

CS224n: NLP with Deep Learning

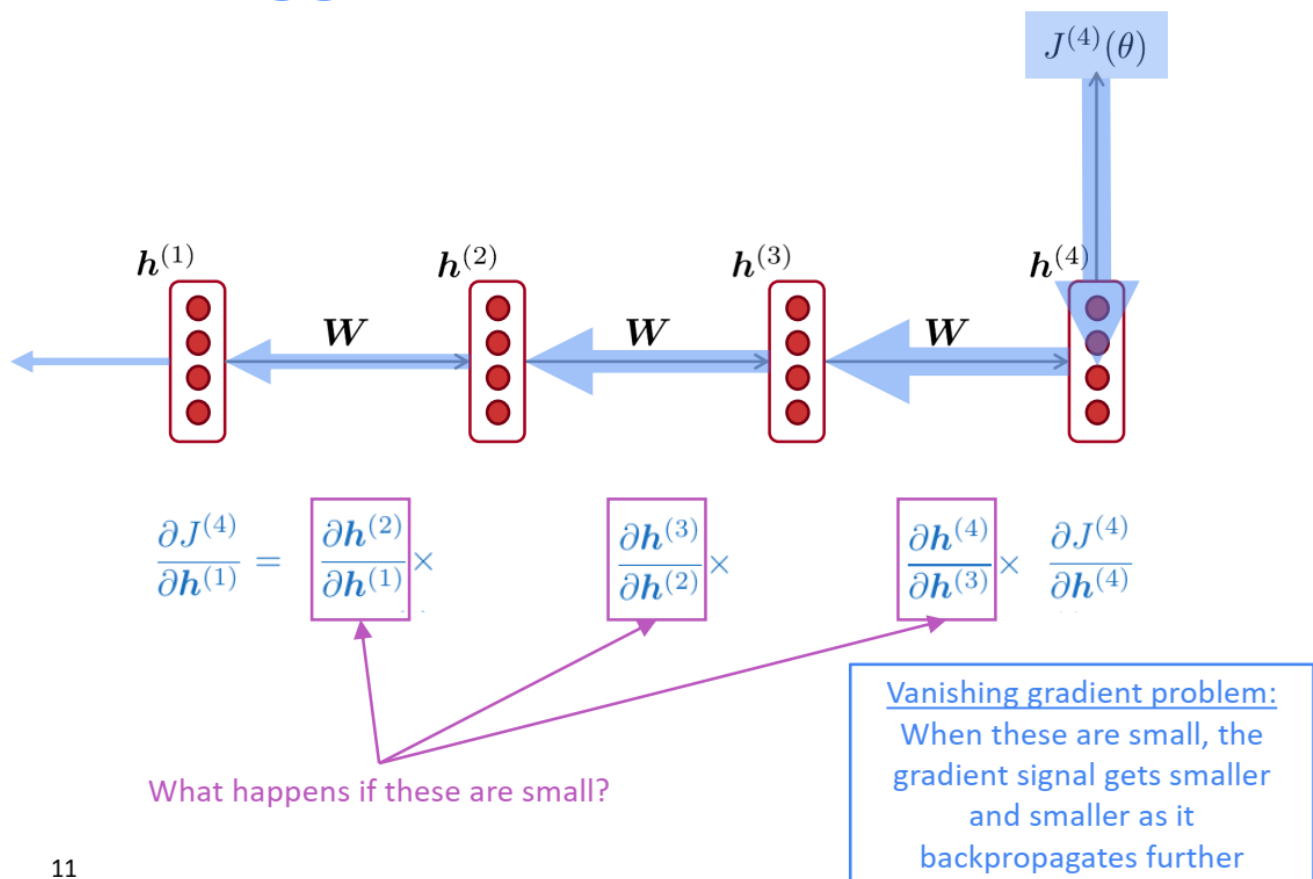
Lecture 7: Vanishing Gradients and Fancy RNN

Vanishing Gradient Problem

Intuition

- When the gradients $\frac{\partial h^{(i)}}{\partial h^{(i-1)}}$ are small, the overall gradient is going to get smaller and smaller, as we go backwards

Vanishing gradient intuition



11

- If the largest eigenvalue of W_h is < 1 :

the gradient will shrink exponentially

→ Vanishing gradient !

Proof Sketch

Vanishing gradient proof sketch

- Recall: $\mathbf{h}^{(t)} = \sigma \left(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_x \mathbf{x}^{(t)} + \mathbf{b}_1 \right)$
- Therefore: $\frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}} = \text{diag} \left(\sigma' \left(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_x \mathbf{x}^{(t)} + \mathbf{b}_1 \right) \right) \mathbf{W}_h$ (chain rule)
- Consider the gradient of the loss $J^{(i)}(\theta)$ on step i , with respect to the hidden state $\mathbf{h}^{(j)}$ on some previous step j .

$$\begin{aligned} \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(j)}} &= \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} \prod_{j < t \leq i} \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}} && \text{(chain rule)} \\ &= \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} \boxed{\mathbf{W}_h^{(i-j)}} \prod_{j < t \leq i} \text{diag} \left(\sigma' \left(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_x \mathbf{x}^{(t)} + \mathbf{b}_1 \right) \right) && \left(\text{value of } \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}} \right) \end{aligned}$$

