# Recurrent Networks –take 3

# Unboundedness in time



$$h = \sigma(W^{(h)}x + b^{(h)})$$
$$y = \text{softmax}(W^{(o)}h + b^{(o)})$$

$$h^{(t)} = \sigma(W^{(h)}x^{(t)} + Uh^{(t-1)} + b^{(h)})$$
$$y^{(t)} = \text{softmax}(W^{(o)}h^{(t)} + b^{(o)})$$

Simple Recurrent Network (SRN/RNN)
(Elman 1990)

# Misbehaving Gradients

$$z^{(t)} = W^{(h)}x^{(t)} + Uh^{(t-1)} + b^{(h)}$$

$$h^{(t)} = \sigma(z^{(t)})$$

$$\hat{y}^{(t)} = \text{softmax}\left(W^{(o)}h^{(t)} + b^{(o)}\right)$$

$$L^{(t)} = - y^{(t)}\log(\hat{y}^{(t)})$$

$$\frac{\partial \mathbf{L}^{(3)}}{\partial \mathbf{U}} = \frac{\partial \mathbf{L}^{(3)}}{\partial \hat{\mathbf{y}}^{(3)}} \frac{\partial \hat{\mathbf{y}}^{(3)}}{\partial \mathbf{h}^{(3)}} \frac{\partial \mathbf{h}^{(3)}}{\partial \mathbf{z}^{(3)}} \frac{\partial \mathbf{z}^{(3)}}{\partial \mathbf{U}}$$

$$\frac{\partial \mathbf{z}^{(3)}}{\partial \mathbf{U}} = \mathbf{U}\frac{\partial \mathbf{h}^{(2)}}{\partial \mathbf{U}} + \frac{\partial \mathbf{U}}{\partial \mathbf{U}}\mathbf{h}^{(2)} \qquad \frac{\partial \sigma(\mathbf{x})}{\partial \mathbf{x}} = \sigma(\mathbf{x})(1 - \sigma(\mathbf{x}))$$

$$= \mathbf{U}\left(\frac{\partial \mathbf{h}^{(2)}}{\partial \mathbf{z}^{(2)}} \frac{\partial \mathbf{z}^{(2)}}{\partial \mathbf{U}}\right) + \mathbf{h}^{(2)}$$

$$= \mathbf{U}\left(\frac{\partial \mathbf{h}^{(2)}}{\partial \mathbf{z}^{(2)}}\left(\mathbf{U}\frac{\partial \mathbf{h}^{(1)}}{\partial \mathbf{U}} + \mathbf{h}^{(1)}\right)\right) + \mathbf{h}^{(2)}$$

$$= \mathbf{U}\left(\frac{\partial \mathbf{h}^{(2)}}{\partial \mathbf{z}^{(2)}}\left(\mathbf{U}\left(\frac{\partial \mathbf{h}^{(1)}}{\partial \mathbf{z}^{(1)}}\left(\mathbf{U}\frac{\partial \mathbf{h}^{(0)}}{\partial \mathbf{U}} + \mathbf{h}^{(0)}\right) + \mathbf{h}^{(1)}\right)\right) + \mathbf{h}^{(2)}$$

# The Consequences of Vanishing Gradients

- **LM task:**
  *When she tried to print her tickets, she found that the printer was out of toner. She went to the stationery store to buy more toner. It was very overpriced. After installing the toner into the printer, she finally printed her _____*

- If the gradient $\frac{\partial L^{(45)}}{\partial x^{(7)}}$ is too small, the network won't learn this (or any other) long-distance dependency!
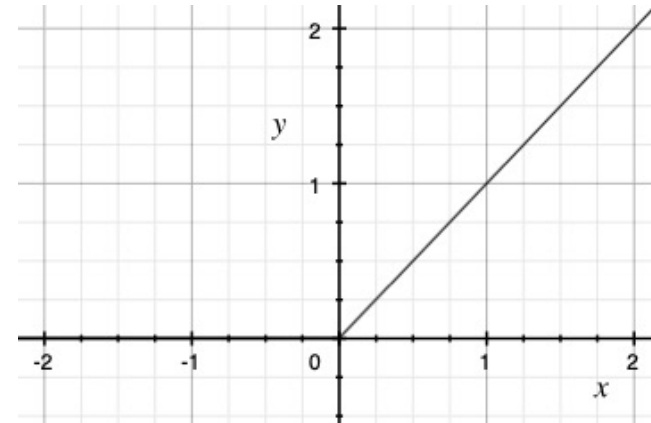
# Gradient solutions

- **Exploding gradients:** gradient clipping

  - If $\|g\| > \text{threshold}, g \longleftarrow \dfrac{\text{threshold}}{\|g\|} g$

