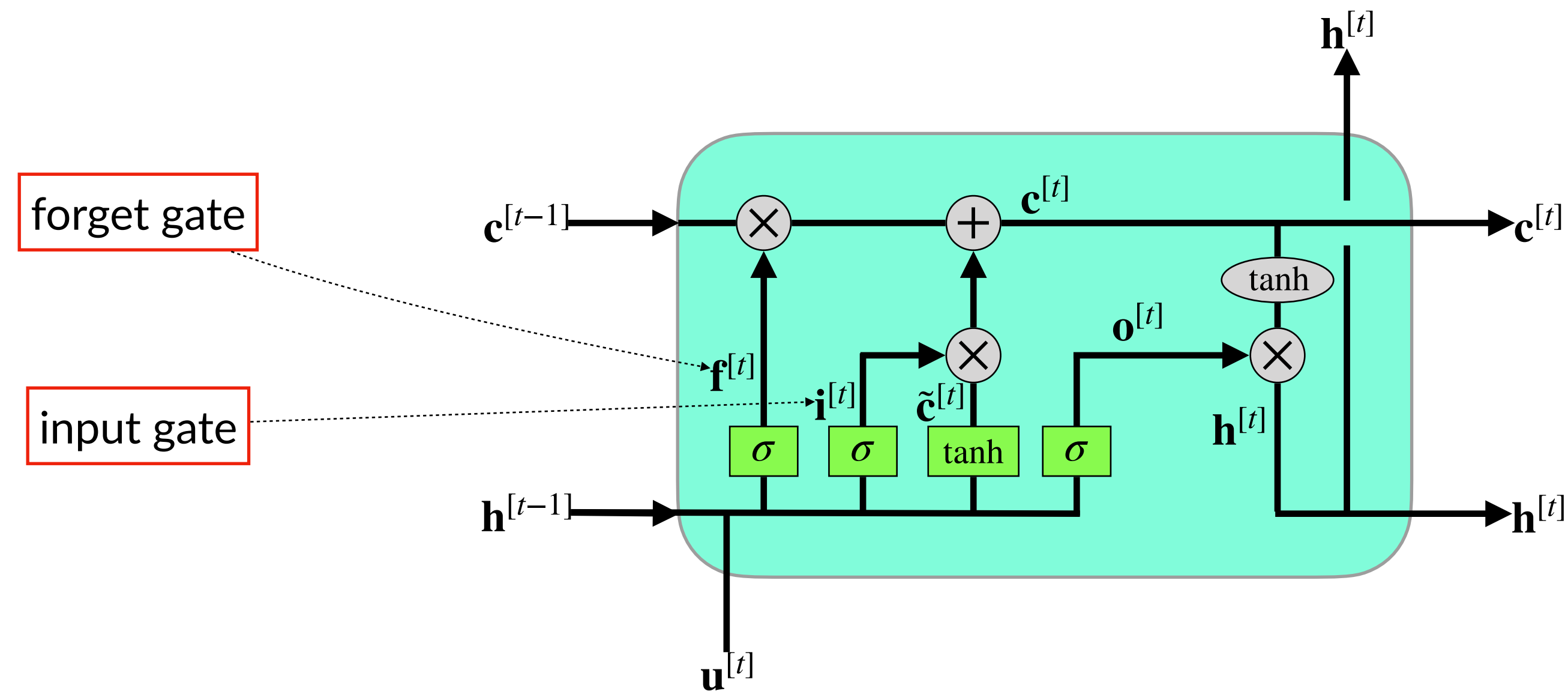If $\mathbf{u}^{[t]} \in \mathbb{R}^m$, how many parameters?

$\mathbf{h}^{[t]}$

## LSTM

$\mathbf{c}^{[t]}$

$\mathbf{u}^{[1]}$ $\mathbf{u}^{[2]}$ $\mathbf{u}^{[3]}$ $\mathbf{u}^{[t]}$

...

another    brick    in    the

all gate values range from $0$ to $1$
given the sigmoid activation ($\sigma$)
$\mathbf{f}^{[t]}, \mathbf{i}^{[t]}, \mathbf{o}^{[t]} \in (0,1)^n$

$\tilde{\mathbf{c}}^{[t]}, \mathbf{h}^{[t]} \in (-1,1)^n$ and $\mathbf{c}^{[t]} \in \mathbb{R}^n$

$\odot$  Hadamard or element-wise product

**New content**: similarly to simple RNN, there is an input sequence $\mathbf{u}^{[1]}, \ldots, \mathbf{u}^{[t]}$ and there is a dependency to the previous hidden state $\mathbf{h}^{[t-1]}$

independent from each other could be computed in parallel

$$\tilde{\mathbf{c}}^{[t]} = \tanh\left(\mathbf{U}_c \cdot \mathbf{u}^{[t]} + \mathbf{W}_c \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_c\right)$$

**Forget gate**: what should be forgotten from the previous cell state; $0 \to 1 \sim$ forget $\to$ keep.

$$\mathbf{f}^{[t]} = \sigma\left(\mathbf{U}_f \cdot \mathbf{u}^{[t]} + \mathbf{W}_f \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_f\right)$$

**Input gate**: what should be kept from the new content?

$$\mathbf{i}^{[t]} = \sigma\left(\mathbf{U}_i \cdot \mathbf{u}^{[t]} + \mathbf{W}_i \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_i\right)$$

**Cell state**: forget some past, keep some present

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

**Output gate**: what parts of the cell state will be passed on to the hidden state

$$\mathbf{o}^{[t]} = \sigma\left(\mathbf{U}_o \cdot \mathbf{u}^{[t]} + \mathbf{W}_o \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_o\right)$$

