

■ CS286 AI for Science and Engineering

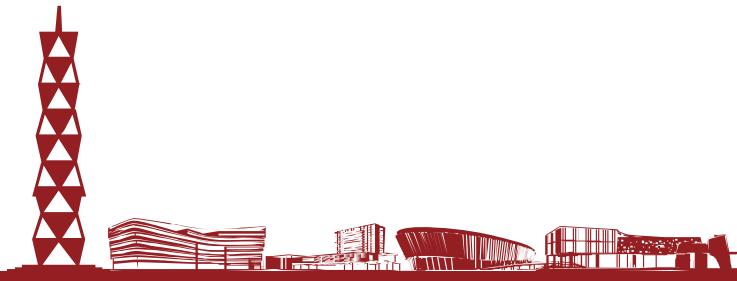
Lecture 8: Recurrent Neural Networks

Xuming He (何旭明)

PhD, Associate Professor

School of Information Science and Technology (SIST), ShanghaiTech University

Fall Semester, 2020





Outline



- Basic Recurrent Neural Networks
- Training Recurrent Neural Networks
- LSTM
- RNN Applications



Acknowledgement: Hugo Larochelle's, Mehryar Mohri@NYU's, Yingyu Liang@Princeton's, Bhiksha Raj@CMU's & Feifei Li@Stanford's course notes

Motivation: Sequence modeling

- Modeling a sequence of tokens
 - Running example: sentences, but it can also be gene or protein sequences, etc.
- Goal: learn/build a good distribution of sentences
- Inputs: a corpus of sentences $\mathbf{s}^{(1)}, \dots, \mathbf{s}^{(N)}$
- Output: a distribution $p(\mathbf{s})$
- Common approach: **maximum likelihood**
 - Assume sentences are independent

$$\max \prod_{i=1}^N p(\mathbf{s}^{(i)})$$





Sequence modeling

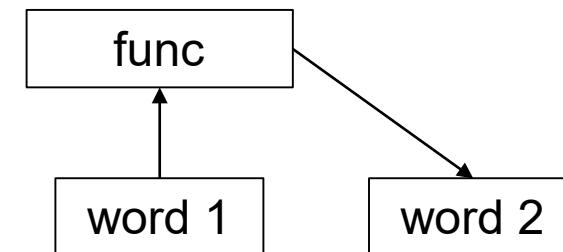
- What is $p(\mathbf{s})$?
- A sentence is a sequence of words w_1, w_2, \dots, w_T .

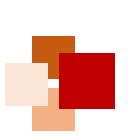
$$p(\mathbf{s}) = p(w_1, \dots, w_T) = p(w_1)p(w_2 | w_1) \cdots p(w_T | w_1, \dots, w_{T-1}).$$

- Essentially aim to predict the next word
- Markovian assumption
 - The distribution over the next word depends on the preceding few words. For example,

$$p(w_t | w_1, \dots, w_{t-1}) = p(w_t | w_{t-3}, w_{t-2}, w_{t-1}).$$

- Autoregressive model
 - Memoryless

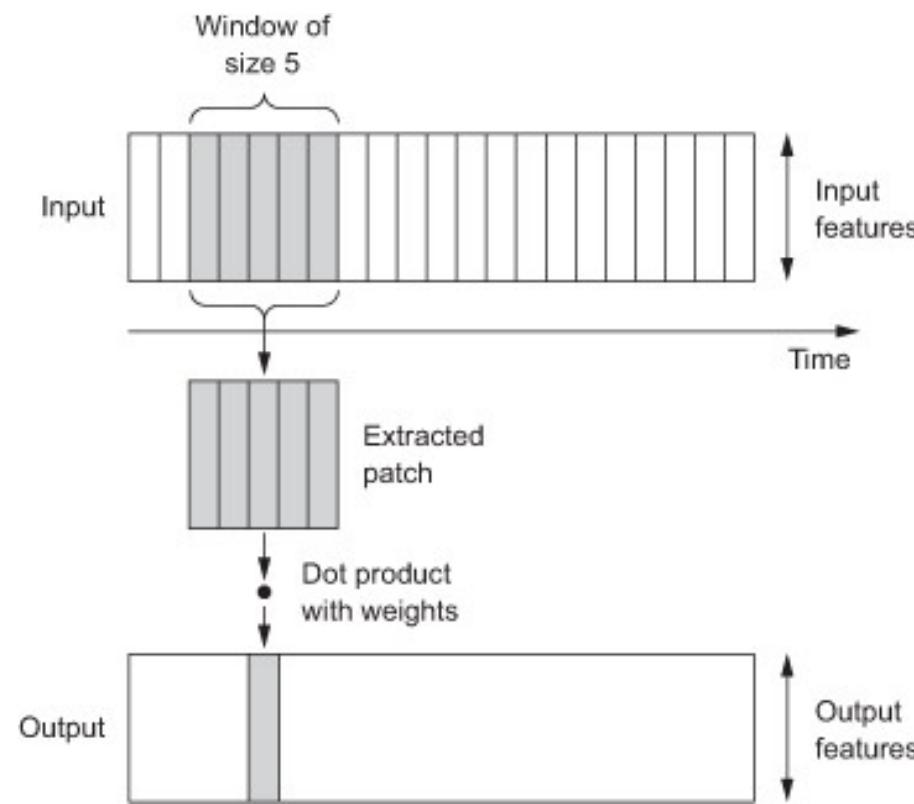




Sequence modeling



- What is $p(s)$?
- A sentence is a sequence of words w_1, w_2, \dots, w_T .
- Can we use CNN?



<https://medium.com/@jon.froland/convolutional-neural-networks-for-sequence-processing-part-1-420dd9b500>





Outline

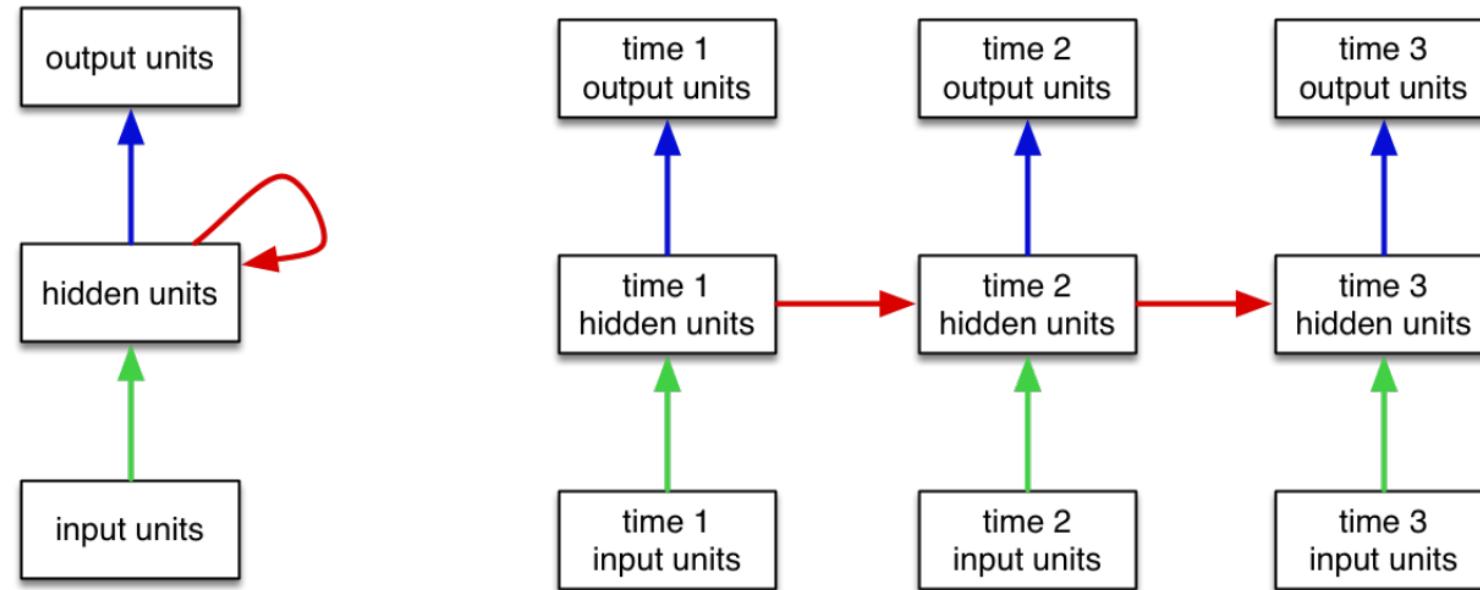


- Basic Recurrent Neural Networks
 - (Vanilla) RNN models
 - Example: language modeling
 - Other RNN problem settings
 - More examples: Neural language models, Image generation



Recurrent Neural Network

- Recurrent Neural Network as a dynamical system with one set of hidden units feeding into themselves
 - The network's graph has self-loops
- The RNN's graph can be unrolled by explicitly representing the units at all time steps
 - The weights and biases are shared



Recurrent Neural Network

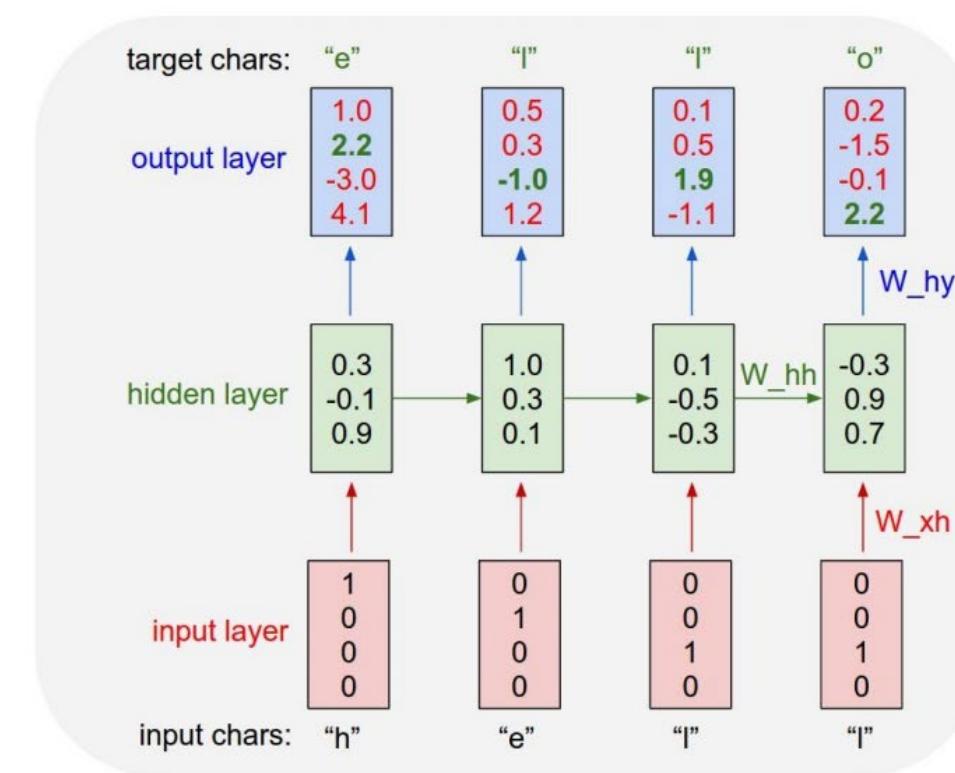


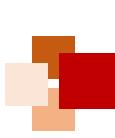
- The RNN's graph can be unrolled by explicitly representing the units at all time steps
 - The weights and biases are shared

**Example:
Character-level
Language Model**

Vocabulary:
[h,e,l,o]

Example training
sequence:
“hello”





Recurrent Neural Network



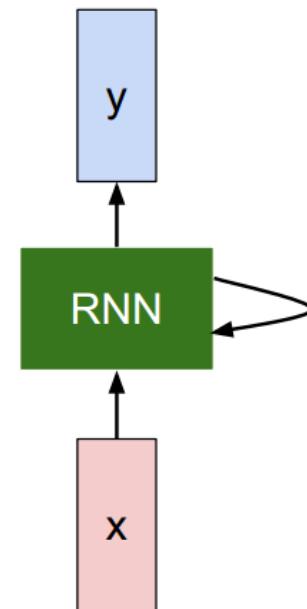
- General formulation

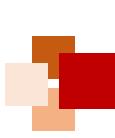
We can process a sequence of vectors \mathbf{x} by applying a **recurrence formula** at every time step:

$$h_t = f_W(h_{t-1}, x_t)$$

new state / old state input vector at some time step

some function with parameters W



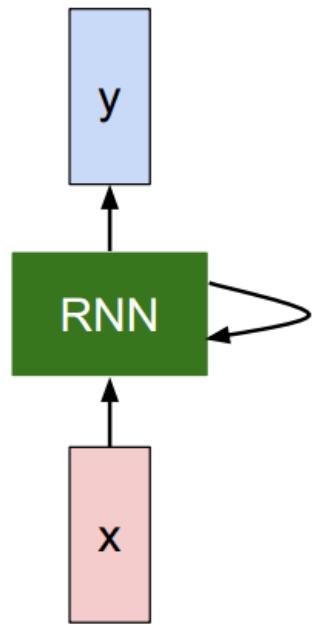


(Vanilla) Recurrent Neural Network



- General formulation

The state consists of a single “*hidden*” vector \mathbf{h} :

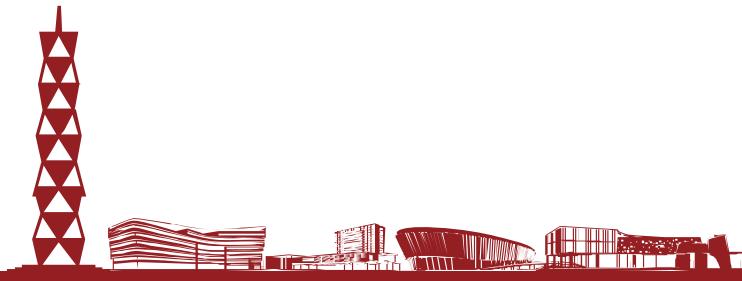


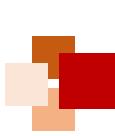
$$h_t = f_W(h_{t-1}, x_t)$$



$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$y_t = W_{hy}h_t$$





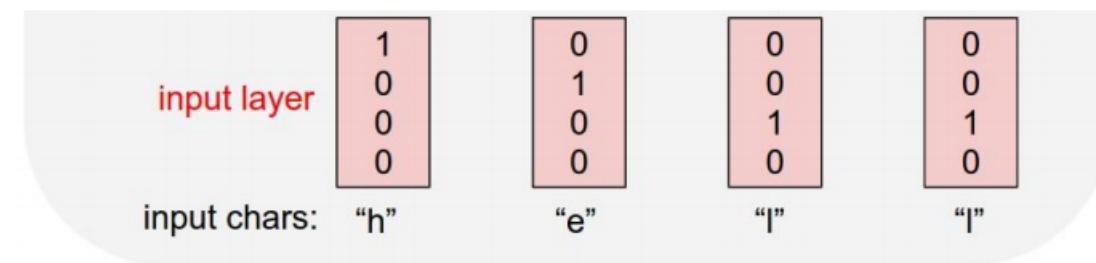
RNN for language modeling



Example: Character-level Language Model

Vocabulary:
[h,e,l,o]

Example training
sequence:
“hello”



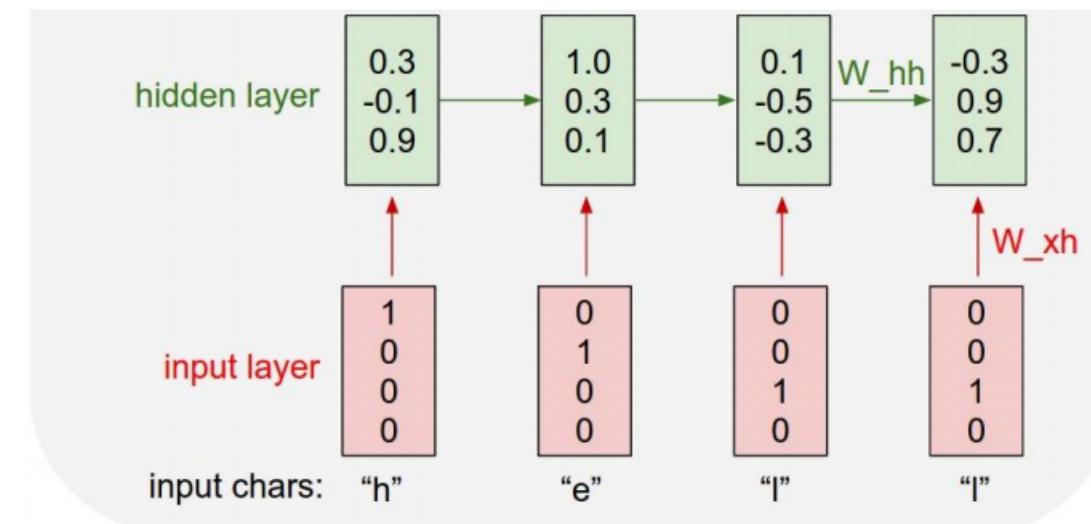
RNN for language modeling

**Example:
Character-level
Language Model**

Vocabulary:
[h,e,l,o]

Example training
sequence:
“hello”