- **Vanishing gradients:**

  - change activation function to RELU

  - gating networks

# Gating Networks

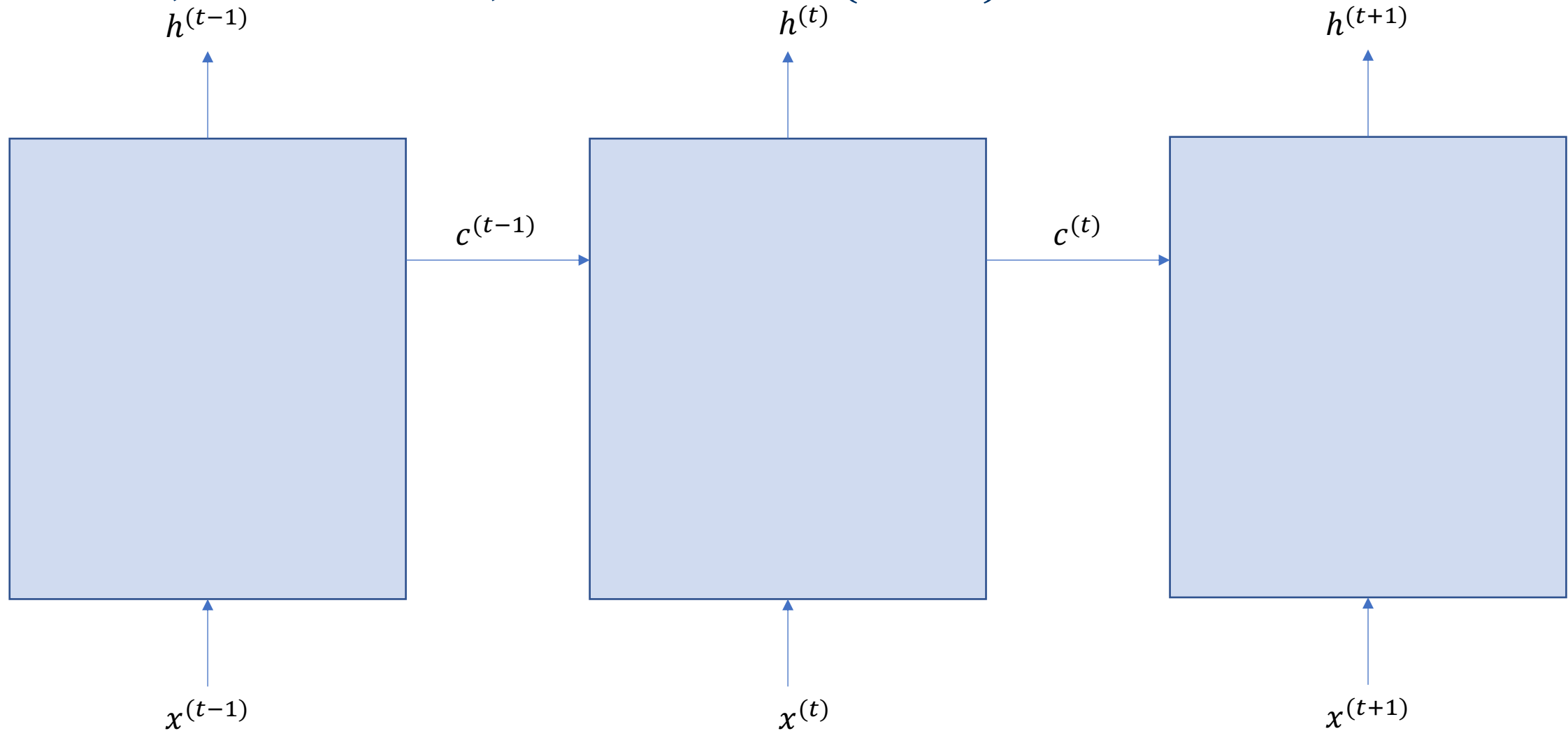- s: storage

- x: input

- g: gating vector

$$\mathbf{s}^{(t)} = \mathbf{x} \odot \mathbf{g} + \mathbf{s}^{(t-1)} \odot (1 - \mathbf{g})$$

$$\begin{bmatrix} 3 \\ 4 \\ 6 \end{bmatrix} \odot \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} + \begin{bmatrix} 8 \\ 9 \\ 2 \end{bmatrix} \odot \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$
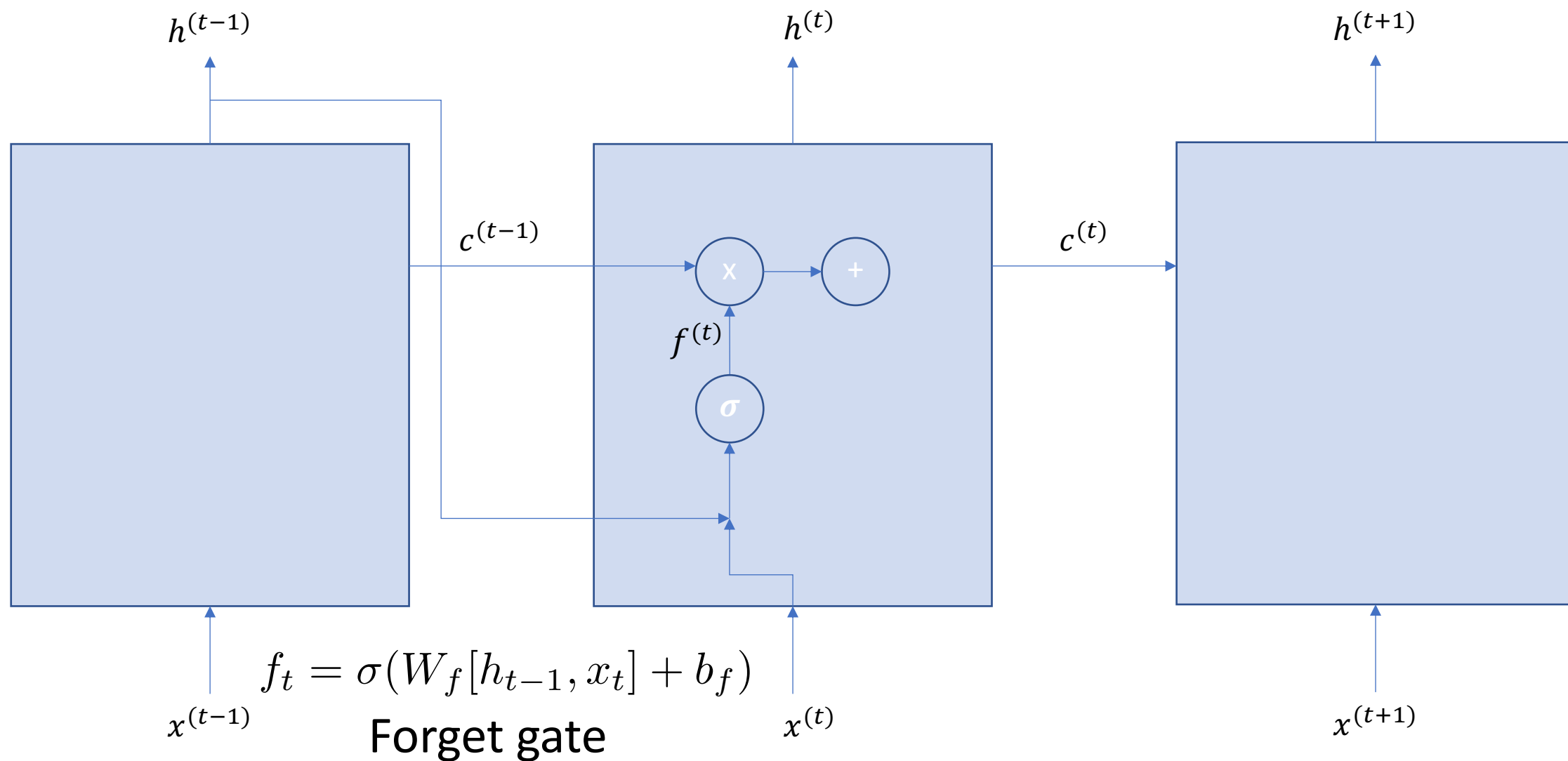
# Long Short-Term Memory (LSTM)

