# CSE 481: NLP Capstone
# Spring 2017

Yejin Choi
University of Washington

# Office Hour News

- Hannah:
  - Wed 2 - 3pm @ CSE 220
- Maarten:
  - Wed 2 - 3pm @ CSE 220
- Yejin:
  - Tue 2pm - 3:30pm
  - Wed 5pm - 5:30pm @ CSE 578
- All:
  - Thu 12pm – 1:25pm @ ??? for some weeks
- Google doc sign up required

| Week | Dates | Topic | Leader |
|---|---|---|---|
| 1 | Mar 28, 30 | Course Overview, Project Pitch, TensorFlow Tutorial | Hannah, Maarten |
| 2 | Apr 4, 6 | Project Proposal Presentations & Discussion | Yejin |
| 3 | Apr 11, 13 | Lecture on Deep Learning & Project Update Meetings | Yejin |
| 4 | Apr 18, 20 | Lecture on Deep Learning & Project Update Meetings | Yejin |
| 5 | Apr 25, 27 | In Class Project Update Presentations! | All Students |
| 6 | May 2, 4 | Lecture on Deep Learning & Project Update Meetings | Yejin |
| 7 | May 9, 11 | Lecture on Deep Learning & Project Update Meetings | Yejin |
| 8 | May 16, 18 | In Class Project *Demo* and Presentations! | All Students |
| 9 | May 23, 25 | Lecture on Deep Learning & Project Update Meetings | Yejin |
| 10 | May 30, Jun 1 | Finale! - final poster presentation & demo @ CSE Atrium | All Students |

# GPU NEWS!

# GPU NEWS!

1. Back in stock!

   – desktop with 2 GPUs can be set up at $4000

2. Microsoft Azure kindly agreed to donate free GPU cycles for the class!!!!!
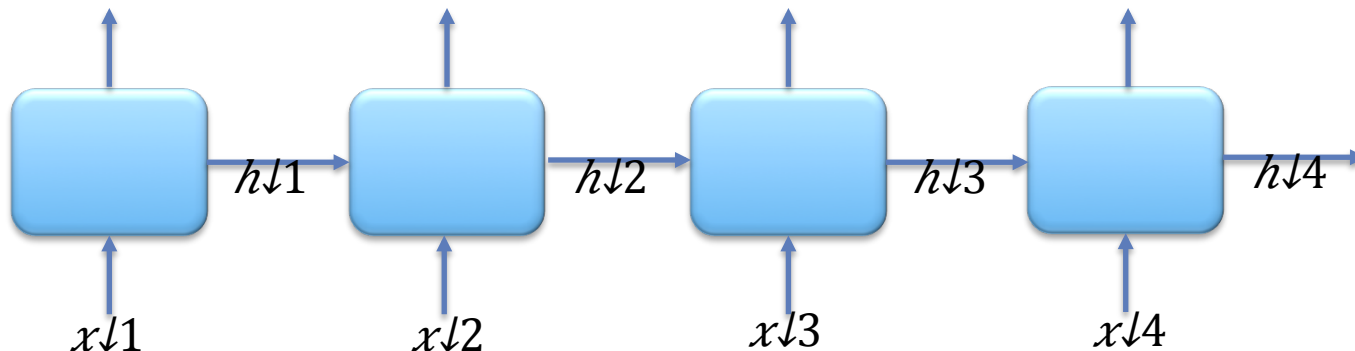
3. You can sign up to Azure today for free $200 credits

# RECURRENT NEURAL NETWORKS

# Recurrent Neural Networks (RNNs)

- Each RNN unit computes a new hidden state using the previous state and a new input $$h_t = f(x_t, h_{t-1})$$
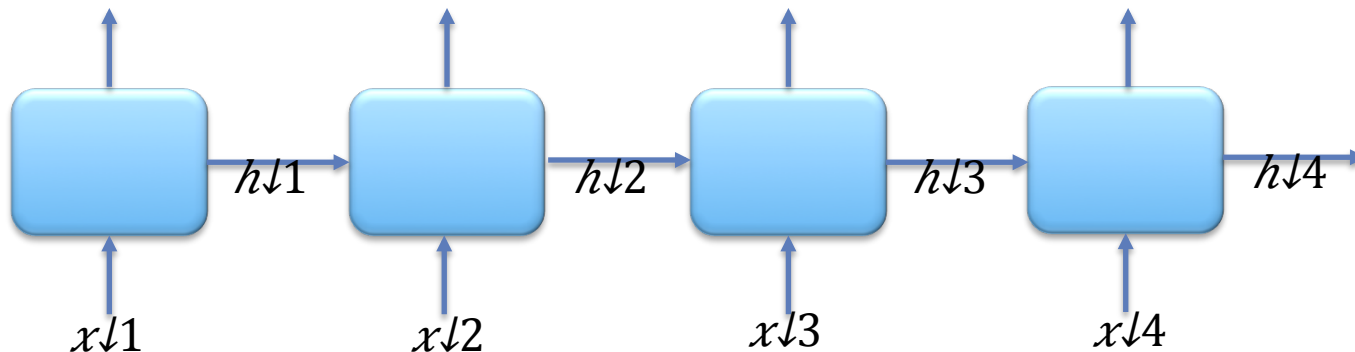- Each RNN unit (optionally) makes an output using the current hidden state $$y_t = \text{softmax}(V h_t)$$

$$h_t \in R^D$$

- Hidden states $h_t \in R^D$ are continuous vectors
  - Can represent very rich information
  - Possibly the entire history from the beginning
- Parameters are shared (tied) across all RNN units (unlike feedforward NNs)

# Recurrent Neural Networks (RNNs)

- Generic RNNs:   $h_t = f(x_t, h_{t-1})$

  $$y_t = \operatorname{softmax}(V h_t)$$

- Vanilla RNN:   $h_t = \tanh(U x_t + W h_{t-1} + b)$

  $$y_t = \operatorname{softmax}(V h_t)$$

# Recurrent Neural Networks (RNNs)

- Generic RNNs:  $h_t = f(x_t, h_{t-1})$

- Vanilla RNNs:  $h_t = \tanh(U x_t + W h_{t-1} + b)$

- LSTMs (Long Short-term Memory Networks):

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

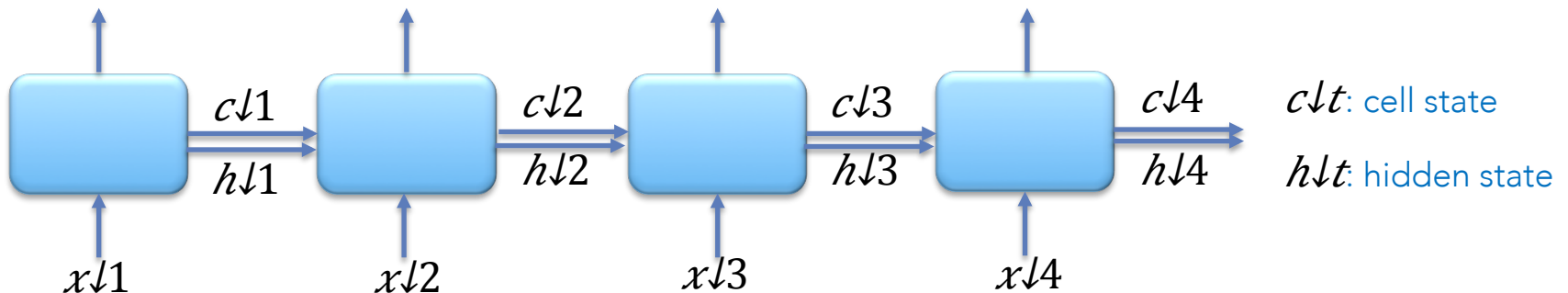$$f_t = \sigma(U^{(f)} x_t + W^{(f)} h_{t-1} + b^{(f)})$$

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

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

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

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

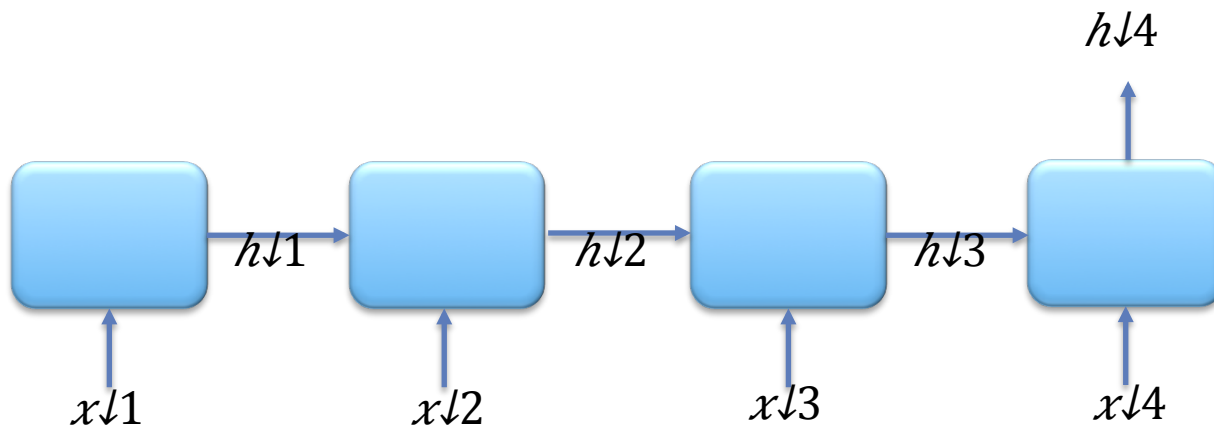There are many known variations to this set of equations!

$c_1$ $\quad$ $c_2$ $\quad$ $c_3$ $\quad$ $c_4$ $\qquad$ $c_t$: cell state

$h_1$ $\quad$ $h_2$ $\quad$ $h_3$ $\quad$ $h_4$ $\qquad$ $h_t$: hidden state

$x_1$ $\qquad$ $x_2$ $\qquad$ $x_3$ $\qquad$ $x_4$

# Many uses of RNNs
## 1. Classification (seq to one)

- Input: a sequence
- Output: one label (classification)
- Example: sentiment classification

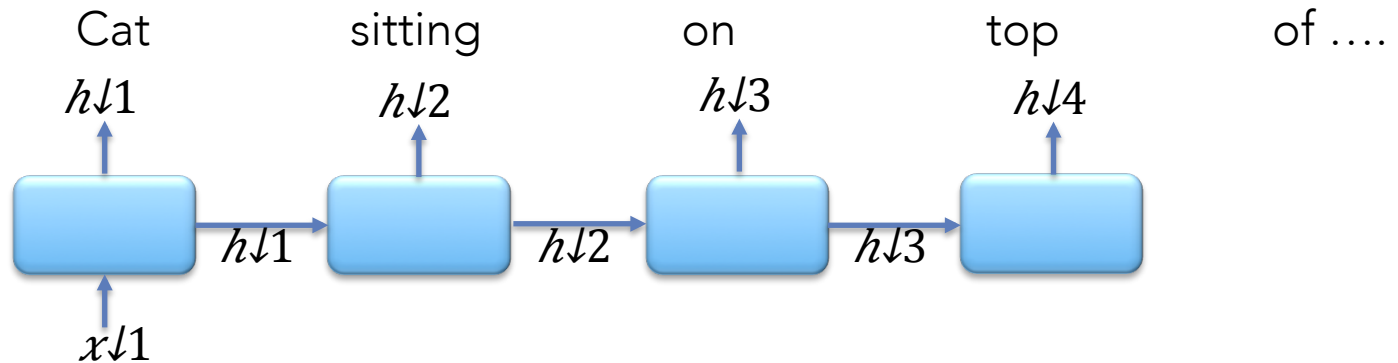$$h_t = f(x_t, h_{t-1})$$
$$y = \mathrm{softmax}(V h_n)$$

# Many uses of RNNs
# 2. one to seq

- Input: one item
- Output: a sequence
- Example: Image captioning

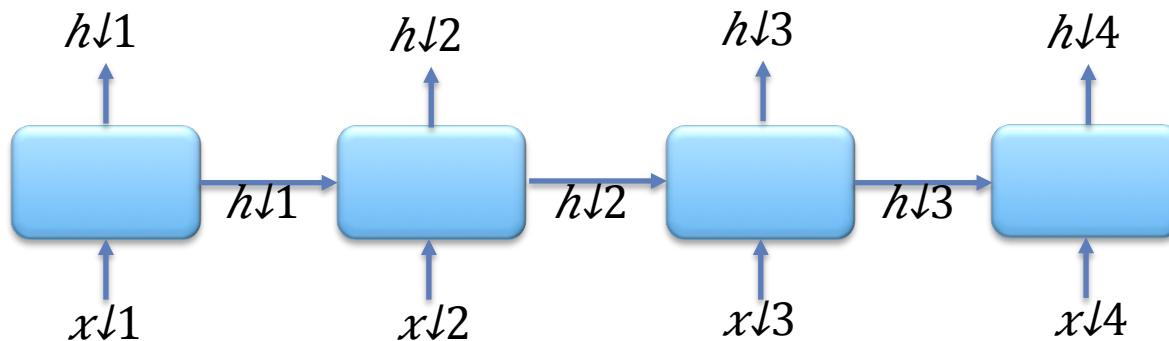$$h_t = f(x_t, h_{t-1})$$
$$y_t = \mathrm{softmax}(V h_t)$$

# Many uses of RNNs
# 3. sequence tagging

- Input: a sequence
- Output: a sequence (of the same length)
- Example: POS tagging, Named Entity Recognition
- How about Language Models?
  - Yes! RNNs can be used as LMs!
  - RNNs make markov assumption: T/F?

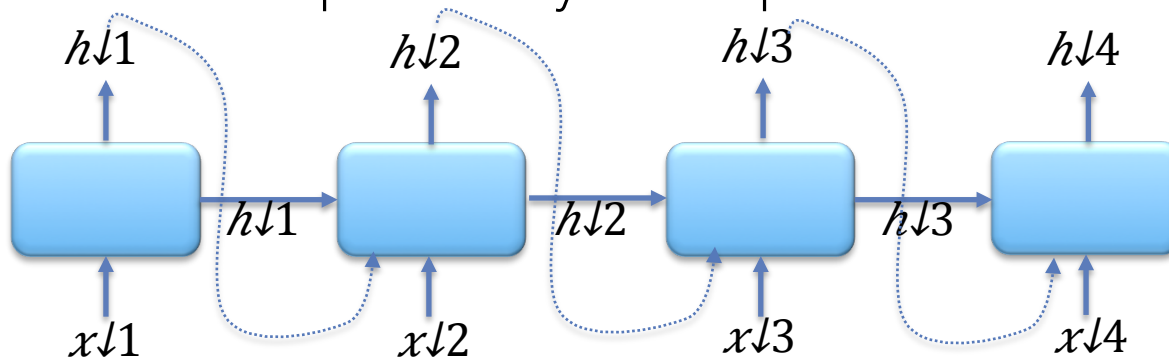$$h_t = f(x_t, h_{t-1})$$
$$y_t = \mathrm{softmax}(V h_t)$$

# Many uses of RNNs
# 4. Language models

- Input: a sequence of words
- Output: one next word
- Output: or a sequence of next words
- During training, x_t is the actual word in the training sentence.
- During testing, x_t is the word predicted from the previous time step.
- Does RNN LMs make Markov assumption?
  - i.e., the next word depends only on the previous N words
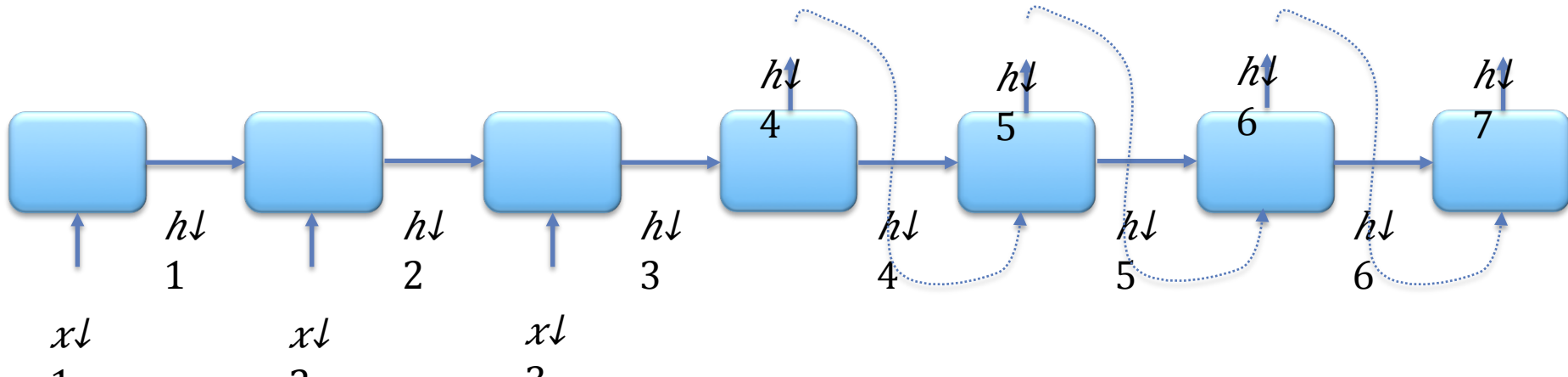
$$h_t = f(x_t, h_{t-1})$$
$$y_t = \text{softmax}(Vh_t)$$

