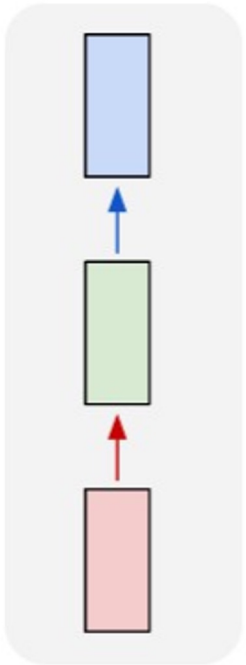# Sequence Modeling (RNNs)

# So far: "Feedforward" Neural Networks
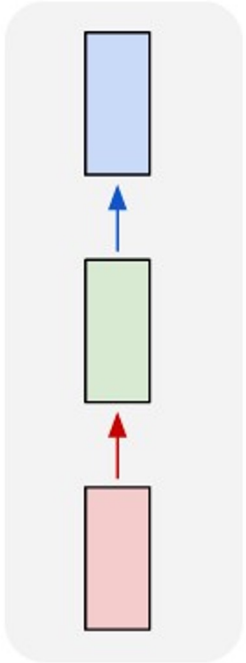
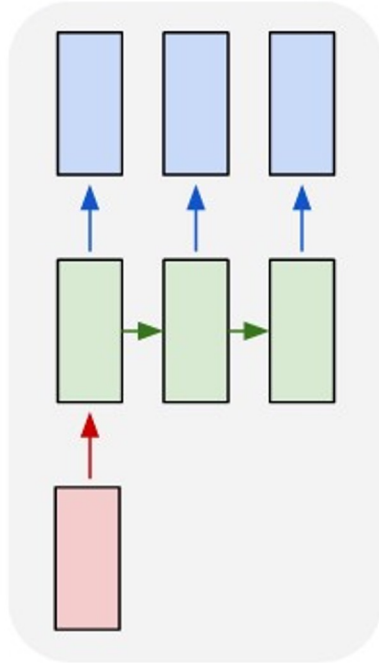one to one



e.g. **Image classification**
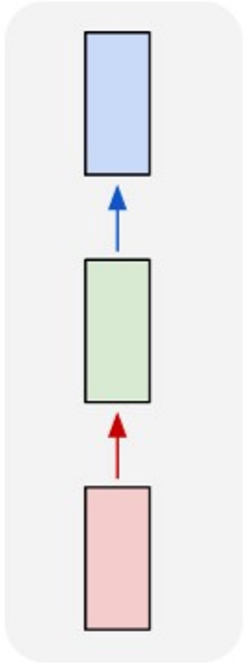Image -> Label

# Recurrent Neural Networks: Process Sequences



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

# Recurrent Neural Networks: Process Sequences

one to one

one to many

many to one

e.g. **Video classification:**
Sequence of images -> label
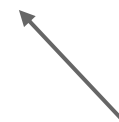
# Recurrent Neural Networks: Process Sequences



e.g. **Machine Translation:**
Sequence of words -> Sequence of words

# Recurrent Neural Networks: Process Sequences



e.g. **Per-frame video classification:**
Sequence of images -> Sequence of labels

# Sequential Processing of Non-Sequential Data

Classify images by taking
a series of "glimpses"



Ba, Mnih, and Kavukcuoglu, "Multiple Object Recognition with Visual Attention", ICLR 2015.
Gregor et al, "DRAW: A Recurrent Neural Network For Image Generation", ICML 2015

Slide from Justin Johnson

# Sequential Processing of Non-Sequential Data

## Generate images one piece at a time!



Gregor et al, "DRAW: A Recurrent Neural Network For Image Generation", ICML 2015

Slide from Justin Johnson

# Recurrent Neural Networks



Key idea: RNNs have an "internal state" that is updated as a sequence is processed

# Recurrent Neural Networks



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    old state    input vector at some time step

some function with parameters W

# Recurrent Neural Networks

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    old state    input vector at
with parameters W    some time step

y

RNN

x

# (Vanilla) Recurrent Neural Networks

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

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



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

$$y_t = W_{hy} h_t + b_y$$

Sometimes called a "Vanilla RNN" or an "Elman RNN" after Prof. Jeffrey Elman

13

# RNN Computational Graph

Initial hidden state
Either set to all 0,
Or learn it

h₀

x₁

# RNN Computational Graph

# RNN Computational Graph

16

# RNN Computational Graph

# RNN Computational Graph

Re-use the same weight matrix at every time-step

# RNN Computational Graph (Many to Many)

# RNN Computational Graph (Many to Many)

# RNN Computational Graph (Many to Many)

# RNN Computational Graph (Many to One)

# RNN Computational Graph (One to Many)

23

# Sequence to Sequence (seq2seq)
# (Many to one) + (One to many)

**Many to one**: Encode input
sequence in a single vector



Sutskever et al, "Sequence to Sequence Learning with Neural Networks", NIPS 2014

Slide from Justin Johnson

24

# Sequence to Sequence (seq2seq)
# (Many to one) + (One to many)

**One to many**: Produce output sequence from single input vector

**Many to one**: Encode input sequence in a single vector



Sutskever et al, "Sequence to Sequence Learning with Neural Networks", NIPS 2014

Slide from Justin Johnson

25

# Example: Language Modeling

Given characters 1, 2, ..., t-1,
model predicts character t

Training sequence: "hello"

Vocabulary: [h, e, l, o]

# Example: Language Modeling

Given characters 1, 2, ..., t-1,
model predicts character t

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

Training sequence: "hello"

Vocabulary: [h, e, l, o]

# Example: Language Modeling

Given characters 1, 2, ..., t-1, model predicts character t

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

Training sequence: "hello"

Vocabulary: [h, e, l, o]

# Example: Language Modeling

Given characters 1, 2, ..., t-1, model predicts character t

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

Training sequence: "hello"

Vocabulary: [h, e, l, o]



Slide from Justin Johnson

# Example: Language Modeling

Given "he", predict "l"

Given characters 1, 2, …, t-1, model predicts character t

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

Training sequence: "hello"

Vocabulary: [h, e, l, o]

30

# Example: Language Modeling

Given "hel", predict "l"

Given characters 1, 2, ..., t-1, model predicts character t

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

Training sequence: "hello"

Vocabulary: [h, e, l, o]

# Example: Language Modeling

Given "hell", predict "o"

Given characters 1, 2, ..., t-1, model predicts character t

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

Training sequence: "hello"

Vocabulary: [h, e, l, o]

# Example: Language Modeling

At test-time, **generate** new text: sample characters one at a time, feed back to model

Training sequence: "hello"

Vocabulary: [h, e, l, o]

# Example: Language Modeling

At test-time, **generate** new text: sample characters one at a time, feed back to model

