*CS 1699: Deep Learning*
# Modeling Sequences: Recurrent Neural Networks, Transformers

Prof. Adriana Kovashka
University of Pittsburgh
March 24, 2020

# Plan for this lecture

- Recurrent neural networks
  - Basics
  - Training (backprop through time, vanishing gradient)
  - Recurrent networks with gates (GRU, LSTM)
- Applications in NLP and vision
  - Image/video captioning
  - Neural machine translation (beam search, attention)
- Transformers
  - Self-attention
  - BERT
  - Cross-modal transformers for VQA and VCR

# Recurrent neural networks

# Some pre-RNN captioning results



This is a picture of one sky, one road and one sheep. The gray sky is over the gray road. The gray sheep is by the gray road.



Here we see one road, one sky and one bicycle. The road is near the blue sky, and near the colorful bicycle. The colorful bicycle is within the blue sky.



This is a picture of two dogs. The first dog is near the second furry dog.

# Results with Recurrent Neural Networks



"man in black shirt is playing guitar."

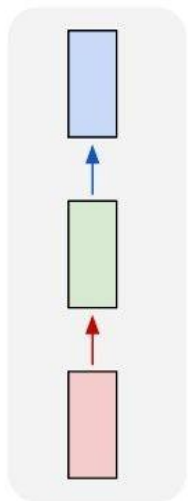"construction worker in orange safety vest is working on road."
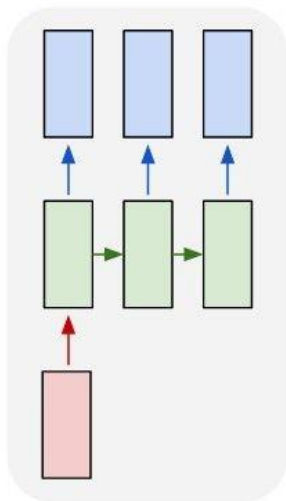
"two young girls are playing with lego toy."

"boy is doing backflip on wakeboard."

Karpathy and Fei-Fei, CVPR 2015

# Recurrent Networks offer a lot of flexibility:

| one to one | one to many | many to one | many to many | many to many |
|---|---|---|---|---|

**vanilla neural networks**

Andrej Karpathy

# Recurrent Networks offer a lot of flexibility:

| one to one | one to many | many to one | many to many | many to many |

e.g. **image captioning**
image -> sequence of words
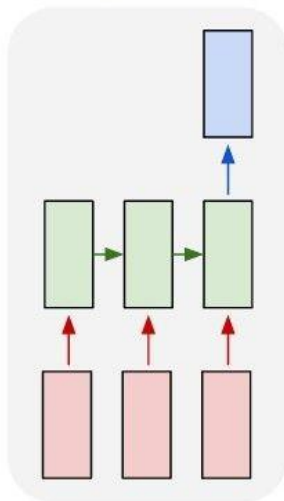
# Recurrent Networks offer a lot of flexibility:



one to one     one to many     many to one     many to many     many to many

e.g. **sentiment classification**
sequence of words -> sentiment

# Recurrent Networks offer a lot of flexibility:



e.g. **machine translation**
seq of words -> seq of words

# Recurrent Networks offer a lot of flexibility:



one to one     one to many     many to one     many to many     many to many
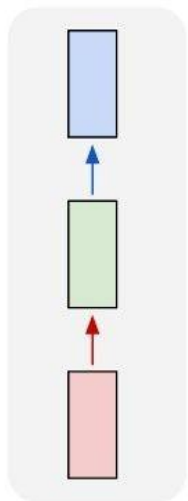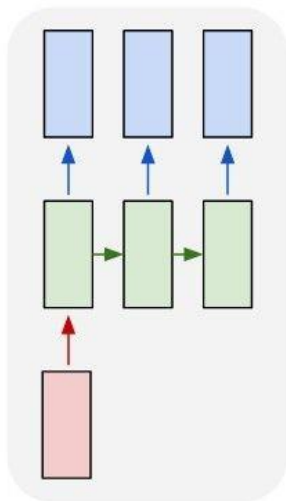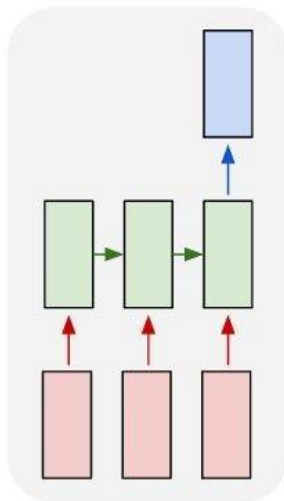
e.g. **video classification on frame level**

# Recurrent Neural Network

# Recurrent Neural Network



y

RNN

x

usually want to output a prediction at some time steps

# Recurrent Neural Network

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

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

new state — some function with parameters W — old state — input vector at some time step

y

RNN

x

# Recurrent Neural Network

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

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

Notice: the same function and the same set
of parameters are used at every time step.



Andrej Karpathy

# (Vanilla) Recurrent Neural Network

The state consists of a single *"hidden"* vector **h**:

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

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

$$y_t = W_{hy} h_t$$

y

RNN

x

# Example

**Character-level language model example**

Vocabulary:
[h,e,l,o]

Example training sequence:
**"hello"**

# Example

**Character-level language model example**

Vocabulary:
[h,e,l,o]

Example training sequence:
**"hello"**

# Example

**Character-level language model example**

Vocabulary: [h,e,l,o]

Example training sequence: **"hello"**

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

# Example

**Character-level language model example**

Vocabulary:
[h,e,l,o]

Example training sequence:
**"hello"**



What do we still need to specify, for this to work?

What kind of loss can we formulate?

# Training a Recurrent Neural Network

- Get a big corpus of text which is a sequence of words $x^{(1)}, \ldots, x^{(T)}$
- Feed into RNN; compute output distribution $\hat{y}^{(t)}$ for *every step t*.
  - i.e. predict probability dist of *every word*, given words so far

- Loss function on step *t* is cross-entropy between predicted probability distribution $\hat{y}^{(t)}$, and true next word $y^{(t)}$ (one-hot); V is vocabulary

$$J^{(t)}(\theta) = CE(y^{(t)}, \hat{y}^{(t)}) = -\sum_{w \in V} y_w^{(t)} \log \hat{y}_w^{(t)} = -\log \hat{y}_{x_{t+1}}^{(t)}$$

- Average this to get overall loss for entire training set:

$$J(\theta) = \frac{1}{T} \sum_{t=1}^{T} J^{(t)}(\theta) = \frac{1}{T} \sum_{t=1}^{T} -\log \hat{y}_{x_{t+1}}^{(t)}$$

# The vanishing/exploding gradient problem

- The error at a time step ideally can tell a previous time step from many steps away to change during backprop
- Multiply the same matrix at each time step during backprop

# The vanishing gradient problem

- Total error is the sum of each error at time steps t

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

- Chain rule:

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

- More chain rule:

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

- Derivative of vector wrt vector is a Jacobian matrix of partial derivatives; norm of this matrix can become very small or very large quickly [Bengio et al 1994, Pascanu et al. 2013], leading to vanishing/exploding gradient

# What uses of language models from everyday life can you think of?

# Now in more detail…

# Language Modeling

- **Language Modeling** is the task of predicting what word comes next.

*the students opened their* _____

- *books*
- *laptops*
- *exams*
- *minds*

- More formally: given a sequence of words $\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \ldots, \boldsymbol{x}^{(t)}$, compute the probability distribution of the next word $\boldsymbol{x}^{(t+1)}$ :

$$P(\boldsymbol{x}^{(t+1)} \mid \boldsymbol{x}^{(t)}, \ldots, \boldsymbol{x}^{(1)})$$

where $\boldsymbol{x}^{(t+1)}$ can be any word in the vocabulary $V = \{\boldsymbol{w}_1, \ldots, \boldsymbol{w}_{|V|}\}$

- A system that does this is called a **Language Model**.

# n-gram Language Models

- First we make a simplifying assumption: $\boldsymbol{x}^{(t+1)}$ depends only on the preceding *n-1* words.

*n*-1 words

$$P(\boldsymbol{x}^{(t+1)}|\boldsymbol{x}^{(t)}, \ldots, \boldsymbol{x}^{(1)}) = P(\boldsymbol{x}^{(t+1)}|\boldsymbol{x}^{(t)}, \ldots, \boldsymbol{x}^{(t-n+2)}) \qquad \text{(assumption)}$$

prob of a n-gram

prob of a (n-1)-gram

$$= \frac{P(\boldsymbol{x}^{(t+1)}, \boldsymbol{x}^{(t)}, \ldots, \boldsymbol{x}^{(t-n+2)})}{P(\boldsymbol{x}^{(t)}, \ldots, \boldsymbol{x}^{(t-n+2)})} \qquad \text{(definition of conditional prob)}$$

- **Question:** How do we get these *n*-gram and (*n*-1)-gram probabilities?
- **Answer:** By counting them in some large corpus of text!

$$\approx \frac{\text{count}(\boldsymbol{x}^{(t+1)}, \boldsymbol{x}^{(t)}, \ldots, \boldsymbol{x}^{(t-n+2)})}{\text{count}(\boldsymbol{x}^{(t)}, \ldots, \boldsymbol{x}^{(t-n+2)})} \qquad \text{(statistical approximation)}$$

# Sparsity Problems with n-gram Language Models

**Sparsity Problem 1**

**Problem:** What if *"students opened their $w$"* never occurred in data? Then $w$ has probability 0!

**(Partial) Solution:** Add small $\delta$ to the count for every $w \in V$. This is called *smoothing*.

$$P(w|\text{students opened their}) = \frac{\text{count}(\text{students opened their } w)}{\text{count}(\text{students opened their})}$$

**Sparsity Problem 2**

**Problem:** What if *"students opened their"* never occurred in data? Then we can't calculate probability for *any* $w$!

**(Partial) Solution:** Just condition on *"opened their"* instead. This is called *backoff*.

**Note:** Increasing *n* makes sparsity problems *worse.* Typically we can't have *n* bigger than 5.

# A fixed-window neural Language Model

output distribution

$$\hat{\boldsymbol{y}} = \mathrm{softmax}(\boldsymbol{U}\boldsymbol{h} + \boldsymbol{b}_2) \in \mathbb{R}^{|V|}$$

books

laptops

a          zoo

$\boldsymbol{U}$

hidden layer

$$\boldsymbol{h} = f(\boldsymbol{W}\boldsymbol{e} + \boldsymbol{b}_1)$$

$\boldsymbol{W}$

concatenated word embeddings

$$\boldsymbol{e} = [\boldsymbol{e}^{(1)}; \boldsymbol{e}^{(2)}; \boldsymbol{e}^{(3)}; \boldsymbol{e}^{(4)}]$$

words / one-hot vectors

$$\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \boldsymbol{x}^{(3)}, \boldsymbol{x}^{(4)}$$

*the*
$\boldsymbol{x}^{(1)}$

*students*
$\boldsymbol{x}^{(2)}$

*opened*
$\boldsymbol{x}^{(3)}$

*their*
$\boldsymbol{x}^{(4)}$

Abigail See

# A fixed-window neural Language Model

**Improvements** over *n*-gram LM:
- No sparsity problem
- Don't need to store all observed *n*-grams

Remaining **problems**:
- Fixed window is too small
- Enlarging window enlarges $W$
- Window can never be large enough!
- $x^{(1)}$ and $x^{(2)}$ are multiplied by completely different weights in $W$. No symmetry in how the inputs are processed.

We need a neural architecture that can process *any length input*

books

laptops

a                    zoo

$U$

$W$

the $x^{(1)}$   students $x^{(2)}$   opened $x^{(3)}$   their $x^{(4)}$

# Recurrent Neural Networks (RNN)

A family of neural architectures

outputs (optional) $\{$    $\hat{y}^{(1)}$    $\hat{y}^{(2)}$    $\hat{y}^{(3)}$    $\hat{y}^{(4)}$    ...

$h^{(1)}$    $h^{(2)}$    $h^{(3)}$    $h^{(4)}$

hidden states $\{$    $W$    $W$    $W$    $W$    ...

input sequence (any length) $\{$    $x^{(1)}$    $x^{(2)}$    $x^{(3)}$    $x^{(4)}$    ...

Abigail See

# A RNN Language Model

$$\hat{y}^{(4)} = P(x^{(5)}|\text{the students opened their})$$

output distribution

$$\hat{y}^{(t)} = \text{softmax}\left(Uh^{(t)} + b_2\right) \in \mathbb{R}^{|V|}$$

hidden states

$$h^{(t)} = \sigma\left(W_h h^{(t-1)} + W_e e^{(t)} + b_1\right)$$

$h^{(0)}$ is the initial hidden state

word embeddings

$$e^{(t)} = Ex^{(t)}$$

words / one-hot vectors
$$x^{(t)} \in \mathbb{R}^{|V|}$$

*books*

*laptops*

a    zoo

$U$

$h^{(0)}$   $h^{(1)}$   $h^{(2)}$   $h^{(3)}$   $h^{(4)}$

$W_h$   $W_h$   $W_h$   $W_h$

$W_e$   $W_e$   $W_e$   $W_e$

$e^{(1)}$   $e^{(2)}$   $e^{(3)}$   $e^{(4)}$

$E$   $E$   $E$   $E$

*the*   *students*   *opened*   *their*
$$x^{(1)} \quad x^{(2)} \quad x^{(3)} \quad x^{(4)}$$

**Note**: this input sequence could be much longer, but this slide doesn't have space!

Abigail See

# A RNN Language Model

RNN **Advantages**:
- Can process any length input
- Computation for step *t* can (in theory) use information from many steps back
- Model size doesn't increase for longer input
- Same weights applied on every timestep, so there is symmetry in how inputs are processed.

RNN **Disadvantages**:
- Recurrent computation is slow
- In practice, difficult to access information from many steps back

$$\hat{y}^{(4)} = P(x^{(5)}|\text{the students opened their})$$



Abigail See

# Recall: Training a RNN Language Model

- Get a big corpus of text which is a sequence of words $x^{(1)}, \ldots, x^{(T)}$
- Feed into RNN-LM; compute output distribution $\hat{y}^{(t)}$ for *every step t.*
  - i.e. predict probability dist of *every word*, given words so far

- Loss function on step *t* is cross-entropy between predicted probability distribution $\hat{y}^{(t)}$, and the true next word $y^{(t)}$ (one-hot for $x^{(t+1)}$):

$$J^{(t)}(\theta) = CE(y^{(t)}, \hat{y}^{(t)}) = -\sum_{w \in V} y_w^{(t)} \log \hat{y}_w^{(t)} = -\log \hat{y}_{x_{t+1}}^{(t)}$$

- Average this to get overall loss for entire training set:

$$J(\theta) = \frac{1}{T} \sum_{t=1}^{T} J^{(t)}(\theta) = \frac{1}{T} \sum_{t=1}^{T} -\log \hat{y}_{x_{t+1}}^{(t)}$$

# Training a RNN Language Model

= negative log prob
of "students"

Loss ⟶ $J^{(1)}(\theta)$ $\quad J^{(2)}(\theta) \quad\quad J^{(3)}(\theta) \quad\quad J^{(4)}(\theta)$

Predicted
prob dists ⟶ $\hat{\boldsymbol{y}}^{(1)}$ $\quad\quad \hat{\boldsymbol{y}}^{(2)} \quad\quad \hat{\boldsymbol{y}}^{(3)} \quad\quad \hat{\boldsymbol{y}}^{(4)}$

$\boldsymbol{U} \quad\quad \boldsymbol{U} \quad\quad \boldsymbol{U} \quad\quad \boldsymbol{U}$

$\boldsymbol{h}^{(0)} \quad\quad \boldsymbol{h}^{(1)} \quad\quad \boldsymbol{h}^{(2)} \quad\quad \boldsymbol{h}^{(3)} \quad\quad \boldsymbol{h}^{(4)}$

$\boldsymbol{W}_h \quad\quad \boldsymbol{W}_h \quad\quad \boldsymbol{W}_h \quad\quad \boldsymbol{W}_h \quad\quad \boldsymbol{W}_h \quad$ ...

$\boldsymbol{W}_e \quad\quad \boldsymbol{W}_e \quad\quad \boldsymbol{W}_e \quad\quad \boldsymbol{W}_e$

$\boldsymbol{e}^{(1)} \quad\quad \boldsymbol{e}^{(2)} \quad\quad \boldsymbol{e}^{(3)} \quad\quad \boldsymbol{e}^{(4)}$

$\boldsymbol{E} \quad\quad \boldsymbol{E} \quad\quad \boldsymbol{E} \quad\quad \boldsymbol{E}$

Corpus ⟶ *the* $\quad$ *students* $\quad$ *opened* $\quad$ *their* $\quad$ *exams* $\quad$ ...
$\boldsymbol{x}^{(1)} \quad\quad \boldsymbol{x}^{(2)} \quad\quad \boldsymbol{x}^{(3)} \quad\quad \boldsymbol{x}^{(4)}$

Abigail See

# Training a RNN Language Model

= negative log prob
of "opened"



Loss $\longrightarrow$ $J^{(1)}(\theta)$ $\boxed{J^{(2)}(\theta)}$ $J^{(3)}(\theta)$ $J^{(4)}(\theta)$

Predicted prob dists $\longrightarrow$ $\hat{\boldsymbol{y}}^{(1)}$ $\hat{\boldsymbol{y}}^{(2)}$ $\hat{\boldsymbol{y}}^{(3)}$ $\hat{\boldsymbol{y}}^{(4)}$

$\boldsymbol{U}$ $\boldsymbol{U}$ $\boldsymbol{U}$ $\boldsymbol{U}$

$\boldsymbol{h}^{(0)}$ $\boldsymbol{h}^{(1)}$ $\boldsymbol{h}^{(2)}$ $\boldsymbol{h}^{(3)}$ $\boldsymbol{h}^{(4)}$

$\boldsymbol{W}_h$ $\boldsymbol{W}_h$ $\boldsymbol{W}_h$ $\boldsymbol{W}_h$ $\boldsymbol{W}_h$ ...

$\boldsymbol{W}_e$ $\boldsymbol{W}_e$ $\boldsymbol{W}_e$ $\boldsymbol{W}_e$

$\boldsymbol{e}^{(1)}$ $\boldsymbol{e}^{(2)}$ $\boldsymbol{e}^{(3)}$ $\boldsymbol{e}^{(4)}$

$\boldsymbol{E}$ $\boldsymbol{E}$ $\boldsymbol{E}$ $\boldsymbol{E}$

Corpus $\longrightarrow$ *the* *students* *opened* *their* *exams* ...
$\boldsymbol{x}^{(1)}$ $\boldsymbol{x}^{(2)}$ $\boldsymbol{x}^{(3)}$ $\boldsymbol{x}^{(4)}$

Abigail See

# Training a RNN Language Model

= negative log prob of "their"

Loss → $J^{(1)}(\theta)$  $J^{(2)}(\theta)$  $\boxed{J^{(3)}(\theta)}$  $J^{(4)}(\theta)$

Predicted prob dists → $\hat{\boldsymbol{y}}^{(1)}$  $\hat{\boldsymbol{y}}^{(2)}$  $\hat{\boldsymbol{y}}^{(3)}$  $\hat{\boldsymbol{y}}^{(4)}$

$\boldsymbol{U}$  $\boldsymbol{U}$  $\boldsymbol{U}$  $\boldsymbol{U}$

$\boldsymbol{h}^{(0)}$  $\boldsymbol{h}^{(1)}$  $\boldsymbol{h}^{(2)}$  $\boldsymbol{h}^{(3)}$  $\boldsymbol{h}^{(4)}$

$\boldsymbol{W}_h$  $\boldsymbol{W}_h$  $\boldsymbol{W}_h$  $\boldsymbol{W}_h$  $\boldsymbol{W}_h$  ...

$\boldsymbol{W}_e$  $\boldsymbol{W}_e$  $\boldsymbol{W}_e$  $\boldsymbol{W}_e$

$\boldsymbol{e}^{(1)}$  $\boldsymbol{e}^{(2)}$  $\boldsymbol{e}^{(3)}$  $\boldsymbol{e}^{(4)}$

$\boldsymbol{E}$  $\boldsymbol{E}$  $\boldsymbol{E}$  $\boldsymbol{E}$

Corpus →  *the* $\boldsymbol{x}^{(1)}$  *students* $\boldsymbol{x}^{(2)}$  *opened* $\boldsymbol{x}^{(3)}$  *their* $\boldsymbol{x}^{(4)}$  *exams*  ...

Abigail See

# Training a RNN Language Model



= negative log prob
of "exams"

Loss ⟶ $J^{(1)}(\theta)$    $J^{(2)}(\theta)$    $J^{(3)}(\theta)$    $\boxed{J^{(4)}(\theta)}$

Predicted prob dists ⟶ $\hat{\boldsymbol{y}}^{(1)}$   $\hat{\boldsymbol{y}}^{(2)}$   $\hat{\boldsymbol{y}}^{(3)}$   $\hat{\boldsymbol{y}}^{(4)}$

$\boldsymbol{U}$   $\boldsymbol{U}$   $\boldsymbol{U}$   $\boldsymbol{U}$

$\boldsymbol{h}^{(0)}$   $\boldsymbol{h}^{(1)}$   $\boldsymbol{h}^{(2)}$   $\boldsymbol{h}^{(3)}$   $\boldsymbol{h}^{(4)}$

$\boldsymbol{W}_h$   $\boldsymbol{W}_h$   $\boldsymbol{W}_h$   $\boldsymbol{W}_h$   $\boldsymbol{W}_h$   ...

$\boldsymbol{W}_e$   $\boldsymbol{W}_e$   $\boldsymbol{W}_e$   $\boldsymbol{W}_e$

$\boldsymbol{e}^{(1)}$   $\boldsymbol{e}^{(2)}$   $\boldsymbol{e}^{(3)}$   $\boldsymbol{e}^{(4)}$

$\boldsymbol{E}$   $\boldsymbol{E}$   $\boldsymbol{E}$   $\boldsymbol{E}$

Corpus ⟶ *the*   *students*   *opened*   *their*   *exams*   ...

$\boldsymbol{x}^{(1)}$   $\boldsymbol{x}^{(2)}$   $\boldsymbol{x}^{(3)}$   $\boldsymbol{x}^{(4)}$

Abigail See

# Training a RNN Language Model

Loss ⟶ $J^{(1)}(\theta)$ + $J^{(2)}(\theta)$ + $J^{(3)}(\theta)$ + $J^{(4)}(\theta)$ + ... = $J(\theta) = \dfrac{1}{T}\displaystyle\sum_{t=1}^{T} J^{(t)}(\theta)$

Predicted prob dists ⟶ $\hat{\boldsymbol{y}}^{(1)}$ $\hat{\boldsymbol{y}}^{(2)}$ $\hat{\boldsymbol{y}}^{(3)}$ $\hat{\boldsymbol{y}}^{(4)}$

$\boldsymbol{U}$ $\boldsymbol{U}$ $\boldsymbol{U}$ $\boldsymbol{U}$

$\boldsymbol{h}^{(0)}$ $\boldsymbol{h}^{(1)}$ $\boldsymbol{h}^{(2)}$ $\boldsymbol{h}^{(3)}$ $\boldsymbol{h}^{(4)}$

$\boldsymbol{W}_h$ $\boldsymbol{W}_h$ $\boldsymbol{W}_h$ $\boldsymbol{W}_h$ $\boldsymbol{W}_h$ ...

$\boldsymbol{W}_e$ $\boldsymbol{W}_e$ $\boldsymbol{W}_e$ $\boldsymbol{W}_e$

$\boldsymbol{e}^{(1)}$ $\boldsymbol{e}^{(2)}$ $\boldsymbol{e}^{(3)}$ $\boldsymbol{e}^{(4)}$

$\boldsymbol{E}$ $\boldsymbol{E}$ $\boldsymbol{E}$ $\boldsymbol{E}$

Corpus ⟶ *the*    *students*    *opened*    *their*    *exams*    ...

$\boldsymbol{x}^{(1)}$    $\boldsymbol{x}^{(2)}$    $\boldsymbol{x}^{(3)}$    $\boldsymbol{x}^{(4)}$

Abigail See

# Training a RNN Language Model

- However: Computing loss and gradients across entire corpus $\boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(T)}$ is too expensive!

$$J(\theta) = \frac{1}{T} \sum_{t=1}^{T} J^{(t)}(\theta)$$

- In practice, consider $\boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(T)}$ as a sentence (or a document)

- Recall: Stochastic Gradient Descent allows us to compute loss and gradients for small chunk of data, and update.

- Compute loss $J(\theta)$ for a sentence (actually a batch of sentences), compute gradients and update weights. Repeat.

