

Recurrent  
Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

Deep RNNs

LSTMs

GRUs

NLP

Summary

# Recurrent Architectures

Stephen Scott

(Adapted from Vinod Variyam and Ian Goodfellow)

[sscott2@unl.edu](mailto:sscott2@unl.edu)

# Introduction

Recurrent  
Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

Deep RNNs

LSTMs

GRUs

NLP

Summary

- All our architectures so far work on fixed-sized inputs
- Recurrent neural networks work on variable-length **sequences** of inputs
- E.g., text, biological sequences, video, audio
- Can also try 1D convolutions, but can lose long-term relationships in input
- Especially useful for NLP applications: translation, speech-to-text, sentiment analysis
- Can also **create novel output**: e.g., Shakespearean text, music

# Outline

Recurrent  
Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

Deep RNNs

LSTMs

GRUs

NLP

Summary

- Basic RNNs
- Input/Output Mappings
- Training
- Long short-term memory
- Gated Recurrent Unit
- Applications to natural language processing
  - Embedded representations
  - Bidirectional RNNs
  - Attention mechanisms
  - Transformers
- Summary

# Basic Recurrent Cell

Recurrent  
Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

Deep RNNs

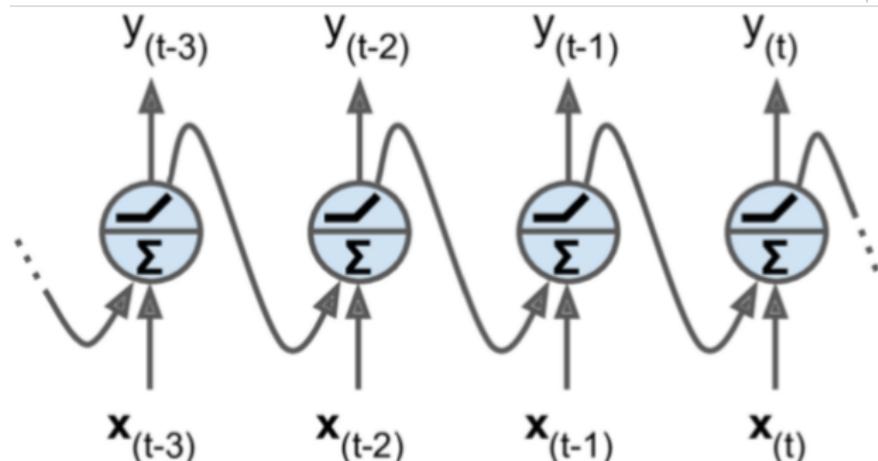
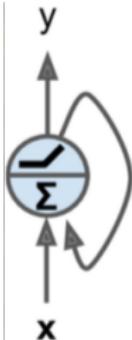
LSTMs

GRUs

NLP

Summary

- A recurrent cell (or recurrent neuron) has connections pointing **backward** as well as forward
- At time step (frame)  $t$ , neuron receives input vector  $x_{(t)}$  as usual, but also receives its own output  $y_{(t-1)}$  from previous step



# Basic Recurrent Layer

Recurrent  
Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

Deep RNNs

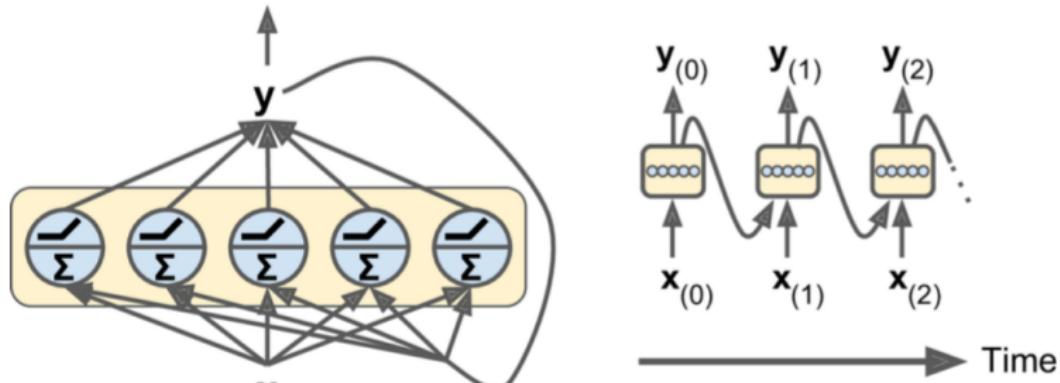
LSTMs

GRUs

NLP

Summary

- Can build a layer of recurrent cells, where each node gets both the vector  $x_{(t)}$  and the vector  $y_{(t-1)}$



# Basic Recurrent Layer

Recurrent  
Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

Deep RNNs

LSTMs

GRUs

NLP

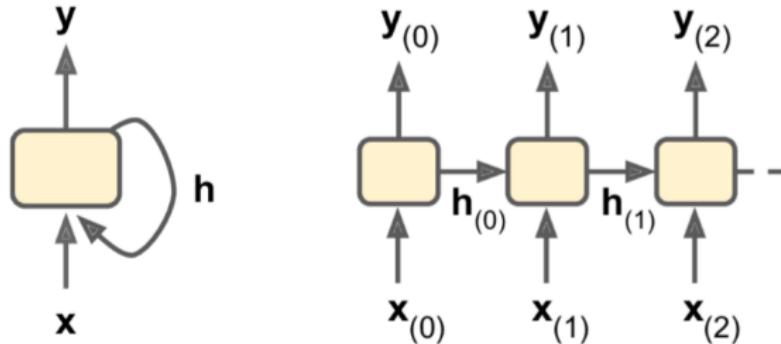
Summary

- Each node in the recurrent layer has independent weights for both  $x_{(t)}$  and  $y_{(t-1)}$
- For a single recurrent node, denote by  $w_x$  and  $w_y$
- For the entire layer, combine into matrices  $W_x$  and  $W_y$
- For activation function  $\phi$  and bias vector  $b$ , output vector is

$$\mathbf{y}_{(t)} = \phi \left( W_x^\top \mathbf{x}_{(t)} + W_y^\top \mathbf{y}_{(t-1)} + \mathbf{b} \right)$$