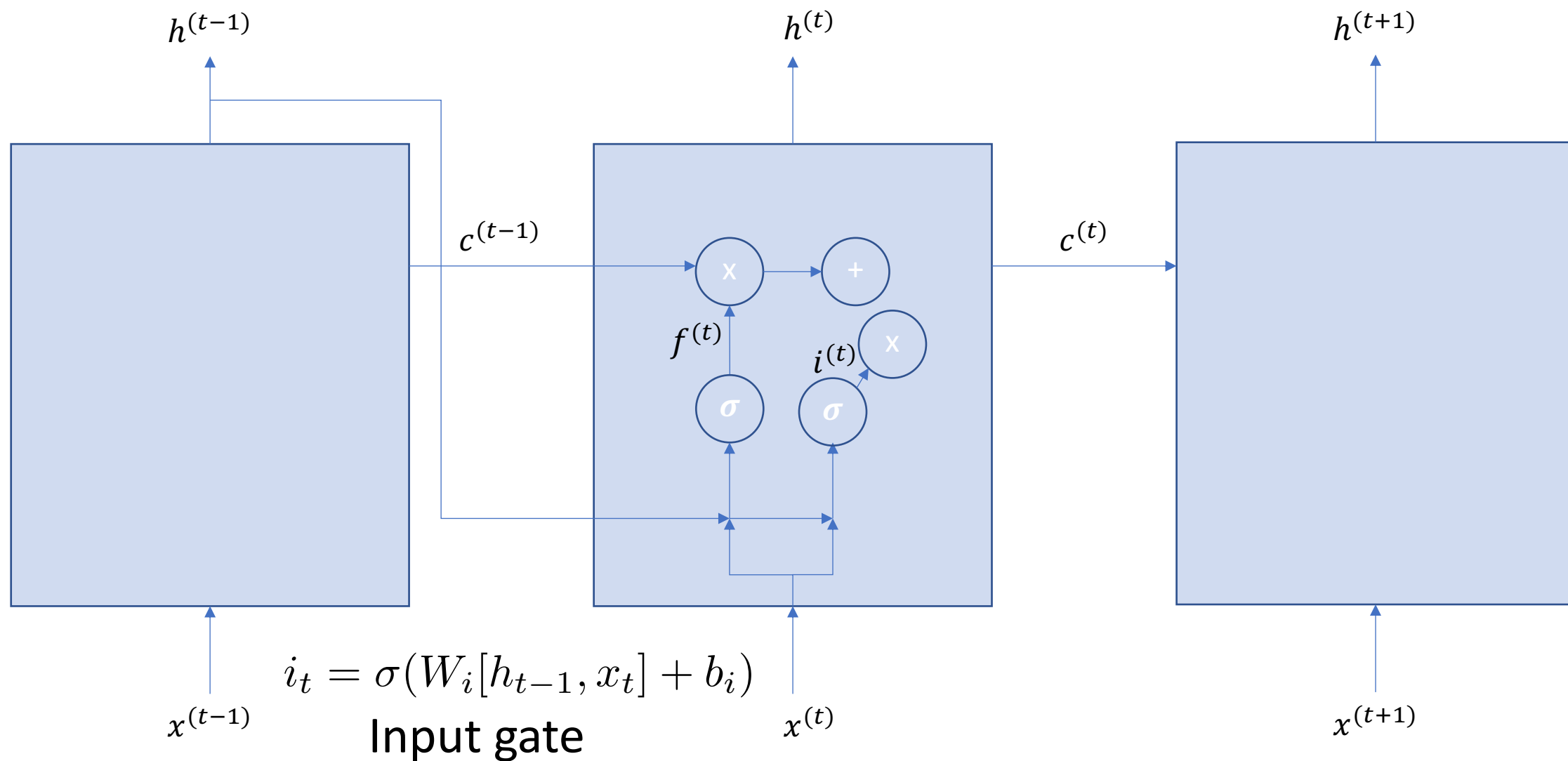## Hochreiter and Schmidhuber (1997)
## Gers, Schmidhuber, and Cummins (2000)