Abigail See

# Backpropagation for RNNs



**Question:** What's the derivative of $J^{(t)}(\theta)$ w.r.t. the repeated weight matrix $\boldsymbol{W}_h$ ?

**Answer:** $\dfrac{\partial J^{(t)}}{\partial \boldsymbol{W}_h} = \displaystyle\sum_{i=1}^{t} \dfrac{\partial J^{(t)}}{\partial \boldsymbol{W}_h}\bigg|_{(i)}$

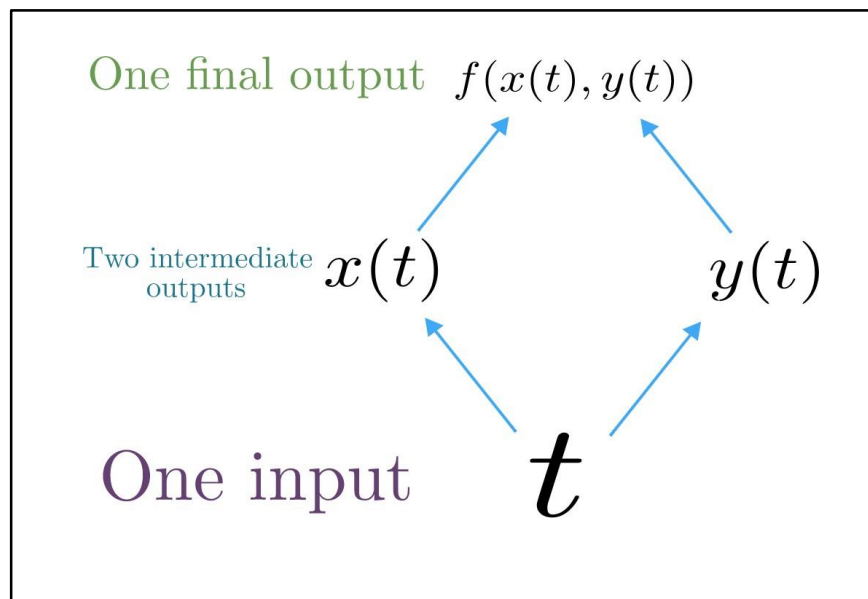"The gradient w.r.t. a repeated weight is the sum of the gradient w.r.t. each time it appears"

**Why?**

Abigail See

# Multivariable Chain Rule

- Given a multivariable function $f(x, y)$, and two single variable functions $x(t)$ and $y(t)$, here's what the multivariable chain rule says:

$$\underbrace{\frac{d}{dt} f(x(t), y(t))}_{\text{Derivative of composition function}} = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt}$$

One final output $f(x(t), y(t))$

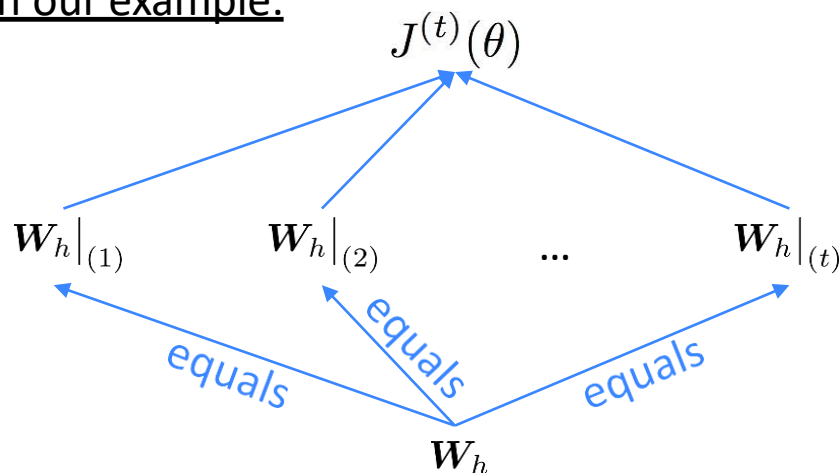Two intermediate outputs $x(t)$ $\qquad$ $y(t)$

One input $t$

Abigail See

# Backpropagation for RNNs: Proof sketch

- Given a multivariable function $f(x, y)$, and two single variable functions $x(t)$ and $y(t)$, here's what the multivariable chain rule says:

$$\underbrace{\frac{d}{dt} f(x(t), y(t))}_{\text{Derivative of composition function}} = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt}$$

In our example:

$$J^{(t)}(\theta)$$

$$\left. W_h \right|_{(1)} \quad \left. W_h \right|_{(2)} \quad \dots \quad \left. W_h \right|_{(t)}$$

equals    equals    equals

$$W_h$$

Apply the multivariable chain rule:

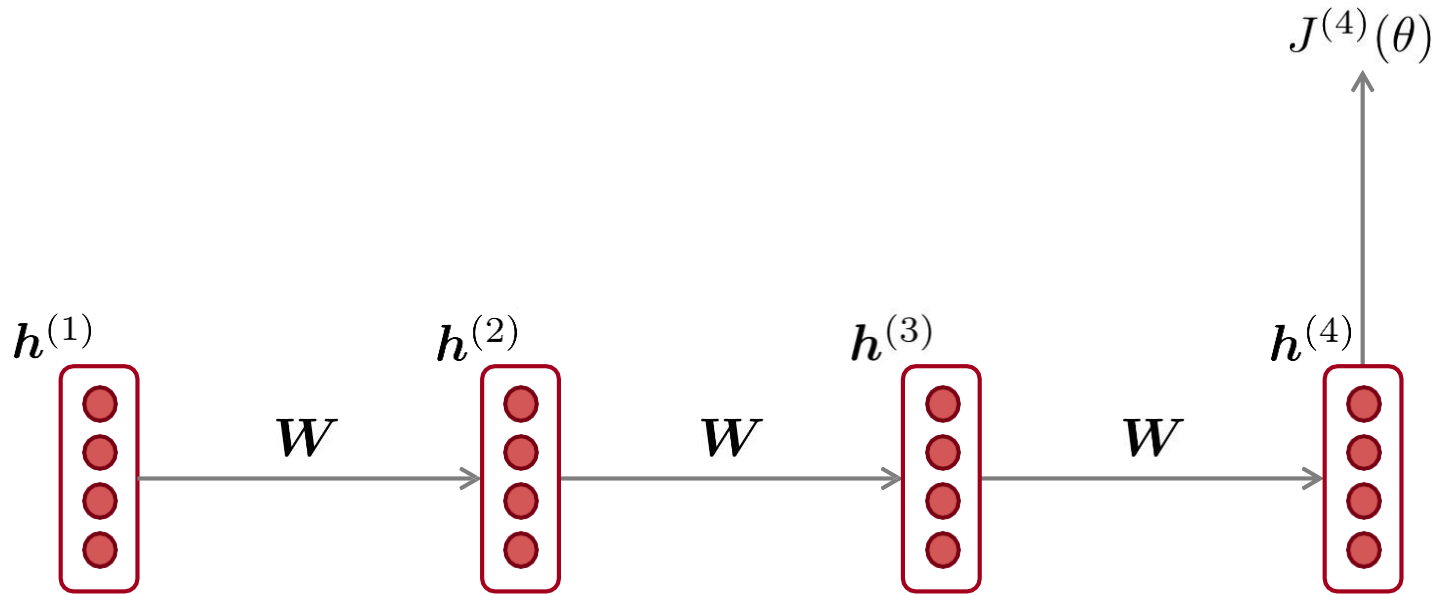$$= 1$$

$$\frac{\partial J^{(t)}}{\partial W_h} = \sum_{i=1}^{t} \left. \frac{\partial J^{(t)}}{\partial W_h} \right|_{(i)} \boxed{\frac{\left. \partial W_h \right|_{(i)}}{\partial W_h}}$$

$$= \sum_{i=1}^{t} \left. \frac{\partial J^{(t)}}{\partial W_h} \right|_{(i)}$$

Abigail See

# Backpropagation for RNNs



$$\frac{\partial J^{(t)}}{\partial \boldsymbol{W_h}} = \boxed{\sum_{i=1}^{t} \frac{\partial J^{(t)}}{\partial \boldsymbol{W_h}}\bigg|_{(i)}}$$
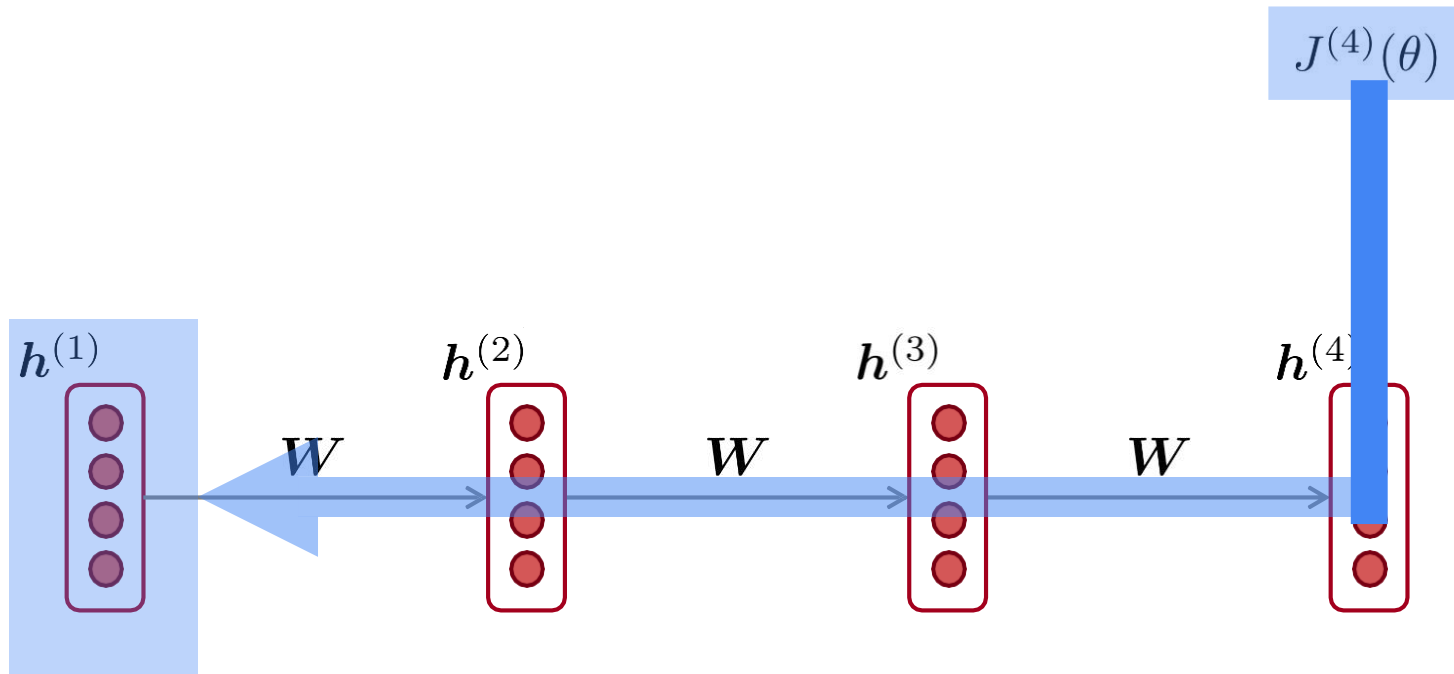
**Question:** How do we calculate this?

**Answer:** Backpropagate over timesteps $i=t,\ldots,0$, summing gradients as you go.
This algorithm is called **"backpropagation through time"**

Abigail See

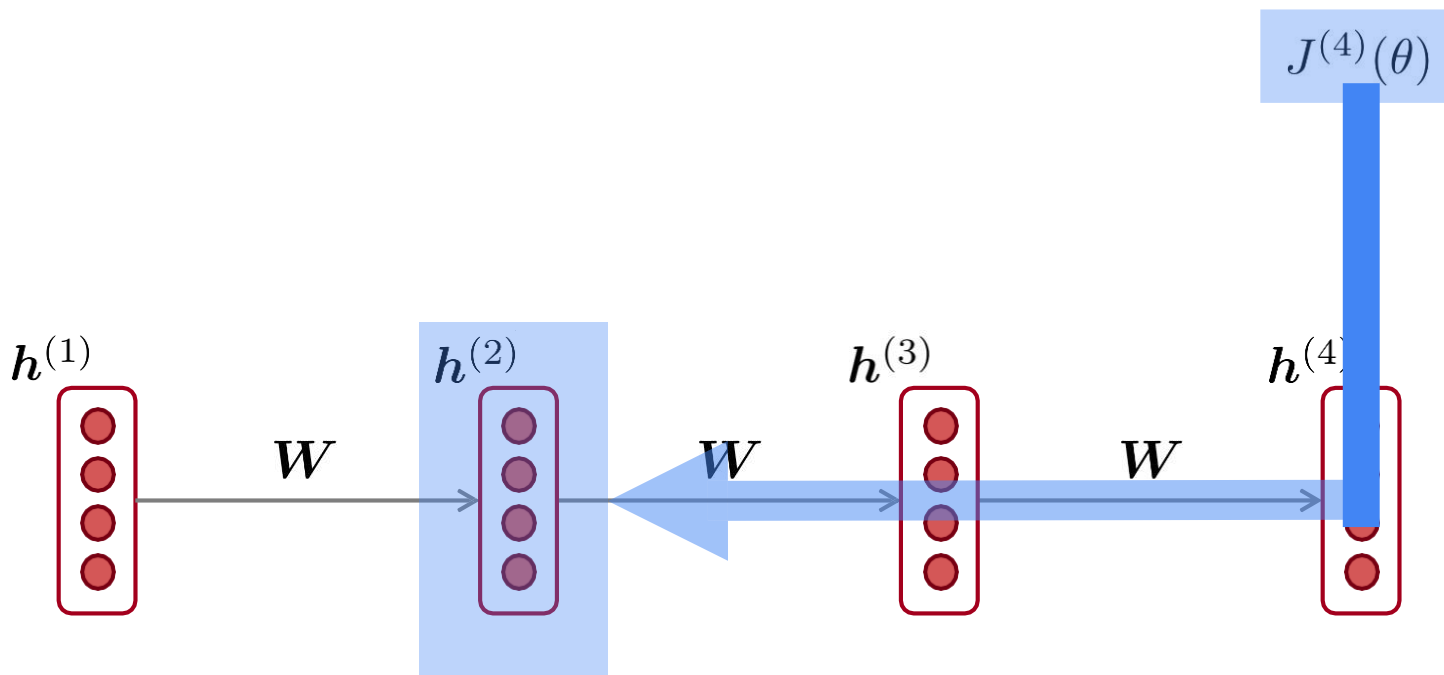# Vanishing gradient intuition

$$J^{(4)}(\theta)$$

$$\boldsymbol{h}^{(1)} \qquad \boldsymbol{W} \qquad \boldsymbol{h}^{(2)} \qquad \boldsymbol{W} \qquad \boldsymbol{h}^{(3)} \qquad \boldsymbol{W} \qquad \boldsymbol{h}^{(4)}$$

# Vanishing gradient intuition

$$J^{(4)}(\theta)$$

$$\boldsymbol{h}^{(1)} \qquad \boldsymbol{h}^{(2)} \qquad \boldsymbol{h}^{(3)} \qquad \boldsymbol{h}^{(4)}$$

$$\boldsymbol{W} \qquad \boldsymbol{W} \qquad \boldsymbol{W}$$

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

Abigail See

# Vanishing gradient intuition

$$J^{(4)}(\theta)$$

$$\boldsymbol{h}^{(1)} \qquad \boldsymbol{h}^{(2)} \qquad \boldsymbol{h}^{(3)} \qquad \boldsymbol{h}^{(4)}$$

$$\boldsymbol{W} \qquad \boldsymbol{W} \qquad \boldsymbol{W}$$

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

chain rule!

Abigail See

# Vanishing gradient intuition



$$\frac{\partial J^{(4)}}{\partial \boldsymbol{h}^{(1)}} = \frac{\partial \boldsymbol{h}^{(2)}}{\partial \boldsymbol{h}^{(1)}} \times \qquad \frac{\partial \boldsymbol{h}^{(3)}}{\partial \boldsymbol{h}^{(2)}} \times \frac{\partial J^{(4)}}{\partial \boldsymbol{h}^{(3)}}$$

chain rule!

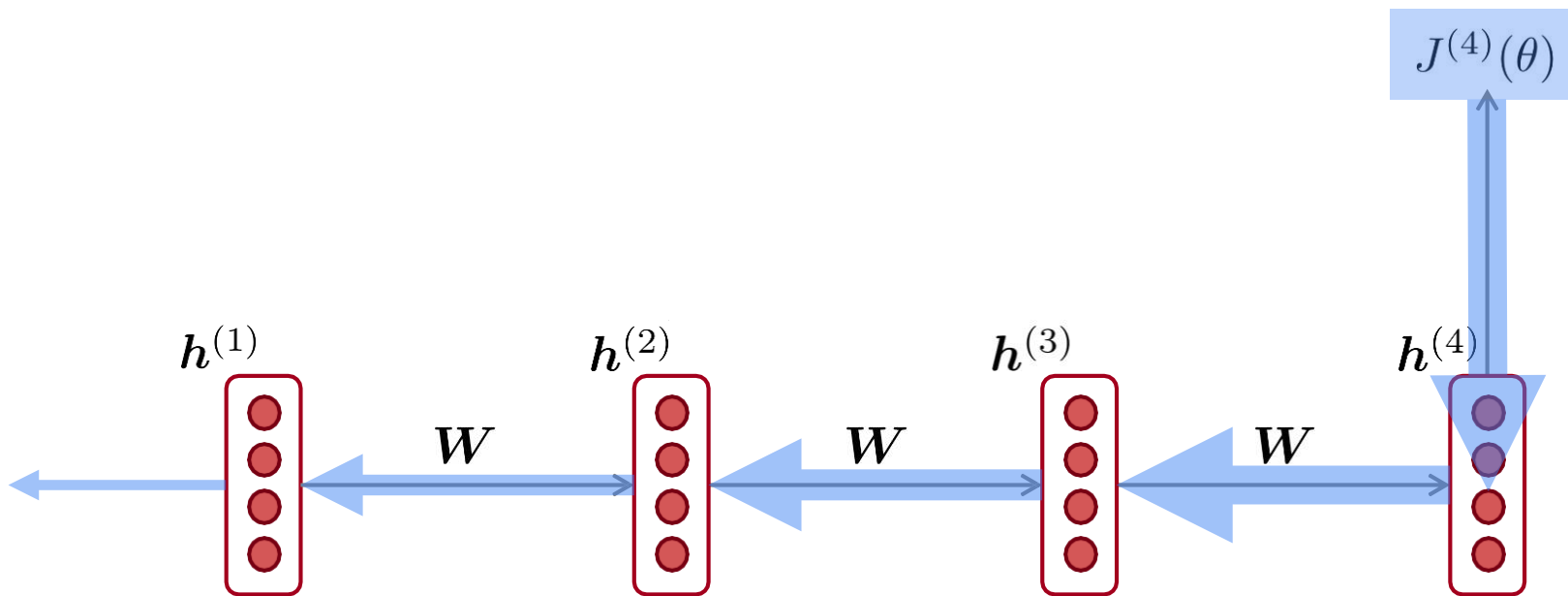# Vanishing gradient intuition



$$\frac{\partial J^{(4)}}{\partial h^{(1)}} = \quad \frac{\partial h^{(2)}}{\partial h^{(1)}} \times \qquad \frac{\partial h^{(3)}}{\partial h^{(2)}} \times \qquad \frac{\partial h^{(4)}}{\partial h^{(3)}} \times \quad \frac{\partial J^{(4)}}{\partial h^{(4)}}$$

chain rule!

# Vanishing gradient intuition



$$J^{(4)}(\theta)$$

$$\boldsymbol{h}^{(1)} \qquad \boldsymbol{W} \qquad \boldsymbol{h}^{(2)} \qquad \boldsymbol{W} \qquad \boldsymbol{h}^{(3)} \qquad \boldsymbol{W} \qquad \boldsymbol{h}^{(4)}$$

$$\frac{\partial J^{(4)}}{\partial \boldsymbol{h}^{(1)}} = \frac{\partial \boldsymbol{h}^{(2)}}{\partial \boldsymbol{h}^{(1)}} \times \frac{\partial \boldsymbol{h}^{(3)}}{\partial \boldsymbol{h}^{(2)}} \times \frac{\partial \boldsymbol{h}^{(4)}}{\partial \boldsymbol{h}^{(3)}} \times \frac{\partial J^{(4)}}{\partial \boldsymbol{h}^{(4)}}$$

What happens if these are small?

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

Abigail See

# Vanishing gradient proof sketch

- Recall: $\quad h^{(t)} = \sigma\left(W_h h^{(t-1)} + W_x x^{(t)} + b_1\right)$

- Therefore: $\quad \dfrac{\partial h^{(t)}}{\partial h^{(t-1)}} = \text{diag}\left(\sigma'\left(W_h h^{(t-1)} + W_x x^{(t)} + b_1\right)\right) W_h \quad$ (chain rule)

- Consider the gradient of the loss $J^{(i)}(\theta)$ on step *i*, with respect to the hidden state $h^{(j)}$ on some previous step *j*.

$$\frac{\partial J^{(i)}(\theta)}{\partial h^{(j)}} = \frac{\partial J^{(i)}(\theta)}{\partial h^{(i)}} \prod_{j < t \leq i} \frac{\partial h^{(t)}}{\partial h^{(t-1)}} \qquad \text{(chain rule)}$$

$$= \frac{\partial J^{(i)}(\theta)}{\partial h^{(i)}} \boxed{W_h^{(i-j)}} \prod_{j < t \leq i} \text{diag}\left(\sigma'\left(W_h h^{(t-1)} + W_x x^{(t)} + b_1\right)\right) \qquad \left(\text{value of } \frac{\partial h^{(t)}}{\partial h^{(t-1)}}\right)$$

If $W_h$ is small, then this term gets vanishingly small as *i* and *j* get further apart

**Source**: "On the difficulty of training recurrent neural networks", Pascanu et al, 2013. http://proceedings.mlr.press/v28/pascanu13.pdf
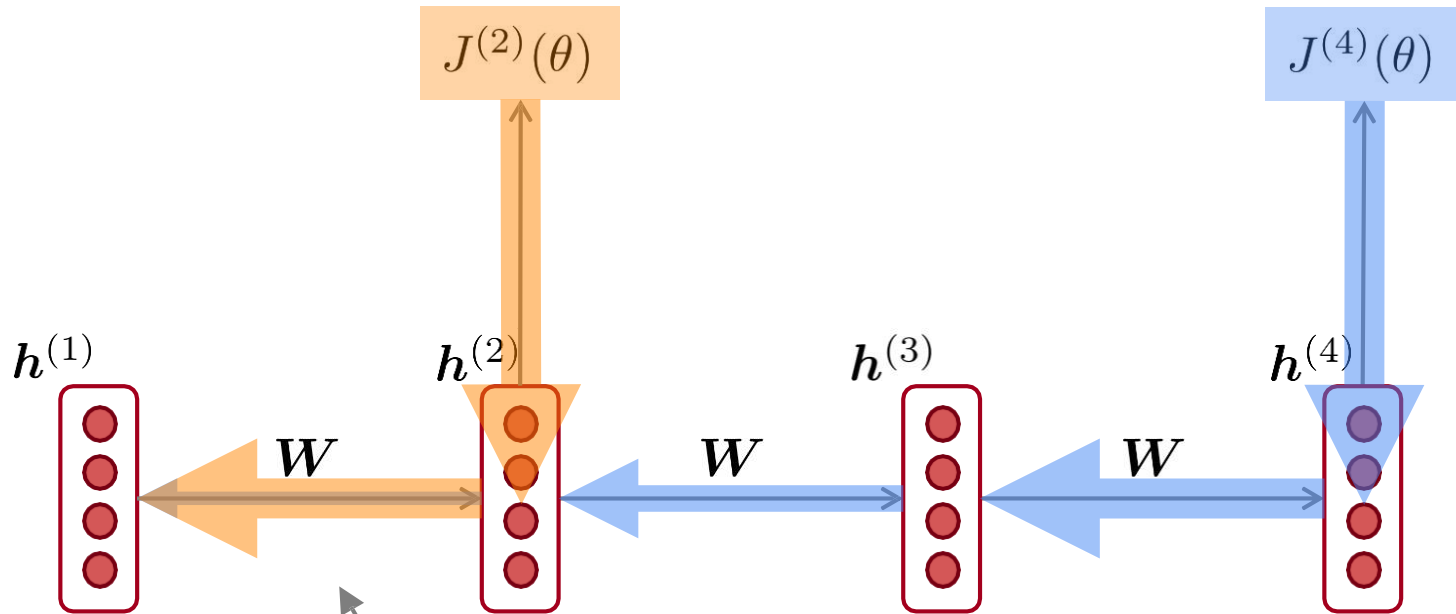
Abigail See

# Vanishing gradient proof sketch

- Consider matrix L2 norms:

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

- Pascanu et al showed that that if the largest eigenvalue of $W_h$ is less than 1, then the gradient $\left\| \frac{\partial J^{(i)}(\theta)}{\partial \boldsymbol{h}^{(j)}} \right\|$ will shrink exponentially
  - Here the bound is 1 because we have sigmoid nonlinearity

- There's a similar proof relating a largest eigenvalue >1 to exploding gradients

Abigail See

# Why is vanishing gradient a problem?

$J^{(2)}(\theta)$

$J^{(4)}(\theta)$

$\boldsymbol{h}^{(1)}$       $\boldsymbol{h}^{(2)}$       $\boldsymbol{h}^{(3)}$       $\boldsymbol{h}^{(4)}$

$\boldsymbol{W}$       $\boldsymbol{W}$       $\boldsymbol{W}$
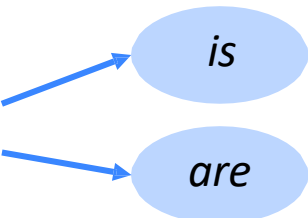
Gradient signal from faraway is lost because it's much smaller than gradient signal from close-by.