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

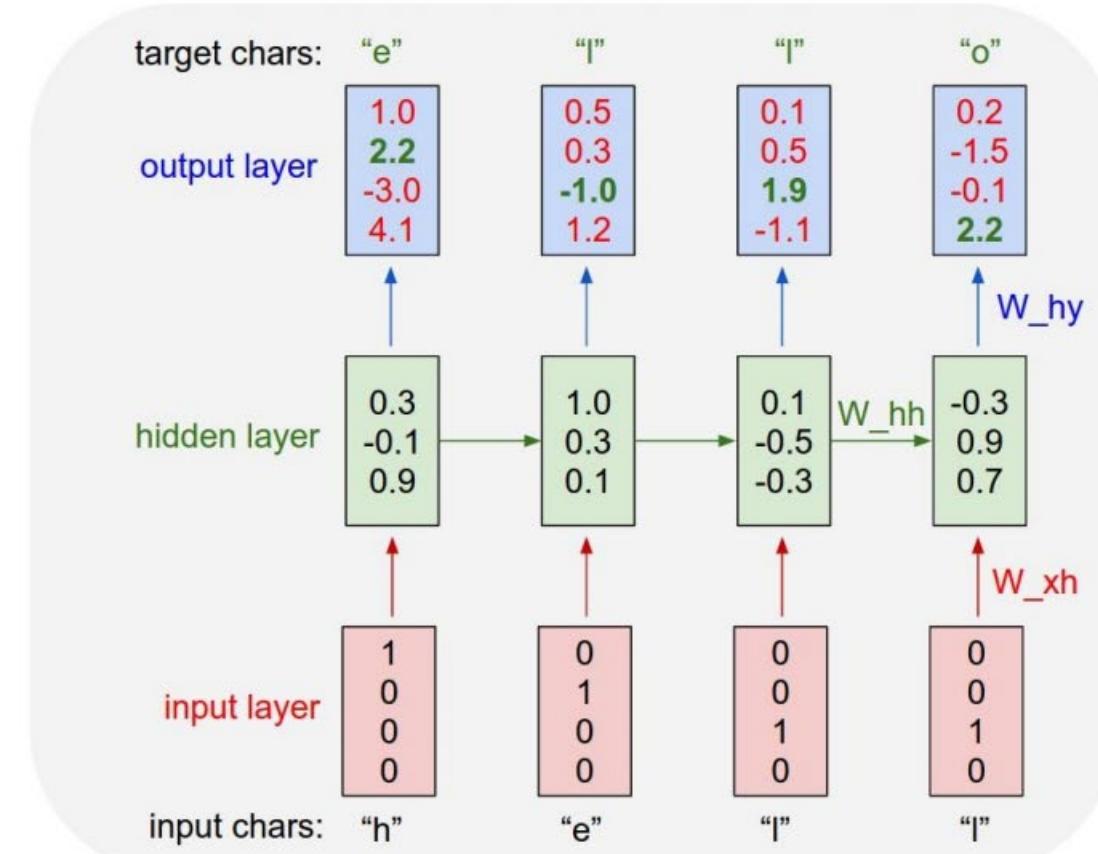


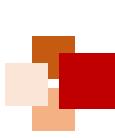
RNN for language modeling

Example: Character-level Language Model

Vocabulary:
[h,e,l,o]

Example training
sequence:
“hello”





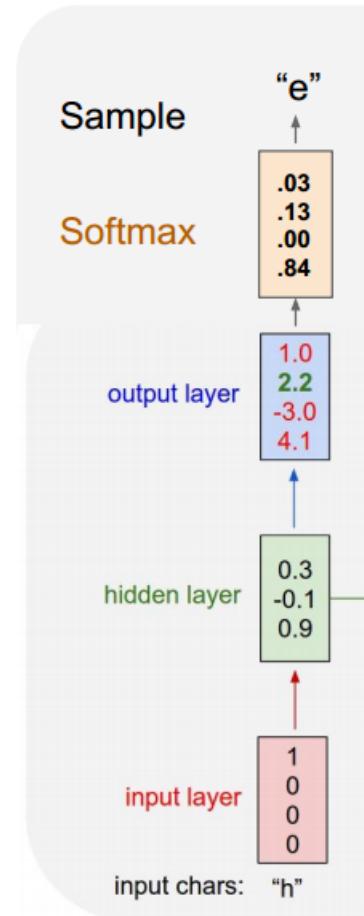
RNN for language modeling

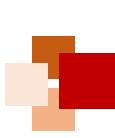


Example: Character-level Language Model Sampling

Vocabulary:
[h,e,l,o]

At test-time sample
characters one at a time,
feed back to model





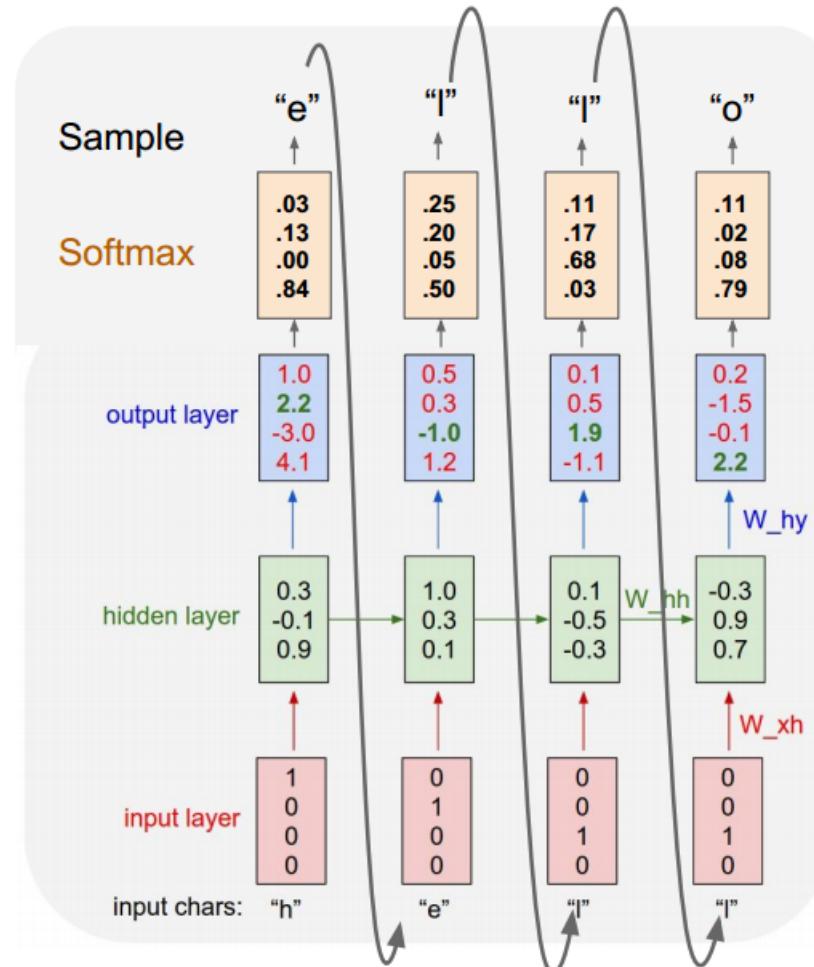
RNN for language modeling



**Example:
Character-level
Language Model
Sampling**

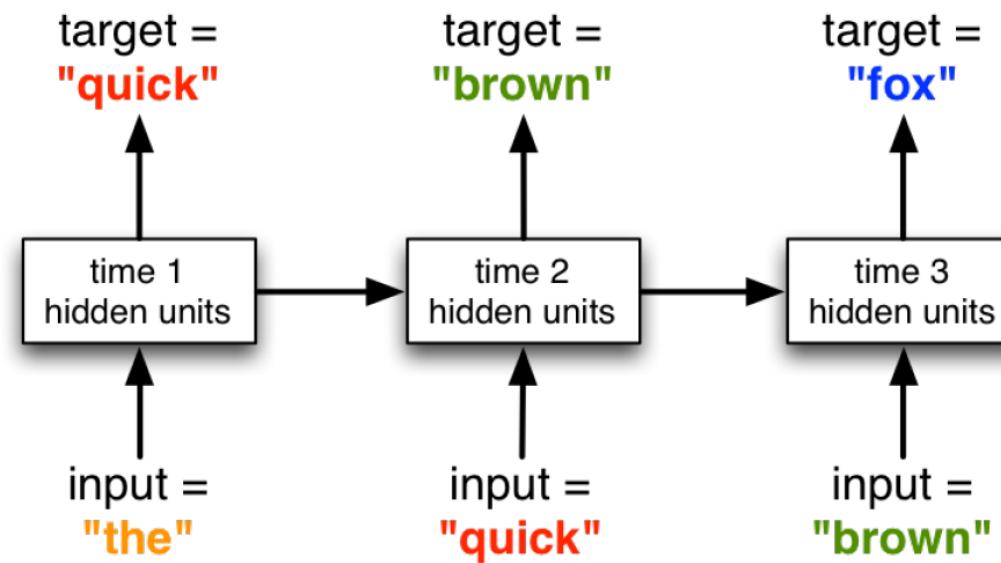
Vocabulary:
[h,e,l,o]

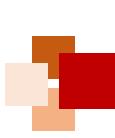
At test-time sample
characters one at a time,
feed back to model



RNN for language modeling

- Modeling at word level
 - Each word is represented as an indicator vector
 - The model predicts a distribution over words

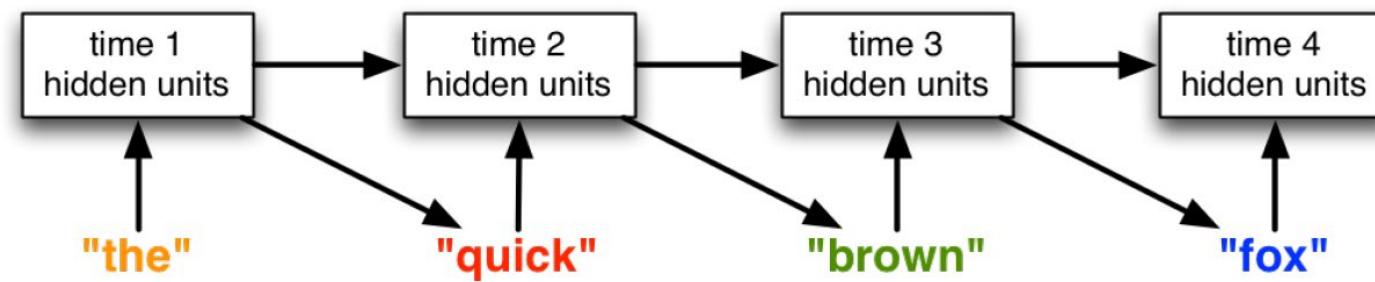




RNN for language modeling

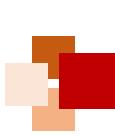


- Generating from a RNN language model
 - The outputs are fed back to the network



- Training time: the inputs are the token from the training set (**teacher forcing**).





RNN for language modeling



at first:

tyntd-iafhatawiaoihrdemot lytdws e ,tfti, astai f ogoh eoase rrranbyne 'nhthnee e
plia tkldrgd t o idoe ns,smtt h ne etie h,hregtrs nigtike,aoaenns lng

↓ train more

"Tmont thithey" fomesscerliund
Keushey. Thom here
sheulke, anmerenith ol sivh I lalterthend Bleipile shuw fil on aseterlome
coaniogennc Phe lism thond hon at. MeiDimorotion in ther thize."

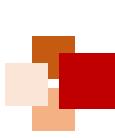
↓ train more

Aftair fall unsuch that the hall for Prince Velzonski's that me of
her hearly, and behs to so arwage fiving were to it beloge, pavu say falling misfort
how, and Gogition is so overelical and ofter.

↓ train more

"Why do what that day," replied Natasha, and wishing to himself the fact the
princess, Princess Mary was easier, fed in had oftened him.
Pierre aking his soul came to the packs and drove up his father-in-law women.

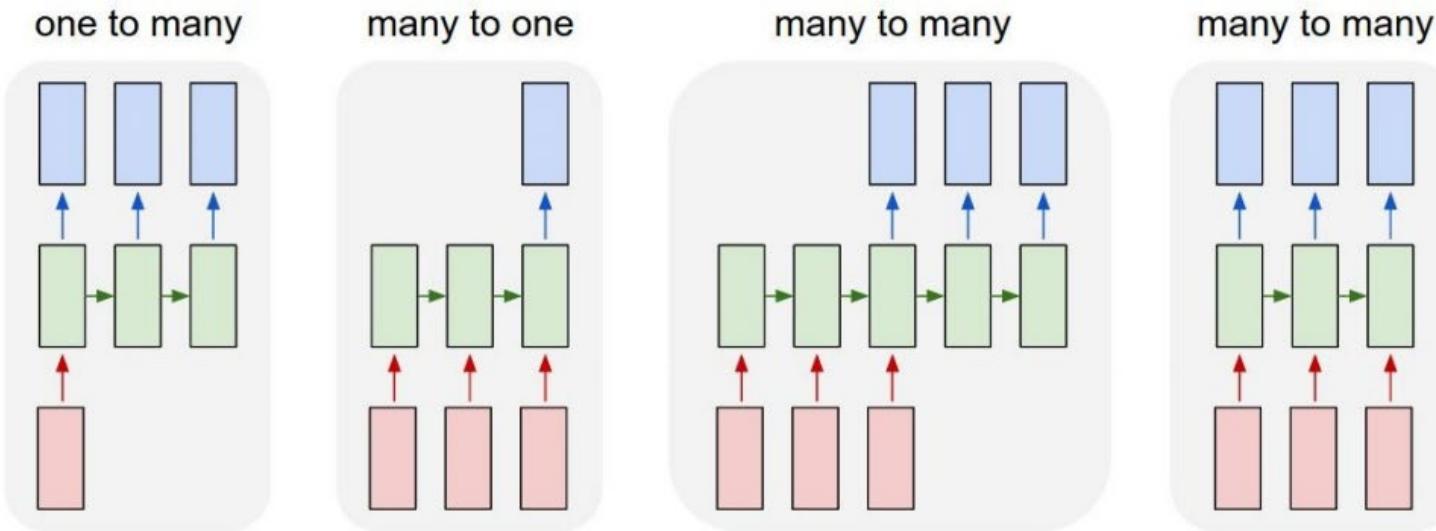




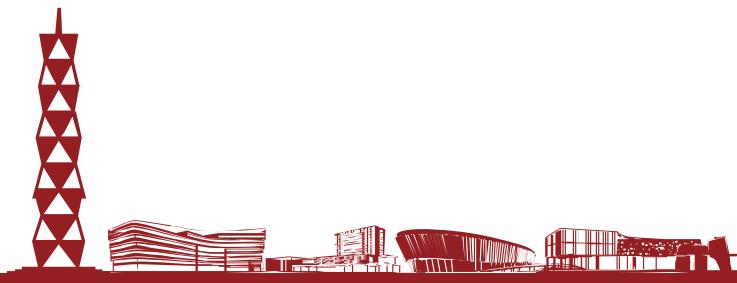
Recurrent Neural Network



- Recurrent Neural Networks: model variants



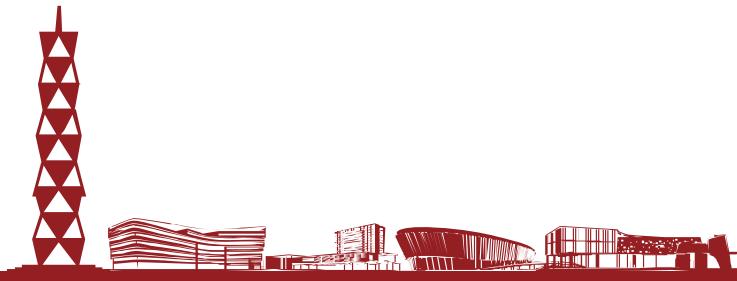
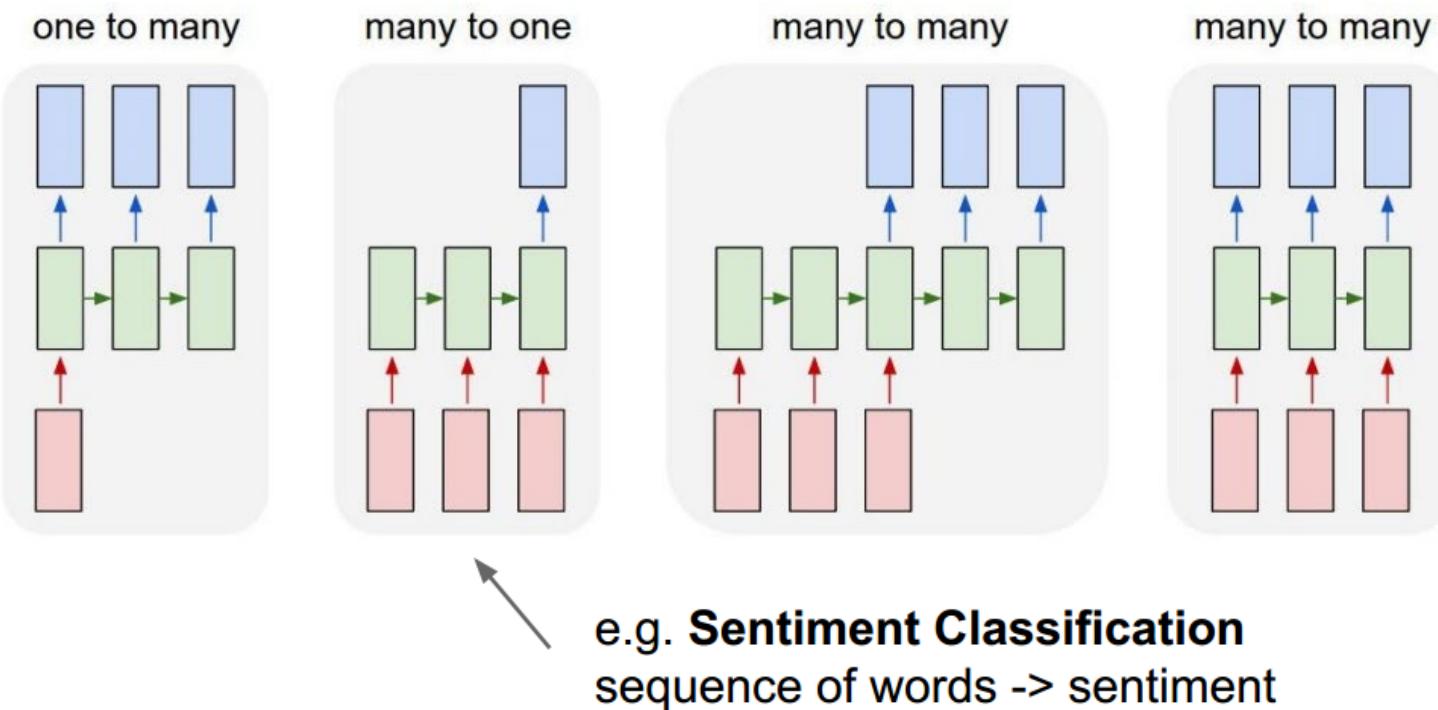
e.g. **Image Captioning**
image -> sequence of words



