

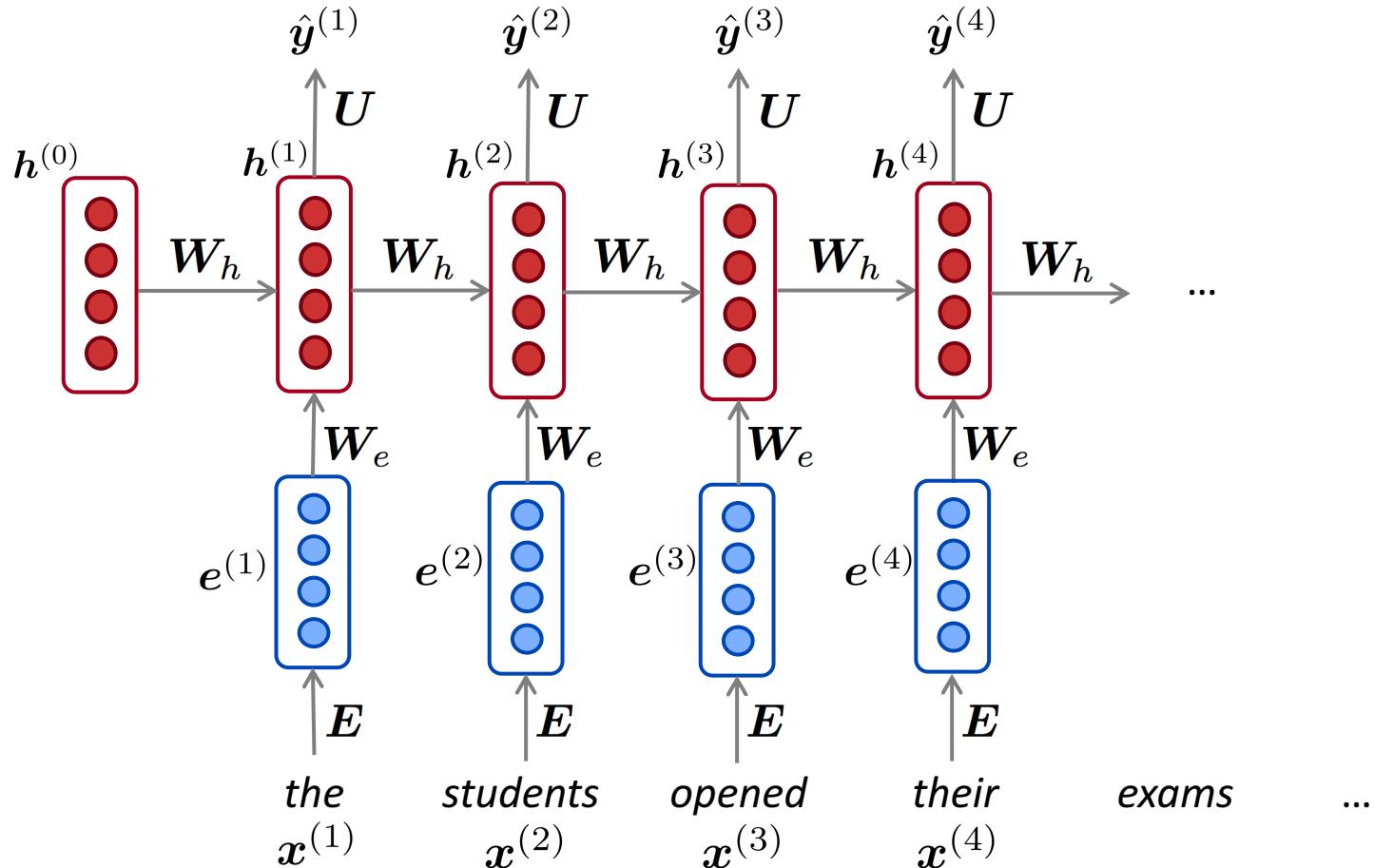
This lecture



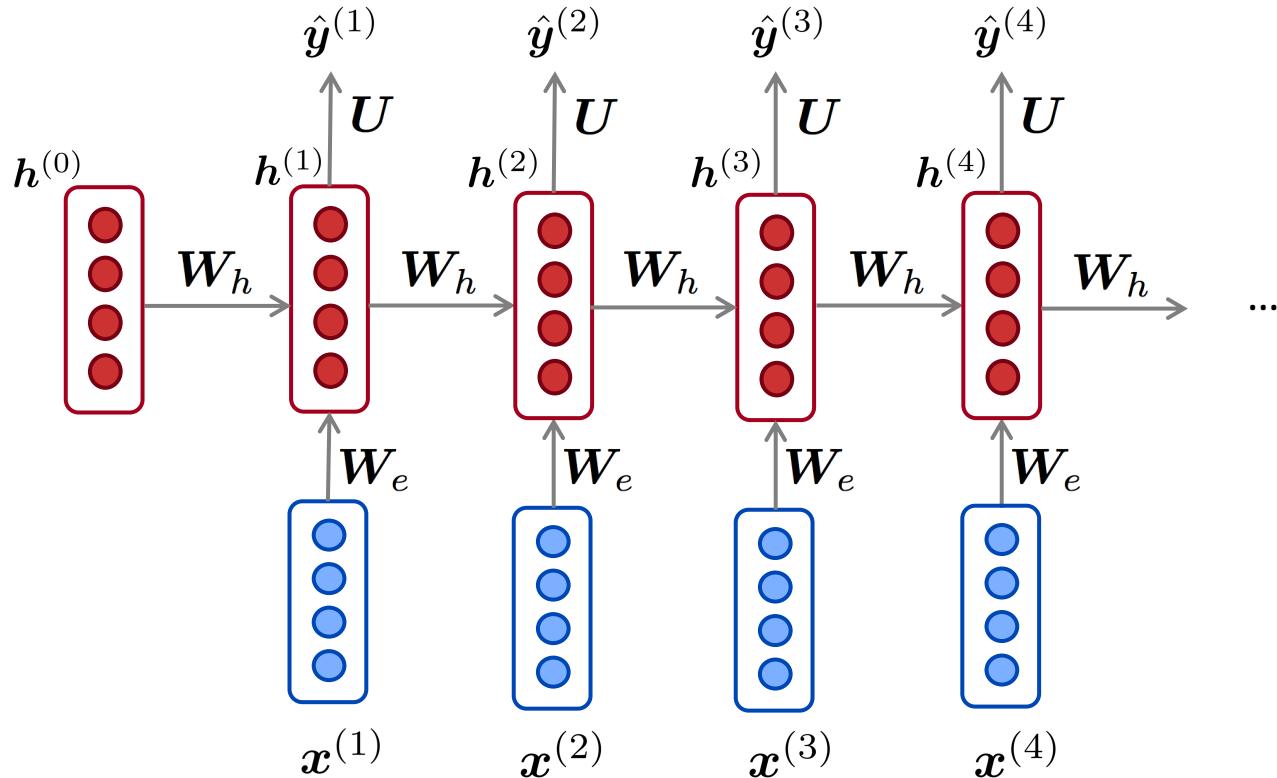
- Vanishing Gradient problem
- Fancy RNNs:
 - GRU
 - LSTM (!)
 - Bidirectional
 - Multi-layer