# Memory and State

- Since a node's output depends on its past, it can be thought of having **memory** or **state**
- State at time  $t$  is  $\mathbf{h}_{(t)} = f(\mathbf{h}_{(t-1)}, \mathbf{x}_{(t)})$  and output  $\mathbf{y}_{(t)} = g(\mathbf{h}_{(t-1)}, \mathbf{x}_{(t)})$
- State could be the same as the output, or separate
- Can think of  $\mathbf{h}_{(t)}$  as storing important information about input sequence
- Analogous to convolutional outputs summarizing important image features



# Input/Output Mappings

## Sequence to Sequence

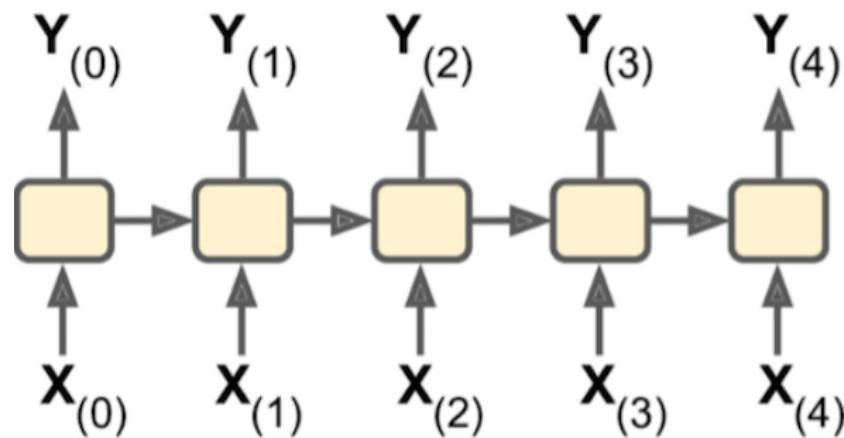
Recurrent  
Architectures  
Stephen Scott

Introduction  
Basic Idea  
I/O Mappings

Training  
Deep RNNs  
LSTMs  
GRUs  
NLP  
Summary

Many ways to employ this basic architecture:

- **Sequence to sequence:** Input is a sequence and output is a sequence
- E.g., series of stock predictions, one day in advance

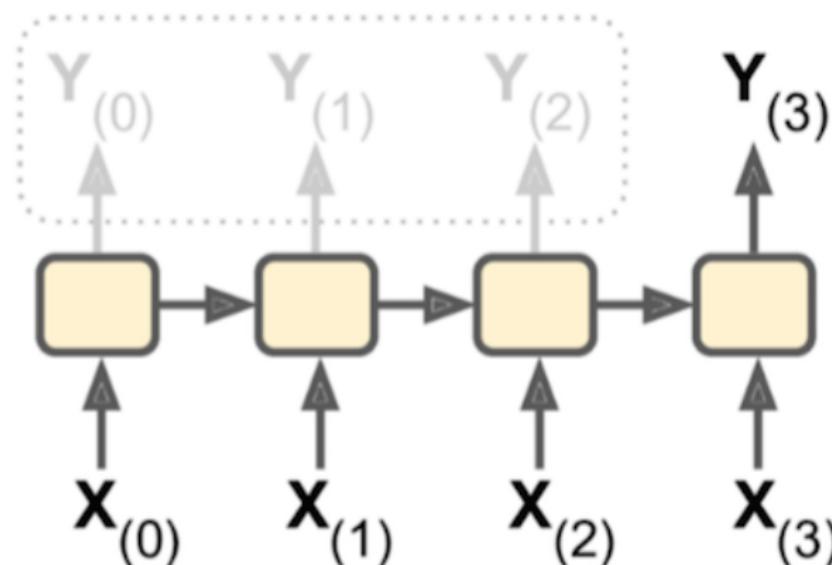


# Input/Output Mappings

## Sequence to Vector

- **Sequence to vector:** Input is sequence and output a vector/score/classification
- E.g., sentiment score of movie review

### Ignored outputs



# Input/Output Mappings

## Vector to Sequence

Recurrent  
Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

Deep RNNs

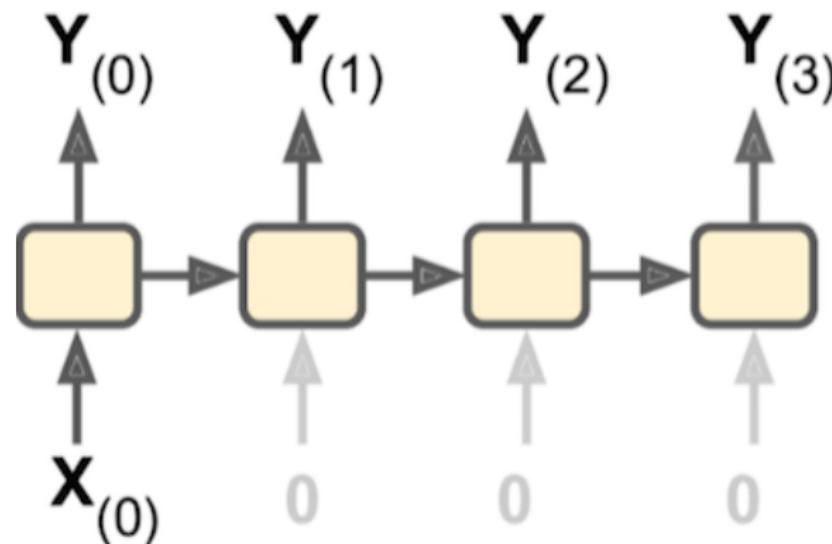
LSTMs

GRUs

NLP

Summary

- **Vector to sequence:** Input is a single vector (other times zeroes or repeated) and output is a sequence
- E.g., image to caption