# Long Short-Term Memory (LSTM)



$$f_t = \sigma(W_f[h_{t-1}, x_t] + b_f)$$

Forget gate

# Long Short-Term Memory (LSTM)



$$i_t = \sigma(W_i[h_{t-1}, x_t] + b_i)$$

Input gate

# Long Short-Term Memory (LSTM)



$$\tilde{c}_t = \tanh(W_c[h_{t-1}, x_t] + b_c)$$

Candidate cell state

# Long Short-Term Memory (LSTM)



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

$h^{(t-1)}$

$h^{(t)}$

$h^{(t+1)}$

$c^{(t-1)}$

$c^{(t)}$

$f^{(t)}$

$i^{(t)}$

$\tilde{c}^{(t)}$

$x^{(t-1)}$

$x^{(t)}$

$x^{(t+1)}$

# Long Short-Term Memory (LSTM)

Output gate

$$o_t = \sigma(W_o[h_{t-1}, x_t] + b_o)$$

# Long Short-Term Memory (LSTM)

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

# Long Short-Term Memory (LSTM)

$$f_t = \sigma(W_f[h_{t-1}, x_t] + b_f)$$

$$i_t = \sigma(W_i[h_{t-1}, x_t] + b_i)$$

$$o_t = \sigma(W_o[h_{t-1}, x_t] + b_o)$$

$$\tilde{c}_t = \tanh(W_c[h_{t-1}, x_t] + b_c)$$

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

Important point:

$\dfrac{\partial h^{(t)}}{\partial x^{(t-k)}}$ involves the application of the chain rule across only two activation functions

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

# Gated Recurrent Unit (GRU)
## Cho et al. (2014)

$$z^{(t)} = \sigma(W_z[h^{(t-1)}, x^{(t)}] + b_z)$$

$$h^{(t)} = (1 - z^{(t)}) \odot h^{(t-1)} + z^{(t)} \odot \tilde{h}^{(t)}$$
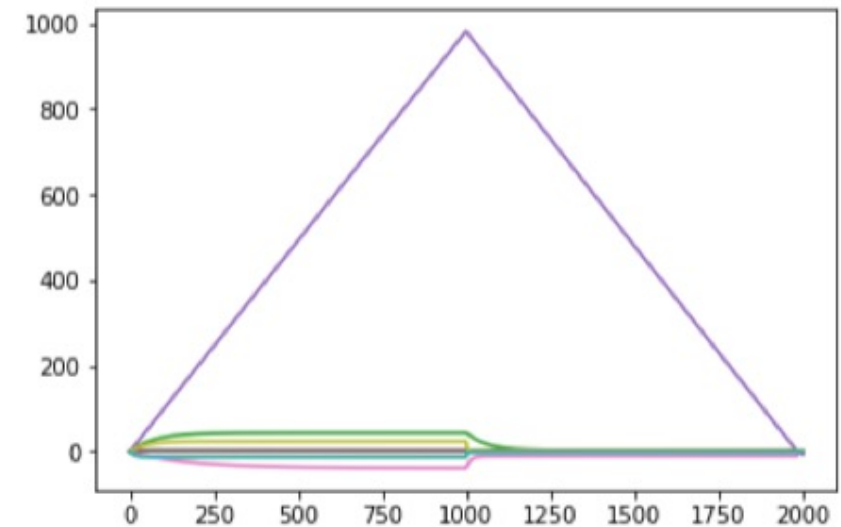
$$\tilde{h}^{(t)} = \tanh(W[r^{(t)} \odot h^{(t-1)}, x^{(t)}])$$

$$r^{(t)} = \sigma(W_r[h^{(t-1)}, x^{(t)}] + b_r)$$

# LSTM-GRU differences

- Weiss, Goldberg and Yahav (2017): train 10d hidden unit networks on acceptance task for $a^n b^n$ (up to $n = 100$)

- Generalization results:

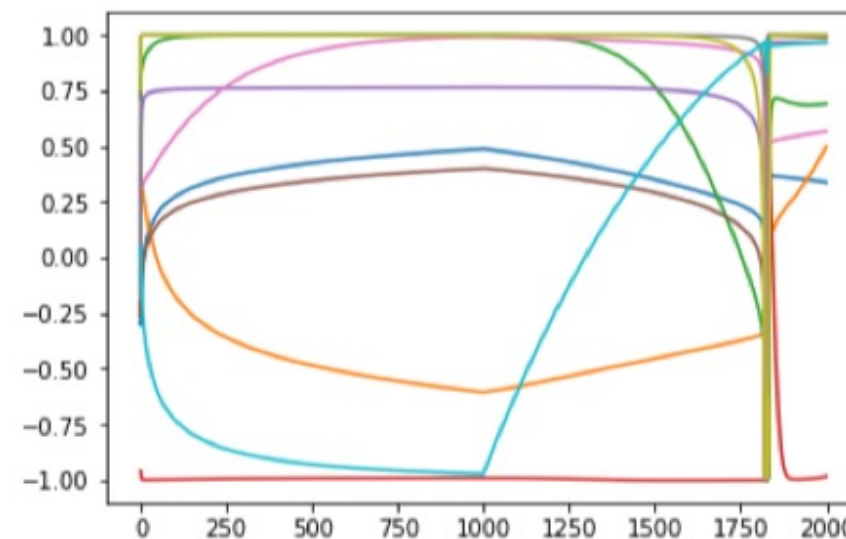| LSTM | $n = 256$ |
|---|---|
| GRU | $n > 37$ accepts $a^n b^{n+1}$<br>$n > 97$ accepts $a^n b^{n+2}$ |