Recurrent Neural Network

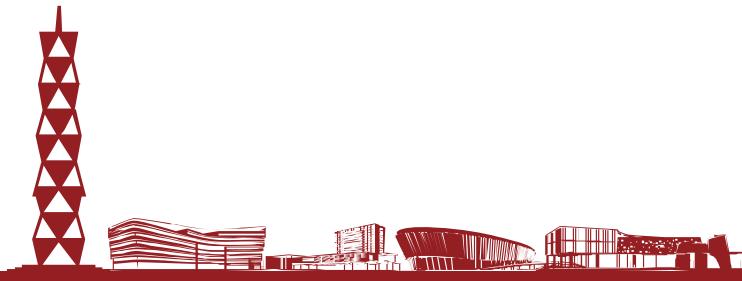
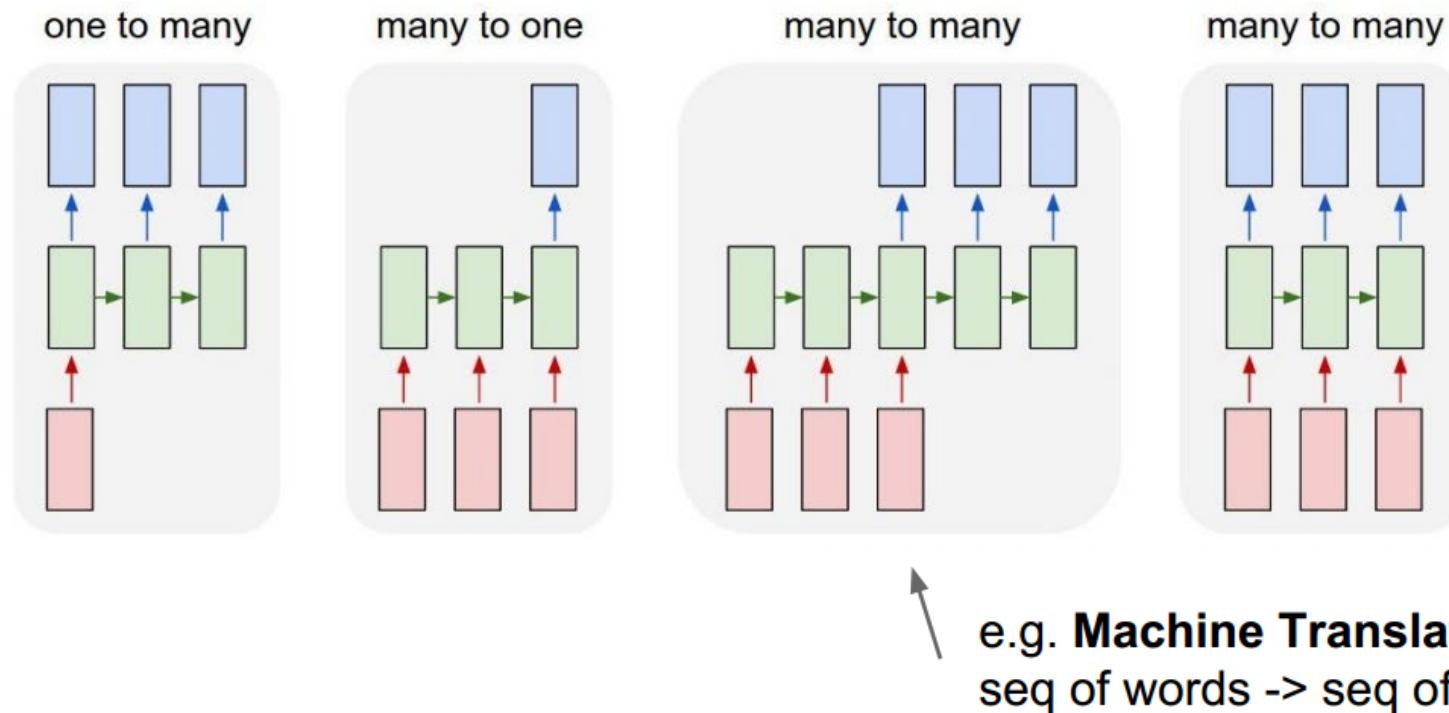


- Recurrent Neural Networks: model variants



Recurrent Neural Network

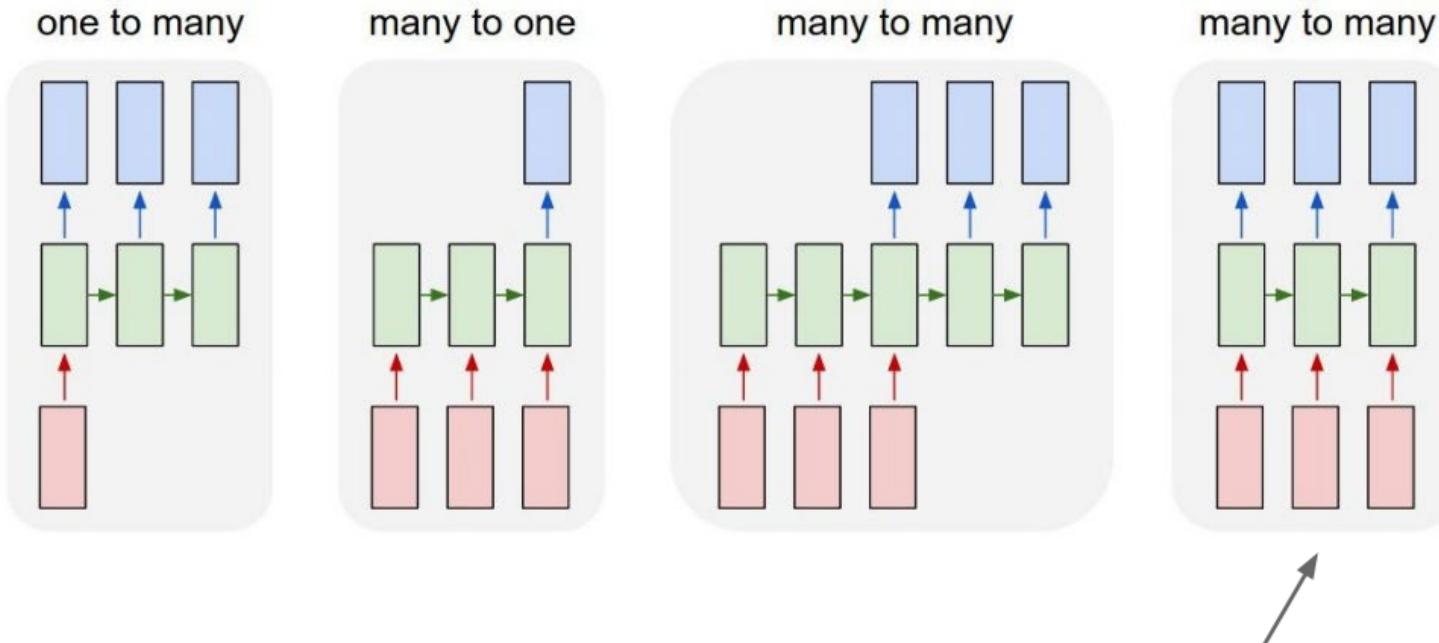
- Recurrent Neural Networks: model variants



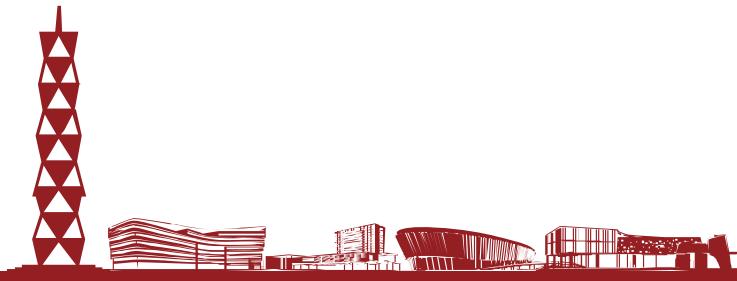
Recurrent Neural Network

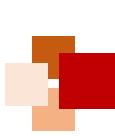


- Recurrent Neural Networks: model variants



e.g. Video classification on frame level



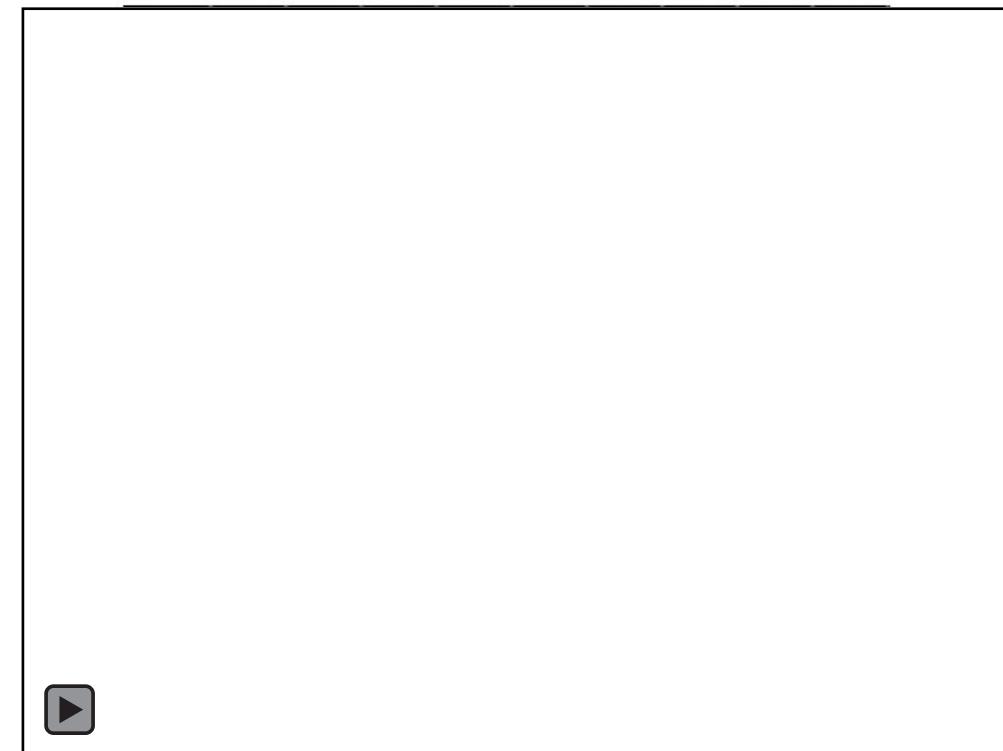


Recurrent Neural Network

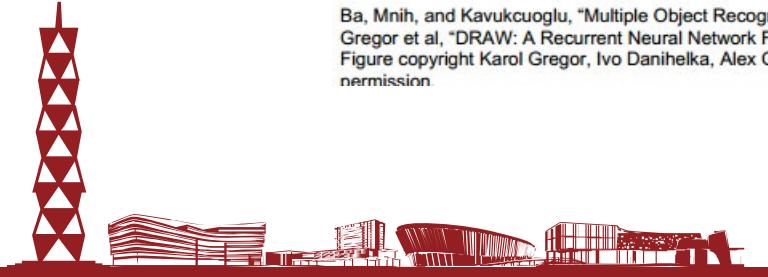


- Sequential Processing of Non-Sequence Data

Classify images by taking a series of “glimpses”



Ba, Mnih, and Kavukcuoglu, "Multiple Object Recognition with Visual Attention", ICLR 2015.
Gregor et al, "DRAW: A Recurrent Neural Network For Image Generation", ICML 2015
Figure copyright Karol Gregor, Ivo Danihelka, Alex Graves, Danilo Jimenez Rezende, and Daan Wierstra, 2015. Reproduced with permission



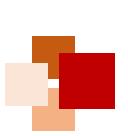


Outline



- Training Recurrent Neural Networks
 - Backpropagation through time
 - Gradient problems in training RNNs

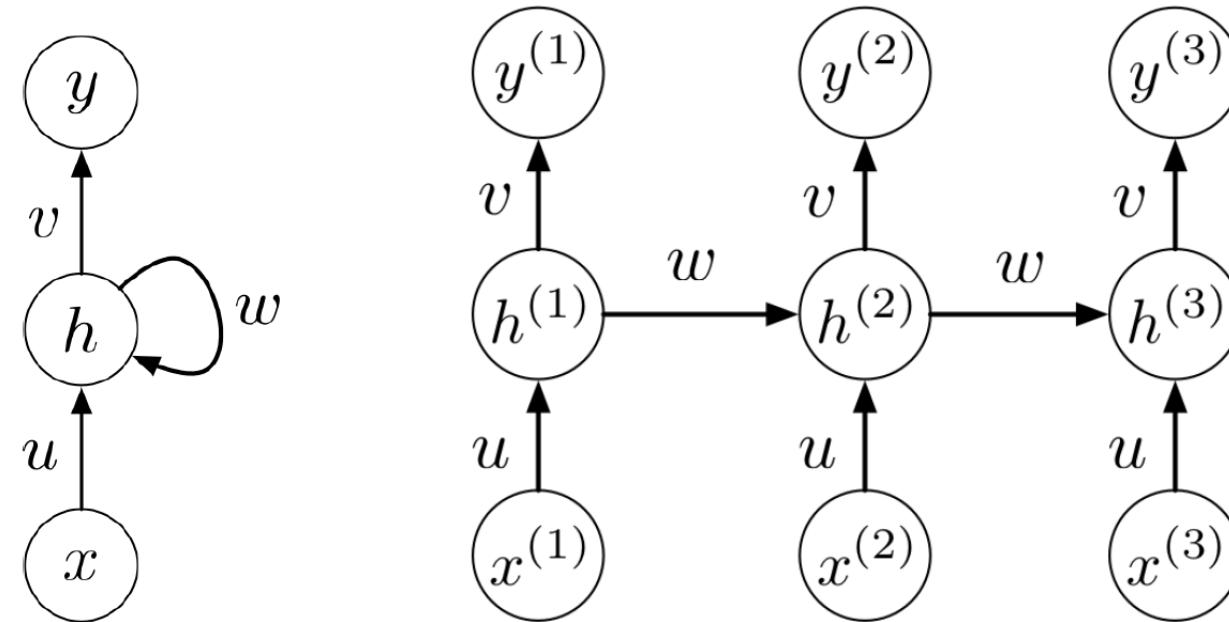




BPTT example



- A simple network
 - Everything is scalar
 - Unrolled computation graph with shared parameters

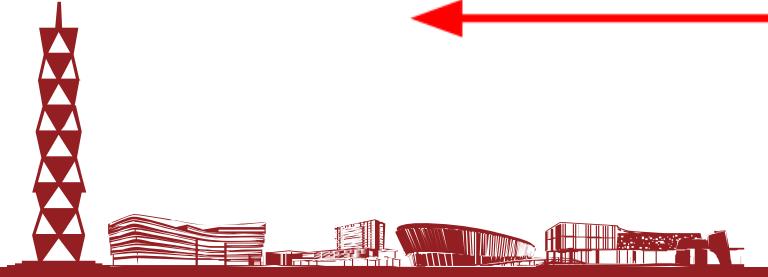
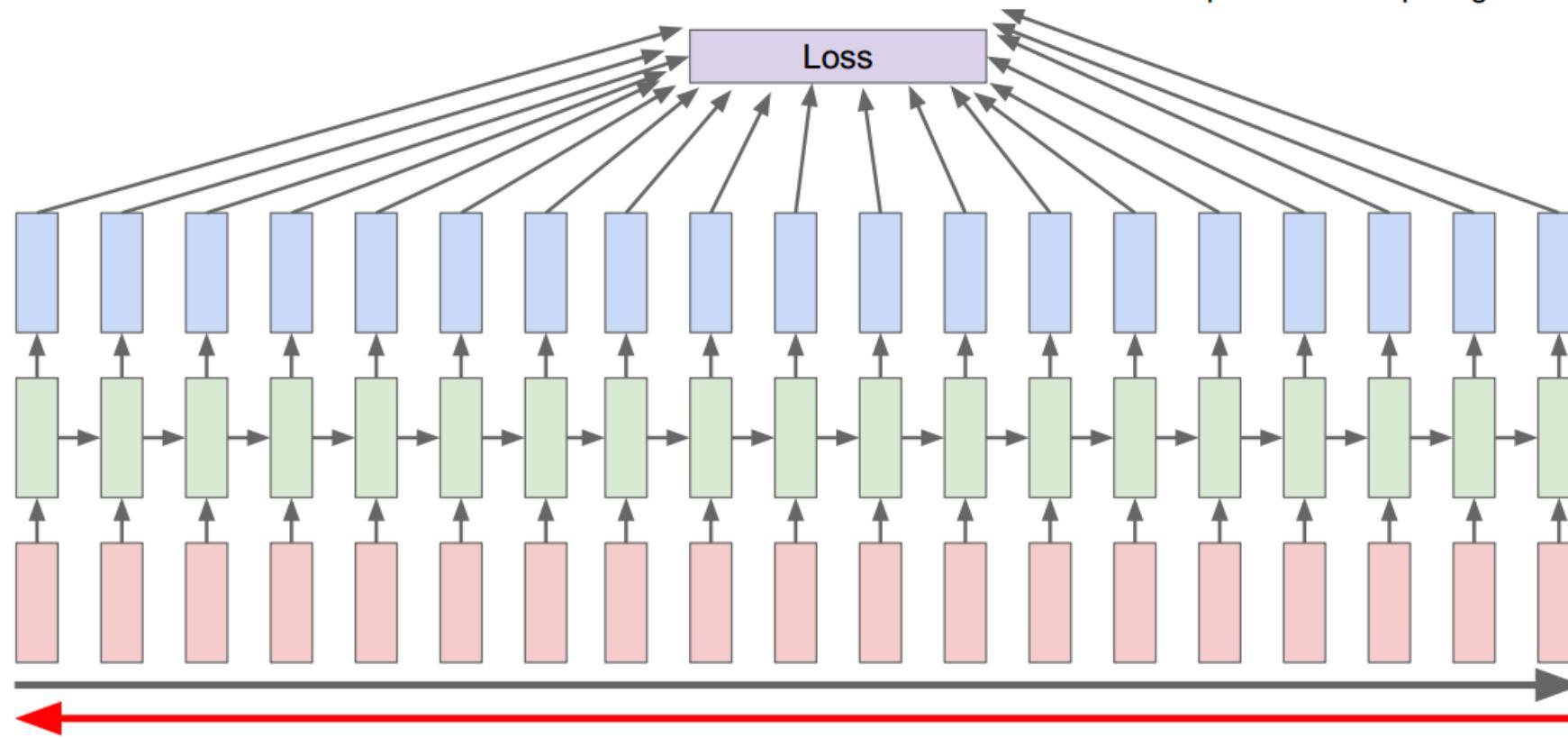