# Input/Output Mappings

## Encoder-Decoder Architecture

Recurrent  
Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

Deep RNNs

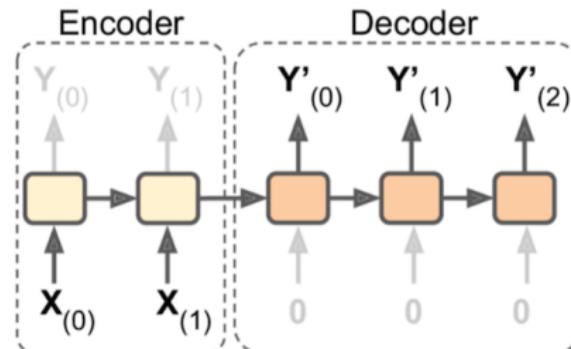
LSTMs

GRUs

NLP

Summary

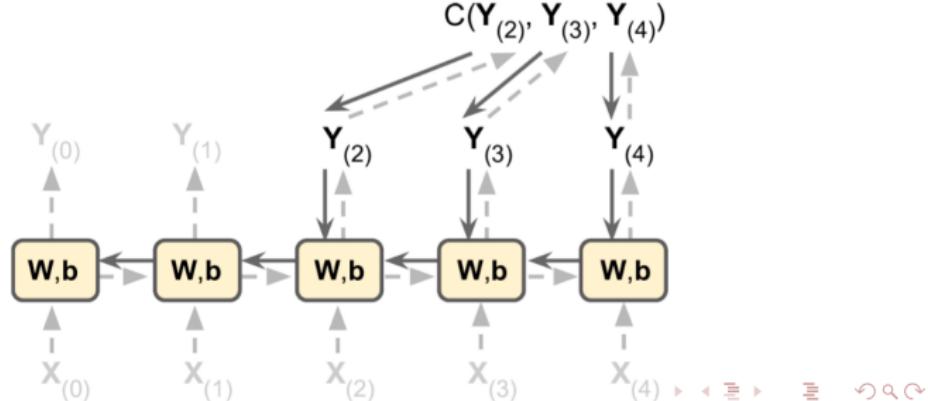
- Encoder-decoder: Sequence-to-vector (**encoder**) followed by vector-to-sequence (**decoder**)
- Input sequence  $(x_1, \dots, x_T)$  yields hidden outputs  $(h_1, \dots, h_T)$ , then mapped to **context vector**  $c = f(h_1, \dots, h_T)$
- Decoder output  $y_{t'}$  depends on previously output  $(y_1, \dots, y_{t'-1})$  and  $c$
- Example application: **neural machine translation**



# Training

## Backpropagation Through Time (BPTT)

- Unroll through time and use BPTT
- Forward pass mini-batch of sequences through unrolled network yields output sequence  $Y_{(0)}, \dots, Y_{(T)}$
- Output sequence evaluated using **sequence loss**  $C(Y_{(0)}, Y_{(1)}, \dots, Y_{(T)})$ , e.g., weighted sum of cross-entropies
- Gradients propagated backward through unrolled network (summing over all time steps), and parameters



# Training Issues

Recurrent Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

Deep RNNs

LSTMs

GRUs

NLP

Summary

- BPTT means that gradient is flowing through longer paths in graph ⇒ **exploding or vanishing gradients**
  - Can happen with any network, but RNNs very susceptible
  - **Clipping** gradients to range  $[-1, +1]$  can mitigate explosions

# Deep RNNs

Recurrent  
Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

Deep RNNs

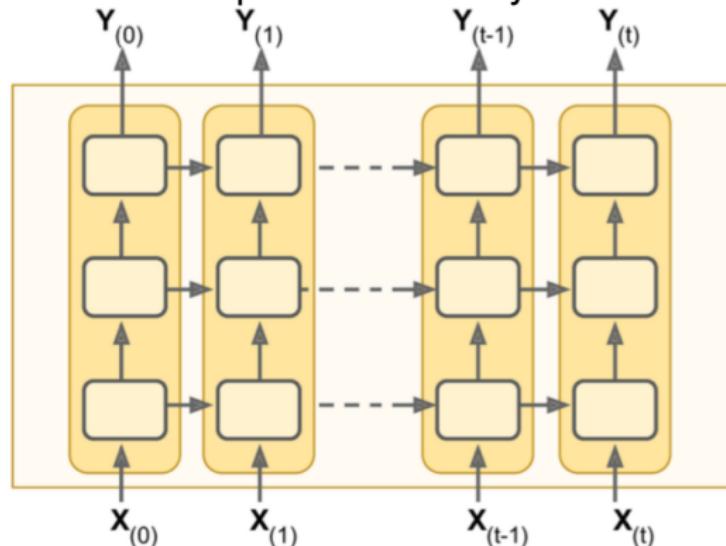
LSTMs

GRUs

NLP

Summary

A **deep RNN** has multiple recurrent layers stacked



# Training over Many Time Steps

Recurrent  
Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

Deep RNNs

LSTMs

GRUs

NLP

Summary

- Vanishing and exploding gradients can be a problem with RNNs, like with other deep networks
  - Can maybe address with gradient clipping, etc.
- Can still suffer from long training times with long input sequences
  - **Truncated backpropagation through time** addresses this by limiting the number of unrolled training steps
  - Inhibits ability to learn long-term patterns
- In general, also have problem of first inputs of sequence have diminishing impact as sequence grows
  - E.g., first few words of long text sequence
- Goal: Introduce **long-term memory** to RNNs
- Allow a network to **accumulate** information about the past, but also decide when to **forget** information

# Long Short-Term Memory

Hochreiter and Schmidhuber (1997)

Recurrent  
Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

Deep RNNs

LSTMs

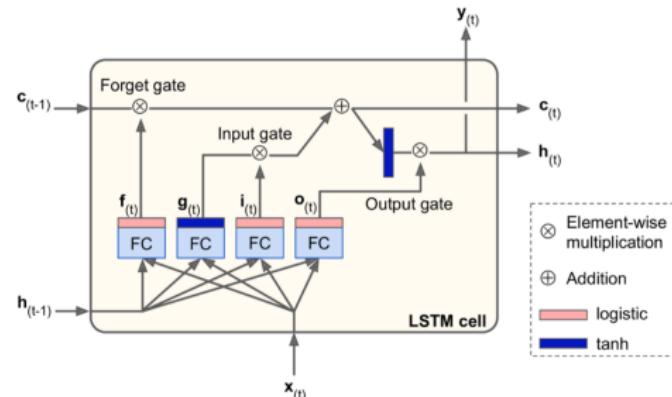
GRUs

NLP

Summary

- Vector  $h_{(t)}$  = **short-term state**,  $c_{(t)}$  = **long-term state**

- At time  $t$ , some memories from  $c_{(t-1)}$  are forgotten in the **forget gate** and new ones (from **input gate**) added