**Hidden state**: maintains the part of the cell state that will pass from the output gate

$$\mathbf{h}^{[t]} = \mathbf{o}^{[t]} \odot \tanh\left(\mathbf{c}^{[t]}\right)$$

If $\mathbf{u}^{[t]} \in \mathbb{R}^m$, how many parameters?    $= 4 \cdot n \cdot (m + n + 1)$

element-wise operation

$$\tilde{\mathbf{c}}^{[t]} = \tanh\big(\mathbf{U}_c \cdot \mathbf{u}^{[t]} + \mathbf{W}_c \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_c\big)$$

$$\mathbf{f}^{[t]} = \sigma\Big(\mathbf{U}_f \cdot \mathbf{u}^{[t]} + \mathbf{W}_f \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_f\Big)$$
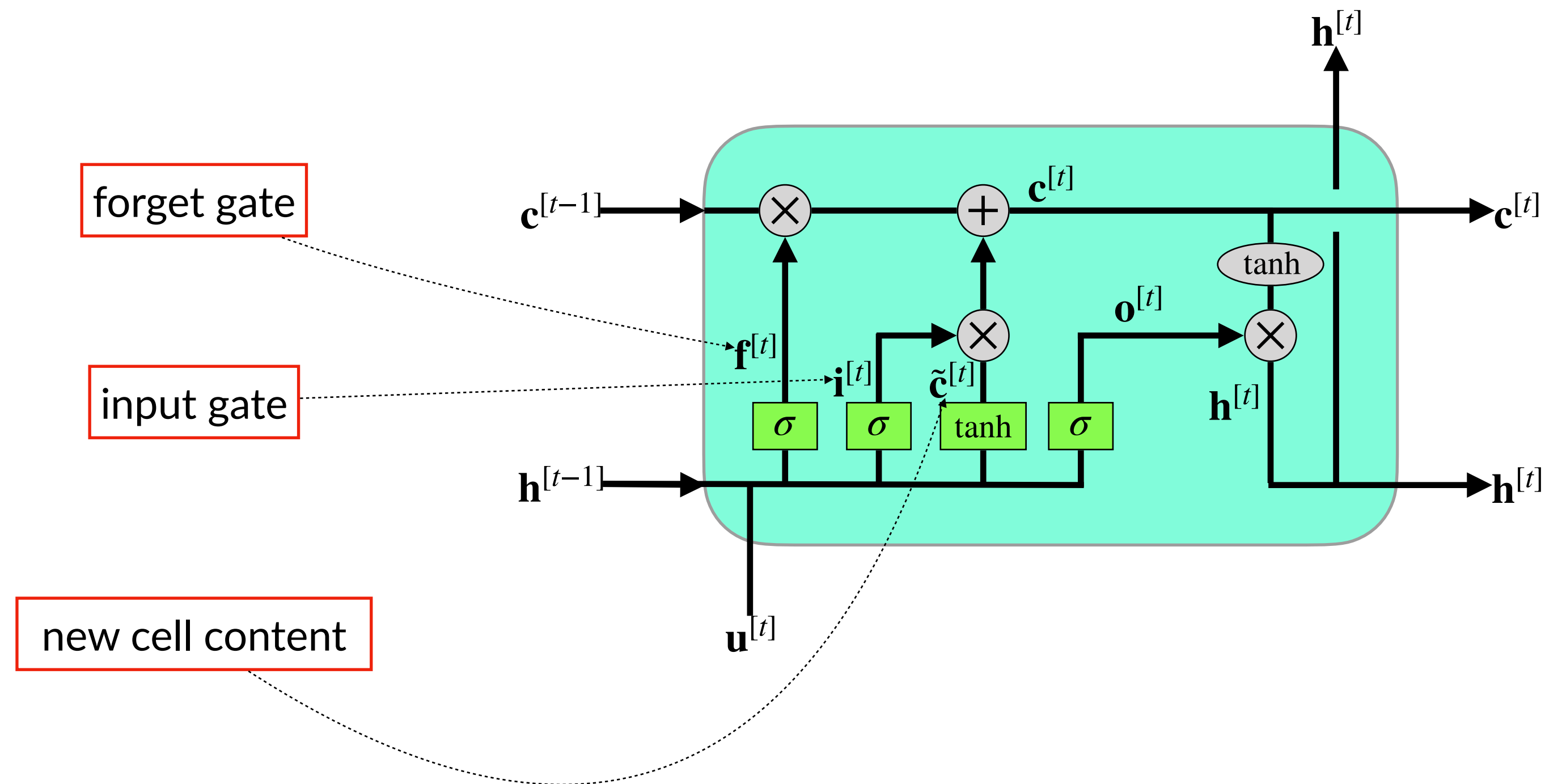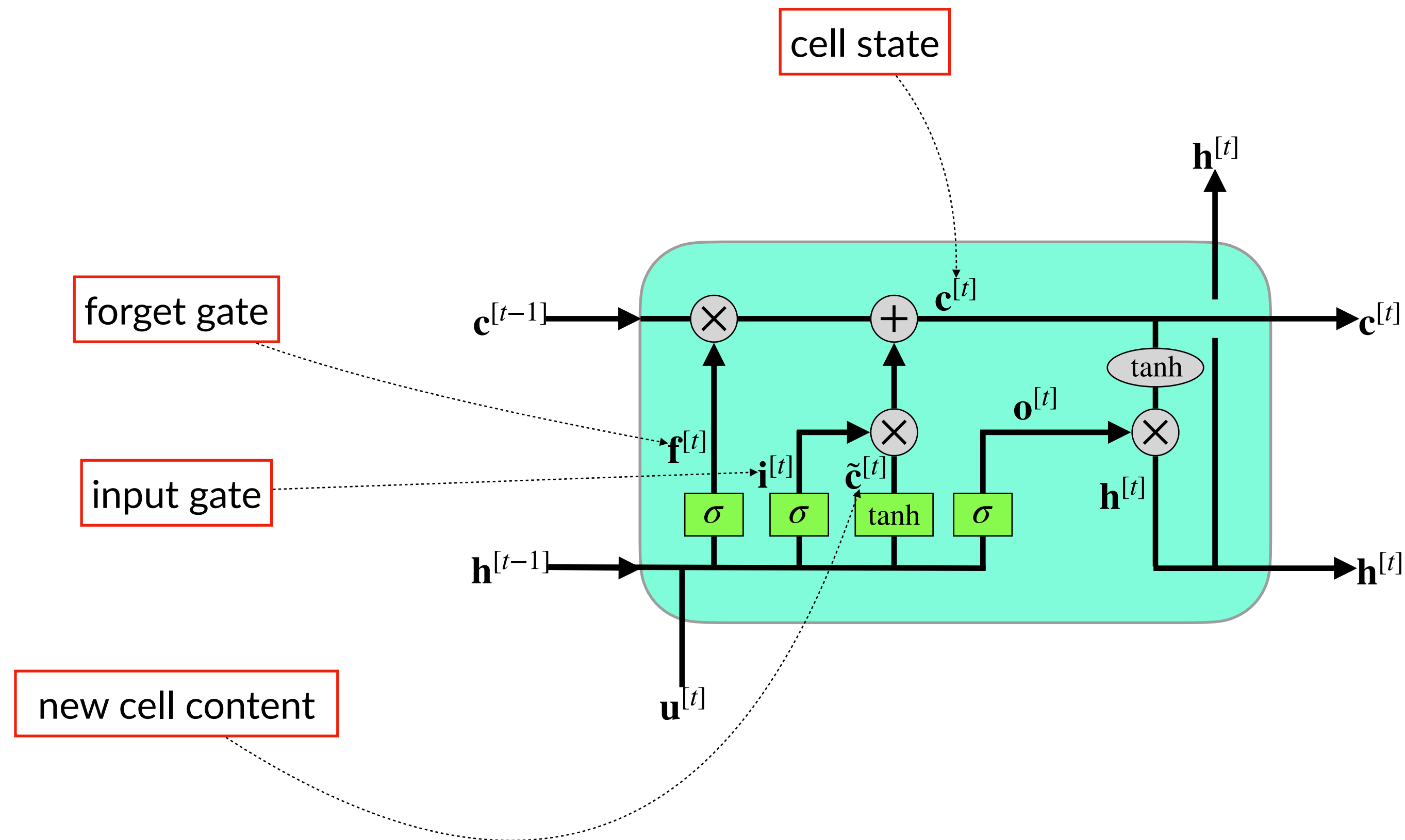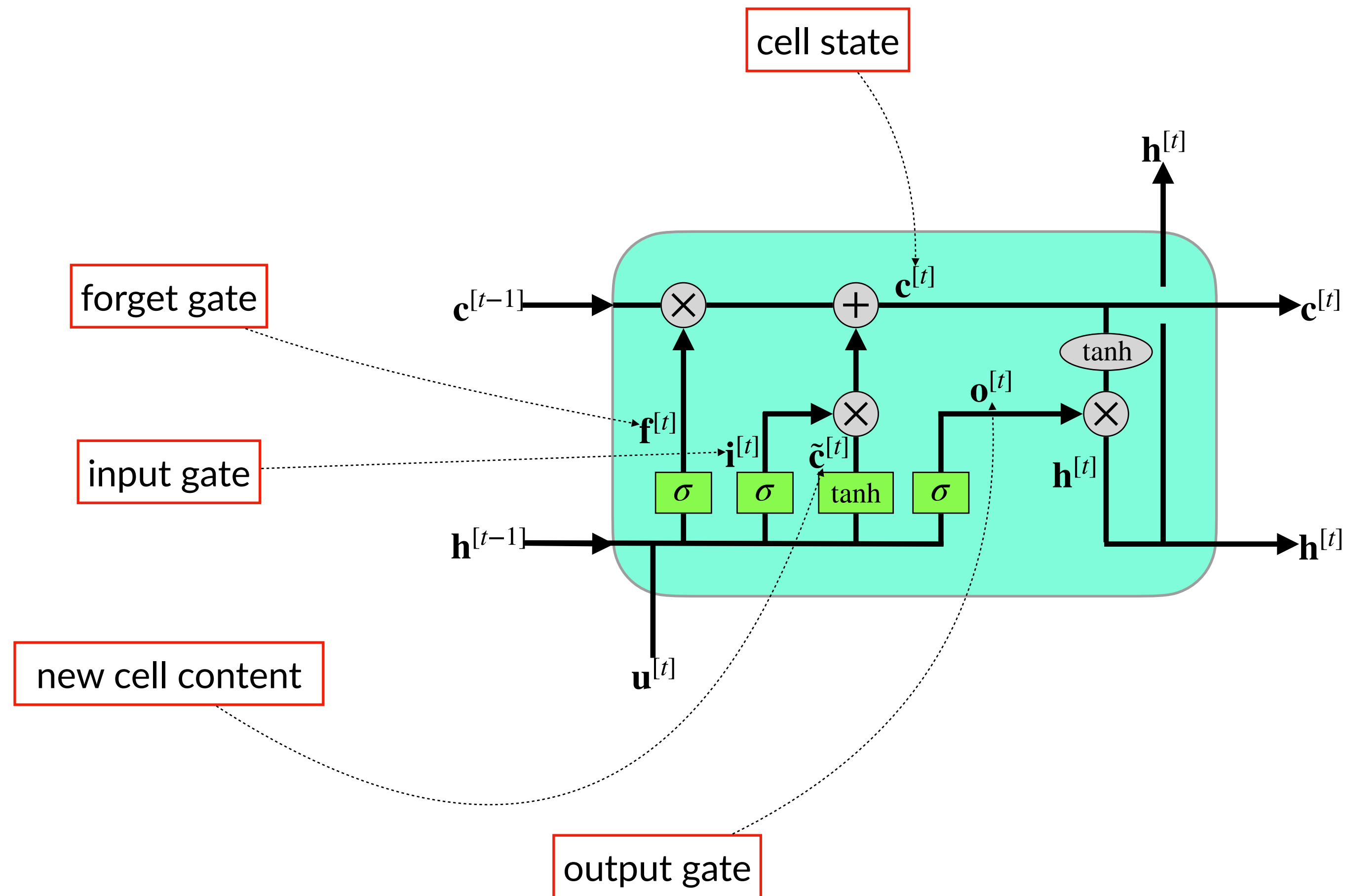
$$\mathbf{i}^{[t]} = \sigma\big(\mathbf{U}_i \cdot \mathbf{u}^{[t]} + \mathbf{W}_i \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_i\big)$$

