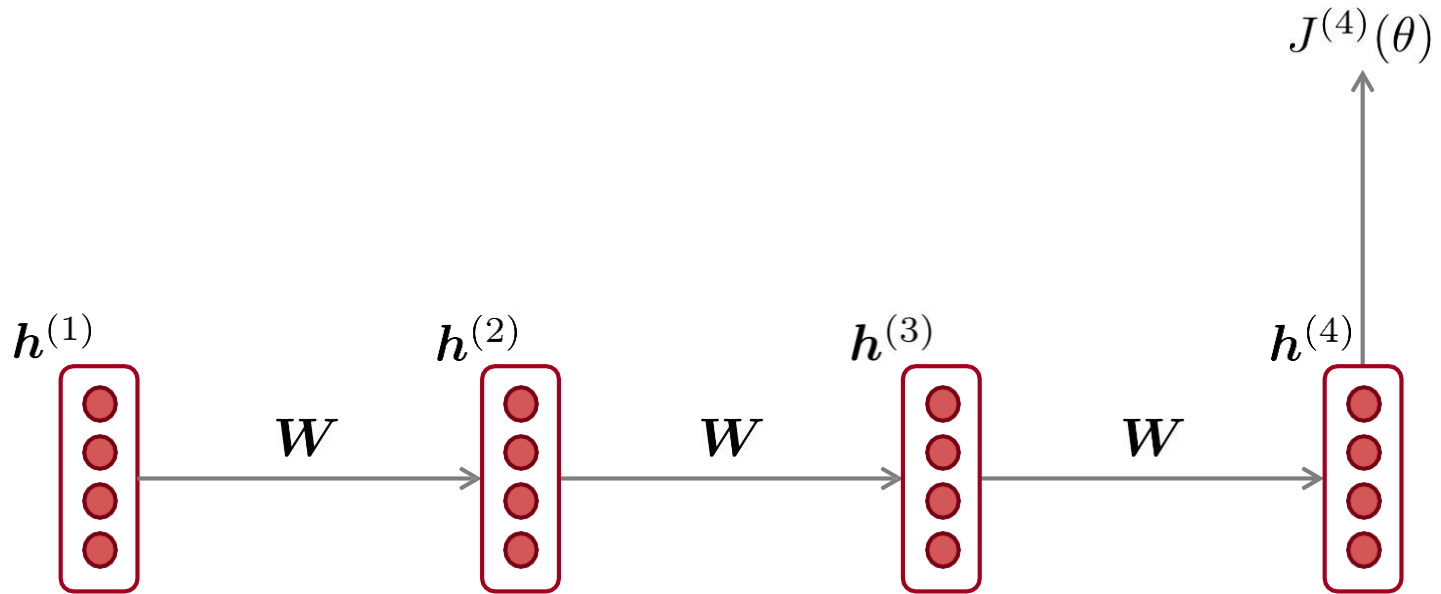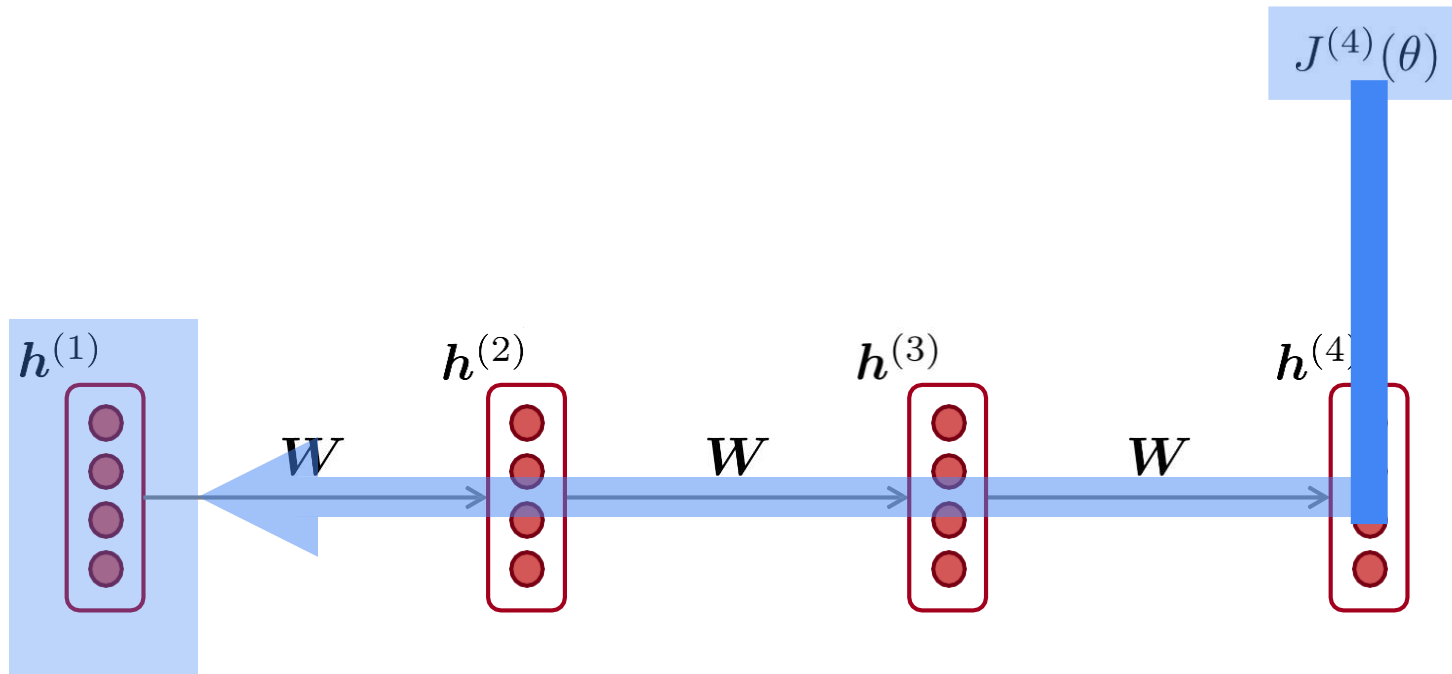# Recurrent Neural Networks (Part II)