If \mathbf{W}_h is small, then this term gets vanishingly small as i and j get further apart

Vanishing gradient proof sketch

- Consider matrix L2 norms:

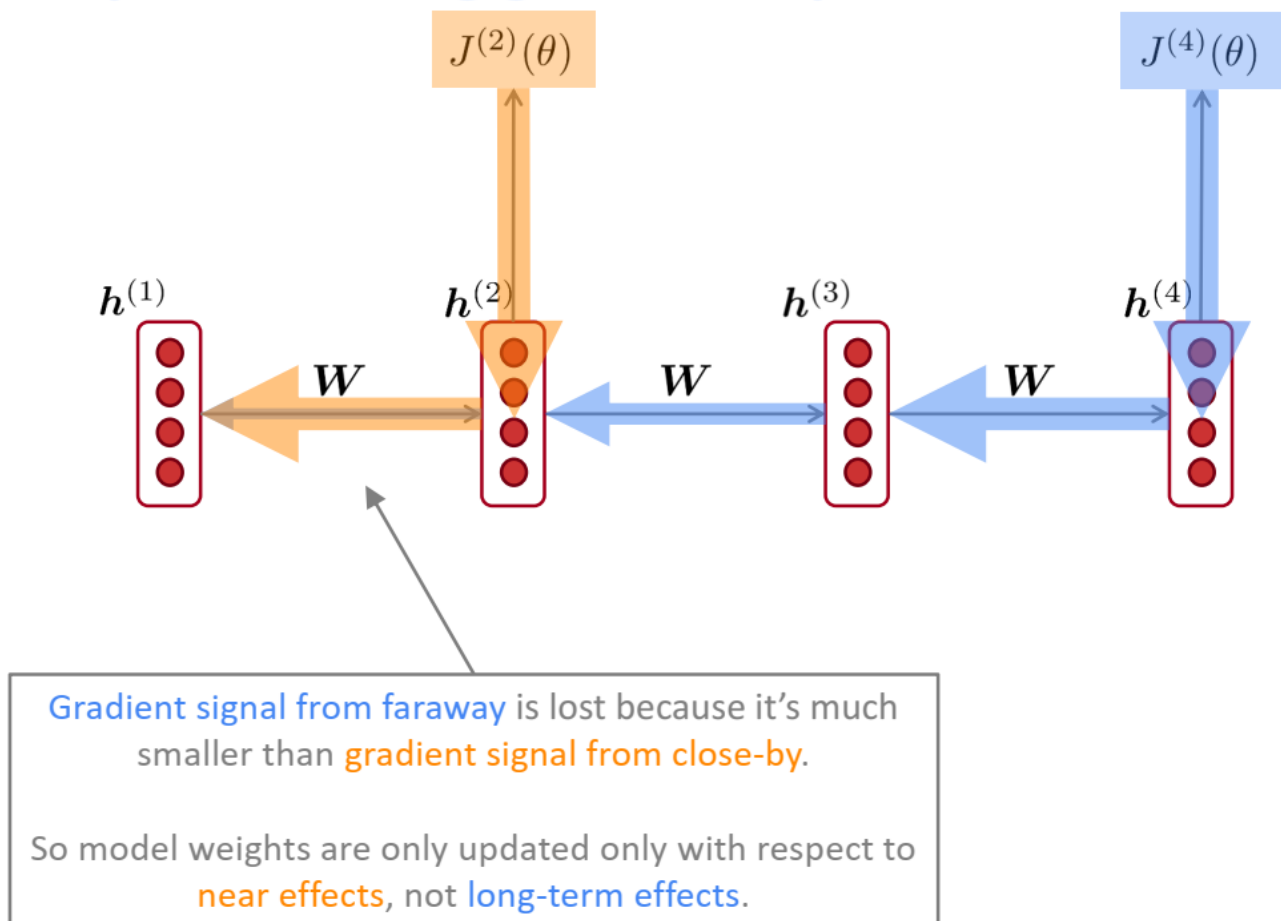
$$\left\| \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(j)}} \right\| \leq \left\| \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} \right\| \|\mathbf{W}_h\|^{(i-j)} \prod_{j < t \leq i} \left\| \text{diag} \left(\sigma' \left(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_x \mathbf{x}^{(t)} + \mathbf{b}_1 \right) \right) \right\|$$

- Pascanu et al showed that that if the **largest eigenvalue** of \mathbf{W}_h is **less than 1**, then the gradient $\left\| \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(j)}} \right\|$ will **shrink** exponentially
 - Here the bound is 1 because we have sigmoid nonlinearity
- There's a similar proof relating a **largest eigenvalue >1** to **exploding gradients**

But why is this a problem?

- The gradient from steps far away will be much smaller than gradients from close-by
- The weights will only be updated with respect to close effects, not long-term !

Why is vanishing gradient a problem?