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

$$\mathbf{o}^{[t]} = \sigma\big(\mathbf{U}_o \cdot \mathbf{u}^{[t]} + \mathbf{W}_o \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_o\big)$$

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

$$\tilde{\mathbf{c}}^{[t]} = \tanh\big(\mathbf{U}_c \cdot \mathbf{u}^{[t]} + \mathbf{W}_c \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_c\big)$$

$$\mathbf{f}^{[t]} = \sigma\Big(\mathbf{U}_f \cdot \mathbf{u}^{[t]} + \mathbf{W}_f \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_f\Big)$$
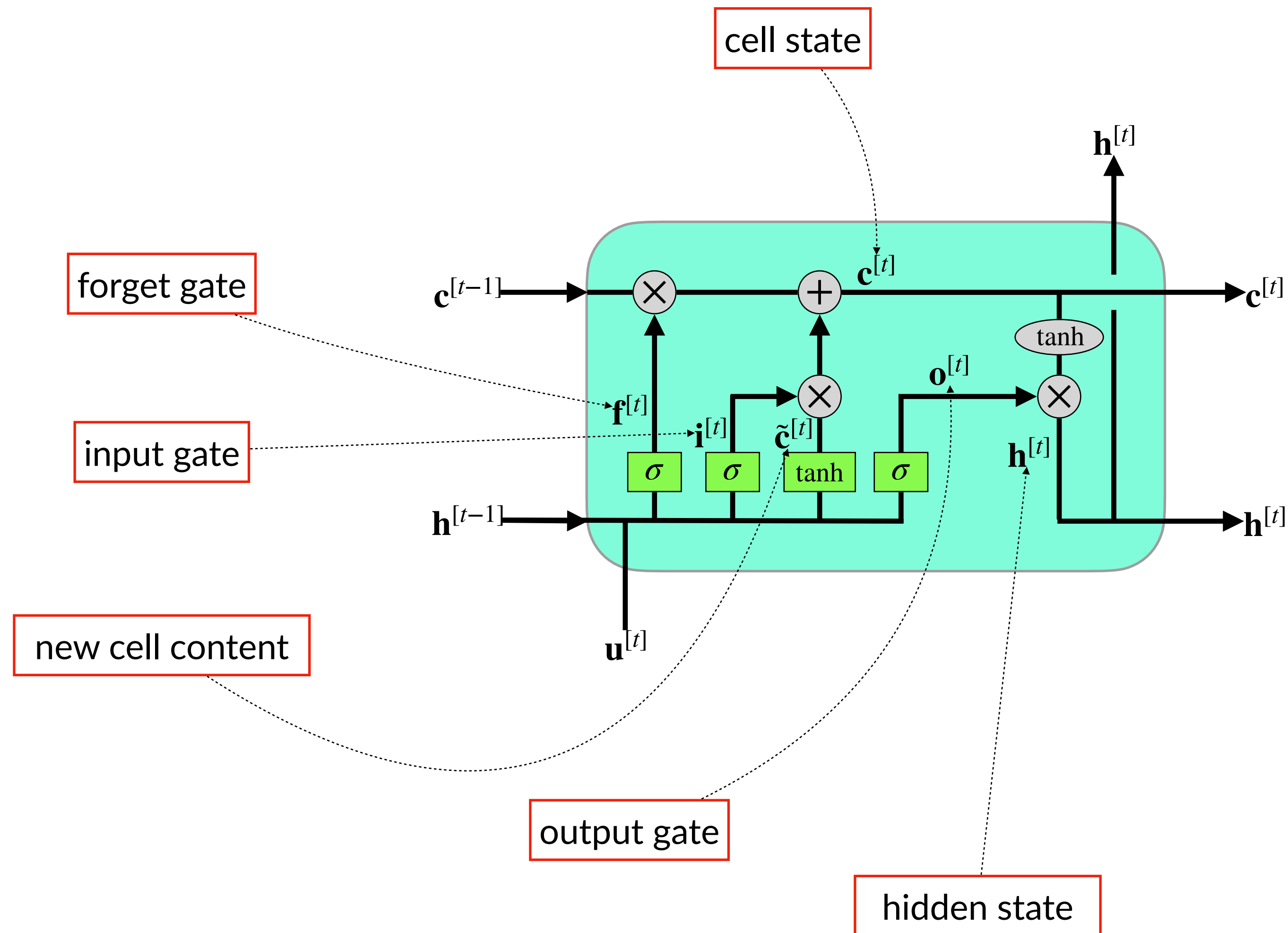
$$\mathbf{i}^{[t]} = \sigma\big(\mathbf{U}_i \cdot \mathbf{u}^{[t]} + \mathbf{W}_i \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_i\big)$$

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

$$\mathbf{o}^{[t]} = \sigma\big(\mathbf{U}_o \cdot \mathbf{u}^{[t]} + \mathbf{W}_o \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_o\big)$$

