



# Generating text with a RNN Language Model

Let's have some fun!

- You can train a RNN-LM on any kind of text, then generate text in that style.
- RNN-LM trained on *Harry Potter*:



“Sorry,” Harry shouted, panicking—“I’ll leave those brooms in London, are they?”

“No idea,” said Nearly Headless Nick, casting low close by Cedric, carrying the last bit of treacle Charms, from Harry’s shoulder, and to answer him the common room perched upon it, four arms held a shining knob from when the spider hadn’t felt it seemed. He reached the teams too.

Source: <https://medium.com/deep-writing/harry-potter-written-by-artificial-intelligence-8a9431803da6>



# Generating text with a RNN Language Model

Let's have some fun!

- You can train a RNN-LM on any kind of text, then generate text in that style.
- RNN-LM trained on **recipes**:

Title: CHOCOLATE RANCH BARBECUE

Categories: Game, Casseroles, Cookies, Cookies

Yield: 6 Servings

2 tb Parmesan cheese -- chopped

1 c Coconut milk

3 Eggs, beaten



Place each pasta over layers of lumps. Shape mixture into the moderate oven and simmer until firm. Serve hot in bodied fresh, mustard, orange and cheese.

Combine the cheese and salt together the dough in a large skillet; add the ingredients and stir in the chocolate and pepper.

Source: <https://gist.github.com/nylki/1efbaa36635956d35bcc>



# Generating text with a RNN Language Model

Let's have some fun!

- You can train a RNN-LM on any kind of text, then generate text in that style.
- RNN-LM trained on paint color names:

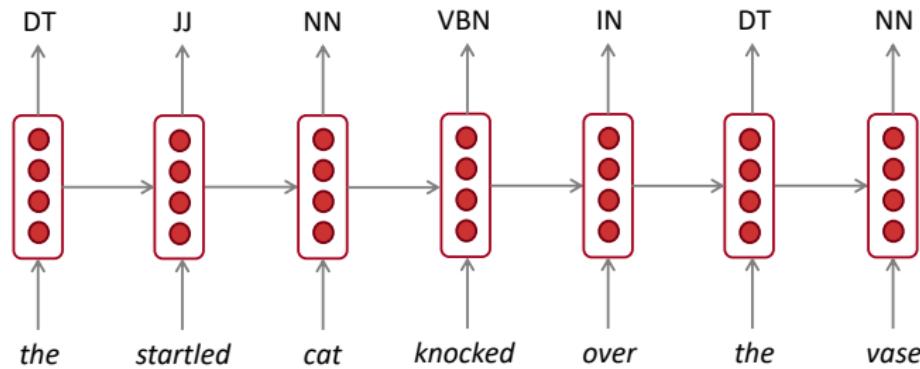
	Ghasty Pink 231 137 165
	Power Gray 151 124 112
	Navel Tan 199 173 140
	Bock Coe White 221 215 236
	Horble Gray 178 181 196
	Homestar Brown 133 104 85
	Snader Brown 144 106 74
	Golder Craam 237 217 177
	Hurky White 232 223 215
	Burf Pink 223 173 179
	Rose Hork 230 215 198
	Sand Dan 201 172 143
	Grade Bat 48 94 83
	Light Of Blast 175 150 147
	Grass Bat 176 99 108
	Sindis Poop 204 205 194
	Dope 219 209 179
	Testing 156 101 106
	Stoner Blue 152 165 159
	Burble Simp 226 181 132
	Stanky Bean 197 162 171
	Turdly 190 164 116

This is an example of a character-level RNN-LM (predicts what character comes next)



# RNNs can be used for tagging

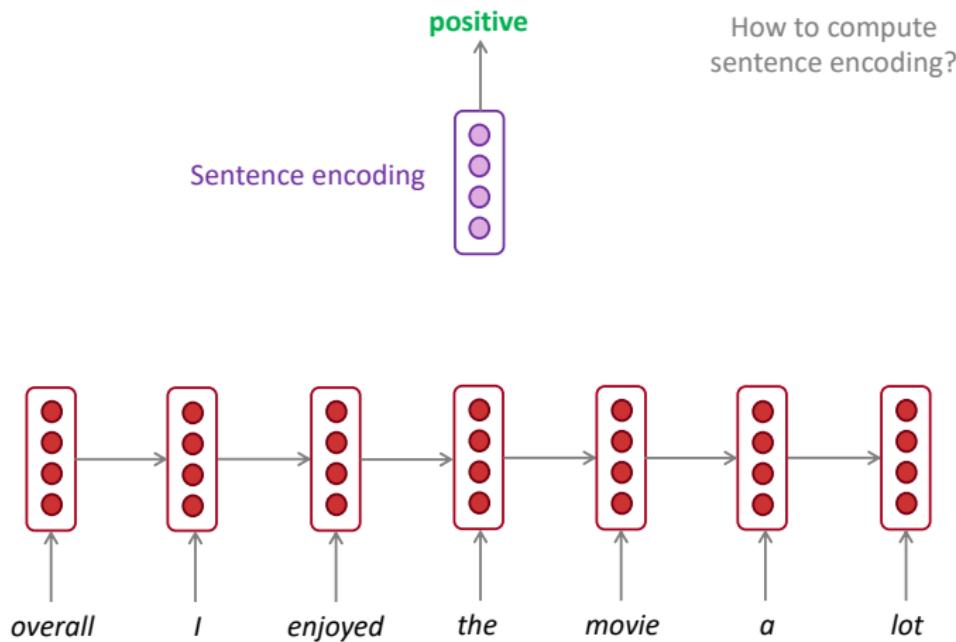
e.g. part-of-speech tagging, named entity recognition