Training sequence: "hello"

Vocabulary: [h, e, l, o]



Sample

Softmax

| |
|---|
| .03 |
| .13 |
| .00 |
| .84 |

output layer

| |
|---|
| 1.0 |
| 2.2 |
| -3.0 |
| 4.1 |

hidden layer

| |
|---|
| 0.3 |
| -0.1 |
| 0.9 |

input layer

| | |
|---|---|
| 1 | 0 |
| 0 | 1 |
| 0 | 0 |
| 0 | 0 |

input chars:   "h"      "e"

34

# Example: Language Modeling

At test-time, **generate** new
text: sample characters one
at a time, feed back to model

Training sequence: "hello"

Vocabulary: [h, e, l, o]

35

# Example: Language Modeling

At test-time, **generate** new text: sample characters one at a time, feed back to model

Training sequence: "hello"

Vocabulary: [h, e, l, o]

36

# Example: Language Modeling

So far: encode inputs as **one-hot-vector**

$$[w_{11}\ w_{12}\ w_{13}\ w_{14}]\ [1]\qquad [w_{11}]$$
$$[w_{21}\ w_{22}\ w_{23}\ w_{14}]\ [0]\ =\ [w_{21}]$$
$$[w_{31}\ w_{32}\ w_{33}\ w_{14}]\ [0]\qquad [w_{31}]$$
$$[0]$$

Matrix multiply with a one-hot vector just extracts a column from the weight matrix. Often extract this into a separate **embedding** layer

# Example: Language Modeling

So far: encode inputs as **one-hot-vector**

$[w_{11}\ w_{12}\ w_{13}\ w_{14}]\ [1]\qquad [w_{11}]$
$[w_{21}\ w_{22}\ w_{23}\ w_{14}]\ [0]\ =\ [w_{21}]$
$[w_{31}\ w_{32}\ w_{33}\ w_{14}]\ [0]\qquad [w_{31}]$
$\qquad\qquad\qquad\qquad\quad [0]$

Matrix multiply with a one-hot vector just extracts a column from the weight matrix. Often extract this into a separate **embedding** layer

# Backpropagation Through Time

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



Slide from Justin Johnson

39

# Backpropagation Through Time

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

Problem: Takes a lot of memory for long sequences!



Slide from Justin Johnson

40

# Truncated Backpropagation Through Time

Loss

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

41

# Truncated Backpropagation Through Time



Loss

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

# Truncated Backpropagation Through Time

43

# Example: Image Captioning



Figure from Karpathy et a, "Deep Visual-Semantic Alignments for Generating Image Descriptions", CVPR 2015

Mao et al, "Explain Images with Multimodal Recurrent Neural Networks", NeurIPS 2014 Deep Learning and Representation Workshop
Karpathy and Fei-Fei, "Deep Visual-Semantic Alignments for Generating Image Descriptions", CVPR 2015
Vinyals et al, "Show and Tell: A Neural Image Caption Generator", CVPR 2015
Donahue et al, "Long-term Recurrent Convolutional Networks for Visual Recognition and Description", CVPR 2015
Chen and Zitnick, "Learning a Recurrent Visual Representation for Image Caption Generation", CVPR 2015

Slide from Justin Johnson

# Example: Image Captioning



**Recurrent Neural Network**

**Convolutional Neural Network**

Slide from Justin Johnson

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

**Transfer learning**: Take CNN trained on ImageNet, chop off last layer

Slide from Justin Johnson

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>

Slide from Justin Johnson

68

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

**Before:**

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

**Now:**

$$\tanh(\textcolor{red}{W_{hh}h_{t-1}} + \textcolor{blue}{W_{xh}x_t} + \textcolor{magenta}{W_{ih}v} + b_h)$$

$W_{ih}$

y0

h0

x0

<START>

Slide from Justin Johnson

69

**Before:**

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

**Now:**

$$\tanh(\boldsymbol{W_{hh}h_{t-1}} + \boldsymbol{W_{xh}x_t} + \boldsymbol{W_{ih}v} + b_h)$$

Sample word and copy to input

Slide from Justin Johnson

Before:

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

Now:

$$\tanh(W_{hh}h_{t-1} + W_{xh}x_t + W_{ih}v + b_h)$$

$W_{ih}$

Sample word and copy to input

Slide from Justin Johnson

72

**Before:**

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

**Now:**

$$\tanh(\boldsymbol{W_{hh}h_{t-1}} + \boldsymbol{W_{xh}x_t} + \boldsymbol{W_{ih}v} + b_h)$$

$\boldsymbol{W_{ih}}$

Sample word and copy to input

Before:

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

Now:

$$\tanh(\boldsymbol{W_{hh}h_{t-1}} + \boldsymbol{W_{xh}x_t} + \boldsymbol{W_{ih}v} + b_h)$$

74

# Image Captioning: Example Results

*A cat sitting on a suitcase on the floor*



*A cat is sitting on a tree branch*



*A dog is running in the grass with a frisbee*



*A white teddy bear sitting in the grass*



*Two people walking on the beach with surfboards*



*A tennis player in action on the court*



*Two giraffes standing in a grassy field*



*A man riding a dirt bike on a dirt track*

Slide from Justin Johnson

# Image Captioning: Failure Cases

*A woman is holding a cat in her hand*



*A person holding a computer mouse on a desk*



*A woman standing on a beach holding a surfboard*



*A bird is perched on a tree branch*



*A man in a baseball uniform throwing a ball*

# Vanilla RNN Gradient Flow



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

$$= \tanh\left((W_{hh} \quad W_{hx})\begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} + b_h\right)$$

$$= \tanh\left(W\begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} + b_h\right)$$

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013

Slide from Justin Johnson

# Vanilla RNN Gradient Flow

Backpropagation from
$h_t$ to $h_{t-1}$ multiplies by W
(actually $W_{hh}^T$)



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

$$= \tanh\left((W_{hh} \quad W_{hx})\begin{pmatrix}h_{t-1}\\x_t\end{pmatrix} + b_h\right)$$

$$= \tanh\left(W\begin{pmatrix}h_{t-1}\\x_t\end{pmatrix} + b_h\right)$$

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013

Slide from Justin Johnson

78

# Vanilla RNN Gradient Flow



Computing gradient of
$h_0$ involves many
factors of W
(and repeated tanh)

# Vanilla RNN Gradient Flow



Computing gradient of $h_0$ involves many factors of W (and repeated tanh)

Largest singular value > 1:
**Exploding gradients**

Largest singular value < 1:
**Vanishing gradients**

# Vanilla RNN Gradient Flow



Computing gradient of $h_0$ involves many factors of W (and repeated tanh)