(a) $a^n b^n$-LSTM on $a^{1000} b^{1000}$

# LSTM-GRU differences

- Weiss, Goldberg and Yahav (2017): train 10d hidden unit networks on acceptance task for $a^n b^n$ (up to $n = 100$)

- Generalization results:

| LSTM | $n = 256$ |
|------|-----------|
| GRU | $n > 37$ accepts $a^n b^{n+1}$ <br> $n > 97$ accepts $a^n b^{n+2}$ |



(c) $a^n b^n$-GRU on $a^{1000} b^{1000}$

# LSTM-GRU differences

LSTM

GRU

Output of sigmoid: [0,1]

Output of tanh: [-1,1]

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

Output of sigmoid: [0,1]

1-output of sigmoid:
convex combination of
the present and the past

LSTM cell states can count!

GRU hidden states can't!

# Vanishing Gradients Everywhere!

- We also find vanishing gradients in deep MLPs

  - Chain rule + activation functions will result in the gradient getting ever smaller as it propogates backward.

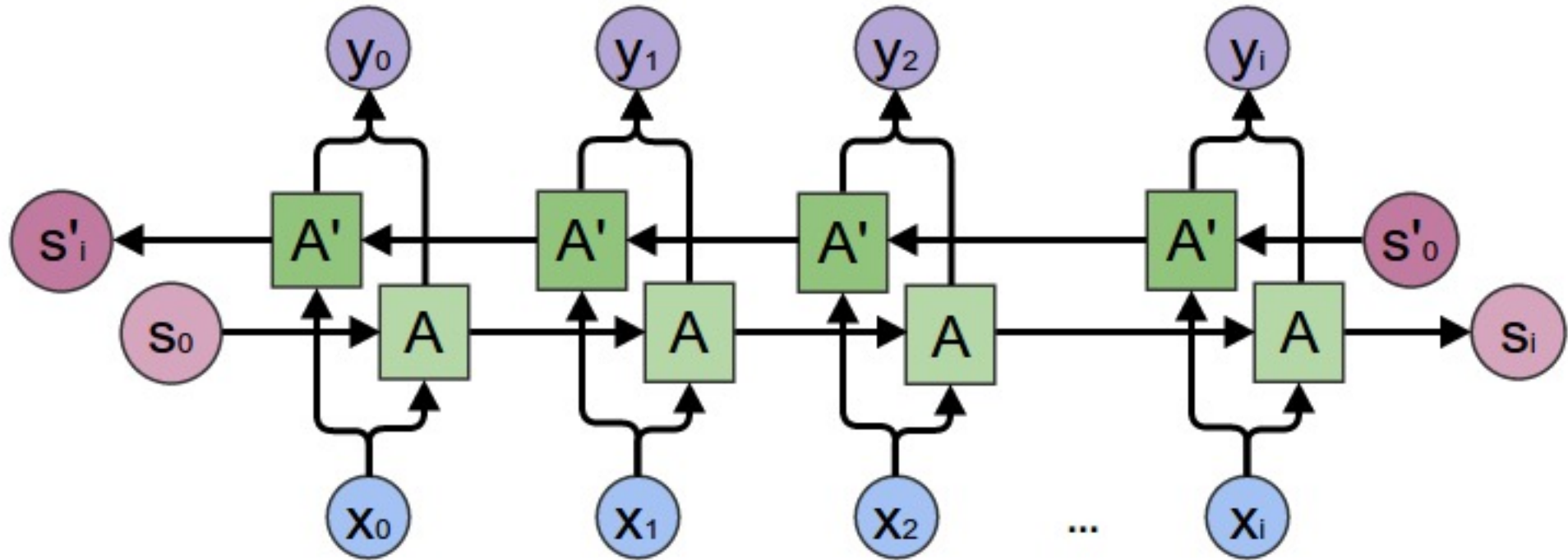  - He et al.'s 2015 proposal: Residual (skip) connections