Recurrent Neural Network



Backpropagation through time

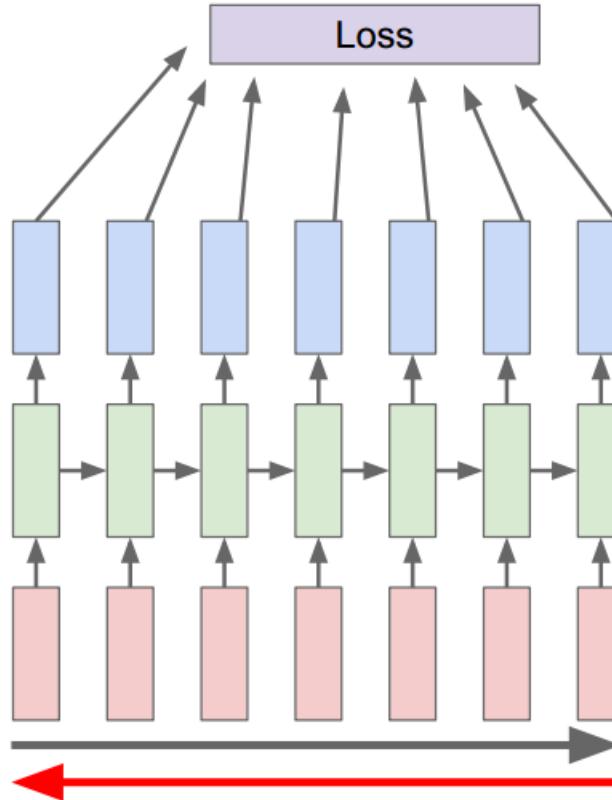
Forward through entire sequence to compute loss, then backward through entire sequence to compute gradient



Recurrent Neural Network



Truncated Backpropagation through time



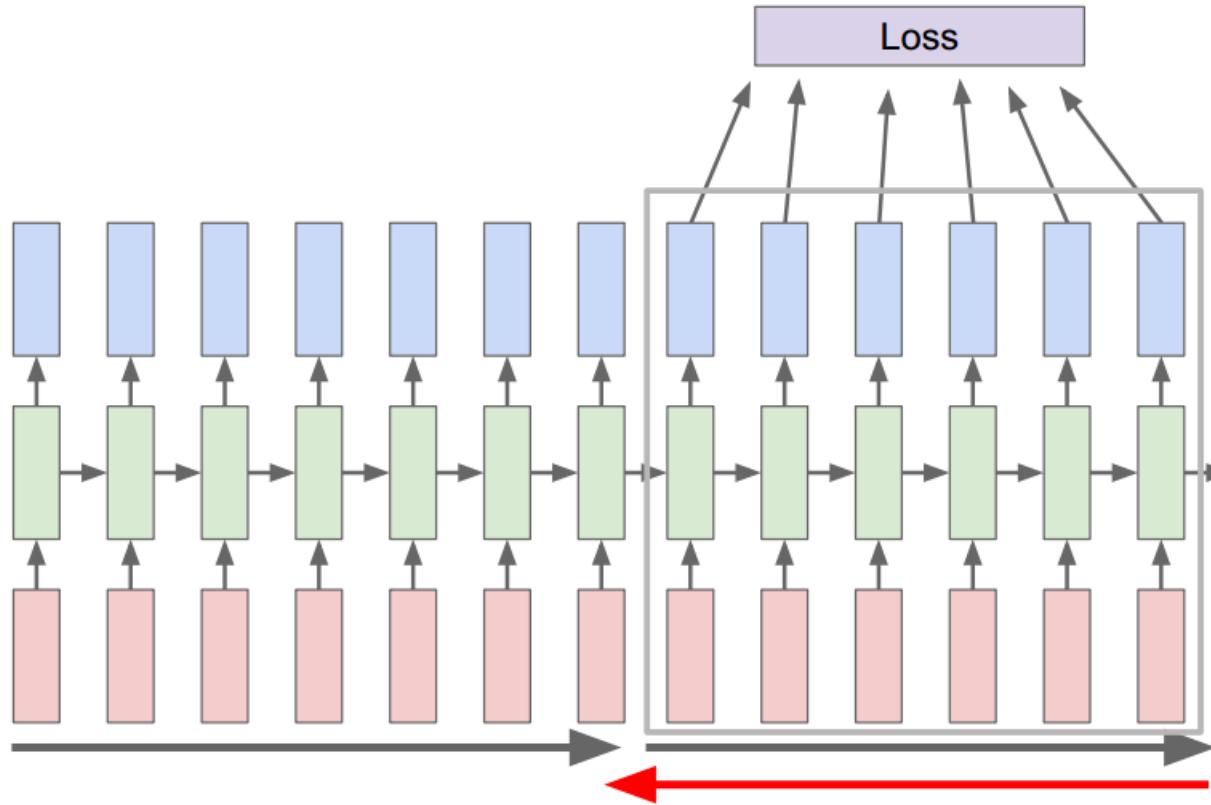
Run forward and backward
through chunks of the
sequence instead of whole
sequence



Recurrent Neural Network



Truncated Backpropagation through time



Carry hidden states
forward in time forever,
but only backpropagate
for some smaller
number of steps



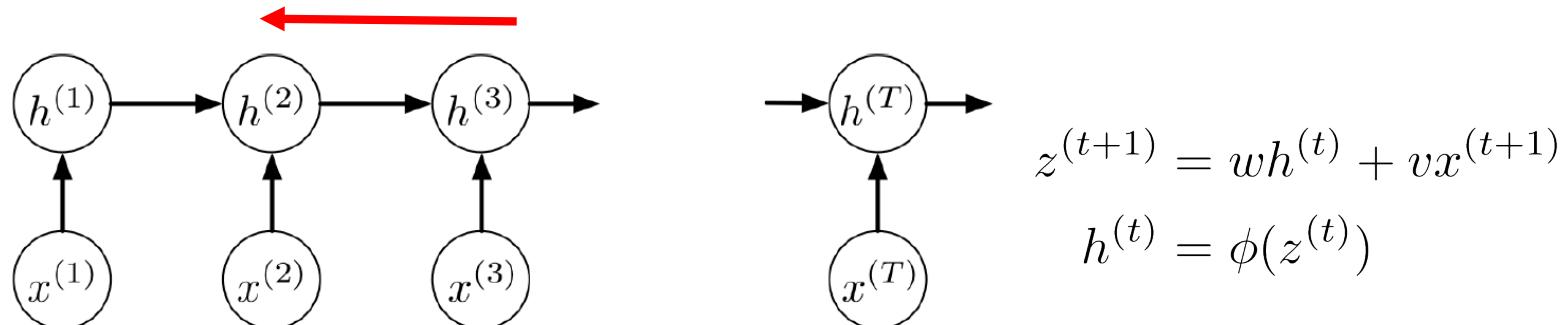
Challenges in training RNNs

- RNN
 - BP through time is used to compute the gradient descent update
- Problems
 - The updates are mathematically correct, but gradient descent fails because the gradients explode or vanish
 - This limits the scope of the dependencies over time



Why gradients explode or vanish

- Motivating example:
 - Consider a univariate version of the vanilla RNNs



Backprop updates:

$$\overline{h^{(t)}} = \overline{z^{(t+1)}} w$$

$$\overline{z^{(t)}} = \overline{h^{(t)}} \phi'(z^{(t)})$$

Applying this recursively:

$$\overline{h^{(1)}} = \underbrace{w^{T-1} \phi'(z^{(2)}) \cdots \phi'(z^{(T)})}_{\text{the Jacobian } \partial h^{(T)} / \partial h^{(1)}} \overline{h^{(T)}}$$

With linear activations:

$$\partial h^{(T)} / \partial h^{(1)} = w^{T-1}$$

Exploding:

$$w = 1.1, T = 50 \Rightarrow \frac{\partial h^{(T)}}{\partial h^{(1)}} = 117.4$$

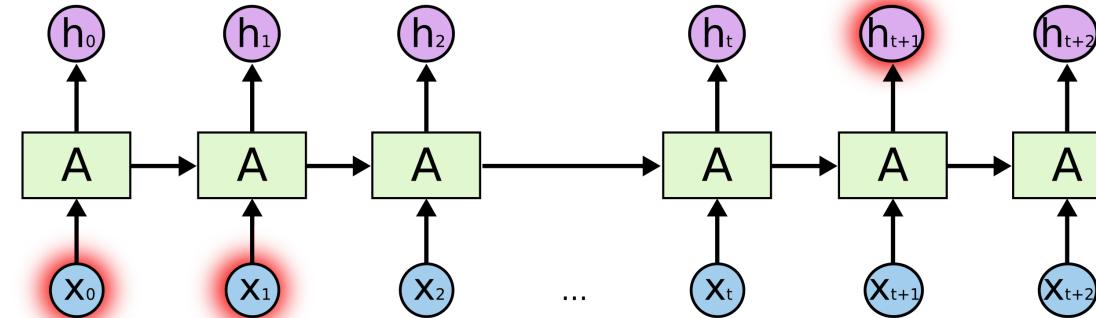
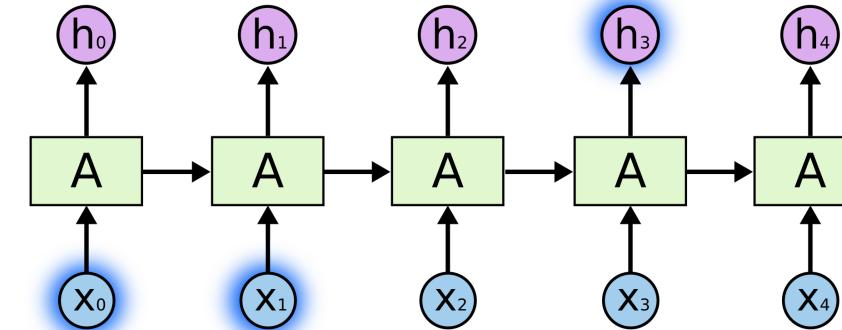
Vanishing:

$$w = 0.9, T = 50 \Rightarrow \frac{\partial h^{(T)}}{\partial h^{(1)}} = 0.00515$$



Vanilla RNN

- Difficulty in modeling long-term dependency



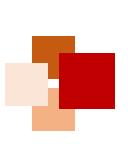


Outline



- LSTM
 - Long-Term Short Term Memory (LSTM)
 - RNNs with LSTM





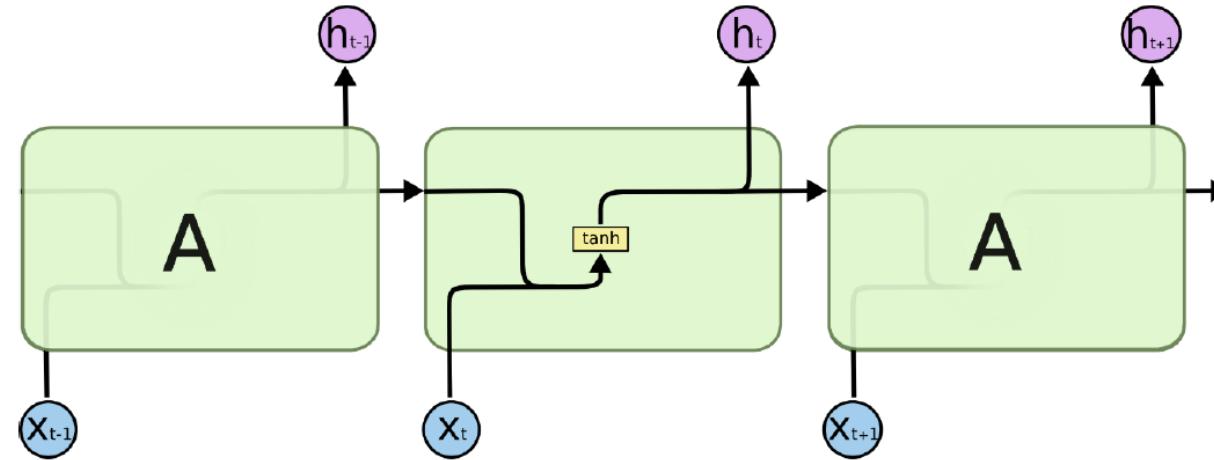
Long-term Short Term Memory

- Replacing a vanilla RNN neuron by the **LSTM unit**
- Why it is called LSTM
 - A network's activations are its **short-term memory** and its weights are its **long-term memory**
 - The LSTM architecture wants the short-term memory to last for a long time period
- Key idea
 - Composed of **memory cells** which have controllers that decide when to store or forget information



Standard RNN

- Recall

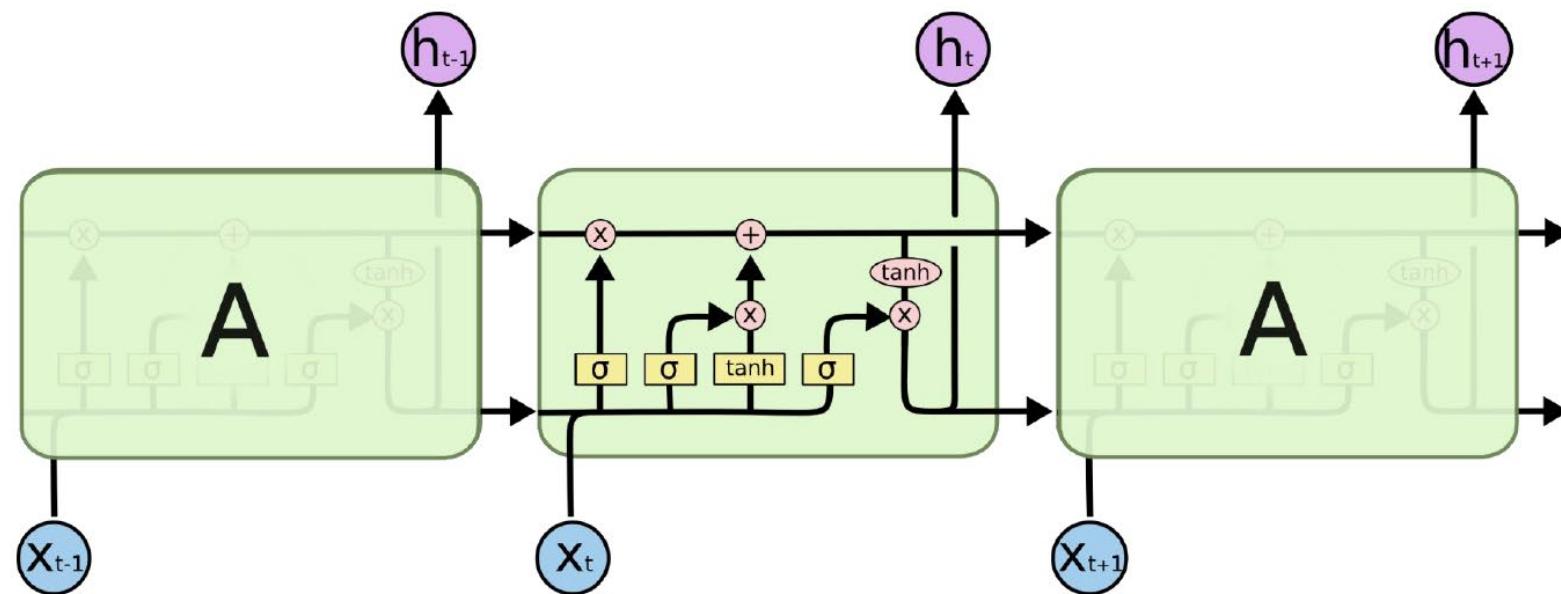


- Each recurrent neuron receives past outputs and current input
- Pass through a tanh function