Largest singular value > 1:
**Exploding gradients**

Largest singular value < 1:
**Vanishing gradients**

**Gradient clipping**: Scale gradient if its norm is too big

```
grad_norm = np.sum(grad * grad)
if grad_norm > threshold:
    grad *= (threshold / grad_norm)
```

# Vanilla RNN Gradient Flow



Computing gradient of $h_0$ involves many factors of W (and repeated tanh)

Largest singular value > 1:
**Exploding gradients**

Largest singular value < 1:
**Vanishing gradients** $\rightarrow$ **Change RNN architecture!**

# Long Short Term Memory (LSTM)

**Vanilla RNN**

$$h_t = \tanh\left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} + b_h\right)$$

Hochreiter and Schmidhuber, "Long Short Term Memory", Neural Computation 1997

# Long Short Term Memory (LSTM)

**Vanilla RNN**

$$h_t = \tanh\left(W\begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} + b_h\right)$$

**LSTM**

$$\begin{pmatrix} i_t \\ f_t \\ o_t \\ g_t \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix}\left(W\begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} + b_h\right)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$

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

Hochreiter and Schmidhuber, "Long Short Term Memory", Neural Computation 1997

Slide from Justin Johnson

84

# Long Short Term Memory (LSTM)

**Vanilla RNN**

$$h_t = \tanh\left(W\begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} + b_h\right)$$

Two vectors at each timestep:

Cell state: $c_t \in \mathbb{R}^H$

Hidden state: $h_t \in \mathbb{R}^H$

**LSTM**

$$\begin{pmatrix} i_t \\ f_t \\ o_t \\ g_t \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix}\left(W\begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} + b_h\right)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$

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

Hochreiter and Schmidhuber, "Long Short Term Memory", Neural Computation 1997

Slide from Justin Johnson

85

# Long Short Term Memory (LSTM)

**Vanilla RNN**

$$h_t = \tanh\left(W\begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} + b_h\right)$$

Compute four "gates" per timestep:

Input gate: $i_t \in \mathbb{R}^H$

Forget gate: $f_t \in \mathbb{R}^H$

Output gate: $o_t \in \mathbb{R}^H$

"Gate?" gate: $g_t \in \mathbb{R}^H$

**LSTM**

$$\begin{pmatrix} i_t \\ f_t \\ o_t \\ g_t \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix}\left(W\begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} + b_h\right)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$

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

Hochreiter and Schmidhuber, "Long Short Term Memory", Neural Computation 1997

Slide from Justin Johnson

86

# Long Short Term Memory (LSTM)

**i**: <u>Input gate</u>, whether to write to cell

**f**: <u>Forget gate</u>, Whether to erase cell

**o**: <u>Output gate</u>, How much to reveal cell

**g**: <u>Gate gate</u> (?), How much to write to cell

Input vector (x)



x

h

W

Previous hidden state (h)

4h x 2h

sigmoid → i

sigmoid → f

sigmoid → o

tanh → g

4h

4*h

$$\begin{pmatrix} i_t \\ f_t \\ o_t \\ g_t \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} \left( W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} + b_h \right)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$

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

# Long Short Term Memory (LSTM)



$$\begin{pmatrix} i_t \\ f_t \\ o_t \\ g_t \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} \left( W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} + b_h \right)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$

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

# Long Short Term Memory (LSTM): Gradient Flow



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

$$\begin{pmatrix} i_t \\ f_t \\ o_t \\ g_t \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} \left( W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} + b_h \right)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$
$$h_t = o_t \odot \tanh(c_t)$$

# Long Short Term Memory (LSTM): Gradient Flow

## Uninterrupted gradient flow!

# Long Short Term Memory (LSTM): Gradient Flow

## Uninterrupted gradient flow!



Similar to ResNet!

# Long Short Term Memory (LSTM): Gradient Flow

## Uninterrupted gradient flow!



Similar to ResNet!

In between: **Highway Networks**

$$g_t = F(x, W_t)$$
$$y_t = g_t \odot H(x, W_h) + (1 - g_t) \odot x_t$$

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

Slide from Justin Johnson

# Single-Layer RNNs

$$h_t = \tanh\left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} + b_h\right)$$

LSTM:

$$\begin{pmatrix} i_t \\ f_t \\ o_t \\ g_t \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} \left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} + b_h\right)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$

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

93

# Mutilayer RNNs

$$h_t^\ell = \tanh\left(W\begin{pmatrix} h_{t-1}^\ell \\ h_t^{\ell-1} \end{pmatrix} + b_h^\ell\right)$$

LSTM:

$$\begin{pmatrix} i_t^\ell \\ f_t^\ell \\ o_t^\ell \\ g_t^\ell \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix}\left(W\begin{pmatrix} h_{t-1}^\ell \\ h_t^{\ell-1} \end{pmatrix} + b_h^\ell\right)$$

$$c_t^\ell = f_t^\ell \odot c_{t-1}^\ell + i_t^\ell \odot g_t^\ell$$

$$h_t^\ell = o_t^\ell \odot \tanh(c_t^\ell)$$

depth

**Two-layer RNN**: Pass hidden states from one RNN as inputs to another RNN



time

Slide from Justin Johnson

94

# Mutilayer RNNs

$$h_t^\ell = \tanh\left(W\begin{pmatrix} h_{t-1}^\ell \\ h_t^{\ell-1} \end{pmatrix} + b_h^\ell\right)$$

LSTM:

$$\begin{pmatrix} i_t^\ell \\ f_t^\ell \\ o_t^\ell \\ g_t^\ell \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix}\left(W\begin{pmatrix} h_{t-1}^\ell \\ h_t^{\ell-1} \end{pmatrix} + b_h^\ell\right)$$

$$c_t^\ell = f_t^\ell \odot c_{t-1}^\ell + i_t^\ell \odot g_t^\ell$$

$$h_t^\ell = o_t^\ell \odot \tanh(c_t^\ell)$$

**Three-layer RNN**



Slide from Justin Johnson

# Other RNN Variants

**Gated Recurrent Unit (GRU)**

Cho et al "Learning phrase representations using RNN encoder-decoder for statistical machine translation", 2014

$$r_t = \sigma(W_{xr}x_t + W_{hr}h_{t-1} + b_r)$$
$$z_t = \sigma(W_{xz}x_t + W_{hz}h_{t-1} + b_z)$$
$$\tilde{h}_t = \tanh(W_{xh}x_t + W_{hh}(r_T \odot h_{t-1}) + b_h)$$
$$h_t = z_t \odot h_{t-1} + (1 - z_t) \odot \tilde{h}_t$$

# Sequence-to-Sequence with RNNs

**Input**: Sequence $x_1, \ldots x_T$
**Output**: Sequence $y_1, \ldots, y_{T'}$

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



Sutskever et al, "Sequence to sequence learning with neural networks", NeurIPS 2014