- Result sent out as  $c_{(t)}$
- $h_{(t)}$  (and  $y_{(t)}$ ) comes from processing long-term state in **output gate**

# Long Short-Term Memory

Hochreiter and Schmidhuber (1997)

Recurrent  
Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

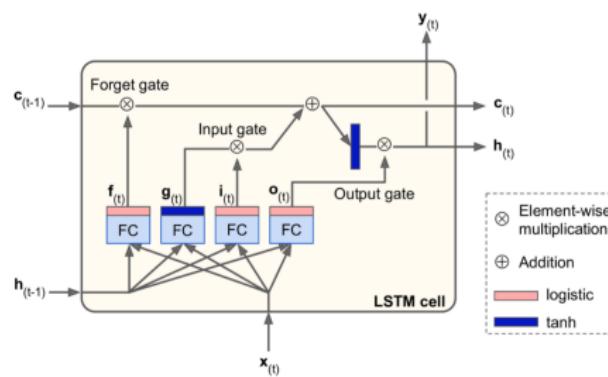
Deep RNNs

LSTMs

GRUs

NLP

Summary



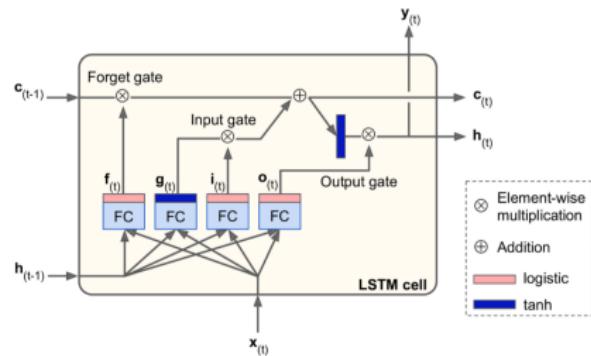
- $g_{(t)}$  combines input  $x_{(t)}$  with state  $h_{(t-1)}$
- $f_{(t)}, i_{(t)}, o_{(t)}$  are **gate controllers**
- $f_{(t)} \in [0, 1]^n$  controls forgetting of  $c_{(t-1)}$
- $i_{(t)}$  controls remembering of  $g_{(t)}$

- $o_{(t)}$  controls what of  $c_{(t)}$  goes to output and  $h_{(t)}$
- Output depends on long- and short-term memory
- Network learns what to remember long-term based on  $x_{(t)}$  and  $h_{(t-1)}$

# Long Short-Term Memory

- $\mathbf{i}_{(t)} = \sigma(W_{xi}^\top \mathbf{x}_{(t)} + W_{hi}^\top \mathbf{h}_{(t-1)} + \mathbf{b}_i)$
  - $\mathbf{f}_{(t)} = \sigma(W_{xf}^\top \mathbf{x}_{(t)} + W_{hf}^\top \mathbf{h}_{(t-1)} + \mathbf{b}_f)$
  - $\mathbf{o}_{(t)} = \sigma(W_{xo}^\top \mathbf{x}_{(t)} + W_{ho}^\top \mathbf{h}_{(t-1)} + \mathbf{b}_o)$
  - $\mathbf{g}_{(t)} = \tanh(W_{xg}^\top \mathbf{x}_{(t)} + W_{hg}^\top \mathbf{h}_{(t-1)} + \mathbf{b}_g)$

- $\mathbf{c}_{(t)} = \mathbf{f}_{(t)} \otimes \mathbf{c}_{(t-1)} + \mathbf{i}_{(t)} \otimes \mathbf{g}_{(t)}$
  - $\mathbf{y}_{(t)} = \mathbf{h}_{(t)} = \mathbf{o}_{(t)} \otimes \tanh(\mathbf{c}_{(t)})$



- Can add **peephole connection**: Let  $c_{(t-1)}$  affect  $f_{(t)}$  and  $i_{(t)}$  and  $c_{(t)}$  affect  $o_{(t)}$

# Gated Recurrent Unit

Recurrent  
Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

Deep RNNs

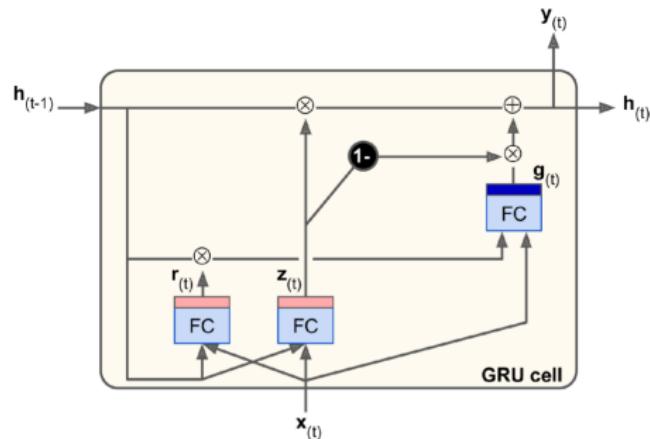
LSTMs

GRUs

NLP

Summary

- Simplified LSTM
- Merge  $c_{(t)}$  into  $\mathbf{h}_{(t)}$
- Merge  $f_{(t)}$  and  $i_{(t)}$  into  $\mathbf{z}_{(t)}$ 
  - $\mathbf{z}_{(t),i} = 0 \Rightarrow$  forget  $\mathbf{h}_{(t-1),i}$  and add in  $\mathbf{g}_{(t),i}$



- $\mathbf{o}_{(t)}$  replaced by  $r_{(t)} \Rightarrow$  forget part of  $\mathbf{h}_{(t-1)}$  when computing  $\mathbf{g}_{(t)}$