# Many uses of RNNs
## 5. seq2seq (aka "encoder-decoder")

- Input: a sequence
- Output: a sequence (of *different* length)
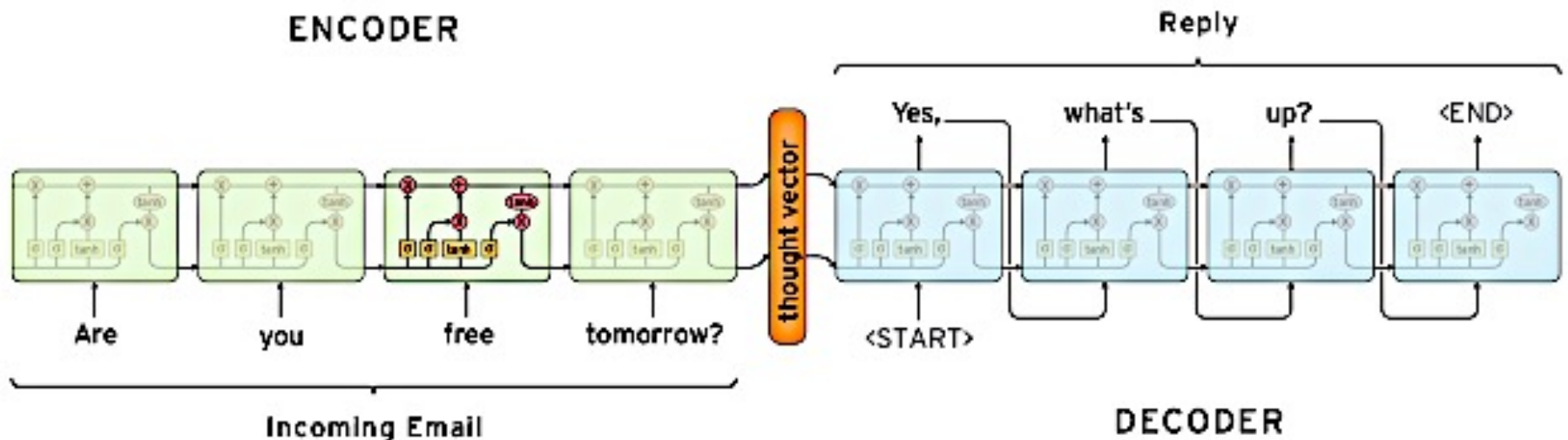- Examples?

$$h_t = f(x_t, h_{t-1})$$
$$y_t = \text{softmax}(V h_t)$$

# Many uses of RNNs
# 4. seq2seq (aka "encoder-decoder")

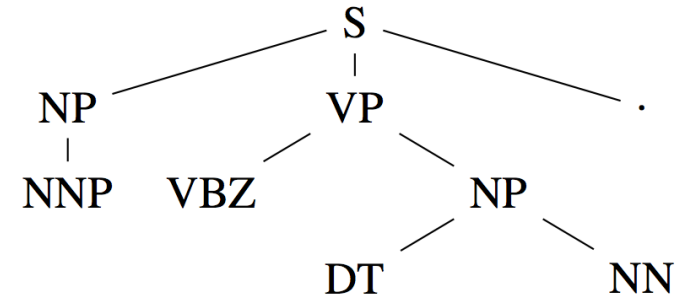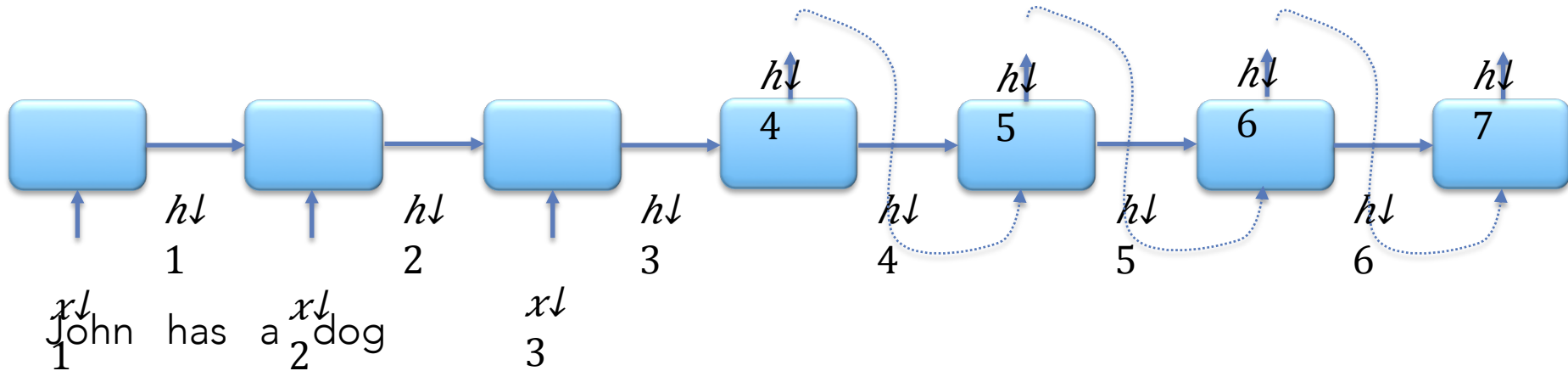- Conversation and Dialogue
- Machine Translation

# Many uses of RNNs
## 4. seq2seq (aka "encoder-decoder")

Parsing!
- *"Grammar as Foreign Language"* (Vinyals et al., 2015)
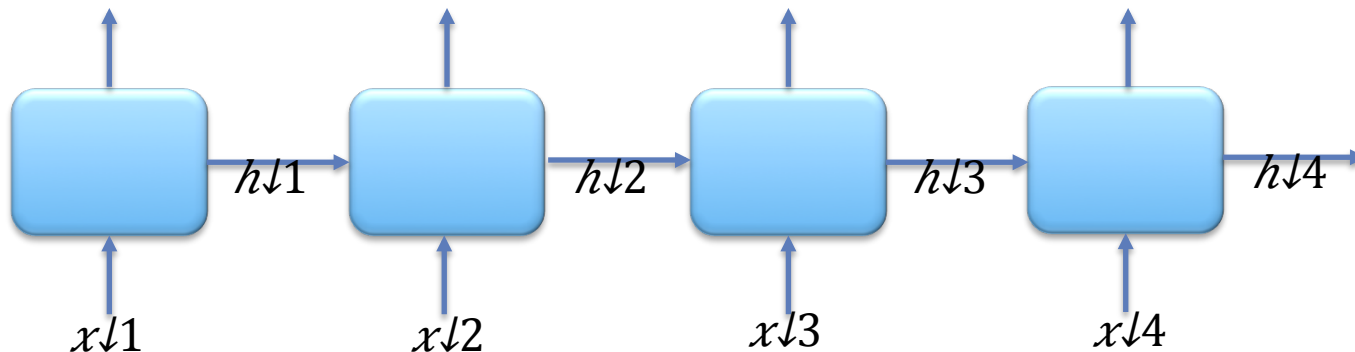
(S (NP NNP )$_{NP}$ (VP VBZ (NP DT NN )$_{NP}$ )$_{VP}$ . )$_S$

# Recurrent Neural Networks (RNNs)

- Generic RNNs:
$$h_t = f(x_t, h_{t-1})$$
$$y_t = \text{softmax}(V h_t)$$

- Vanilla RNN:
$$h_t = \tanh(U x_t + W h_{t-1} + b)$$
$$y_t = \text{softmax}(V h_t)$$

# Recurrent Neural Networks (RNNs)

- Generic RNNs: $h_t = f(x_t, h_{t-1})$

- Vanilla RNNs: $h_t = \tanh(Ux_t + Wh_{t-1} + b)$

- LSTMs (Long Short-term Memory Networks):

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

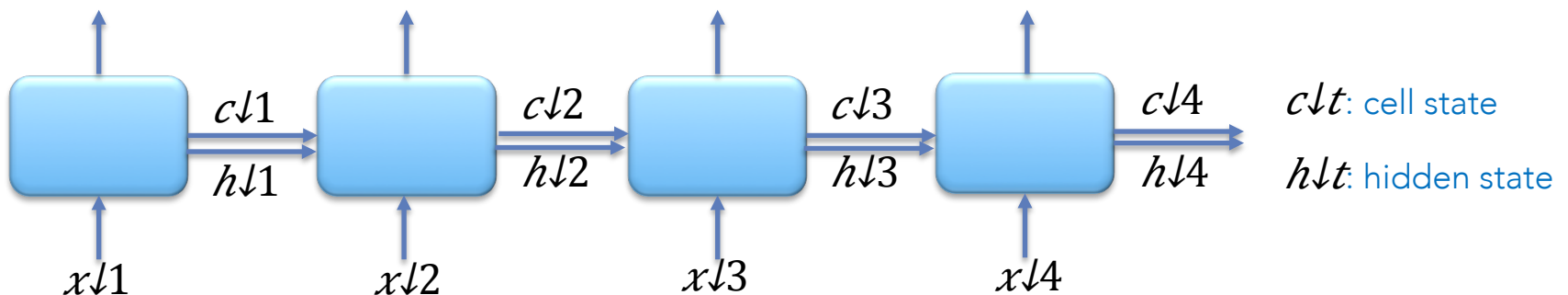$$f_t = \sigma(U^{(f)}x_t + W^{(f)}h_{t-1} + b^{(f)})$$

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

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

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

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

There are many known variations to this set of equations!



$c_t$: cell state

$h_t$: hidden state

# LSTMS (LONG SHORT-TERM MEMORY NETWORKS)



Figure by Christopher Olah (colah.github.io)

# LSTMS (LONG SHORT-TERM MEMORY NETWORKS

sigmoid: [0,1]

Forget gate: forget the past or not

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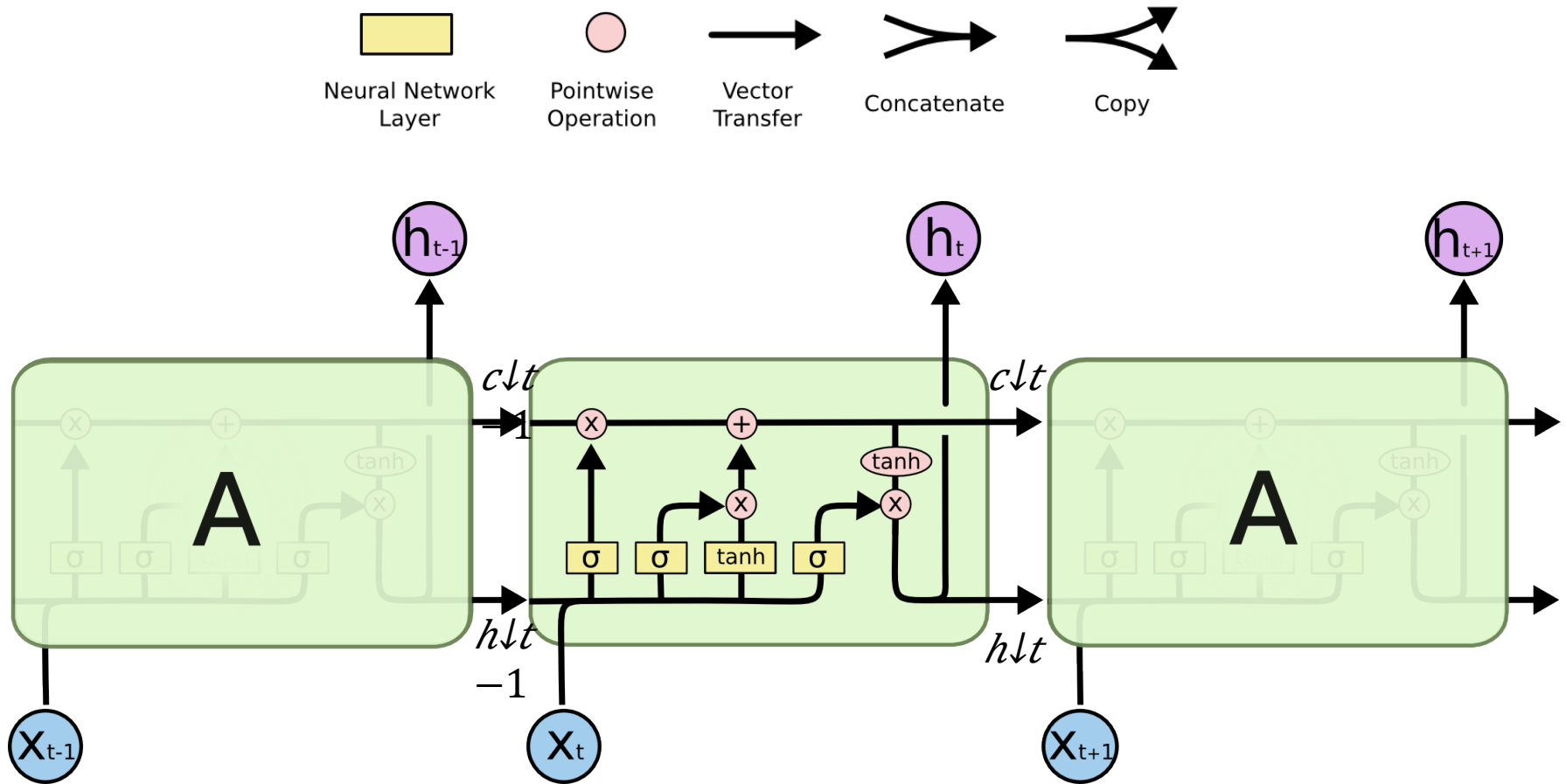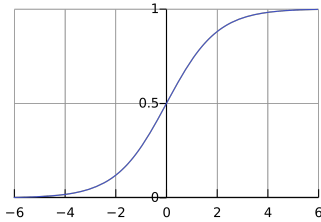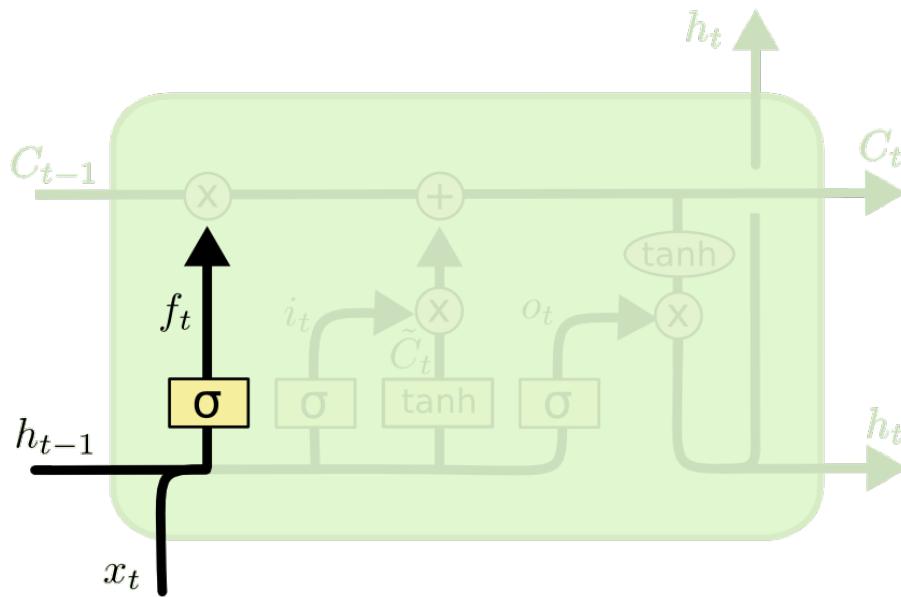


Figure by Christopher Olah (colah.github.io)

# LSTMS (LONG SHORT-TERM MEMORY NETWORKS)

sigmoid:
[0,1]



tanh:
[-1,1]



Forget gate: forget the past or not
$$f_t = \sigma(U^{(f)}x_t + W^{(f)}h_{t-1} + b^{(f)})$$

Input gate: use the input or not
$$i_t = \sigma(U^{(i)}x_t + W^{(i)}h_{t-1} + b^{(i)})$$

New cell content (temp):
$$\tilde{c}_t = \tanh(U^{(c)}x_t + W^{(c)}h_{t-1} + b^{(c)})$$



Figure by Christopher Olah (colah.github.io)
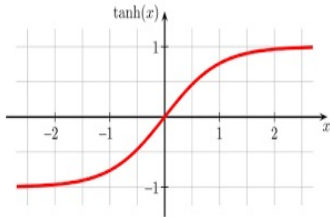
# LSTMS (LONG SHORT-TERM MEMORY NETWORKS)

sigmoid:
[0,1]

tanh:
[-1,1]

Forget gate: forget the past or not
$$f_t = \sigma(U^{(f)}x_t + W^{(f)}h_{t-1} + b^{(f)})$$

Input gate: use the input or not
$$i_t = \sigma(U^{(i)}x_t + W^{(i)}h_{t-1} + b^{(i)})$$

New cell content (temp):
$$\tilde{c}_t = \tanh(U^{(c)}x_t + W^{(c)}h_{t-1} + b^{(c)})$$

New cell content:
  - mix old cell with the new temp cell
$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$$

Figure by Christopher Olah (colah.github.io)

# LSTMS (LONG SHORT-TERM MEMORY NETWORKS

Output gate: output from the new cell or not

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

Hidden state:

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

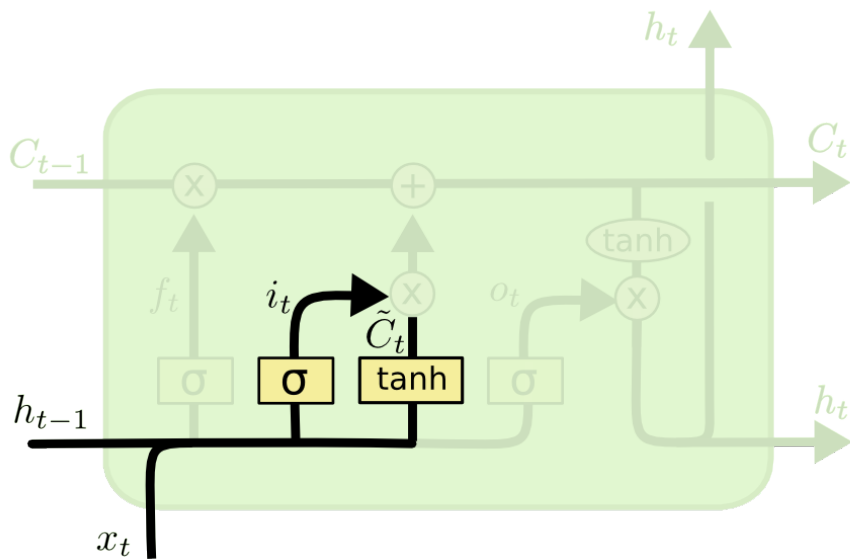Forget gate: forget the past or not

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

Input gate: use the input or not

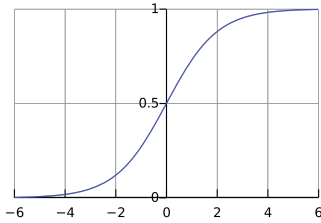$$i_t = \sigma(U^{(i)}x_t + W^{(i)}h_{t-1} + b^{(i)})$$