RNN Refresher

- Multiply the same matrix at each time step during forward prop



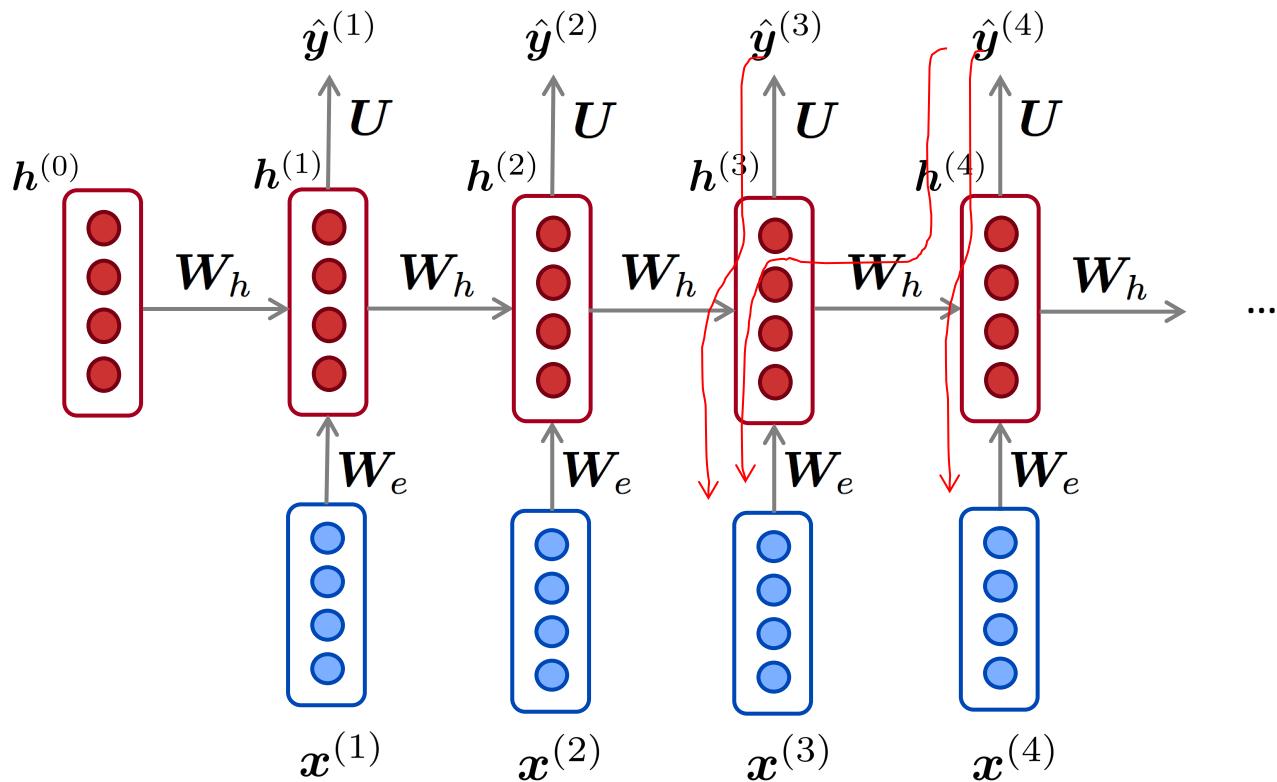
Simplify inputs to just x (usually word vectors)



- Ideally inputs from many time steps ago can modify output y
- Take $\frac{\partial E_2}{\partial W}$ for an example RNN with 2 time steps! Insightful!

The vanishing/exploding gradient problem

- Multiply the same matrix at each time step during backprop



The vanishing gradient problem - Details

- Similar but simpler RNN formulation:

$$\begin{aligned} h_t &= Wf(h_{t-1}) + W^{(hx)}x_{[t]} \\ \hat{y}_t &= W^{(S)}f(h_t) \end{aligned}$$

- Total error is the sum of each error (aka cost function, aka J in previous lectures when it was cross entropy error, could be other cost functions too), but at time steps t

$$\frac{\partial E}{\partial W} = \sum_{t=1}^T \frac{\partial E_t}{\partial W}$$

- Hardcore chain rule application:

$$\frac{\partial E_t}{\partial W} = \sum_{k=1}^t \frac{\partial E_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \frac{\partial h_t}{\partial h_k} \frac{\partial h_k}{\partial W}$$

The vanishing gradient problem - Details

- Similar to backprop but less efficient formulation
- Useful for analysis we'll look at:

$$\frac{\partial E_t}{\partial W} = \sum_{k=1}^t \frac{\partial E_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \boxed{\frac{\partial h_t}{\partial h_k}} \frac{\partial h_k}{\partial W}$$

We'll show this can quickly become very small or very large

- Remember: $h_t = Wf(h_{t-1}) + W^{(hx)}x_{[t]}$
- More chain rule, remember:

$$\frac{\partial h_t}{\partial h_k} = \prod_{j=k+1}^t \frac{\partial h_j}{\partial h_{j-1}}$$

- Each partial is a Jacobian:

$$\frac{d\mathbf{f}}{d\mathbf{x}} = \begin{bmatrix} \frac{\partial \mathbf{f}}{\partial x_1} & \dots & \frac{\partial \mathbf{f}}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