So model weights are only updated only with respect to near effects, not long-term effects.

Abigail See

# Effect of vanishing gradient on RNN-LM

- **LM task:** *When she tried to print her tickets, she found that the printer was out of toner. She went to the stationery store to buy more toner. It was very overpriced. After installing the toner into the printer, she finally printed her* _____

- To learn from this training example, the RNN-LM needs to model the dependency between *"tickets"* on the 7[th] step and the target word *"tickets"* at the end.

- But if gradient is small, the model can't learn this dependency
  - So the model is unable to predict similar long-distance dependencies at test time

Abigail See

# Effect of vanishing gradient on RNN-LM

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

  *is*

  *are*

- **Correct answer**: *The writer of the books is planning a sequel*

- **Syntactic** **recency**: *The writer of the books is*　　　　　　(correct)

- **Sequential** **recency:** *The writer of the books are*　　　　(incorrect)

- Due to vanishing gradient, RNN-LMs are better at learning from sequential recency than syntactic recency, so they make this type of error more often than we'd like [Linzen et al 2016]

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

Abigail See

# Why is <u>exploding</u> gradient a problem?

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

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

- This can cause bad updates: we take too large a step and reach a bad parameter configuration (with large loss)

- In the worst case, this will result in Inf or NaN in your network (then you have to restart training from an earlier checkpoint)

Abigail See

# Gradient clipping: solution for exploding gradient

- Gradient clipping: if the norm of the gradient is greater than some threshold, scale it down before applying SGD update

**Algorithm 1** Pseudo-code for norm clipping

$\hat{\mathbf{g}} \leftarrow \frac{\partial \mathcal{E}}{\partial \theta}$

**if** $\|\hat{\mathbf{g}}\| \geq threshold$ **then**

$\quad \hat{\mathbf{g}} \leftarrow \frac{threshold}{\|\hat{\mathbf{g}}\|} \hat{\mathbf{g}}$

**end if**

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

**Source**: "On the difficulty of training recurrent neural networks", Pascanu et al, 2013.  http://proceedings.mlr.press/v28/pascanu13.pdf

Abigail See

# RNNs with Gates

# How to fix vanishing gradient problem?

- The main problem is that *it's too difficult for the RNN to learn to preserve information over many timesteps.*

- In a vanilla RNN, the hidden state is constantly being rewritten

$$h^{(t)} = \sigma \left( W_h h^{(t-1)} + W_x x^{(t)} \right)$$

- How about a RNN with separate memory?

Richard Socher

# Gated Recurrent Units (GRUs)

- More complex hidden unit computation in recurrence!

- Introduced by Cho et al. 2014

- Main ideas:

  - keep around memories to capture long distance dependencies

  - allow error messages to flow at different strengths depending on the inputs

# Gated Recurrent Units (GRUs)

- Standard RNN computes hidden layer at next time step directly: $h_t = f\left(W^{(hh)}h_{t-1} + W^{(hx)}x_t\right)$

- 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

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

Richard Socher

# Gated Recurrent Units (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:  $\tilde{h}_t = \tanh\left(Wx_t + r_t \circ Uh_{t-1}\right)$
  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:  $h_t = z_t \circ h_{t-1} + (1 - z_t) \circ \tilde{h}_t$

# Gated Recurrent Units (GRUs)



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

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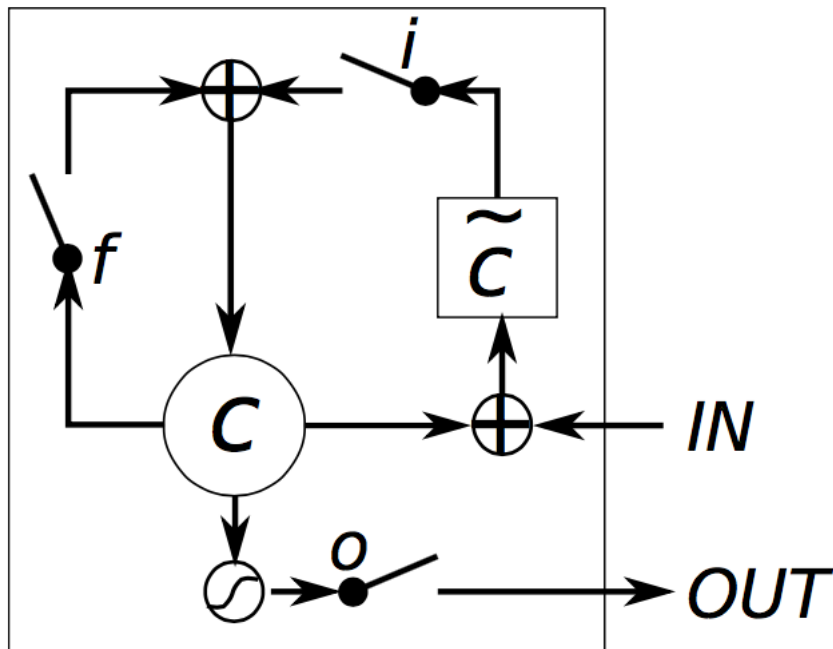
$$\tilde{h}_t = \tanh\left(Wx_t + r_t \circ Uh_{t-1}\right)$$

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

Final memory $h_{t-1}$ $h_t$

Memory (reset) $\tilde{h}_{t-1}$ $\tilde{h}_t$

Update gate $z_{t-1}$ $z_t$

Reset gate $r_{t-1}$ $r_t$

Input: $x_{t-1}$ $x_t$

Richard Socher

# Gated Recurrent Units (GRUs)

- If reset **r** is close to 0, ignore previous hidden state:  Allows model to drop information that is irrelevant in the future

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

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

$$\tilde{h}_t = \tanh\left(Wx_t + r_t \circ Uh_{t-1}\right)$$

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

- If update **z** is close to 1, can copy information through many time steps, i.e. copy paste state: Less vanishing gradient!

- Units with short-term dependencies often have reset gates (r) very active; ones with long-term dependencies have active update gates (z)

# Long-short-term-memories (LSTMs)

- Proposed by Hochreiter and Schmidhuber in 1997

- We can make the units even more complex

- Allow each time step to modify

  - Input gate (current cell matters) $\quad i_t = \sigma\left(W^{(i)} x_t + U^{(i)} h_{t-1}\right)$

  - Forget (gate 0, forget past) $\quad f_t = \sigma\left(W^{(f)} x_t + U^{(f)} h_{t-1}\right)$

  - Output (how much cell is exposed) $o_t = \sigma\left(W^{(o)} x_t + U^{(o)} h_{t-1}\right)$

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

- Final memory cell: $\quad c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$

- Final hidden state: $\quad h_t = o_t \circ \tanh(c_t)$

# Long-short-term-memories (LSTMs)



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

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

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

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

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

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

Intuition: memory cells can keep information intact, unless inputs makes them forget it or overwrite it with new input

Cell can decide to output this information or just store it

Richard Socher, figure from wildml.com

# Review on your own: Gated Recurrent Units (GRU)

- Proposed by Cho et al. in 2014 as a simpler alternative to the LSTM.
- On each timestep $t$ we have input $\boldsymbol{x}^{(t)}$ and hidden state $\boldsymbol{h}^{(t)}$ (no cell state).

**Update gate:** controls what parts of hidden state are updated vs preserved

**Reset gate:** controls what parts of previous hidden state are used to compute new content

$$\boldsymbol{u}^{(t)} = \sigma\left(\boldsymbol{W}_u \boldsymbol{h}^{(t-1)} + \boldsymbol{U}_u \boldsymbol{x}^{(t)} + \boldsymbol{b}_u\right)$$

$$\boldsymbol{r}^{(t)} = \sigma\left(\boldsymbol{W}_r \boldsymbol{h}^{(t-1)} + \boldsymbol{U}_r \boldsymbol{x}^{(t)} + \boldsymbol{b}_r\right)$$

**New hidden state content:** reset gate selects useful parts of prev hidden state. Use this and current input to compute new hidden content.

$$\tilde{\boldsymbol{h}}^{(t)} = \tanh\left(\boldsymbol{W}_h (\boldsymbol{r}^{(t)} \circ \boldsymbol{h}^{(t-1)}) + \boldsymbol{U}_h \boldsymbol{x}^{(t)} + \boldsymbol{b}_h\right)$$

$$\boldsymbol{h}^{(t)} = (1 - \boldsymbol{u}^{(t)}) \circ \boldsymbol{h}^{(t-1)} + \boldsymbol{u}^{(t)} \circ \tilde{\boldsymbol{h}}^{(t)}$$

**Hidden state:** update gate simultaneously controls what is kept from previous hidden state, and what is updated to new hidden state content

**How does this solve vanishing gradient?**
GRU makes it easier to retain info  long-term (e.g. by setting update gate to 0)

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

Abigail See

# Review on your own: Long Short-Term Memory (LSTM)

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

**Forget gate:** controls what is kept vs forgotten, from previous cell state

**Input gate:** controls what parts of the new cell content are written to cell

**Output gate:** controls what parts of cell are output to hidden state

**Sigmoid function**: all gate values are between 0 and 1

$$f^{(t)} = \sigma\left(W_f h^{(t-1)} + U_f x^{(t)} + b_f\right)$$

$$i^{(t)} = \sigma\left(W_i h^{(t-1)} + U_i x^{(t)} + b_i\right)$$

$$o^{(t)} = \sigma\left(W_o h^{(t-1)} + U_o x^{(t)} + b_o\right)$$

All these are vectors of same length $n$

**New cell content:** this is the new content to be written to the cell

**Cell state**: erase ("forget") some content from last cell state, and write ("input") some new cell content

**Hidden state**: read ("output") some content from the cell

$$\tilde{c}^{(t)} = \tanh\left(W_c h^{(t-1)} + U_c x^{(t)} + b_c\right)$$

$$c^{(t)} = f^{(t)} \circ c^{(t-1)} + i^{(t)} \circ \tilde{c}^{(t)}$$

$$h^{(t)} = o^{(t)} \circ \tanh c^{(t)}$$

Gates are applied using element-wise product

Abigail See

# Review on your own: Long Short-Term Memory (LSTM)

You can think of the LSTM equations visually like this:



Write some new cell content

Forget some cell content

Output some cell content to the hidden state

Compute the forget gate

Compute the input gate

Compute the new cell content

Compute the output gate

Neural Network Layer | Pointwise Operation | Vector Transfer | Concatenate | Copy

Abigail See

# LSTM vs GRU

- Researchers have proposed many gated RNN variants, but LSTM and GRU are the most widely-used

- The biggest difference is that GRU is quicker to compute and has fewer parameters

- There is no conclusive evidence that one consistently performs better than the other

- LSTM is a good default choice (especially if your data has particularly long dependencies, or you have lots of training data)

- Rule of thumb: start with LSTM, but switch to GRU if you want something more efficient

# LSTMs: real-world success

- In 2013-2015, LSTMs started achieving state-of-the-art results
  - Successful tasks include: handwriting recognition, speech recognition, machine translation, parsing, image captioning
  - LSTM became the dominant approach

- Now (2019), other approaches (e.g. Transformers) have become more dominant for certain tasks.
  - For example in **WMT** (a MT conference + competition):
  - In WMT 2016, the summary report contains "RNN" 44 times
  - In WMT 2018, the report contains "RNN" 9 times and "Transformer" 63 times

**Source:** "Findings of the 2016 Conference on Machine Translation (WMT16)", Bojar et al. 2016, http://www.statmt.org/wmt16/pdf/W16-2301.pdf
**Source:** "Findings of the 2018 Conference on Machine Translation (WMT18)", Bojar et al. 2018, http://www.statmt.org/wmt18/pdf/WMT028.pdf
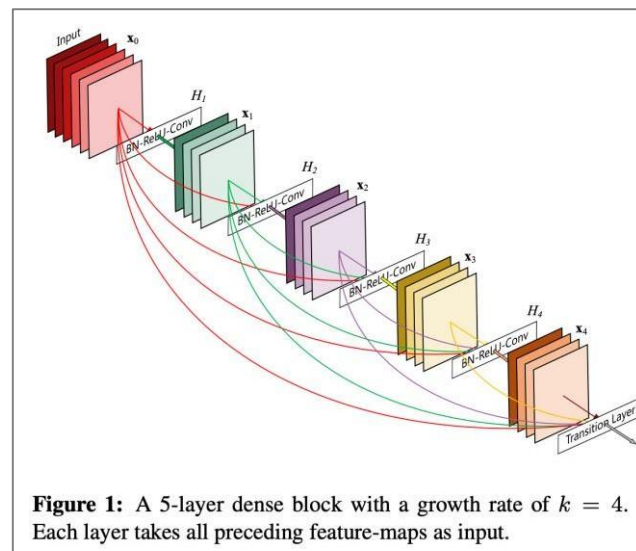
Abigail See

# Is vanishing/exploding gradient just a RNN problem?

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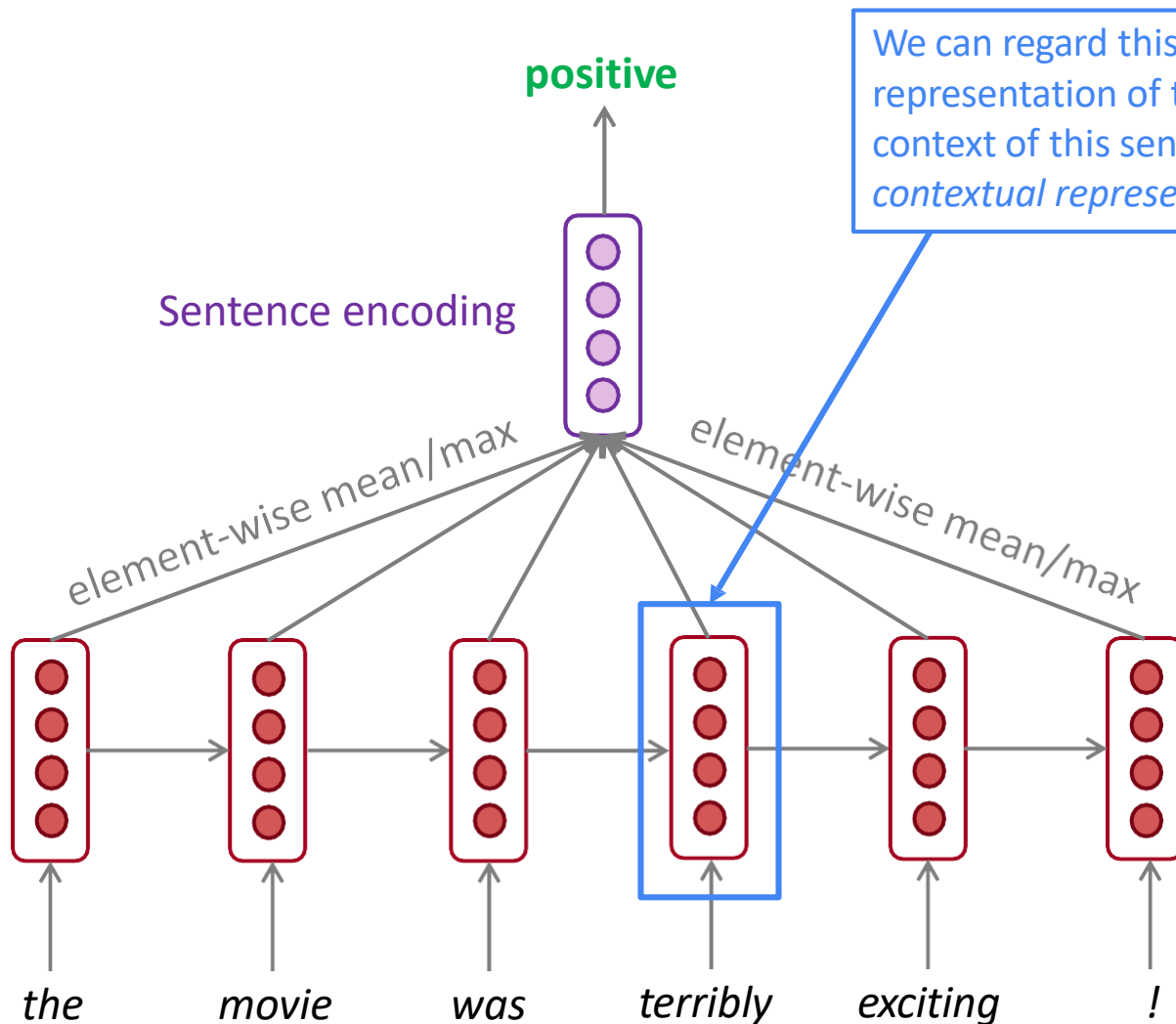
For example:

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



Figure 2. Residual learning: a building block.

"Deep Residual Learning for Image Recognition", He et al, 2015. https://arxiv.org/pdf/1512.03385.pdf

Abigail See

# Is vanishing/exploding gradient just a RNN problem?

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

For example:

- Dense connections aka "DenseNet"
- Directly connect everything to everything!



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

"Densely Connected Convolutional Networks", Huang et al, 2017. https://arxiv.org/pdf/1608.06993.pdf

Abigail See

# Bidirectional RNNs: motivation

Task: Sentiment Classification



positive

Sentence encoding

We can regard this hidden state as a representation of the word *"terribly"* in the context of this sentence. We call this a *contextual representation.*

element-wise mean/max

element-wise mean/max

These contextual representations only contain information about the *left* context (e.g. *"the movie was"*).

**What about *right* context?**

In this example, *"exciting"* is in the right context and this modifies the meaning of *"terribly"* (from negative to positive)

the     movie     was     terribly     exciting     !

Abigail See

# Bidirectional RNNs

This contextual representation of "terribly" has both left and right context!

Concatenated hidden states

Backward RNN

Forward RNN

*the*   *movie*   *was*   *terribly*   *exciting*   *!*

# Bidirectional RNNs

On timestep $t$:

This is a general notation to mean "compute one forward step of the RNN" – it could be a vanilla, LSTM or GRU computation.

Forward RNN $\quad \overrightarrow{\boldsymbol{h}}^{(t)} = \boxed{\text{RNN}_{\text{FW}}}(\overrightarrow{\boldsymbol{h}}^{(t-1)}, \boldsymbol{x}^{(t)})$

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

Generally, these two RNNs have separate weights

Concatenated hidden states $\quad \boxed{\boldsymbol{h}^{(t)}} = [\overrightarrow{\boldsymbol{h}}^{(t)}; \overleftarrow{\boldsymbol{h}}^{(t)}]$

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

Abigail See

# Multi-layer RNNs

RNN layer 3

RNN layer 2

RNN layer 1

the       movie       was       terribly       exciting       !

# Evaluating Language Models

- The standard evaluation metric for Language Models is perplexity.

$$\text{perplexity} = \prod_{t=1}^{T} \left( \frac{1}{P_{\text{LM}}(\boldsymbol{x}^{(t+1)} \mid \boldsymbol{x}^{(t)}, \ldots, \boldsymbol{x}^{(1)})} \right)^{1/T}$$

Normalized by number of words

Inverse probability of corpus, according to Language Model

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

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

**Lower** perplexity is better!

# Recap thus far

- **<u>Language Model</u>**: A system that predicts the next word

- **<u>Recurrent Neural Network</u>**: A family of neural networks that:
    - Take sequential input of any length
    - Apply the same weights on each step
    - Can optionally produce output on each step

- Vanishing gradient problem: what it is, why it happens, and why it's bad for RNNs

- LSTMs and GRUs: more complicated RNNs that use gates to control information flow; they are more resilient to vanishing gradients

# Plan for this lecture

- Recurrent neural networks
  - Basics
  - Training (backprop through time, vanishing gradient)
  - Recurrent networks with gates (GRU, LSTM)
- Applications in NLP and vision
  - Image/video captioning
  - Neural machine translation (beam search, attention)
- Transformers
  - Self-attention
  - BERT
  - Cross-modal transformers for VQA and VCR

# Applications

# Why should we care about Language Modeling?

- Language Modeling is a benchmark task that helps us measure our progress on understanding language

- Language Modeling is a subcomponent of many NLP tasks, especially those involving generating text or estimating the probability of text:

  - Predictive typing
  - Speech recognition
  - Handwriting recognition
  - Spelling/grammar correction
  - Authorship identification
  - Machine translation
  - Summarization
  - Dialogue
  - etc.

Abigail See

# Generating text with a RNN Language Model

You can use a RNN Language Model to generate text by repeated sampling. Sampled output is next step's input.

# Generating text with a RNN Language Model

- Let's have some fun!

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

- RNN-LM trained on Obama speeches:

*The United States will step up to the cost of a new challenges of the American people that will share the fact that we created the problem. They were attacked and so that they have to say that all the task of the final days of war that I will not be able to get this done.*

Abigail See

# Generating text with a RNN Language Model

- Let's have some fun!
- You can train a RNN-LM on any kind of text, then generate text in that style.
- RNN-LM trained on *Harry Potter*:



"Sorry," Harry shouted, panicking—"I'll leave those brooms in London, are they?"

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

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

Abigail See

# Generating text with a RNN Language Model

- Let's have some fun!

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

- RNN-LM trained on paint color names:

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

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

Abigail See

# Generating poetry with RNNs



### Sonnet 116 – Let me not …

*by William Shakespeare*

Let me not to the marriage of true minds
     Admit impediments. Love is not love
Which alters when it alteration finds,
     Or bends with the remover to remove:
O no! it is an ever–fixed mark
     That looks on tempests and is never shaken;
It is the star to every wandering bark,
     Whose worth's unknown, although his height be taken.
Love's not Time's fool, though rosy lips and cheeks
     Within his bending sickle's compass come:
Love alters not with his brief hours and weeks,
     But bears it out even to the edge of doom.
If this be error and upon me proved,
     I never writ, nor no man ever loved.

# Generating poetry with RNNs

at first:

```
tyntd-iafhatawiaoihrdemot  lytdws  e ,tfti, astai f ogoh eoase rrranbyne 'nhthnee e
plia tklrgd 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 shuwy 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.
```

More info: http://karpathy.github.io/2015/05/21/rnn-effectiveness/

Andrej Karpathy

# Generating poetry with RNNs

```
PANDARUS:
Alas, I think he shall be come approached and the day
When little srain would be attain'd into being never fed,
And who is but a chain and subjects of his death,
I should not sleep.

Second Senator:
They are away this miseries, produced upon my soul,
Breaking and strongly should be buried, when I perish
The earth and thoughts of many states.

DUKE VINCENTIO:
Well, your wit is in the care of side and that.

Second Lord:
They would be ruled after this chamber, and
my fair nues begun out of the fact, to be conveyed,
Whose noble souls I'll have the heart of the wars.

Clown:
Come, sir, I will make did behold your worship.

VIOLA:
I'll drink it.
```

```
VIOLA:
Why, Salisbury must find his flesh and thought
That which I am not aps, not a man and in fire,
To show the reining of the raven and the wars
To grace my hand reproach within, and not a fair are hand,
That Caesar and my goodly father's world;
When I was heaven of presence and our fleets,
We spare with hours, but cut thy council I am great,
Murdered and by thy master's ready there
My power to give thee but so much as hell:
Some service in the noble bondman here,
Would show him to her wine.

KING LEAR:
O, if you were a feeble sight, the courtesy of your law,
Your sight and several breath, will wear the gods
With his heads, and my hands are wonder'd at the deeds,
So drop upon your lordship's head, and your opinion
Shall be against your honour.
```

# Generating textbooks with RNNs

open source textbook on algebraic geometry



Latex source

# Generating textbooks with RNNs

For $\bigoplus_{n=1,\dots,m}$ where $\mathcal{L}_{m_\bullet} = 0$, hence we can find a closed subset $\mathcal{H}$ in $\mathcal{H}$ and any sets $\mathcal{F}$ on $X$, $U$ is a closed immersion of $S$, then $U \to T$ is a separated algebraic space.