# RNNs can be used for sentence classification

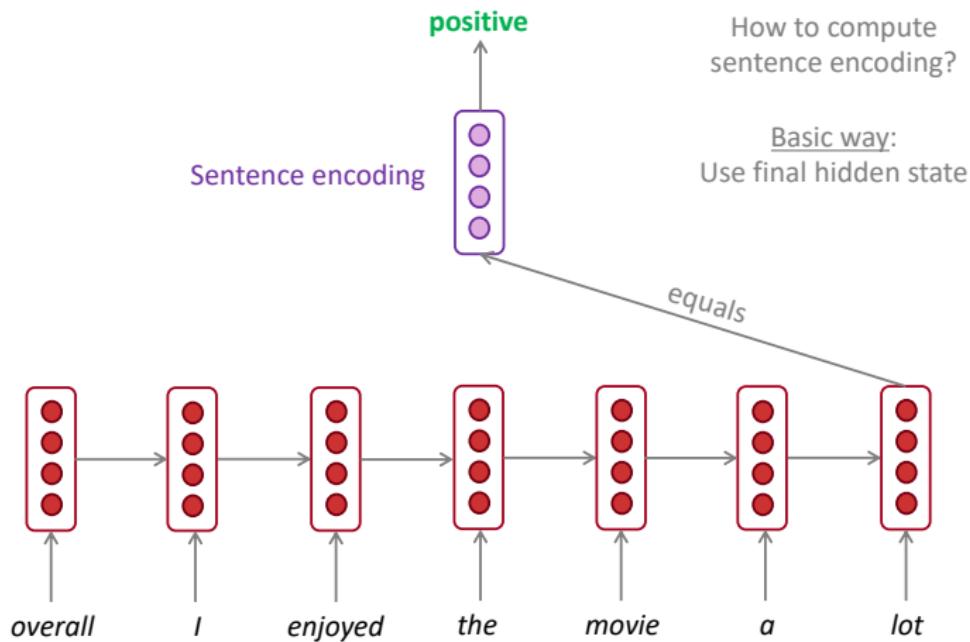
e.g. sentiment classification





# RNNs can be used for sentence classification

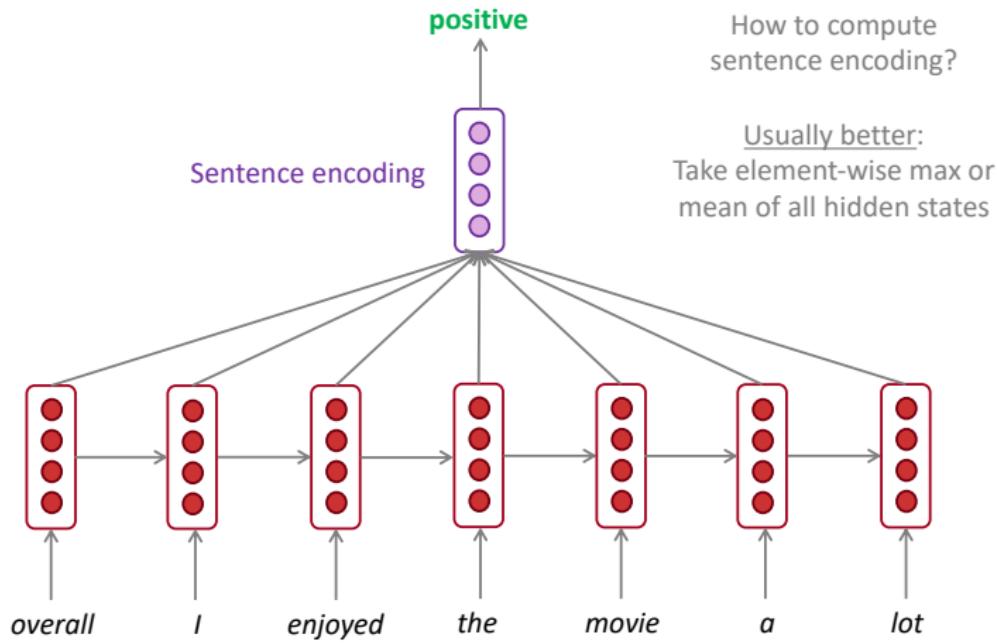
e.g. sentiment classification





# RNNs can be used for sentence classification

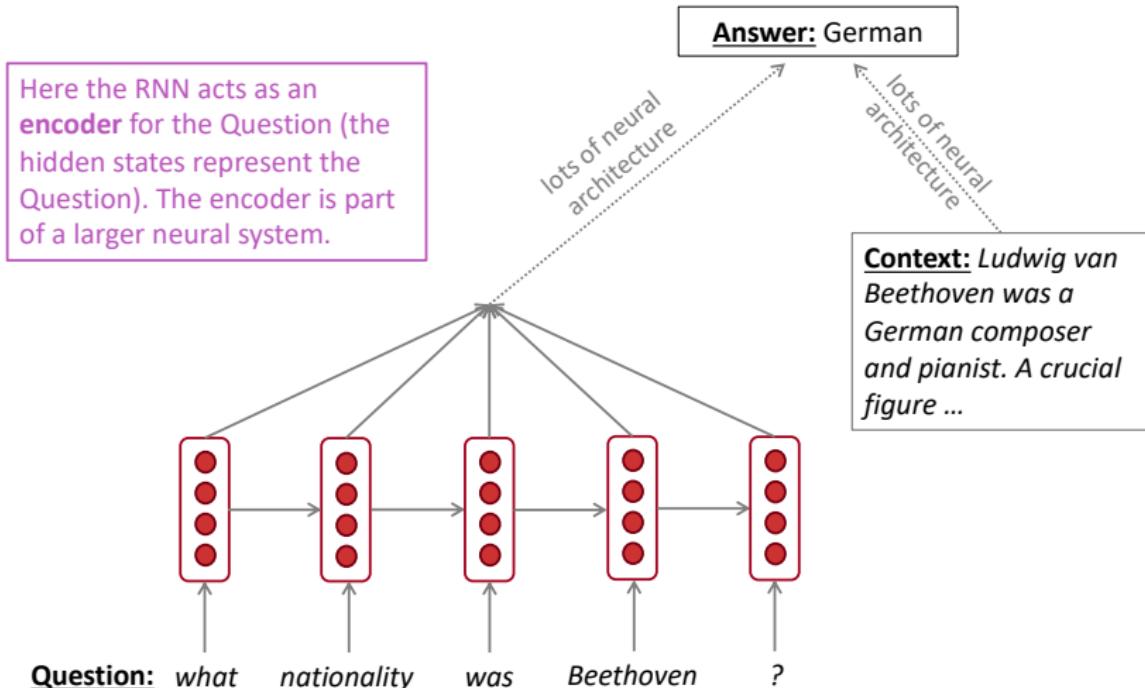
e.g. sentiment classification





# RNNs can be used as an encoder module

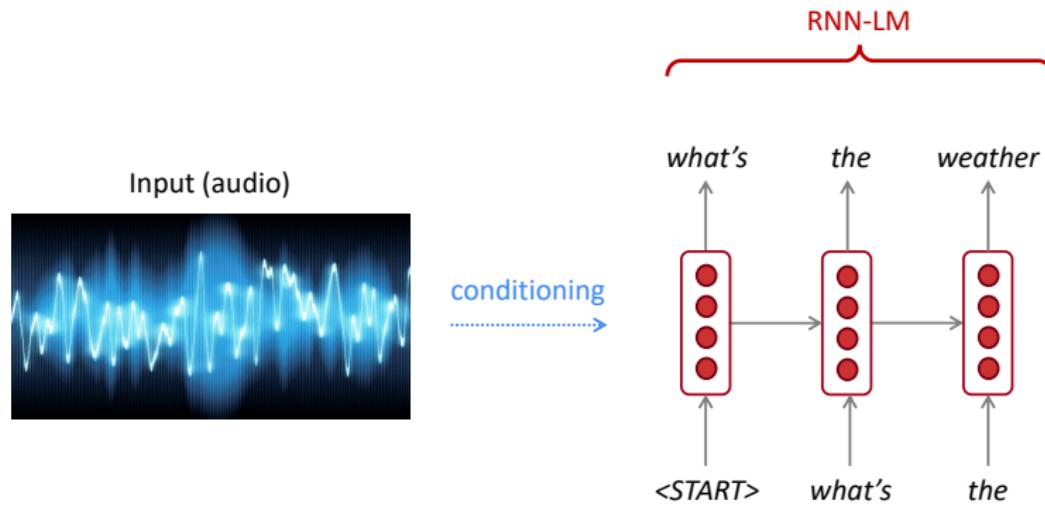
e.g. question answering, machine translation, *many other tasks!*





# RNN-LMs can be used to generate text

e.g. speech recognition, machine translation, summarization



This is an example of a *conditional language model*.

We'll see Machine Translation in much more detail later.



## A note on terminology

The RNN described in this lecture = simple/vanilla/Elman RNN



**Next lecture:** You will learn about other RNN flavors

like **GRU**



and **LSTM**



and multi-layer RNNs



**By the end of the course:** You will understand phrases like  
“stacked bidirectional LSTM with residual connections and self-attention”





# Content

3

## Language Models Based on Recurrent Neural Networks

- Recurrent Neural Networks (RNNs)
- Training a neural language model
- Other applications of neural language models
- **Evaluation of language models**
- Gradient vanishing and exploding
- Long Short Term Memory (LSTM)
- Gated Recurrent Units (GRUs)
- Vanishing and exploding gradient for other neural networks
- Bi-directional RNNs and multi-layer RNNs



## Evaluating Language Models

- The standard **evaluation metric** for Language Models is **perplexity**.