# (Issues with RNNs, Adv. RNNs)

# Vanishing gradient intuition
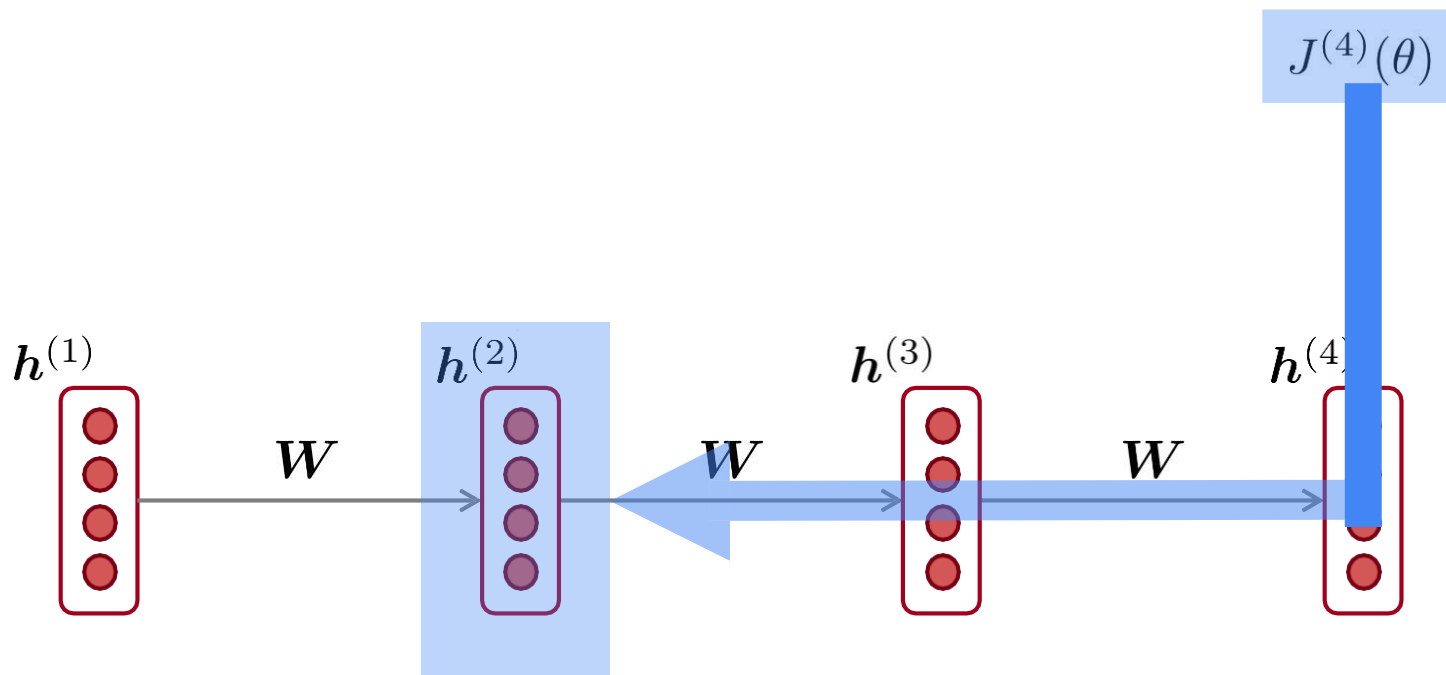
# Vanishing gradient intuition



$$\frac{\partial J^{(4)}}{\partial \boldsymbol{h}^{(1)}} = \quad ?$$

# Vanishing gradient intuition

$$J^{(4)}(\theta)$$

$$h^{(1)} \quad\quad W \quad\quad h^{(2)} \quad\quad W \quad\quad h^{(3)} \quad\quad W \quad\quad h^{(4)}$$



$$\frac{\partial J^{(4)}}{\partial \boldsymbol{h}^{(1)}} = \frac{\partial \boldsymbol{h}^{(2)}}{\partial \boldsymbol{h}^{(1)}} \times \frac{\partial J^{(4)}}{\partial \boldsymbol{h}^{(2)}}$$

chain rule!

# Vanishing gradient intuition

$$J^{(4)}(\theta)$$

$$h^{(1)} \qquad\qquad h^{(2)} \qquad\qquad h^{(3)} \qquad\qquad h^{(4)}$$
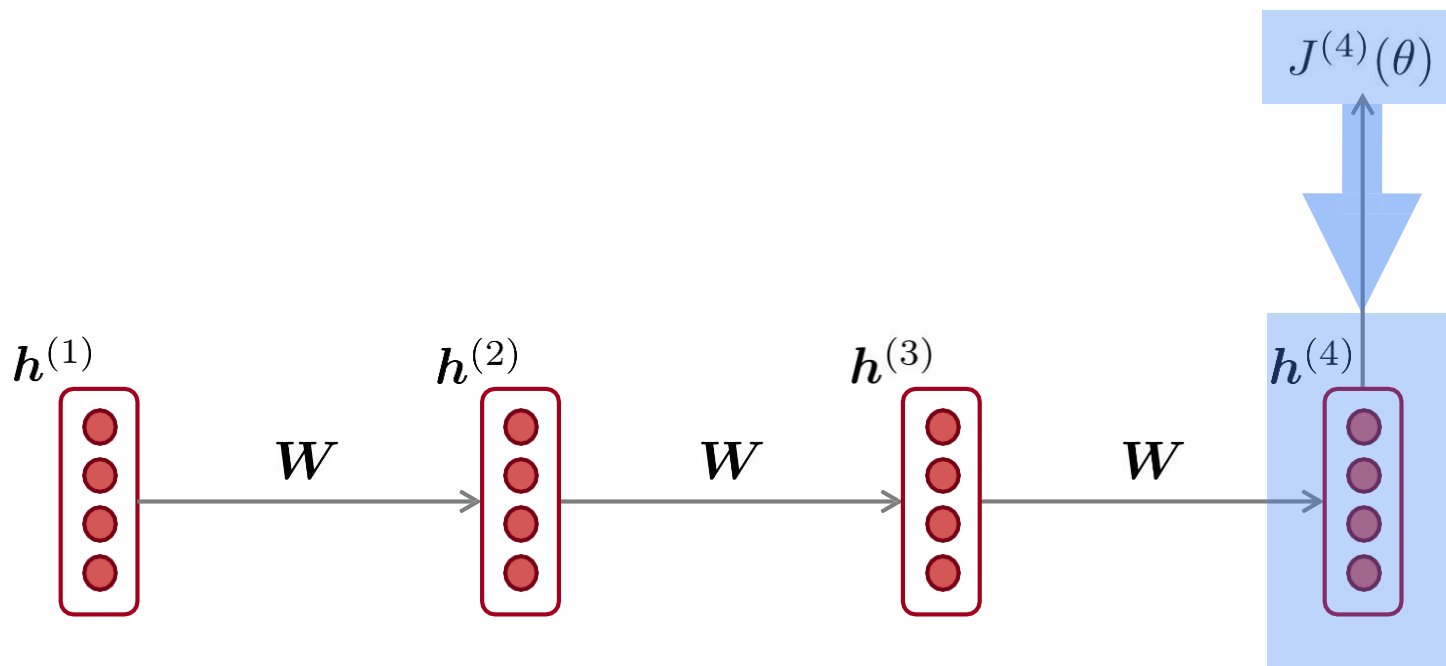
$$W \qquad\qquad W \qquad\qquad W$$