*Proof.* Proof of (1). It also start we get

$$S = \mathrm{Spec}(R) = U \times_X U \times_X U$$

and the comparicoly in the fibre product covering we have to prove the lemma generated by $\coprod Z \times_U U \to V$. Consider the maps $M$ along the set of points $Sch_{fppf}$ and $U \to U$ is the fibre category of $S$ in $U$ in Section, **??** and the fact that any $U$ affine, see Morphisms, Lemma **??**. Hence we obtain a scheme $S$ and any open subset $W \subset U$ in $Sh(G)$ such that $\mathrm{Spec}(R') \to S$ is smooth or an

$$U = \bigcup U_i \times_{S_i} U_i$$

which has a nonzero morphism we may assume that $f_i$ is of finite presentation over $S$. We claim that $\mathcal{O}_{X,x}$ is a scheme where $x, x', s'' \in S'$ such that $\mathcal{O}_{X,x'} \to \mathcal{O}'_{X',x'}$ is separated. By Algebra, Lemma **??** we can define a map of complexes $\mathrm{GL}_{S'}(x'/S'')$ and we win.

To prove study we see that $\mathcal{F}|_U$ is a covering of $\mathcal{X}'$, and $\mathcal{T}_i$ is an object of $\mathcal{F}_{X/S}$ for $i > 0$ and $\mathcal{F}_p$ exists and let $\mathcal{F}_i$ be a presheaf of $\mathcal{O}_X$-modules on $\mathcal{C}$ as a $\mathcal{F}$-module. In particular $\mathcal{F} = U/\mathcal{F}$ we have to show that

$$\widetilde{M}^\bullet = \mathcal{I}^\bullet \otimes_{\mathrm{Spec}(k)} \mathcal{O}_{S,s} - i_X^{-1}\mathcal{F})$$

is a unique morphism of algebraic stacks. Note that

$$\mathrm{Arrows} = (Sch/S)_{fppf}^{opp}, (Sch/S)_{fppf}$$

and

$$V = \Gamma(S,\mathcal{O}) \longmapsto (U, \mathrm{Spec}(A))$$

is an open subset of $X$. Thus $U$ is affine. This is a continuous map of $X$ is the inverse, the groupoid scheme $S$.

*Proof.* See discussion of sheaves of sets. $\square$

The result for prove any open covering follows from the less of Example **??**. It may replace $S$ by $X_{spaces,\acute{e}tale}$ which gives an open subspace of $X$ and $T$ equal to $S_{Zar}$, see Descent, Lemma **??**. Namely, by Lemma **??** we see that $R$ is geometrically regular over $S$.

---

**Lemma 0.1.** *Assume (3) and (3) by the construction in the description.*

*Suppose $X = \lim |X|$ (by the formal open covering $X$ and a single map $\underline{\mathrm{Proj}}_X(\mathcal{A}) = \mathrm{Spec}(B)$ over $U$ compatible with the complex*

$$\mathrm{Set}(\mathcal{A}) = \Gamma(X, \mathcal{O}_{X,\mathcal{O}_X}).$$

*When in this case of to show that $\mathcal{Q} \to \mathcal{C}_{Z/X}$ is stable under the following result in the second conditions of (1), and (3). This finishes the proof. By Definition **??** (without element is when the closed subschemes are catenary. If $T$ is surjective we may assume that $T$ is connected with residue fields of $S$. Moreover there exists a closed subspace $Z \subset X$ of $X$ where $U$ in $X'$ is proper (some defining as a closed subset of the uniqueness it suffices to check the fact that the following theorem*

   (1) *$f$ is locally of finite type. Since $S = \mathrm{Spec}(R)$ and $Y = \mathrm{Spec}(R)$.*

*Proof.* This is form all sheaves of sheaves on $X$. But given a scheme $U$ and a surjective étale morphism $U \to X$. Let $U \cap U = \coprod_{i=1,\dots,n} U_i$ be the scheme $X$ over $S$ at the schemes $X_i \to X$ and $U = \lim_i X_i$. $\square$

The following lemma surjective restrocomposes of this implies that $\mathcal{F}_{x_0} = \mathcal{F}_{x_0} = \mathcal{F}_{\mathcal{X},\dots,0}$.

**Lemma 0.2.** *Let $X$ be a locally Noetherian scheme over $S$, $E = \mathcal{F}_{X/S}$. Set $\mathcal{I} = \mathcal{J}_1 \subset \mathcal{I}'_n$. Since $\mathcal{I}^n \subset \mathcal{I}^n$ are nonzero over $i_0 \leq \mathfrak{p}$ is a subset of $\mathcal{J}_{n,0} \circ \overline{A}_2$ works.*

**Lemma 0.3.** *In Situation **??**. Hence we may assume $\mathfrak{q}' = 0$.*

*Proof.* We will use the property we see that $\mathfrak{p}$ is the mext functor (**??**). On the other hand, by Lemma **??** we see that

$$D(\mathcal{O}_{X'}) = \mathcal{O}_X(D)$$

where $K$ is an $F$-algebra where $\delta_{n+1}$ is a scheme over $S$. $\square$

Andrej Karpathy

# Generating textbooks with RNNs



*Proof.* Omitted. □

**Lemma 0.1.** *Let $\mathcal{C}$ be a set of the construction.*
*Let $\mathcal{C}$ be a gerber covering. Let $\mathcal{F}$ be a quasi-coherent sheaves of $\mathcal{O}$-modules. We have to show that*

$$\mathcal{O}_{\mathcal{O}_X} = \mathcal{O}_X(\mathcal{L})$$

.

*Proof.* This is an algebraic space with the composition of sheaves $\mathcal{F}$ on $X_{\acute{e}tale}$ we have

$$\mathcal{O}_X(\mathcal{F}) = \{morph_1 \times_{\mathcal{O}_X} (\mathcal{G}, \mathcal{F})\}$$

where $\mathcal{G}$ defines an isomorphism $\mathcal{F} \to \mathcal{F}$ of $\mathcal{O}$-modules. □

**Lemma 0.2.** *This is an integer $\mathcal{Z}$ is injective.*

*Proof.* See Spaces, Lemma ??. □

**Lemma 0.3.** *Let $S$ be a scheme. Let $X$ be a scheme and $X$ is an affine open covering. Let $\mathcal{U} \subset \mathcal{X}$ be a canonical and locally of finite type. Let $X$ be a scheme. Let $X$ be a scheme which is equal to the formal complex.*

*The following to the construction of the lemma follows.*

*Let $X$ be a scheme. Let $X$ be a scheme covering. Let*

$$b: X \to Y' \to Y \to Y \to Y' \times_X Y \to X.$$

*be a morphism of algebraic spaces over $S$ and $Y$.*

*Proof.* Let $X$ be a nonzero scheme of $X$. Let $X$ be an algebraic space. Let $\mathcal{F}$ be a quasi-coherent sheaf of $\mathcal{O}_X$-modules. The following are equivalent

(1) $\mathcal{F}$ is an algebraic space over $S$.
(2) If $X$ is an affine open covering.

Consider a common structure on $X$ and $X$ the functor $\mathcal{O}_X(U)$ which is locally of finite type. □

This since $\mathcal{F} \in \mathcal{F}$ and $x \in \mathcal{G}$ the diagram

$$\mathrm{Spec}(K_\psi) \qquad \mathrm{Mor}_{Sets} \quad \mathrm{d}(\mathcal{O}_{X_{X/k}}, \mathcal{G})$$

is a limit. Then $\mathcal{G}$ is a finite type and assume $S$ is a flat and $\mathcal{F}$ and $\mathcal{G}$ is a finite type $f_*$. This is of finite type diagrams, and

- the composition of $\mathcal{G}$ is a regular sequence,
- $\mathcal{O}_{X'}$ is a sheaf of rings. □

*Proof.* We have see that $X = \mathrm{Spec}(R)$ and $\mathcal{F}$ is a finite type representable by algebraic space. The property $\mathcal{F}$ is a finite morphism of algebraic stacks. Then the cohomology of $X$ is an open neighbourhood of $U$. □

*Proof.* This is clear that $\mathcal{G}$ is a finite presentation, see Lemmas ??.
A *reduced above* we conclude that $U$ is an open covering of $\mathcal{C}$. The functor $\mathcal{F}$ is a "field

$$\mathcal{O}_{X,x} \longrightarrow \mathcal{F}_{\overline{x}} \quad -1(\mathcal{O}_{X_{\acute{e}tale}}) \longrightarrow \mathcal{O}_{X_i}^{-1} \mathcal{O}_{X_\lambda}(\mathcal{O}_{X_\eta}^{\overline{v}})$$

is an isomorphism of covering of $\mathcal{O}_{X_i}$. If $\mathcal{F}$ is the unique element of $\mathcal{F}$ such that $X$ is an isomorphism.
The property $\mathcal{F}$ is a disjoint union of Proposition ?? and we can filtered set of presentations of a scheme $\mathcal{O}_X$-algebra with $\mathcal{F}$ are opens of finite type over $S$.
If $\mathcal{F}$ is a scheme theoretic image points. □

If $\mathcal{F}$ is a finite direct sum $\mathcal{O}_{X_\lambda}$ is a closed immersion, see Lemma ??. This is a sequence of $\mathcal{F}$ is a similar morphism.

Andrej Karpathy

# Generating code with RNNs

```c
static void do_command(struct seq_file *m, void *v)
{
  int column = 32 << (cmd[2] & 0x80);
  if (state)
    cmd = (int)(int_state ^ (in_8(&ch->ch_flags) & Cmd) ? 2 : 1);
  else
    seq = 1;
  for (i = 0; i < 16; i++) {
    if (k & (1 << 1))
      pipe = (in_use & UMXTHREAD_UNCCA) +
        ((count & 0x00000000ffffff8) & 0x000000f) << 8;
    if (count == 0)
      sub(pid, ppc_md.kexec_handle, 0x20000000);
    pipe_set_bytes(i, 0);
  }
  /* Free our user pages pointer to place camera if all dash */
  subsystem_info = &of_changes[PAGE_SIZE];
  rek_controls(offset, idx, &soffset);
  /* Now we want to deliberately put it to device */
  control_check_polarity(&context, val, 0);
  for (i = 0; i < COUNTER; i++)
    seq_puts(s, "policy ");
}
```

Generated
C code

# Image Captioning



CVPR 2015:
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 Captioning



**Recurrent Neural Network**

**Convolutional Neural Network**

# Image Captioning



test image

| image |
|---|

| conv-64 |
|---|
| conv-64 |
| maxpool |

| conv-128 |
|---|
| conv-128 |
| maxpool |

| conv-256 |
|---|
| conv-256 |
| maxpool |

| conv-512 |
|---|
| conv-512 |
| maxpool |

| conv-512 |
|---|
| conv-512 |
| maxpool |

| FC-4096 |
|---|
| FC-4096 |
| FC-1000 |
| softmax |



test image

Andrej Karpathy

image

conv-64
conv-64
maxpool

conv-128
conv-128
maxpool

conv-256
conv-256
maxpool

conv-512
conv-512
maxpool

conv-512
conv-512
maxpool

FC-4096
FC-4096

X

test image

Andrej Karpathy

# Image Captioning



test image

| |
|---|
| image |
| conv-64 |
| conv-64 |
| maxpool |
| conv-128 |
| conv-128 |
| maxpool |
| conv-256 |
| conv-256 |
| maxpool |
| conv-512 |
| conv-512 |
| maxpool |
| conv-512 |
| conv-512 |
| maxpool |
| FC-4096 |
| FC-4096 |

x0
<START>

<START>

# Image Captioning



test image

**before:**

$h = \tanh(W_{xh} * x + W_{hh} * h)$

**now:**

$h = \tanh(W_{xh} * x + W_{hh} * h + W_{ih} * im)$

image
conv-64
conv-64
maxpool

conv-128
conv-128
maxpool

conv-256
conv-256
maxpool

conv-512
conv-512
maxpool

conv-512
conv-512
maxpool

FC-4096
FC-4096

im

**Wih**

y0

h0

x0
<START>

<START>

Andrej Karpathy

# Image Captioning



test image

sample!

<START>

# Image Captioning



test image

<START>

Andrej Karpathy

# Image Captioning



test image

sample!

y0   y1

h0 → h1

x0
<START>

straw    hat

<START>

image
conv-64
conv-64
maxpool
conv-128
conv-128
maxpool
conv-256
conv-256
maxpool
conv-512
conv-512
maxpool
conv-512
conv-512
maxpool
FC-4096
FC-4096

# Image Captioning



test image

image

conv-64
conv-64
maxpool

conv-128
conv-128
maxpool

conv-256
conv-256
maxpool

conv-512
conv-512
maxpool

conv-512
conv-512
maxpool

FC-4096
FC-4096

y0  y1  y2

h0 → h1 → h2

x0
<START>  straw  hat

<START>

# Image Captioning



test image

Caption generated:
"straw hat"

sample
<END> token
=> finish.

<START>

# Image Captioning



"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."

Andrej Karpathy

# Video Captioning



Key Insight:

Generate feature representation of the video and "decode" it to a sentence

Venugopalan et al., "Translating Videos to Natural Language using Deep Recurrent Neural Networks", NAACL-HTL 2015

# Video Captioning

Input Video     Convolutional Net     Recurrent Net     Output



$$\frac{1}{n}\sum$$

LSTM → LSTM → *A*

LSTM → LSTM → *boy*

LSTM → LSTM → *is*

LSTM → LSTM → *playing*

LSTM → LSTM → *golf*

LSTM → LSTM → *<EOS>*

Mean across
all frames

Venugopalan et al., "Translating Videos to Natural Language using Deep Recurrent Neural Networks", NAACL-HTL 2015

# Video Captioning



FGM: A person is dancing with the person on the stage.

YT: A group of men are riding the forest.

I+V: **A group of people are dancing.**

GT: Many men and women are dancing in the street.



FGM: A person is cutting a potato in the kitchen.

YT: A man is slicing a tomato.

I+V: **A man is slicing a carrot.**

GT: A man is slicing carrots.



FGM: A person is walking with a person in the forest.

YT: A monkey is walking.

I+V: **A bear is eating a tree.**

GT: Two bear cubs are digging into dirt and plant matter at the base of a tree.



FGM: A person is riding a horse on the stage.

YT: A group of playing are playing in the ball.

I+V: **A basketball player is playing**.

GT: Dwayne wade does a fancy layup in an allstar game.

Venugopalan et al., "Translating Videos to Natural Language using Deep Recurrent Neural Networks", NAACL-HTL 2015

# Video Captioning



| | | | | | | |
|---|---|---|---|---|---|---|
| English Sentence | → | RNN encoder | ○○○ | RNN decoder | → French Sentence | [Sutskever et al. NIPS'14] |
| 📷 | → | Encode | ○○○ | RNN decoder | → Sentence | [Donahue et al. CVPR'15] [Vinyals et al. CVPR'15] |
| ▶ | → | Encode | ○○○ | RNN decoder | → Sentence | [Venugopalan et. al. NAACL'15] |
| ▶ | | RNN encoder | ○○○ | RNN decoder | → Sentence | [Venugopalan et. al. ICCV'15] (this work) |

3

Venugopalan et al., "Sequence to Sequence - Video to Text", ICCV 2015

# Video Captioning

S2VT Overview



Now decode it to a sentence!

Encoding stage

Decoding stage

A man is talking ...

Venugopalan et al., "Sequence to Sequence - Video to Text", ICCV 2015

# Neural Machine Translation?

- Neural Machine Translation (NMT) is a way to do Machine Translation with a *single neural network*

- The neural network architecture is called sequence-to-sequence (aka seq2seq) and it involves *two* RNNs.

Abigail See

# Neural Machine Translation (NMT)

The sequence-to-sequence model



Encoding of the source sentence.
Provides initial hidden state
for Decoder RNN.

Target sentence (output)

he    hit    me    with    a    pie    <END>

Encoder RNN

argmax    argmax    argmax    argmax    argmax    argmax    argmax

il    a    m'    entarté

<START>    he    hit    me    with    a    pie

Decoder RNN

Source sentence (input)

Encoder RNN produces
an encoding of the
source sentence.

Decoder RNN is a Language Model that generates
target sentence, *conditioned on encoding*.

Note: This diagram shows **test time** behavior:
decoder output is fed in ·······> as next step's input

Abigail See

# Greedy decoding

- We saw how to generate (or "decode") the target sentence by taking argmax on each step of the decoder



- This is greedy decoding (take most probable word on each step)
- **Problems with this method?**

# Problems with greedy decoding

- Greedy decoding has no way to undo decisions!
  - Input: *il am'entarté*        *(he hit me with a pie)*
  - →*he* _____
  - →*he hit* _____
  - →*he hit* <span style="color:red">*a*</span> _____        <span style="color:red">(whoops! no going back now…)</span>

- How to fix this?

Abigail See

# Exhaustive search decoding

- Ideally we want to find a (length *T*) translation *y* that maximizes

$$P(y|x) = P(y_1|x)\,P(y_2|y_1,x)\,P(y_3|y_1,y_2,x)\ldots,P(y_T|y_1,\ldots,y_{T-1},x)$$
$$= \prod_{t=1}^{T} P(y_t|y_1,\ldots,y_{t-1},x)$$

- We could try computing all possible sequences  *y*
  - This means that on each step *t* of the decoder, we're tracking $V^T$ possible  partial translations, where *V*  is vocabulary size
  - This $O(V^T)$ complexity is far too expensive!

# Beam search decoding

- <u>Core idea:</u> On each step of decoder, keep track of the *k* most  probable partial translations (*hypotheses*)
  - *k* is the beam size (in practice around 5 to 10)

- A hypothesis $y_1, \cdots, y_t$ has a score which is its log probability:

$$\text{score}(y_1, \ldots, y_t) = \log P_{\text{LM}}(y_1, \ldots, y_t | x) = \sum_{i=1}^{t} \log P_{\text{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$$

  - Scores are all negative, and higher score is better
  - We search for high-scoring hypotheses, tracking top *k* on each step

- Beam search is not guaranteed to find optimal solution
- But much more efficient than exhaustive  search!

Abigail See

# Beam search decoding: example

Beam size = k = 2. Blue numbers = $\text{score}(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{\text{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$

<START>

Calculate prob
dist of next word

# Beam search decoding: example

Beam size = k = 2. Blue numbers = $\mathrm{score}(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{\mathrm{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$

-0.7 = log P$_{LM}$(*he*|*<START>*)

*he*

*<START>*

*I*

-0.9 = log P$_{LM}$(*I*|*<START>*)

Take top *k* words
and compute scores

Abigail See

# Beam search decoding: example

Beam size = k = 2. Blue numbers = $\mathrm{score}(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{\mathrm{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$

-1.7 = log P$_{\text{LM}}$(*hit*|*<START> he*) + -0.7

-0.7

*he*

*hit*

*struck*

-2.9 = log P$_{\text{LM}}$(*struck*|*<START> he*) + -0.7

*<START>*

-1.6 = log P$_{\text{LM}}$(*was*|*<START> I*) + -0.9

*was*

*I*

*got*

-0.9

-1.8 = log P$_{\text{LM}}$(*got*|*<START> I*) + -0.9

For each of the *k* hypotheses, find
top *k* next words and calculate scores

Abigail See

# Beam search decoding: example

Beam size = k = 2. Blue numbers = $\text{score}(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{\text{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$



-1.7
-0.7

hit

he

struck

-2.9

&lt;START&gt;

-1.6

was

I

got

-0.9

-1.8

Of these $k^2$ hypotheses,
just keep $k$ with highest scores

Abigail See

# Beam search decoding: example

Beam size = k = 2. Blue numbers = $\mathrm{score}(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{\mathrm{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$

-2.8 = log P$_{LM}$(*a*|*<START> he hit*) + -1.7

-1.7

-0.7

*hit*

*a*

*he*

*struck*

*me*

-2.9

-2.5 = log P$_{LM}$(*me*|*<START> he hit*) + -1.7

*<START>*

-2.9 = log P$_{LM}$(*hit*|*<START> I was*) + -1.6

-1.6

*hit*

*was*

*I*

*struck*

*got*

-0.9

-1.8

-3.8 = log P$_{LM}$(*struck*|*<START> I was*) + -1.6

For each of the *k* hypotheses, find top *k* next words and calculate scores
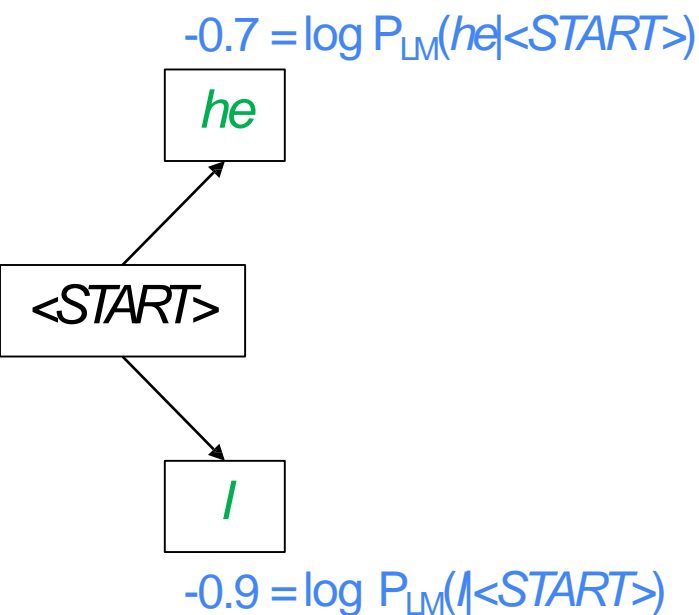
Abigail See

# Beam search decoding: example

Beam size = k = 2. Blue numbers = $\text{score}(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{\text{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$



-2.8
*a*

-1.7
*hit*

-0.7
*he*

*struck*
-2.9

*me*
-2.5

<START>

-2.9
*hit*

-1.6
*was*

*I*
-0.9

*got*
-1.8

*struck*
-3.8

Of these *k²* hypotheses,
just keep *k* with highest scores

Abigail See

# Beam search decoding: example

Beam size = k = 2. Blue numbers = $\mathrm{score}(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{\mathrm{LM}}(y_i|y_1, \ldots, y_{i-1}, x)$



For each of the *k* hypotheses, find
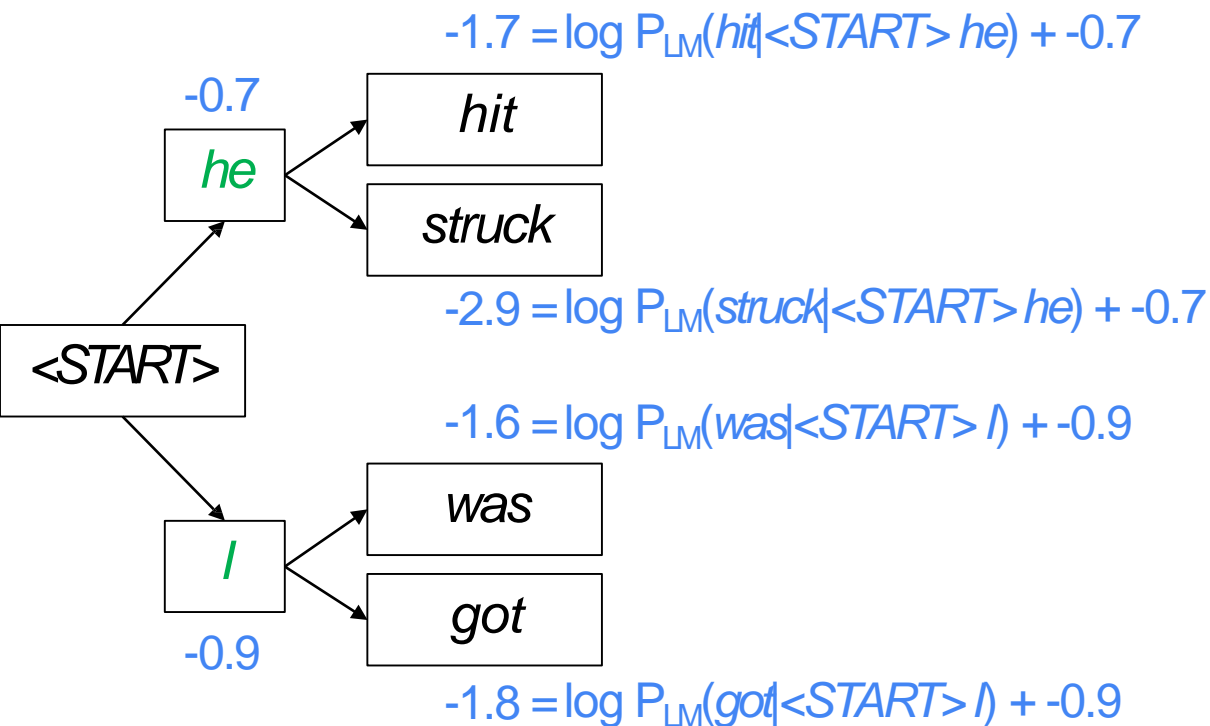top *k* next words and calculate scores