The vanishing gradient problem - Details

- Analyzing the norms of the Jacobians yields:

$$\left\| \frac{\partial h_j}{\partial h_{j-1}} \right\| \leq \|W^T\| \|\text{diag}[f'(h_{j-1})]\| \leq \beta_W \beta_h$$

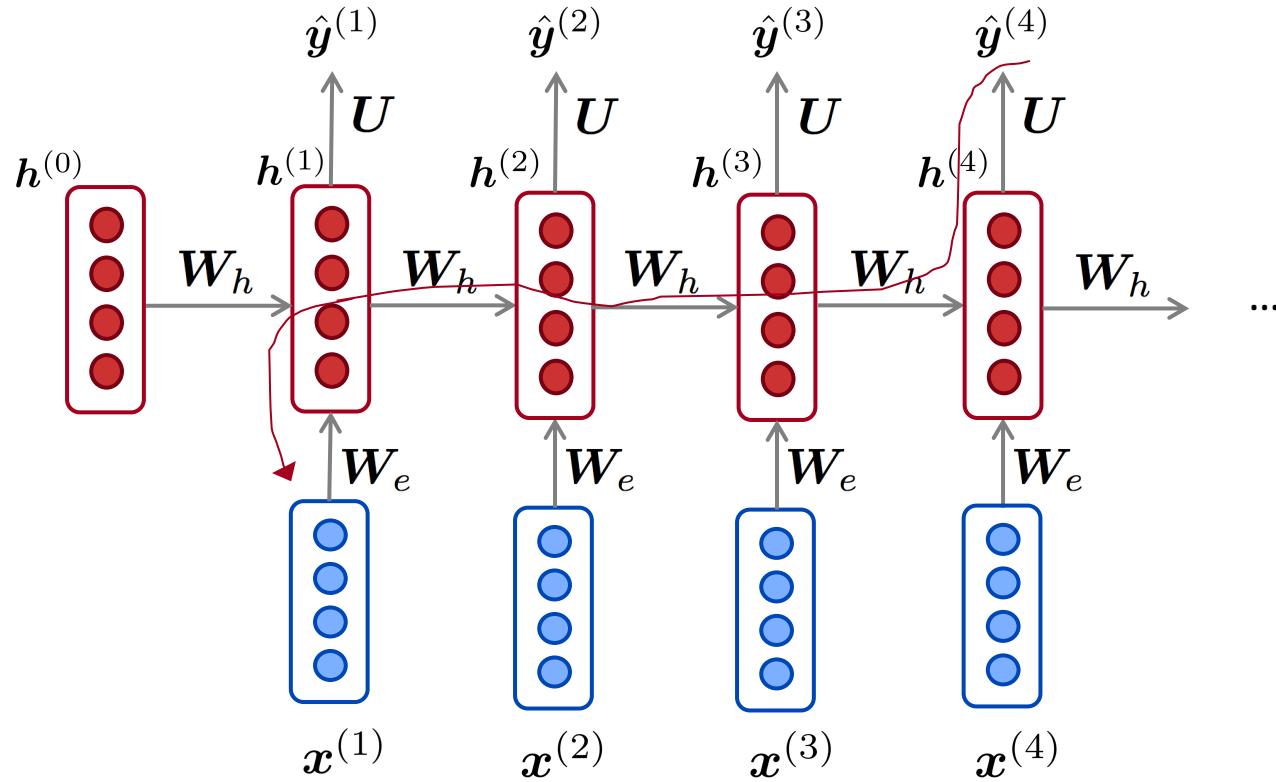
- Where we defined β 's as upper bounds of the norms
- The gradient is a product of Jacobian matrices, each associated with a step in the forward computation.

$$\left\| \frac{\partial h_t}{\partial h_k} \right\| = \left\| \prod_{j=k+1}^t \frac{\partial h_j}{\partial h_{j-1}} \right\| \leq (\beta_W \beta_h)^{t-k}$$

- This can become **very small or very large quickly** [Bengio et al 1994], and the locality assumption of gradient descent breaks down. → **Vanishing or exploding gradient**

Why is the vanishing gradient a problem?

- Ideally, the error E^t on step t can flow backwards, via backprop, and allow the weights on a previous timestep (maybe many timesteps ago) to change.



Why is the vanishing gradient a problem?

1. Gradients can be seen as a measure of influence of the past on the future
2. How does the perturbation at time t affect predictions at $t+n$?

Why is the vanishing gradient a problem?

When we only observe

$$\left\| \frac{\partial h_t}{\partial h_k} \right\| = \left\| \prod_{j=k+1}^t \frac{\partial h_j}{\partial h_{j-1}} \right\| \leq (\beta_W \beta_h)^{t-k} \quad \text{going to 0}$$

We cannot tell whether

1. No dependency between t and $t+n$ in data, or
2. Wrong configuration of parameters

The vanishing gradient problem for language models

- The vanishing gradient problem can cause problems for RNN Language Models:
- When predicting the next word, information from many time steps in the past is not taken into consideration.
- Example:

Jane walked into the room. John walked in too. It was late in the day. Jane said hi to _____

Trick for exploding gradient: clipping trick

- The solution first introduced by Mikolov is to clip gradients so that their norm has some maximum value.