Slide from Justin Johnson          97

# Sequence-to-Sequence with RNNs

**Input**: Sequence $x_1, \ldots x_T$
**Output**: Sequence $y_1, \ldots, y_{T'}$

From final hidden state predict:
**Initial decoder state** $s_0$
**Context vector** c (often $c=h_T$)

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



Sutskever et al, "Sequence to sequence learning with neural networks", NeurIPS 2014

Slide from Justin Johnson

# Sequence-to-Sequence with RNNs

**Input**: Sequence $x_1, \ldots x_T$

**Output**: Sequence $y_1, \ldots, y_{T'}$

**Decoder**: $s_t = g_U(y_{t-1}, s_{t-1}, c)$

estamos

From final hidden state predict:
**Initial decoder state** $s_0$
**Context vector** $c$ (often $c=h_T$)

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



| $h_1$ | $h_2$ | $h_3$ | $h_4$ | | $s_0$ | | $s_1$ |
| $x_1$ | $x_2$ | $x_3$ | $x_4$ | | $c$ | | $y_0$ |

$y_1$

we    are    eating    bread                    [START]

Sutskever et al, "Sequence to sequence learning with neural networks", NeurIPS 2014

Slide from Justin Johnson

99

# Sequence-to-Sequence with RNNs

**Input**: Sequence $x_1, \dots x_T$

**Output**: Sequence $y_1, \dots, y_{T'}$

**Decoder:** $s_t = g_U(y_{t-1}, s_{t-1}, c)$

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

From final hidden state predict:
**Initial decoder state** $s_0$
**Context vector** $c$ (often $c=h_T$)



Sutskever et al, "Sequence to sequence learning with neural networks", NeurIPS 2014

Slide from Justin Johnson

100

# Sequence-to-Sequence with RNNs

**Input**: Sequence $x_1, \ldots x_T$

**Output**: Sequence $y_1, \ldots, y_{T'}$

**Decoder:** $s_t = g_U(y_{t-1}, s_{t-1}, c)$

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

From final hidden state predict:
**Initial decoder state** $s_0$
**Context vector** $c$ (often $c=h_T$)



Sutskever et al, "Sequence to sequence learning with neural networks", NeurIPS 2014

Slide from Justin Johnson

101

# Sequence-to-Sequence with RNNs

**Input**: Sequence $x_1, \ldots x_T$

**Output**: Sequence $y_1, \ldots, y_{T'}$

**Decoder:** $s_t = g_U(y_{t-1}, s_{t-1}, c)$

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

From final hidden state predict:
**Initial decoder state** $s_0$
**Context vector** c (often c=$h_T$)



estamos    comiendo    pan    [STOP]

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

$h_1$  $h_2$  $h_3$  $h_4$

$s_0$

$s_1$  $s_2$  $s_3$  $s_4$

c

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

$y_0$  $y_1$  $y_2$  $y_3$

we    are    eating    bread

[START]    estamos    comiendo    pan

**Problem: Input sequence bottlenecked through fixed-sized vector. What if T=1000?**

Sutskever et al, "Sequence to sequence learning with neural networks", NeurIPS 2014

102

# Sequence-to-Sequence with RNNs

**Input**: Sequence $x_1, \ldots x_T$

**Output**: Sequence $y_1, \ldots, y_{T'}$

**Decoder:** $s_t = g_U(y_{t-1}, s_{t-1}, c)$

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

From final hidden state predict:
**Initial decoder state** $s_0$
**Context vector** c (often $c=h_T$)

estamos    comiendo    pan    [STOP]



**Problem: Input sequence bottlenecked through fixed-sized vector. What if T=1000?**

**Idea: use new context vector at each step of decoder!**

Sutskever et al, "Sequence to sequence learning with neural networks", NeurIPS 2014

# Sequence-to-Sequence with RNNs **and Attention**

**Input**: Sequence $x_1, \ldots x_T$
**Output**: Sequence $y_1, \ldots, y_{T'}$

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

From final hidden state:
**Initial decoder state** $s_0$

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

Slide from Justin Johnson

104

# Sequence-to-Sequence with RNNs **and Attention**

Compute (scalar) **alignment scores**
$e_{t,i} = f_{att}(s_{t-1}, h_i)$     ($f_{att}$ is an MLP)

From final hidden state:
**Initial decoder state** $s_0$



we     are     eating     bread

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

Slide from Justin Johnson

# Sequence-to-Sequence with RNNs **and Attention**



Compute (scalar) **alignment scores**
$e_{t,i} = f_{att}(s_{t-1}, h_i)$ ($f_{att}$ is an MLP)

Normalize alignment scores
to get **attention weights**
$0 < a_{t,i} < 1$ $\sum_i a_{t,i} = 1$

From final hidden state:
**Initial decoder state** $s_0$

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

# Sequence-to-Sequence with RNNs **and Attention**



Compute (scalar) **alignment scores**
$e_{t,i} = f_{att}(s_{t-1}, h_i)$      ($f_{att}$ is an MLP)

Normalize alignment scores to get **attention weights**
$0 < a_{t,i} < 1$    $\sum_i a_{t,i} = 1$

Compute context vector as linear combination of hidden states
$c_t = \sum_i a_{t,i} h_i$

Use context vector in decoder: $s_t = g_U(y_{t-1}, s_{t-1}, c_t)$

From final hidden state:
**Initial decoder state** $s_0$

**This is all differentiable! Do not supervise attention weights – backprop through everything**

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

Slide from Justin Johnson

107

# Sequence-to-Sequence with RNNs **and Attention**



Compute (scalar) **alignment scores**
$e_{t,i} = f_{att}(s_{t-1}, h_i)$      ($f_{att}$ is an MLP)

Normalize alignment scores
to get **attention weights**
$0 < a_{t,i} < 1$    $\sum_i a_{t,i} = 1$

Compute context vector as linear
combination of hidden states
$c_t = \sum_i a_{t,i} h_i$

Use context vector in
decoder: $s_t = g_U(y_{t-1}, s_{t-1}, c_t)$

From final hidden state:
**Initial decoder state** $s_0$

**Intuition**: Context vector
<u>attends</u> to the relevant
part of the input sequence
*"estamos" = "we are"*
so maybe $a_{11}=a_{12}=0.45$,
$a_{13}=a_{14}=0.05$

**This is all differentiable! Do not
supervise attention weights –
backprop through everything**

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

Slide from Justin Johnson

108

# Sequence-to-Sequence with RNNs

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

Slide from Justin Johnson

109

# Sequence-to-Sequence with RNNs **and Attention**



Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

Slide from Justin Johnson

# Sequence-to-Sequence with RNNs **and Attention**



**Intuition**: Context vector attends to the relevant part of the input sequence *"comiendo" = "eating"* so maybe $a_{21}=a_{24}=0.05$, $a_{22}=0.1$, $a_{23}=0.8$