14

Remark:

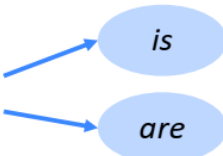


- We are studying $\frac{\partial J}{\partial h^{(i-1)}}$ because we have to calculate it to get $\frac{\partial J}{\partial W}$

Effect on RNN-LM

- Incapability of learning long-distance dependencies
- Syntactic recency (correct):
Making dependencies to words close **in their syntax**
- Sequential recency (incorrect):
Making dependencies to words close **spatially**

Unfortunately, RNN-LMs are better at learning sequential recency than syntactic recency:
(see example below)

Effect of vanishing gradient on RNN-LM

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

17

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

Exploding Gradient Problem

- If our gradient is too big, the update term for our parameter becomes too big
- Take update steps too large, and get stuck in a local optimum
- Get **Inf** or **Nan** values

Solution: Gradient clipping

- Scale down gradients that are higher than some threshold

LSTM

- Reason to exist = To solve the vanishing gradient problem
- A RNN with ability to learn / preserve information from many timesteps ago

What about a RNN with separate memory?

Structure

For each step t :

- a hidden state $h^{(t)}$: $(n, 1)$
- a cell state $c^{(t)}$: $(n, 1)$
Stores long-term information

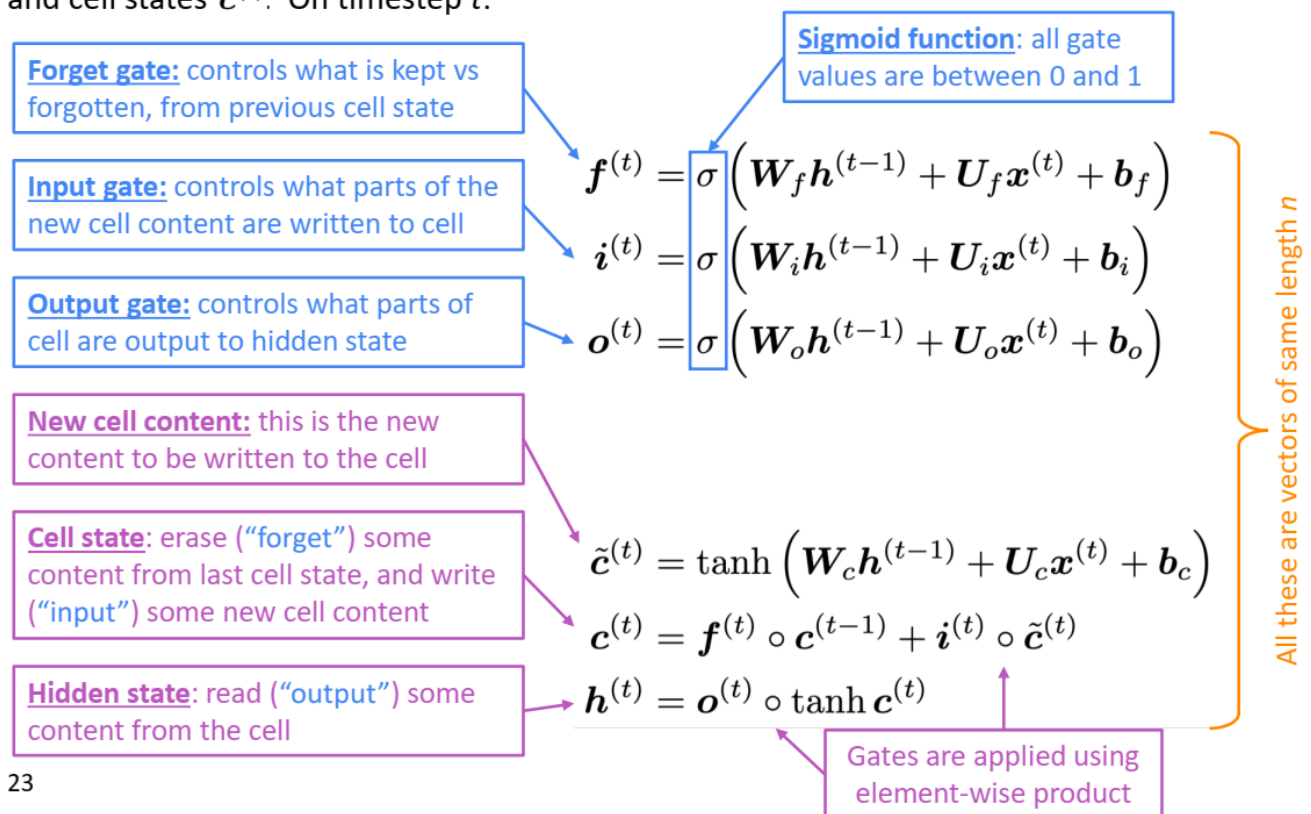
LSTM can erase, write and read information from cell !

This is done using gates of length n , whose values are in $[0, 1]$, thanks to our **sigmoid** friend:

- 1: open gate
- 0: closed gate

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 :



23

Cell state $c^{(t)}$ is the sum of:

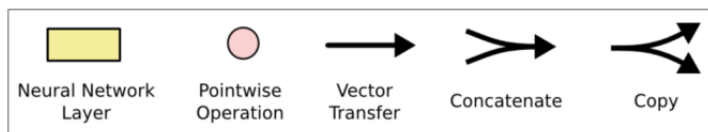
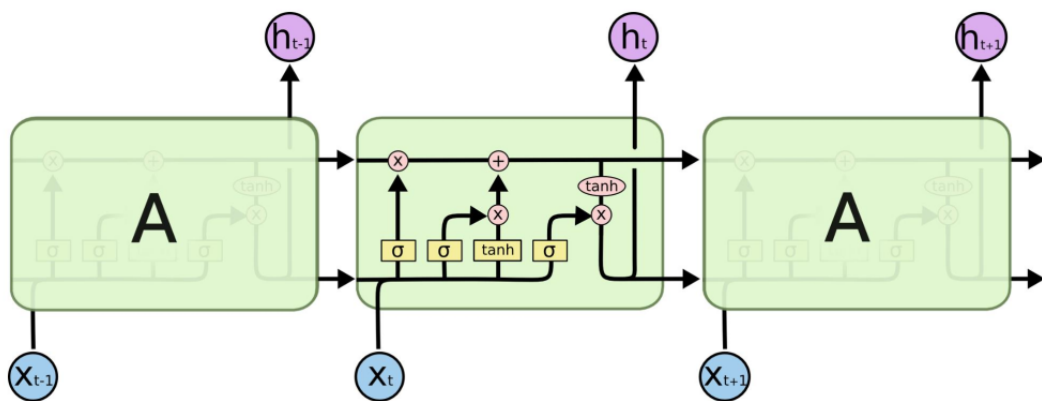
- previous cell state, masked by our forget gate
- our new cell, masked by our input gate

We can think of the hidden states as the output of the RNN

- These cell states ~ general memory, generally not accessible from the outside
- The hidden states will be passed on to the model

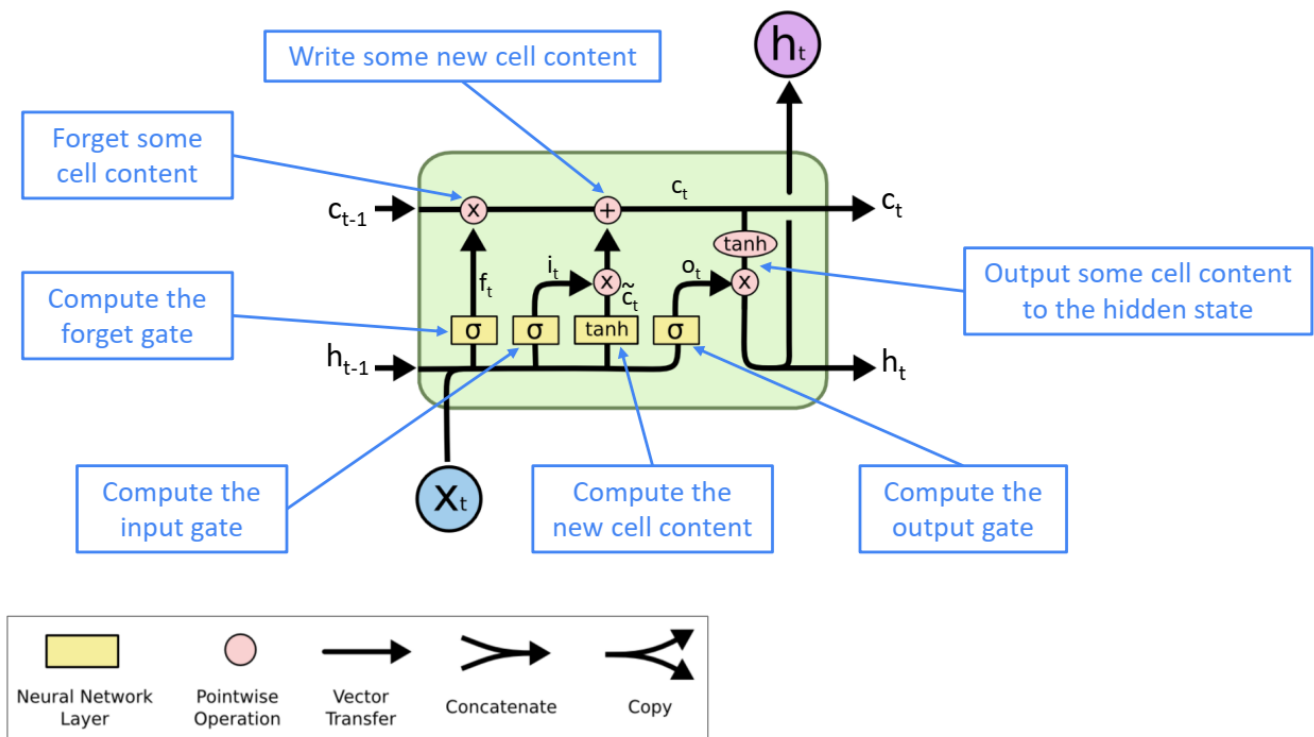
Long Short-Term Memory (LSTM)

You can think of the LSTM equations visually like this:



Long Short-Term Memory (LSTM)

You can think of the LSTM equations visually like this:



25

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

It is possible to set the forget gate so that the LSTM remembers everything from each time step

/!\ We could still have vanishing/exploding gradients, even when using LSTMs !

- With RNN, the hidden states were a bottleneck:
all gradients have to pass through them !!
- With LSTM, we can view the cell as a shortcut connection
There is a potential route within the cell that doesn't make the gradient vanish !

Applications

- Handwriting recognition
- Speech recognition
- Machine translation
- Parsing
- Image captioning

However, more recently, Transformers seem to have taken over !!