$$\text{perplexity} = \prod_{t=1}^T \left( \frac{1}{P_{\text{LM}}(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)})} \right)^{1/T}$$

Inverse probability of corpus, according to Language Model

Normalized by  
number of words

- This is equal to the exponential of the cross-entropy loss  $J(\theta)$ :

$$= \prod_{t=1}^T \left( \frac{1}{\hat{\mathbf{y}}_{\mathbf{x}_{t+1}}^{(t)}} \right)^{1/T} = \exp \left( \frac{1}{T} \sum_{t=1}^T -\log \hat{\mathbf{y}}_{\mathbf{x}_{t+1}}^{(t)} \right) = \exp(J(\theta))$$

**Lower perplexity is better!**



# RNNs have greatly improved perplexity

*n*-gram model →  
↓  
Increasingly complex RNNs

Model	Perplexity
Interpolated Kneser-Ney 5-gram (Chelba et al., 2013)	67.6
RNN-1024 + MaxEnt 9-gram (Chelba et al., 2013)	51.3
RNN-2048 + BlackOut sampling (Ji et al., 2015)	68.3
Sparse Non-negative Matrix factorization (Shazeer et al., 2015)	52.9
LSTM-2048 (Jozefowicz et al., 2016)	43.7
2-layer LSTM-8192 (Jozefowicz et al., 2016)	30
Ours small (LSTM-2048)	43.9
Ours large (2-layer LSTM-2048)	39.8

Perplexity improves  
(lower is better) ↓

Source: <https://research.fb.com/building-an-efficient-neural-language-model-over-a-billion-words/>



# Content

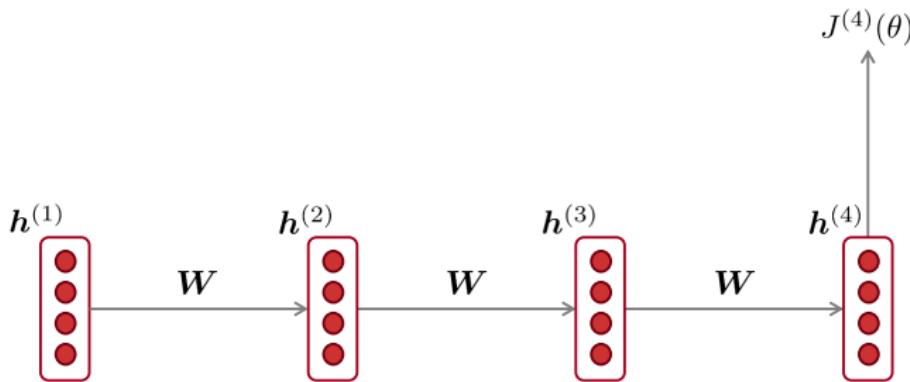
3

## Language Models Based on Recurrent Neural Networks

- Recurrent Neural Networks (RNNs)
- Training a neural language model
- Other applications of neural language models
- Evaluation of language models
- **Gradient vanishing and exploding**
- Long Short Term Memory (LSTM)
- Gated Recurrent Units (GRUs)
- Vanishing and exploding gradient for other neural networks
- Bi-directional RNNs and multi-layer RNNs