$$\frac{\partial J^{(4)}}{\partial h^{(1)}} = \quad \frac{\partial h^{(2)}}{\partial h^{(1)}} \times \qquad \frac{\partial h^{(3)}}{\partial h^{(2)}} \times \quad \frac{\partial J^{(4)}}{\partial h^{(3)}}$$
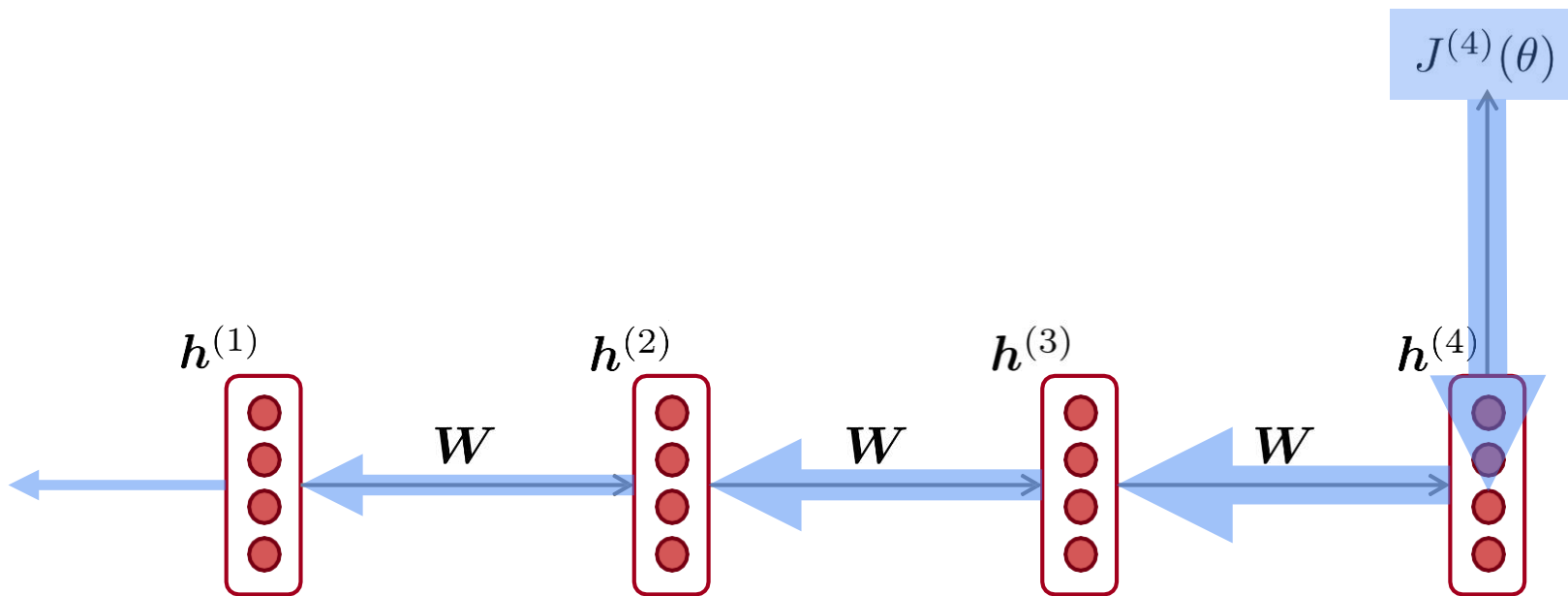
chain rule!

# Vanishing gradient intuition

$$J^{(4)}(\theta)$$

$$\boldsymbol{h}^{(1)} \qquad \boldsymbol{h}^{(2)} \qquad \boldsymbol{h}^{(3)} \qquad \boldsymbol{h}^{(4)}$$

$$\boldsymbol{W} \qquad \boldsymbol{W} \qquad \boldsymbol{W}$$

$$\frac{\partial J^{(4)}}{\partial \boldsymbol{h}^{(1)}} = \quad \frac{\partial \boldsymbol{h}^{(2)}}{\partial \boldsymbol{h}^{(1)}} \times \qquad \frac{\partial \boldsymbol{h}^{(3)}}{\partial \boldsymbol{h}^{(2)}} \times \qquad \frac{\partial \boldsymbol{h}^{(4)}}{\partial \boldsymbol{h}^{(3)}} \times \quad \frac{\partial J^{(4)}}{\partial \boldsymbol{h}^{(4)}}$$

chain rule!

# Vanishing gradient intuition



$J^{(4)}(\theta)$

$$\boldsymbol{h}^{(1)} \qquad \boldsymbol{h}^{(2)} \qquad \boldsymbol{h}^{(3)} \qquad \boldsymbol{h}^{(4)}$$

$$\boldsymbol{W} \qquad \boldsymbol{W} \qquad \boldsymbol{W}$$

$$\frac{\partial J^{(4)}}{\partial \boldsymbol{h}^{(1)}} = \boxed{\frac{\partial \boldsymbol{h}^{(2)}}{\partial \boldsymbol{h}^{(1)}}} \times \boxed{\frac{\partial \boldsymbol{h}^{(3)}}{\partial \boldsymbol{h}^{(2)}}} \times \boxed{\frac{\partial \boldsymbol{h}^{(4)}}{\partial \boldsymbol{h}^{(3)}}} \times \frac{\partial J^{(4)}}{\partial \boldsymbol{h}^{(4)}}$$

What happens if these are small?