Gated Recurrent Units (GRU)

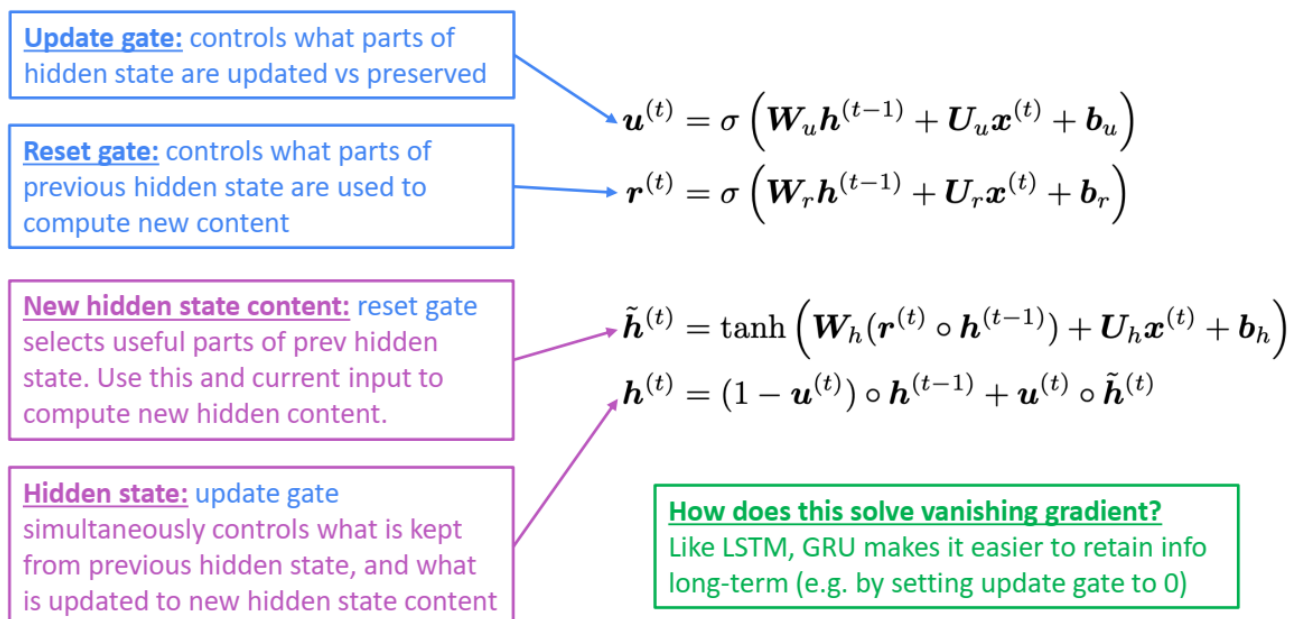
- Simpler alternative to LSTMs
- We have no cell state, only hidden states
- But we still have our cool gates!

Gates

- Update gate:
controls what is updated or preserved
~ input & forget gates for LSTM
- Reset gate:
selects which parts of the previous state are useful
- New hidden state = combination of previous hidden state, and computed new content

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 $\mathbf{x}^{(t)}$ and hidden state $\mathbf{h}^{(t)}$ (no cell state).



28

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

GRU compared to LSTM:

- Quicker
- Less parameters
- Similar performances

Other fixes

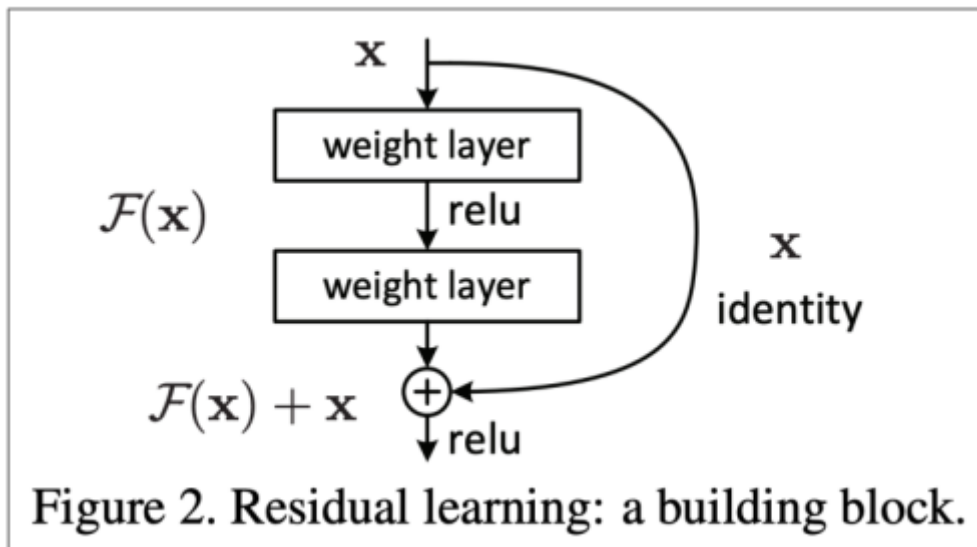
Gradient clipping

(See above)