Abigail See

# Beam search decoding: example

Beam size = k = 2. Blue numbers = $\text{score}(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{\text{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$



-0.7 → he
-0.9 → I

he → hit (-1.7)
he → struck (-2.9)

I → was (-1.6)
I → got (-1.8)

hit → a (-2.8)
hit → me (-2.5)

was → hit (-2.9)
was → struck (-3.8)

a → tart (-4.0)
a → pie (-3.4)

me → with (-3.3)
me → on (-3.5)

Of these $k^2$ hypotheses, just keep $k$ with highest scores

Abigail See

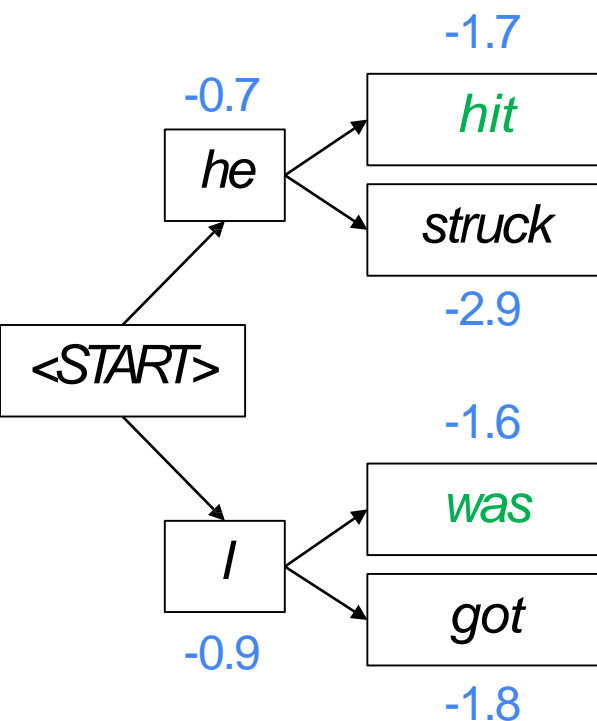# Beam search decoding: example

Beam size = k = 2. Blue numbers = $\text{score}(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{\text{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$



For each of the *k* hypotheses, find top *k* next words and calculate scores
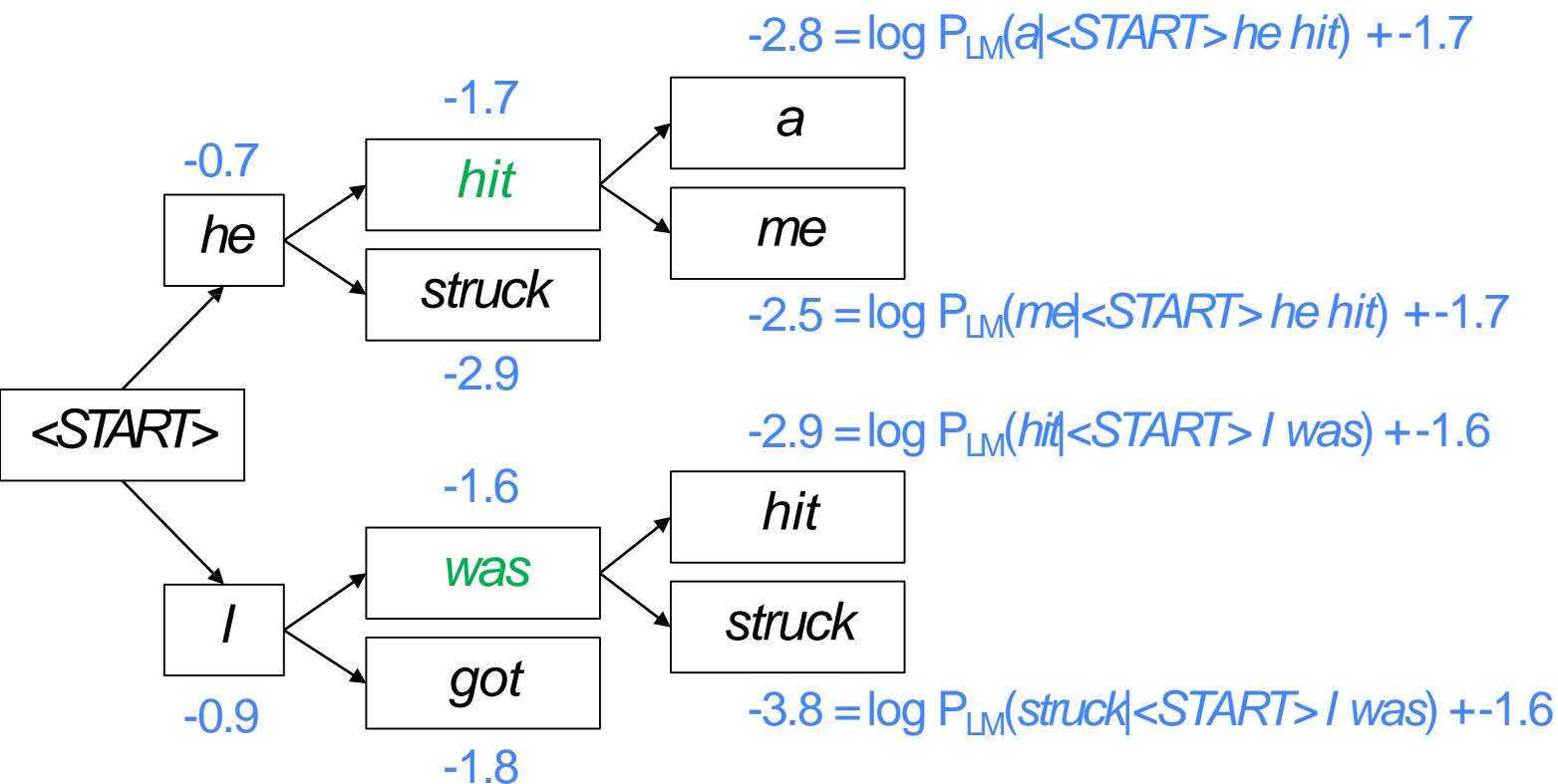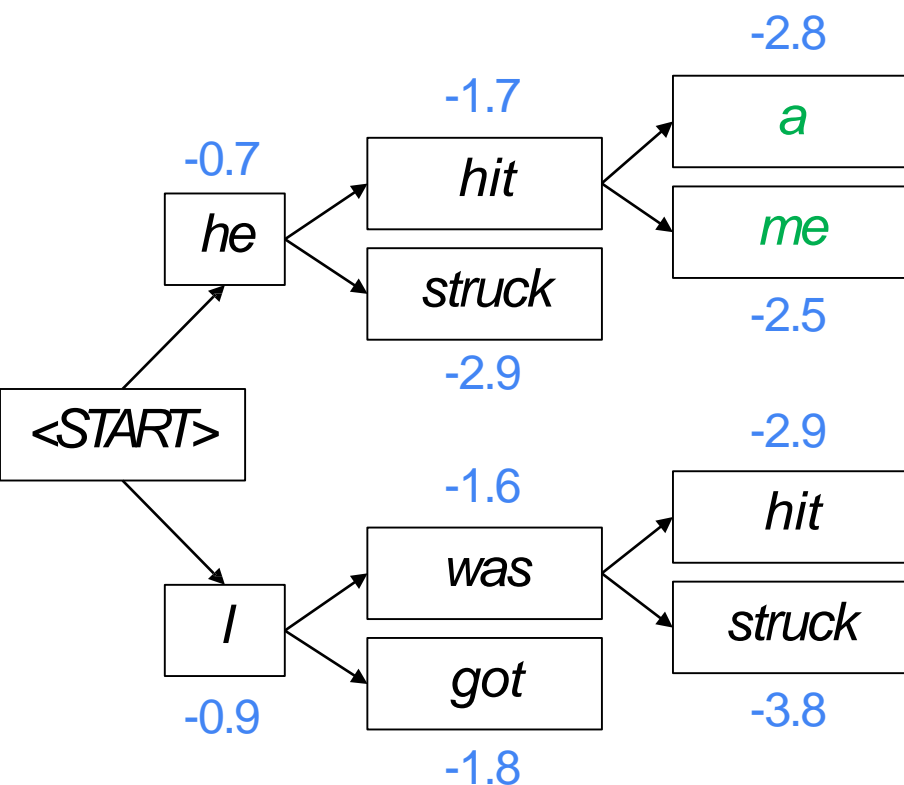
# Beam search decoding: example

Beam size = k = 2. Blue numbers = $\text{score}(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{\text{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$



-4.0    -4.8

| tart | | in |

-2.8

| pie | | with |
-3.4    -4.5

-1.7

| hit | | a |

-0.7

| he | | struck | | me |
-2.9    -2.5

-3.3    -3.7

| with | | a |

| on | | one |
-3.5    -4.3

<START>

-1.6    -2.9

| was | | hit |

-0.9

| I | | got | | struck |
-1.8    -3.8

Abigail See

# Beam search decoding: example

Beam size = k = 2. Blue numbers = $\mathrm{score}(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{\mathrm{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$



For each of the *k* hypotheses, find
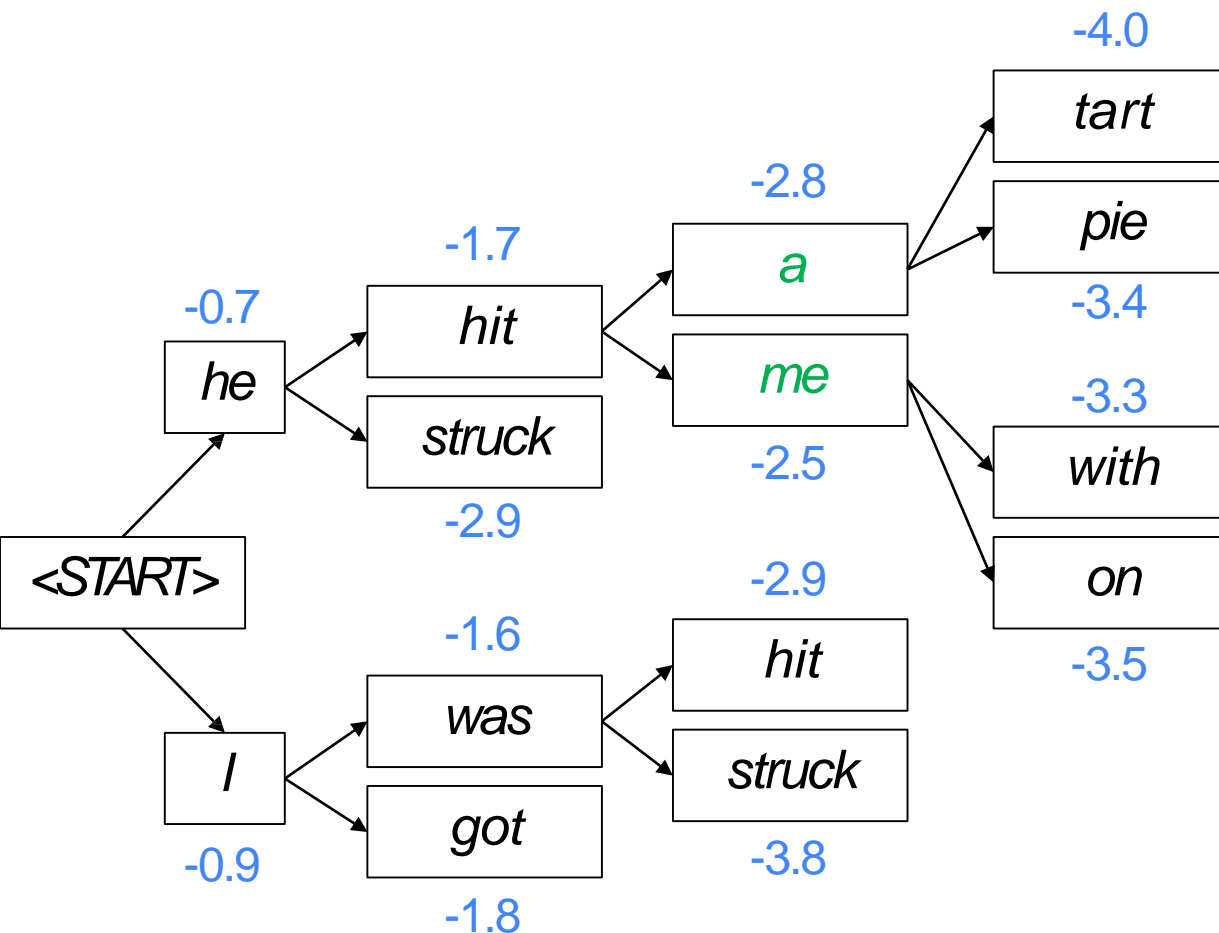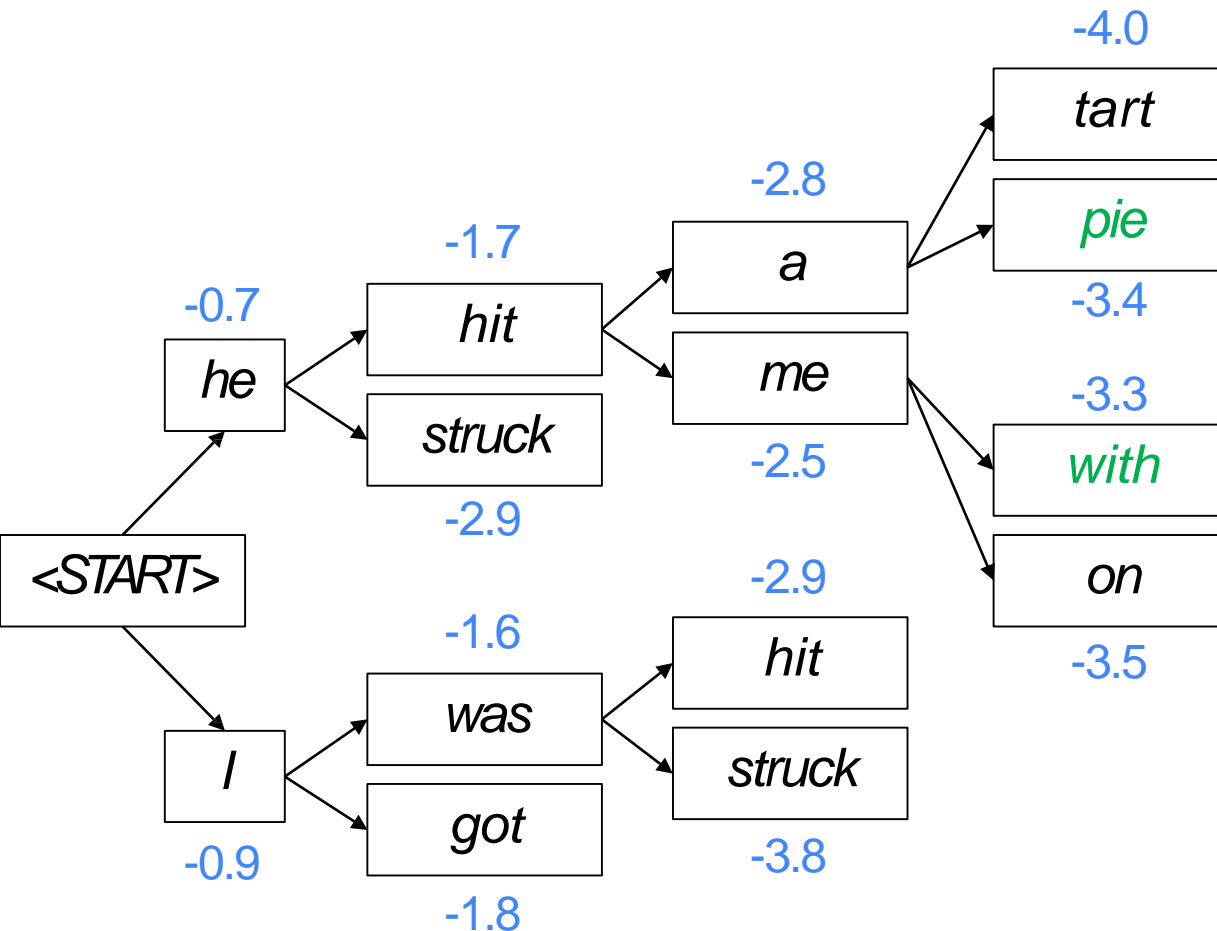top *k* next words and calculate scores

# Beam search decoding: example

Beam size = k = 2. Blue numbers = $\mathrm{score}(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{\mathrm{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$

-4.0
tart

-4.8
in

-2.8
a

pie
-3.4

with
-4.5

-1.7
hit

me
-2.5

-4.3
*pie*

-0.7
he

struck
-2.9

with
-3.3

a
-3.7

tart

on
-3.5

one
-4.3

-4.6

<START>

-2.9
hit

pie
-5.0

-1.6
was

struck
-3.8

tart

I

got
-1.8

-0.9

-5.3

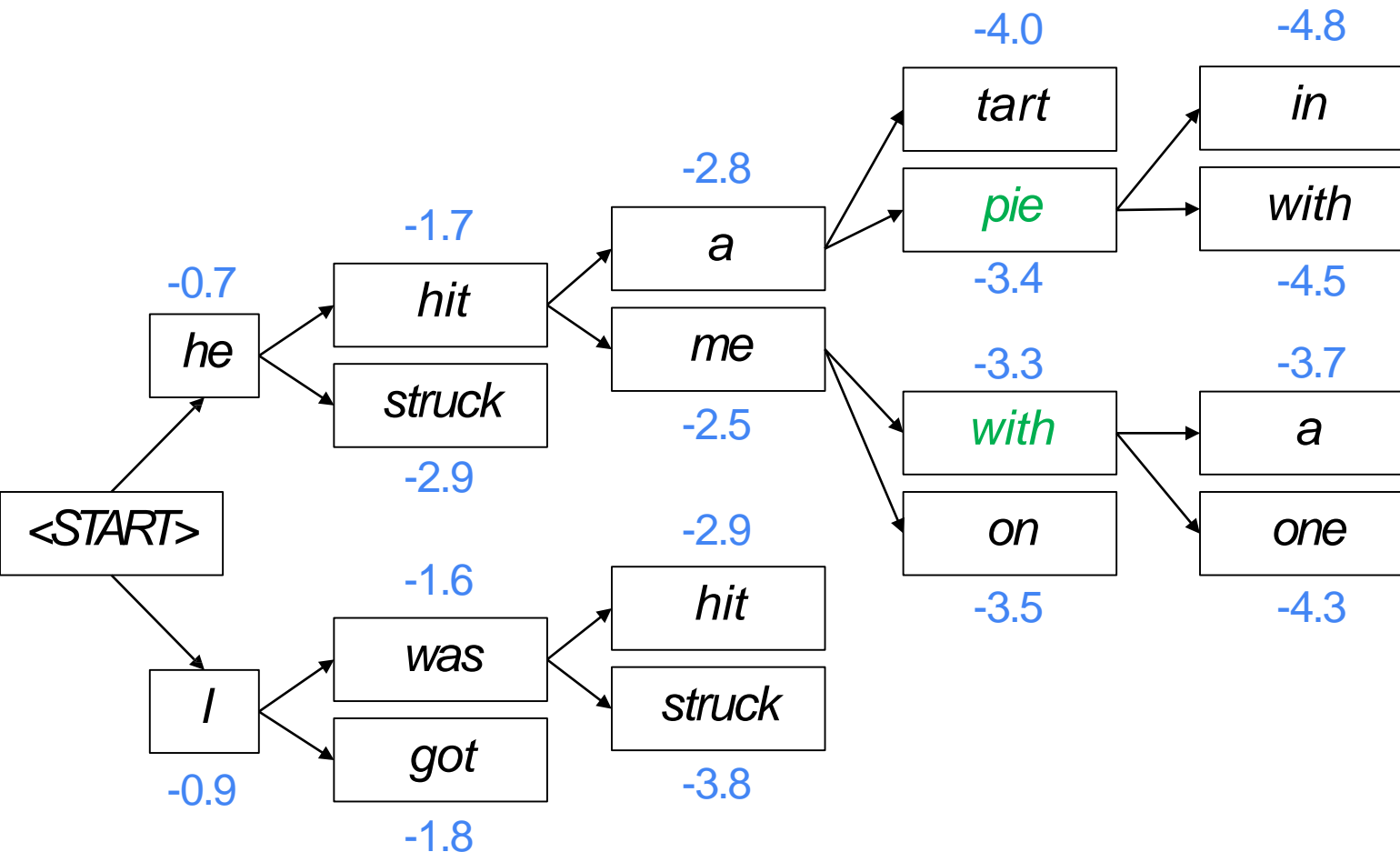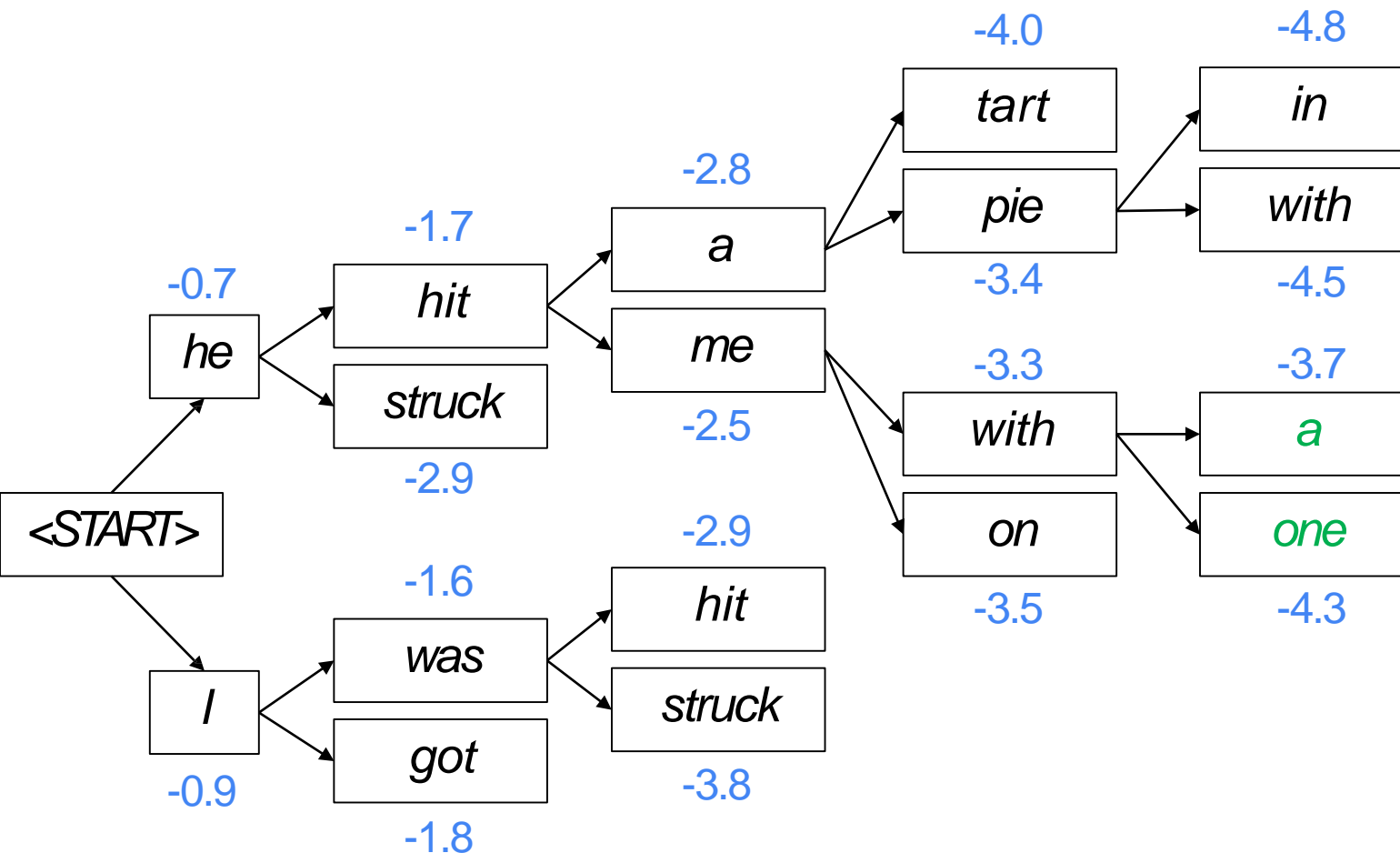This is the top-scoring hypothesis!