New cell content (temp):

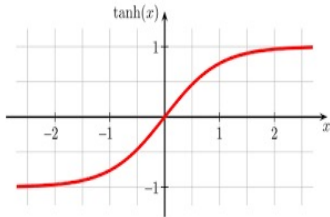$$\tilde{c}_t = \tanh(U^{(c)}x_t + W^{(c)}h_{t-1} + b^{(c)})$$

New cell content:
  - mix old cell with the new temp cell

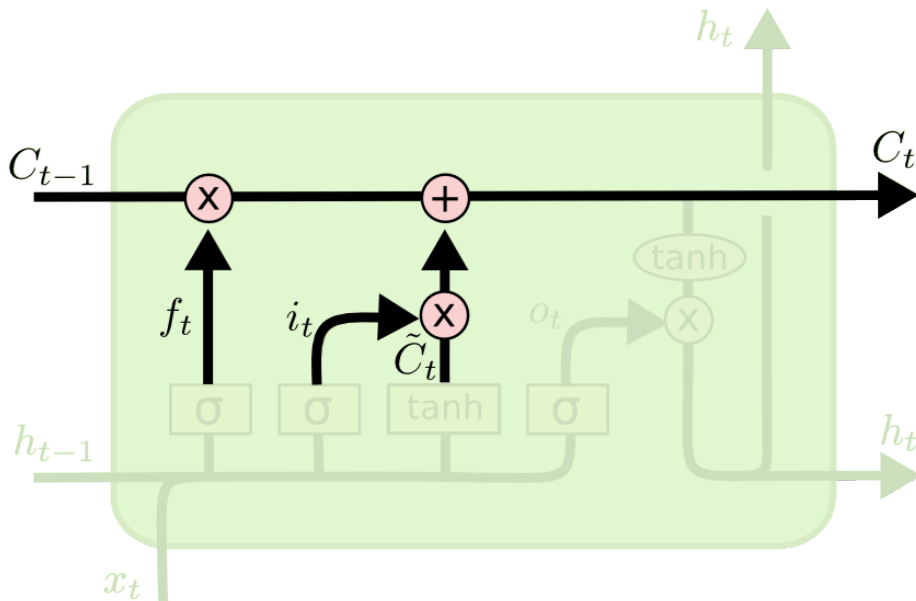$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$$



Figure by Christopher Olah (colah.github.io)

# LSTMS (LONG SHORT-TERM MEMORY NETWORKS

Forget gate: forget the past or not

Input gate: use the input or not

Output gate: output from the new cell or not

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

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

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

---

New cell content (temp):

New cell content:
   - mix old cell with the new temp cell

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

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

Hidden state:

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

# 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.

Example from Graves 2012

# Preservation of gradient information by LSTM



- 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.

Example from Graves 2012

# Recurrent Neural Networks (RNNs)

- Generic RNNs: $h_t = f(x_t, h_{t-1})$
- Vanilla RNNs: $h_t = \tanh(U x_t + W h_{t-1} + b)$
- GRUs (Gated Recurrent Units):

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

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

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

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

Z: Update gate
R: Reset gate

Less parameters than LSTMs. Easier to train for comparable performance!



$h{\downarrow}1$　$h{\downarrow}2$　$h{\downarrow}3$　$h{\downarrow}4$

$x{\downarrow}1$　$x{\downarrow}2$　$x{\downarrow}3$　$x{\downarrow}4$

# Gates

- Gates contextually control information flow
- Open/close with sigmoid
- In LSTMs and GRUs, they are used to (contextually) maintain longer term history

# Bi-directional RNNs

Outputs $\quad \cdots \quad y_{t-1} \qquad y_t \qquad y_{t+1} \quad \cdots$

Backward Layer $\quad \overleftarrow{h}_{t-1} \qquad \overleftarrow{h}_t \qquad \overleftarrow{h}_{t+1}$

Forward Layer $\quad \overrightarrow{h}_{t-1} \qquad \overrightarrow{h}_t \qquad \overrightarrow{h}_{t+1}$

Inputs $\quad \cdots \quad x_{t-1} \qquad x_t \qquad x_{t+1} \quad \cdots$

- Can incorporate context from both directions
- Generally improves over uni-directional RNNs

# Google NMT (Oct 2016)

# Recursive Neural Networks

- Sometimes, inference over a tree structure makes more sense than sequential structure
- An example of compositionality in ideological bias detection (red → conservative, blue → liberal, gray → neutral) in which modifier phrases and punctuation cause polarity switches at higher levels of the parse tree



They dubbed it the " death tax " and created a big lie about its adverse effects on small businesses

Example from Iyyer et al., 2014

# Recursive Neural Networks

- NNs connected as a tree
- Tree structure is fixed a priori
- Parameters are shared, similarly as RNNs

# Tree LSTMs

$y_1$    $y_2$    $y_3$    $y_4$

$x_1$    $x_2$    $x_3$    $x_4$

$y_1$

$y_2$    $x_1$    $y_3$

$x_2$    $y_4$    $y_6$

$x_4$    $x_5$    $x_6$

Figure 1: **Top:** A chain-structured LSTM network. **Bottom:** A tree-structured LSTM network with arbitrary branching factor.

- Are tree LSTMs more expressive than sequence LSTMs?

- I.e., *recursive* vs *recurrent*

- When Are Tree Structures Necessary for Deep Learning of Representations? Jiwei Li, Minh-Thang Luong, Dan Jurafsky and Eduard Hovy. EMNLP, 2015.

# Neural Probabilistic Language Model (Bengio 2003)

$i$-th output $= P(w_t = i \mid context)$

softmax

most computation here

tanh

$C(w_{t-n+1})$     $C(w_{t-2})$     $C(w_{t-1})$

Table look-up in $C$

Matrix $C$

shared parameters across words

index for $w_{t-n+1}$     index for $w_{t-2}$     index for $w_{t-1}$

# Neural Probabilistic Language Model (Bengio 2003)



$i$-th output $= P(w_t = i \mid context)$

softmax

most computation here

tanh

$C(w_{t-n+1})$  $C(w_{t-2})$  $C(w_{t-1})$

Table look–up in $C$

Matrix $C$ shared parameters across words

index for $w_{t-n+1}$   index for $w_{t-2}$   index for $w_{t-1}$

- Each word prediction is a separate feed forward neural network
- Feedforward NNLM is a Markovian language model
- Dashed lines show optional direct connections

$$NN_{DMLP1}(\mathbf{x}) = [\tanh(\mathbf{x}\mathbf{W}^1 + \mathbf{b}^1), \mathbf{x}]W^2 + \mathbf{b}^2$$

▶ $\mathbf{W}^1 \in \mathbb{R}^{d_{\text{in}} \times d_{\text{hid}}}, \mathbf{b}^1 \in \mathbb{R}^{1 \times d_{\text{hid}}}$; first affine transformation

▶ $\mathbf{W}^2 \in \mathbb{R}^{(d_{\text{hid}} + d_{\text{in}}) \times d_{\text{out}}}, \mathbf{b}^2 \in \mathbb{R}^{1 \times d_{\text{out}}}$; second affine transformation

# LEARNING: BACKPROPAGATION

# Error Backpropagation

- Model parameters: $\vec{\theta} = \{w_{ij}^{(1)}, w_{jk}^{(2)}, w_{kl}^{(3)}\}$

  for brevity: $\vec{\theta} = \{w_{ij}, w_{jk}, w_{kl}\}$

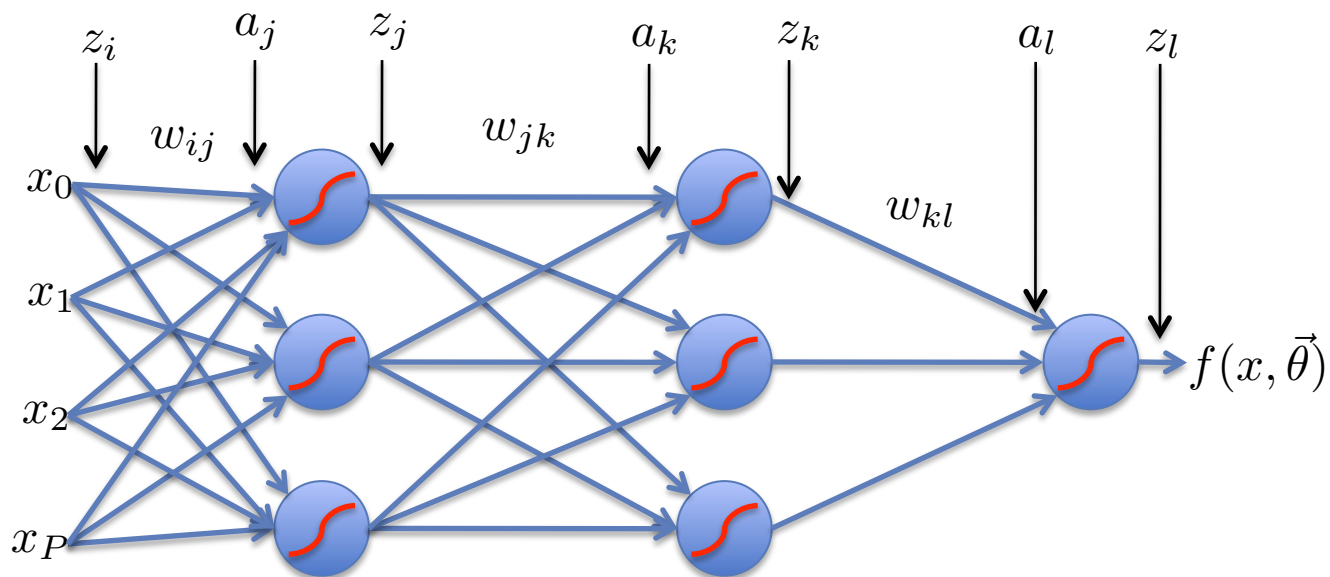# Learning: Gradient Descent

$$w_{ij}^{t+1} \quad = \quad w_{ij}^{t} - \eta \frac{\partial R}{w_{ij}}$$

$$w_{jk}^{t+1} \quad = \quad w_{jk}^{t} - \eta \frac{\partial R}{w_{kl}}$$

$$w_{kl}^{t+1} \quad = \quad w_{kl}^{t} - \eta \frac{\partial R}{w_{kl}}$$

# Backpropagation

- Starts with a forward sweep to compute all the intermediate function $z_i$ values

- Through backprop, computes the partial derivatives recursively $\delta_j$ $\dfrac{\partial R}{\partial w_{ij}}$

- A form of dynamic programming

  - Instead of considering exponentially many paths between a weight w_ij and the final loss (risk), store and reuse intermediate results.

- A type of automatic differentiation. (there are other variants e.g., recursive differentiation only through forward propagation



Forward    Inputs    Outputs

Gradient

# Backpropagation

Primary Interface Language

- TensorFlow (https://www.tensorflow.org/)
- Torch (http://torch.ch/)
- Theano (http://deeplearning.net/software/theano/)
- CNTK (https://github.com/Microsoft/CNTK)
- cnn (https://github.com/clab/cnn)
- Caffe (http://caffe.berkeleyvision.org/)

- Python
- Lua
- Python
- C++
- C++
- C++



Forward

Inputs

Outputs

Gradient

# Cross Entropy Loss (aka log loss, logistic loss)

- Cross Entropy

$$H(p, q) = -\sum_y p(y) \log q(y)$$

Predicted prob

True prob

- Related quantities

  $$H(p) = \sum_y p(y) \log p(y)$$

  - Entropy

  - KL divergence (the distance between two distributions p and q)

    $$D_{KL}(p||q) = \sum_y p(y) \log \frac{p(y)}{q(y)}$$

    $$H(p, q) = E_p[-\log q] = H(p) + D_{KL}(p||q)$$

- Use Cross Entropy for models that should have more probabilistic flavor (e.g., language models)

- Use Mean Squared Error loss for models that focus on correct/incorrect predictions

  $$\text{MSE} = \frac{1}{2}(y - f(x))^2$$

# RNN Learning: Backprop Through Time (BPTT)

- Similar to backprop with non-recurrent NNs
- But unlike feedforward (non-recurrent) NNs, each unit in the computation graph repeats the exact same parameters…
- Backprop gradients of the parameters of each unit as if they are different parameters
- When updating the parameters using the gradients, use the average gradients throughout the entire chain of units.

# LEARNING: TRAINING DEEP NETWORKS

# Vanishing / exploding Gradients

- Deep networks are hard to train
- Gradients go through multiple layers
- The multiplicative effect tends to lead to *exploding* or *vanishing* gradients
- Practical solutions w.r.t.
  - network architecture
  - numerical operations

# Vanishing / exploding Gradients

- Practical solutions w.r.t. network architecture
  - Add skip connections to reduce distance
    - Residual networks, highway networks, …
  - Add gates (and memory cells) to allow longer term memory
    - LSTMs, GRUs, memory networks, …

# Gradients of deep networks

$$NN_{layer}(\mathbf{x}) = \text{ReLU}(\mathbf{x}\mathbf{W}^1 + \mathbf{b}^1)$$

$\mathbf{h}_n$

$|$

$\mathbf{h}_{n-1}$

$|$

$\cdots$
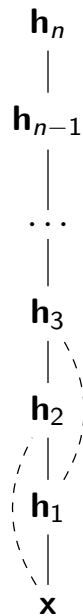
$|$

$\mathbf{h}_2$

$|$

$\mathbf{h}_1$

$|$

$\mathbf{x}$

▶ Can have similar issues with vanishing gradients.

$$\frac{\partial L}{\partial h_{n-1,j_{n-1}}} = \sum_{j_n} \mathbf{1}(h_{n,j_n} > 0) W_{j_{n-1},j_n} \frac{\partial L}{\partial h_{n,j_n}}$$

Diagram borrowed from Alex Rush

# Effects of Skip Connections on Gradients

- Thought Experiment: Additive Skip-Connections

$$NN_{s/1}(\mathbf{x}) = \frac{1}{2}\,\text{ReLU}(\mathbf{x}\mathbf{W}^1 + \mathbf{b}^1) + \frac{1}{2}\mathbf{x}$$

$$\frac{\partial L}{\partial h_{n-1,j_{n-1}}} = \frac{1}{2}\left(\sum_{j_n}\mathbf{1}(h_{n,j_n} > 0)\,W_{j_{n-1},j_n}\frac{\partial L}{\partial h_{n,j_n}}\right) +$$

$$\frac{1}{2}\left(h_{n-1,j_{n-1}}\frac{\partial L}{\partial h_{n,j_{n-1}}}\right)$$

$\mathbf{h}_n$

$\mathbf{h}_{n-1}$

$\ldots$

$\mathbf{h}_3$

$\mathbf{h}_2$

$\mathbf{h}_1$

$\mathbf{x}$

Diagram borrowed from Alex Rush
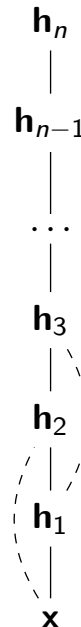
# Effects of Skip Connections on Gradients

- Thought Experiment: Dynamic Skip-Connections

$$NN_{sl2}(\mathbf{x}) = (1-t)\,\mathrm{ReLU}(\mathbf{x}\mathbf{W}^1 + \mathbf{b}^1) + t\mathbf{x}$$

$$t = \sigma(\mathbf{x}\mathbf{W}^t + b^t)$$

$$\mathbf{W}^1 \in \mathbb{R}^{d_{\mathrm{hid}} \times d_{\mathrm{hid}}}$$

$$\mathbf{W}^t \in \mathbb{R}^{d_{\mathrm{hid}} \times 1}$$

$$\mathbf{h}_n$$
$$|$$
$$\mathbf{h}_{n-1}$$
$$|$$
$$\ldots$$
$$|$$
$$\mathbf{h}_3$$
$$\mathbf{h}_2$$
$$\mathbf{h}_1$$
$$\mathbf{x}$$

# Highway Network (Srivastava et al., 2015)

- A plain feedforward neural network:
$$\mathbf{y} = H(\mathbf{x}, \mathbf{W_H}).$$

  – H is a typical affine transformation followed by a non-linear activation
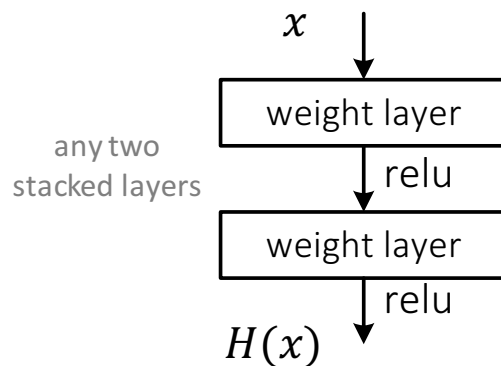
- Highway network:
$$\mathbf{y} = H(\mathbf{x}, \mathbf{W_H}) \cdot T(\mathbf{x}, \mathbf{W_T}) + \mathbf{x} \cdot C(\mathbf{x}, \mathbf{W_C}).$$
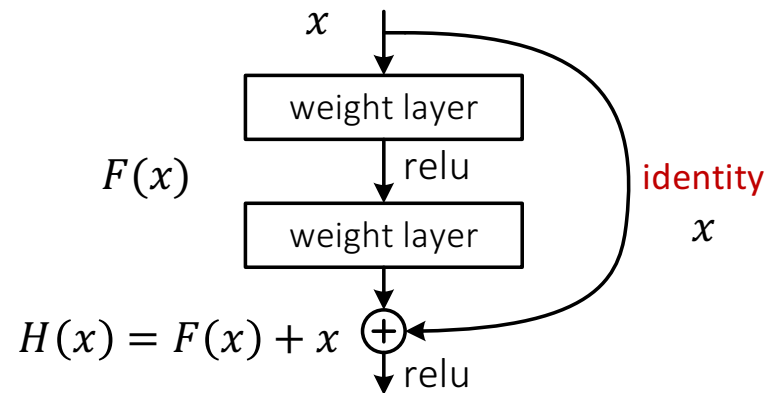
  – T is a "transform gate"
  – C is a "carry gate"
  – Often C = 1 – T for simplicity