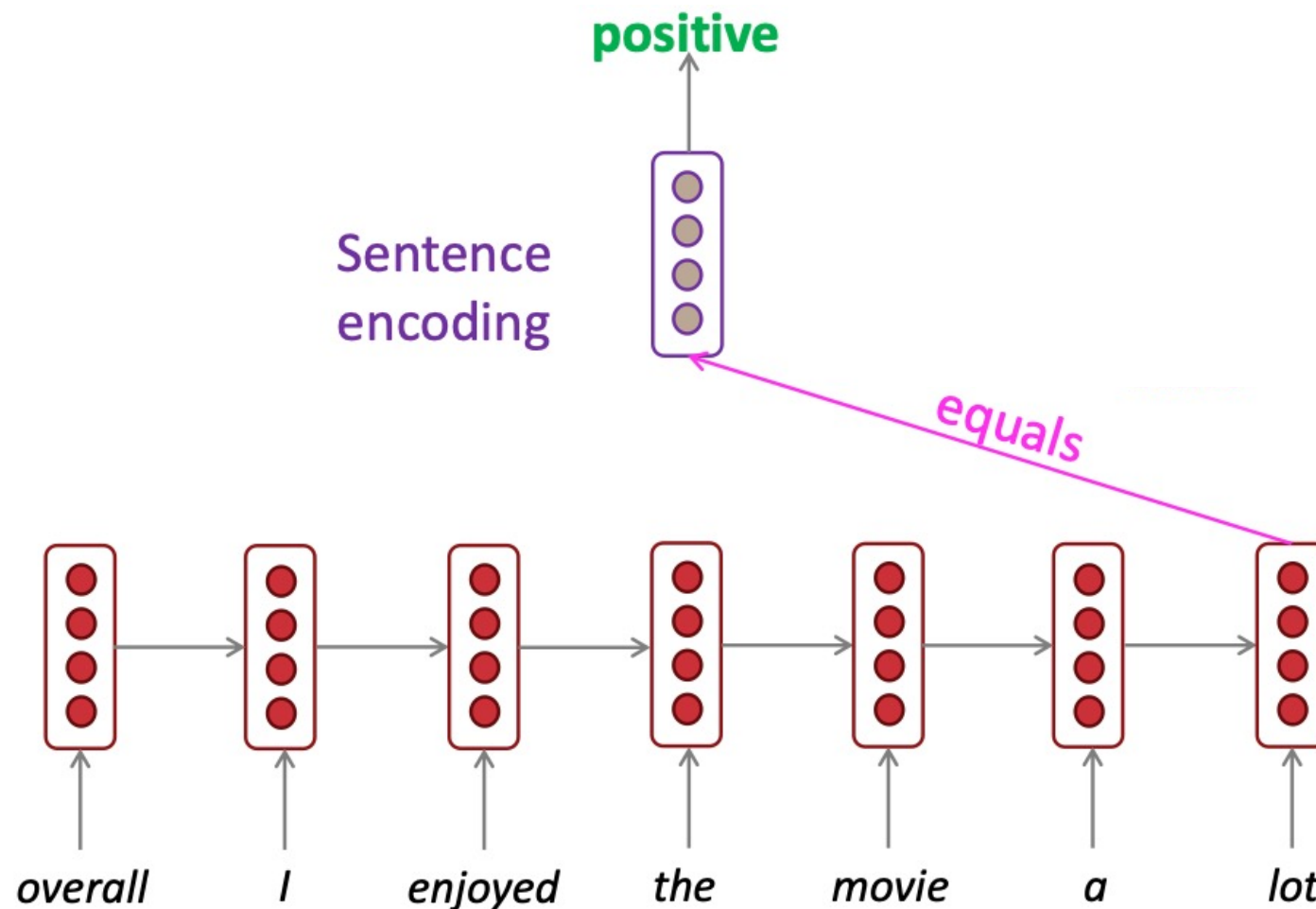| method | top-1 err. | top-5 err. |
|---|---|---|
| VGG [40] (ILSVRC'14) | - | 8.43[†] |
| GoogLeNet [43] (ILSVRC'14) | - | 7.89 |
| VGG [40] (v5) | 24.4 | 7.1 |
| PReLU-net [12] | 21.59 | 5.71 |
| BN-inception [16] | 21.99 | 5.81 |
| ResNet-34 B | 21.84 | 5.71 |
| ResNet-34 C | 21.53 | 5.60 |
| ResNet-50 | 20.74 | 5.25 |
| ResNet-101 | 19.87 | 4.60 |
| ResNet-152 | **19.38** | **4.49** |