# Gated Recurrent Unit

Recurrent  
Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

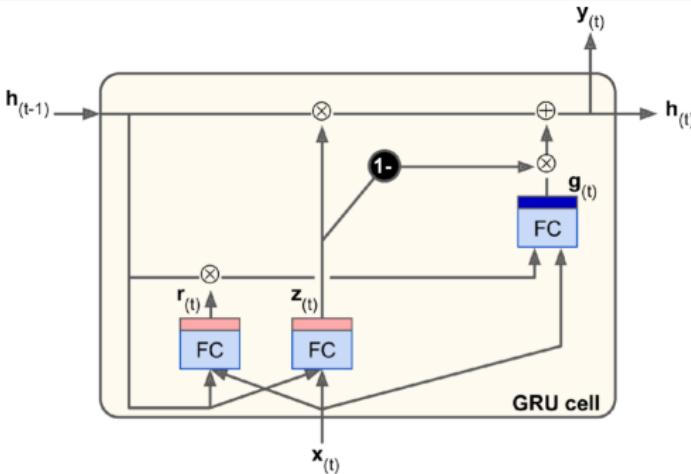
Deep RNNs

LSTMs

GRUs

NLP

Summary



- $z_{(t)} = \sigma (W_{xz}^\top x_{(t)} + W_{hz}^\top h_{(t-1)} + b_z)$
- $r_{(t)} = \sigma (W_{xr}^\top x_{(t)} + W_{hr}^\top h_{(t-1)} + b_r)$
- $g_{(t)} = \tanh (W_{xg}^\top x_{(t)} + W_{hg}^\top (r_{(t)} \otimes h_{(t-1)}) + b_g)$
- $y_{(t)} = h_{(t)} = z_{(t)} \otimes h_{(t-1)} + (1 - z_{(t)}) \otimes g_{(t)}$

# Application: Natural Language Processing

## Introduction

Recurrent  
Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

Deep RNNs

LSTMs

GRUs

NLP

Embeddings

Architecture

Bidirectional

Attention

Transformer

GPT

BERT

Summary

- Common application of sequential models such as RNNs
- Example problems:
  - Sentiment analysis
  - Information extraction
  - Question answering
  - Summarization
  - **Machine translation**
- Some tasks involve classification, some text generation

# Application: Natural Language Processing

## Outline

Recurrent  
Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

Deep RNNs

LSTMs

GRUs

NLP

Embeddings

Architecture

Bidirectional

Attention

Transformer

GPT

BERT

Summary

- Embedded representations of words
- Neural machine translation (NMT) via encoder-decoder architecture
- Bidirectional RNNs
- Attention mechanisms
- Transformer architectures

# Embedded Representations of Words

Recurrent  
Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

Deep RNNs

LSTMs

GRUs

NLP

Embeddings

Architecture

Bidirectional

Attention

Transformer

GPT

BERT

Summary

- To process text with RNN, need **numeric** representation of words
- Simple approach: **Sparse vectors**, using a **one-hot** vector per word
  - $|V|$ -dimensional vector, where  $V$  = vocabulary
  - Represent document with disjunction of its words' vectors ⇒ **bag-of-words**
- **Issues with this:**
  - Very high-dimensional inputs ( $|V| > 10^5$ )
  - All word vectors orthogonal ⇒ don't cluster
- Alternative: **Dense vector** representation where similar words near each other
- Train embedding as part of RNN or use pre-trained one
  - E.g., **word2vec**

# NMT via Encoder-Decoder Architecture

## Training

Recurrent  
Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

Deep RNNs

LSTMs

GRUs

NLP

Embeddings

Architecture

Bidirectional

Attention

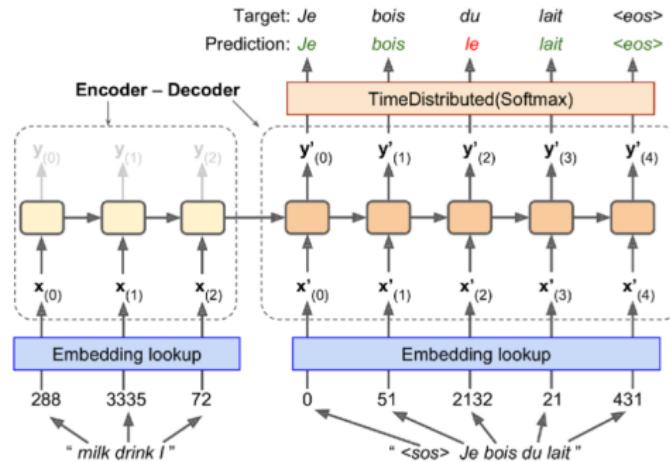
Transformer

GPT

BERT

Summary

- Pre-trained word embeddings fed into input
  - Both input and translation (shifted back one time step)
  - Each time step is classification  $\Rightarrow$  use cross-entropy
- Encoder maps word sequence to vector, decoder maps to translation via softmax distribution