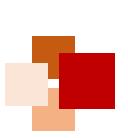
Long Short Term Memory(LSTM)

- LSTM uses multiplicative gates that decide if something is important or not



Hochreiter and Schmidhuber, "Long Short Term Memory", Neural Computation

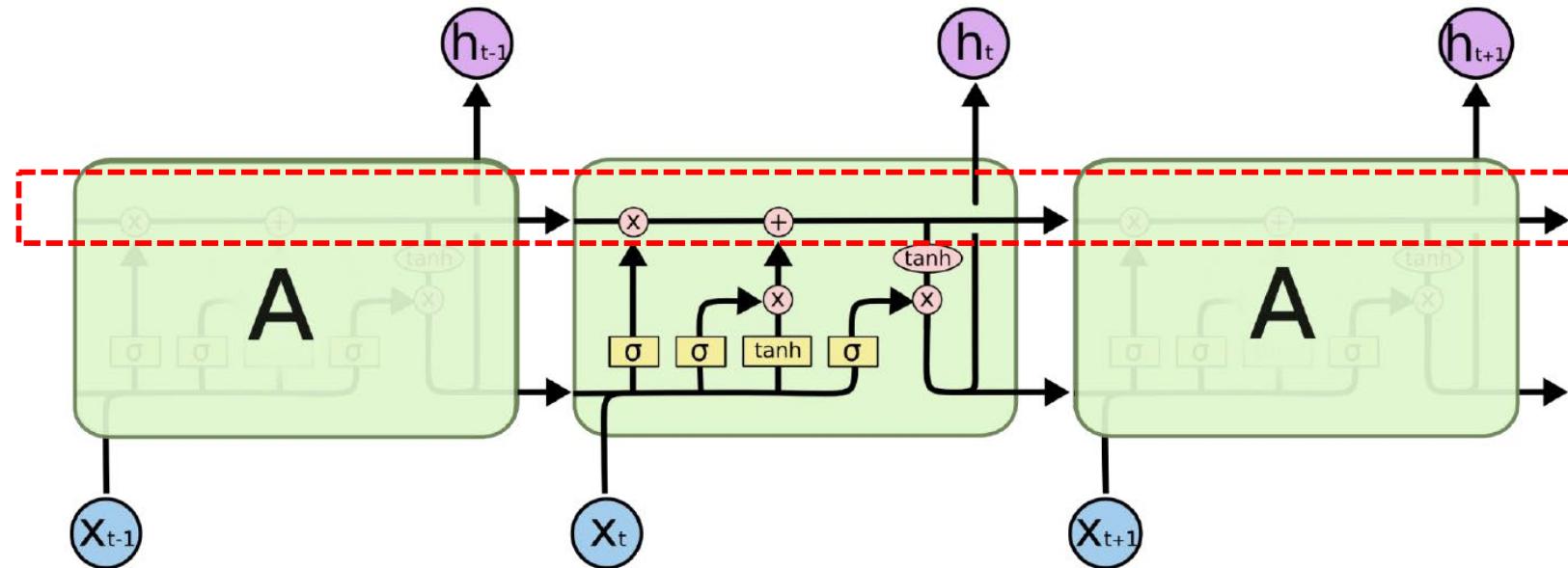




Long Short Term Memory(LSTM)

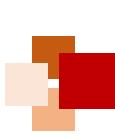


- Key component: a remembered cell state



Hochreiter and Schmidhuber, "Long Short Term Memory", Neural Computation

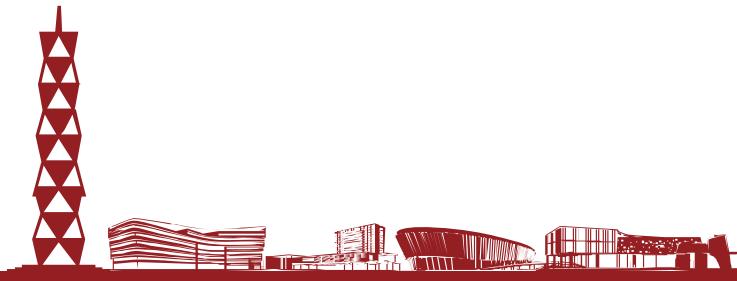
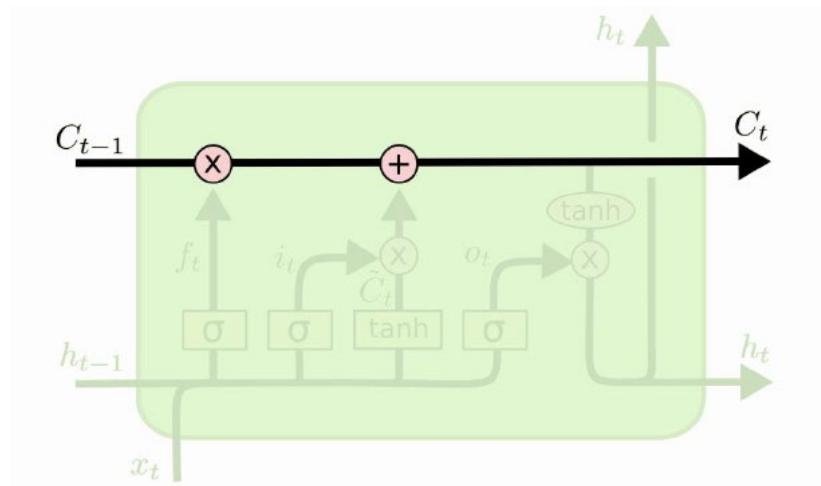


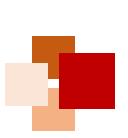


LSTM: cell state



- A linear history
 - Carries information through
 - Only affected by a gate and addition of current information, which is also gated

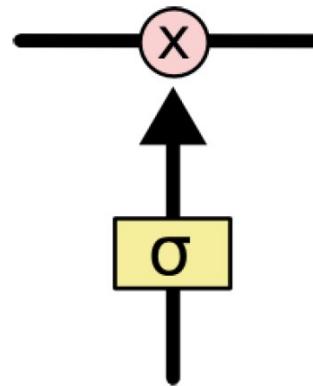




LSTM: gates



- Gates are simple sigmoid units with output range in (0,1)
- Controls how much of the information will be let through



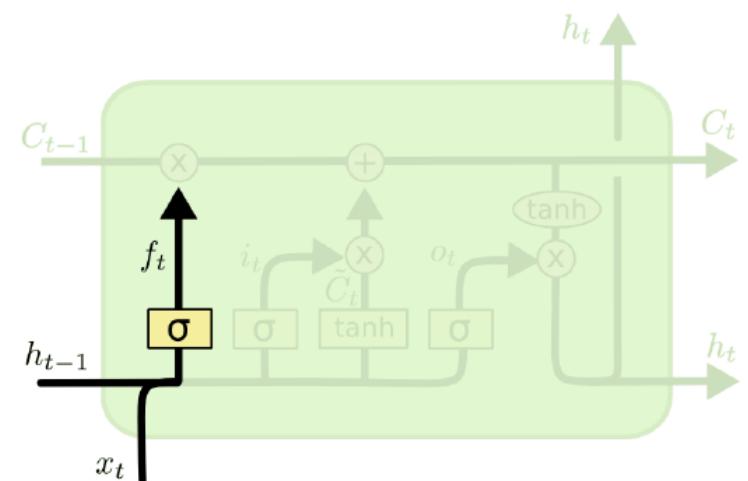
- Three gates
 - Forget gate
 - Input gate
 - Output gate



LSTM: forget gate

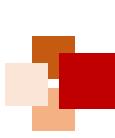


- The first gate determines whether to carry over the history or to forget it
 - Soft decision: how much of the history C_{t-1} to carry over
 - Determined by the current input x_t and the previous state h_{t-1}



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

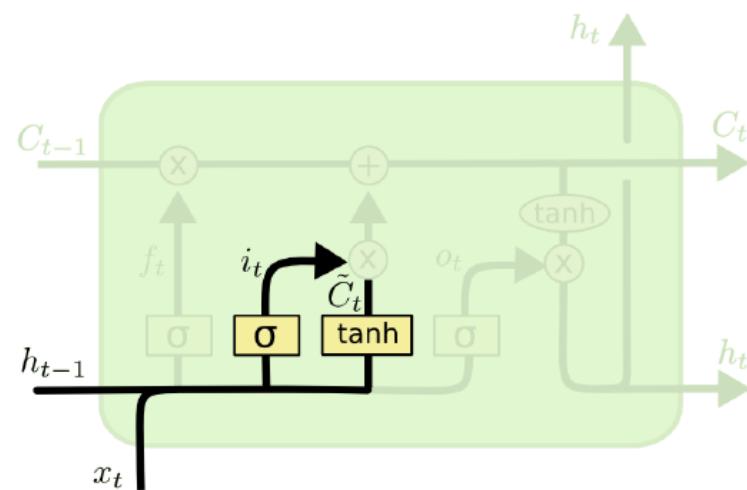




LSTM: input gate



- The second gate has two parts
 - A gate that decides if it is worth remembering
 - A nonlinear transformation that extracts new and interesting information from the input
 - Both use the current input and the previous state



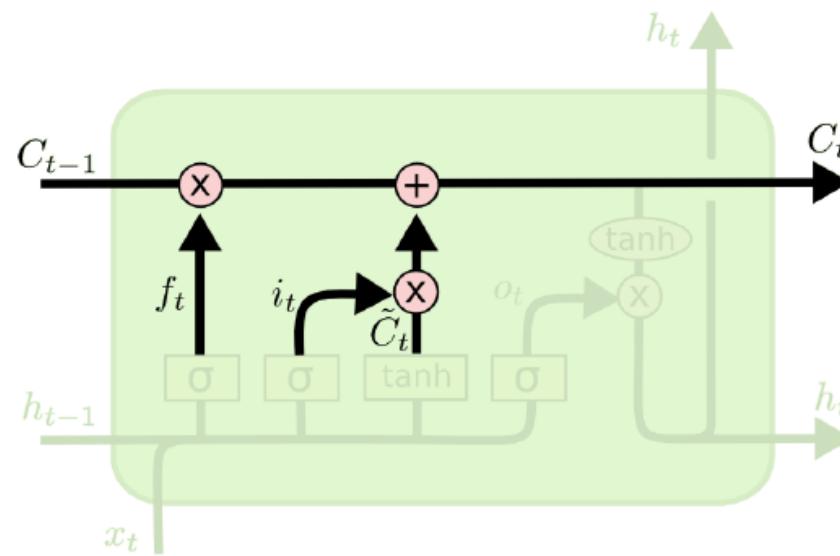
$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$



LSTM: Memory cell update

- The output of the second part is added into the current memory cell
 - The input and state jointly decide how much of history info is kept and how much of embedded input info is added
 - A dynamic mixture at each time step

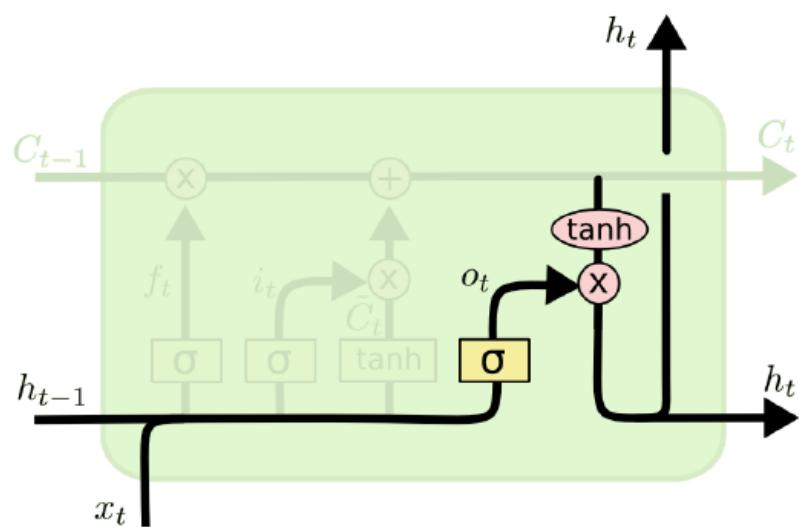


$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

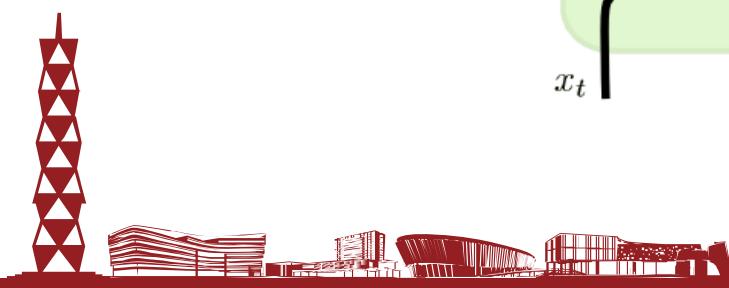


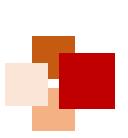
LSTM: Output gate

- The third gate is the output gate
 - To decide if the memory cell contents are worth reporting at this time using the current input and previous state
- The output of the cell or the state
 - A nonlinear transform of the cell values
 - Compress it with tanh to make it in (-1,1)



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

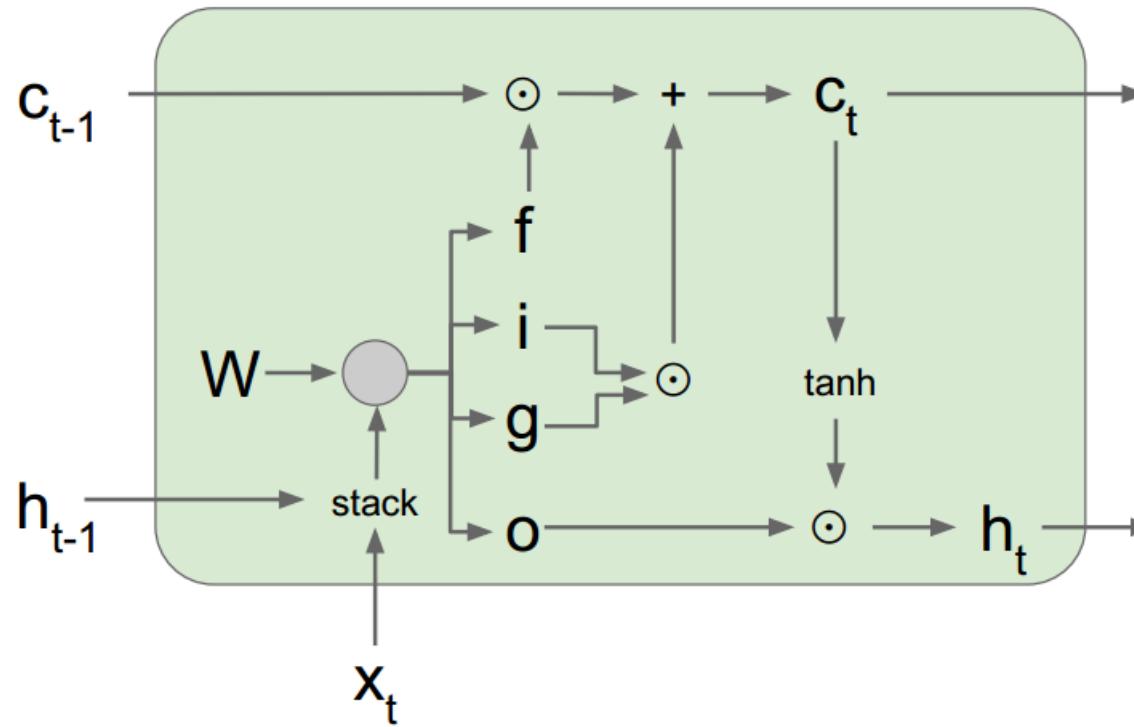




Long Short Term Memory(LSTM)



[Hochreiter et al., 1997]

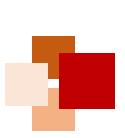


$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

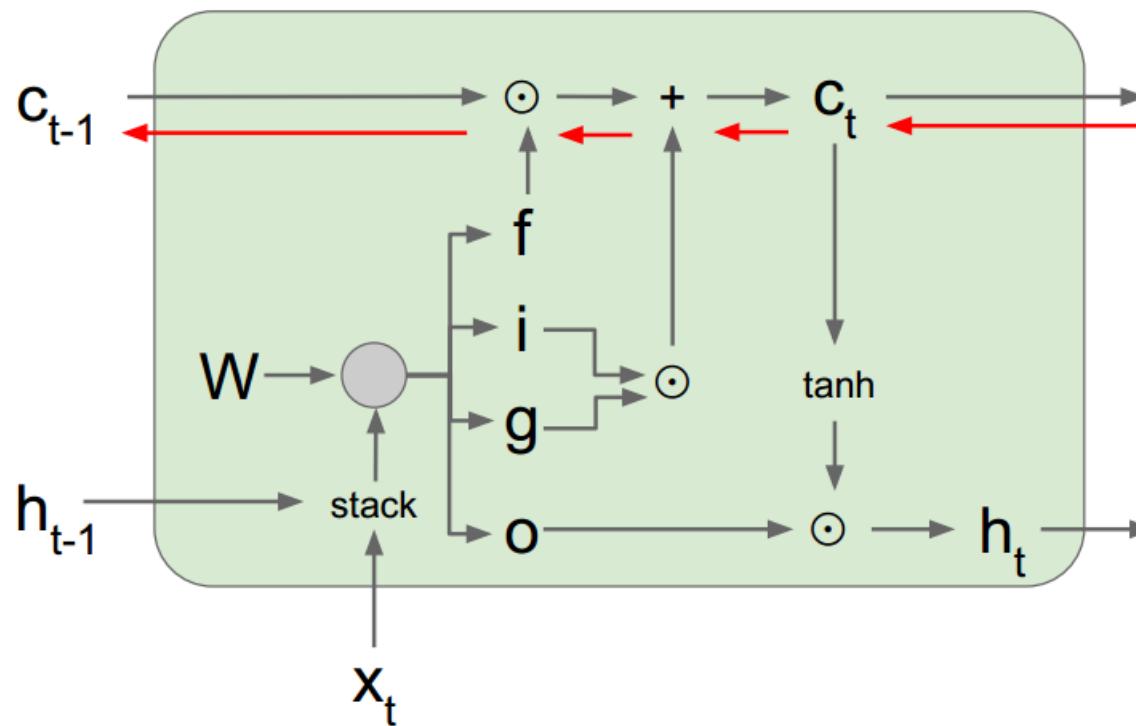




LSTM: Backpropagation



[Hochreiter et al., 1997]