Repeat: Use $s_1$ to compute new context vector $c_2$

Use $c_2$ to compute $s_2$, $y_2$

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

Slide from Justin Johnson

111

# Sequence-to-Sequence with RNNs **and Attention**

**Use a different context vector in each timestep of decoder**
- Input sequence not bottlenecked through single vector
- At each timestep of decoder, context vector "looks at" different parts of the input sequence



Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

Slide from Justin Johnson

# Sequence-to-Sequence with RNNs **and Attention**

**Example**: English to French translation

**Input**: "The agreement on the European Economic Area was signed in August 1992."

**Output**: "L'accord sur la zone économique européenne a été signé en août 1992."

Visualize attention weights $a_{t,i}$



Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

# Sequence-to-Sequence with RNNs **and Attention**

**Example**: English to French translation

**Input**: "**The agreement on the** European Economic Area was signed **in August 1992**."

**Output**: "**L'accord sur la** zone économique européenne a été signé **en août 1992**."

Visualize attention weights $a_{t,i}$

**Diagonal attention means words correspond in order**

**Diagonal attention means words correspond in order**



Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

Slide from Justin Johnson

114

# Sequence-to-Sequence with RNNs **and Attention**

**Example**: English to French translation

**Input**: "**The agreement on the European Economic Area** was signed **in August 1992**."

**Output**: "**L'accord sur la zone économique européenne** a été signé **en août 1992**."

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

Visualize attention weights $a_{t,i}$

Diagonal attention means words correspond in order

Attention figures out different word orders

Diagonal attention means words correspond in order

# Sequence-to-Sequence with RNNs **and Attention**

Visualize attention weights $a_{t,i}$

**Example**: English to French translation

**Input**: "**The agreement on the European Economic Area** **was signed** **in August 1992**."

**Output**: "**L'accord sur la zone économique européenne** **a été signé** **en août 1992**."

Diagonal attention means words correspond in order

Attention figures out different word orders

Verb conjugation

Diagonal attention means words correspond in order



Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

Slide from Justin Johnson

116

# Sequence-to-Sequence with RNNs **and Attention**

The decoder doesn't use the fact that $h_i$ form an ordered sequence – it just treats them as an unordered set $\{h_i\}$

Can use similar architecture given any set of input hidden vectors $\{h_i\}$!



Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

Slide from Justin Johnson

117

# Attention



$$\boldsymbol{y_i} = \sum_j w_{ij}\boldsymbol{x_{ij}}$$

$$w_{ij} = softmax_j(\boldsymbol{x}_i^T \boldsymbol{x}_j / \sqrt{d_k})$$

$$w_{ij} = \frac{e^{\boldsymbol{x}_i^T \boldsymbol{x}_j}}{\sum_j e^{\boldsymbol{x}_i^T \boldsymbol{x}_j}}$$

Source: http://peterbloem.nl/blog/transformers  See also: Attention is all you need

# Attention (with key, query and value)



$$y_i = \sum_j w_{ij} W_v x_{ij}$$

$$w_{ij} = softmax_j((W_q x_i)^T W_k x_j / \sqrt{d_k})$$

Source: http://peterbloem.nl/blog/transformers  See also: Attention is all you need

# The Transformer

$$\boxed{x_1} \quad \boxed{x_2} \quad \boxed{x_3} \quad \boxed{x_4}$$

Vaswani et al, "Attention is all you need", NeurIPS 2017

# The Transformer

All vectors interact
with each other



Vaswani et al, "Attention is all you need", NeurIPS 2017

Slide from Justin Johnson

121

# The Transformer

Residual connection

All vectors interact
with each other



Self-Attention

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

Vaswani et al, "Attention is all you need", NeurIPS 2017

Slide from Justin Johnson

# The Transformer

Recall **Layer Normalization**:

Given $h_1, ..., h_N$     (Shape: D)

scale: $\gamma$                (Shape: D)

shift: $\beta$                 (Shape: D)

$\mu_i = (\sum_j h_{i,j})/D$         (scalar)

$\sigma_i = (\sum_j (h_{i,j} - \mu_i)^2/D)^{1/2}$   (scalar)

$z_i = (h_i - \mu_i) / \sigma_i$

$y_i = \gamma * z_i + \beta$

Ba et al, 2016

Residual connection

All vectors interact with each other



Vaswani et al, "Attention is all you need", NeurIPS 2017

Slide from Justin Johnson

# The Transformer

Recall **Layer Normalization**:
Given $h_1, ..., h_N$ (Shape: D)
scale: $\gamma$ (Shape: D)
shift: $\beta$ (Shape: D)
$\mu_i = (\sum_j h_{i,j})/D$ (scalar)
$\sigma_i = (\sum_j (h_{i,j} - \mu_i)^2/D)^{1/2}$ (scalar)
$z_i = (h_i - \mu_i) / \sigma_i$
$y_i = \gamma * z_i + \beta$

Ba et al, 2016

MLP independently
on each vector

Residual connection

All vectors interact
with each other



Vaswani et al, "Attention is all you need", NeurIPS 2017

Slide from Justin Johnson

# The Transformer

Recall **Layer Normalization**:
Given $h_1, ..., h_N$    (Shape: D)
scale: $\gamma$              (Shape: D)
shift: $\beta$              (Shape: D)
$\mu_i = (\sum_j h_{i,j})/D$         (scalar)
$\sigma_i = (\sum_j (h_{i,j} - \mu_i)^2/D)^{1/2}$  (scalar)
$z_i = (h_i - \mu_i) / \sigma_i$
$y_i = \gamma * z_i + \beta$