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

$$\tilde{\mathbf{c}}^{[t]} = \tanh\big(\mathbf{U}_c \cdot \mathbf{u}^{[t]} + \mathbf{W}_c \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_c\big)$$

$$\mathbf{f}^{[t]} = \sigma\Big(\mathbf{U}_f \cdot \mathbf{u}^{[t]} + \mathbf{W}_f \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_f\Big)$$
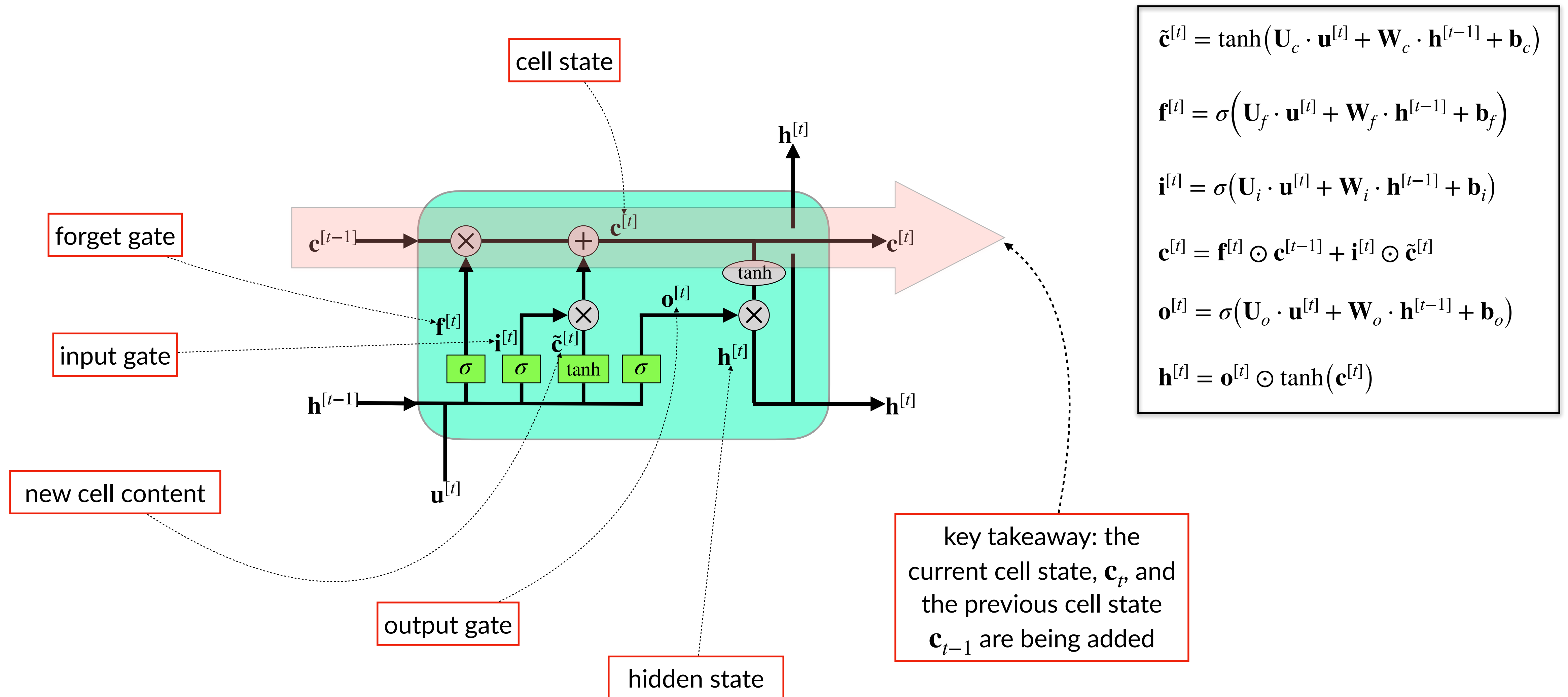
$$\mathbf{i}^{[t]} = \sigma\big(\mathbf{U}_i \cdot \mathbf{u}^{[t]} + \mathbf{W}_i \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_i\big)$$

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

$$\mathbf{o}^{[t]} = \sigma\big(\mathbf{U}_o \cdot \mathbf{u}^{[t]} + \mathbf{W}_o \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_o\big)$$
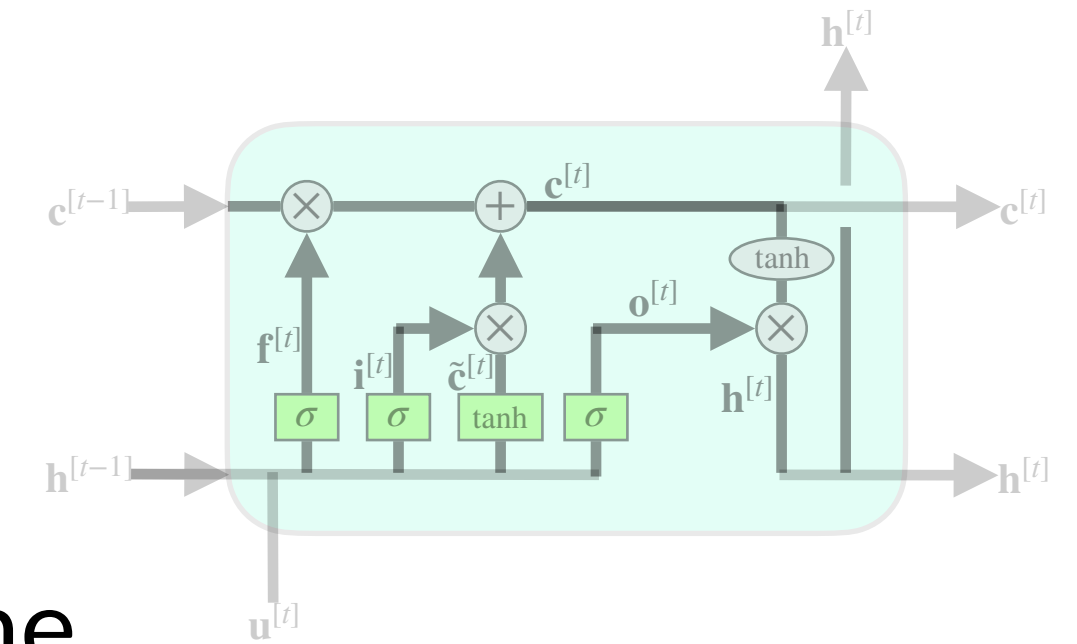
$$\mathbf{h}^{[t]} = \mathbf{o}^{[t]} \odot \tanh\big(\mathbf{c}^{[t]}\big)$$

$$\tilde{\mathbf{c}}^{[t]} = \tanh\left(\mathbf{U}_c \cdot \mathbf{u}^{[t]} + \mathbf{W}_c \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_c\right)$$

$$\mathbf{f}^{[t]} = \sigma\left(\mathbf{U}_f \cdot \mathbf{u}^{[t]} + \mathbf{W}_f \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_f\right)$$

$$\mathbf{i}^{[t]} = \sigma\left(\mathbf{U}_i \cdot \mathbf{u}^{[t]} + \mathbf{W}_i \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_i\right)$$