Backpropagation from c_t to
 c_{t-1} only elementwise
multiplication by f , no matrix
multiply by W

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

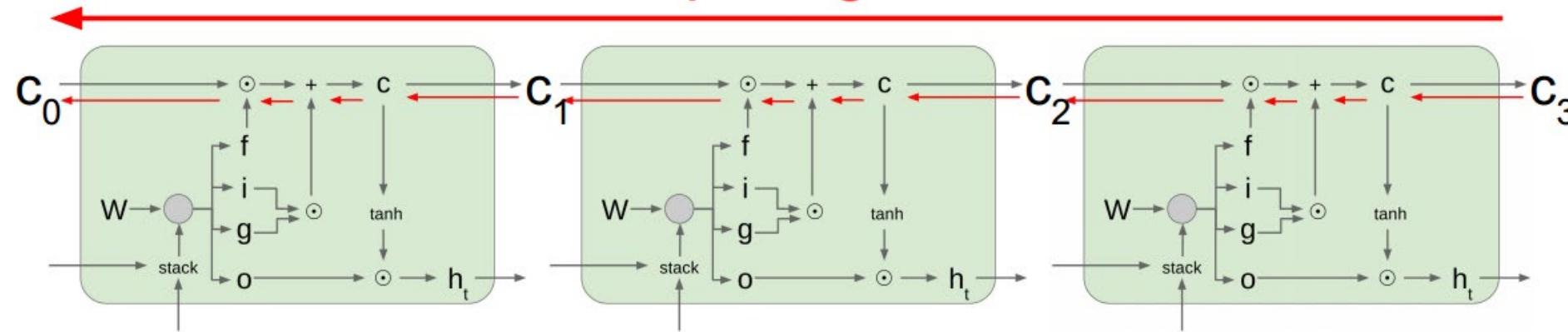
$$h_t = o \odot \tanh(c_t)$$

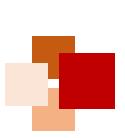


LSTM: Backpropagation



Uninterrupted gradient flow!

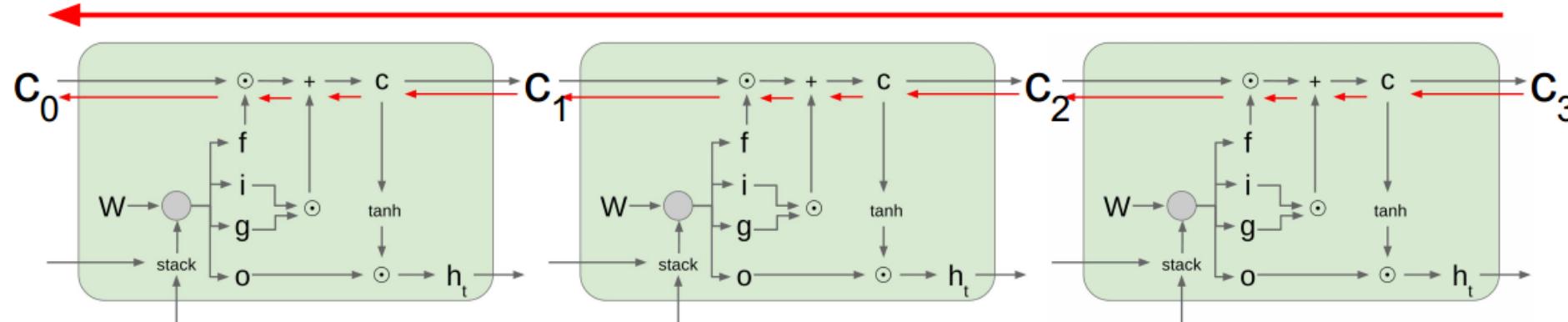




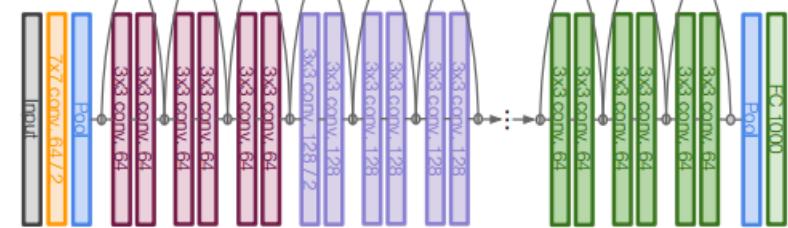
LSTM: Backpropagation



Uninterrupted gradient flow!



Similar to ResNet!

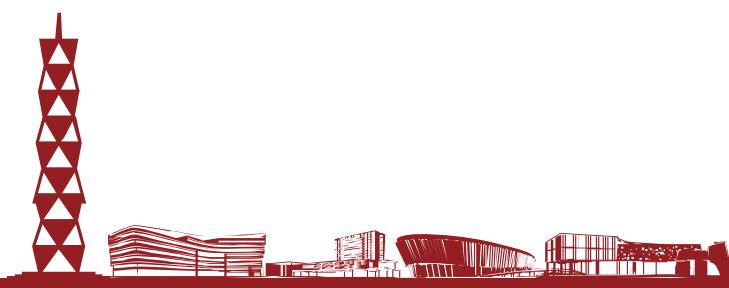


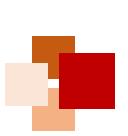
In between:
Highway Networks

$$g = T(x, W_T)$$

$$y = g \odot H(x, W_H) + (1 - g) \odot x$$

Srivastava et al, "Highway Networks",
ICML DL Workshop 2015





Multi-Layer RNNs



Multilayer RNNs

$$h_t^l = \tanh W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

$h \in \mathbb{R}^n$ $W^l [n \times 2n]$

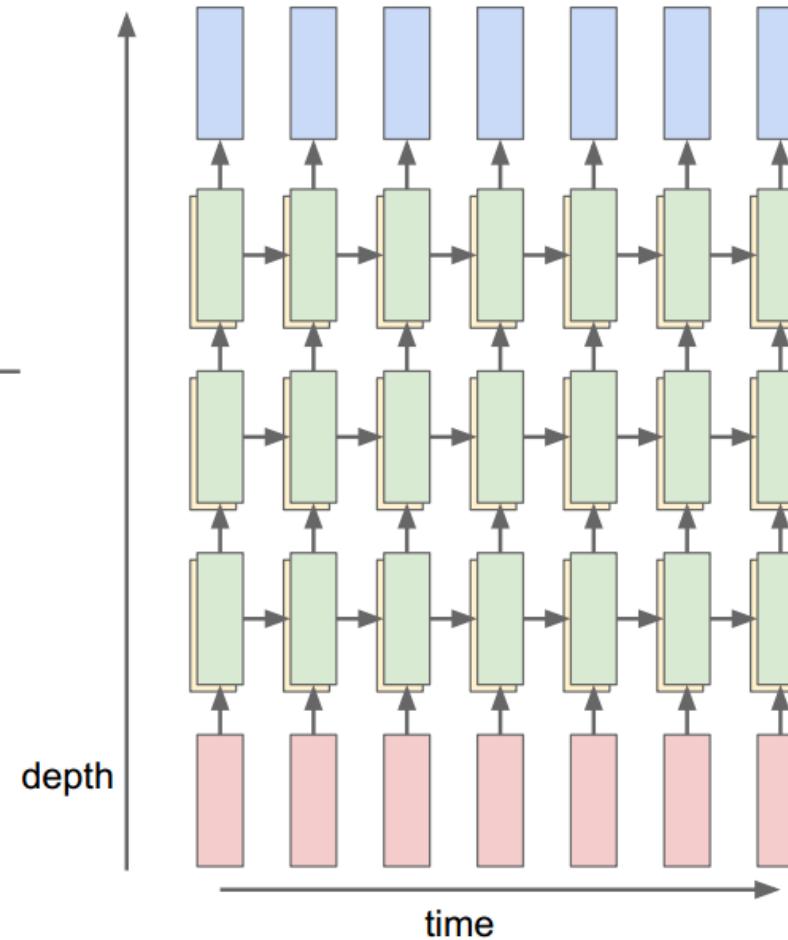
LSTM:

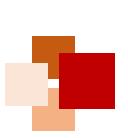
$$W^l [4n \times 2n]$$

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \tanh \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

$$c_t^l = f \odot c_{t-1}^l + i \odot g$$

$$h_t^l = o \odot \tanh(c_t^l)$$

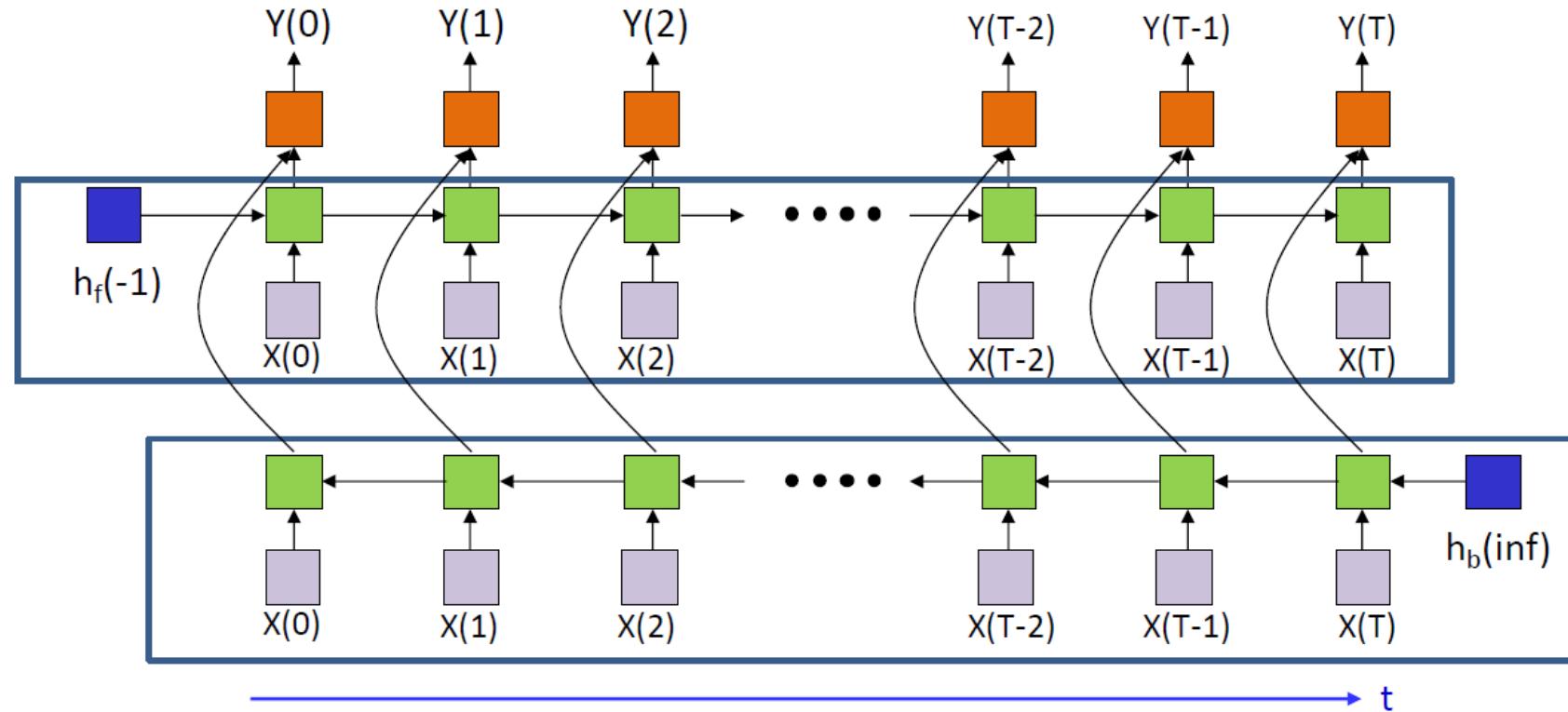




Bidirectional LSTM



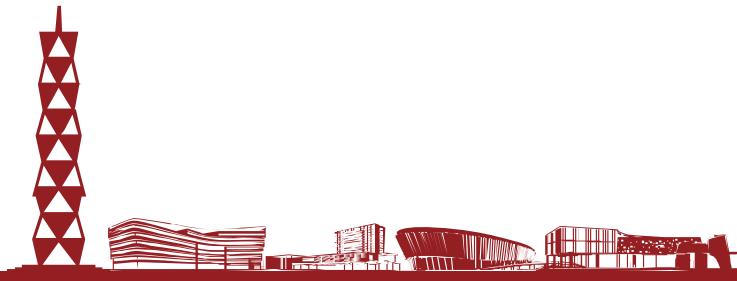
- Two opposite directions
 - Noncausal but complementary global context
 - Can have multiple layers of LSTM units in either direction





Outline

- RNN Applications
 - Sequence-to-sequence model: NMT
 - Image-to-sequence model: Image Captioning

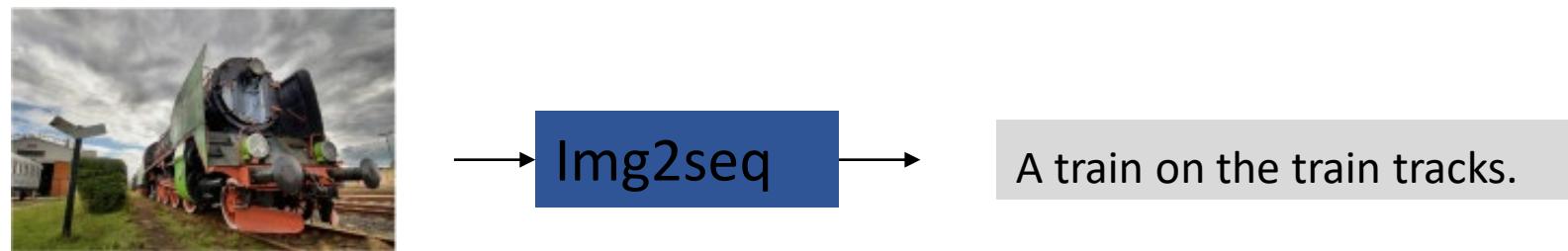


X-to-Sequence problems

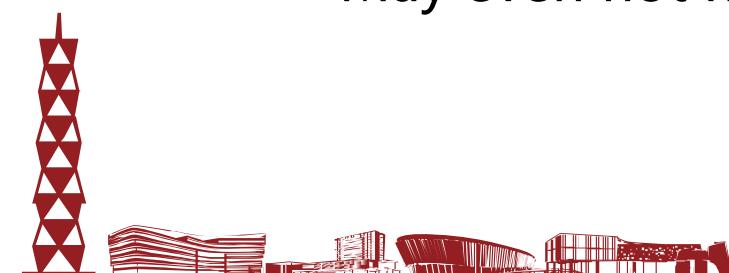
- Sequence or non-sequence in, sequence comes out
 - Machine translation



- Image caption generation



- No notion of “synchrony” between input and output
 - May even not have a notion of “alignment”



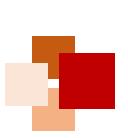
Sequence-to-sequence problem

- Task definition



- A sequence X_1, \dots, X_N goes in
- A different sequence Y_1, \dots, Y_M comes out
- Example: machine translation
 - The output is in a different language
- Example: dialog
 - “I have a problem” -> “How may I help you”