Table 4. Error rates (%) of **single-model** results on the ImageNet validation set (except [†] reported on the test set).
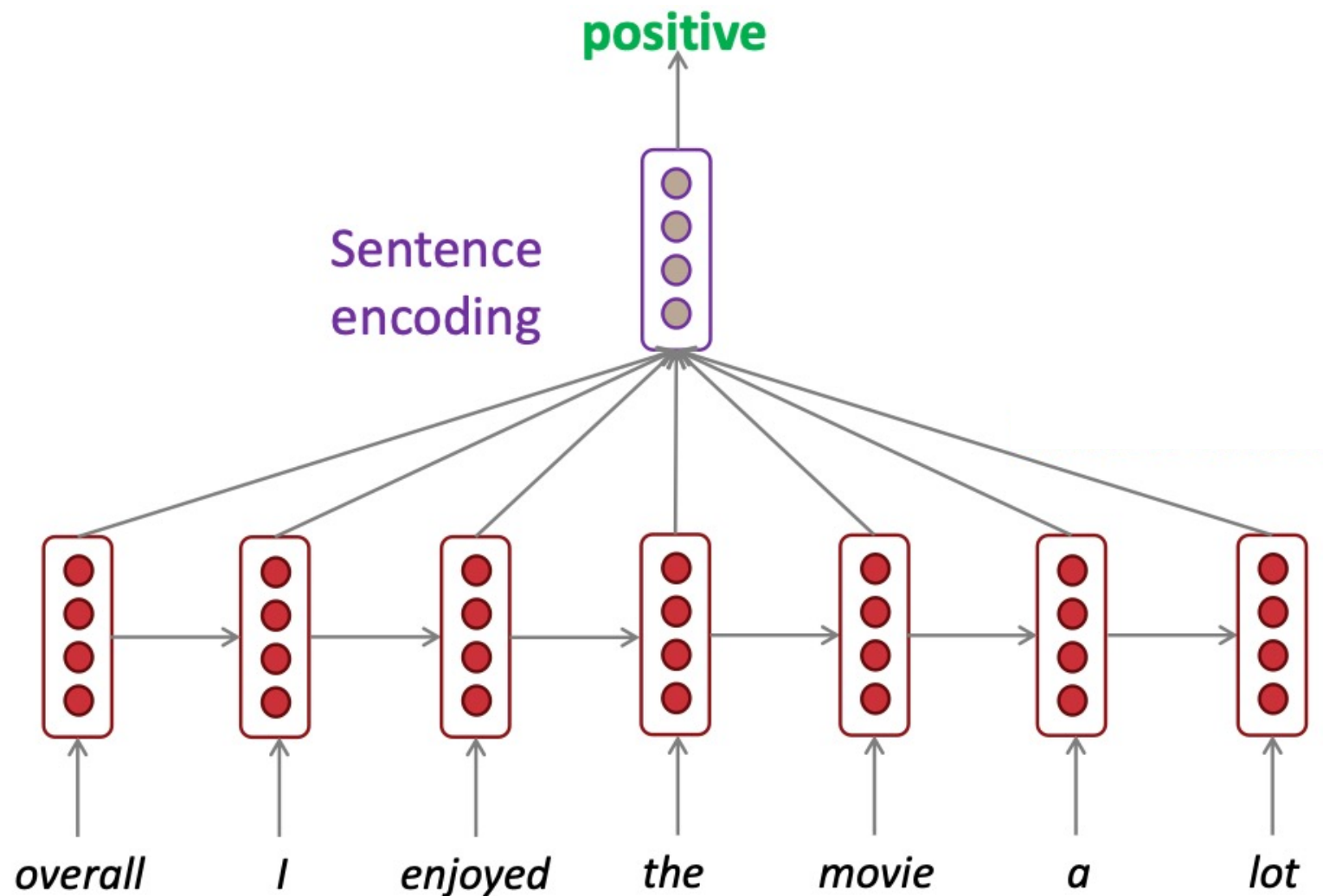
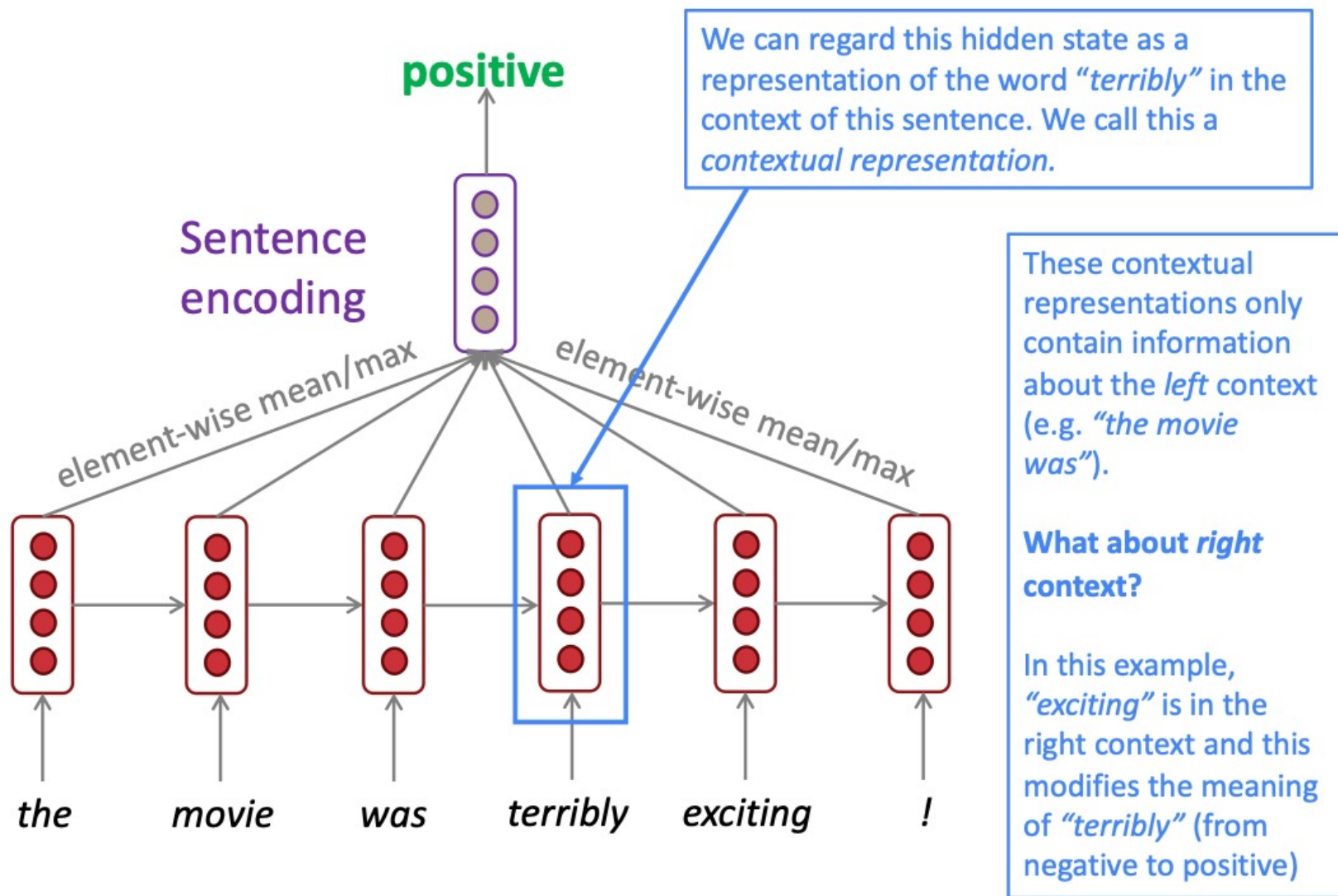# RNN Architectural variations