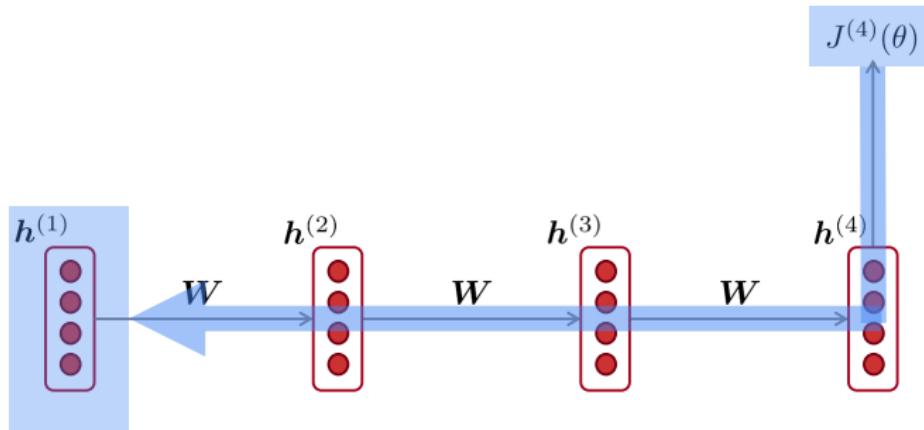


## Vanishing gradient intuition





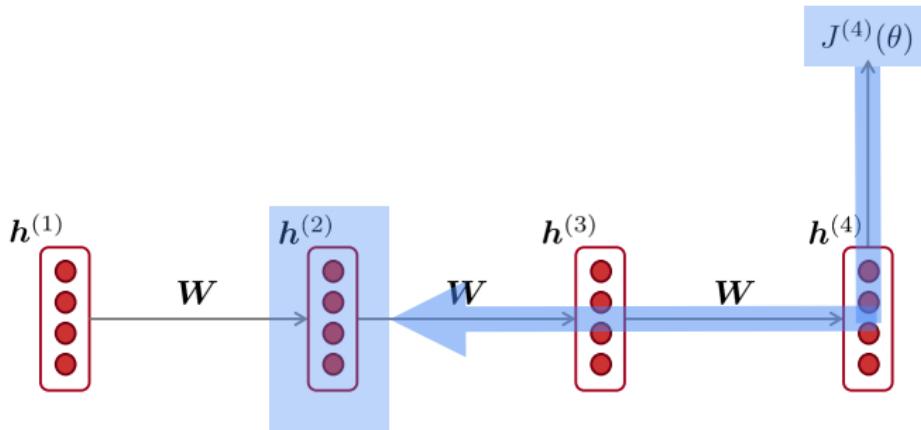
## Vanishing gradient intuition



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



## Vanishing gradient intuition

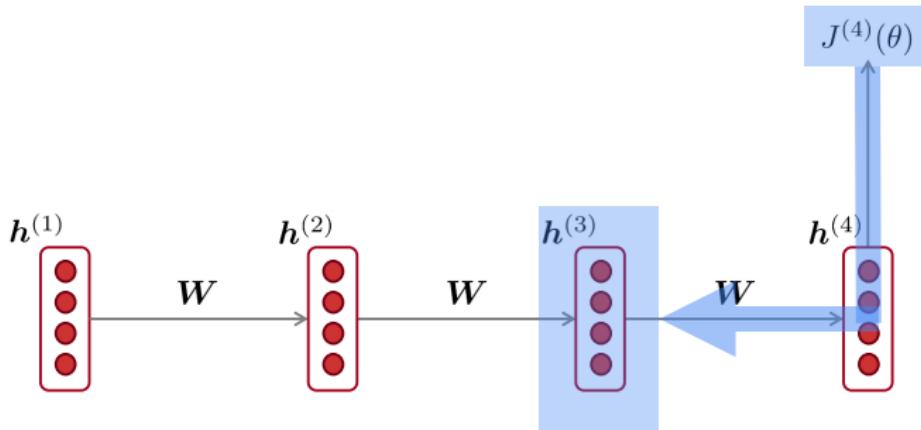


$$\frac{\partial J^{(4)}}{\partial \mathbf{h}^{(1)}} = \frac{\partial \mathbf{h}^{(2)}}{\partial \mathbf{h}^{(1)}} \times \frac{\partial J^{(4)}}{\partial \mathbf{h}^{(2)}}$$

chain rule!



# Vanishing gradient intuition



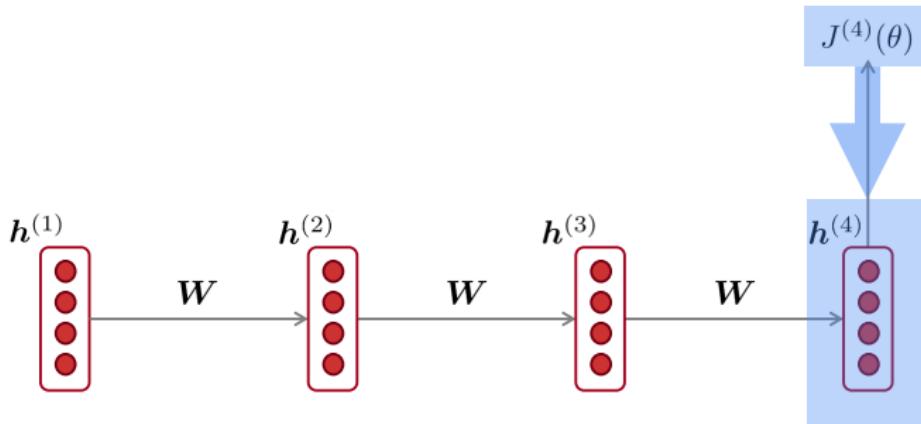
$$\frac{\partial J^{(4)}}{\partial \mathbf{h}^{(1)}} = \frac{\partial \mathbf{h}^{(2)}}{\partial \mathbf{h}^{(1)}} \times$$

$$\frac{\partial \mathbf{h}^{(3)}}{\partial \mathbf{h}^{(2)}} \times \frac{\partial J^{(4)}}{\partial \mathbf{h}^{(3)}}$$

chain rule!



## Vanishing gradient intuition



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

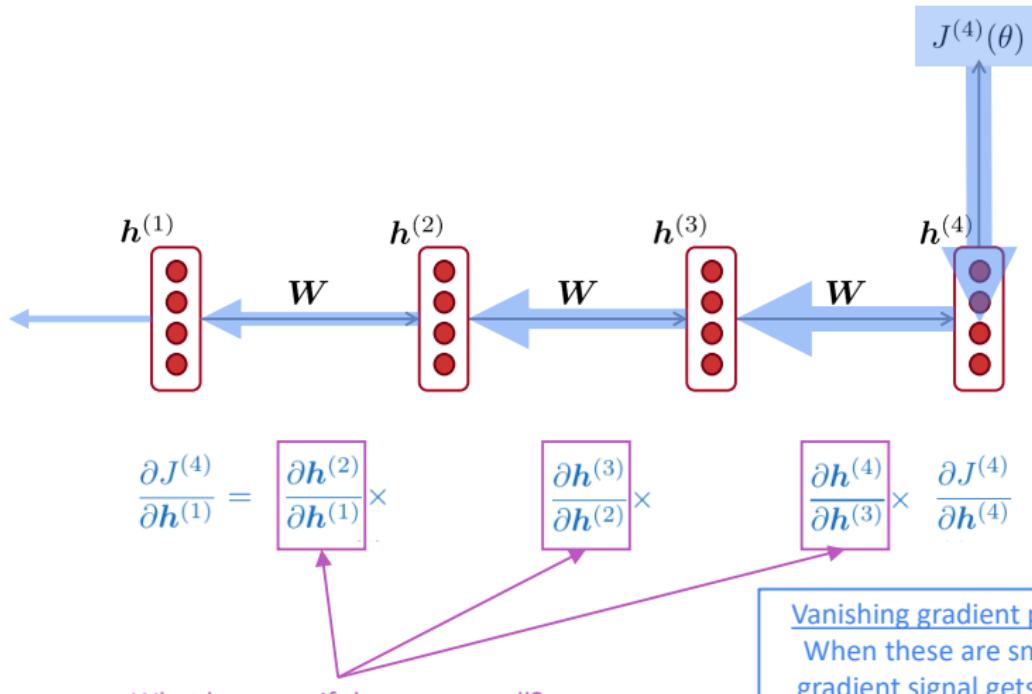
$$\frac{\partial \mathbf{h}^{(3)}}{\partial \mathbf{h}^{(2)}} \times$$

$$\frac{\partial \mathbf{h}^{(4)}}{\partial \mathbf{h}^{(3)}} \times \frac{\partial J^{(4)}}{\partial \mathbf{h}^{(4)}}$$

chain rule!