Algorithm 1 Pseudo-code for norm clipping the gradients whenever they explode

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

- Makes a big difference in RNNs and many other unstable models

Main solution for better RNNs: Better Units

- The main solution to the Vanishing Gradient Problem is to use a more complex hidden unit computation in recurrence!
- Gated Recurrent Units (GRU) introduced by [Cho et al. 2014] and LSTMs [Hochreiter & Schmidhuber, 1999]
- Main ideas:
 - keep around memories to capture long distance dependencies
 - allow error messages to flow at different strengths depending on the inputs

GRUs

- Standard RNN computes hidden layer at next time step directly:
$$h_t = f \left(W^{(hh)} h_{t-1} + W^{(hx)} x_t \right)$$
- Simplified from last lecture's:
$$h^{(t)} = \sigma \left(W_h h^{(t-1)} + W_e e^{(t)} + b_1 \right)$$
- GRU first computes an update **gate** (another layer) based on current input word vector and hidden state

$$z_t = \sigma \left(W^{(z)} x_t + U^{(z)} h_{t-1} \right)$$

- Compute reset gate similarly but with different weights

$$r_t = \sigma \left(W^{(r)} x_t + U^{(r)} h_{t-1} \right)$$

GRUs

- Update gate

$$z_t = \sigma \left(W^{(z)} x_t + U^{(z)} h_{t-1} \right)$$

- Reset gate

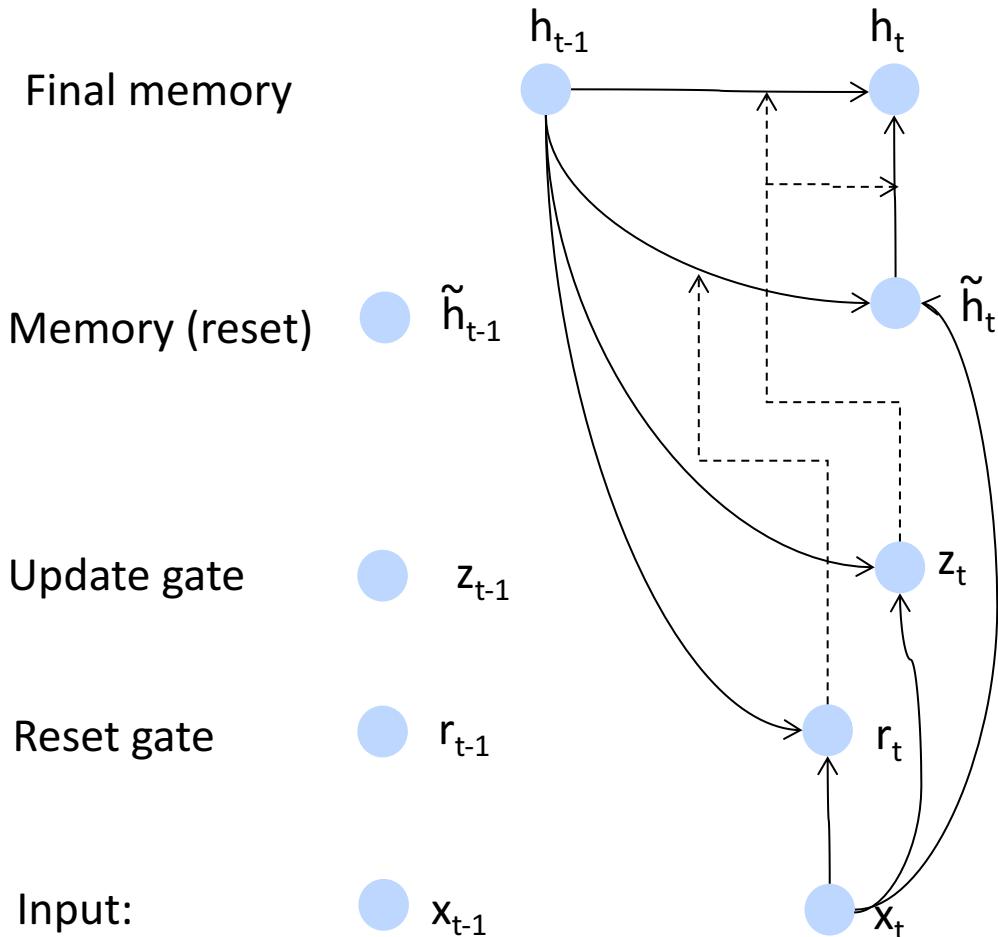
$$r_t = \sigma \left(W^{(r)} x_t + U^{(r)} h_{t-1} \right)$$

- New memory content: $\tilde{h}_t = \tanh(W x_t + r_t \circ U h_{t-1})$

If reset gate unit is ~ 0 , then this ignores previous memory and only stores the new word information

- Final memory at time step combines current and previous time steps: $h_t = z_t \circ h_{t-1} + (1 - z_t) \circ \tilde{h}_t$

GRU illustration



$$z_t = \sigma \left(W^{(z)} x_t + U^{(z)} h_{t-1} \right)$$
$$r_t = \sigma \left(W^{(r)} x_t + U^{(r)} h_{t-1} \right)$$

$$\tilde{h}_t = \tanh \left(W x_t + r_t \circ U h_{t-1} \right)$$

$$h_t = z_t \circ h_{t-1} + (1 - z_t) \circ \tilde{h}_t$$

GRU intuition

- If reset is close to 0, ignore previous hidden state
→ Allows model to drop information that is irrelevant in the future
- Update gate z controls how much of past state should matter now.
 - If z close to 1, then we can copy information in that unit through many time steps! **Less vanishing gradient!**
- Units with short-term dependencies often have reset gates very active