**Vanishing gradient problem:**
When these are small, the gradient signal gets smaller and smaller as it backpropagates further

# Vanishing gradient proof sketch (linear case)

- Recall: $h^{(t)} = \sigma\left(W_h h^{(t-1)} + W_x x^{(t)} + b_1\right)$

- What if $\sigma$ were the identity function, $\sigma(x) = x$ ?

$$\frac{\partial h^{(t)}}{\partial h^{(t-1)}} = \text{diag}\left(\sigma'\left(W_h h^{(t-1)} + W_x x^{(t)} + b_1\right)\right) W_h \qquad \text{(chain rule)}$$

$$= I\ W_h = W_h$$

- Consider the gradient of the loss $J^{(i)}(\theta)$ on step $i$, with respect to the hidden state $h^{(j)}$ on some previous step $j$. Let $\ell = i - j$

$$\frac{\partial J^{(i)}(\theta)}{\partial h^{(j)}} = \frac{\partial J^{(i)}(\theta)}{\partial h^{(i)}} \prod_{j < t \leq i} \frac{\partial h^{(t)}}{\partial h^{(t-1)}} \qquad \text{(chain rule)}$$

$$= \frac{\partial J^{(i)}(\theta)}{\partial h^{(i)}} \prod_{j < t \leq i} W_h = \frac{\partial J^{(i)}(\theta)}{\partial h^{(i)}} \boxed{W_h^{\ell}} \qquad \left(\text{value of } \frac{\partial h^{(t)}}{\partial h^{(t-1)}}\right)$$

If $W_h$ is "small", then this term gets exponentially problematic as $\ell$ becomes large

**Source**: "On the difficulty of training recurrent neural networks", Pascanu et al, 2013. http://proceedings.mlr.press/v28/pascanu13.pdf
(and supplemental materials), at http://proceedings.mlr.press/v28/pascanu13-supp.pdf

# Vanishing gradient proof sketch (linear case)

- What's wrong with $W_h^\ell$ ?

sufficient but not necessary

- Consider if the eigenvalues of $W_h$ are all less than 1:

$$\lambda_1, \lambda_2, \ldots, \lambda_n < 1$$
$$q_1, q_2, \cdots, q_n \text{ (eigenvectors)}$$

- We can write $\frac{\partial J^{(i)}(\theta)}{\partial h^{(i)}} W_h^\ell$ using the eigenvectors of $W_h$ as a basis:

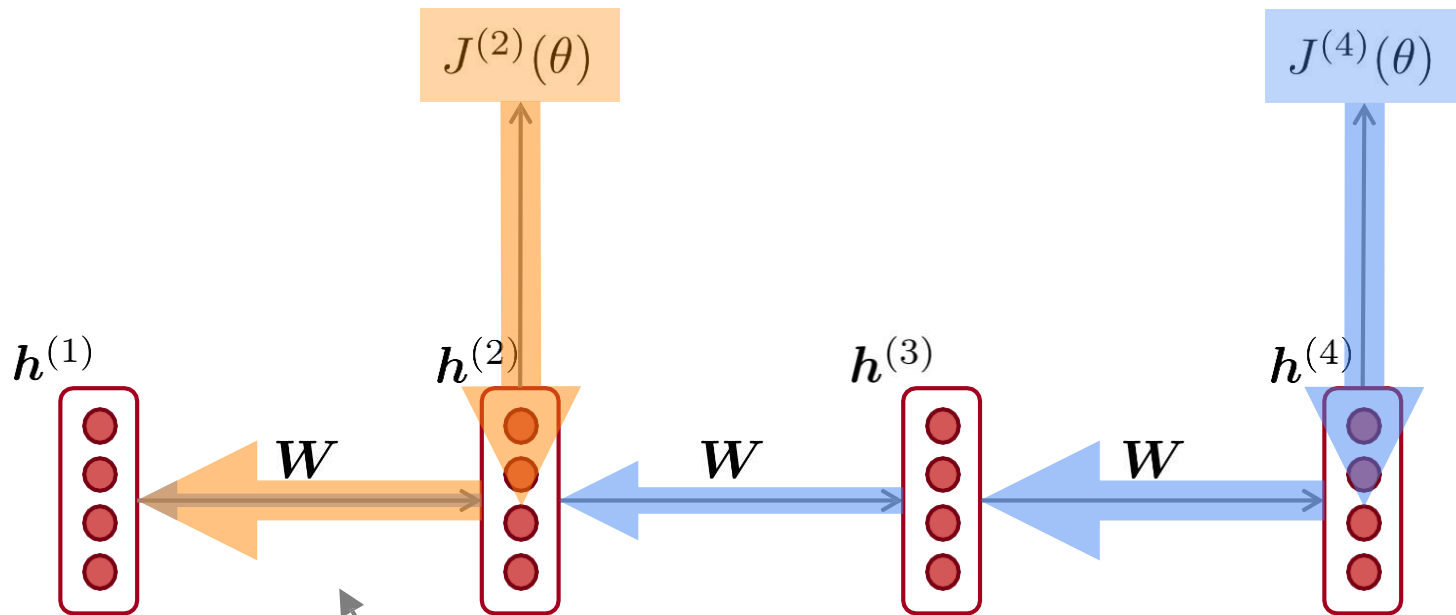$$\frac{\partial J^{(i)}(\theta)}{\partial h^{(i)}} W_h^\ell = \sum_{i=1}^n c_i \boxed{\lambda_i^\ell} q_i \approx 0 \text{ (for large } \ell)$$

Approaches 0 as $\ell$ grows, so gradient vanishes

- What about nonlinear activations $\sigma$ (i.e., what we use?)
  - Pretty much the same thing, except the proof requires $\lambda_i < \gamma$
    for some $\gamma$ dependent on dimensionality and $\sigma$

Source: "On the difficulty of training recurrent neural networks", Pascanu et al, 2013. http://proceedings.mlr.press/v28/pascanu13.pdf
(and supplemental materials), at http://proceedings.mlr.press/v28/pascanu13-supp.pdf

# Why is vanishing gradient a problem?

$J^{(2)}(\theta)$

$J^{(4)}(\theta)$

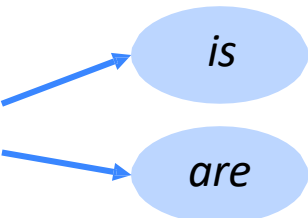$h^{(1)}$  $h^{(2)}$  $h^{(3)}$  $h^{(4)}$

$W$  $W$  $W$

Gradient signal from faraway is lost because it's much smaller than gradient signal from close-by.