Ba et al, 2016

Residual connection

MLP independently
on each vector

Residual connection

All vectors interact
with each other



Vaswani et al, "Attention is all you need", NeurIPS 2017

Slide from Justin Johnson

# The Transformer

Recall **Layer Normalization**:
Given $h_1, ..., h_N$    (Shape: D)
scale: $\gamma$               (Shape: D)
shift: $\beta$              (Shape: D)
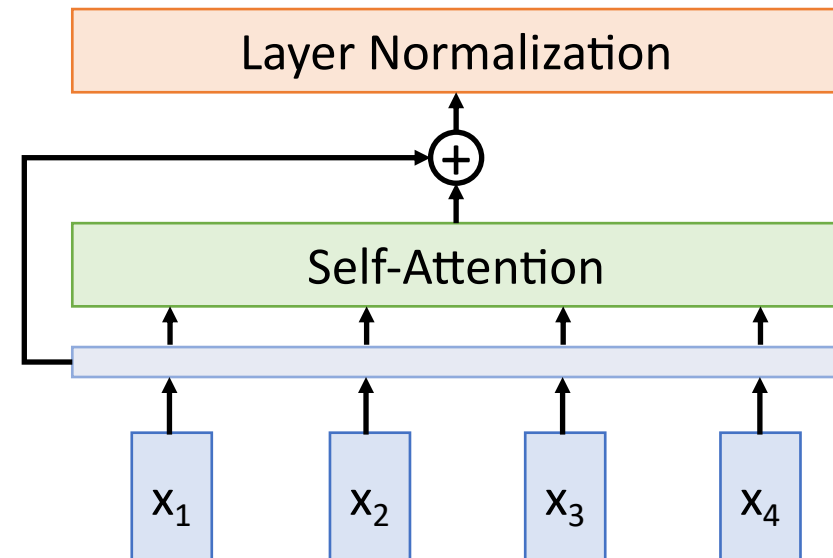$\mu_i = (\sum_j h_{i,j})/D$       (scalar)
$\sigma_i = (\sum_j (h_{i,j} - \mu_i)^2/D)^{1/2}$  (scalar)
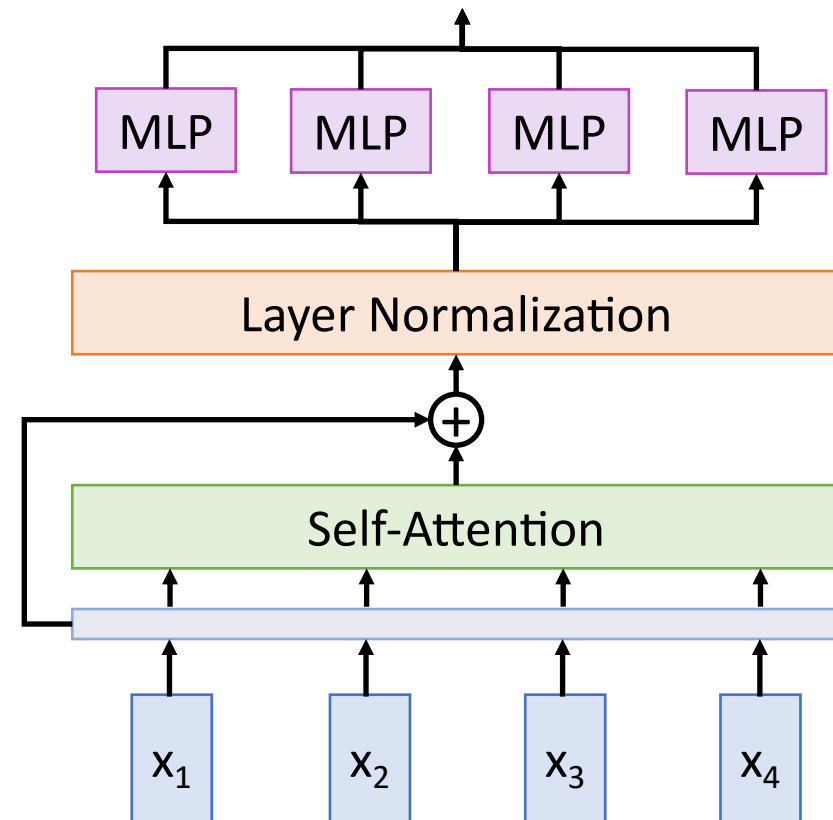$z_i = (h_i - \mu_i) / \sigma_i$
$y_i = \gamma * z_i + \beta$

Ba et al, 2016

Residual connection

MLP independently
on each vector

Residual connection

All vectors interact
with each other



Vaswani et al, "Attention is all you need", NeurIPS 2017

Slide from Justin Johnson

# The Transformer

**Transformer Block:**
**Input**: Set of vectors x
**Output**: Set of vectors y

Self-attention is the only interaction between vectors!

Layer norm and MLP work independently per vector

Highly scalable, highly parallelizable

Vaswani et al, "Attention is all you need", NeurIPS 2017

# Post-Norm Transformer



**Layer normalization** is **after** residual connections

Vaswani et al, "Attention is all you need", NeurIPS 2017

Slide from Justin Johnson

# Pre-Norm Transformer

**Layer normalization** is **inside** residual connections

Gives more stable training, commonly used in practice



Baevski & Auli, "Adaptive Input Representations for Neural Language Modeling", arXiv 2018

Slide from Justin Johnson

# The Transformer



**Transformer Block:**
**Input**: Set of vectors x
**Output**: Set of vectors y

Self-attention is the only interaction between vectors!