$$z_t = \sigma \left(W^{(z)} x_t + U^{(z)} h_{t-1} \right)$$

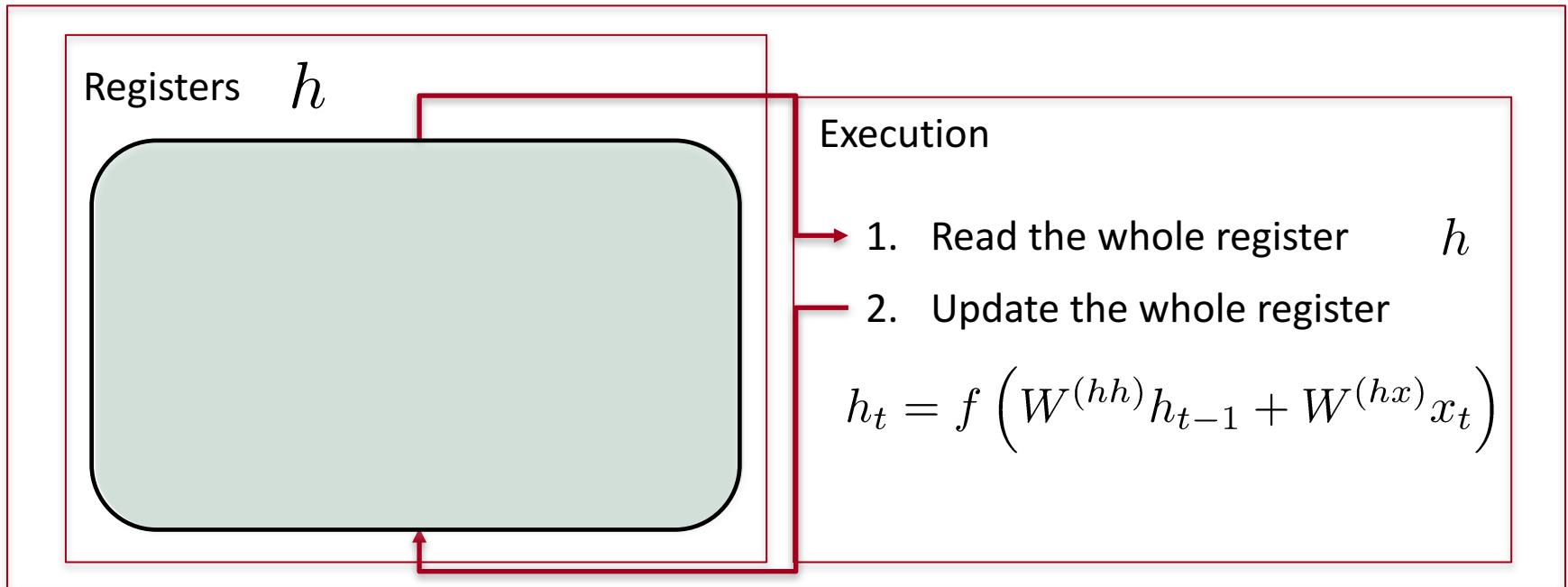
$$r_t = \sigma \left(W^{(r)} x_t + U^{(r)} h_{t-1} \right)$$

$$\tilde{h}_t = \tanh (W x_t + r_t \circ U h_{t-1})$$

$$h_t = z_t \circ h_{t-1} + (1 - z_t) \circ \tilde{h}_t$$

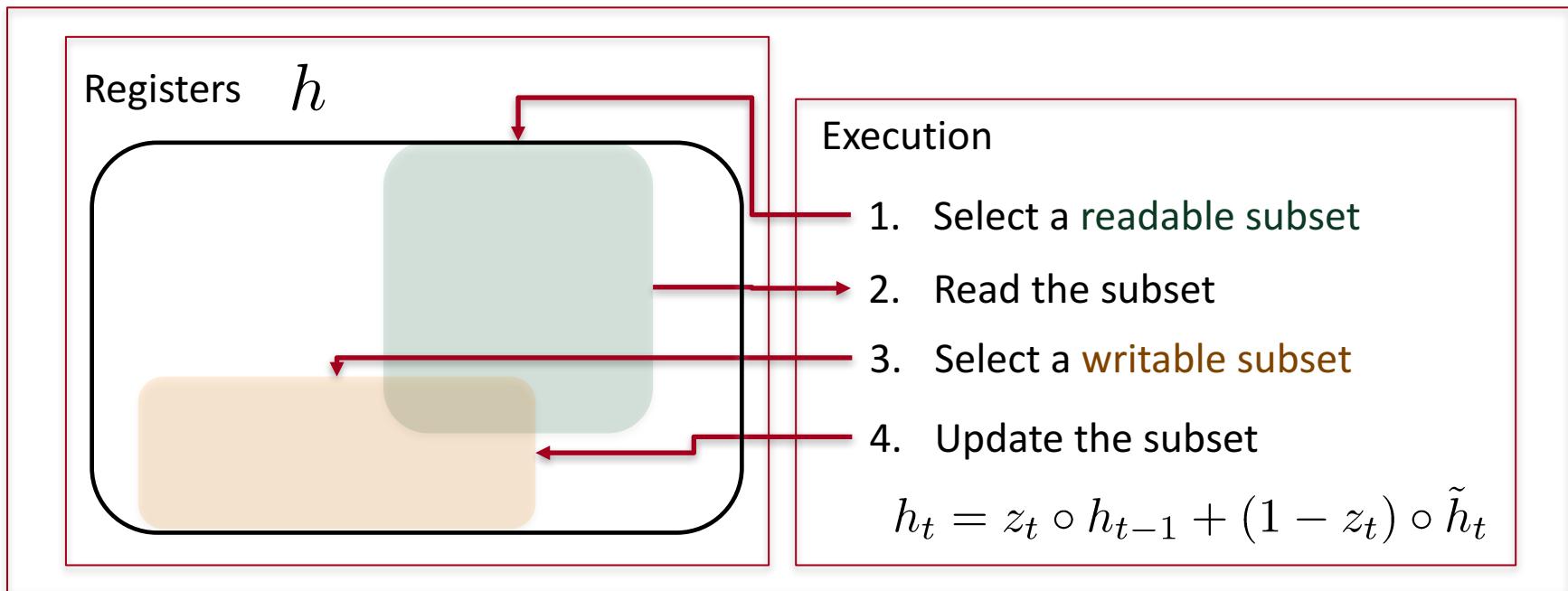
GRU Comparison to Standard tanh-RNN

Vanilla RNN ...



GRU Comparison to Standard tanh-RNN

GRU ...



Gated recurrent units are much more versatile and adaptive in which elements of the hidden vector h they update!

Long-short-term-memories (LSTMs)

- LSTM is even more complex than GRU
- Allow each time step to modify
 - Input gate (current cell matters)
 - Forget (gate 0, forget past)
 - Output (how much cell is exposed)
 - New memory cell
- Final memory cell:
- Final hidden state:

$$i_t = \sigma \left(W^{(i)}x_t + U^{(i)}h_{t-1} \right)$$

$$f_t = \sigma \left(W^{(f)}x_t + U^{(f)}h_{t-1} \right)$$

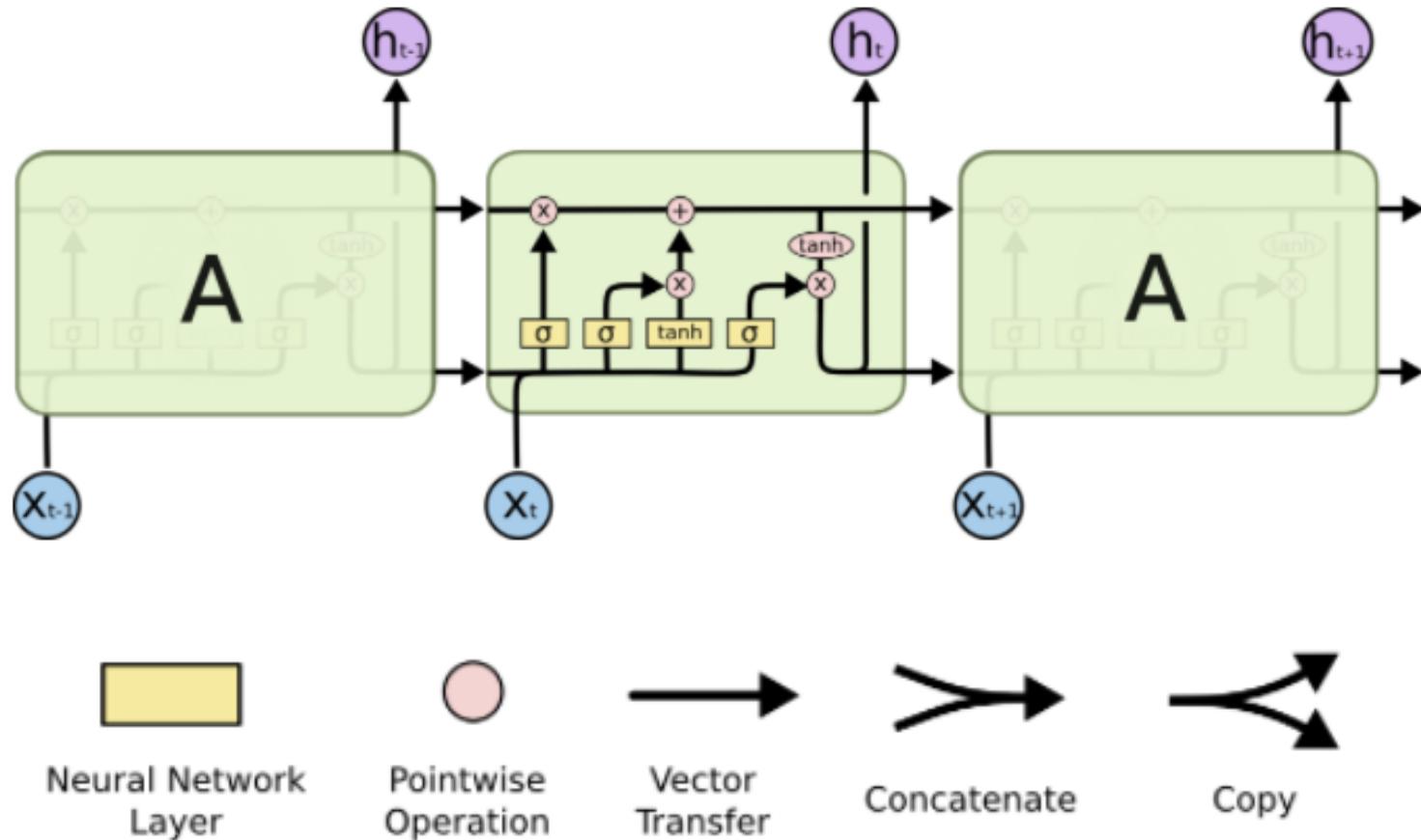
$$o_t = \sigma \left(W^{(o)}x_t + U^{(o)}h_{t-1} \right)$$