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

$$\mathbf{o}^{[t]} = \sigma\left(\mathbf{U}_o \cdot \mathbf{u}^{[t]} + \mathbf{W}_o \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_o\right)$$

$$\mathbf{h}^{[t]} = \mathbf{o}^{[t]} \odot \tanh\left(\mathbf{c}^{[t]}\right)$$

cell state

forget gate

input gate

new cell content

$\mathbf{c}^{[t-1]}$    $\mathbf{c}^{[t]}$    $\mathbf{c}^{[t]}$

$\mathbf{h}^{[t]}$

$\mathbf{f}^{[t]}$   $\mathbf{i}^{[t]}$   $\tilde{\mathbf{c}}^{[t]}$   $\mathbf{o}^{[t]}$

$\sigma$   $\sigma$   tanh   $\sigma$

$\mathbf{h}^{[t-1]}$    $\mathbf{h}^{[t]}$    $\mathbf{h}^{[t]}$

$\mathbf{u}^{[t]}$

$$\tilde{\mathbf{c}}^{[t]} = \tanh\left(\mathbf{U}_c \cdot \mathbf{u}^{[t]} + \mathbf{W}_c \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_c\right)$$

$$\mathbf{f}^{[t]} = \sigma\left(\mathbf{U}_f \cdot \mathbf{u}^{[t]} + \mathbf{W}_f \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_f\right)$$

$$\mathbf{i}^{[t]} = \sigma\left(\mathbf{U}_i \cdot \mathbf{u}^{[t]} + \mathbf{W}_i \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_i\right)$$

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

$$\mathbf{o}^{[t]} = \sigma\left(\mathbf{U}_o \cdot \mathbf{u}^{[t]} + \mathbf{W}_o \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_o\right)$$

$$\mathbf{h}^{[t]} = \mathbf{o}^{[t]} \odot \tanh\left(\mathbf{c}^{[t]}\right)$$

$$\tilde{\mathbf{c}}^{[t]} = \tanh\left(\mathbf{U}_c \cdot \mathbf{u}^{[t]} + \mathbf{W}_c \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_c\right)$$

$$\mathbf{f}^{[t]} = \sigma\left(\mathbf{U}_f \cdot \mathbf{u}^{[t]} + \mathbf{W}_f \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_f\right)$$

$$\mathbf{i}^{[t]} = \sigma\left(\mathbf{U}_i \cdot \mathbf{u}^{[t]} + \mathbf{W}_i \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_i\right)$$

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

$$\mathbf{o}^{[t]} = \sigma\left(\mathbf{U}_o \cdot \mathbf{u}^{[t]} + \mathbf{W}_o \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_o\right)$$

$$\mathbf{h}^{[t]} = \mathbf{o}^{[t]} \odot \tanh\left(\mathbf{c}^{[t]}\right)$$

$$\tilde{\mathbf{c}}^{[t]} = \tanh\big(\mathbf{U}_c \cdot \mathbf{u}^{[t]} + \mathbf{W}_c \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_c\big)$$

$$\mathbf{f}^{[t]} = \sigma\Big(\mathbf{U}_f \cdot \mathbf{u}^{[t]} + \mathbf{W}_f \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_f\Big)$$

$$\mathbf{i}^{[t]} = \sigma\big(\mathbf{U}_i \cdot \mathbf{u}^{[t]} + \mathbf{W}_i \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_i\big)$$

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

$$\mathbf{o}^{[t]} = \sigma\big(\mathbf{U}_o \cdot \mathbf{u}^{[t]} + \mathbf{W}_o \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_o\big)$$

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

$$\tilde{\mathbf{c}}^{[t]} = \tanh\big(\mathbf{U}_c \cdot \mathbf{u}^{[t]} + \mathbf{W}_c \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_c\big)$$

$$\mathbf{f}^{[t]} = \sigma\Big(\mathbf{U}_f \cdot \mathbf{u}^{[t]} + \mathbf{W}_f \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_f\Big)$$

$$\mathbf{i}^{[t]} = \sigma\big(\mathbf{U}_i \cdot \mathbf{u}^{[t]} + \mathbf{W}_i \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_i\big)$$

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

$$\mathbf{o}^{[t]} = \sigma\big(\mathbf{U}_o \cdot \mathbf{u}^{[t]} + \mathbf{W}_o \cdot \mathbf{h}^{[t-1]} + \mathbf{b}_o\big)$$

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

cell state

forget gate

input gate

new cell content

output gate

hidden state

key takeaway: the current cell state, $\mathbf{c}_t$, and the previous cell state $\mathbf{c}_{t-1}$ are being added

▸ LSTM can preserve information over many time steps using its gates

▸ LSTM: If the forget gate value is set to $\mathbf{f}_i^{[t]} = 1$ for a cell dimension $i$ and the corresponding input gate value $\mathbf{i}_i^{[t]} = 0$, then the cell value from the previous time step, $\mathbf{c}_i^{[t-1]}$, is maintained intact

▸ Simple RNN: much harder to maintain previous state information given at least an entire row of the recurrent matrix $\mathbf{W}_h$ should be set to $1$ which in turn will invalidate the entire RNN rationale: $\mathbf{h}_j^{[t]} \propto \mathbf{W}_h[j, :] \cdot \mathbf{h}^{[t-1]}$