Layer norm and MLP work independently per vector

Highly scalable, highly parallelizable

A **Transformer** is a sequence of transformer blocks
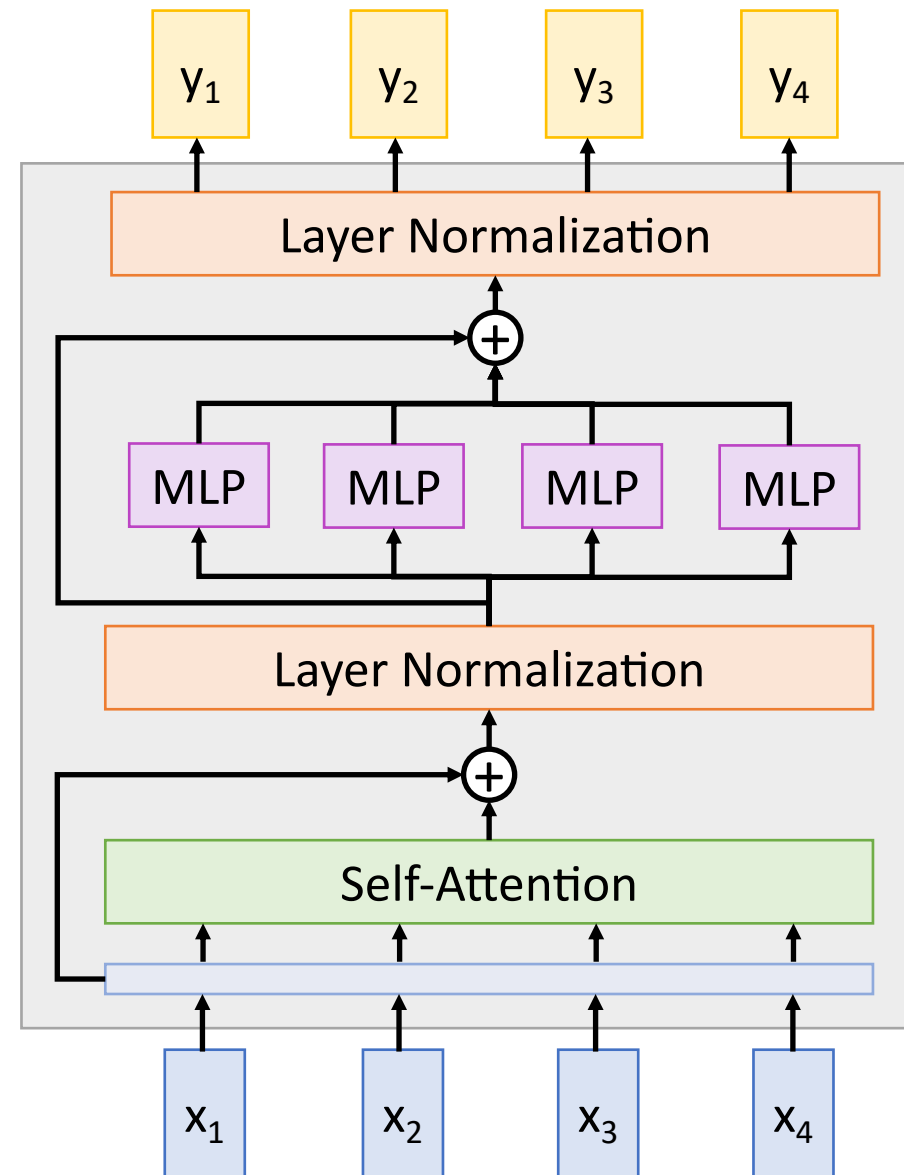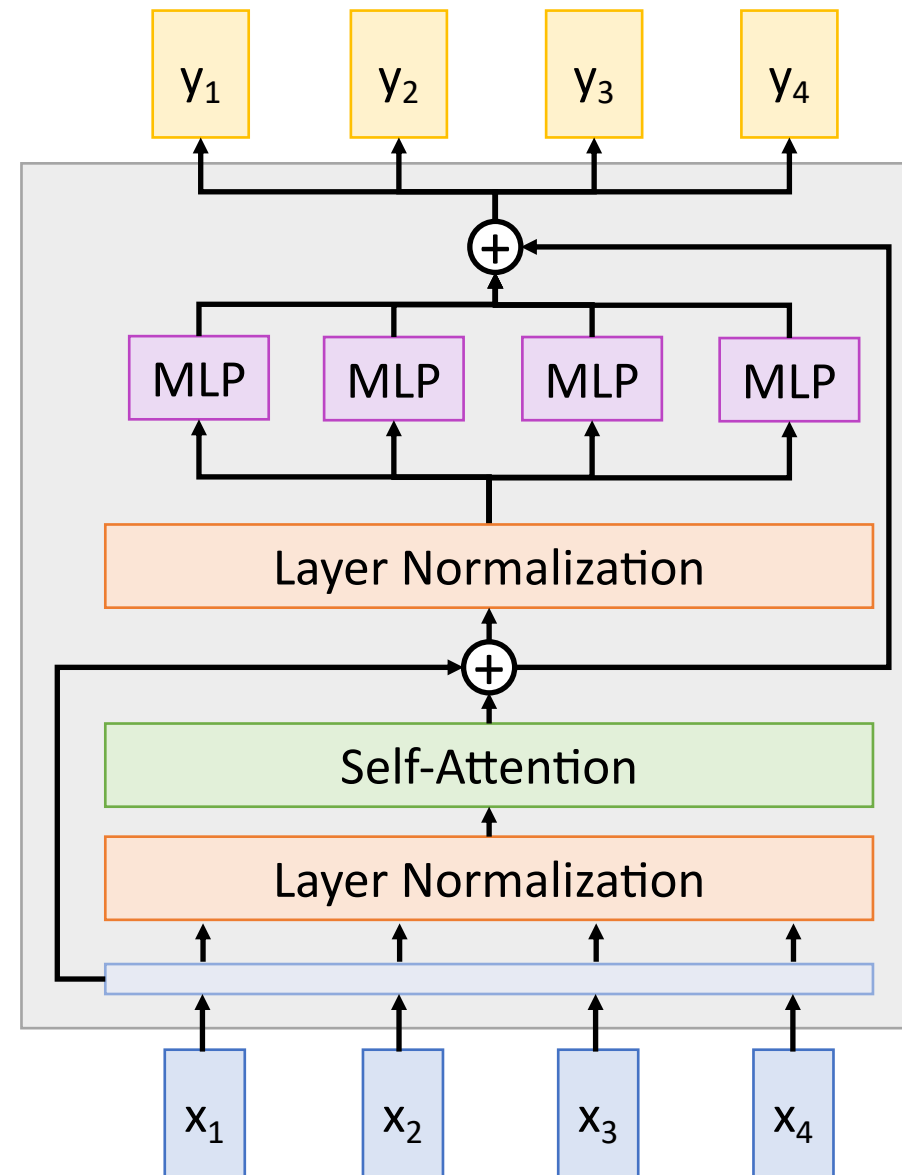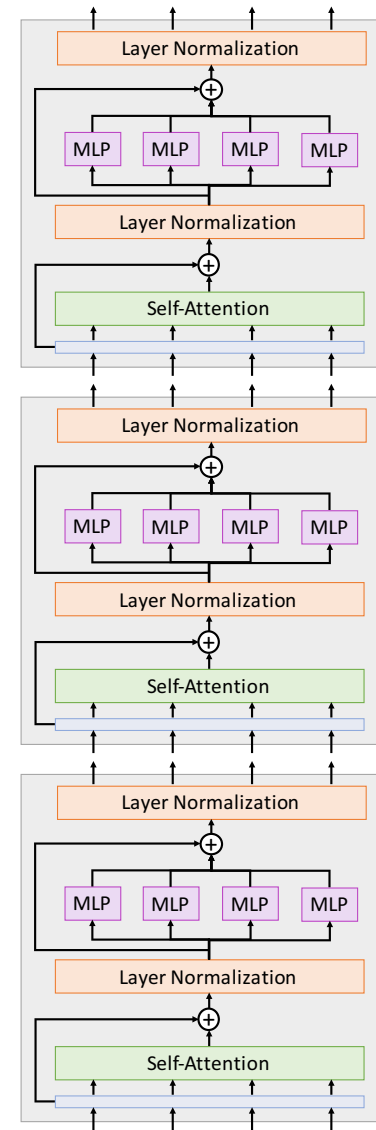
Vaswani et al:
12 blocks, $D_Q$=512, 6 heads

Vaswani et al, "Attention is all you need", NeurIPS 2017

Slide from Justin Johnson

# The Transformer

Table 2: The Transformer achieves better BLEU scores than previous state-of-the-art models on the English-to-German and English-to-French newstest2014 tests at a fraction of the training cost.