## Vanishing gradient intuition



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



## Vanishing gradient proof sketch

- Recall:  $\mathbf{h}^{(t)} = \sigma(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_x \mathbf{x}^{(t)} + \mathbf{b}_1)$
- 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  $j$ .

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

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:

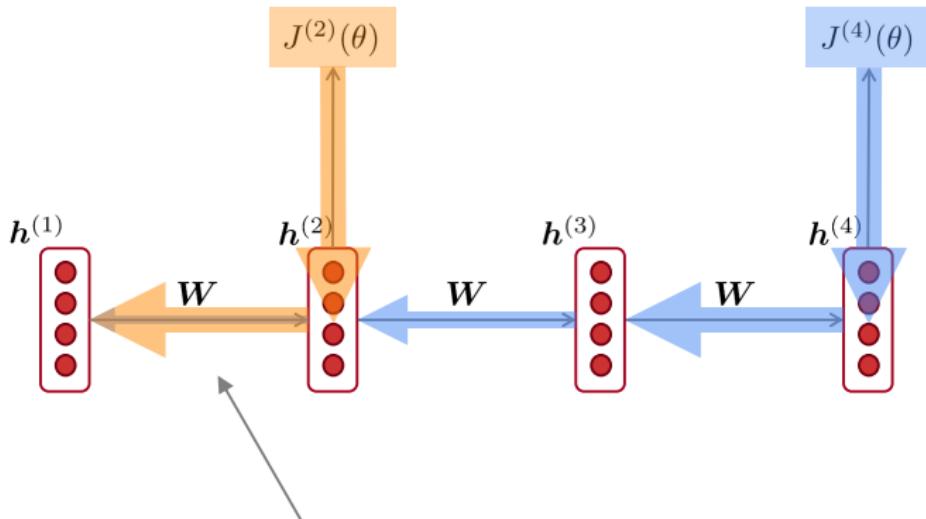
$$\left\| \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(j)}} \right\| \leq \left\| \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} \right\| \|\mathbf{W}_h\|^{(i-j)} \prod_{j < t \leq i} \left\| \text{diag} \left( \sigma' \left( \mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_x \mathbf{x}^{(t)}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

- Pascanu et al showed that 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





# Why is vanishing gradient a problem?



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.



## Why is vanishing gradient a problem?

- Another explanation: Gradient can be viewed as a measure of *the effect of the past on the future*
- If the gradient becomes vanishingly small over longer distances (step  $t$  to step  $t+n$ ), then we can't tell whether:
  1. There's **no dependency** between step  $t$  and  $t+n$  in the data
  2. We have **wrong parameters** to capture the true dependency between  $t$  and  $t+n$



## 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 7<sup>th</sup> 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* \_\_
- 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]





## Why is exploding gradient a problem?

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

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

learning rate

- 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

- Gradient clipping: 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

---

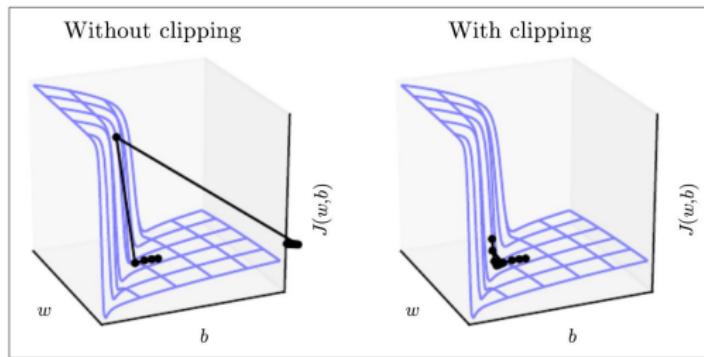
```
hat{g} ← ∂E/∂θ
if ||hat{g}|| ≥ threshold then
    hat{g} ← threshold / ||hat{g}|| * hat{g}
end if
```

---

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



# Gradient clipping: solution for exploding gradient



- This shows the loss surface of a simple RNN (hidden state is a scalar not a vector)
- The “cliff” is dangerous because it has steep gradient
- On the left, gradient descent takes two very big steps due to steep gradient, resulting in climbing the cliff then shooting off to the right (both bad updates)
- On the right, gradient clipping reduces the size of those steps, so effect is less drastic



## 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*

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

- How about a RNN with separate *memory*?



# Content

3

## Language Models Based on Recurrent Neural Networks

- Recurrent Neural Networks (RNNs)
- Training a neural language model
- Other applications of neural language models
- Evaluation of language models
- Gradient vanishing and exploding
- **Long Short Term Memory (LSTM)**
- Gated Recurrent Units (GRUs)
- Vanishing and exploding gradient for other neural networks
- Bi-directional RNNs and multi-layer RNNs