# Residual Networks



- Plaint net
- <span style="color:red">Residual</span> net

any two stacked layers

$x$

weight layer

relu

weight layer

relu

$H(x)$

$x$

weight layer

relu

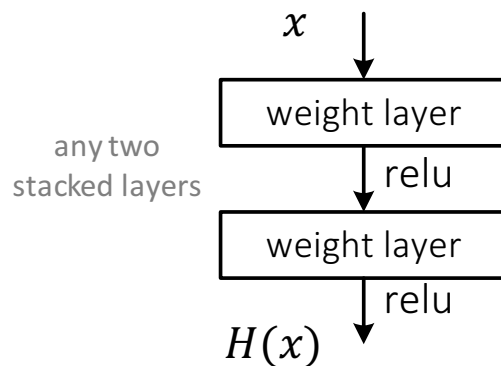$F(x)$

weight layer
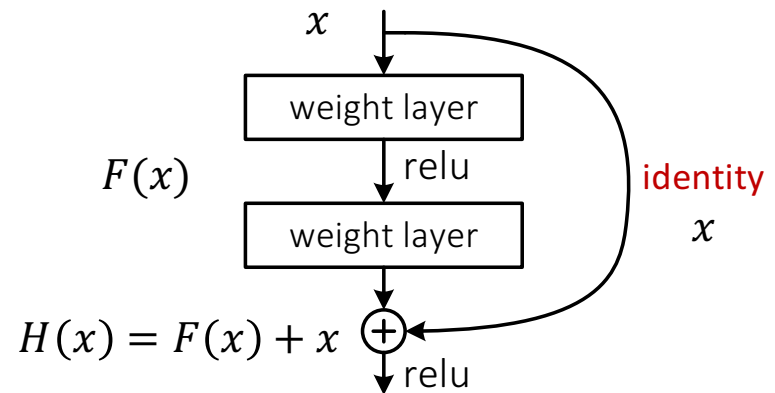
identity
$x$

$H(x) = F(x) + x$ $\oplus$

relu

- ResNet (He et al. 2015): first very deep (152 layers) network successfully trained for object recognition

# Residual Networks

• Plaint net

• **Residual** net

any two
stacked layers

$x$

weight layer

relu

weight layer

relu

$H(x)$

$x$

weight layer

relu

weight layer

$F(x)$

identity

$x$

$H(x) = F(x) + x$ $\oplus$

relu
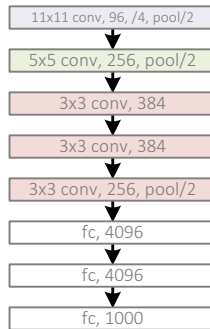
- F(x) is a residual mapping with respect to identity
- Direct input connection +x leads to a nice property w.r.t. back propagation --- more direct influence from the final loss to any deep layer
- In contrast, LSTMs & Highway networks allow for long distance input connection only through "gates".
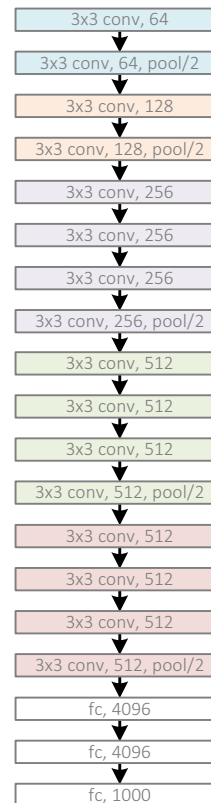
# Residual Networks

## Revolution of Depth

AlexNet, 8 layers
(ILSVRC 2012)

| 11x11 conv, 96, /4, pool/2 |
| 5x5 conv, 256, pool/2 |
| 3x3 conv, 384 |
| 3x3 conv, 384 |
| 3x3 conv, 256, pool/2 |
| fc, 4096 |
| fc, 4096 |
| fc, 1000 |

VGG, 19 layers
(ILSVRC 2014)

| 3x3 conv, 64 |
| 3x3 conv, 64, pool/2 |
| 3x3 conv, 128 |
| 3x3 conv, 128, pool/2 |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| 3x3 conv, 256, pool/2 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512, pool/2 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512, pool/2 |
| fc, 4096 |
| fc, 4096 |
| fc, 1000 |

GoogleNet, 22 layers
(ILSVRC 2014)

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016

52

# Residual Networks

## Revolution of Depth

AlexNet, 8 layers
(ILSVRC 2012)

VGG, 19 layers
(ILSVRC 2014)
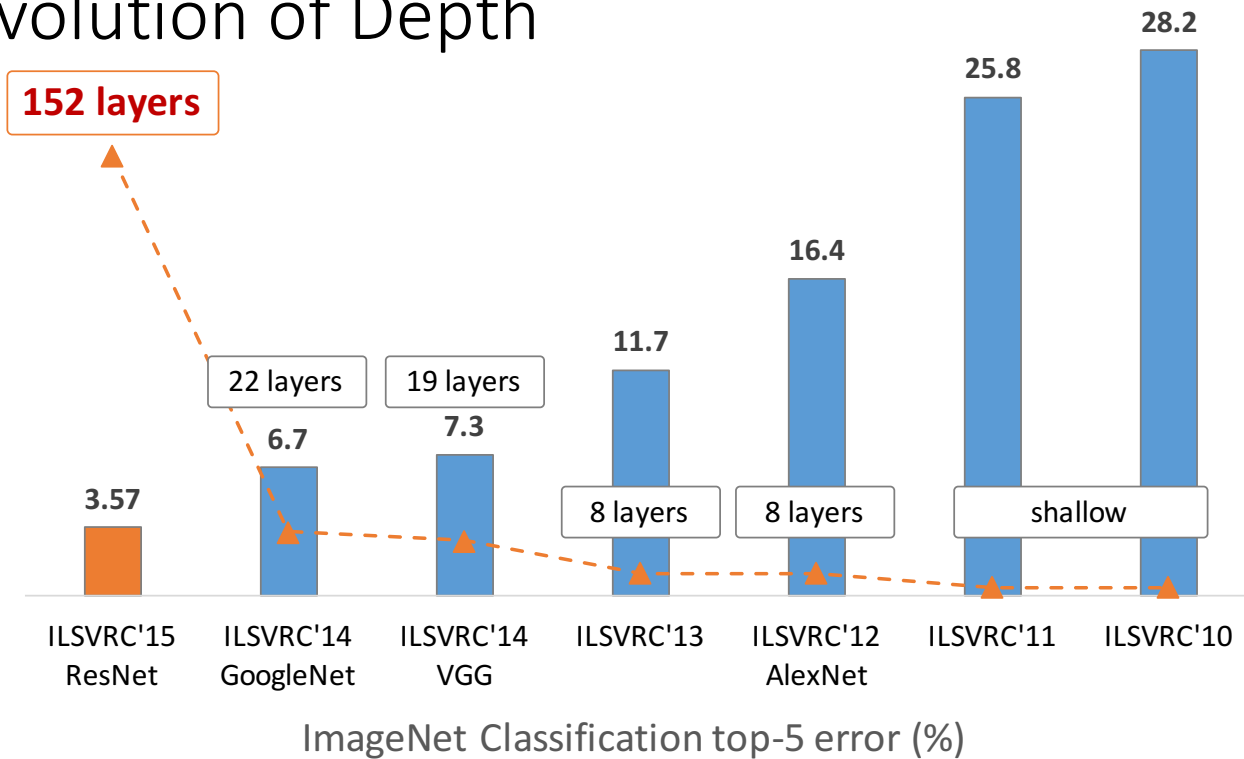
ResNet, 152 layers
(ILSVRC 2015)

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

# Residual Networks

## Revolution of Depth



**152 layers**

28.2

25.8

16.4

22 layers    19 layers

11.7

3.57

6.7    7.3

8 layers    8 layers    shallow

ILSVRC'15    ILSVRC'14    ILSVRC'14    ILSVRC'13    ILSVRC'12    ILSVRC'11    ILSVRC'10
ResNet       GoogleNet    VGG                       AlexNet

ImageNet Classification top-5 error (%)

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

# Highway Network (Srivastava et al., 2015)

- A plain feedforward neural network:
$$\mathbf{y} = H(\mathbf{x}, \mathbf{W_H}).$$

  – H is a typical affine transformation followed by a non-linear activation

- Highway network:
$$\mathbf{y} = H(\mathbf{x}, \mathbf{W_H}) \cdot T(\mathbf{x}, \mathbf{W_T}) + \mathbf{x} \cdot C(\mathbf{x}, \mathbf{W_C}).$$

  – T is a "transform gate"
  – C is a "carry gate"
  – Often C = 1 – T for simplicity

# @Schmidhubered

# Vanishing / exploding Gradients
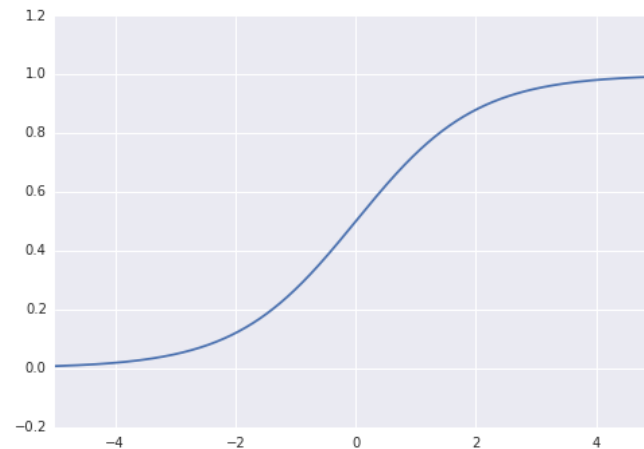
- Practical solutions w.r.t. numerical operations
  - Gradient Clipping: bound gradients by a max value

  - Gradient Normalization: renormalize gradients when they are above a fixed norm
  - Careful initialization, smaller learning rates
  - Avoid saturating nonlinearities (like tanh, sigmoid)
    - ReLU or hard-tanh instead

# Sigmoid

- Often used for gates
- Pro: neuron-like, differentiable
- Con: gradients saturate to zero almost everywhere except x near zero => vanishing gradients
- Batch normalization helps

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$
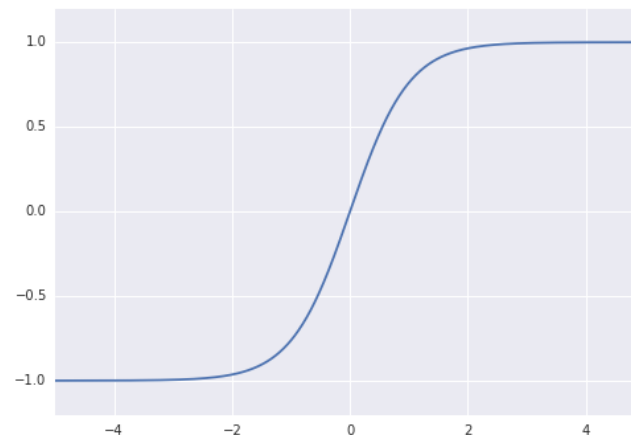
$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$

# Tanh

- Often used for hidden states & cells in RNNs, LSTMs
- Pro: differentiable, often converges faster than sigmoid
- Con: gradients easily saturate to zero => vanishing gradients

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$
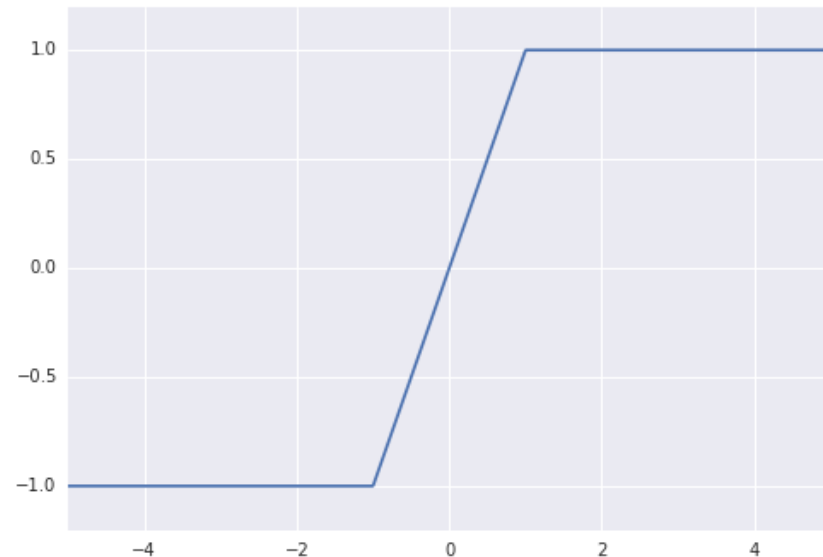
$$\tanh'(x) = 1 - \tanh^2(x)$$

$$\tanh(x) = 2\sigma(2x) - 1$$

# Hard Tanh

- Pro: computationally cheaper
- Con: saturates to zero easily, doesn't differentiate at 1, -1

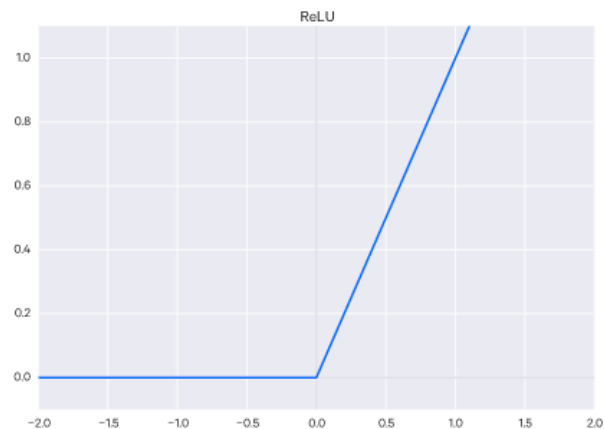$$\mathrm{hardtanh}(t) = \begin{cases} -1 & t < -1 \\ t & -1 \le t \le 1 \\ 1 & t > 1 \end{cases}$$

# ReLU

- Pro: doesn't saturate for x > 0, computationally cheaper, induces sparse NNs

- Con: non-differentiable at 0

- Used widely in deep NN, but not as much in RNNs

- We informally use subgradients:

$$\mathrm{ReLU}(x) = \max(0, x)$$

$$\frac{d\,\mathrm{ReLU}(x)}{dx} = \begin{cases} 1 & x > 0 \\ 0 & x < 0 \\ 1 \text{ or } 0 & o.w \end{cases}$$

# Vanishing / exploding Gradients

- Practical solutions w.r.t. numerical operations
  - Gradient Clipping: bound gradients by a max value
  - Gradient Normalization: renormalize gradients when they are above a fixed norm
  - Careful initialization, smaller learning rates
  - Avoid saturating nonlinearities (like tanh, sigmoid)
    - ReLU or hard-tanh instead
  - Batch Normalization: add intermediate input normalization layers

# Batch Normalization

**Input:** Values of $x$ over a mini-batch: $\mathcal{B} = \{x_{1...m}\}$;
Parameters to be learned: $\gamma, \beta$

**Output:** $\{y_i = \text{BN}_{\gamma,\beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^{m} x_i \qquad \text{// mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^{m} (x_i - \mu_{\mathcal{B}})^2 \qquad \text{// mini-batch variance}$$

$$\widehat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \qquad \text{// normalize}$$

$$y_i \leftarrow \gamma \widehat{x}_i + \beta \equiv \text{BN}_{\gamma,\beta}(x_i) \qquad \text{// scale and shift}$$

# Regularization

- Regularization by objective term

$$\mathcal{L}(\theta) = \sum_{i=1}^{n} \max\{0, 1 - (\hat{y}_c - \hat{y}_{c'})\} + \lambda ||\theta||^2$$

  - Modify loss with L1 or L2 norms

- Less depth, smaller hidden states, early stopping

- Dropout
  - Randomly delete parts of network during training
  - Each node (and its corresponding incoming and outgoing edges) dropped with a probability p
  - P is higher for internal nodes, lower for input nodes
  - The full network is used for testing
  - Faster training, better results
  - Vs. Bagging

# Convergence of backprop

- Without non-linearity or hidden layers, learning is convex optimization
  - Gradient descent reaches **global minima**
- Multilayer neural nets (with nonlinearity) are **not convex**
  - Gradient descent gets stuck in local minima
  - Selecting number of hidden units and layers =  fuzzy process
  - NNs have made a HUGE comeback in the last few years
    - Neural nets are back with a new name
      - Deep belief networks
      - Huge error reduction when trained with lots of data on GPUs