▸ Depends on the task, but say an RNN can model ~10 time steps accurately, then an LSTM can probably capture ~100 time steps

## Recipe RNN LM output

**Input:** "Fish and chips"

---

**Name:** Fish and chips with Broccoli and Salad of Creamy Thyme Broth

**Ingredients:**
- 1 cup frozen peas, thawed
- 1/4 cup chopped fresh cilantro leaves
- 1 tablespoon finely chopped fresh dill
- 1/2 cup sugar
- 1/2 cup corn tortillas
- 1 cup shredded smoked mozzarella or parmesan cheese
- 1/2 cup white wine
- 1 cup chicken broth
- Salt and pepper

**Instructions:** *Season salad with salt and pepper. In a large saute pan over medium-high heat, cook poblano pepper for 1 minute. Add broccoli rabe, spring onions, thyme, and bay leaves and sprinkle with salt and pepper to taste. Cook until vegetables are soft, about 10 minutes. Add the spinach and stir until completely melted. Add sugar and simmer until sauce thickens, about 1 minute. Remove from heat and stir in lemon juice. Serve with steamed roasted garlic bread.*

Certain LSTM cells "learn" to have larger values…

*towards the end of a line*



*inside if statements*



Source: karpathy.github.io/2015/05/21/rnn-effectiveness/

*COMP0087 - Recurrent Neural Networks*

Certain LSTM cells "learn" to have larger values…

*when the code expression's depth increases*

*inside comments or double quotes*

Source: karpathy.github.io/2015/05/21/rnn-effectiveness/

e.g. tasks like **part-of-speech** (POS) tagging and **named entity recognition** (NER)

noun    noun    verb    determiner    adjective    noun

Windows    Me    was    a    great    failure

negative

combine all hidden
states (element-
wise mean or max)

e.g. *text / sentence,*
*sentiment classification*

Windows    Me    was    a    great    failure

Windows   Me   was   a   great   failure

"forward" RNN

Windows    Me    was    a    great    failure

"backward" RNN

"forward" RNN

Windows    Me    was    a    great    failure

Hidden state
via concatenation
has context from
both directions

"backward" RNN

"forward" RNN

Windows    Me    was    a    great    failure

*"great"* product vs. *"great"* failure

Hidden state via concatenation has context from both directions

"backward" RNN

"forward" RNN

Windows   Me   was   a   great   failure

Hidden state via concatenation has context from both directions

"backward" RNN

"forward" RNN

Windows    Me    was    a    great    failure

hidden state of the bidirectional RNN

$$\mathbf{h}^{[t]} = \left[ \overrightarrow{\mathbf{h}}^{[t]} ; \overleftarrow{\mathbf{h}}^{[t]} \right]$$

$$\overleftarrow{\mathbf{h}}^{[t]} = \mathrm{RNN}_{\mathrm{B}} \left( \overleftarrow{\mathbf{h}}^{[t+1]}, \mathbf{u}^{[t]} \right)$$
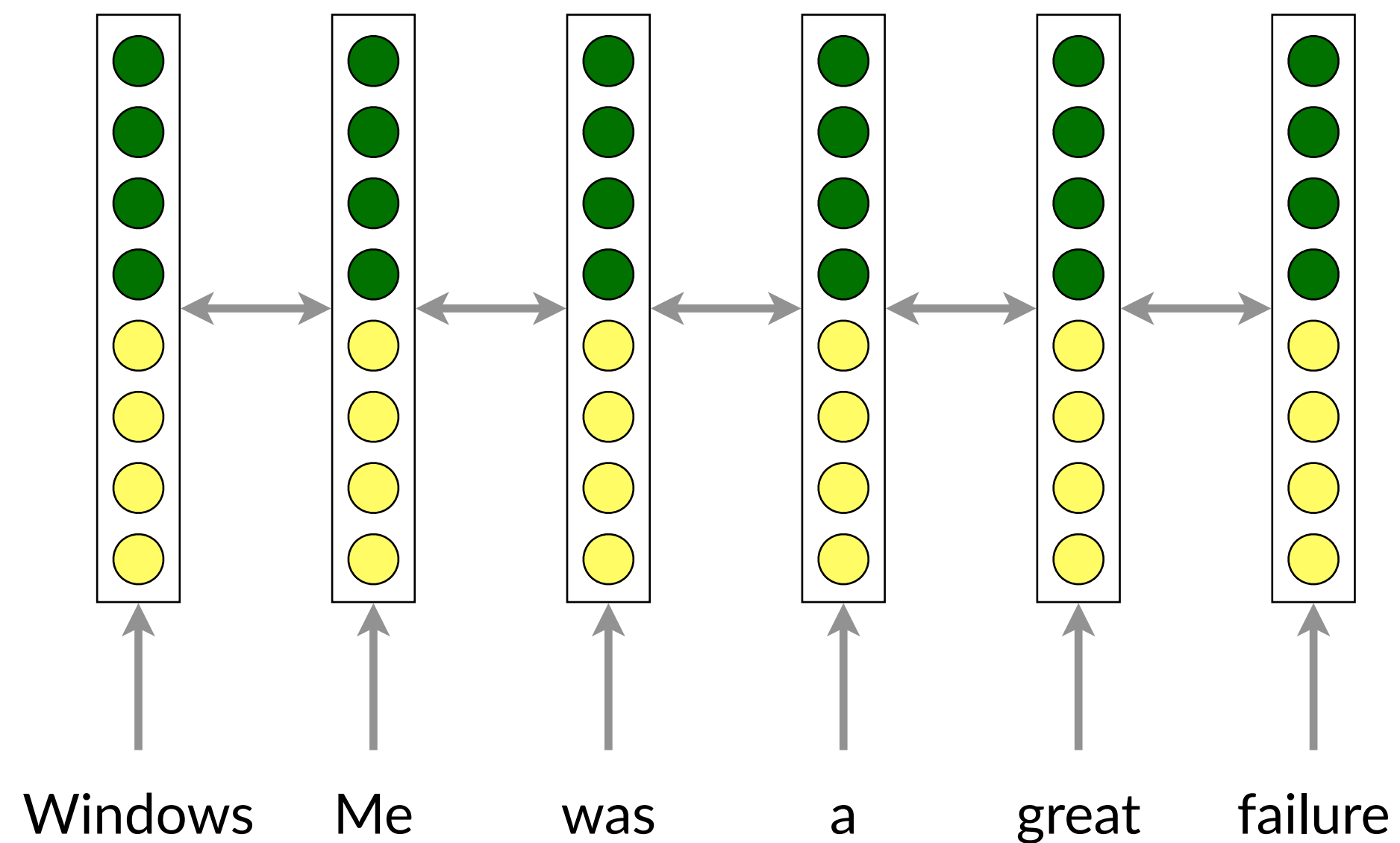
$$\overrightarrow{\mathbf{h}}^{[t]} = \mathrm{RNN}_{\mathrm{F}} \left( \overrightarrow{\mathbf{h}}^{[t-1]}, \mathbf{u}^{[t]} \right)$$