# Long Short-Term Memory (LSTM)

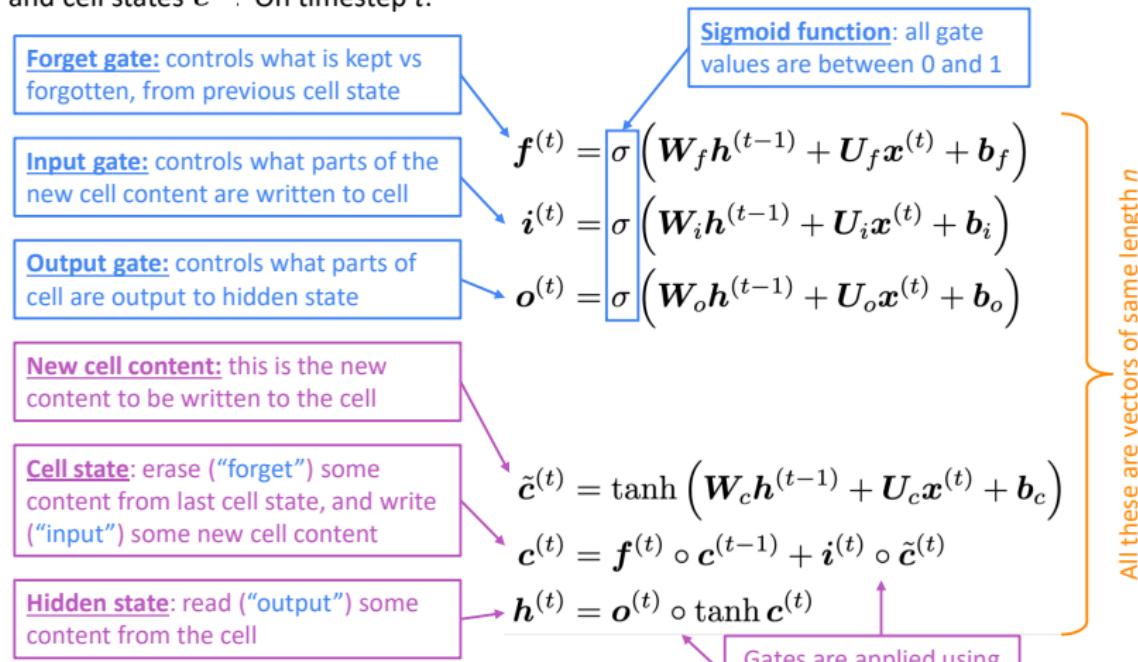
- A type of RNN proposed by Hochreiter and Schmidhuber in 1997 as a solution to the vanishing gradients problem.
- On step  $t$ , there is a hidden state  $h^{(t)}$  and a cell state  $c^{(t)}$ 
  - Both are vectors length  $n$
  - The cell stores long-term information
  - The LSTM can erase, write and 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 (LSTM)

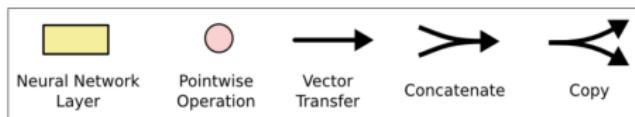
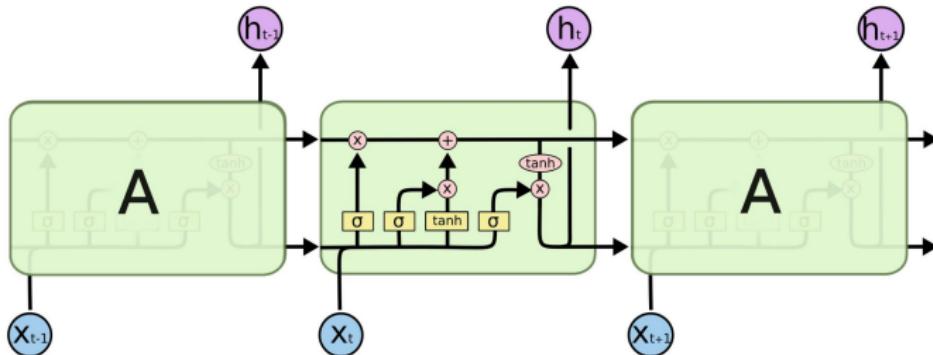
We have a sequence of inputs  $x^{(t)}$ , and we will compute a sequence of hidden states  $h^{(t)}$  and cell states  $c^{(t)}$ . On timestep  $t$ :





# Long Short-Term Memory (LSTM)

You can think of the LSTM equations visually like this:



24

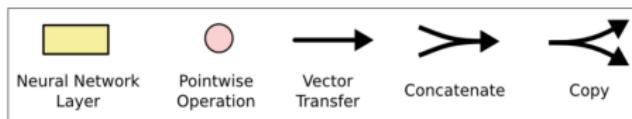
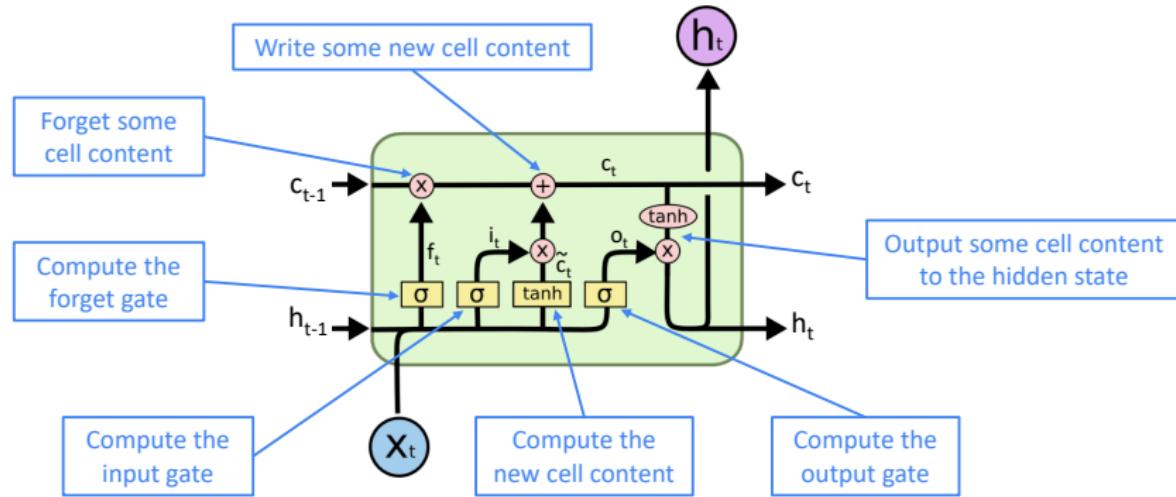
Source: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Christopher Manning, Natural Language Processing with Deep Learning, Standford U. CS224n



# Long Short-Term Memory (LSTM)

You can think of the LSTM equations visually like this:





## 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



## LSTMs: real-world success

- In 2013-2015, LSTMs started achieving state-of-the-art results
  - Successful tasks include: handwriting recognition, speech recognition, machine translation, parsing, image captioning
  - LSTM became the dominant approach
- Now (2019), other approaches (e.g. Transformers) have become more dominant for certain tasks.
  - For example in WMT (a MT conference + competition):
  - In WMT 2016, the summary report contains "RNN" 44 times
  - In WMT 2018, the report contains "RNN" 9 times and "Transformer" 63 times

Source: "Findings of the 2016 Conference on Machine Translation (WMT16)", Bojar et al. 2016, <http://www.statmt.org/wmt16/pdf/W16-2301.pdf>

Source: "Findings of the 2018 Conference on Machine Translation (WMT18)", Bojar et al. 2018, <http://www.statmt.org/wmt18/pdf/WMT028.pdf>