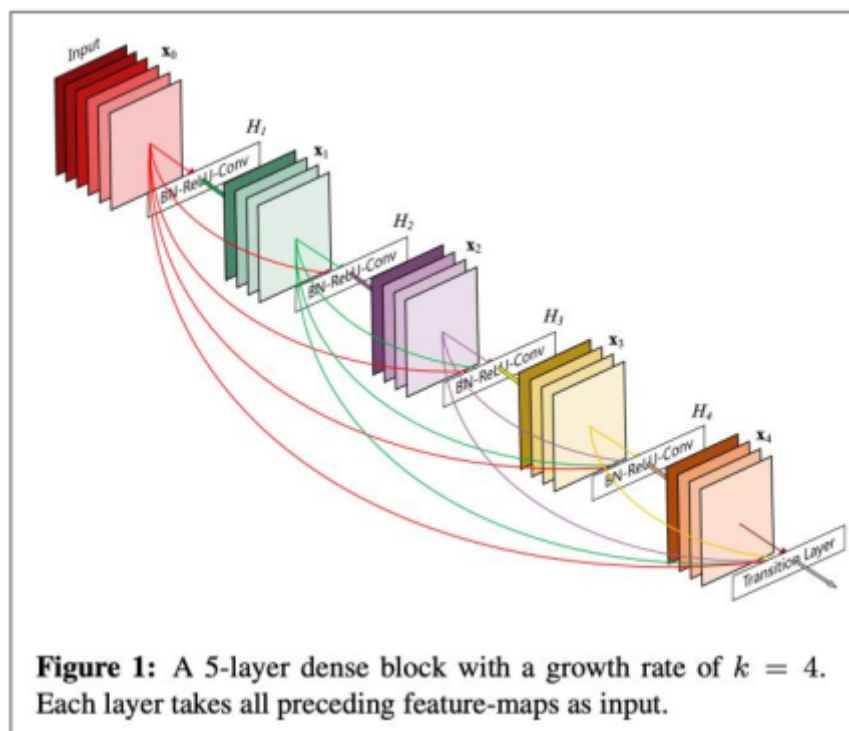
Skip connections

- Create direct connections from lower layers to newer layers
→ Allowing the gradients to flow more easily

Example: ResNet



- DenseNet:
Going even further: connect every layer to each other!



Recap 1

Recap

- Today we've learnt:
 - **Vanishing gradient problem**: what it is, why it happens, and why it's bad for RNNs
 - **LSTMs and GRUs**: more complicated RNNs that use gates to control information flow; they are more resilient to vanishing gradients
 - Remainder of this lecture:
 - **Bidirectional** RNNs
 - **Multi-layer** RNNs
- } Both of these are pretty simple
-

More fancy RNN variants

Bidirectional RNN

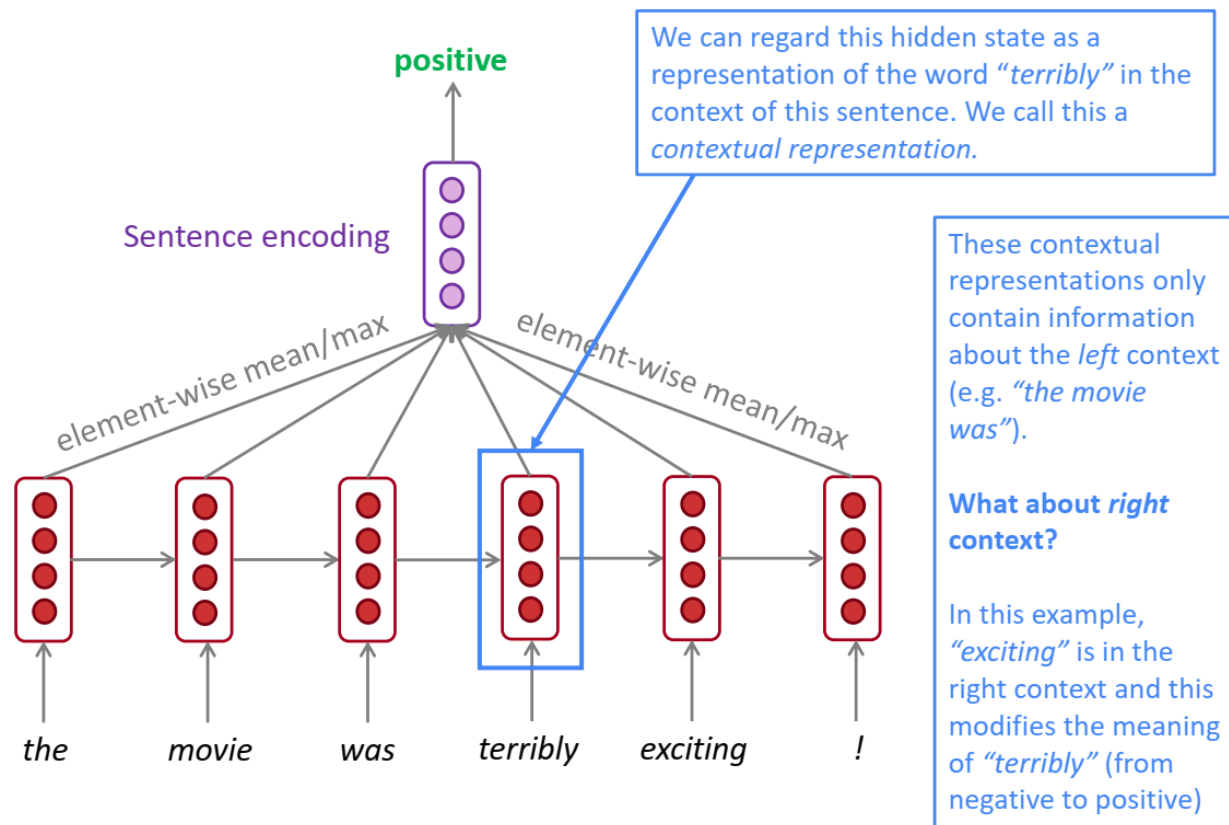
Motivation

- Our contextual representation only gets information from the left context !