# RECAP

# Vanishing / exploding Gradients

- Deep networks are hard to train
- Gradients go through multiple layers
- The multiplicative effect tends to lead to *exploding* or *vanishing* gradients
- Practical solutions w.r.t.
  – network architecture
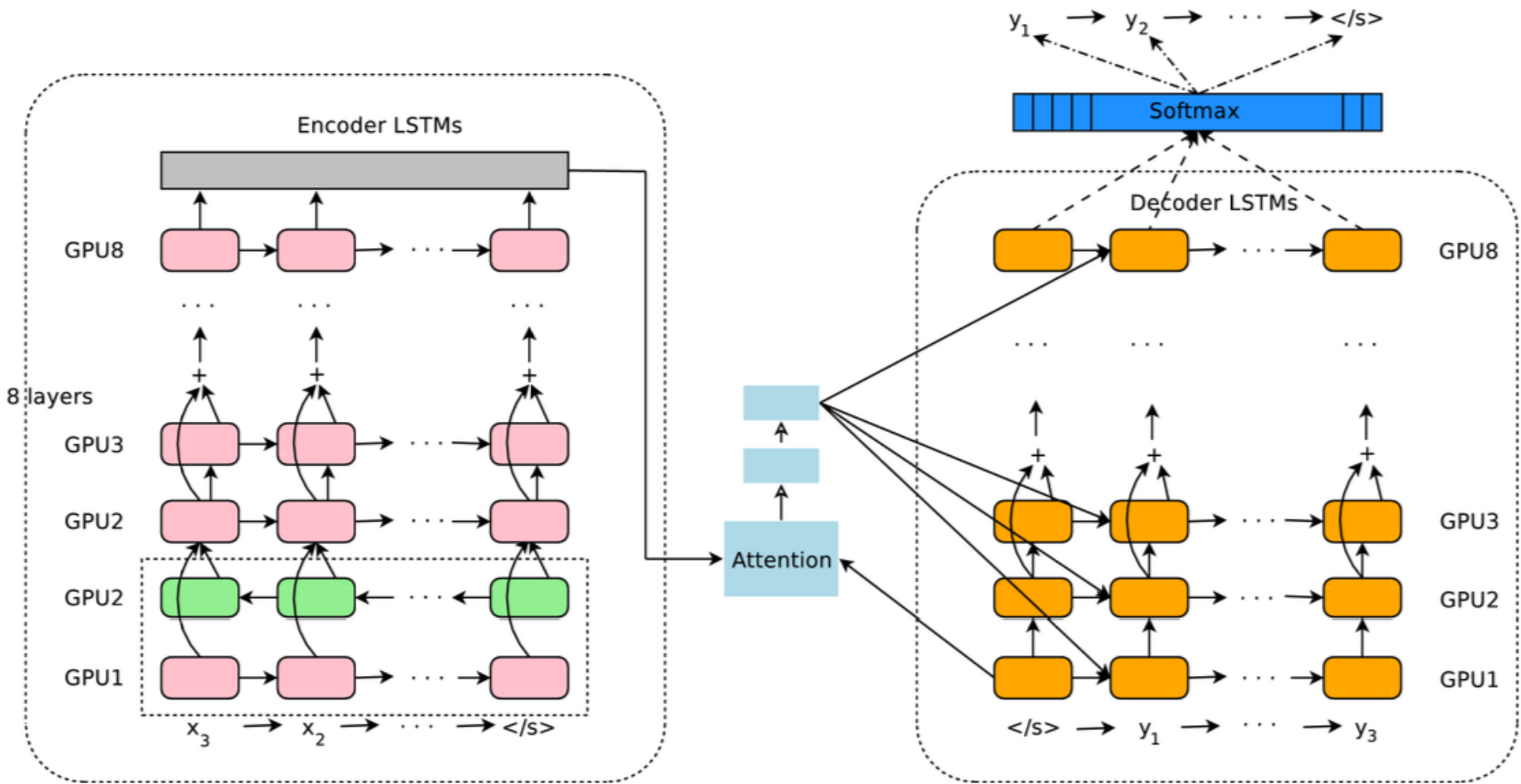  – numerical operations

# Vanishing / exploding Gradients

- Practical solutions w.r.t. network architecture
  - Add skip connections to reduce distance
    - Residual networks, highway networks, …
  - Add gates (and memory cells) to allow longer term memory
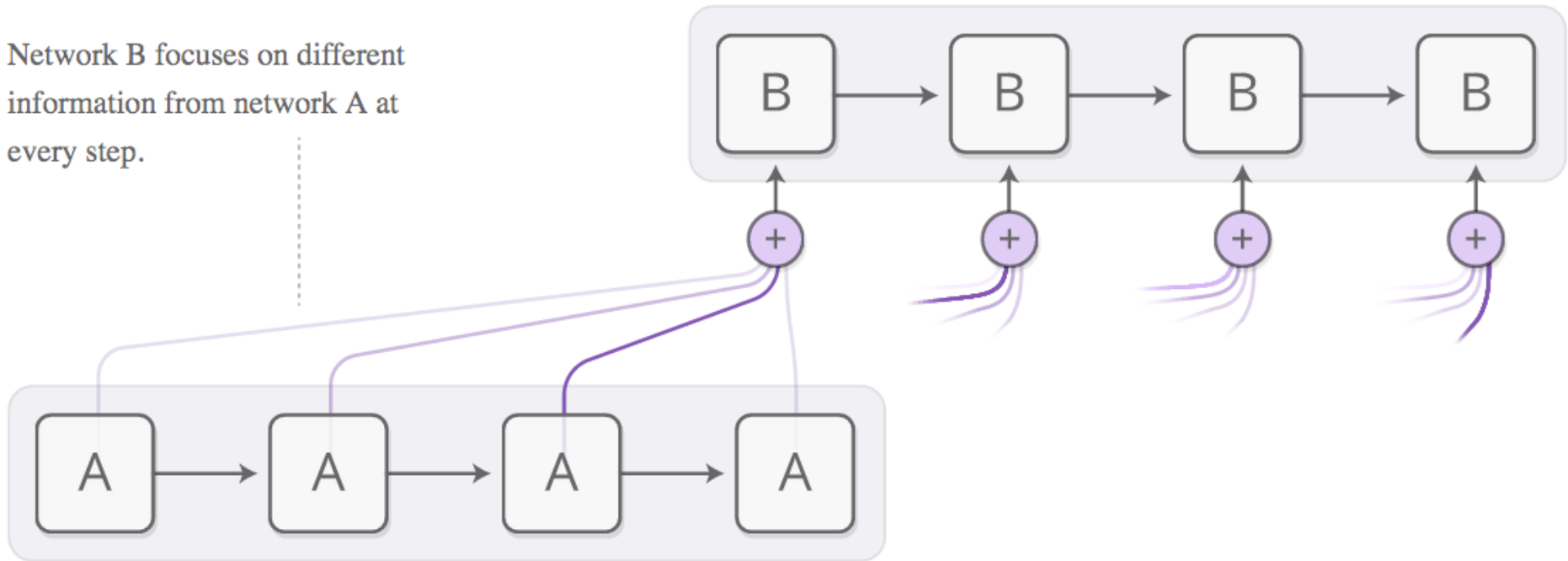    - LSTMs, GRUs, memory networks, …
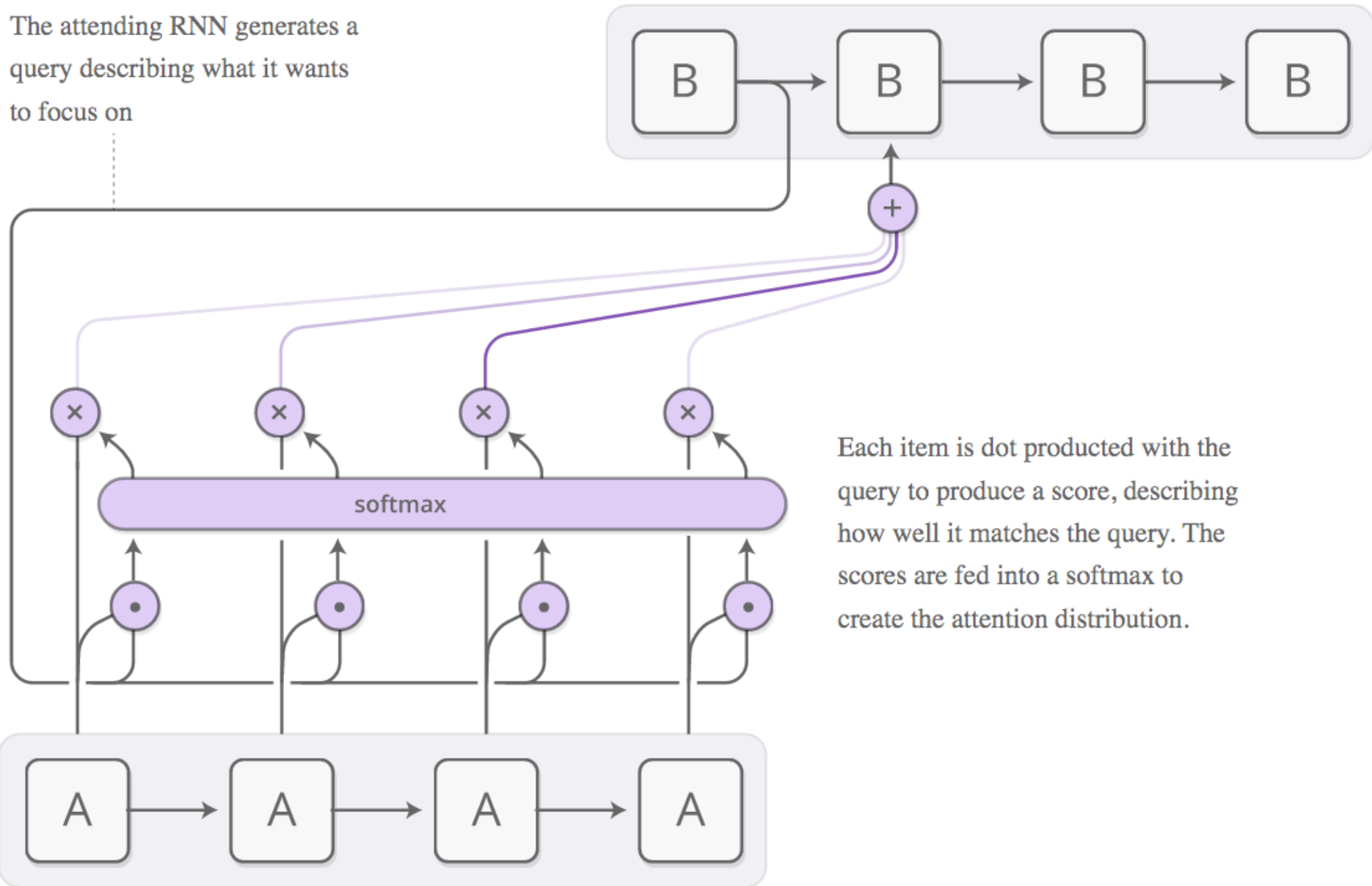
# seq2seq (aka "encoder-decoder")



LSTM Encoder

LSTM Decoder

# Google NMT (Oct 2016)

# ATTENTION!

# Seq-to-Seq with Attention

Network B focuses on different information from network A at every step.

The attending RNN generates a query describing what it wants to focus on

B → B → B → B

softmax

A → A → A → A

Each item is dot producted with the query to produce a score, describing how well it matches the query. The scores are fed into a softmax to create the attention distribution.

# Trial: Hard Attention

- At each step generating the target word $\mathbf{s}_i^t$
- Compute the best alignment to the source word $\mathbf{s}_j^s$
- And incorporate the source word to generate the target word

$$w_{i+1}^t = \operatorname{argmax}_w O(w, s_{i+1}^t, s_j^s)$$
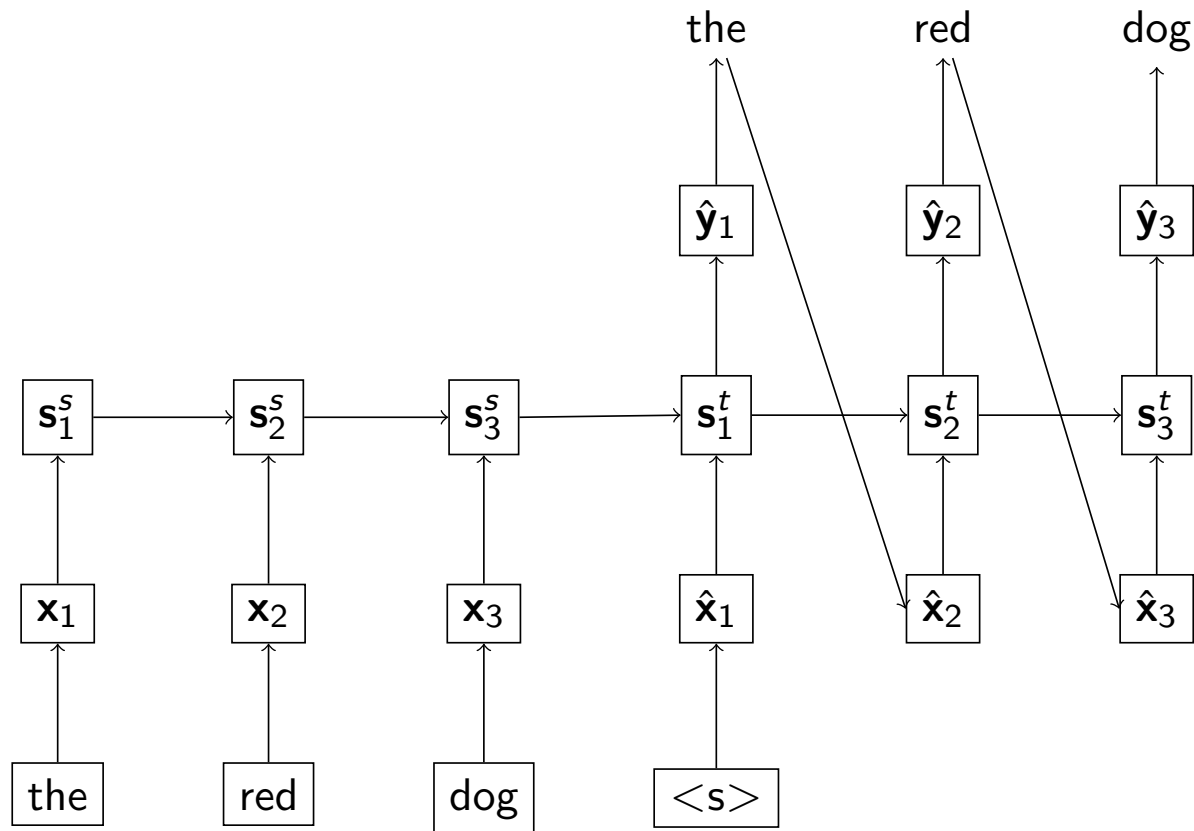
- Contextual *hard* alignment. How?

$$z_j = \tanh([s_i^t, s_j^s]W + b)$$

$$j = \operatorname{argmax}_j z_j$$

- Problem?

# Encoder – Decoder Architecture

Sequence-to-Sequence

# Attention: Soft Alignments

- At each step generating the target word $\mathbf{s}_i^t$
- Compute the attention $\mathbf{c}$ to the source sequence $\mathbf{s}^s$
- And incorporate the attention to generate the target word

$$w_{i+1}^t = \operatorname{argmax}_w O(w, s_{i+1}^t, c)$$

- Contextual *attention* as *soft* alignment. How?

$$z_j = \tanh([s_i^t, s_j^s]W + b)$$

$$\alpha = \operatorname{softmax}(z)$$

$$c = \sum_j \alpha_j s_j^s$$

  – Step-1: compute the attention weights
  – Step-2: compute the attention vector as interpolation

# Attention function parameterization

- Feedforward NNs

$$z_j = \tanh([s_i^t; s_j^s]W + b)$$
$$z_j = \tanh([s_i^t; s_j^s; s_i^t \circ s_j^s]W + b)$$