# Content

3

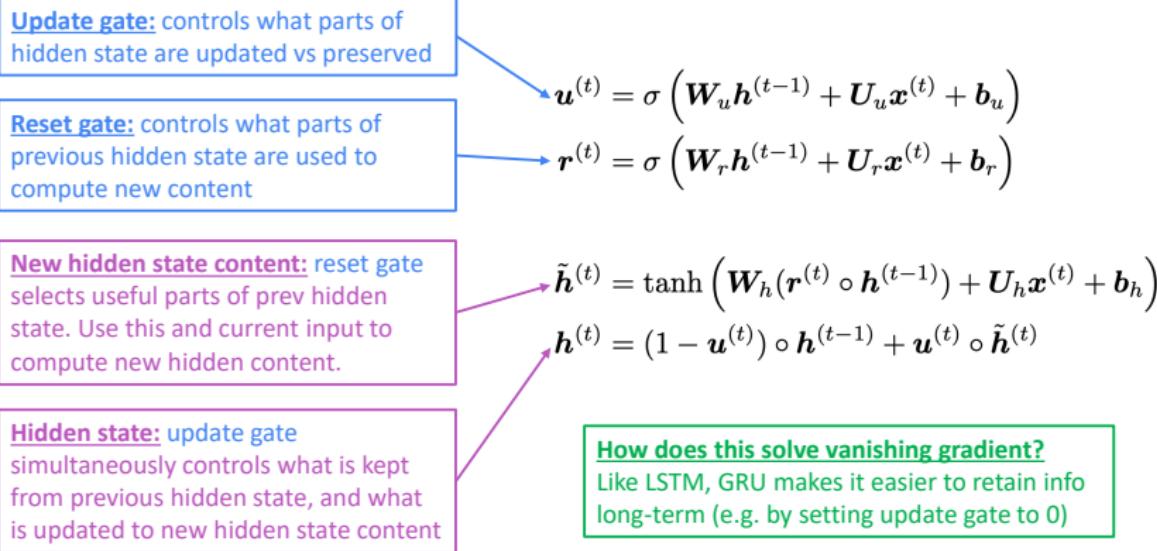
## Language Models Based on Recurrent Neural Networks

- Recurrent Neural Networks (RNNs)
- Training a neural language model
- Other applications of neural language models
- Evaluation of language models
- Gradient vanishing and exploding
- Long Short Term Memory (LSTM)
- **Gated Recurrent Units (GRUs)**
- Vanishing and exploding gradient for other neural networks
- Bi-directional RNNs and multi-layer RNNs



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





## 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



# Content

3

## Language Models Based on Recurrent Neural Networks

- Recurrent Neural Networks (RNNs)
- Training a neural language model
- Other applications of neural language models
- Evaluation of language models
- Gradient vanishing and exploding
- Long Short Term Memory (LSTM)
- Gated Recurrent Units (GRUs)
- **Vanishing and exploding gradient for other neural networks**
- Bi-directional RNNs and multi-layer RNNs

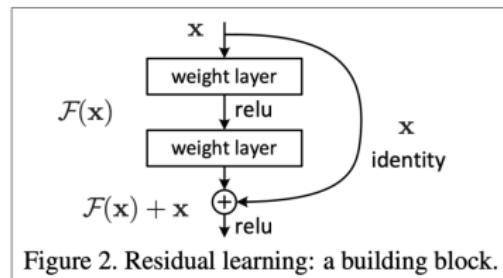


# Is vanishing/exploding gradient just a RNN problem?

- No! It can be a problem for all neural architectures (including **feed-forward** and **convolutional**), especially **deep** ones.
  - Due to chain rule / choice of nonlinearity function, gradient can become vanishingly small as it backpropagates
  - Thus lower layers are learnt very slowly (hard to train)
  - Solution: lots of new deep feedforward/convolutional architectures that add more direct connections (thus allowing the gradient to flow)

For example:

- Residual connections** aka “ResNet”
- Also known as **skip-connections**
- The **identity connection** preserves information by default
- This makes **deep** networks much easier to train





# Is vanishing/exploding gradient just a RNN problem?

- No! It can be a problem for all neural architectures (including **feed-forward** and **convolutional**), especially **deep** ones.
  - Due to chain rule / choice of nonlinearity function, gradient can become vanishingly small as it backpropagates
  - Thus lower layers are learnt very slowly (hard to train)
  - Solution: lots of new deep feedforward/convolutional architectures that add more direct connections (thus allowing the gradient to flow)

For example:

- Dense connections** aka “DenseNet”
- Directly connect everything to everything!

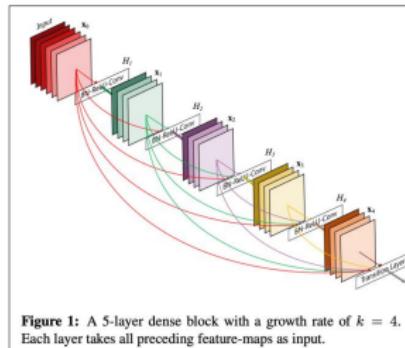


Figure 1: A 5-layer dense block with a growth rate of  $k = 4$ . Each layer takes all preceding feature-maps as input.





## Is vanishing/exploding gradient just a RNN problem?

- No! It can be a problem for all neural architectures (including **feed-forward** and **convolutional**), especially **deep** ones.
  - Due to chain rule / choice of nonlinearity function, gradient can become vanishingly small as it backpropagates
  - Thus lower layers are learnt very slowly (hard to train)
  - Solution: lots of new deep feedforward/convolutional architectures that **add more direct connections** (thus allowing the gradient to flow)

For example:

- **Highway connections** aka “HighwayNet”
- Similar to residual connections, but the identity connection vs the transformation layer is controlled by a **dynamic gate**
- Inspired by LSTMs, but applied to deep feedforward/convolutional networks





## Is vanishing/exploding gradient just a RNN problem?

- No! It can be a problem for all neural architectures (including **feed-forward** and **convolutional**), especially **deep** ones.
  - Due to chain rule / choice of nonlinearity function, gradient can become vanishingly small as it backpropagates
  - Thus lower layers are learnt very slowly (hard to train)
  - Solution: lots of new deep feedforward/convolutional architectures that **add more direct connections** (thus allowing the gradient to flow)



## Is vanishing/exploding gradient just a RNN problem?

- No! It can be a problem for all neural architectures (including **feed-forward** and **convolutional**), especially **deep** ones.
  - Due to chain rule / choice of nonlinearity function, gradient can become vanishingly small as it backpropagates
  - Thus lower layers are learnt very slowly (hard to train)
  - Solution: lots of new deep feedforward/convolutional architectures that **add more direct connections** (thus allowing the gradient to flow)
- Conclusion: Though vanishing/exploding gradients are a general problem, **RNNs are particularly unstable** due to the repeated multiplication by the **same** weight matrix [Bengio et al, 1994]