We might want to get context coming from the right direction as well !

Bidirectional RNNs: motivation

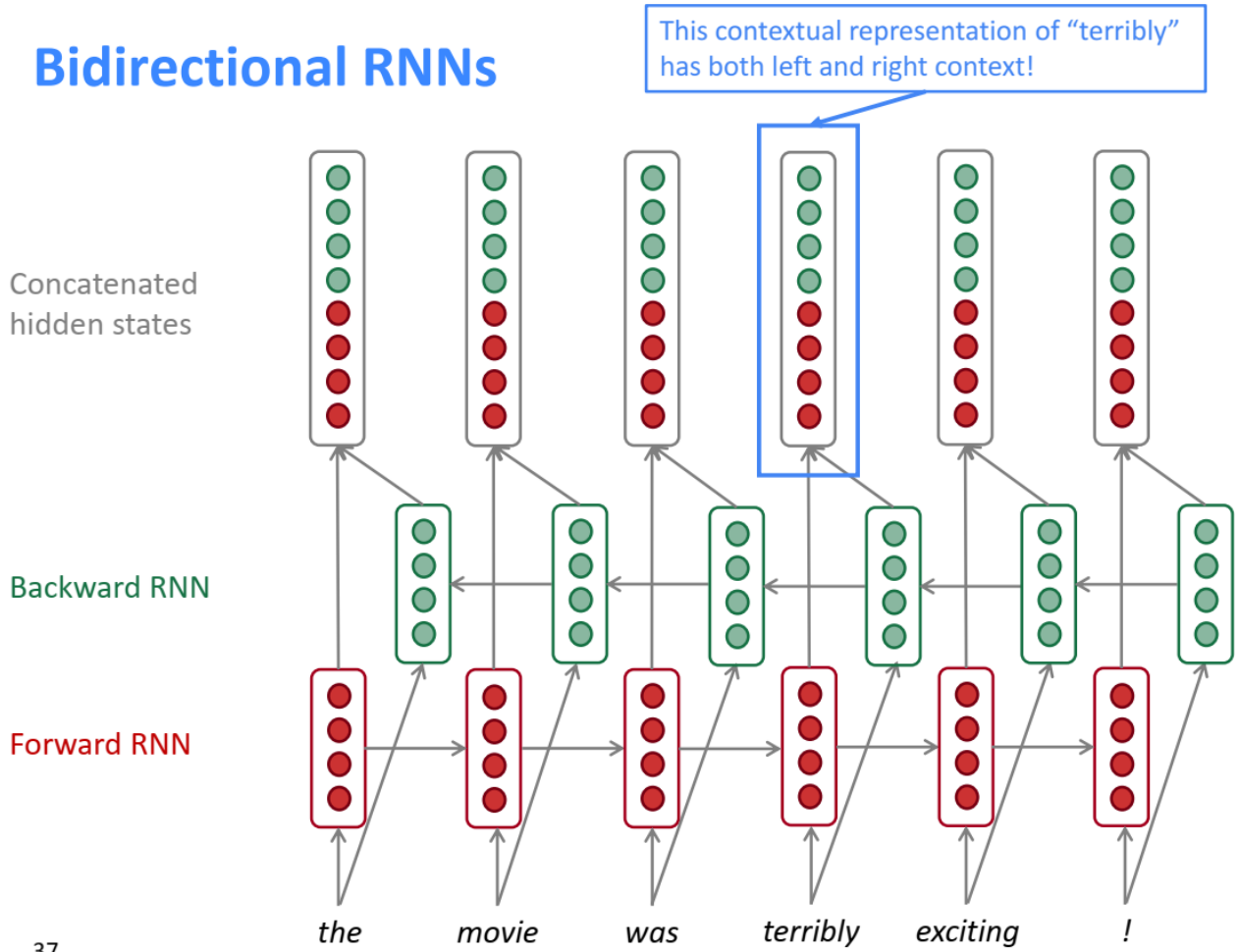
Task: Sentiment Classification



36

- We have 2 parallel RNN: one for each direction
- Hidden states = [Hidden state from left, Hidden state from right]

Bidirectional RNNs



Notations

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 $\vec{h}^{(t)} = \text{RNN}_{\text{FW}}(\vec{h}^{(t-1)}, x^{(t)})$