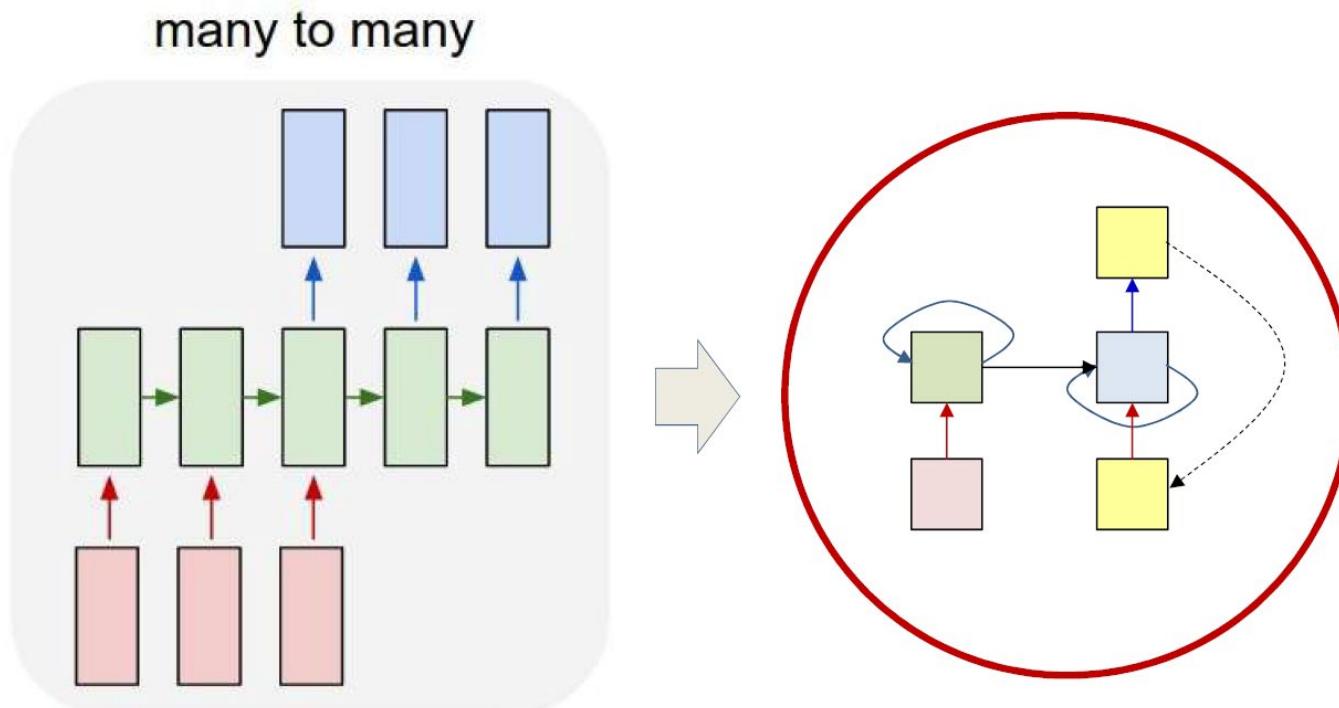


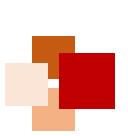


Modeling the problem



- Delayed sequence to sequence
 - Delayed self-referencing sequence-to-sequence

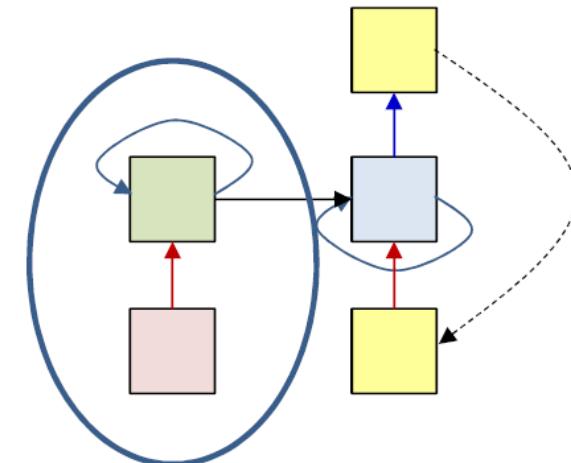
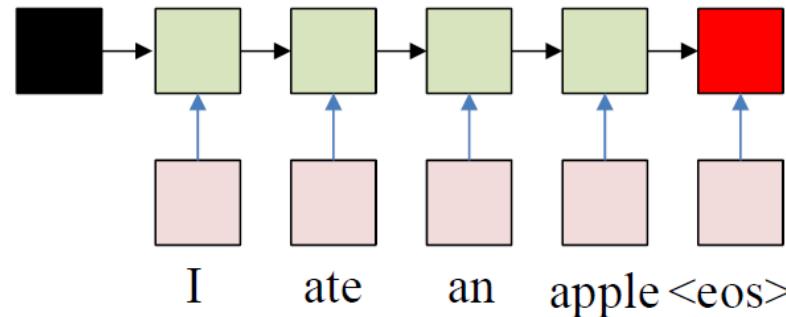


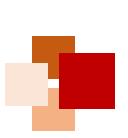


The “simple” translation model



- The input sequence feeds into an recurrent structure
 - The input sequence is terminated by an explicit <eos> symbol
 - The hidden activation at the <eos> “stores” all information about the sentence

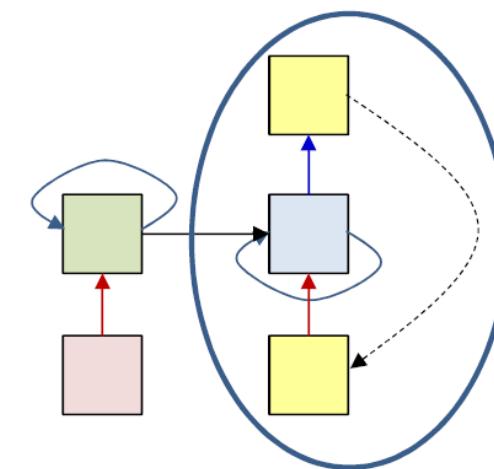
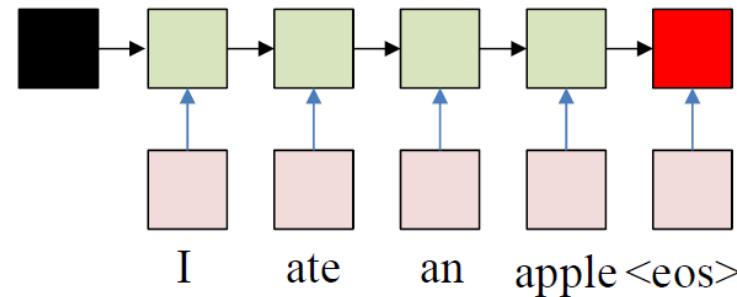


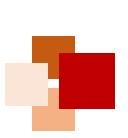


The “simple” translation model



- Subsequently a second RNN uses the hidden activation as initial state to produce a sequence of outputs
 - The output at each time becomes the input at the next time
 - Output production continues until an <eos> is produced

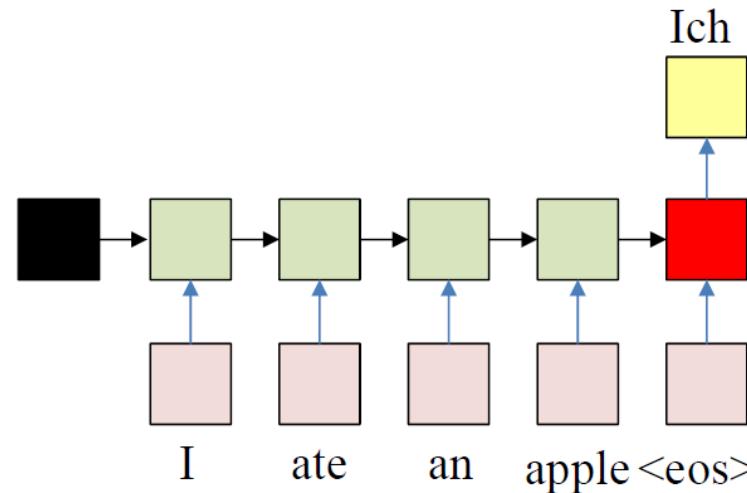


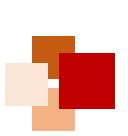


The “simple” translation model



- Subsequently a second RNN uses the hidden activation as initial state to produce a sequence of outputs
 - The output at each time becomes the input at the next time
 - Output production continues until an <eos> is produced

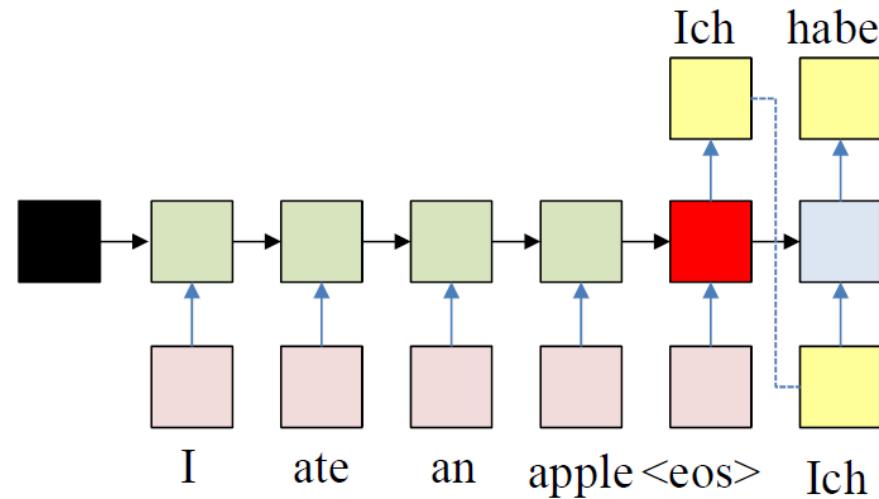


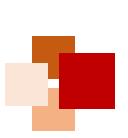


The “simple” translation model



- Subsequently a second RNN uses the hidden activation as initial state to produce a sequence of outputs
 - The output at each time becomes the input at the next time
 - Output production continues until an <eos> is produced

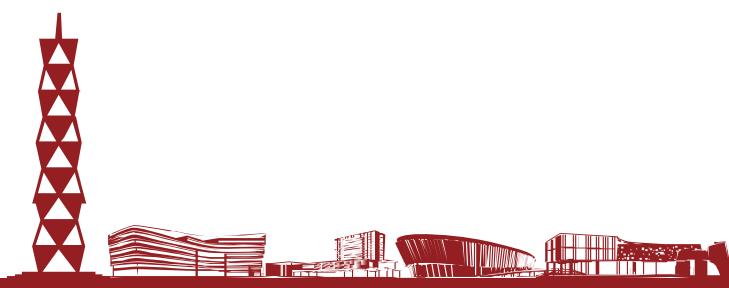
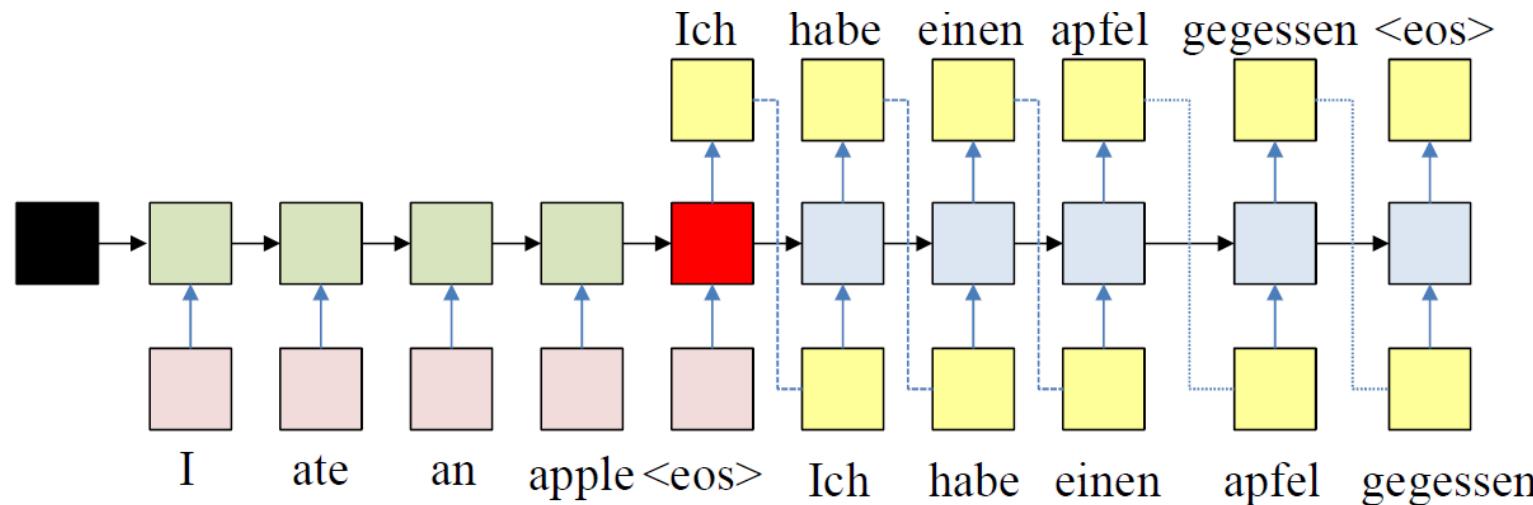


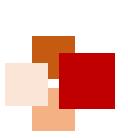


The “simple” translation model



- Subsequently a second RNN uses the hidden activation as initial state to produce a sequence of outputs
 - The output at each time becomes the input at the next time
 - Output production continues until an <eos> is produced