# Beam search decoding: example

Beam size = k = 2. Blue numbers = $\mathrm{score}(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{\mathrm{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$



Backtrack to obtain the full hypothesis

Abigail See

# Beam search decoding: finishing up

- We have our list of completed hypotheses.
- How to select top one with highest score?

- Each hypothesis $y_1, \ldots, y_t$ on our list has a score

$$\text{score}(y_1, \ldots, y_t) = \log P_{\text{LM}}(y_1, \ldots, y_t | x) = \sum_{i=1}^{t} \log P_{\text{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$$

- <u>Problem with this:</u> longer hypotheses have lower scores

- <u>Fix:</u> Normalize by length. Use this to select top one instead:

$$\frac{1}{t} \sum_{i=1}^{t} \log P_{\text{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$$

Abigail See

# How do we evaluate Machine Translation?

**BLEU** (**Bil**ingual **E**valuation **U**nderstudy)

- BLEU compares the <u>machine-written translation</u> to one or several <u>human-written translation</u>(s), and computes a similarity score based on:
  - *n*-gram precision (usually for 1, 2, 3 and 4-grams)
  - Plus a penalty for too-short system translations

- BLEU is useful but imperfect
  - There are many valid ways to translate a sentence
  - So a good translation can get a poor BLEU score because it has low *n*-gram overlap with the human translation ☹

**Source:** "BLEU: a Method for Automatic Evaluation of Machine Translation", Papineni et al, 2002.

Abigail See

# MT progress overtime

[Edinburgh En-De WMT newstest2013 Cased BLEU; NMT 2015 from U. Montréal]

**Source**: http://www.meta-net.eu/events/meta-forum-2016/slides/09_sennrich.pdf

# NMT: the biggest success story of NLP Deep Learning

Neural Machine Translation went from a fringe research activity in **2014** to the leading standard method in **2016**

- **2014**: First seq2seq paper published

- **2016**: Google Translate switches from SMT to NMT

- This is amazing!
  - **SMT** systems, built by hundreds of engineers over many years, outperformed by NMT systems trained by a handful of engineers in a few months

# So is Machine Translation solved?

- **Nope!**
- Many difficulties remain:
  - Out-of-vocabulary words
  - Domain mismatch between train and test data
  - Maintaining context over longer text
  - Low-resource language pairs

**Further reading:** "*Has AI surpassed humans at translation? Not even close!*"
https://www.skynettoday.com/editorials/state_of_nmt

Abigail See

# So is Machine Translation solved?

- **Nope!**
- Using common sense is still hard



Open in Google Translate                    Feedback

**?**

Abigail See

# So is Machine Translation solved?

- **Nope!**
- NMT picks up biases in training data



Didn't specify gender

Abigail See

# NMT research continues

NMT is the **flagship task** for NLP Deep Learning

- NMT research has pioneered many of the recent innovations of NLP Deep Learning

- In **2019**: NMT research continues to thrive
  - Researchers have found *many, many* improvements to the "vanilla" seq2seq NMT system
  - But one improvement is so integral that it is the new vanilla…

# ATTENTION

Abigail See

# Sequence-to-sequence: the bottleneck problem



Encoding of the source sentence.

Target sentence (output)

Encoder RNN

Decoder RNN

he    hit    me    with    a    pie    <END>

il    a    m'    entarté

<START>    he    hit    me    with    a    pie

Source sentence (input)

Problems with this architecture?

Abigail See

# Sequence-to-sequence: the bottleneck problem

Encoding of the source sentence. This needs to capture *all information* about the source sentence. Information bottleneck!

Target sentence (output)

he    hit    me    with    a    pie    <END>

Encoder RNN

Decoder RNN

il    a    m'    entarté

<START>    he    hit    me    with    a    pie

Source sentence (input)

# Attention

- **Attention** provides a solution to the bottleneck problem.

- <u>Core idea</u>: on each step of the decoder, use *direct connection to the encoder* to *focus on a particular part* of the source sequence

- First we will show via diagram (no equations), then we will show with equations

# Sequence-to-sequence with attention



dot product

Attention scores

Encoder RNN

Decoder RNN

*il*   *a*   *m'*   *entarté*   *<START>*

Source sentence (input)

Abigail See

# Sequence-to-sequence with attention



dot product

Attention scores

Encoder RNN

Decoder RNN

*il*  *a*  *m'*  *entarté*  *<START>*

Source sentence (input)

Abigail See

# Sequence-to-sequence with attention



dot product

Attention scores

Encoder RNN

Decoder RNN

il     a     m'     entarté          <START>

Source sentence (input)

Abigail See

# Sequence-to-sequence with attention



dot product

Attention scores

Encoder RNN

Decoder RNN

il    a    m'    entarté    <START>

Source sentence (input)

Abigail See

# Sequence-to-sequence with attention



On this decoder timestep, we're mostly focusing on the first encoder hidden state ("he")

Take softmax to turn the scores into a probability distribution

Attention distribution

Attention scores

Encoder RNN

Decoder RNN

il    a    m'    entarté        <START>

Source sentence (input)

Abigail See

# Sequence-to-sequence with attention



Use the attention distribution to take a **weighted sum** of the encoder hidden states.

The attention output mostly contains information from the hidden states that received high attention.

Attention output

Attention distribution

Attention scores

Encoder RNN

Decoder RNN

*il    a    m'    entarté*    <START>

Source sentence (input)

Abigail See

# Sequence-to-sequence with attention



Attention output

Attention distribution

Attention scores

Encoder RNN

Decoder RNN

he

$y_1$

Concatenate attention output with decoder hidden state, then use to compute $y_1$ as before

il    a    m'    entarté    <START>

Source sentence (input)

Abigail See

# Sequence-to-sequence with attention



Attention output

Attention distribution

Attention scores

Encoder RNN

Decoder RNN

*il    a    m'    entarté*

*<START>    he*

*hit*

$y_2$

Sometimes we take the attention output from the previous step, and also feed it into the decoder (along with the usual decoder input).

Source sentence (input)

Abigail See

# Attention: in equations

- We have encoder hidden states $h_1, \ldots, h_N \in \mathbb{R}^h$

- On timestep $t$, we have decoder hidden state $s_t \in \mathbb{R}^h$

- We get the attention scores $e^t$ for this step:

$$e^t = [s_t^T h_1, \ldots, s_t^T h_N] \in \mathbb{R}^N$$

- We take softmax to get the attention distribution $\alpha^t$ for this step (this is a probability distribution and sums to 1)

$$\alpha^t = \mathrm{softmax}(e^t) \in \mathbb{R}^N$$

- We use $\alpha^t$ to take a weighted sum of the encoder hidden states to get the attention output $a_t$

$$a_t = \sum_{i=1}^{N} \alpha_i^t h_i \in \mathbb{R}^h$$

- Finally we concatenate the attention output $a_t$ with the decoder hidden state $s_t$ and proceed as in the non-attention seq2seq model

$$[a_t; s_t] \in \mathbb{R}^{2h}$$

Abigail See

# Attention is great

- Attention significantly improves NMT performance
  - It's very useful to allow decoder to focus on certain parts of the source
- Attention solves the bottleneck problem
  - Attention allows decoder to look directly at source; bypass bottleneck
- Attention helps with vanishing gradient problem
  - Provides shortcut to faraway states
- Attention provides some interpretability
  - By inspecting attention distribution, we can see what the decoder was focusing on
  - We get (soft) alignment for free!
  - This is cool because we never explicitly trained an alignment system
  - The network just learned alignment by itself

# Attention is a *general* Deep Learning technique

- We've seen that attention is a great way to improve the sequence-to-sequence model for Machine Translation.

- <u>However</u>: You can use attention in many architectures  (not just seq2seq) and many tasks (not just MT)

- **More general definition of attention**:
  - Given a set of vector *values*, and a vector *query*, **attention** is a  technique to compute a weighted sum of the values,  dependent on the query.

- We sometimes say that the query *attends to* the values.

- For example, in the seq2seq + attention model, each decoder hidden state (query) *attends to* all the encoder hidden states (values).

# Plan for this lecture

- Recurrent neural networks
  - Basics
  - Training (backprop through time, vanishing gradient)
  - Recurrent networks with gates (GRU, LSTM)
- Applications in NLP and vision
  - Image/video captioning
  - Neural machine translation (beam search, attention)
- Transformers
  - Self-attention
  - BERT
  - Cross-modal transformers for VQA and VCR

# Transformers
## (meaning representation through context, representation learning, unsupervised learning)

# How do we represent the meaning of a word?

Definition: **meaning** (Webster dictionary)

- the idea that is represented by a word, phrase, etc.

- the idea that a person wants to express by using words, signs, etc.

- the idea that is expressed in a work of writing, art, etc.

Commonest linguistic way of thinking of meaning:

signifier (symbol) $\iff$ signified (idea or thing)

= denotational semantics

# How do we have usable meaning in a computer?

Common solution: Use e.g. WordNet, a thesaurus containing lists of **synonym sets** and **hypernyms** ("is a" relationships).

*e.g. synonym sets containing "good":*

```
from nltk.corpus import wordnet as wn
poses = { 'n':'noun', 'v':'verb', 's':'adj (s)', 'a':'adj', 'r':'adv'}
for synset in wn.synsets("good"):
    print("{}: {}".format(poses[synset.pos()],
              ", ".join([l.name() for l in synset.lemmas()])))
```

noun: good
noun: good, goodness
noun: good, goodness
noun: commodity, trade_good, good
adj: good
adj (sat): full, good  adj:
good
adj (sat): estimable, good, honorable, respectable  adj (sat):
beneficial, good
adj (sat): good
adj (sat): good, just, upright
…
adverb: well, good
adverb: thoroughly, soundly, good

*e.g. hypernyms of "panda":*

```
from nltk.corpus import wordnet as wn
panda = wn.synset("panda.n.01")  hyper =
lambda s: s.hypernyms()
list(panda.closure(hyper))
```

[Synset('procyonid.n.01'),
Synset('carnivore.n.01'),
Synset('placental.n.01'),
Synset('mammal.n.01'),
Synset('vertebrate.n.01'),
Synset('chordate.n.01'),
Synset('animal.n.01'),
Synset('organism.n.01'),
Synset('living_thing.n.01'),
Synset('whole.n.02'),
Synset('object.n.01'),
Synset('physical_entity.n.01'),
Synset('entity.n.01')]

Christopher Manning

# Problems with resources like WordNet

- Great as a resource but missing nuance
  - e.g. "proficient" is listed as a synonym for "good".
    This is only correct in some contexts.

- Missing new meanings of words
  - e.g., wicked, badass, nifty, wizard, genius, ninja, bombest
  - Impossible to keep up-to-date!

- Subjective

- Requires human labor to create and adapt

- Can't compute accurate word similarity

Christopher Manning

# Representing words as discrete symbols

In traditional NLP, we regard words as discrete symbols:
hotel, conference, motel – a localist representation

Means one 1, the rest 0s

Words can be represented by one-hot vectors:

motel = [0 0 0 0 0 0 0 0 0 0 1 0 0 0]
hotel = [0 0 0 0 0 0 0 1 0 0 0 0 0 0]

Vector dimension = number of words in vocab (e.g. 500,000)

# Problem with words as discrete symbols

**Example:** in web search, if user searches for "Seattle motel", we would like to match documents containing "Seattle hotel".

But:

motel = [0 0 0 0 0 0 0 0 0 0 1 0 0 0 0]
hotel = [0 0 0 0 0 0 0 1 0 0 0 0 0 0 0]

These two vectors are orthogonal.

There is no natural notion of **similarity** for one-hot vectors!

**Solution:**

- Could try to rely on WordNet's list of synonyms to get similarity?
  - But it is well-known to fail badly: incompleteness, etc.
- **Instead: learn to encode similarity in the vectors themselves**

# Representing words by their context

- Distributional semantics: **A word's meaning is given by the words that frequently appear close-by**

  - *"You shall know a word by the company it keeps"* (J. R. Firth 1957)

  - One of the most successful ideas of modern statistical NLP!

- When a word *w* appears in a text, its **context** is the set of words that appear nearby (within a fixed-size window).

- Use the many contexts of *w* to build up a representation of *w*

| | | |
|---|---|---|
| …government debt problems turning into | **banking** | crises as happened in 2009… |
| …saying that Europe needs unified | **banking** | regulation to replace the hodgepodge… |
| …India has just given its | **banking** | system a shot in the arm… |

These context words will represent *banking*

Christopher Manning

# Word vectors

We will build a dense vector for each word, chosen so that it is similar to vectors of words that appear in similar contexts

$$
banking = \begin{pmatrix} 0.286 \\ 0.792 \\ -0.177 \\ -0.107 \\ 0.109 \\ -0.542 \\ 0.349 \\ 0.271 \end{pmatrix}
$$

Note: word vectors are sometimes called word embeddings or word representations. They are a distributed representation.

# Word meaning as a neural word vector - visualization

$$expect = \begin{pmatrix} 0.286 \\ 0.792 \\ -0.177 \\ -0.107 \\ 0.109 \\ -0.542 \\ 0.349 \\ 0.271 \\ 0.487 \end{pmatrix}$$
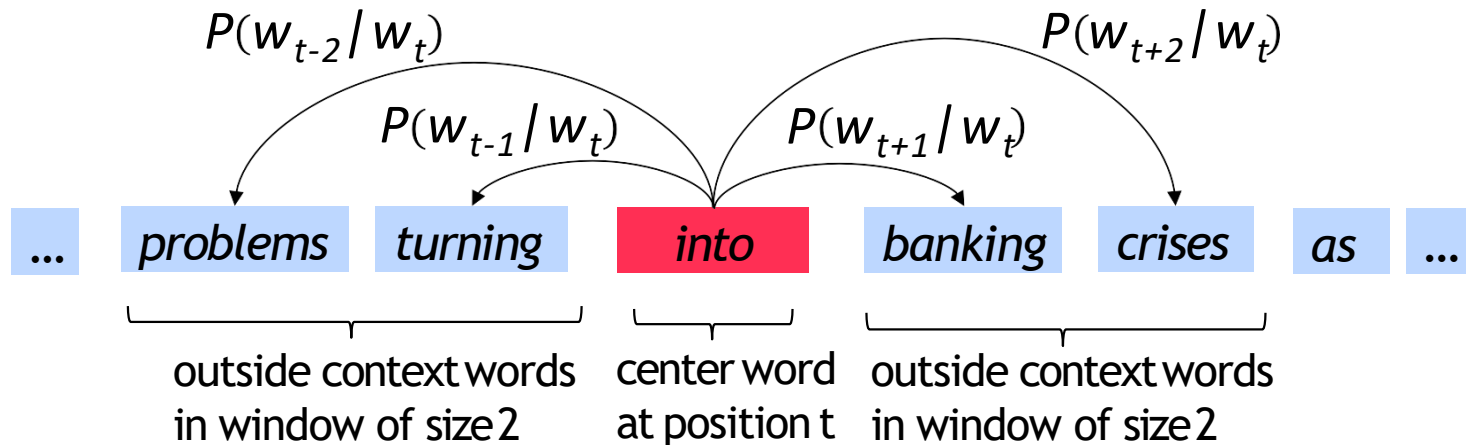
# 3. Word2vec: Overview

Word2vec (Mikolov et al. 2013) is a framework for learning word vectors

Idea:

- We have a large corpus of text
- Every word in a fixed vocabulary is represented by a vector
- Go through each position $t$ in the text, which has a center word $c$ and context ("outside") words $o$
- Use the similarity of the word vectors for $c$ and $o$ to calculate the probability of $o$ given $c$ (or vice versa)
- Keep adjusting the word vectors to maximize this probability

Christopher Manning

# Word2Vec Overview

- Example windows and process for computing $P(w_{t+j}|w_t)$



$P(w_{t-2}|w_t)$  $P(w_{t+2}|w_t)$

$P(w_{t-1}|w_t)$  $P(w_{t+1}|w_t)$

... | *problems* | *turning* | *into* | *banking* | *crises* | *as* | ...

outside context words in window of size 2    center word at position t    outside context words in window of size 2

# Word2Vec Overview

- Example windows and process for computing $P(w_{t+j} | w_t)$

$$P(w_{t-2} | w_t) \qquad\qquad\qquad P(w_{t+2} | w_t)$$

$$P(w_{t-1} | w_t) \qquad\qquad P(w_{t+1} | w_t)$$

| ... | *problems* | *turning* | *into* | *banking* | *crises* | *as* | ... |

outside context words
in window of size 2

center word
at position t

outside context words
in window of size 2

Christopher Manning

# Word2vec: objective function

For each position t = 1, ... , T, predict context words within a window of fixed size $m$, given center word $w_j$.

$$\text{Likelihood} = L(\theta) = \prod_{t=1}^{T} \prod_{\substack{-m \leq j \leq m \\ j \neq 0}} P(w_{t+j} \mid w_t; \theta)$$

$\theta$ is all variables to be optimized

sometimes called *cost* or *loss* function

The objective function is the (average) negative log likelihood:

$$J(\theta) = -\frac{1}{T} \log L(\theta) = -\frac{1}{T} \sum_{t=1}^{T} \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \log P(w_{t+j} \mid w_t; \theta)$$

Minimizing objective function ⟺ Maximizing predictive accuracy

# Word2vec: objective function

- We want to minimize the objective function:

$$J(\theta) = -\frac{1}{T}\log L(\theta) = -\frac{1}{T}\sum_{t=1}^{T}\sum_{\substack{-m \le j \le m \\ j \neq 0}}\log P\left(w_{t+j} \mid w_t; \theta\right)$$

- <u>Question:</u> How to calculate $P(w_{t+j} | w_t ; \theta)$?

- <u>Answer:</u> We will *use two* vectors per word *w*:

  - $v_w$ when *w* is a center word

  - $u_w$ when *w* is a context word

- Then for a center word *c* and a context word *o*:

$$P(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V}\exp(u_w^T v_c)}$$

Christopher Manning

# Word2vec: prediction function

Exponentiation makes anything positive

Dot product compares similarity of $o$ and $c$.
$u^T v = u.v = \sum_{i=1}^{n} u_i v_i$
Larger dot product = larger probability

$$P(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$$

Normalize over entire vocabulary to give probability distribution

- This is an example of the **softmax function** $\mathbb{R}^n \rightarrow \mathbb{R}^n$

$$\text{softmax}(x_i) = \frac{\exp(x_i)}{\sum_{j=1}^{n} \exp(x_j)} = p_i$$

- The softmax function maps arbitrary values to a probability distribution $p_i$
  - "max" because amplifies probability of largest $x_i$
  - "soft" because still assigns some probability to smaller $x_i$
  - Frequently used in Deep Learning

Christopher Manning

# Peters et al. (2018): ELMo: Embeddings from Language Models

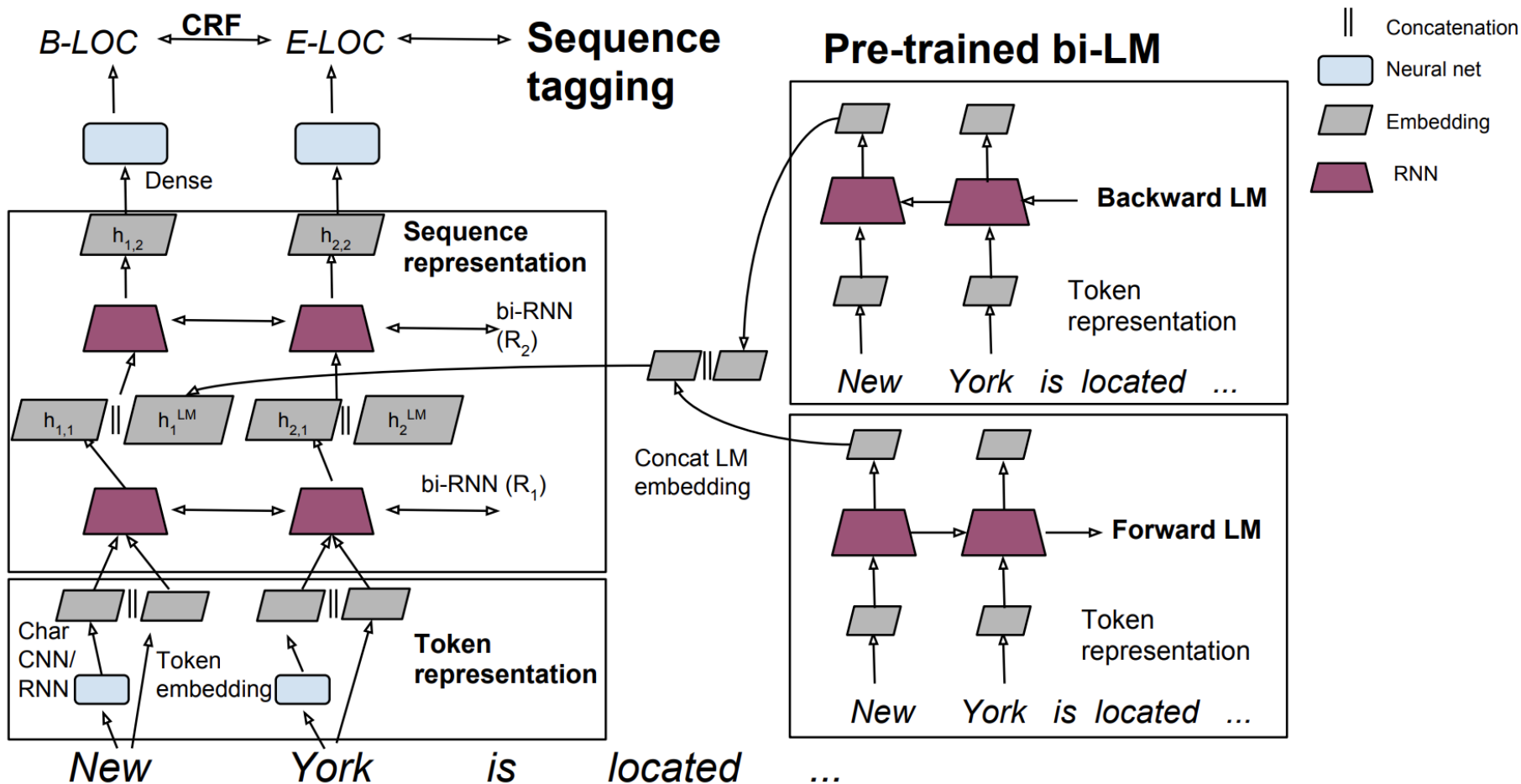Deep contextualized word representations. NAACL 2018. https://arxiv.org/abs/1802.05365

- Breakout version of **word token vectors** or **contextual word vectors**
- Learn word token vectors using long contexts not context windows (here, whole sentence, could be longer)
- Learn a deep Bi-NLM and use all its layers in prediction

# Peters et al. (2018): ELMo: Embeddings from Language Models

- Train a bidirectional LM

- Aim at performant but not overly large LM:

  - Use 2 biLSTM layers

  - Use character CNN to build initial word representation

  - User 4096 dim hidden/cell LSTM states with 512 dim projections to next input

  - Use a residual connection

  - Tie parameters of token input and output (softmax) and tie these between forward and backward LMs

# ELMo used in a sequence tagger



$$\mathbf{h}_{k,1} = [\overrightarrow{\mathbf{h}}_{k,1}; \overleftarrow{\mathbf{h}}_{k,1}; \mathbf{h}_k^{LM}].$$

Christopher Manning, figure from https://tsenghungchen.github.io/posts/elmo/, paper at https://arxiv.org/pdf/1802.05365.pdf

# ELMo results: Great for all tasks

| TASK | PREVIOUS SOTA | | OUR BASELINE | ELMo + BASELINE | INCREASE (ABSOLUTE/ RELATIVE) |
|---|---|---|---|---|---|
| SQuAD | Liu et al. (2017) | 84.4 | 81.1 | 85.8 | 4.7 / 24.9% |
| SNLI | Chen et al. (2017) | 88.6 | 88.0 | $88.7 \pm 0.17$ | 0.7 / 5.8% |
| SRL | He et al. (2017) | 81.7 | 81.4 | 84.6 | 3.2 / 17.2% |
| Coref | Lee et al. (2017) | 67.2 | 67.2 | 70.4 | 3.2 / 9.8% |
| NER | Peters et al. (2017) | $91.93 \pm 0.19$ | 90.15 | $92.22 \pm 0.10$ | 2.06 / 21% |
| SST-5 | McCann et al. (2017) | 53.7 | 51.4 | $54.7 \pm 0.5$ | 3.3 / 6.8% |

Christopher Manning

# ELMo: Weighting of layers

- The two biLSTM NLM layers have differentiated uses/meanings
  - Lower layer is better for lower-level syntax, etc.
    - Part-of-speech tagging, syntactic dependencies, NER
  - Higher layer is better for higher-level semantics
    - Sentiment, Semantic role labeling, question answering, SNLI

# Let's scale it up!



ULMfit
Jan 2018
Training:  1
GPU day

GPT
June 2018
Training
240 GPU days

BERT
Oct 2018
Training
256 TPU days
~320–560
GPU days

GPT-2
Feb 2019
Training
~2048 TPU v3
days according to
a reddit thread

# GPT-2 language model cherry-picked output

**SYSTEM PROMPT (HUMAN-WRITTEN)**

*In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains. Even more surprising to the researchers was the fact that the unicorns spoke perfect English.*

**MODEL COMPLETION (MACHINE-WRITTEN, 10 TRIES)**

The scientist named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown to science.