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

} Generally, these two RNNs have separate weights

Concatenated hidden states $h^{(t)} = [\vec{h}^{(t)}; \overleftarrow{h}^{(t)}]$

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

38

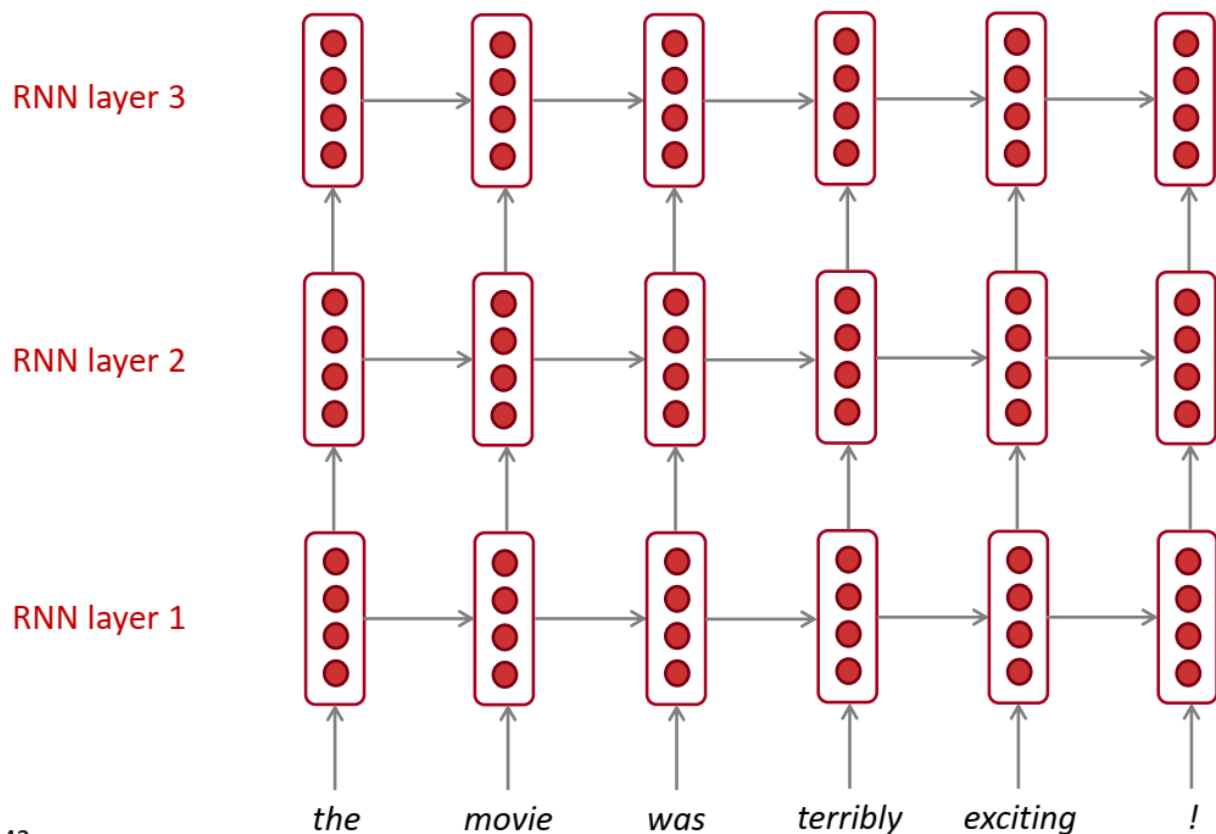
- The 2 RNNs are trained together, not separately
- Not useful for Language Modelling (as the Right context is missing, by definition!)

Multi-layer RNN

- Using multiple RNNs
- Allows to compute more complex representations
- Hidden states of layer i are the inputs to layer $i + 1$

Multi-layer RNNs

The hidden states from RNN layer i are the inputs to RNN layer $i+1$



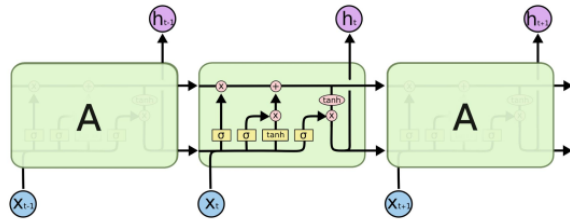
42

- RNNs have to be computed sequentially
→ Expensive to compute, can't go too deep

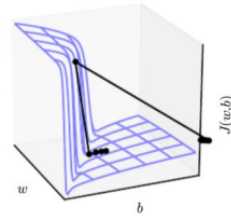
Recap 2: Practical Takeaways

In summary

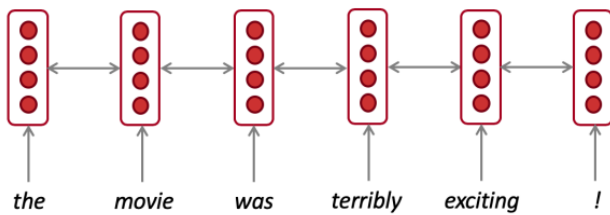
Lots of new information today! What are the **practical takeaways**?



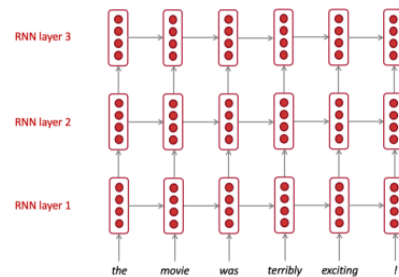
1. LSTMs are powerful but GRUs are faster



2. Clip your gradients



3. Use bidirectionality when possible



4. Multi-layer RNNs are powerful, but you might need skip/dense-connections if it's deep