different weights

Windows    Me    was    a    great    failure
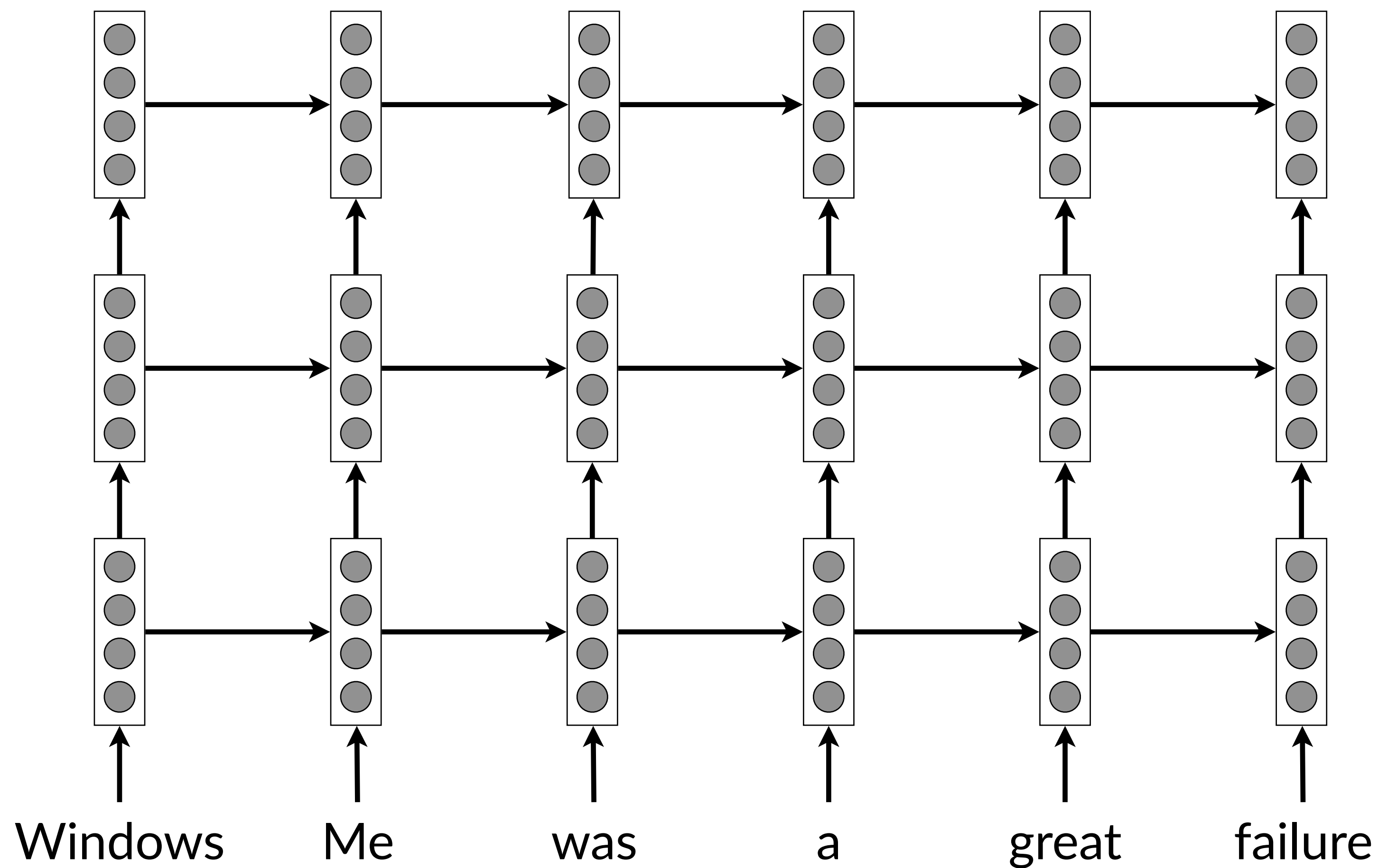
▶ Bidirectional RNNs are very effective in sequence classification

▶ They requires access to the entire sequence, i.e. not necessarily great for language models (*text generators*)

▶ Bidirectional NNs are strong predictors, i.e. BERT: Bidirectional Encoder Representations from Transformers
  aclanthology.org/N19-1423.pdf

*the output of one RNN layer (hidden state) becomes the input to the next*

Windows    Me    was    a    great    failure

*the output of one RNN layer (hidden state) becomes the input to the next*

*higher-level features*

*lower-level features*

Windows    Me    was    a    great    failure

▶ Monday, March 18 (*last week*)

▶ Self-invited "guest" lecture on "*Modelling infectious disease prevalence using web search activity*"