$$\tilde{c}_t = \tanh \left(W^{(c)}x_t + U^{(c)}h_{t-1} \right)$$

$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$$

$$h_t = o_t \circ \tanh(c_t)$$

Some visualizations



By Chris Olah: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

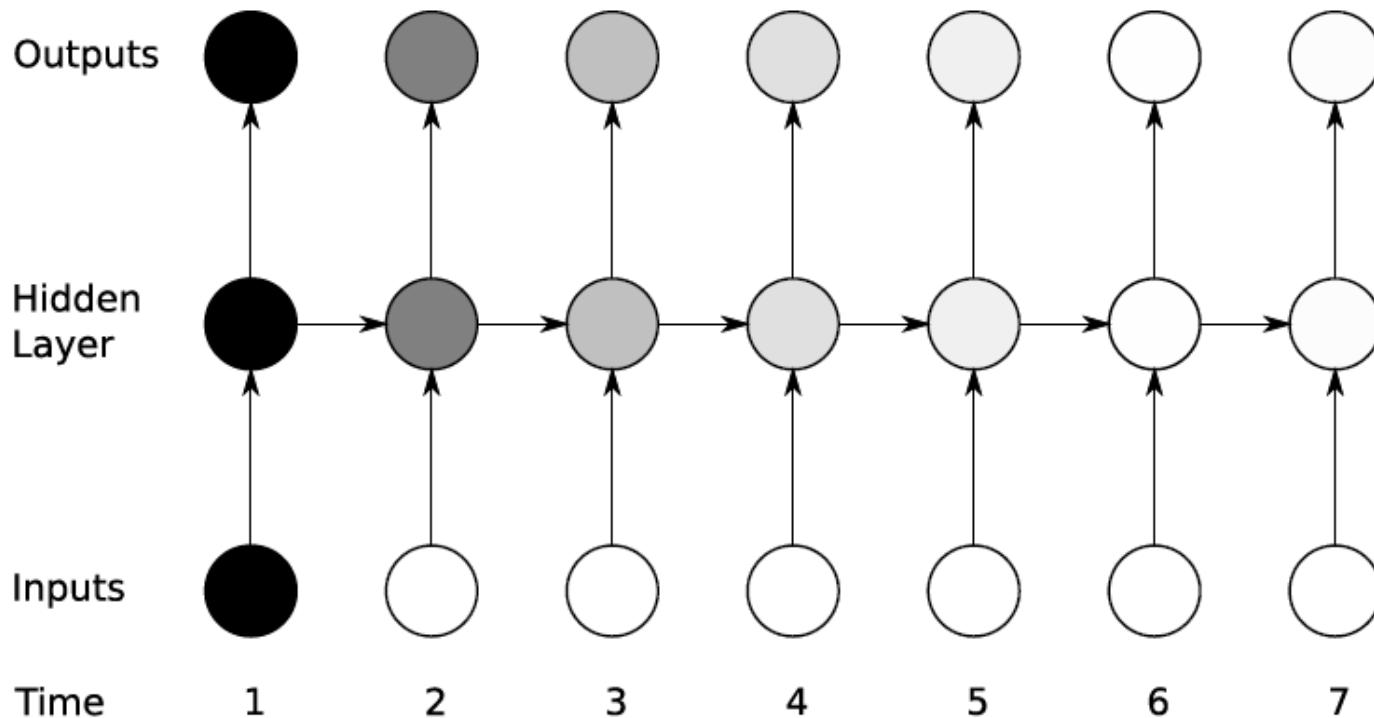


Figure 4.1: **The vanishing gradient problem for RNNs.** The shading of the nodes in the unfolded network indicates their sensitivity to the inputs at time one (the darker the shade, the greater the sensitivity). The sensitivity decays over time as new inputs overwrite the activations of the hidden layer, and the network ‘forgets’ the first inputs.

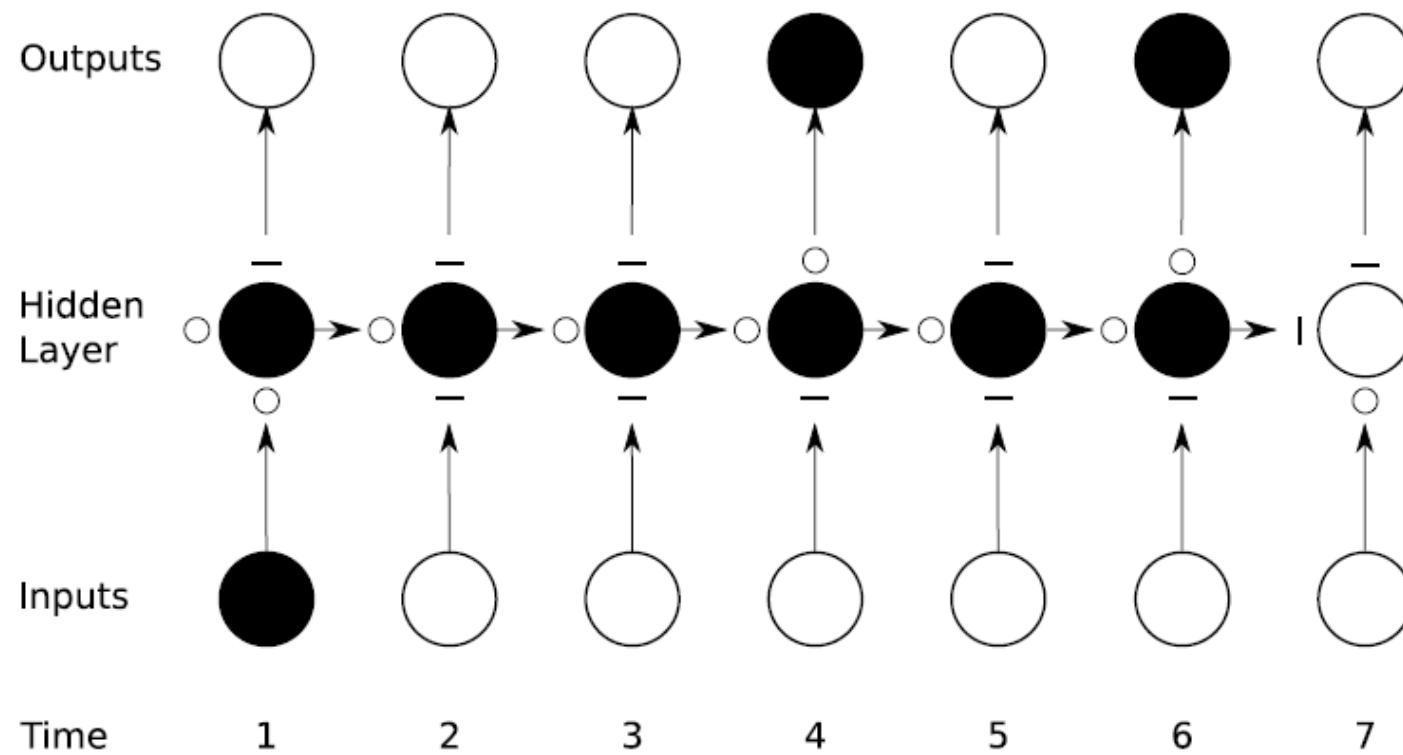


Figure 4.4: **Preservation of gradient information by LSTM.** As in Figure 4.1 the shading of the nodes indicates their sensitivity to the inputs at time one; in this case the black nodes are maximally sensitive and the white nodes are entirely insensitive. The state of the input, forget, and output gates are displayed below, to the left and above the hidden layer respectively. For simplicity, all gates are either entirely open ('O') or closed ('—'). The memory cell 'remembers' the first input as long as the forget gate is open and the input gate is closed. The sensitivity of the output layer can be switched on and off by the output gate without affecting the cell.

LSTMs are a great default for all sequence problems

- Very powerful, especially when stacked and made even deeper (each hidden layer is already computed by a deep internal network)
- Most useful if you have lots and lots of data