| Model | BLEU | | Training Cost (FLOPs) | |
|---|---|---|---|---|
| | EN-DE | EN-FR | EN-DE | EN-FR |
| ByteNet [18] | 23.75 | | | |
| Deep-Att + PosUnk [39] | | 39.2 | | $1.0 \cdot 10^{20}$ |
| GNMT + RL [38] | 24.6 | 39.92 | $2.3 \cdot 10^{19}$ | $1.4 \cdot 10^{20}$ |
| ConvS2S [9] | 25.16 | 40.46 | $9.6 \cdot 10^{18}$ | $1.5 \cdot 10^{20}$ |
| MoE [32] | 26.03 | 40.56 | $2.0 \cdot 10^{19}$ | $1.2 \cdot 10^{20}$ |
| Deep-Att + PosUnk Ensemble [39] | | 40.4 | | $8.0 \cdot 10^{20}$ |
| GNMT + RL Ensemble [38] | 26.30 | 41.16 | $1.8 \cdot 10^{20}$ | $1.1 \cdot 10^{21}$ |
| ConvS2S Ensemble [9] | 26.36 | **41.29** | $7.7 \cdot 10^{19}$ | $1.2 \cdot 10^{21}$ |
| Transformer (base model) | 27.3 | 38.1 | $\mathbf{3.3 \cdot 10^{18}}$ | |
| Transformer (big) | **28.4** | **41.8** | $2.3 \cdot 10^{19}$ | |

Vaswani et al, "Attention is all you need", NeurIPS 2017

131

# The Transformer: Transfer Learning



"ImageNet Moment for Natural Language Processing"

**Pretraining**:
Download a lot of text from the internet
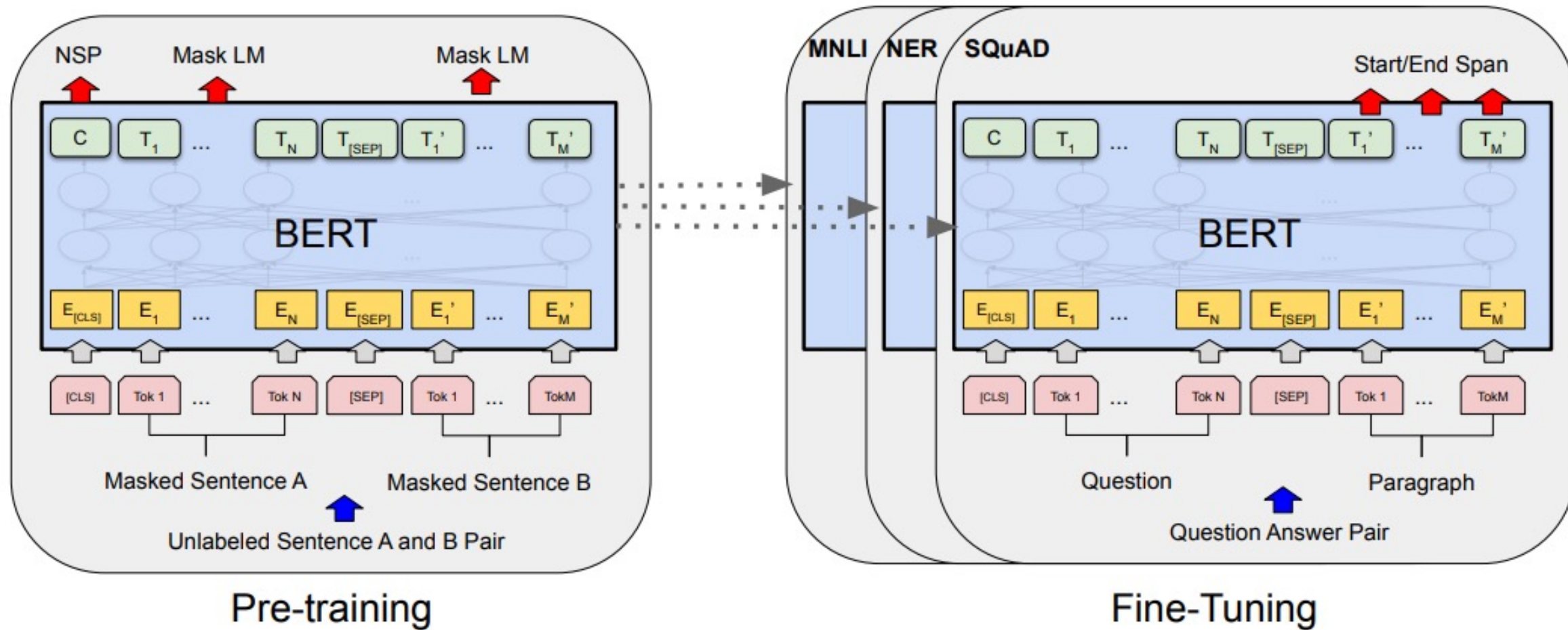
Train a giant Transformer model for language modeling

**Finetuning:**
Fine-tune the Transformer on your own NLP task

Devlin et al, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding", EMNLP 2018

132

# The Transformer: Transfer Learning



Devlin et al, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding", EMNLP 2018

# The Transformer: Transfer Learning

| System | MNLI-(m/mm) 392k | QQP 363k | QNLI 108k | SST-2 67k | CoLA 8.5k | STS-B 5.7k | MRPC 3.5k | RTE 2.5k | **Average** - |
|---|---|---|---|---|---|---|---|---|---|
| Pre-OpenAI SOTA | 80.6/80.1 | 66.1 | 82.3 | 93.2 | 35.0 | 81.0 | 86.0 | 61.7 | 74.0 |
| BiLSTM+ELMo+Attn | 76.4/76.1 | 64.8 | 79.8 | 90.4 | 36.0 | 73.3 | 84.9 | 56.8 | 71.0 |
| OpenAI GPT | 82.1/81.4 | 70.3 | 87.4 | 91.3 | 45.4 | 80.0 | 82.3 | 56.0 | 75.1 |
| BERT$_{BASE}$ | 84.6/83.4 | 71.2 | 90.5 | 93.5 | 52.1 | 85.8 | 88.9 | 66.4 | 79.6 |
| BERT$_{LARGE}$ | **86.7/85.9** | **72.1** | **92.7** | **94.9** | **60.5** | **86.5** | **89.3** | **70.1** | **82.1** |

Table 1: GLUE Test results, scored by the evaluation server (https://gluebenchmark.com/leaderboard). The number below each task denotes the number of training examples. The "Average" column is slightly different than the official GLUE score, since we exclude the problematic WNLI set.[8] BERT and OpenAI GPT are single-model, single task. F1 scores are reported for QQP and MRPC, Spearman correlations are reported for STS-B, and accuracy scores are reported for the other tasks. We exclude entries that use BERT as one of their components.

Devlin et al, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding", EMNLP 2018