So model weights are only updated only with respect to near effects, not long-term effects.

# Effect of vanishing gradient on RNN-LM

- **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 _____*

- To learn from this training example, the RNN-LM needs to model the dependency between *"tickets"* on the 7th step and the target word *"tickets"* at the end.

- But if gradient is small, the model can't learn this dependency
  - So the model is unable to predict similar long-distance dependencies at test time

# Effect of vanishing gradient on RNN-LM

- **LM task:** *The writer of the books* ____

  *is*

  *are*

- **Correct answer**: *The writer of the books is planning a sequel*

- **Syntactic recency:** *The writer of the books is*              (correct)

- **Sequential recency:** *The writer of the books are*              (incorrect)

- Due to vanishing gradient, RNN-LMs are better at learning from sequential recency than syntactic recency, so they make this type of error more often than we'd like [Linzen et al 2016]

"Assessing the Ability of LSTMs to Learn Syntax-Sensitive Dependencies", Linzen et al, 2016. https://arxiv.org/pdf/1611.01368.pdf

# Why is <u>exploding</u> gradient a problem?

- If the gradient becomes too big, then the SGD update step becomes too big:

$$\theta^{new} = \theta^{old} - \overbrace{\alpha}^{\text{learning rate}} \underbrace{\nabla_\theta J(\theta)}_{\text{gradient}}$$

- This can cause bad updates: we take too large a step and reach a bad parameter configuration (with large loss)

- In the worst case, this will result in Inf or NaN in your network (then you have to restart training from an earlier checkpoint)
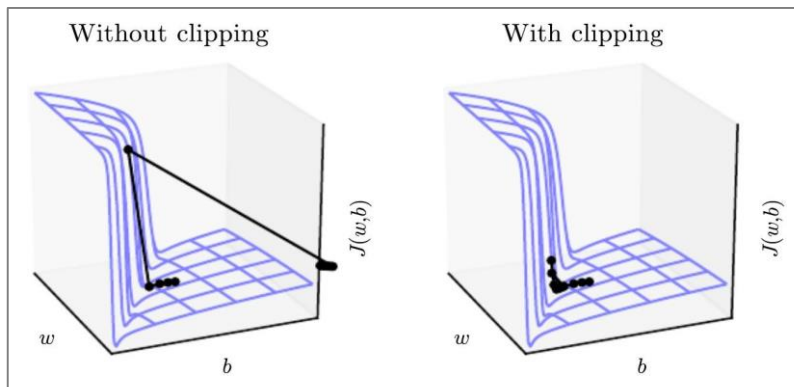
# Gradient clipping: solution for exploding gradient

- <u>Gradient clipping</u>: if the norm of the gradient is greater than some threshold, scale it down before applying SGD update

**Algorithm 1** Pseudo-code for norm clipping

$$\hat{\mathbf{g}} \leftarrow \frac{\partial \mathcal{E}}{\partial \theta}$$
**if** $\|\hat{\mathbf{g}}\| \geq threshold$ **then**
$$\hat{\mathbf{g}} \leftarrow \frac{threshold}{\|\hat{\mathbf{g}}\|} \hat{\mathbf{g}}$$
**end if**

- <u>Intuition</u>: take a step in the same direction, but a smaller step



**Source**: "On the difficulty of training recurrent neural networks", Pascanu et al, 2013. http://proceedings.mlr.press/v28/pascanu13.pdf

# How to fix vanishing gradient problem?

- The main problem is that *it's too difficult for the RNN to learn to preserve information over many timesteps*.

- In a vanilla RNN, the hidden state is constantly being rewritten

$$\boldsymbol{h}^{(t)} = \sigma \left( \boldsymbol{W}_h \boldsymbol{h}^{(t-1)} + \boldsymbol{W}_x \boldsymbol{x}^{(t)} + \boldsymbol{b} \right)$$

- How about a RNN with separate memory?

# Long Short-Term Memory (LSTM)

- A type of RNN proposed by Hochreiter and Schmidhuber in 1997 as a solution to the vanishing gradients problem.

- On step $t$, there is a hidden state $\boldsymbol{h}^{(t)}$ *and* a cell state $\boldsymbol{c}^{(t)}$
  - Both are vectors length $n$
  - The cell stores long-term information
  - The LSTM can erase, write and read information from the cell

- The selection of which information is erased/written/read is controlled by three corresponding gates
  - The gates are also vectors length $n$
  - On each timestep, each element of the gates can be open (1), closed (0), or somewhere in-between.
  - The gates are dynamic: their value is computed based on the current context

"Long short-term memory", Hochreiter and Schmidhuber, 1997. https://www.bioinf.jku.at/publications/older/2604.pdf

# Long Short-Term Memory (LSTM)

We have a sequence of inputs $x^{(t)}$, and we will compute a sequence of hidden states $h^{(t)}$ and cell states $c^{(t)}$. On timestep $t$:

**Sigmoid function**: all gate values are between 0 and 1

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

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

**Output gate:** controls what parts of cell are output to hidden state
**(READ)**

**New cell content:** this is the new content to be written to the cell

**Cell state**: erase ("forget") some content from last cell state, and write ("input") some new cell content

**Hidden state**: read ("output") some content from the cell

$$f^{(t)} = \sigma\left(W_f h^{(t-1)} + U_f x^{(t)} + b_f\right)$$

$$i^{(t)} = \sigma\left(W_i h^{(t-1)} + U_i x^{(t)} + b_i\right)$$

$$o^{(t)} = \sigma\left(W_o h^{(t-1)} + U_o x^{(t)} + b_o\right)$$

$$\tilde{c}^{(t)} = \tanh\left(W_c h^{(t-1)} + U_c x^{(t)} + b_c\right)$$

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

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

All these are vectors of same length $n$

Gates are applied using element-wise product

# Long Short-Term Memory (LSTM)

You can think of the LSTM equations visually like this:

**Source:** http://colah.github.io/posts/2015-08-Understanding-LSTMs/

# Long Short-Term Memory (LSTM)

You can think of the LSTM equations visually like this:

$$f^{(t)} = \sigma\left(W_f h^{(t-1)} + U_f x^{(t)} + b_f\right)$$

$$i^{(t)} = \sigma\left(W_i h^{(t-1)} + U_i x^{(t)} + b_i\right)$$

$$o^{(t)} = \sigma\left(W_o h^{(t-1)} + U_o x^{(t)} + b_o\right)$$

$$\tilde{c}^{(t)} = \tanh\left(W_c h^{(t-1)} + U_c x^{(t)} + b_c\right)$$

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

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

Write some new cell content

Forget some cell content

Output some cell content to the hidden state

Compute the forget gate

Compute the input gate

Compute the new cell content

Compute the output gate

$c_{t-1}$  $c_t$  $c_t$  $h_t$  $f_t$  $i_t$  $\tilde{c}_t$  $o_t$  $\sigma$  $\sigma$  tanh  $\sigma$  tanh  $h_{t-1}$  $h_t$  $x_t$

Neural Network Layer  Pointwise Operation  Vector Transfer  Concatenate  Copy

19

# How does LSTM solve vanishing gradients?

- The LSTM architecture makes it easier for the RNN to preserve information over many timesteps
    - e.g. if the forget gate is set to remember everything on every timestep, then the info in the cell is preserved indefinitely
    - By contrast, it's harder for vanilla RNN to learn a recurrent weight matrix $W_h$ that preserves info in hidden state

- LSTM doesn't *guarantee* that there is no vanishing/exploding gradient, but it does provide an easier way for the model to learn long-distance dependencies

# Gated Recurrent Units (GRU)

- Proposed by Cho et al. in 2014 as a simpler alternative to the LSTM.
- On each timestep $t$ we have input $\boldsymbol{x}^{(t)}$ and hidden state $\boldsymbol{h}^{(t)}$ (no cell state).

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

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

$$\boldsymbol{u}^{(t)} = \sigma\left(\boldsymbol{W}_u \boldsymbol{h}^{(t-1)} + \boldsymbol{U}_u \boldsymbol{x}^{(t)} + \boldsymbol{b}_u\right)$$

$$\boldsymbol{r}^{(t)} = \sigma\left(\boldsymbol{W}_r \boldsymbol{h}^{(t-1)} + \boldsymbol{U}_r \boldsymbol{x}^{(t)} + \boldsymbol{b}_r\right)$$

**New hidden state content:** reset gate selects useful parts of prev hidden state. Use this and current input to compute new hidden content.

$$\tilde{\boldsymbol{h}}^{(t)} = \tanh\left(\boldsymbol{W}_h(\boldsymbol{r}^{(t)} \circ \boldsymbol{h}^{(t-1)}) + \boldsymbol{U}_h \boldsymbol{x}^{(t)} + \boldsymbol{b}_h\right)$$

$$\boldsymbol{h}^{(t)} = (1 - \boldsymbol{u}^{(t)}) \circ \boldsymbol{h}^{(t-1)} + \boldsymbol{u}^{(t)} \circ \tilde{\boldsymbol{h}}^{(t)}$$

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

**How does this solve vanishing gradient?**
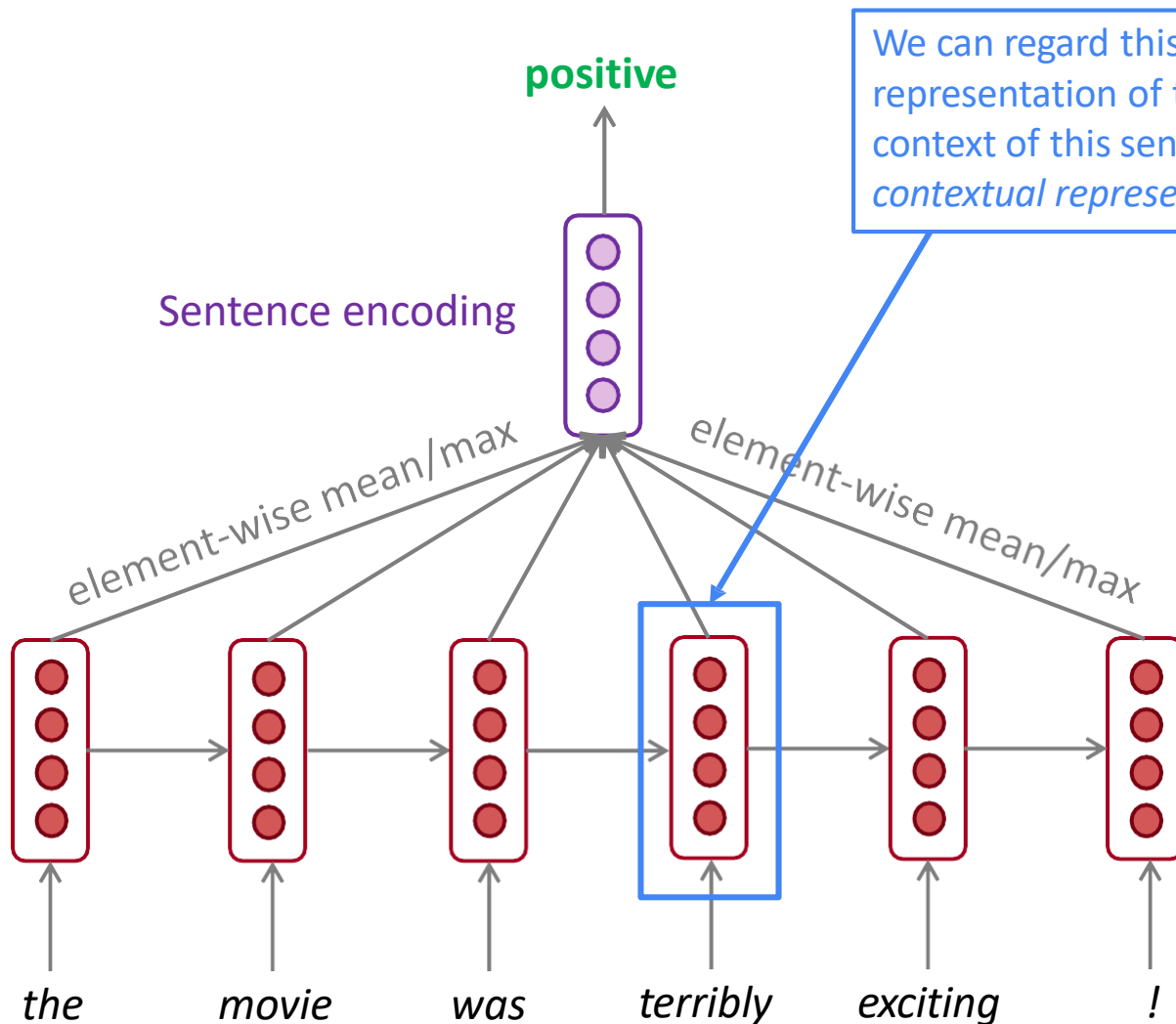Like LSTM, GRU makes it easier to retain info long-term (e.g. by setting update gate to 0)

"Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation", Cho et al. 2014, https://arxiv.org/pdf/1406.1078v3.pdf

# LSTM vs GRU

- Researchers have proposed many gated RNN variants, but LSTM and GRU are the most widely-used

- The biggest difference is that GRU is quicker to compute and has fewer parameters

- There is no conclusive evidence that one consistently performs better than the other

- LSTM is a good default choice (especially if your data has particularly long dependencies, or you have lots of training data)

- Rule of thumb: start with LSTM, but switch to GRU if you want something more efficient

# Bidirectional RNNs: motivation

Task: Sentiment Classification

positive

Sentence encoding

element-wise mean/max

element-wise mean/max

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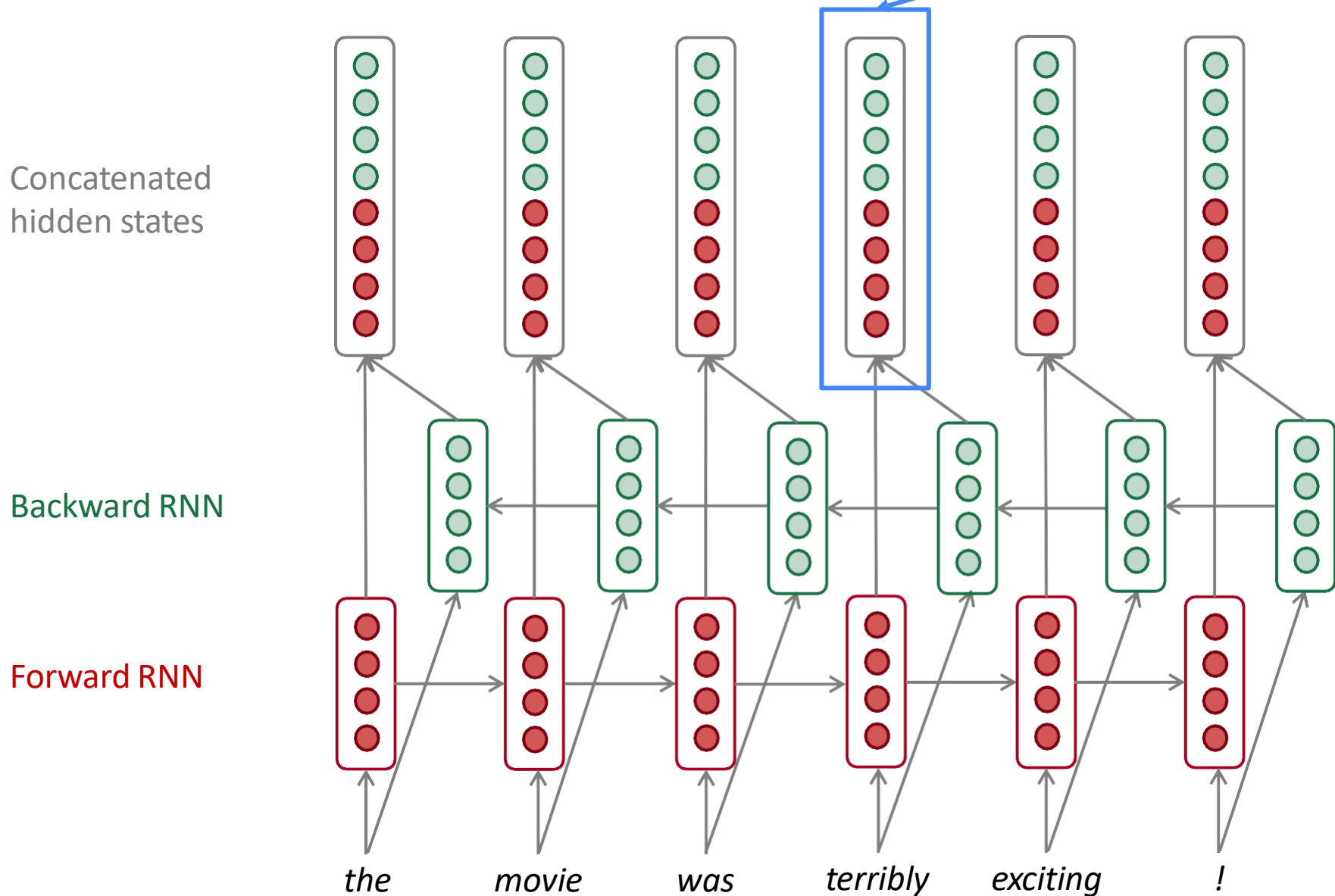
These contextual representations only contain information about the *left* context (e.g. *"the movie was"*).

**What about *right* context?**

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

the       movie       was       terribly       exciting       !

# Bidirectional RNNs

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

Concatenated hidden states

Backward RNN

Forward RNN

*the*  *movie*  *was*  *terribly*  *exciting*  *!*

# Bidirectional RNNs

On timestep $t$:

This is a general notation to mean "compute one forward step of the RNN" – it could be a vanilla, LSTM or GRU computation.