# NMT via Encoder-Decoder Architecture

## Translating

Recurrent  
Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

Deep RNNs

LSTMs

GRUs

NLP

Embeddings

Architecture

Bidirectional

Attention

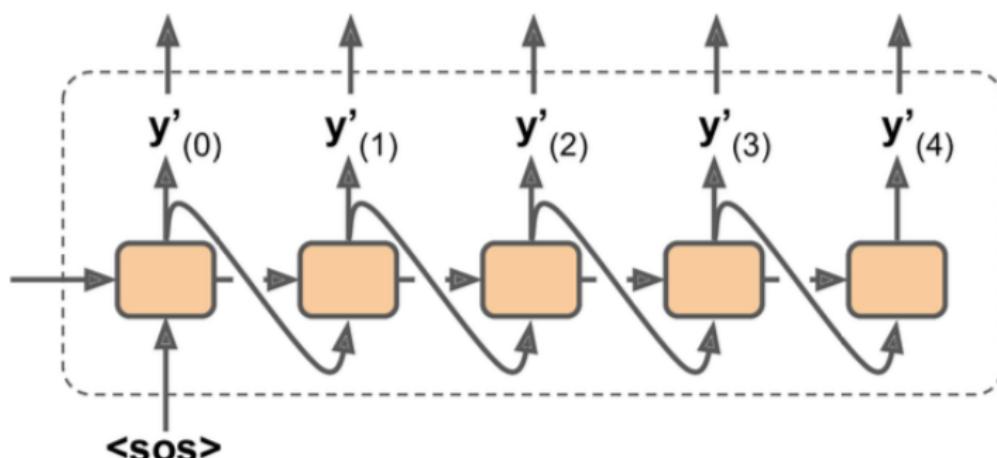
Transformer

GPT

BERT

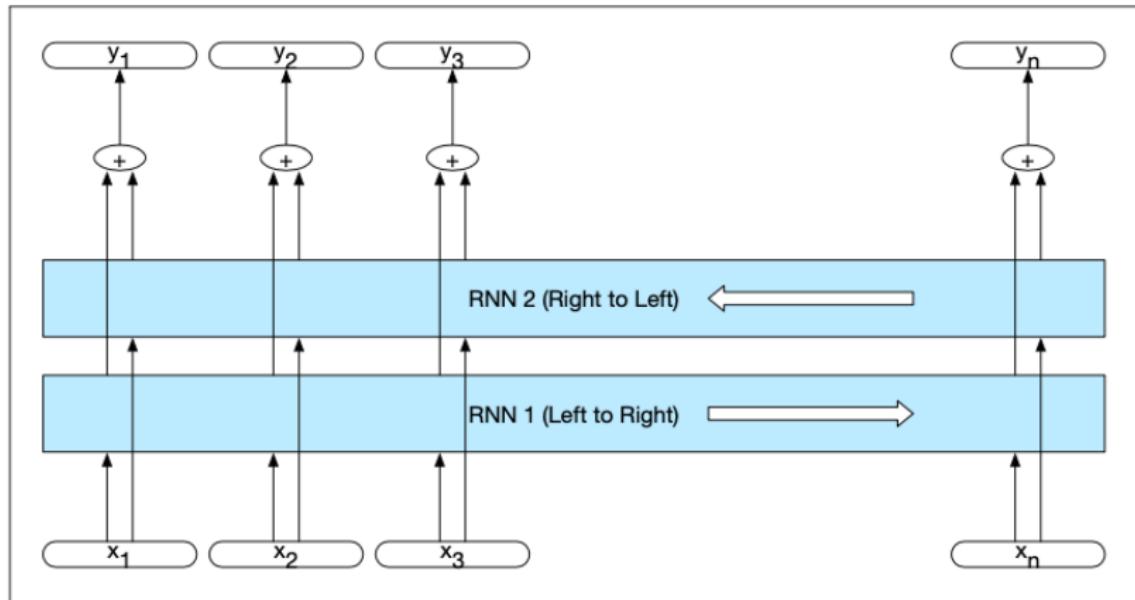
Summary

- After training, do translation by feeding previous translated word  $y'_{(t-1)}$  to decoder



# Bidirectional RNNs

- In addition to working from start to end, can simultaneously process from end towards start
- Two separate RNNs run in opposite directions and outputs are concatenated, added, etc.



Recurrent  
Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

Deep RNNs

LSTMs

GRUs

NLP

Embeddings

Architecture

Bidirectional

Attention

Transformer

GPT

BERT

Summary

# Attention Mechanisms

Recurrent  
Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

Deep RNNs

LSTMs

GRUs

NLP

Embeddings

Architecture

Bidirectional

Attention

Transformer

GPT

BERT

Summary

- Encoder-decoder works through an **embedded space** like an autoencoder, so can represent the entire input as an embedded vector prior to decoding
- Issue:** Need to ensure that the context vector fed into decoder is sufficiently large in dimension to represent context required
- Can address this representation problem via **attention mechanism**
  - Encodes input sequence into a **vector sequence** rather than single vector
  - As it decodes translation, decoder focuses on relevant subset of the vectors