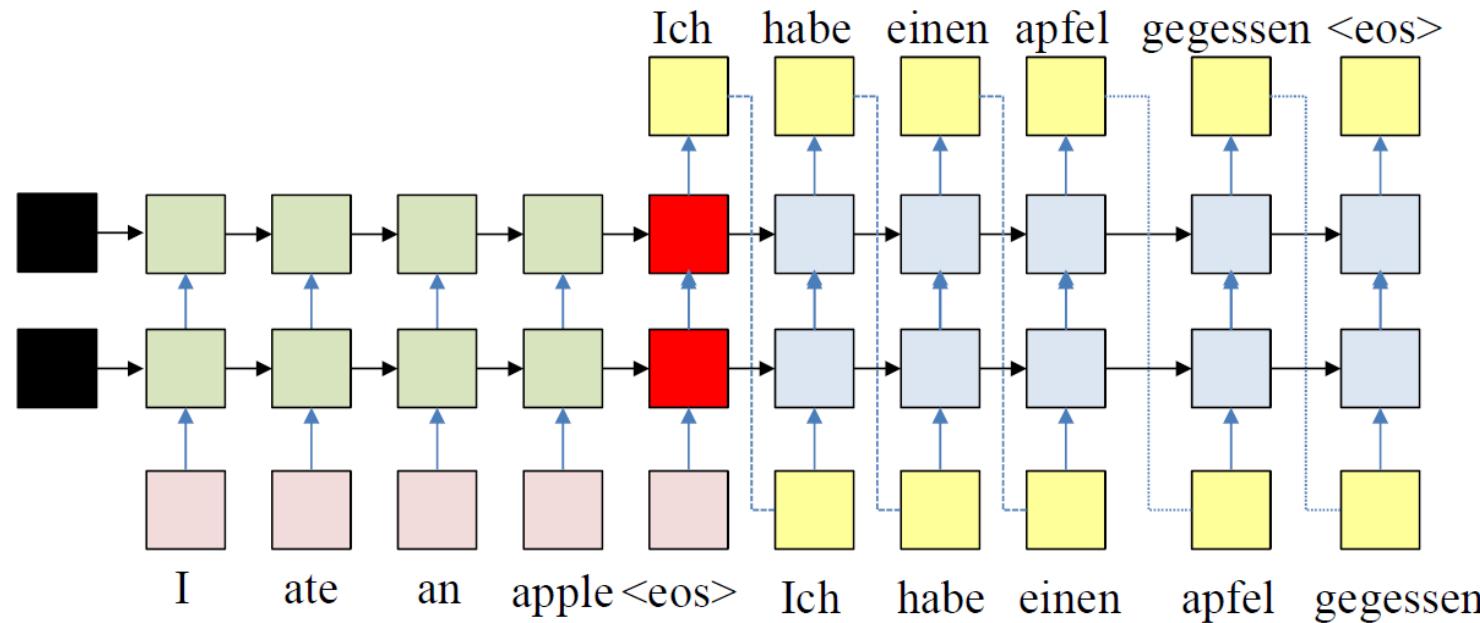


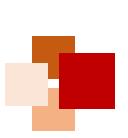


The “simple” translation model



- Such an architecture can be generalized to more layers

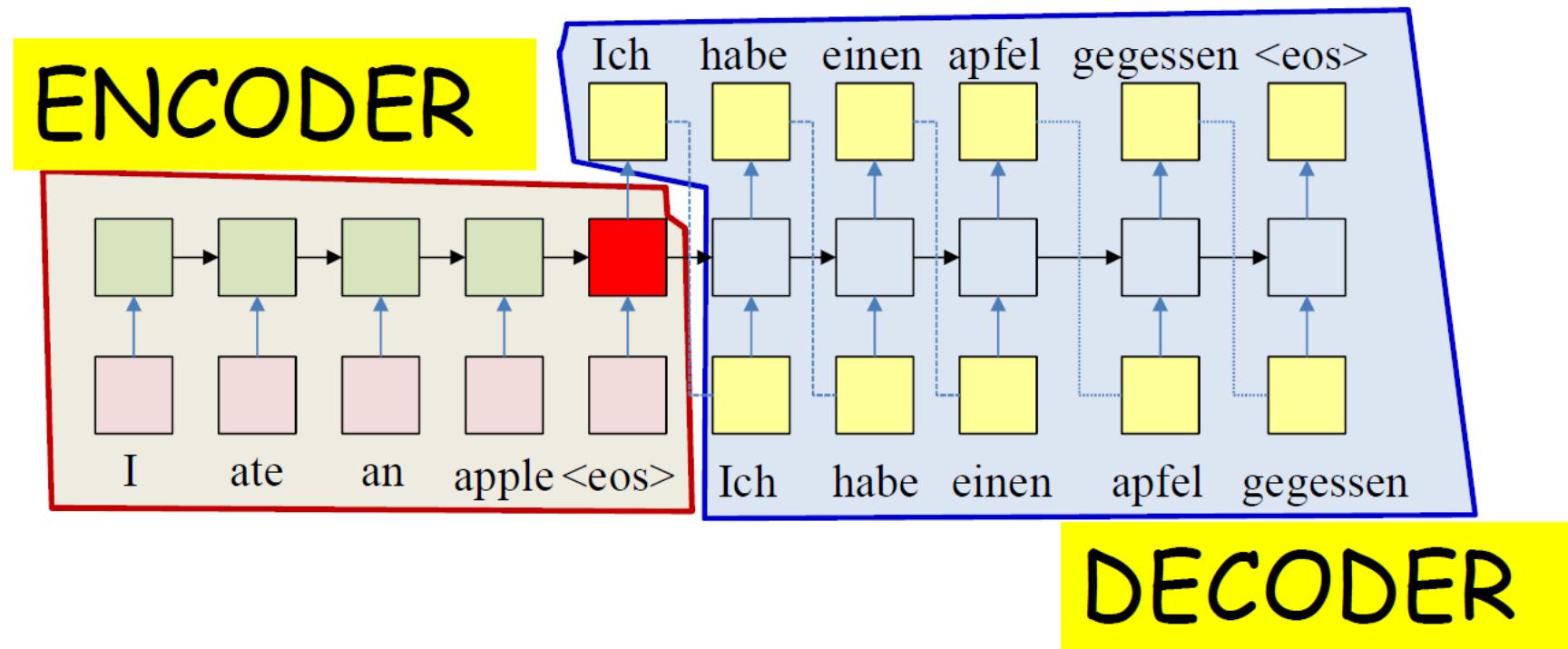


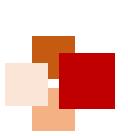


The “simple” translation model



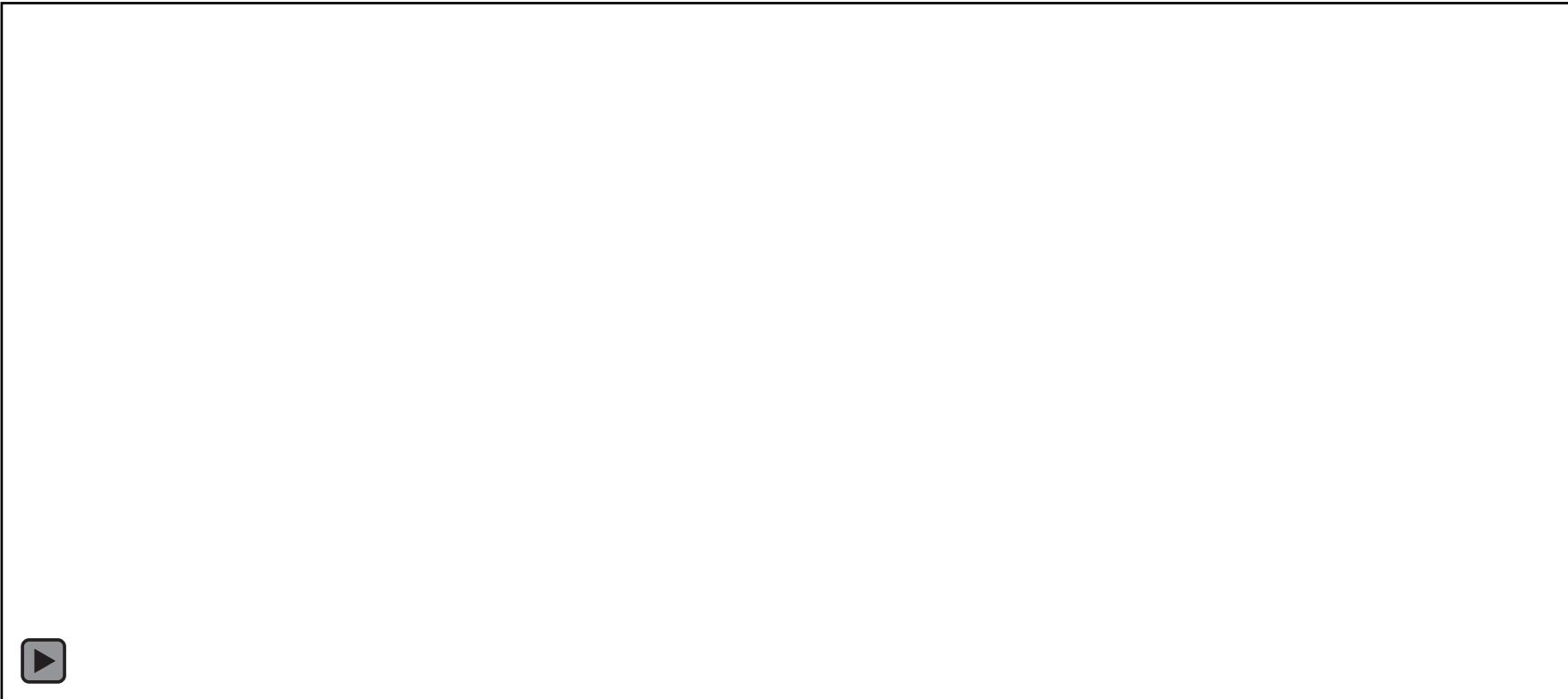
- This is also referred to as an **encoder-decoder** structure





The “simple” translation model

- This is also referred to as an **encoder-decoder** structure

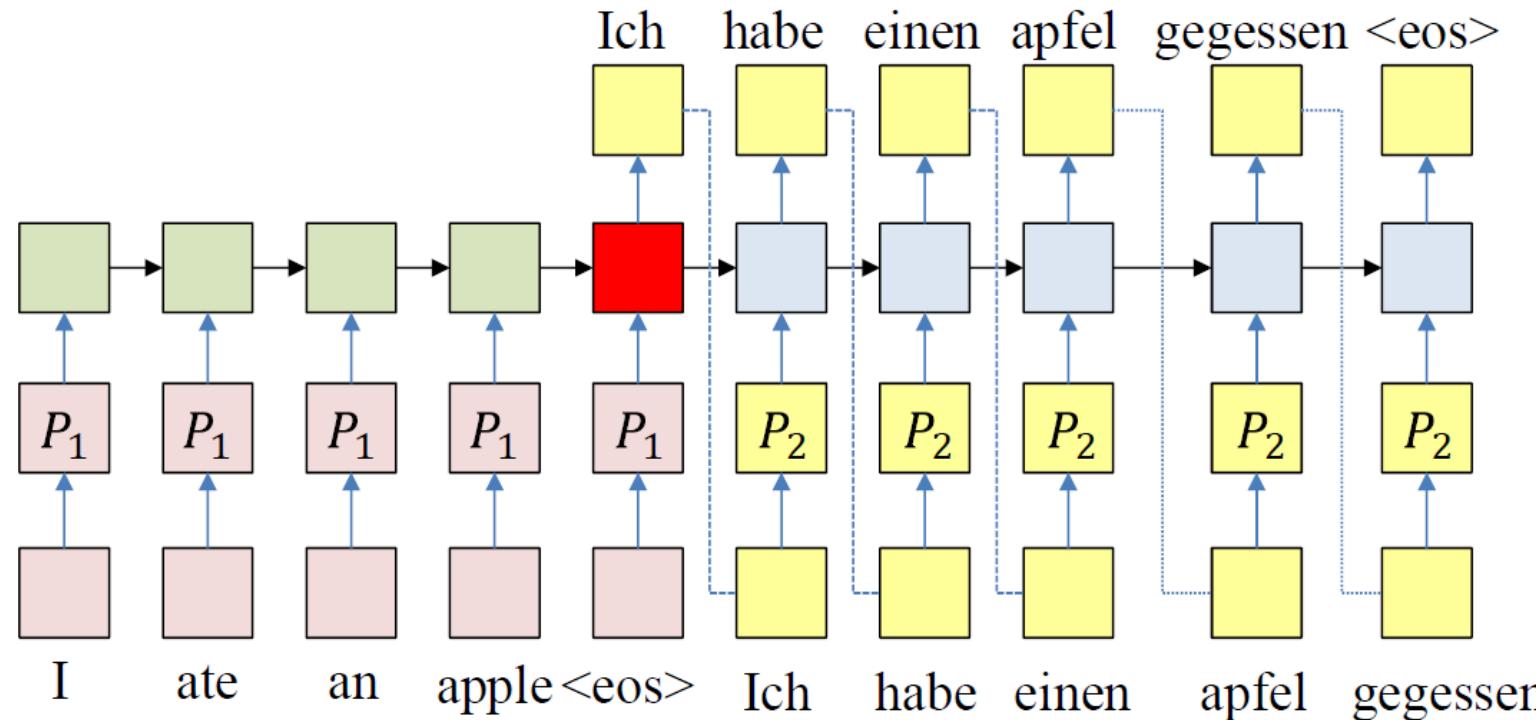


<https://jamalmar.github.io/visualizing-neural-machine-translation-mechanics-of-seq2seq-models-with-attention/>



The “simple” translation model

- A more detailed look
 - Word embedding can be incorporated
 - And will be learned along with the rest of the network

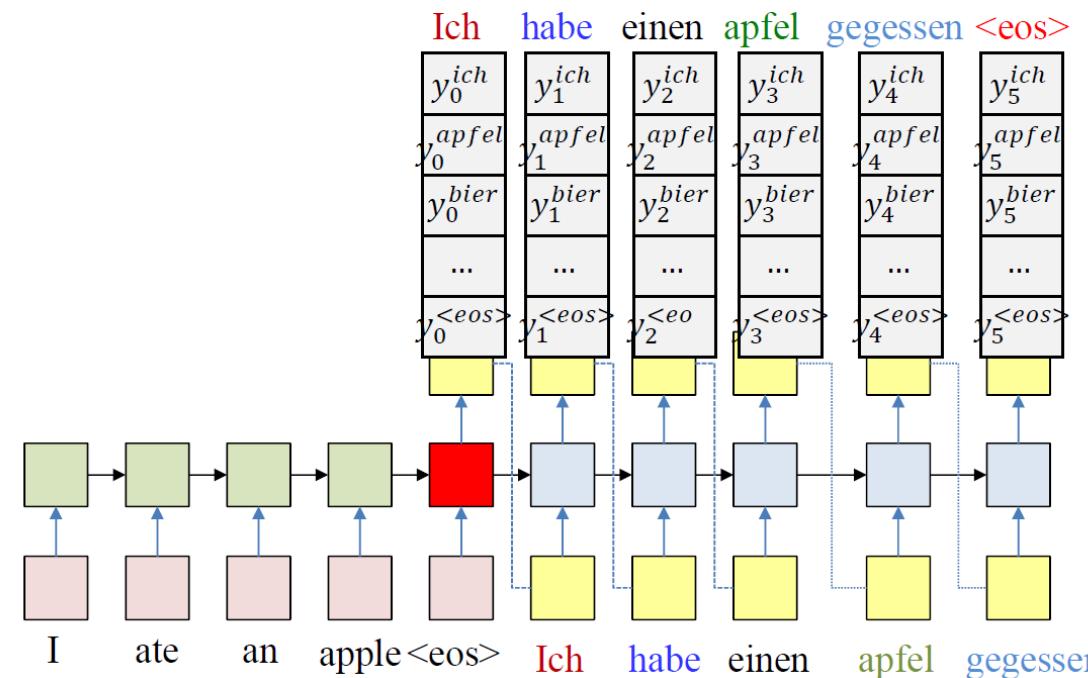


Prediction of translation model

- At each time k , the network produces a probability distribution over the output vocabulary

$$y_k^w = P(O_k = w | O_{k-1}, \dots, O_1, I_1, \dots, I_N)$$

- At each time, a word is drawn from the output distribution
- The drawn word is provided as input to the next time, until <eos>



Prediction of translation model

- At each time k , the network produces a probability distribution over the output vocabulary

$$y_k^w = P(O_k = w | O_{k-1}, \dots, O_1, I_1, \dots, I_N)$$

- At each time, a word is drawn from the output distribution
- The drawn word is provided as input to the next time, until <eos>
- For a given output sequence O_1, \dots, O_L , its probability is

$$P(O_1, \dots, O_L | I_1, \dots, I_N) = y_1^{O_1} y_2^{O_2} \cdots y_L^{O_L}$$

- The objective of drawing: produce the most probable output

$$\arg \max_{O_1, \dots, O_L} y_1^{O_1} y_2^{O_2} \cdots y_L^{O_L}$$



Human machine conversation example

- Task: Human text in, machine response out
- Trained on human-human conversations

Machine: *what is the error that you are running please*

Human: *i am seeing an error related to vpn*

Machine: *what is the error message that you are getting when connecting to vpn using network connect ?*

Human: *connection refused or something like that*

Machine: *may i know the version of network connect you connect ?*

Human: *i am not sure i know that*

Machine: *is the network connect prompting for any user certificate*

“A neural conversational model”, Orin Vinyals and Quoc Le

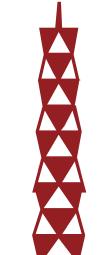


Image-to-sequence problem

- Generating image captions

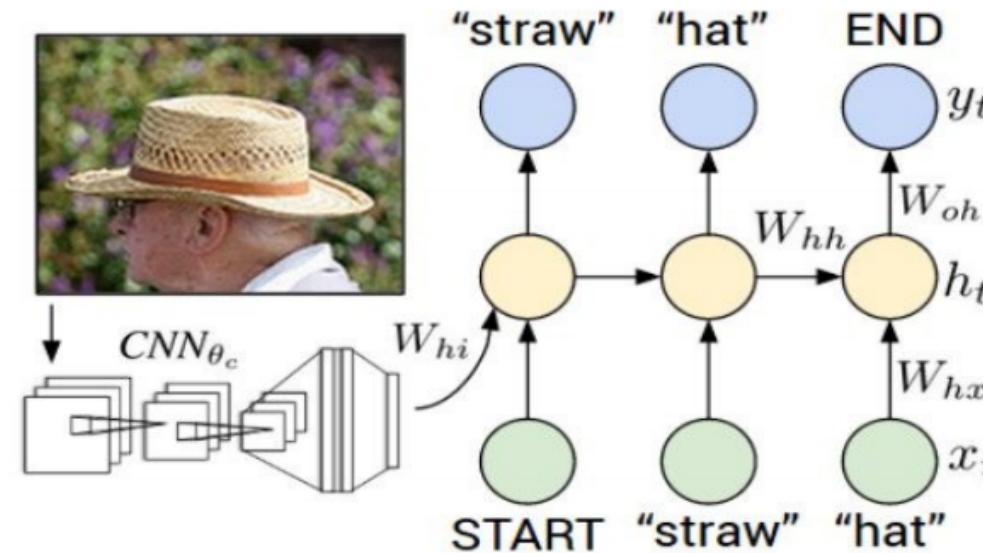


Figure from Karpathy et al, "Deep Visual-Semantic Alignments for Generating Image Descriptions", CVPR 2015; figure copyright IEEE, 2015.
Reproduced for educational purposes.

Explain Images with Multimodal Recurrent Neural Networks, Mao et al.

Deep Visual-Semantic Alignments for Generating Image Descriptions, Karpathy and Fei-Fei

Show and Tell: A Neural Image Caption Generator, Vinyals et al.

Long-term Recurrent Convolutional Networks for Visual Recognition and Description, Donahue et al.

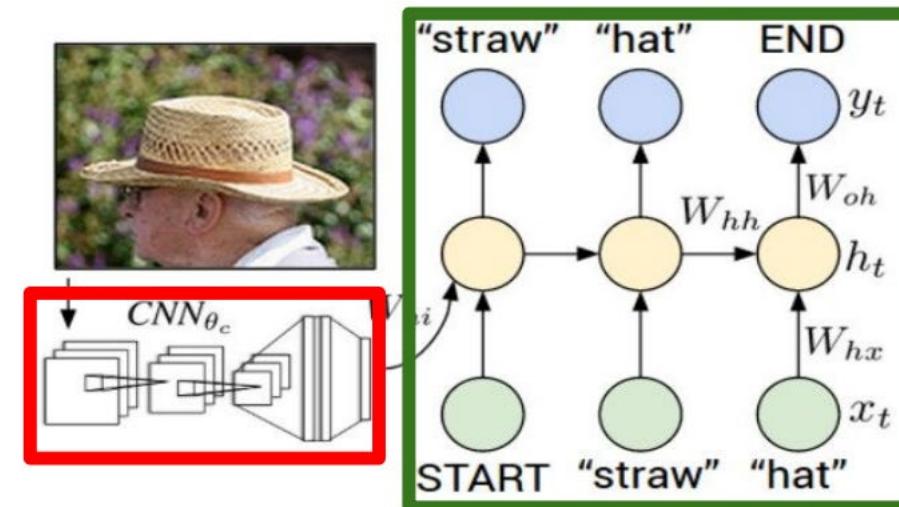
Learning a Recurrent Visual Representation for Image Caption Generation, Chen and Zitnick



Image-to-sequence problem

- Initial state is produced by a state-of-the-art CNN-based image classification system
- Subsequent model is the decoder end of a seq-to-seq model

Recurrent Neural Network



Convolutional Neural Network



Application: Image Captioning

- Example Results



"man in black shirt is playing guitar."



"construction worker in orange safety vest is working on road."



"two young girls are playing with lego toy."



"boy is doing backflip on wakeboard."



"a young boy is holding a baseball bat."



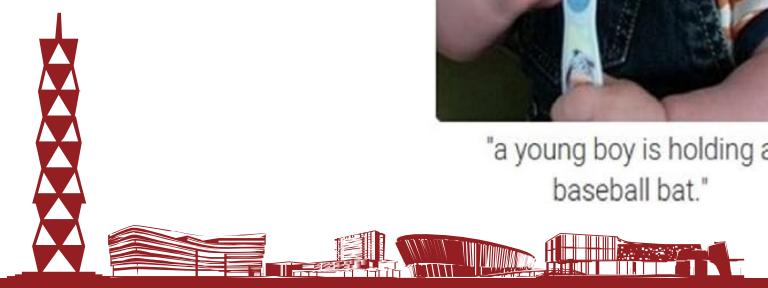
"a cat is sitting on a couch with a remote control."



"a woman holding a teddy bear in front of a mirror."



"a horse is standing in the middle of a road."





Summary



- In this lecture we have learned
 - Basic Recurrent Neural Networks
 - Training issues and LSTMs
 - RNN Applications
- For details, read the following resources
 - https://d2l.ai/chapter_recurrent-neural-networks/index.html
 - https://d2l.ai/chapter_recurrent-modern/index.html
 - https://d2l.ai/chapter_attention-mechanisms/index.html