Forward RNN $\quad \overrightarrow{\boldsymbol{h}}^{(t)} = \boxed{\text{RNN}_{\text{FW}}}(\overrightarrow{\boldsymbol{h}}^{(t-1)}, \boldsymbol{x}^{(t)})$

Backward RNN $\quad \overleftarrow{\boldsymbol{h}}^{(t)} = \text{RNN}_{\text{BW}}(\overleftarrow{\boldsymbol{h}}^{(t+1)}, \boldsymbol{x}^{(t)})$

Generally, these two RNNs have separate weights

Concatenated hidden states $\quad \boxed{\boldsymbol{h}^{(t)}} = [\overrightarrow{\boldsymbol{h}}^{(t)}; \overleftarrow{\boldsymbol{h}}^{(t)}]$

We regard this as "the hidden state" of a bidirectional RNN. This is what we pass on to the next parts of the network.

# Bidirectional RNNs: simplified diagram



the    movie    was    terribly    exciting    !

The two-way arrows indicate bidirectionality and the depicted hidden states are assumed to be the concatenated forwards+backwards states.
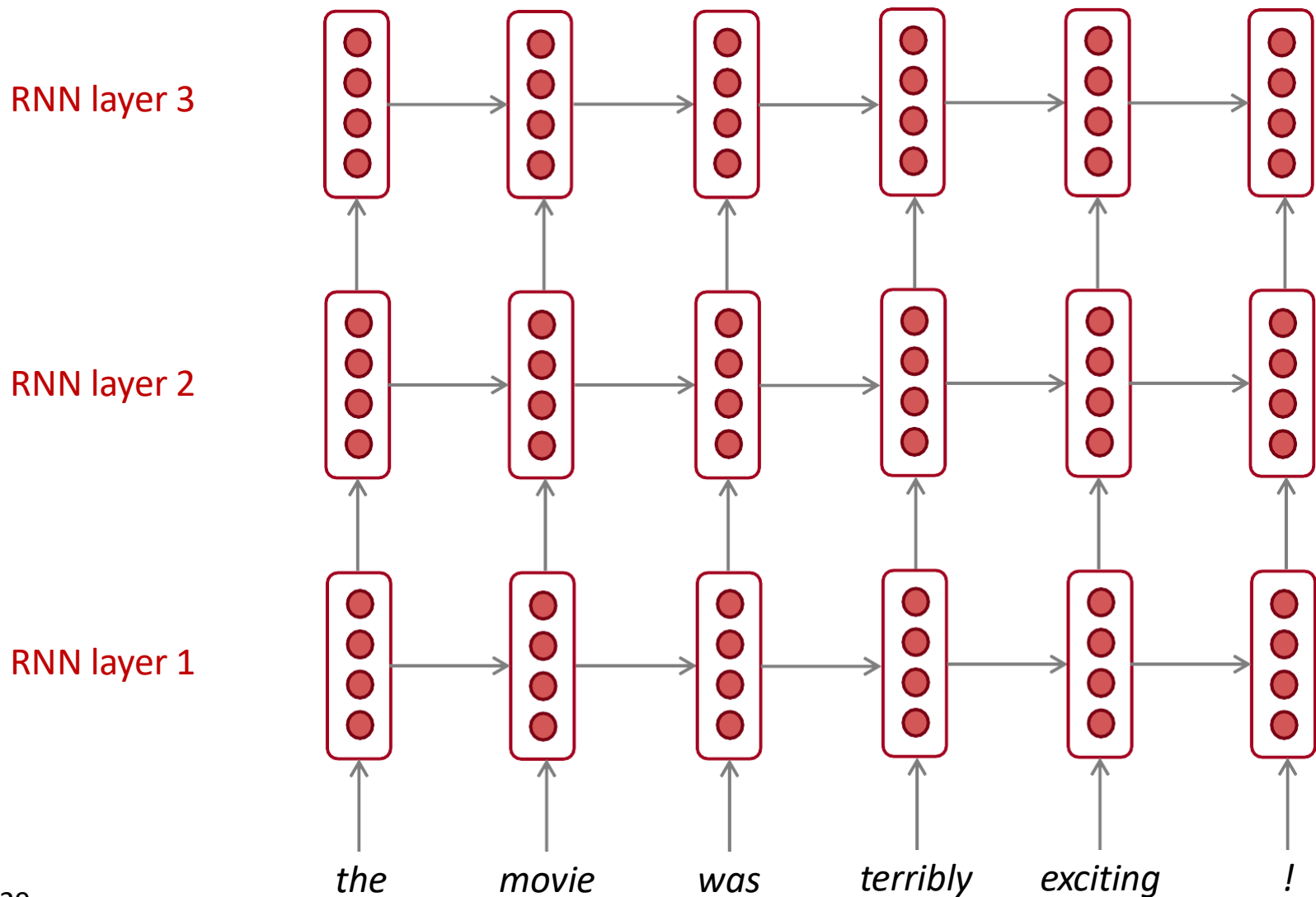
# Bidirectional RNNs

- Note: bidirectional RNNs are only applicable if you have access to the entire input sequence.

  - They are **not** applicable to Language Modeling, because in LM you *only* have left context available.

- If you do have entire input sequence (e.g. any kind of encoding), bidirectionality is powerful (you should use it by default).

- For example, BERT (**Bidirectional** Encoder Representations from Transformers) is a powerful pretrained contextual representation system built on bidirectionality.

  - You will learn more about BERT later in the course!

# Multi-layer RNNs

- RNNs are already "deep" on one dimension (they unroll over many timesteps)

- We can also make them "deep" in another dimension by applying multiple RNNs – this is a multi-layer RNN.

- This allows the network to compute more complex representations
    - The lower RNNs should compute lower-level features and the higher RNNs should compute higher-level features.

- Multi-layer RNNs are also called *stacked RNNs*.

# Multi-layer RNNs

RNN layer 3

RNN layer 2

RNN layer 1

the    movie    was    terribly    exciting    !

29

# Multi-layer RNNs in practice

- High-performing RNNs are often multi-layer (but aren't as deep as convolutional or feed-forward networks)

- For example: In a 2017 paper, Britz et al find that for Neural Machine Translation, 2 to 4 layers is best for the encoder RNN, and 4 layers is best for the decoder RNN
  - However, skip-connections/dense-connections are needed to train deeper RNNs (e.g. 8 layers)

- Transformer-based networks (e.g. BERT) can be up to 24 layers
  - You will learn about Transformers later; they have a lot of skipping-like connections

"Massive Exploration of Neural Machine Translation Architecutres", Britz et al, 2017. https://arxiv.org/pdf/1703.03906.pdf