# Attention Mechanisms

Recurrent  
Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

Deep RNNs

LSTMs

GRUs

NLP

Embeddings

Architecture

Bidirectional

Attention

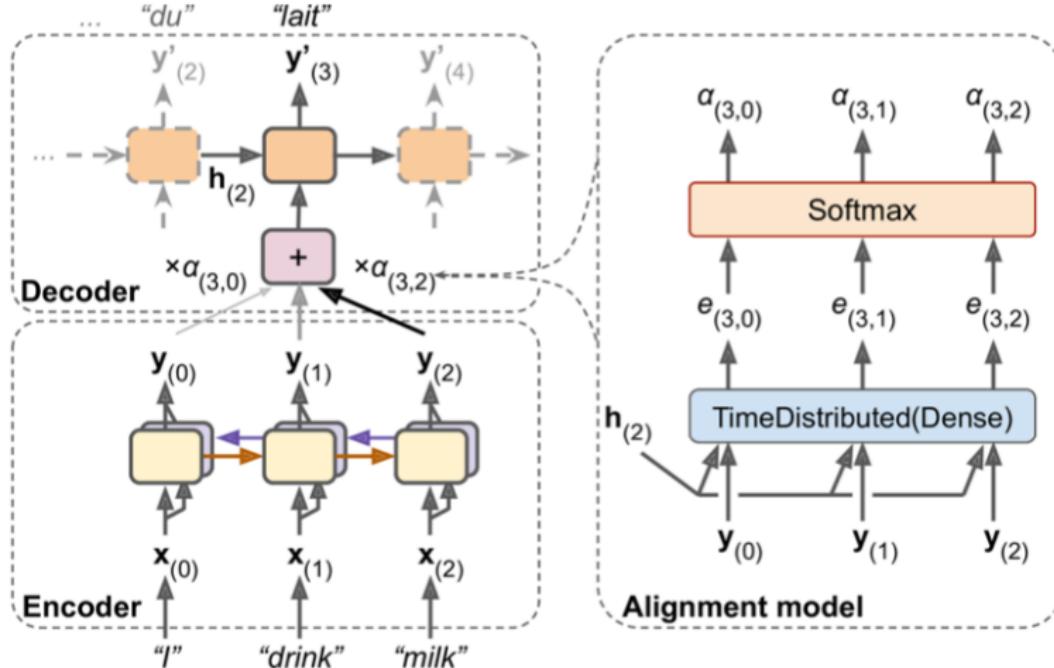
Transformer

GPT

BERT

Summary

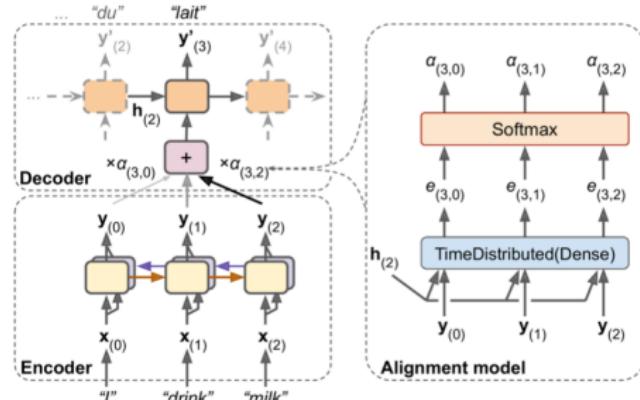
- At time step  $t$ , decoder weights encoder output  $i$  by  $\alpha_{(t,i)}$
- E.g.,  $t = 3$  has  $\alpha_{(3,2)} > \alpha_{(3,0)}, \alpha_{(3,1)}$  since “lait” depends more on “milk” than other words



# Attention Mechanisms

## Alignment Model

- Weights  $\alpha_{(t,i)}$  come from **alignment model** trained at same time as encoder-decoder
- Input is encoder outputs  $(y_{(0)}, y_{(1)}, \dots)$  and decoder hidden state  $(h_{(2)})$
- Output is **energy** score  $e_{(t,i)}$  measuring how well output  $y_{(i)}$  aligns with  $h_{(t-1)}$  (or  $h_{(t)}$ )
  - I.e., how much one should focus on word encoding  $y_{(i)}$  when generating output  $y'_{(t)}$



Recurrent Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

Deep RNNs

LSTMs

GRUs

NLP

Embeddings

Architecture

Bidirectional

Attention

Transformer

GPT

BERT

Summary

## Attention Mechanisms Calculations

## Recurrent Architectures

Stephen Scott

## Introduction

## Basic Idea

## I/O Mappings

## Training

Deep BNNs

| STM

GPUs

NLP

## Embeddings

Architecture

## Bidirectional

## Attention

## Transformer

GPT

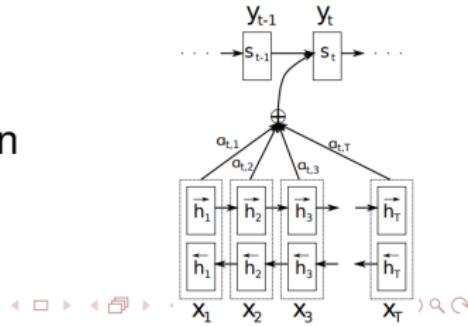
DEW

$$\tilde{\boldsymbol{h}}_{(t)} = \sum_i \alpha_{(t,i)} \boldsymbol{y}_{(i)}$$

$$\alpha_{(t,i)} = \frac{\exp(e_{(t,i)})}{\sum_{i'} \exp(e_{(t,i')})}$$

$$e_{(t,i)} = \begin{cases} \mathbf{h}_{(t)}^\top \mathbf{y}_{(i)} & \text{dot product} \\ \mathbf{h}_{(t)}^\top W \mathbf{y}_{(i)} & \text{bilinear} \\ \mathbf{v}^\top \tanh(W [\mathbf{h}_{(t)}; \mathbf{y}_{(i)}]) & \text{concatenation and activation} \end{cases}$$

$y$  vectors can also be concatenation from bidirectional RNNs



# Attention Mechanisms

Recurrent  
Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

Deep RNNs

LSTMs

GRUs

NLP

Embeddings

Architecture

Bidirectional

Attention

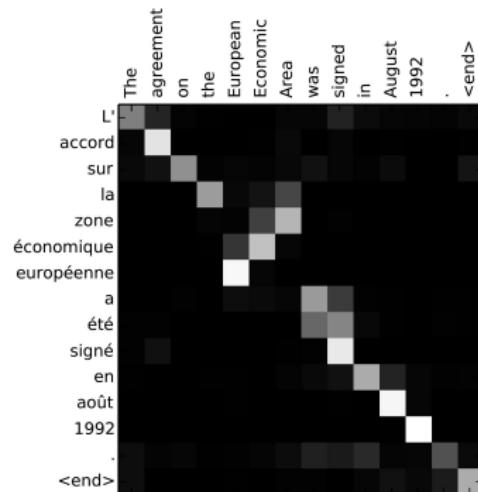
Transformer

GPT

BERT

Summary

- The  $i$ th element of **attention vector**  $\alpha_t$  tells us the probability that target output  $y'_i$  is aligned to (or translated from) input  $x_t$
- Then  $\tilde{h}_i$  is expected annotation over all annotations with probabilities  $\alpha_t$



Recurrent  
Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

Deep RNNs

LSTMs

GRUs

NLP

Embeddings

Architecture

Bidirectional

Attention

Transformer

GPT

BERT

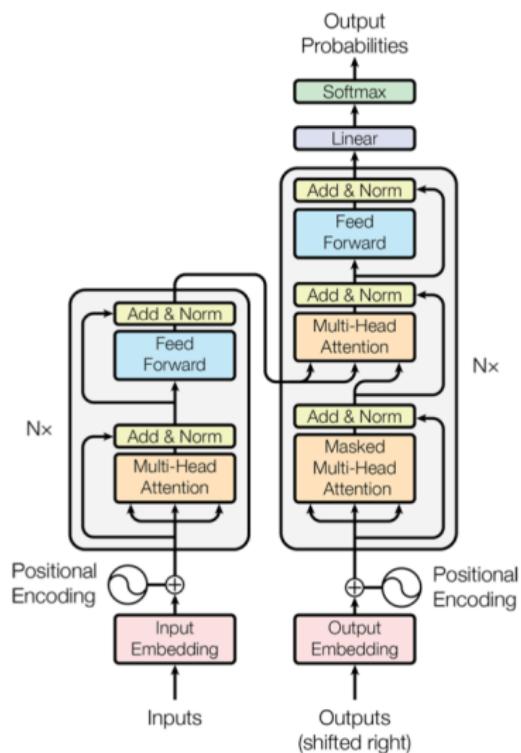
Summary

- State-of-the-art performance in NMT **without** recurrent layers; only attention
  - Feed inputs sequentially as before, but no memory
  - Faster to train
- Basis of cutting-edge approaches such as **BERT** and **GPT-3**