Now, after almost two centuries, the mystery of what sparked this odd phenomenon is finally solved.

Dr. Jorge Pérez, an evolutionary biologist from the University of La Paz, and several companions, were exploring the Andes Mountains when they found a small valley, with no other animals or humans. Pérez noticed that the valley had what appeared to be a natural fountain, surrounded by two peaks of rock and silver snow.

Pérez and the others then ventured further into the valley. …

# Elon Musk's OpenAI builds artificial intelligence so powerful it must be kept locked up for the good of humanity

**Jasper Hamill** Friday 15 Feb 2019 10:06 am

f          y          ●          ❮          272 SHARES

Elon Musk's scientists have announced the creation of a terrifying artificial intelligence that's so smart they refused to release it to the public.

OpenAI's GPT-2 is designed to write just like a human and is an impressive leap forward capable of penning chillingly convincing text.

It was 'trained' by analysing eight million web pages and is capable of writing large tracts based upon a 'prompt' written by a real person.

But the machine mind will not be released in its fully-fledged form because of the risk of it being used for 'malicious purposes' such as generating fake news, impersonating people online, automating the production of spam or churning out 'abusive or faked content to post on social media'.

OpenAI wrote: 'Due to our concerns about malicious applications of the technology, we are not releasing the trained model.

**Elon Musk** ✔
@elonmusk                                    Follow    ⌄

Replying to @georgezachary

To clarify, I've not been involved closely with OpenAI for over a year & don't have mgmt or board oversight

8:19 PM - 16 Feb 2019

**500** Retweets  **14,573** Likes     ●●●●●●●●●

💬 229        ⟲ 500        ♡ 15K        ✉

Christopher Manning

# The Motivation for Transformers

- We want **parallelization** but RNNs are inherently sequential

- Despite GRUs and LSTMs, RNNs still need attention mechanism  to deal with long range dependencies – **path length** between  states grows with sequence otherwise

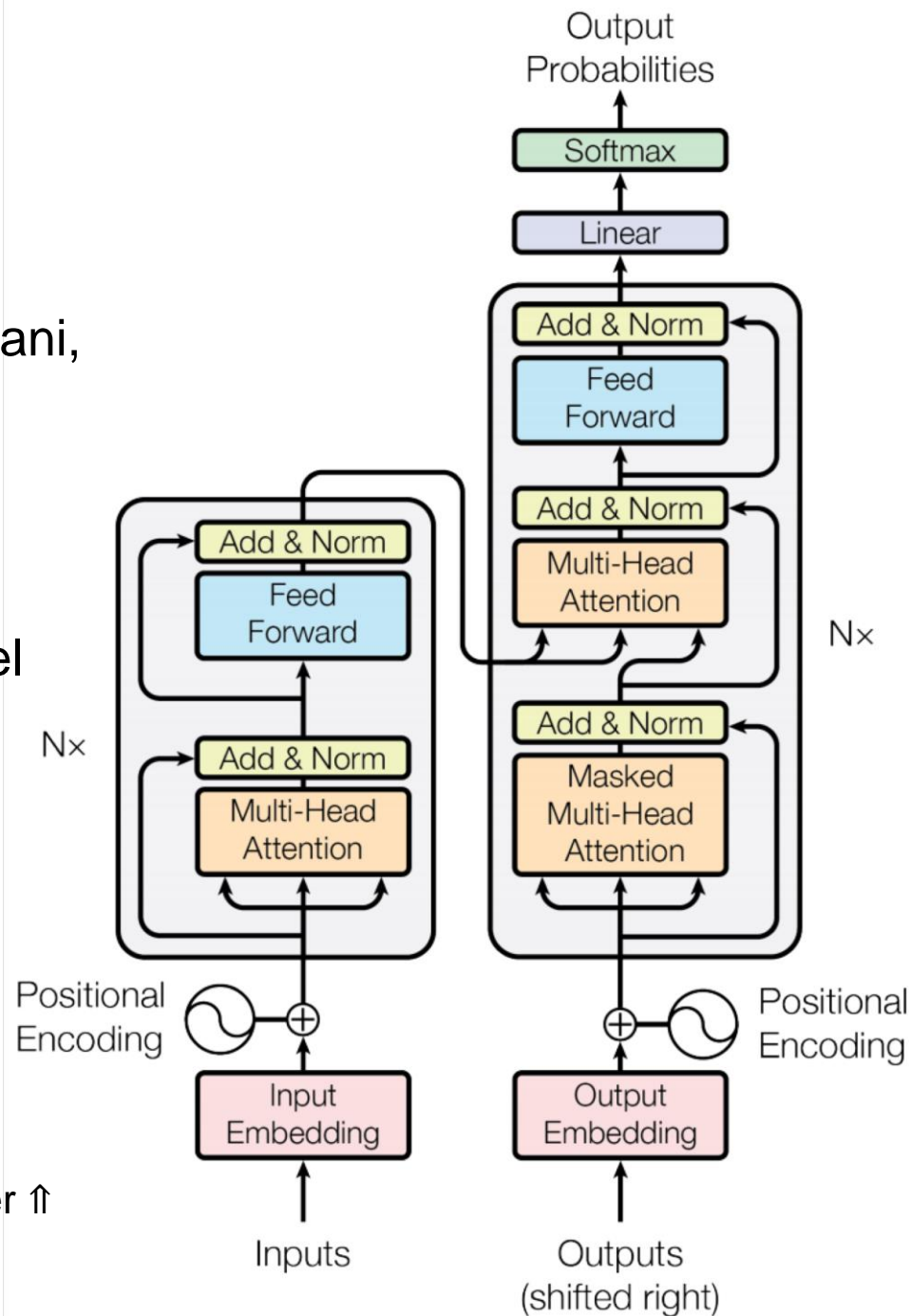- But if **attention** gives us access to any state… maybe we can just  use attention and don't need the RNN?

# Transformer Overview

Attention is all you need. 2017.  Aswani, Shazeer, Parmar, Uszkoreit,  Jones, Gomez, Kaiser, Polosukhin
https://arxiv.org/pdf/1706.03762.pdf

- Non-recurrent sequence-to-sequence encoder-decoder model

- Task: machine translation with parallel corpus

- Predict each translated word

- Final cost/error function is standard cross-entropy error on top of a softmax classifier

This and related figures from paper ⇑

Christopher Manning

# Dot-Product Attention (Extending our previous def.)

- Inputs: a query q and a set of key-value (k-v) pairs to an output
- Query, keys, values, and output are all vectors

- Output is weighted sum of values, where
- Weight of each value is computed by an inner product of query and corresponding key
- Queries and keys have same dimensionality $d_k$ value have $d_v$

$$A(q, K, V) = \sum_i \frac{e^{q \cdot k_i}}{\sum_j e^{q \cdot k_j}} v_i$$

Christopher Manning

# Dot-Product Attention - Matrix notation

- When we have multiple queries q, we stack them in a matrix Q:

$$A(q, K, V) = \sum_i \frac{e^{q \cdot k_i}}{\sum_j e^{q \cdot k_j}} v_i$$

- Becomes: $A(Q, K, V) = softmax(QK^T)V$

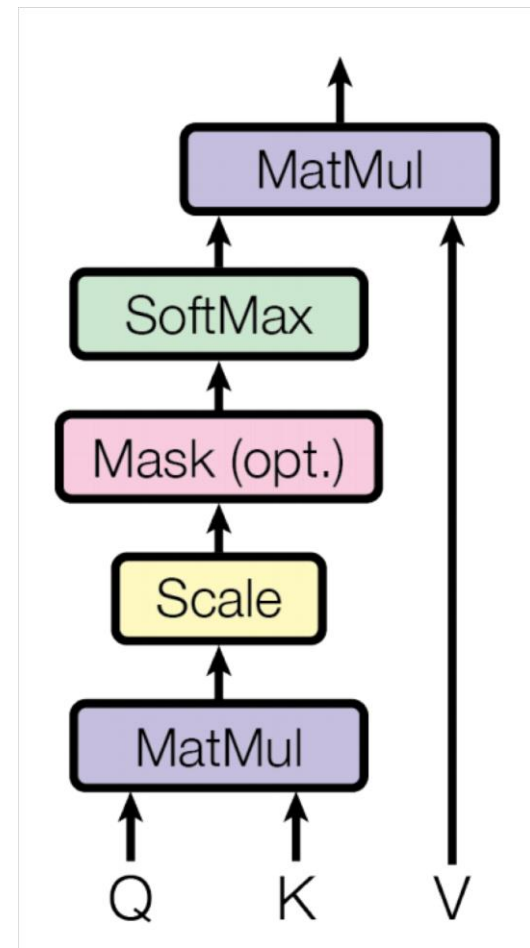$[|Q| \times d_k] \times [d_k \times |K|] \times [|K| \times d_v]$

softmax
row-wise       $= [|Q| \times d_v]$

# Scaled Dot-Product Attention

- Problem: As $d_k$ gets large, the variance of $QK^T$ increases → some values inside the softmax get large → the softmax gets very peaked → hence its gradient gets smaller.

- Solution: Scale by length of query/key vectors:

$$A(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

- The input word vectors are the queries, keys and values
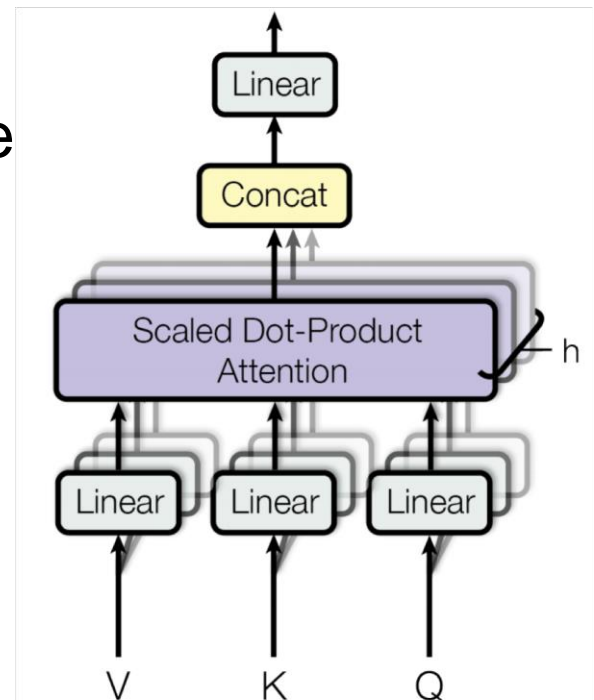- In other words: the word vectors select each other



MatMul

SoftMax

Mask (opt.)

Scale

MatMul

Q      K      V

# Multi-head attention

- Problem with simple self-attention:
- Only one way for words to interact with one-another
- Solution: Multi-head attention
- First map Q, K, V into h=8 many lower dimensional spaces via W matrices
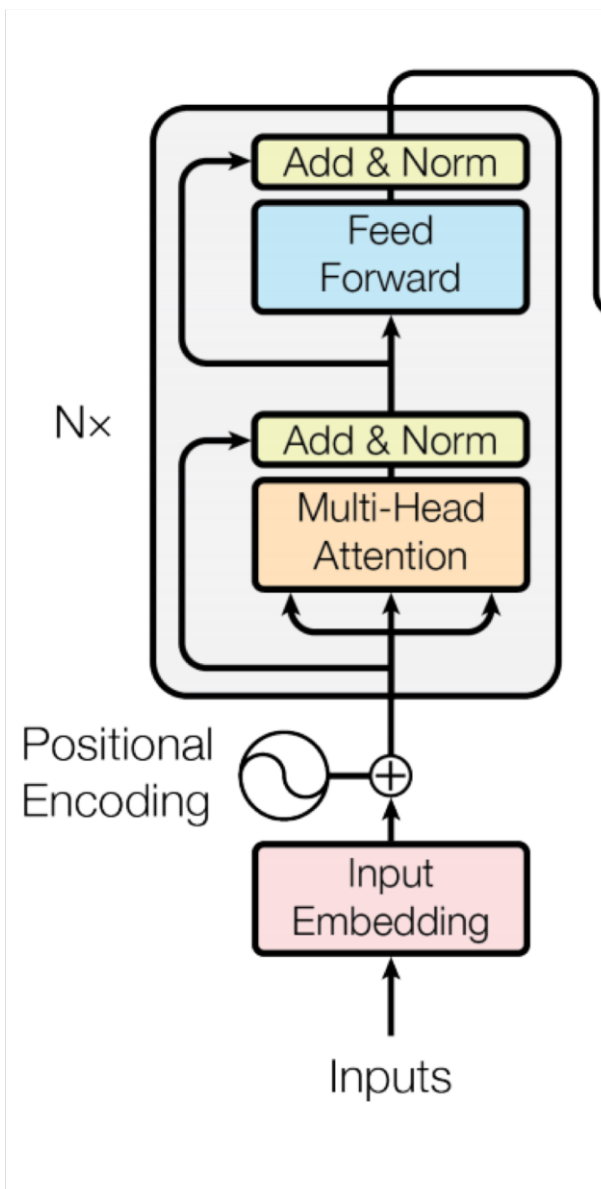- Then apply attention, then concatenate outputs and pipe through linear layer



$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, ..., \text{head}_h)W^O$$
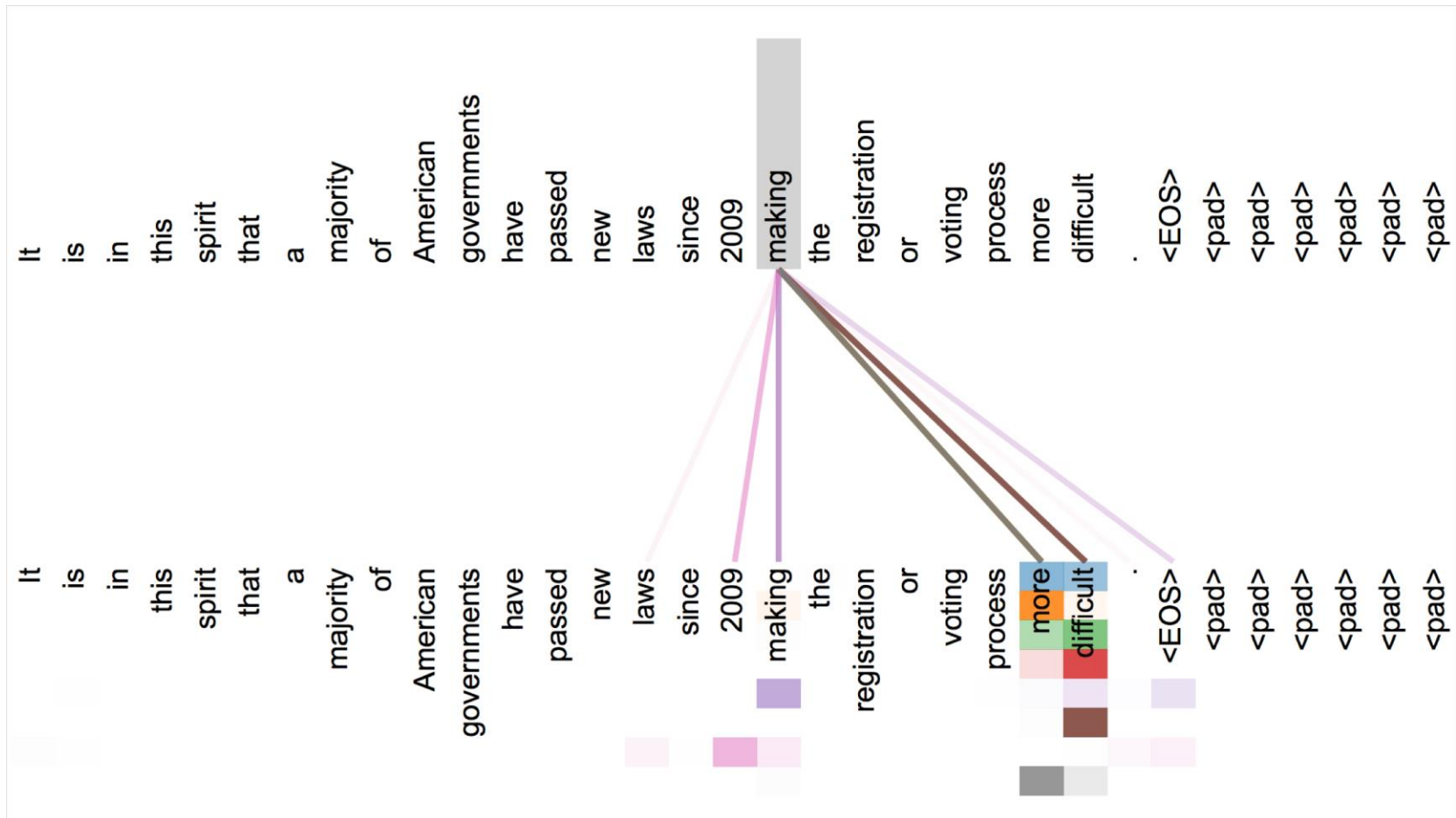$$\text{where head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

# Complete Encoder

- For encoder, at each block, we use the same Q, K and V from the previous layer
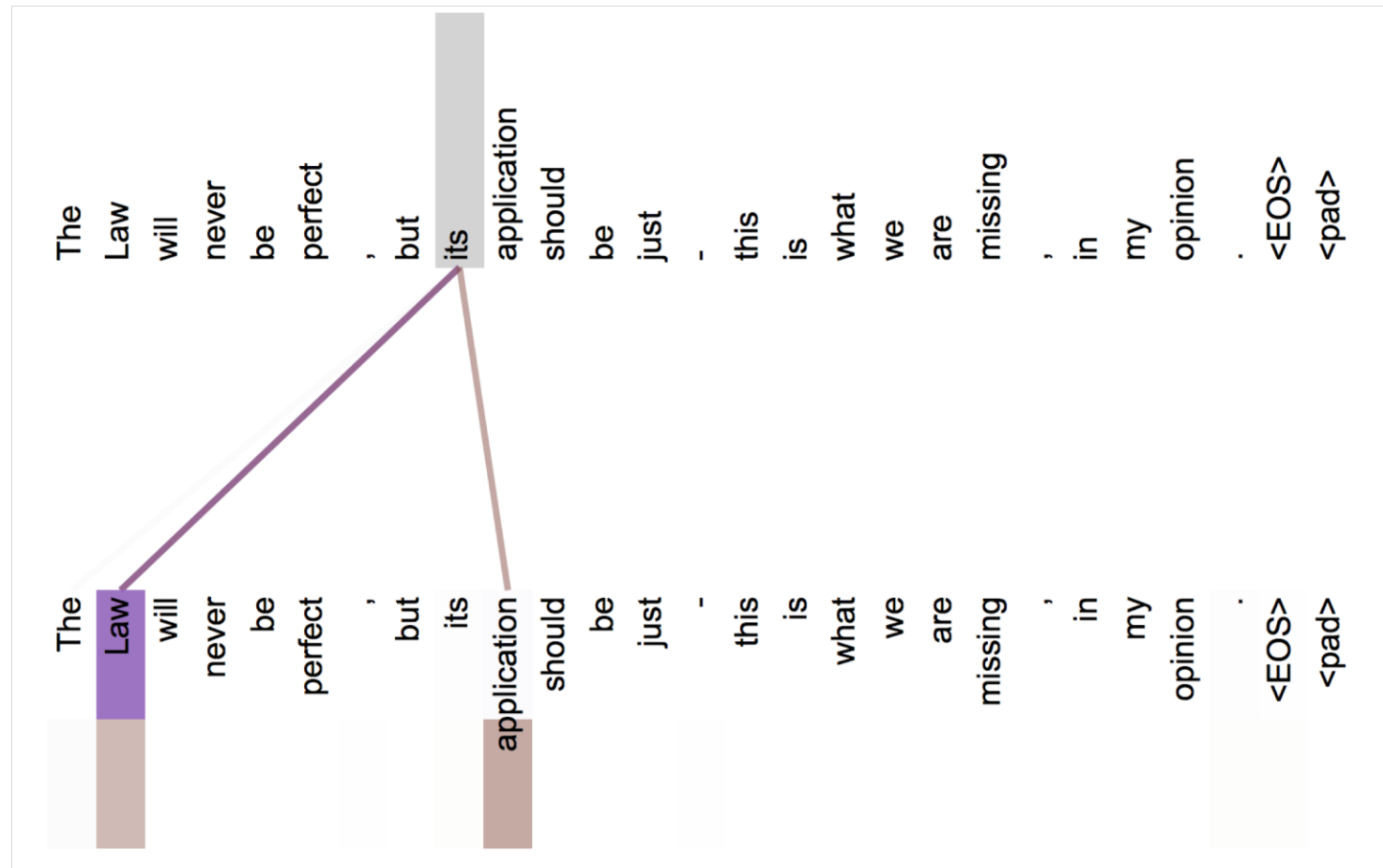
- Blocks are repeated 6 times (in vertical stack)

# Attention visualization in layer 5

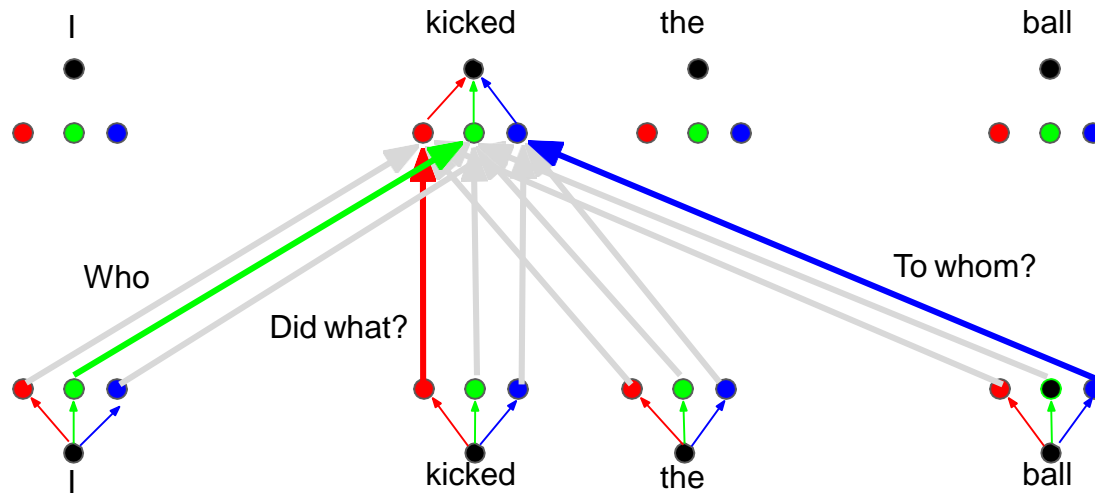- Words start to pay attention to other words in sensible ways

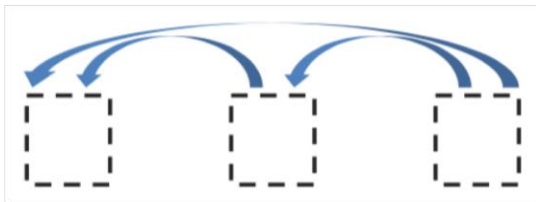# Attention visualization: Implicit anaphora resolution