- Dot product

$$z_j = s_i^t \cdot s_j^s$$

- Cosine similarity

$$z_j = \frac{s_i^t \cdot s_j^s}{||s_i^t||\,||s_j^s||}$$

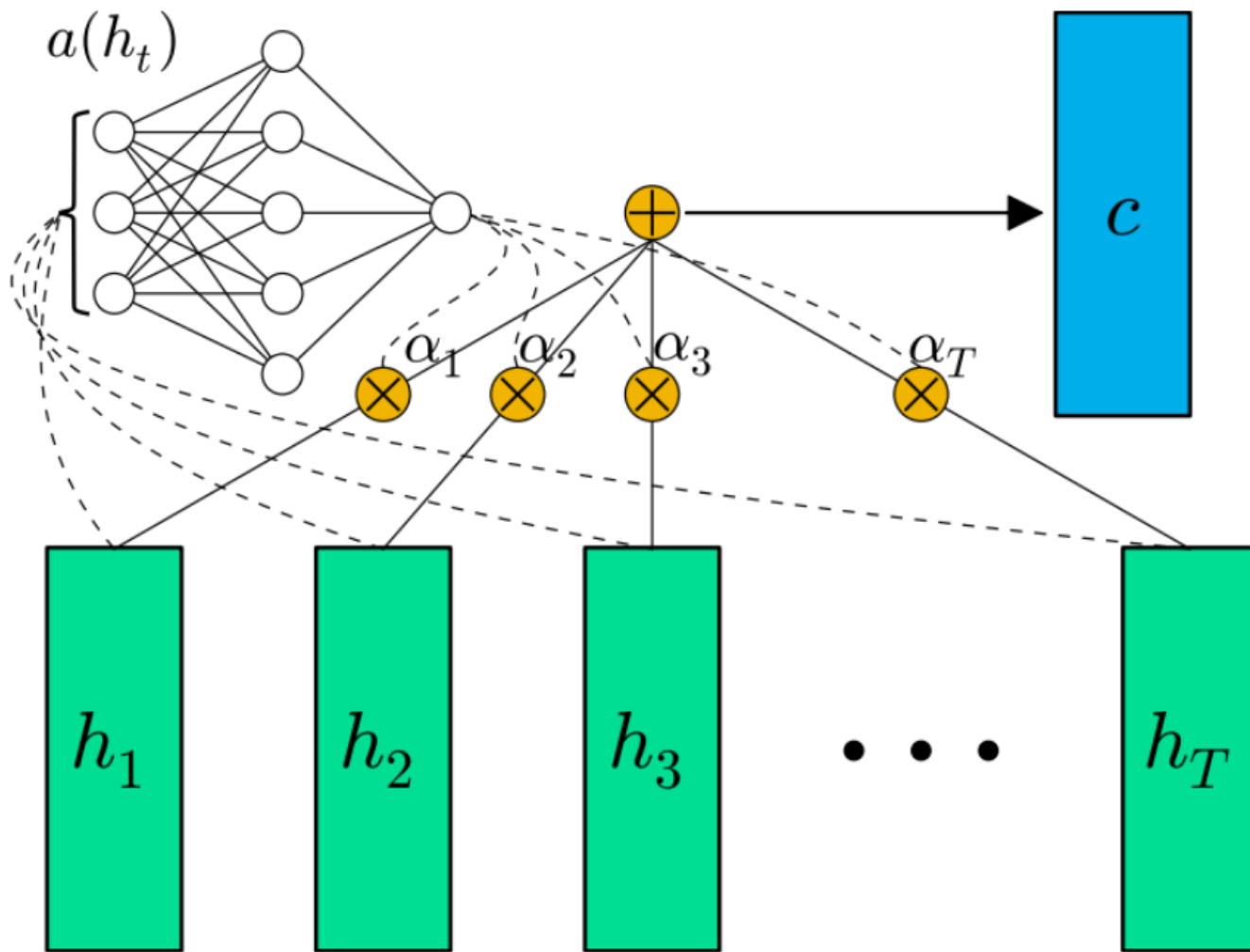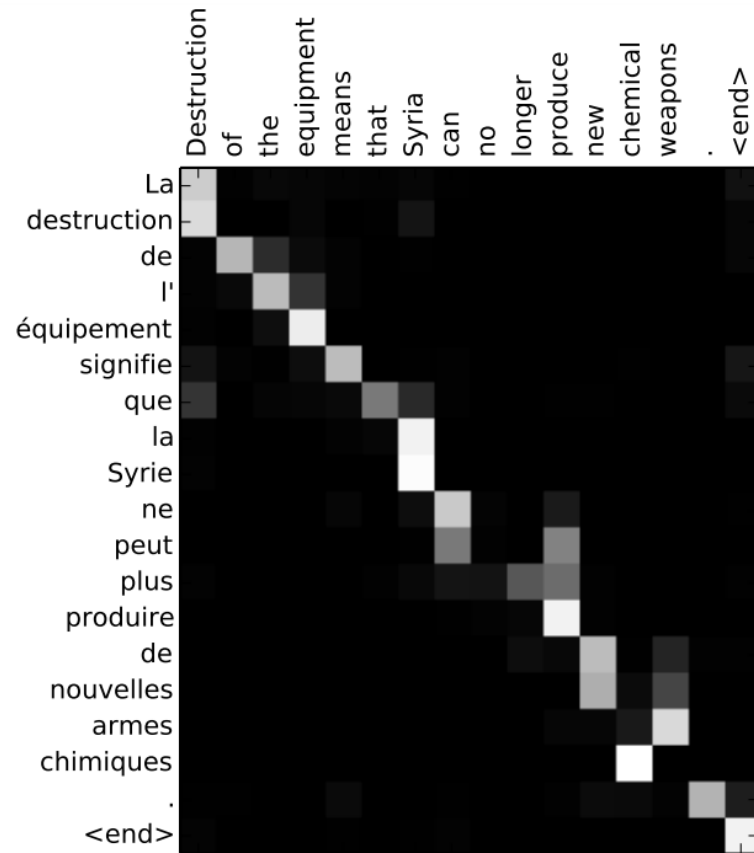- Bi-linear models

$$z_j = {s_i^t}^T W s_j^s$$

Figure 1: Schematic of our proposed "feed-forward" attention mechanism (cf. (Cho, 2015) Figure 1). Vectors in the hidden state sequence $h_t$ are fed into the learnable function $a(h_t)$ to produce a probability vector $\alpha$. The vector $c$ is computed as a weighted average of $h_t$, with weighting given by $\alpha$.

# Learned Attention!

# Qualitative results

*Figure 2.* Attention over time. As the model generates each word, its attention changes to reflect the relevant parts of the image. "soft" (top row) vs "hard" (bottom row) attention. (Note that both models generated the same captions in this example.)
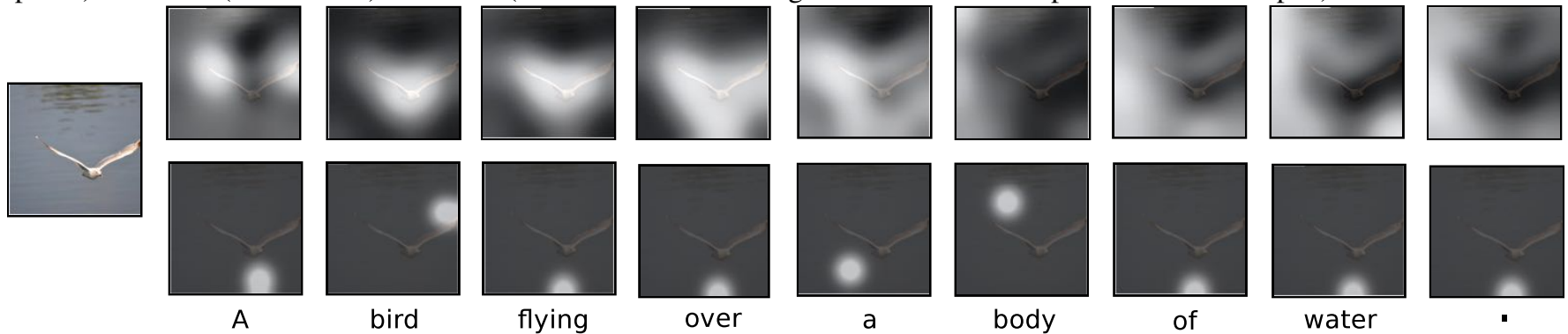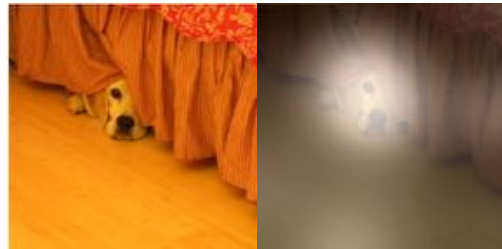


A    A    bird bird    flying flying    over over    a a    body body    of of    water water    . .

*Figure 3.* Examples of attending to the correct object (*white* indicates the attended regions, *underlines* indicated the corresponding word)
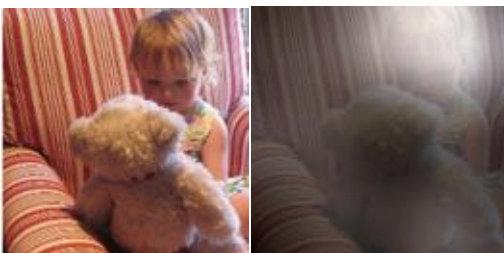


A woman is throwing a frisbee in a park.

A dog is standing on a hardwood floor.

A stop sign is on a road with a mountain in the background.

A little girl sitting on a bed with a teddy bear.

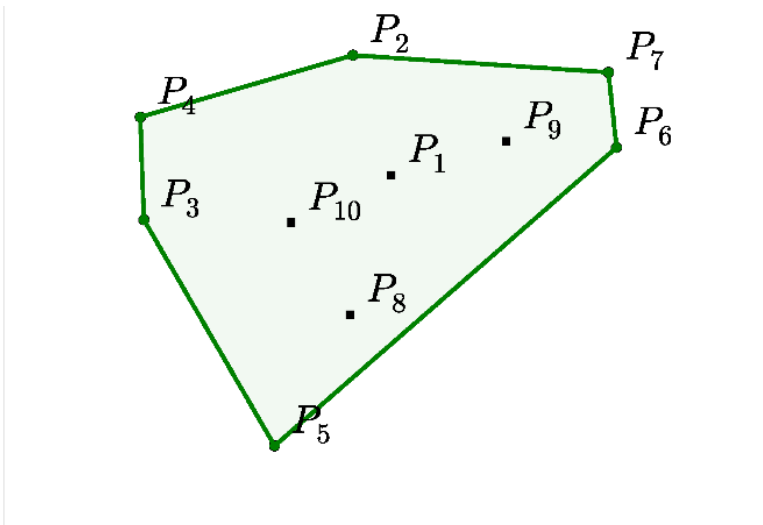A group of people sitting on a boat in the water.

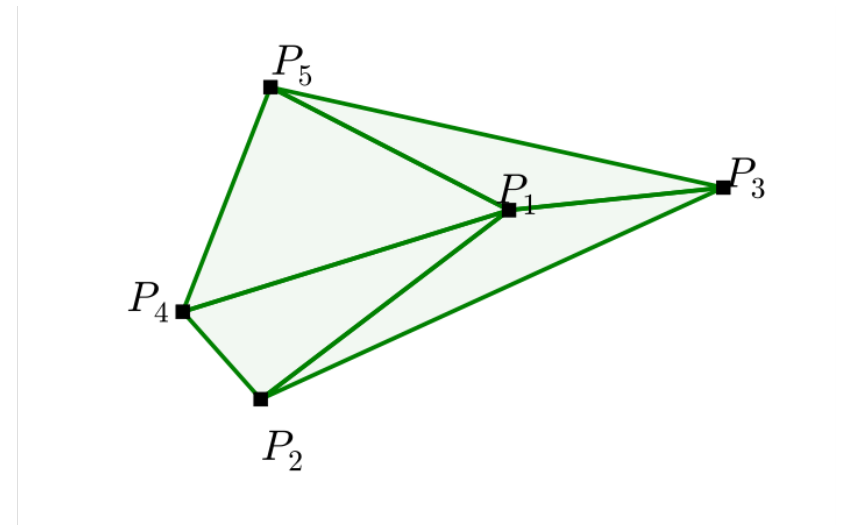A giraffe standing in a forest with trees in the background.

# POINTER NETWORKS

# Convex haul, Delaunay Triangulation, Traveling Salesman

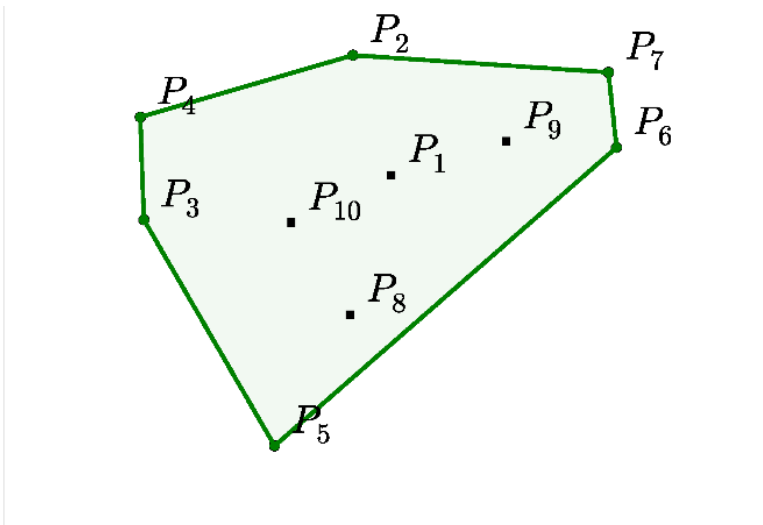## Can we model these problems using seq-to-seq?



(a) Input $\mathcal{P} = \{P_1, \ldots, P_{10}\}$, and the output sequence $\mathcal{C}^{\mathcal{P}} = \{\Rightarrow, 2, 4, 3, 5, 6, 7, 2, \Leftarrow\}$ representing its convex hull.

(b) Input $\mathcal{P} = \{P_1, \ldots, P_5\}$, and the output $\mathcal{C}^{\mathcal{P}} = \{\Rightarrow, (1, 2, 4), (1, 4, 5), (1, 3, 5), (1, 2, 3), \Leftarrow\}$ representing its Delaunay Triangulation.

# Pointer Networks! (Vinyals et al. 2015)

- NNs with attention: content-based attention to input
- Pointer networks: location-based attention to input



(a) Input $\mathcal{P} = \{P_1, \dots, P_{10}\}$, and the output sequence $\mathcal{C}^{\mathcal{P}} = \{\Rightarrow, 2, 4, 3, 5, 6, 7, 2, \Leftarrow\}$ representing its convex hull.
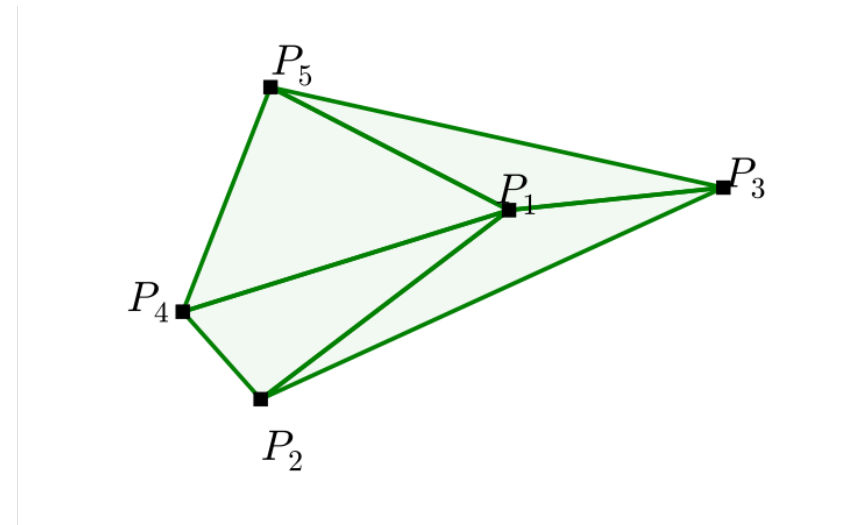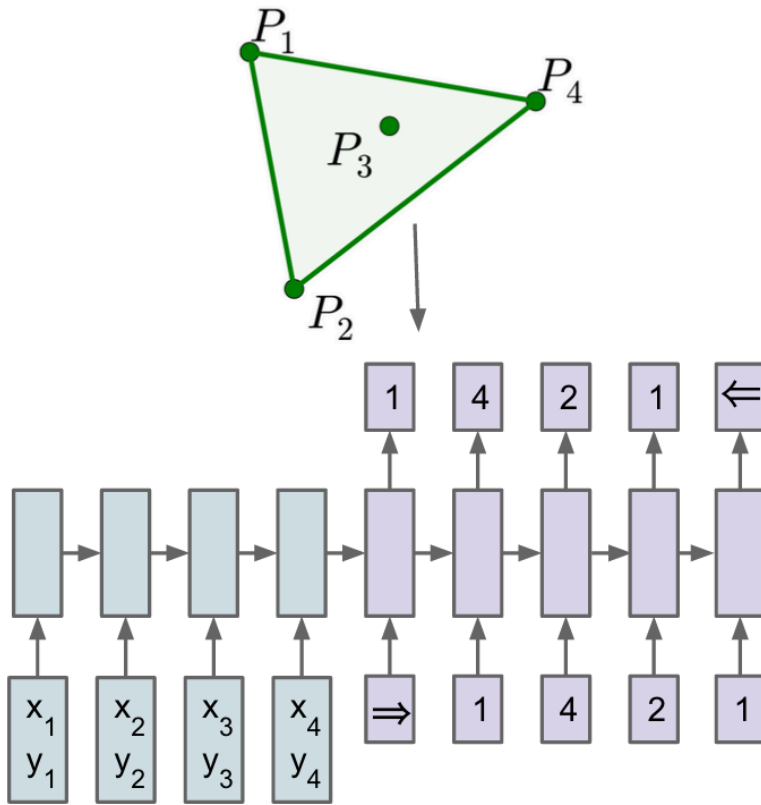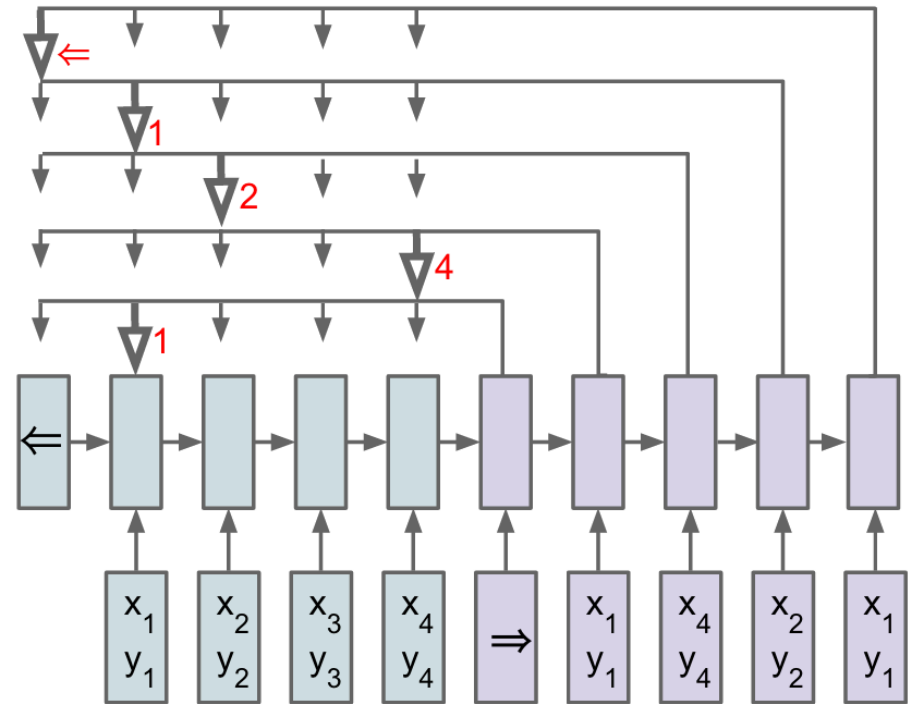
(b) Input $\mathcal{P} = \{P_1, \dots, P_5\}$, and the output $\mathcal{C}^{\mathcal{P}} = \{\Rightarrow, (1, 2, 4), (1, 4, 5), (1, 3, 5), (1, 2, 3), \Leftarrow\}$ representing its Delaunay Triangulation.