# Transformer Architecture (Vaswani et al., 2017)

## Overview

- Encoder on left is stack of  $N$  **multi-head attention** and feed-forward layers (no recurrent)
- Decoder on right takes outputs and encoder outputs and predicts distribution over next word
- Replace memory with **positional encoding**
- Training and inference similar to RNN



# Transformer Architecture (Vaswani et al., 2017)

## Positional Encodings

Recurrent  
Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

Deep RNNs

LSTMs

GRUs

NLP

Embeddings

Architecture

Bidirectional

Attention

Transformer

GPT

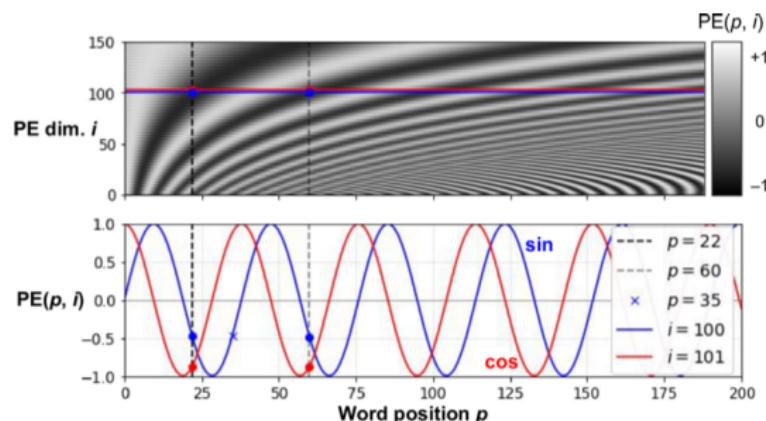
BERT

Summary

- **Positional encoding** vector added to word embedding vector (different sum for same word in different spots)
- Matrix  $P$  such that  $P_{p,i}$  is  $i$ th component of embedding for  $p$ th word in sentence ( $d = 512$ )

$$P_{p,2i} = \sin\left(p/10000^{2i/d}\right) \quad P_{p,2i+1} = \cos\left(p/10000^{2i/d}\right)$$

22nd word positional encoding is vector with  $\approx -0.488$  in 100th position and  $\approx -0.873$  in 101st position



# Transformer Architecture (Vaswani et al., 2017)

## Multi-Head Attention

Recurrent  
Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

Deep RNNs

LSTMs

GRUs

NLP

Embeddings

Architecture

Bidirectional

Attention

Transformer

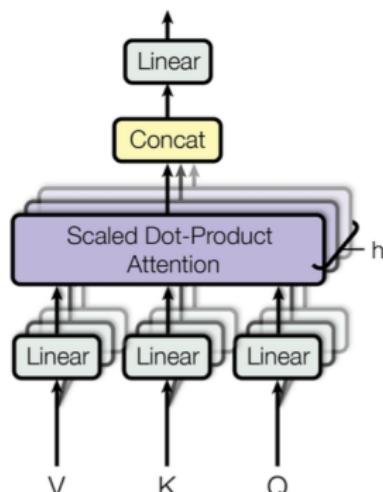
GPT

BERT

Summary

- Generally, attention can be thought of as matching a **query** to a **key**, mapping to a distribution, then scaling a set of **values**
  - E.g., from earlier,  $h_{(t)}$  is query,  $y_{(i)}$  is key

- Work with batches to use **scaled dot-product attention**:



$$\text{Attn}(Q, K, V) = \text{softmax} \left( \frac{QK^\top}{\sqrt{d_k}} \right) V$$

- $d_k$  = number of keys
- Multi-head attention** defines multiple such matrices and concatenates outputs

# Transformer Application: GPT-3

Recurrent  
Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

Deep RNNs

LSTMs

GRUs

NLP

Embeddings

Architecture

Bidirectional

Attention

Transformer

GPT

BERT

Summary

- Initially built via unsupervised pretraining of transformer modules, then fine-tuned on various NLP tasks
- 175 billion parameters (!!!)
- Such realistic generation of fake news articles, etc., creators warn of potential “harmful effects”

# Transformer Application: BERT

Recurrent  
Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

Deep RNNs

LSTMs

GRUs

NLP

Embeddings

Architecture

Bidirectional

Attention

Transformer

GPT

BERT

Summary

## Bidirectional Encoder Representations from Transformers

- An embedding that depends on the **context** of the word in the sentence
  - “We went to the river **bank**”
  - “I need to go to the **bank** to make a deposit”
- In contrast to **context-free models** (word2vec, GloVe)
- BERT trained on entire sentences by **masking** out random words and predicting them
  - ⇒ **Self-supervised training**

# Summary

Recurrent  
Architectures

Stephen Scott

Introduction

Basic Idea

I/O Mappings

Training

Deep RNNs

LSTMs

GRUs

NLP

Summary

- **Recurrent neural networks** process **sequential data**
- At time  $t$  an RNN generates output based on input at time  $t$  and its **state** (or memory) at time  $t - 1$
- RNNs are trained by **unrolling through time** the computation graph and updating with backpropagation
- **LSTMs** and **GRUs** maintain **long-term** memory, which allows for more robust models and easier training
- **Bidirectional** RNNs process from both the beginning of a sequence to the end, and from end to beginning
- Performance is further enhanced via **attention mechanisms** that allow a decoder to focus on different parts of the input at different times
- **Transformer architectures** are non-recurrent, attention-based neural approaches that are the state of the art in NLP