In 5[th] layer. Isolated attentions from just the word 'its' for attention heads 5 and 6.
Note that the attentions are very sharp for this word.
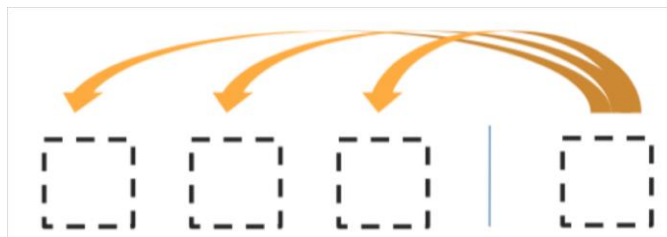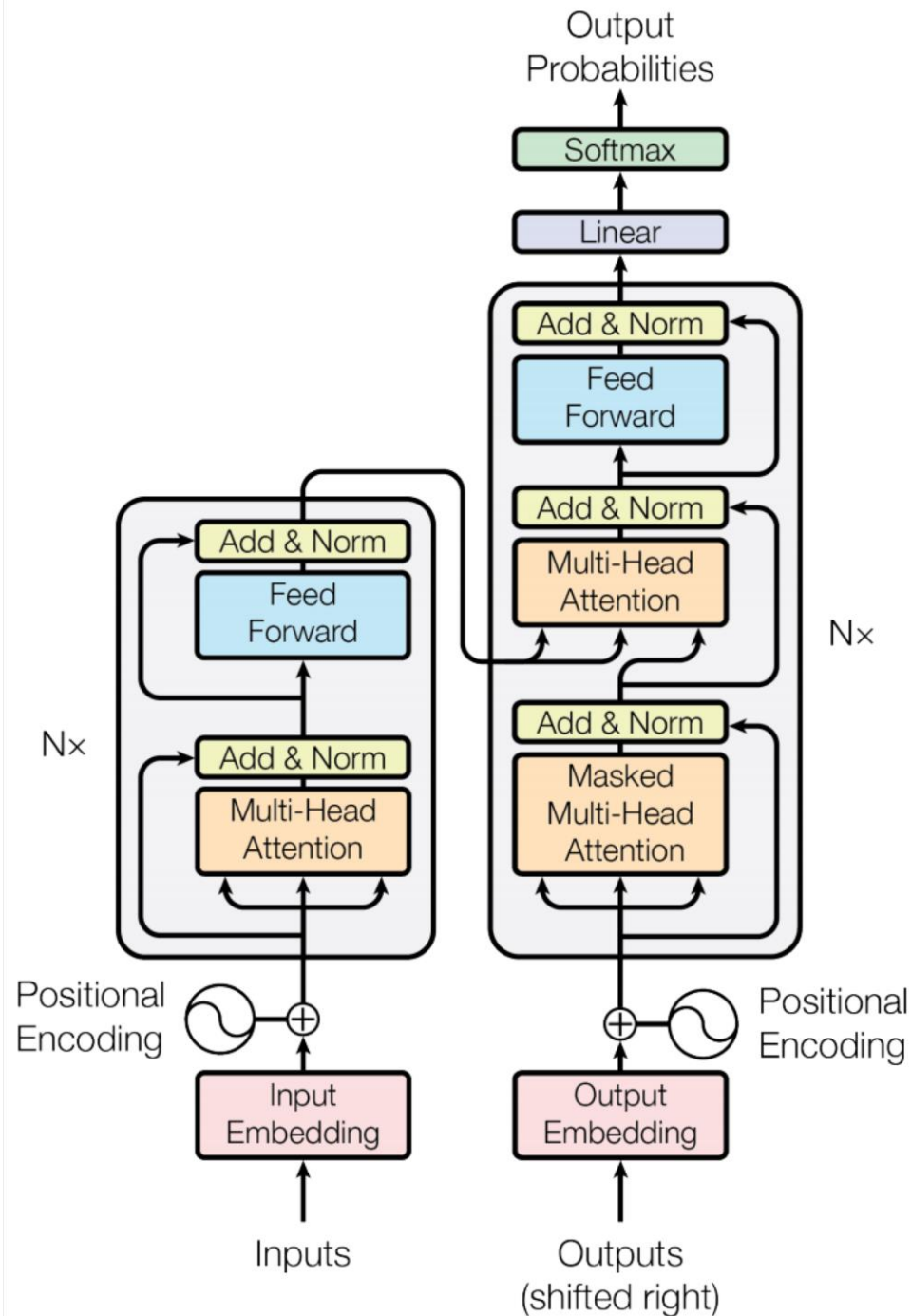
# Parallel attention heads

# Transformer Decoder

- 2 sublayer changes in decoder
- Masked decoder self-attention on previously generated outputs:



- Encoder-Decoder Attention, where queries come from previous decoder layer and keys and values come from output of encoder



Blocks repeated 6 times also

# BERT: Devlin, Chang, Lee, Toutanova (2018)

BERT (Bidirectional Encoder Representations from Transformers):

Pre-training of Deep Bidirectional Transformers for Language  Understanding

Based on slides from Jacob Devlin

Christopher Manning

# BERT: Devlin, Chang, Lee, Toutanova (2018)

- Mask out $k$% of the input words, and then predict the masked words
  - They always use $k$ = 15%

<div align="center">
store          gallon

↑          ↑

the man went to the [MASK] to buy a [MASK] of milk
</div>

- Too little masking: Too expensive to train
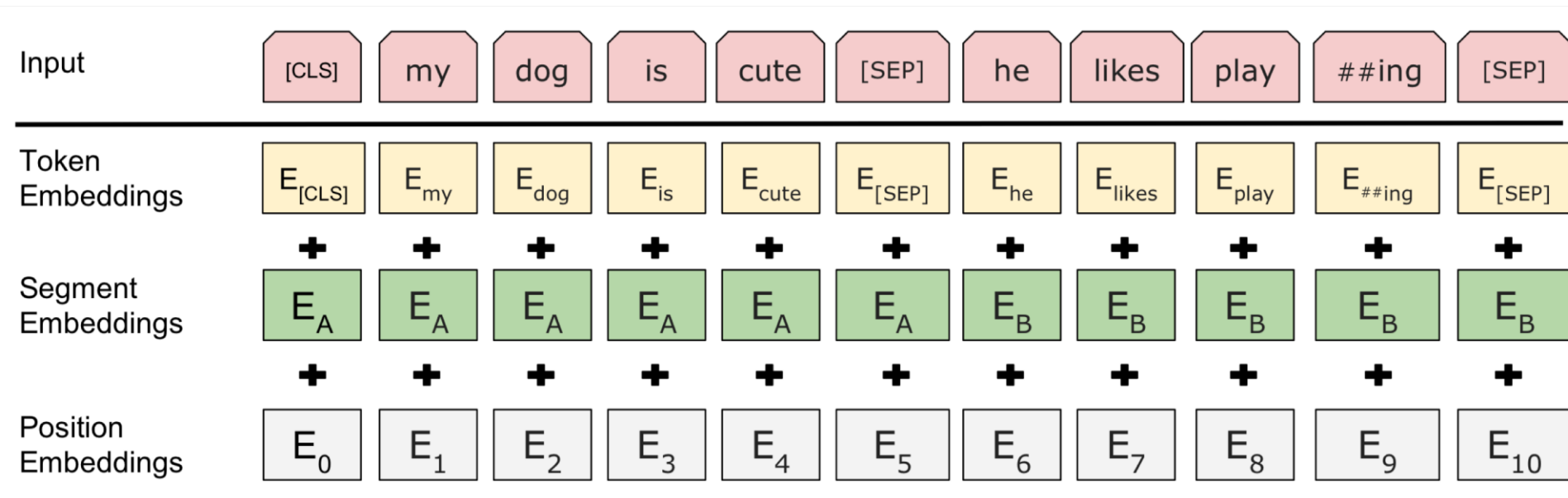- Too much masking: Not enough context

# Additional task: Next sentence prediction

- To learn *relationships* between sentences, predict whether  Sentence B is actual sentence that proceeds Sentence A, or a  random sentence

**Sentence A** = The man went to the store.
**Sentence B** = He bought a gallon of milk.
**Label** = IsNextSentence

**Sentence A** = The man went to the store.
**Sentence B** = Penguins are flightless.
**Label** = NotNextSentence

# BERT sentence pair encoding

| Input | [CLS] | my | dog | is | cute | [SEP] | he | likes | play | ##ing | [SEP] |
|-------|-------|-----|------|-----|------|-------|-----|-------|------|-------|-------|
| Token Embeddings | $E_{[CLS]}$ | $E_{my}$ | $E_{dog}$ | $E_{is}$ | $E_{cute}$ | $E_{[SEP]}$ | $E_{he}$ | $E_{likes}$ | $E_{play}$ | $E_{\#\#ing}$ | $E_{[SEP]}$ |
| | + | + | + | + | + | + | + | + | + | + | + |
| Segment Embeddings | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ |
| | + | + | + | + | + | + | + | + | + | + | + |
| Position Embeddings | $E_0$ | $E_1$ | $E_2$ | $E_3$ | $E_4$ | $E_5$ | $E_6$ | $E_7$ | $E_8$ | $E_9$ | $E_{10}$ |

Token embeddings are word pieces
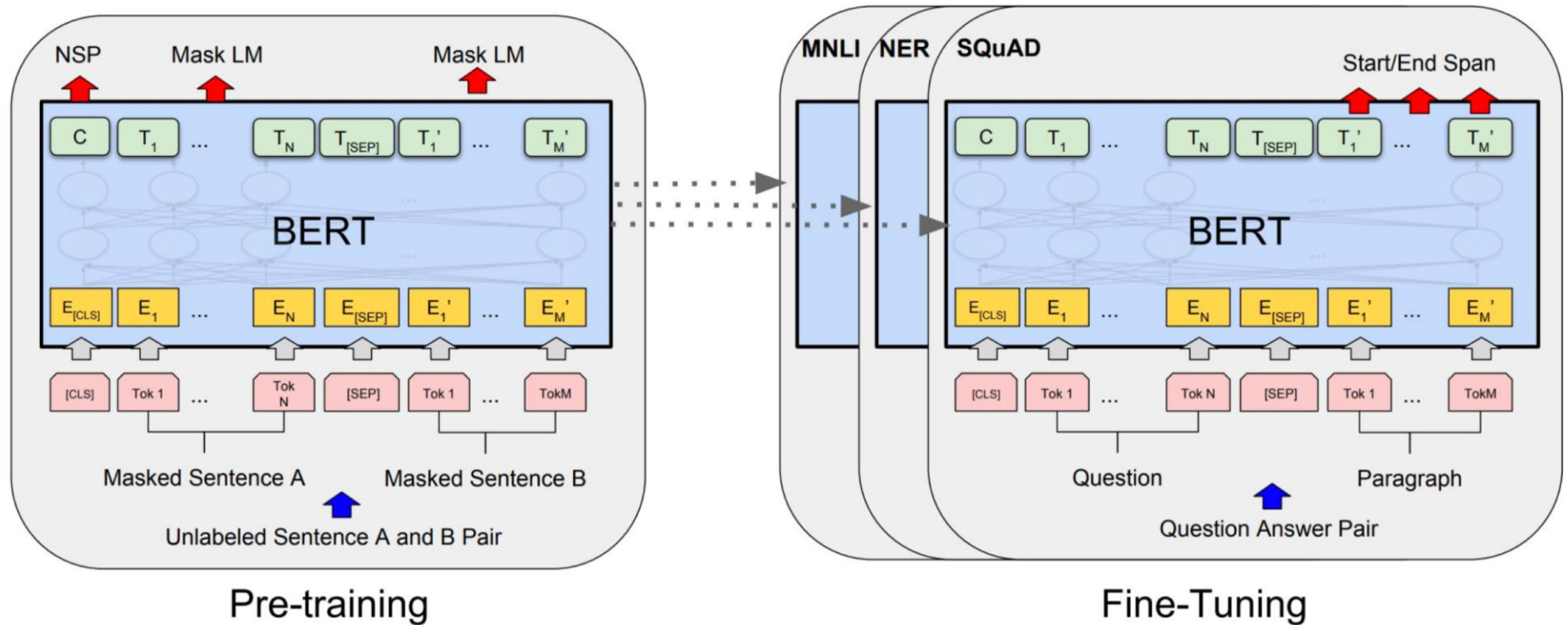Learned segmented embedding represents each sentence
Positional embedding

# BERT model architecture and training

- Transformer encoder (as before)
- Self-attention ⇒ no locality bias
  - Long-distance context has "equal opportunity"
- Single multiplication per layer ⇒ efficiency on GPU/TPU

- Train on Wikipedia + BookCorpus
- Train 2 model sizes:
  - BERT-Base: 12-layer, 768-hidden, 12-head
  - BERT-Large: 24-layer, 1024-hidden, 16-head
- Trained on 4x4 or 8x8 TPU slice for 4 days

Christopher Manning

# BERT model fine tuning

- Simply learn a classifier built on the top layer for each task that you fine tune for

# SQuAD 2.0 leaderboard, 2019-02-07

| Rank | Model | EM | F1 |
|------|-------|-----|-----|
| | Human Performance<br>*Stanford University*<br>(Rajpurkar & Jia et al. '18) | **86.831** | **89.452** |
| 1<br>Jan 15, 2019 | BERT + MMFT + ADA (ensemble)<br>*Microsoft Research Asia* | **85.082** | **87.615** |
| 2<br>Jan 10, 2019 | BERT + Synthetic Self-Training (ensemble)<br>*Google AI Language*<br>https://github.com/google-research/bert | 84.292 | 86.967 |
| 3<br>Dec 13, 2018 | BERT finetune baseline (ensemble)<br>*Anonymous* | 83.536 | 86.096 |
| 4<br>Dec 16, 2018 | Lunet + Verifier + BERT (ensemble)<br>*Layer 6 AI NLP Team* | 83.469 | 86.043 |
| 4<br>Dec 21, 2018 | PAML+BERT (ensemble model)<br>*PINGAN GammaLab* | 83.457 | 86.122 |
| 5<br>Dec 15, 2018 | Lunet + Verifier + BERT (single model)<br>*Layer 6 AI NLP Team* | 82.995 | 86.035 |

Christopher Manning
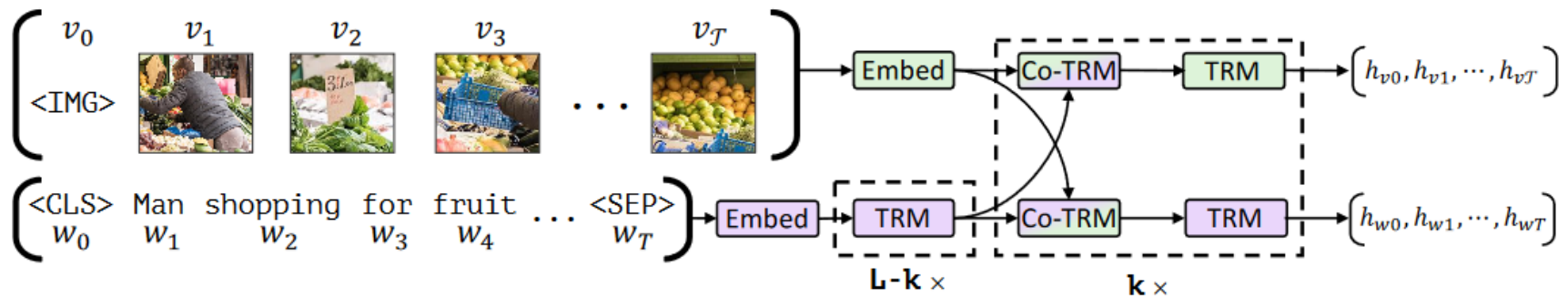
# Cross-modal transformers



Figure 1: Our ViLBERT model consists of two parallel streams for visual (green) and linguistic (purple) processing that interact through novel co-attentional transformer layers. This structure allows for variable depths for each modality and enables sparse interaction through co-attention. Dashed boxes with multiplier subscripts denote repeated blocks of layers.

Lu et al., "ViLBERT: Pretraining Task-Agnostic Visiolinguistic Representations for Vision-and-Language Tasks", NeurIPS 2019
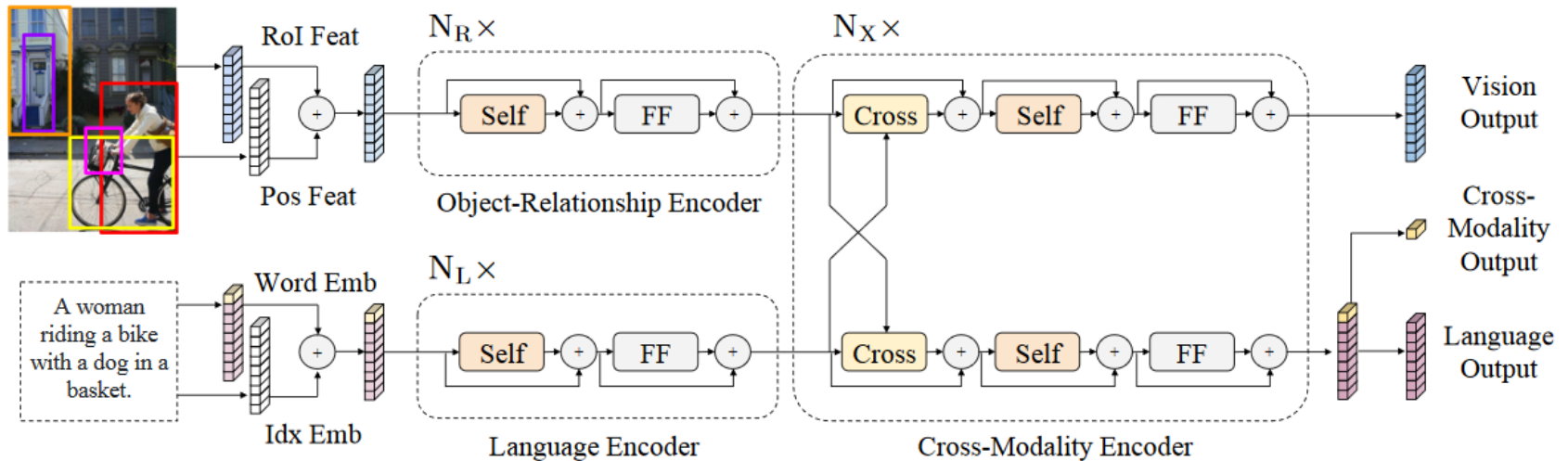
# Cross-modal transformers



Figure 1: The LXMERT model for learning vision-and-language cross-modality representations. 'Self' and 'Cross' are abbreviations for self-attention sub-layers and cross-attention sub-layers, respectively. 'FF' denotes a feed-forward sub-layer.

Tan and Bansal, "LXMERT: Learning Cross-Modality Encoder Representationsfrom Transformers", EMNLP 2019
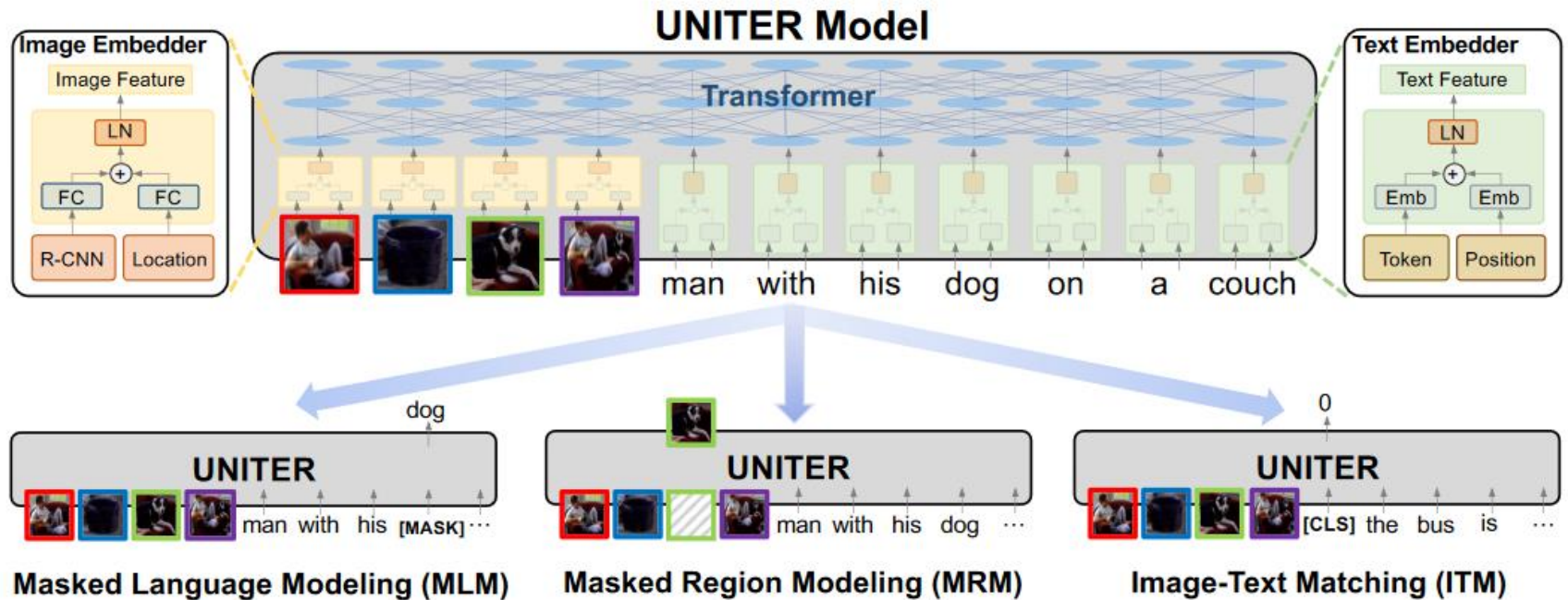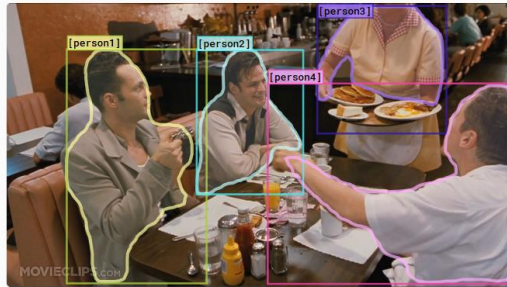
# Cross-modal transformers



Figure 1: Overview of the proposed UNITER model (best viewed in color), consisting of an Image Embedder, a Text Embedder and a multi-layer self-attention Transformer, learned through three pre-training tasks.

Chen et al., "UNITER: Learning UNiversal Image-TExt Representations", arxiv 2019

# Visual Commonsense Reasoning Leaderboard



| Rank | Model | Q->A | QA->R | Q->AR |
|------|-------|------|-------|-------|
| | **Human Performance** <br> *University of Washington* <br><br> (Zellers et al. '18) | 91.0 | 93.0 | 85.0 |
| 📷 <br> September 30, 2019 | **UNITER-large (ensemble)** <br> *MS D365 AI* <br><br> https://arxiv.org /abs/1909.11740 | **79.8** | **83.4** | **66.8** |
| 2 <br> September 23, 2019 | UNITER-large (single model) <br> *MS D365 AI* <br><br> https://arxiv.org /abs/1909.11740 | 77.3 | 80.8 | 62.8 |
| 3 <br> August 9,2019 | ViLBERT (ensemble of 10 models) <br> *Georgia Tech & Facebook AI Research* <br><br> https://arxiv.org /abs/1908.02265 | 76.4 | 78.0 | 59.8 |
| 4 <br> September 23,2019 | VL-BERT (single model) <br> *MSRA & USTC* <br><br> https://arxiv.org /abs/1908.08530 | 75.8 | 78.4 | 59.7 |
| 5 <br> August 9,2019 | ViLBERT (ensemble of 5 models) <br> *Georgia Tech & Facebook AI Research* <br><br> https://arxiv.org /abs/1908.02265 | 75.7 | 77.5 | 58.8 |

https://visualcommonsense.com/leaderboard/

# Additional resource

- Learning about transformers on your own?
  - Key recommended resource:
    - http://nlp.seas.harvard.edu/2018/04/03/attention.html
    - The Annotated Transformer by Sasha Rush
  - An Jupyter Notebook using PyTorch that explains everything!

Christopher Manning

# Recap

- Language modeling is an effective form of unsupervised pretraining for many different supervised tasks
- Attention captures relationships effectively, helps with vanishing gradients
- Attention is cheap to compute and allows better parallelization during training
- Language/sequence models can be extended to settings beyond NLP
- *You will know the meaning of a concept/word/image by the company it keeps*