# Pointer Networks



(a) Sequence-to-Sequence

(b) Ptr-Net

# Pointer Networks

## Attention Mechanism vs Pointer Networks

$$e_{ij} = v_a^\top \tanh\left(W_a s_{i-1} + U_a h_j\right)$$

$$\alpha_{ij} = \frac{\exp\left(e_{ij}\right)}{\sum_{k=1}^{T_x} \exp\left(e_{ik}\right)}$$

$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j$$

Attention mechanism

$$e_{ij} = v_a^\top \tanh\left(W_a s_{i-1} + U_a h_j\right)$$

$$p(C_i \mid C_1, \ldots, C_{i-1}, \mathcal{P}) = \frac{\exp\left(e_{ij}\right)}{\sum_{k=1}^{T_x} \exp\left(e_{ik}\right)}$$

Ptr-Net

Softmax normalizes the vector e$_{ij}$ to be an output distribution over the dictionary of inputs

Diagram borrowed from Keon Kim

# CopyNet (Gu et al. 2016)

- Conversation
  - I: Hello Jack, my name is Chandralekha
  - R: Nice to meet you, Chandralekha
  - I: This new guy doesn't perform exactly as expected.
  - R: what do you mean by "doesn't perform exactly as expected?"
- Translation

# CopyNet (Gu et al. 2016)



**(b) Generate-Mode & Copy-Mode**

Prob(**"Jebara"**) = Prob(**"Jebara"**, g) + Prob(**"Jebara"**, c)

*Softmax*

**Vocabulary**                    **Source**

**M**

**$s_4$**

**(c) State Update**

DNN

Embedding for "Tony"

Selective Read for "Tony"

"Tony"

**M**

$\rho$

**Attentive Read**

hi , Tony Jebara

$s_1$ → $s_2$ → $s_3$ → $s_4$

<eos> hi , Tony

$h_1$ ← $h_2$ ← $h_3$ ← $h_4$ ← $h_5$ ← $h_6$ ← $h_7$ ← $h_8$

hello , my name is Tony Jebara .

**(a) Attention-based Encoder-Decoder (RNNSearch)**

# CopyNet (Gu et al. 2016)

- Key idea: interpolation between generation model & copy model

$$p(y_t|\mathbf{s}_t, y_{t-1}, \mathbf{c}_t, \mathbf{M}) = p(y_t, \mathsf{g}|\mathbf{s}_t, y_{t-1}, \mathbf{c}_t, \mathbf{M})$$
$$+ p(y_t, \mathsf{c}|\mathbf{s}_t, y_{t-1}, \mathbf{c}_t, \mathbf{M}) \quad (4)$$

$$p(y_t, \mathsf{g}|\cdot) = \begin{cases} \frac{1}{Z} e^{\psi_g(y_t)}, & y_t \in \mathcal{V} \\ 0, & y_t \in \mathcal{X} \cap \bar{V} \ (5) \\ \frac{1}{Z} e^{\psi_g(\text{UNK})} & y_t \notin \mathcal{V} \cup \mathcal{X} \end{cases}$$

$$p(y_t, \mathsf{c}|\cdot) = \begin{cases} \frac{1}{Z} \sum_{j:x_j=y_t} e^{\psi_c(x_j)}, & y_t \in \mathcal{X} \\ 0 & \text{otherwise} \end{cases} \ (6)$$

**Generate-Mode:** The same scoring function as in the generic RNN encoder-decoder (Bahdanau et al., 2014) is used, i.e.
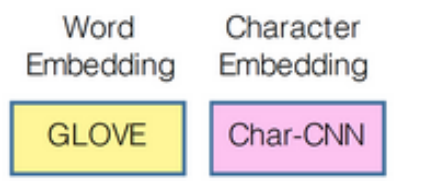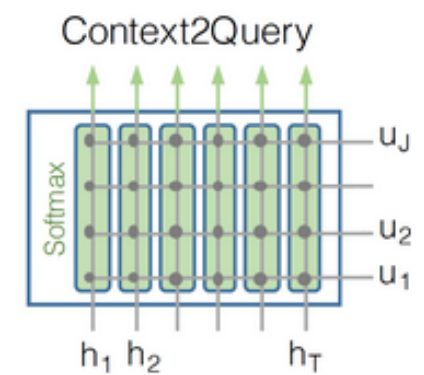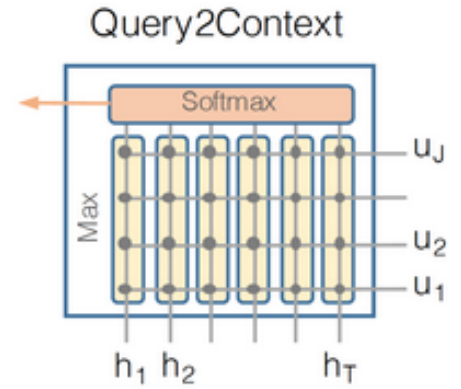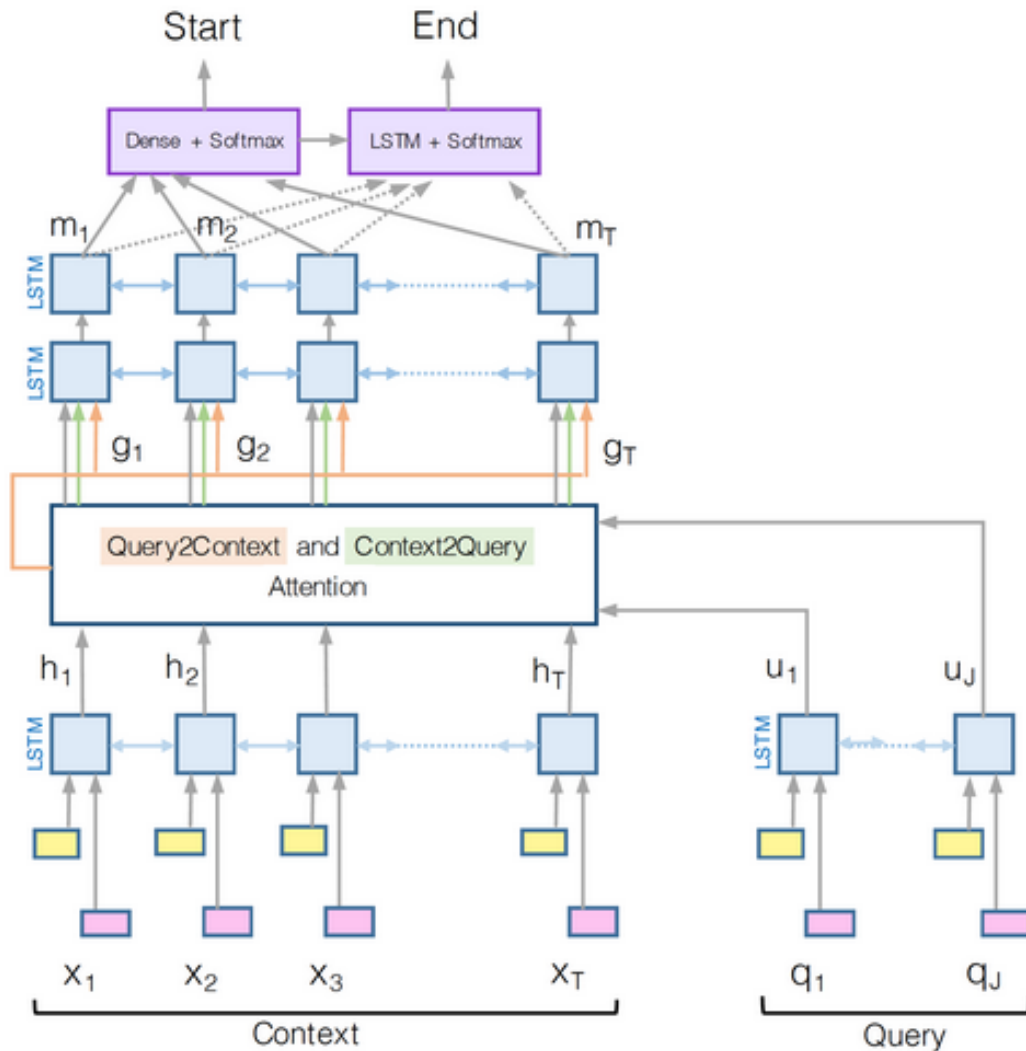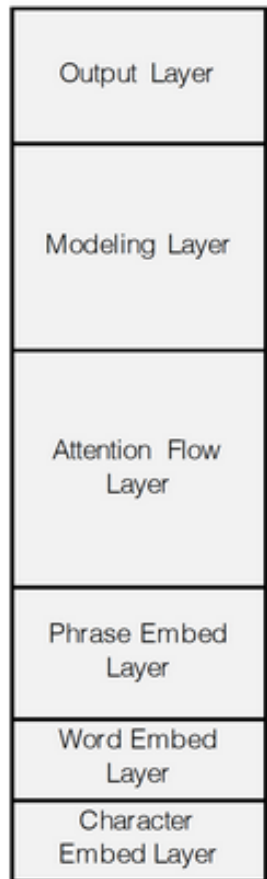
$$\psi_g(y_t = v_i) = \mathbf{v}_i^\top \mathbf{W}_o \mathbf{s}_t, \quad v_i \in \mathcal{V} \cup \text{UNK} \quad (7)$$

where $\mathbf{W}_o \in \mathbb{R}^{(N+1) \times d_s}$ and $\mathbf{v}_i$ is the one-hot indicator vector for $v_i$.

**Copy-Mode:** The score for "copying" the word $x_j$ is calculated as

$$\psi_c(y_t = x_j) = \sigma\left(\mathbf{h}_j^\top \mathbf{W}_c\right) \mathbf{s}_t, \quad x_j \in \mathcal{X} \quad (8)$$
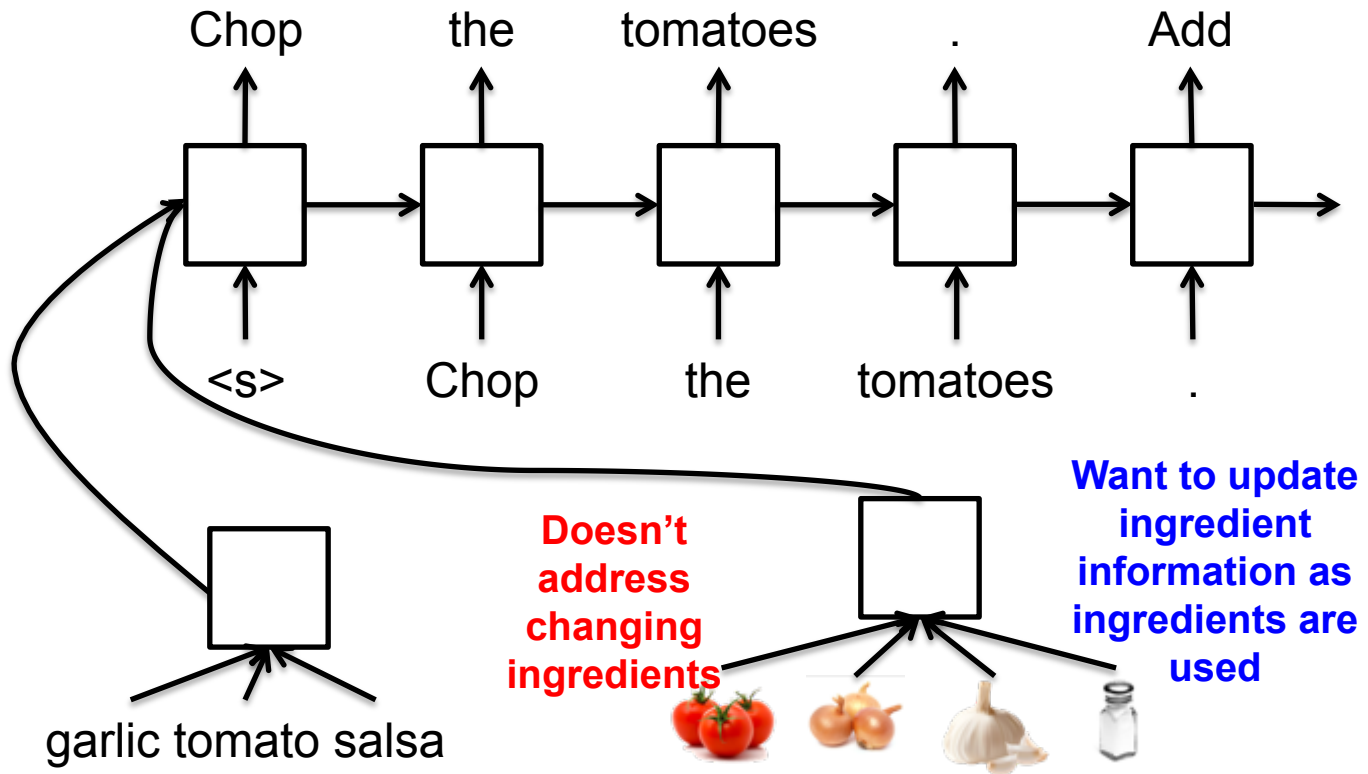
# BiDAF

# NEURAL CHECK LIST

# Neural Checklist Models
## (Kiddon et al., 2016)

- What can we do with gating & attention?
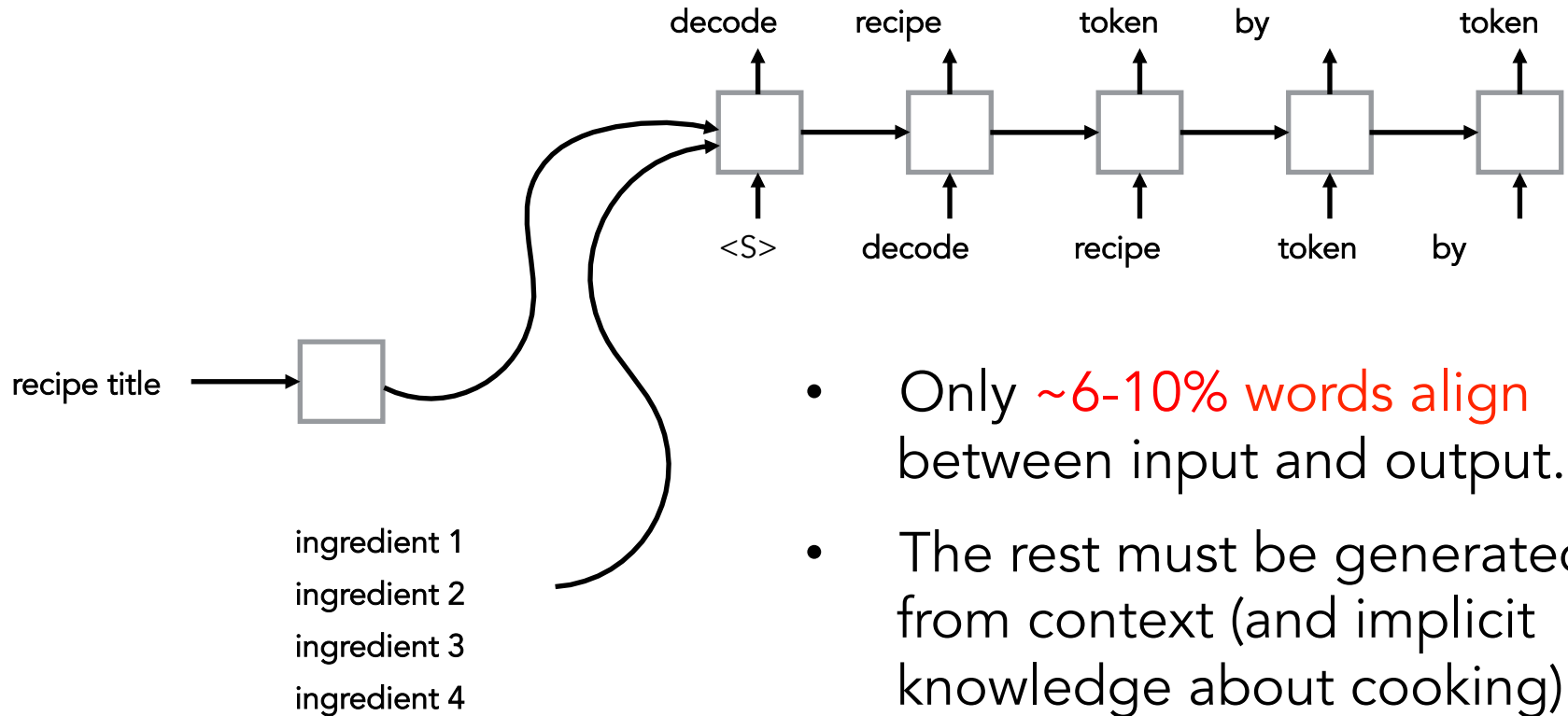
# Encoder--Decoder Architecture

# Encode title - decode recipe

sausage sandwiches  ➡️  Cut each sandwich in halves.
Sandwiches with sandwiches.
Sandwiches, sandwiches, Sandwiches,
sandwiches, sandwiches
sandwiches, sandwiches, sandwiches,
sandwiches, sandwiches, sandwiches, or
sandwiches or triangles, a griddle, each
sandwich.
Top each with a slice of cheese, tomato,
and cheese.
Top with remaining cheese mixture.
Top with remaining cheese.
Broil until tops are bubbly and cheese is
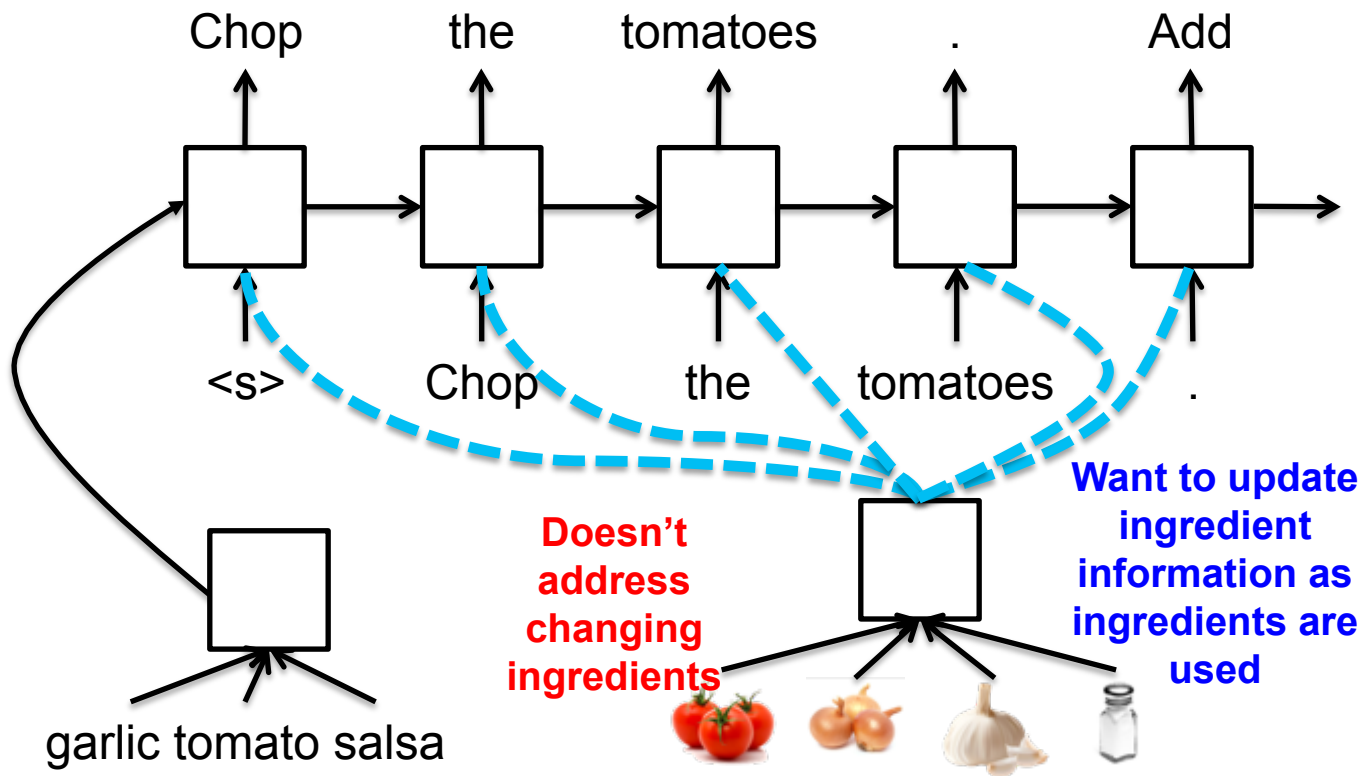melted, about 5 minutes.