- Bidirectional RNNs
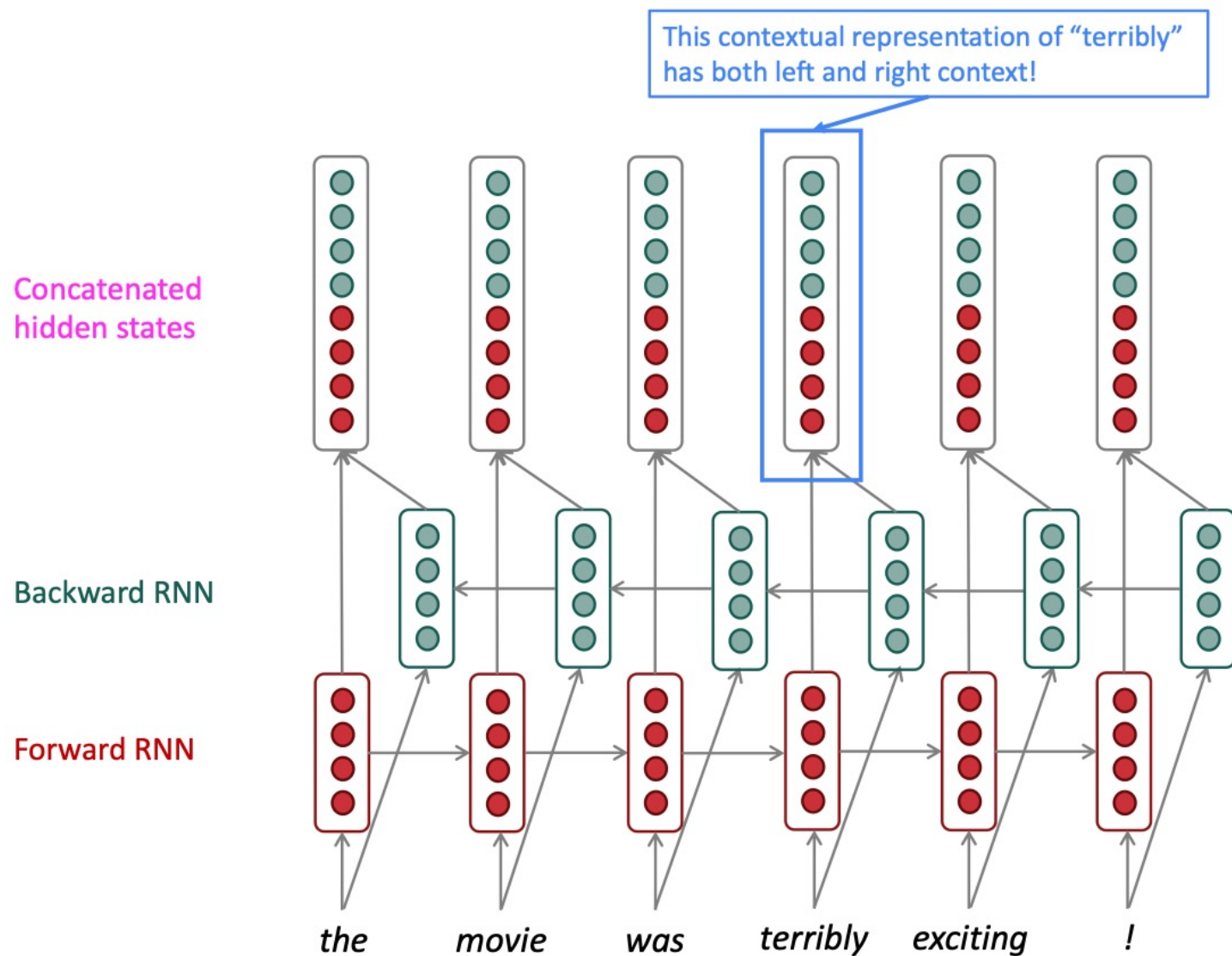
# Bidirectional RNNs in sentiment analysis

# Bidirectional RNNs in sentiment analysis

# Bidirectional RNNs in sentiment analysis



positive

Sentence encoding

element-wise mean/max

element-wise mean/max

the    movie    was    terribly    exciting    !

We can regard this hidden state as a representation of the word *"terribly"* in the context of this sentence. We call this a *contextual representation.*

These contextual representations only contain information about the *left* context (e.g. *"the movie was"*).

**What about *right* context?**

In this example, *"exciting"* is in the right context and this modifies the meaning of *"terribly"* (from negative to positive)

# Bidirectional RNNs in sentiment analysis

This contextual representation of "terribly" has both left and right context!

Concatenated hidden states

Backward RNN

Forward RNN

the   movie   was   terribly   exciting   !

# Architectural variations

- Multilayer RNNs