Deep LSTMs compared to traditional systems 2015

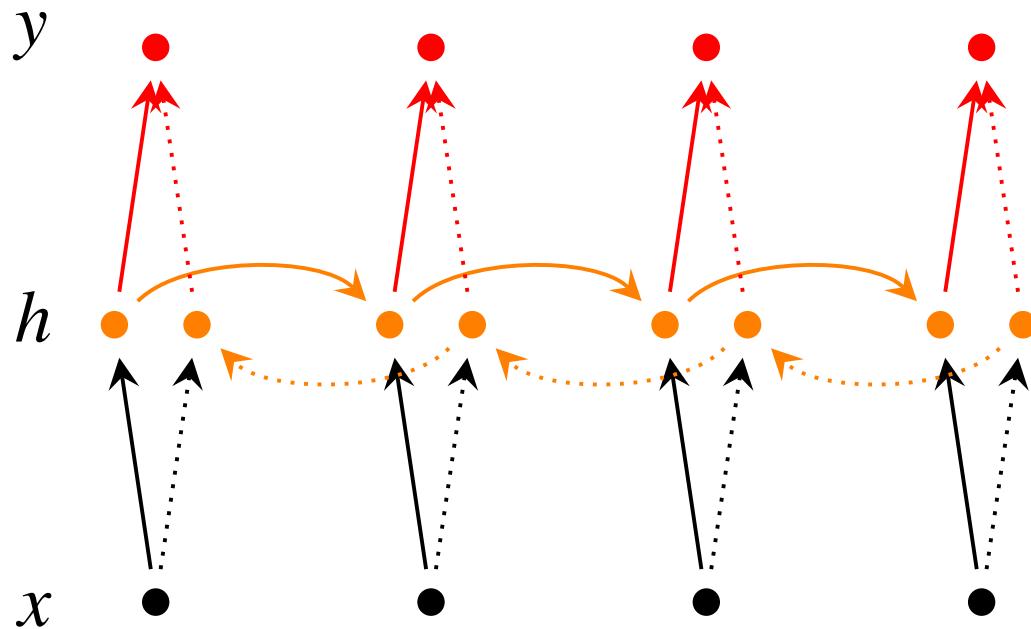
Method	test BLEU score (ntst14)
Bahdanau et al. [2]	28.45
Baseline System [29]	33.30
Single forward LSTM, beam size 12	26.17
Single reversed LSTM, beam size 12	30.59
Ensemble of 5 reversed LSTMs, beam size 1	33.00
Ensemble of 2 reversed LSTMs, beam size 12	33.27
Ensemble of 5 reversed LSTMs, beam size 2	34.50
Ensemble of 5 reversed LSTMs, beam size 12	34.81

Table 1: The performance of the LSTM on WMT'14 English to French test set (ntst14). Note that an ensemble of 5 LSTMs with a beam of size 2 is cheaper than of a single LSTM with a beam of size 12.

Method	test BLEU score (ntst14)
Baseline System [29]	33.30
Cho et al. [5]	34.54
Best WMT'14 result [9]	37.0
Rescoring the baseline 1000-best with a single forward LSTM	35.61
Rescoring the baseline 1000-best with a single reversed LSTM	35.85
Rescoring the baseline 1000-best with an ensemble of 5 reversed LSTMs	36.5
Oracle Rescoring of the Baseline 1000-best lists	~45

Bidirectional RNNs

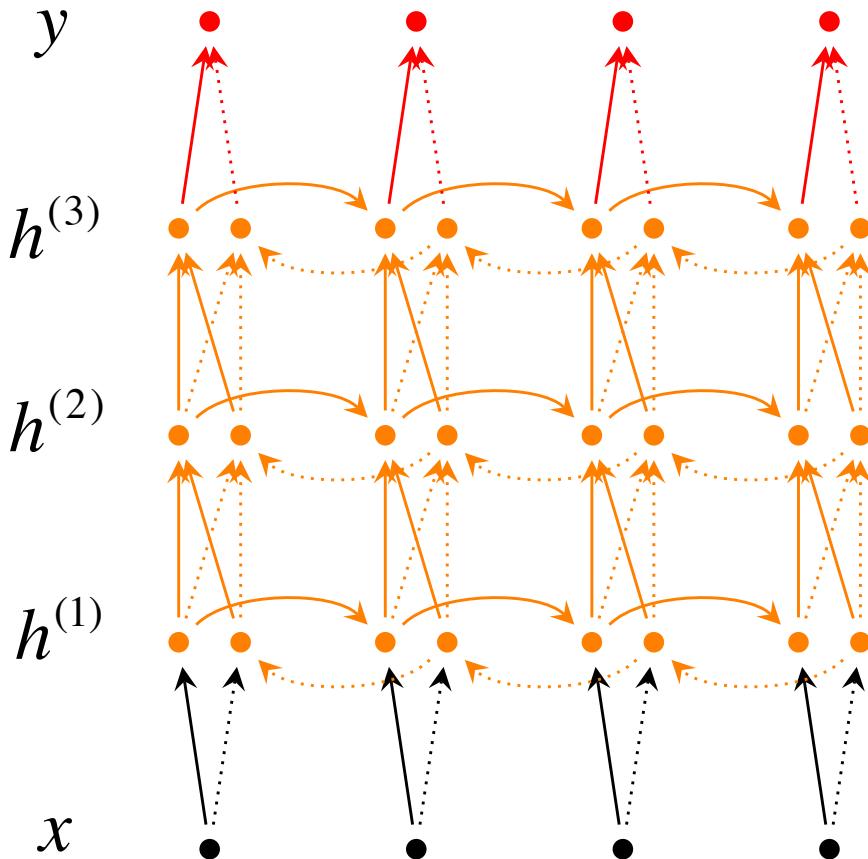
Problem: For classification you want to incorporate information from words both preceding and following



$$\begin{aligned}\vec{h}_t &= f(\vec{W}x_t + \vec{V}\vec{h}_{t-1} + \vec{b}) \\ \overleftarrow{h}_t &= f(\overleftarrow{W}x_t + \overleftarrow{V}\overleftarrow{h}_{t+1} + \overleftarrow{b}) \\ y_t &= g(U[\vec{h}_t; \overleftarrow{h}_t] + c)\end{aligned}$$

$h = [\vec{h}; \overleftarrow{h}]$ now represents (summarizes) the past and future around a single token.

Deep Bidirectional RNNs



$$\begin{aligned}\overrightarrow{h}_t^{(i)} &= f(\overrightarrow{W}^{(i)} \overrightarrow{h}_{t-1}^{(i-1)} + \overrightarrow{V}^{(i)} \overrightarrow{h}_{t-1}^{(i)} + \overrightarrow{b}^{(i)}) \\ \overleftarrow{h}_t^{(i)} &= f(\overleftarrow{W}^{(i)} \overleftarrow{h}_t^{(i-1)} + \overleftarrow{V}^{(i)} \overleftarrow{h}_{t+1}^{(i)} + \overleftarrow{b}^{(i)}) \\ y_t &= g(U[\overrightarrow{h}_t^{(L)}; \overleftarrow{h}_t^{(L)}] + c)\end{aligned}$$

Each memory layer passes an intermediate sequential representation to the next.