# Recipe generation **vs** machine translation

decode    recipe    token    by    token

<S>    decode    recipe    token    by

recipe title

ingredient 1
ingredient 2
ingredient 3
ingredient 4

- Only ~6-10% words align between input and output.

- The rest must be generated from context (and implicit knowledge about cooking)

- Contextual switch between two different input sources

Two input sources

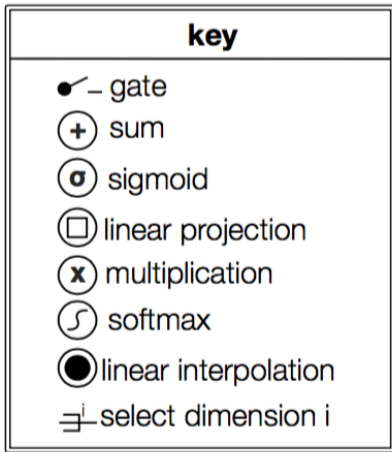# Encoder--Decoder with Attention

Chop    the    tomatoes    .    Add

<s>    Chop    the    tomatoes    .

garlic tomato salsa

**Doesn't address changing ingredients**

**Want to update ingredient information as ingredients are used**

# Neural checklist model

# Let's make salsa!

**Garlic tomato salsa**

tomatoes
onions
garlic
salt

# Neural checklist model

# Neural checklist model

Chop        the        tomatoes        .

| 0.85 |
| 0.10 |
| 0.04 |
| 0.01 |

**non-ingredient**    **new ingredient**

<S>        Chop        the        tomatoes

# Neural checklist model

# Neural checklist model

Add          to          tomatoes                    .

| | |
|---|---|
| 0.94 | 🍅 |
| 0.04 | 🧅 |
| 0.01 | 🧄 |
| 0.01 | 🧂 |

**used ingredient**

.          Add          to          tomatoes

# Checklist is probabilistic

Add      to      tomatoes

| | | |
|---|---|---|
| 0.90 | 🍅 | 0.85 |
| 0.08 | 🧅 | 1.00 |
| 0.01 | 🧄 | 0.04 |
| 0.01 | 🧂 | 0.02 |

**used ingredient**

$$\boldsymbol{\alpha}_t^{new} = \text{new ingredient prob. distribution}$$

$$\mathbf{a}_t^{new} = P(\bullet | \mathbf{h}_t) \cdot \boldsymbol{\alpha}_t^{new}$$

$$\mathbf{a}_{t+1} = \mathbf{a}_t + \mathbf{a}_t^{new}$$

.      Add      to      tomatoes

$\mathbf{a}_t$

$\mathbf{a}_{t+1}$

| 0.85 | 1.00 | 0.04 | 0.02 |
|---|---|---|---|

| 0.85 | 1.00 | 0.04 | 0.02 |
|---|---|---|---|

# Hidden state classifier is soft



tomatoes

Add    to    .

| 0.00 | 🍅 | 0.85 | | 0.90 | 🍅 | 0.85 |
| 0.00 | 🧅 | 1.00 | | 0.08 | 🧅 | 1.00 |
| 0.50 | 🧄 | 0.04 | | 0.01 | 🧄 | 0.04 |
| 0.50 | 🧂 | 0.02 | | 0.01 | 🧂 | 0.02 |

0.01    0.05    0.94

.    Add    to    tomatoes

| 0.85 | 1.00 | 0.04 | 0.02 |

| 0.85 | 1.00 | 0.04 | 0.02 |

# Interpolation



probability distribution over vocabulary

Attention model over available ingredients

Attention model over used ingredients

0.01

0.05

0.94

$$W_o \in \mathbb{R}^{|V| \times k}$$

$$\mathbf{w_t} = \mathrm{softmax}(W_o \mathbf{b_t})$$

$$\mathbf{o}_t = P(\bullet | \mathbf{h}_t) \mathbf{c}_t^{LM}$$
$$+ P(\bullet | \mathbf{h}_t) \mathbf{c}_t^{new}$$
$$+ P(\bullet | \mathbf{h}_t) \mathbf{c}_t^{used}$$

$$\mathbf{c}_t^{LM} = W_h \mathbf{h}_t$$

$$W_h \in \mathbb{R}^{k \times k}$$

# Choose ingredient via attention

$$\boldsymbol{\alpha}_t^{new} = \text{softmax}(\gamma E_t^{new} \mathbf{c}_t^{LM})$$

temperature term

available ingredient
embeddings

content vector from
language model

available ingredient embeddings
$E_t^{new}$

$\boldsymbol{\alpha}_t^{new}$

0.75

hidden state

?

0.01

Attention models for other NLP tasks
MT (Balasubramanian et al. 13,
            Bahdanau et al. 14)
        Sentence summarization (Rush et
        al. 15)
        Machine reading (Cheng et al.
        16)
        Image captioning (Xu et al. 15)

0.24

# Attention-generated embeddings

Can generate an embedding from the attention probabilities

ingredient embeddings
$E^T$

$\boldsymbol{\alpha}_t^{new}$

0.75

attention embedding

0.01

0.24

$\mathbf{c}_t^{new}$

$$\mathbf{c}_t^{new} = E^T \boldsymbol{\alpha}_t^{new}$$

# Discussion Points

- Strength and challenges of deep learning?

  *… what do NNs think about this?*

# Hafez: Neural Sonnet Writer
## (Ghazvininejad et al. 2016)

# Neural Sonnets

Deep Convolution Network
   Outrageous channels on the wrong *connections,*
   An empty space without an open *layer,*
   A closet full of black and blue *extensions,*
   Connections by the *closure operator.*

Theory
   Another way to reach the wrong *conclusion!*
   A vision from a total *transformation,*
   Created by the great *magnetic fusion,*
   Lots of people need an *explanation.*

# Discussion Points

- Strength and challenges of deep learning?

- Representation learning
  - Less efforts on feature engineering (at the cost of more hyperparameter tuning!)
  - In computer vision: NN learned representation is significantly better than human engineered features
  - In NLP: often NN induced representation is concatenated with additional human engineered features.
- Data
  - Most success from massive amount of clean (expensive) data
  - Recent surge of data creation type papers (especially AI challenge type tasks)
  - Which significantly limits the domains & applications
  - Need stronger models for unsupervised & distantly supervised approaches

# Discussion Points

- Strength and challenges of deep learning?

- Architecture
  - allows for flexible, expressive, and creative modeling

- Easier entry to the field
  - Recent breakthrough from engineering advancements than theoretic advancements
  - Several NN platforms, code sharing culture

# Neural Recipe Example #1

In a small bowl , combine the cheese , eggplant , basil , oregano , tomato sauce and onion . Mix well .
Shape mixture into 6 patties , each about 3/4-inch thick.
Place on baking sheet .
Bake at 350 degrees for 30 minutes or until lightly browned .
Southern living magazine , sometime in 1980 .
Typed for you by nancy coleman .

eggplant
cheese cottage
lowfat
chopped onion
bay ground leaf
basil
oregano
tomato sauce
provolone

Cook eggplant in boiling water , covered , for 10 min .
Drain and cut in half lengthwise . scoop out insides leaving 1/2 '' shell . Mash insides with cottage cheese onion , bay leaf , basil , oregano and tomato sauce .
Preheat oven to 350 ^ stuff eggplant halves , place in casserole dish and bake covered for 15 min .
Add a little water to bottom of pan to keep eggplant moist . top with provolone cheese .
Bake 5 more min uncovered 1 serving =

# CONVOLUTION NEURAL NETWORK

Next several slides borrowed from Alex Rush

# Models with Sliding Windows

- Classification/prediction with sliding windows
  - E.g., neural language model
- Feature representations with sliding window
  - E.g., sequence tagging with CRFs or structured perceptron

$$\begin{bmatrix} w_1 & w_2 & w_3 & w_4 & w_5 \end{bmatrix} w_6 \ w_7 \ w_8$$

$$w_1 \begin{bmatrix} w_2 & w_3 & w_4 & w_5 & w_6 \end{bmatrix} w_7 \ w_8$$

$$w_1 \ w_2 \begin{bmatrix} w_3 & w_4 & w_5 & w_6 & w_7 \end{bmatrix} w_8$$

$$\vdots$$

# Sliding Windows w/ Convolution

Let our input be the embeddings of the full sentence, $\mathbf{X} \in \mathbb{R}^{n \times d^0}$

$$\mathbf{X} = [v(w_1), v(w_2), v(w_3), \dots, v(w_n)]$$

Define a window model as $NN_{window} : \mathbb{R}^{1 \times (d_{\mathrm{win}} d^0)} \mapsto \mathbb{R}^{1 \times d_{\mathrm{hid}}}$,

$$NN_{window}(\mathbf{x}_{win}) = \mathbf{x}_{win}\mathbf{W}^1 + \mathbf{b}^1$$

The convolution is defined as $NN_{conv} : \mathbb{R}^{n \times d^0} \mapsto \mathbb{R}^{(n - d_{\mathrm{win}} + 1) \times d_{\mathrm{hid}}}$,

$$NN_{conv}(\mathbf{X}) = \tanh \begin{bmatrix} NN_{window}(\mathbf{X}_{1:d_{\mathrm{win}}}) \\ NN_{window}(\mathbf{X}_{2:d_{\mathrm{win}}+1}) \\ \vdots \\ NN_{window}(\mathbf{X}_{n-d_{\mathrm{win}}:n}) \end{bmatrix}$$

# Pooling Operations

▶ Pooling "over-time" operations $f : \mathbb{R}^{n \times m} \mapsto \mathbb{R}^{1 \times m}$

1. $f_{max}(\mathbf{X})_{1,j} = \max_i X_{i,j}$
2. $f_{min}(\mathbf{X})_{1,j} = \min_i X_{i,j}$
3. $f_{mean}(\mathbf{X})_{1,j} = \sum_i X_{i,j} / n$

$$
f(\mathbf{X}) = \begin{bmatrix} \Downarrow & \Downarrow & \dots \\ \Downarrow & \Downarrow & \dots \\ & \vdots & \\ \Downarrow & \Downarrow & \dots \end{bmatrix} = [\ \dots\ ]
$$

# Convolution + Pooling

$$\hat{y} = \text{softmax}(f_{max}(NN_{conv}(\mathbf{X}))\mathbf{W}^2 + \mathbf{b}^2)$$

▶ $\mathbf{W}^2 \in \mathbb{R}^{d_{\text{hid}} \times d_{\text{out}}}$, $\mathbf{b}^2 \in \mathbb{R}^{1 \times d_{\text{out}}}$

▶ Final linear layer $\mathbf{W}^2$ uses learned window features

# Multiple Convolutions

$$\hat{y} = \text{softmax}([f(NN^1_{conv}(\mathbf{X})), f(NN^2_{conv}(\mathbf{X})), \ldots, f(NN^f_{conv}(\mathbf{X}))]\mathbf{W}^2 + \mathbf{b}^2)$$

- ▶ Concat several convolutions together.

- ▶ Each $NN^1$, $NN^2$, etc uses a different $d_{\text{win}}$

- ▶ Allows for different window-sizes (similar to multiple n-grams)

# Convolution Diagram (kim 2014)



wait for the video and do n't rent it

*n x k* representation of sentence with static and non-static channels

Convolutional layer with multiple filter widths and feature maps

Max-over-time pooling

Fully connected layer with dropout and softmax output

▶ $n = 9$, $d_{\text{hid}} = 4$ , $d_{\text{out}} = 2$

▶ red- $d_{\text{win}} = 2$, blue- $d_{\text{win}} = 3$, (ignore back channel)

# Text Classification (Kim 2014)

| Model | MR | SST-1 | SST-2 | Subj | TREC | CR | MPQA |
|---|---|---|---|---|---|---|---|
| CNN-rand | 76.1 | 45.0 | 82.7 | 89.6 | 91.2 | 79.8 | 83.4 |
| CNN-static | 81.0 | 45.5 | 86.8 | 93.0 | 92.8 | 84.7 | **89.6** |
| CNN-non-static | **81.5** | 48.0 | 87.2 | 93.4 | 93.6 | 84.3 | 89.5 |
| CNN-multichannel | 81.1 | 47.4 | **88.1** | 93.2 | 92.2 | **85.0** | 89.4 |
| RAE (Socher et al., 2011) | 77.7 | 43.2 | 82.4 | – | – | – | 86.4 |
| MV-RNN (Socher et al., 2012) | 79.0 | 44.4 | 82.9 | – | – | – | – |
| RNTN (Socher et al., 2013) | – | 45.7 | 85.4 | – | – | – | – |
| DCNN (Kalchbrenner et al., 2014) | – | 48.5 | 86.8 | – | 93.0 | – | – |
| Paragraph-Vec (Le and Mikolov, 2014) | – | **48.7** | 87.8 | – | – | – | – |

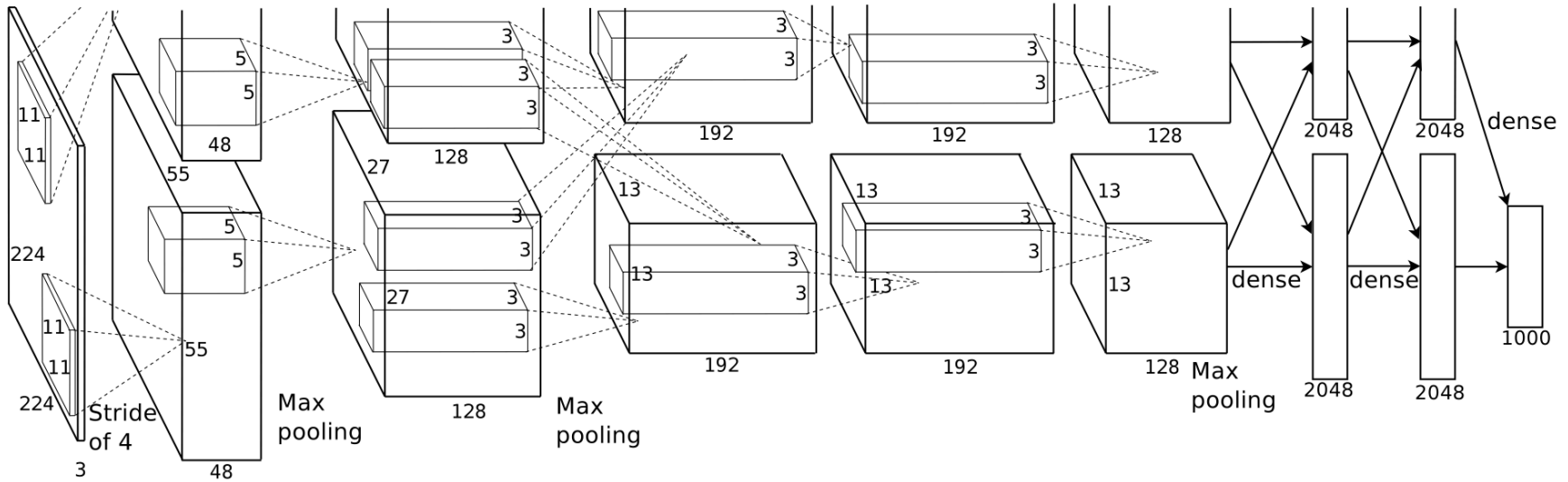# AlexNet (krizhevsky et al., 2012)



Figure 2: An illustration of the architecture of our CNN, explicitly showing the delineation of responsibilities between the two GPUs. One GPU runs the layer-parts at the top of the figure while the other runs the layer-parts at the bottom. The GPUs communicate only at certain layers. The network's input is 150,528-dimensional, and the number of neurons in the network's remaining layers is given by 253,440–186,624–64,896–64,896–43,264–4096–4096–1000.