# Content

3

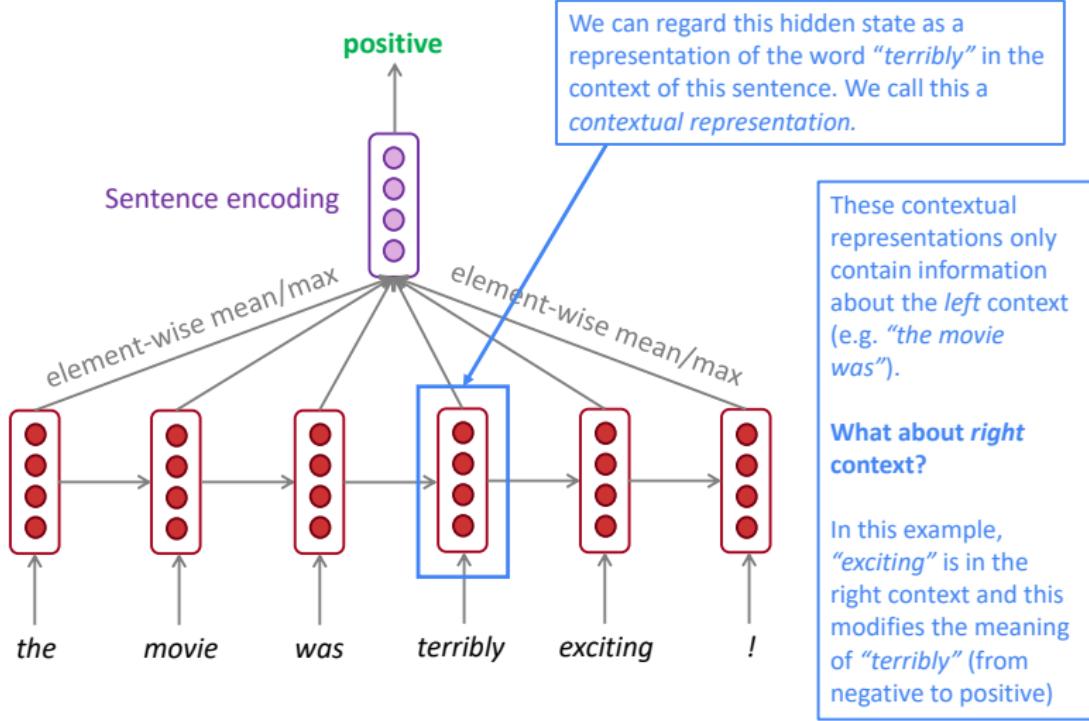
## Language Models Based on Recurrent Neural Networks

- Recurrent Neural Networks (RNNs)
- Training a neural language model
- Other applications of neural language models
- Evaluation of language models
- Gradient vanishing and exploding
- Long Short Term Memory (LSTM)
- Gated Recurrent Units (GRUs)
- Vanishing and exploding gradient for other neural networks
- Bi-directional RNNs and multi-layer RNNs



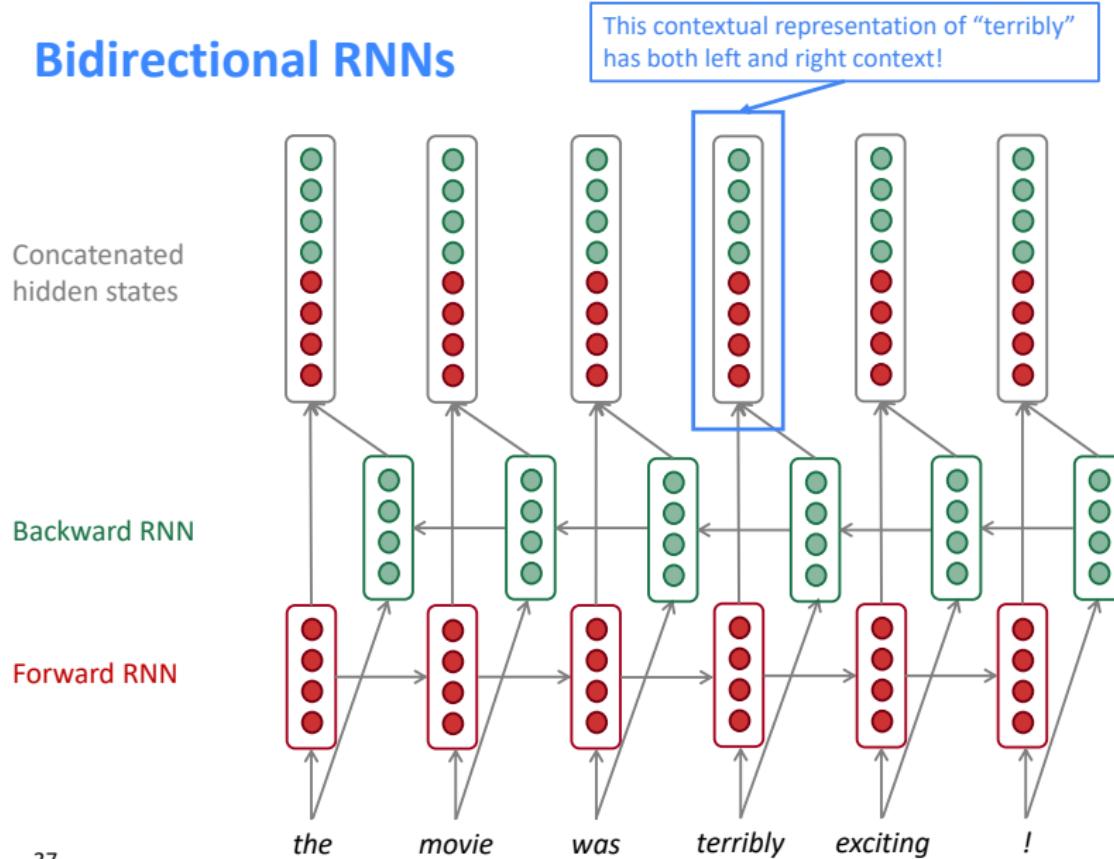
# Bidirectional RNNs: motivation

Task: Sentiment Classification





# Bidirectional RNNs





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

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

Concatenated hidden states

$$\mathbf{h}^{(t)} = [\vec{h}^{(t)}; \overleftarrow{h}^{(t)}]$$

Generally, these two RNNs have separate weights

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