# Copyright Notice

These slides are distributed under the Creative Commons License.

deeplearning.ai

deeplearning.ai

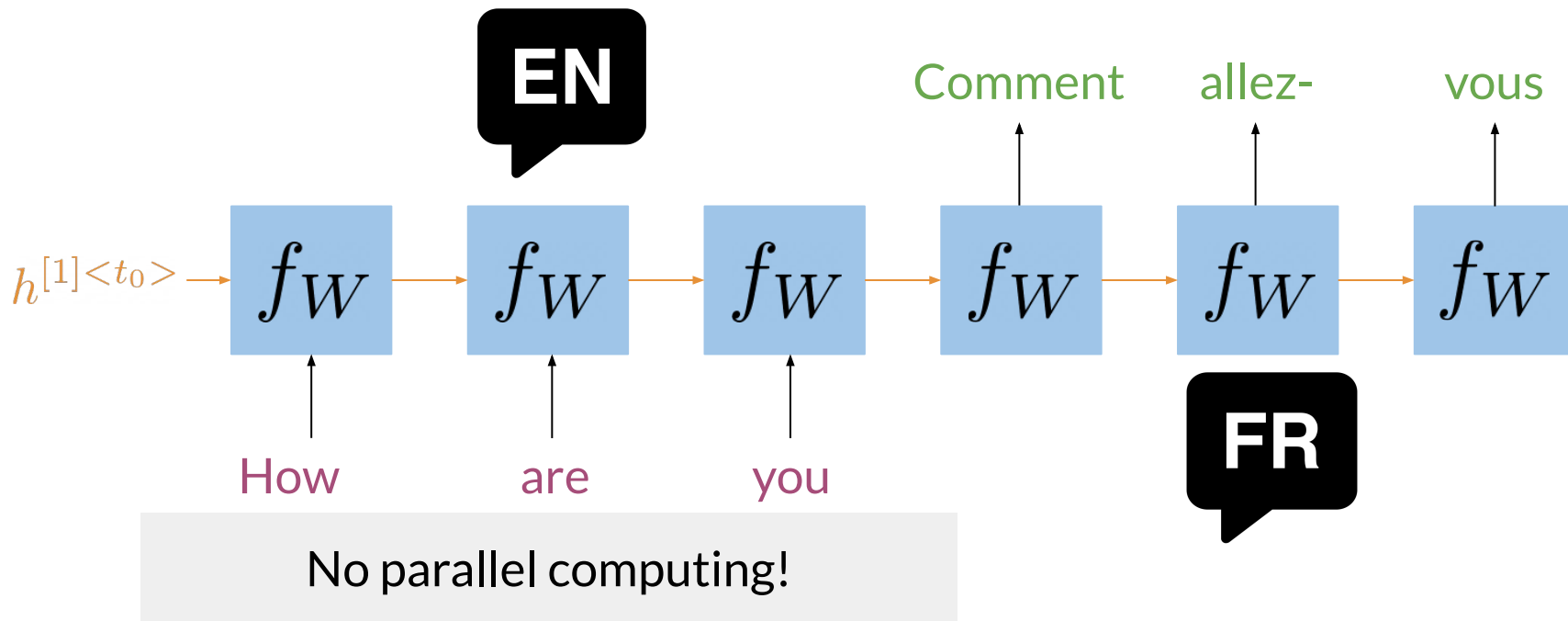# Transformers vs RNNs

# Outline

- Issues with RNNs

- Comparison with Transformers

# Neural Machine Translation



No parallel computing!

# Seq2Seq Architectures



$\hat{y}^{<t_1>}$   $\hat{y}^{<t_2>}$   $\hat{y}^{<T>}$

$h^{[L]<t_0>}$   $f_W^{[L]}$   $f_W^{[L]}$   $f_W^{[L]}$

Loss of information

T sequential steps

Vanishing gradient

$h^{[1]<t_0>}$   $f_W^{[1]}$   $f_W^{[1]}$   $f_W^{[1]}$

$x^{<t_1>}$   $x^{<t_2>}$   $x^{<T>}$

# RNNs vs Transformer: Encoder-Decoder



It's    time    for    tea

Attention
Mechanism

C'est

c

$s_{i-1}$

Decoder

<sos>

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

Encoder

LSTMs

Transformers **don't** use RNNs, such as LSTMs or GRUs

deeplearning.ai

# The Transformer Model

**Attention Is All You Need**

**Ashish Vaswani**[*]
Google Brain
avaswani@google.com

**Noam Shazeer**[*]
Google Brain
noam@google.com

**Niki Parmar**[*]
Google Research
nikip@google.com

**Jakob Uszkoreit**[*]
Google Research
usz@google.com

**Llion Jones**[*]
Google Research
llion@google.com

**Aidan N. Gomez**[*][†]
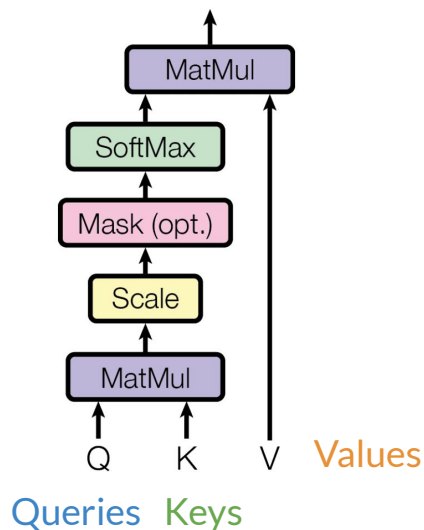University of Toronto
aidan@cs.toronto.edu

**Łukasz Kaiser**[*]
Google Brain
lukaszkaiser@google.com

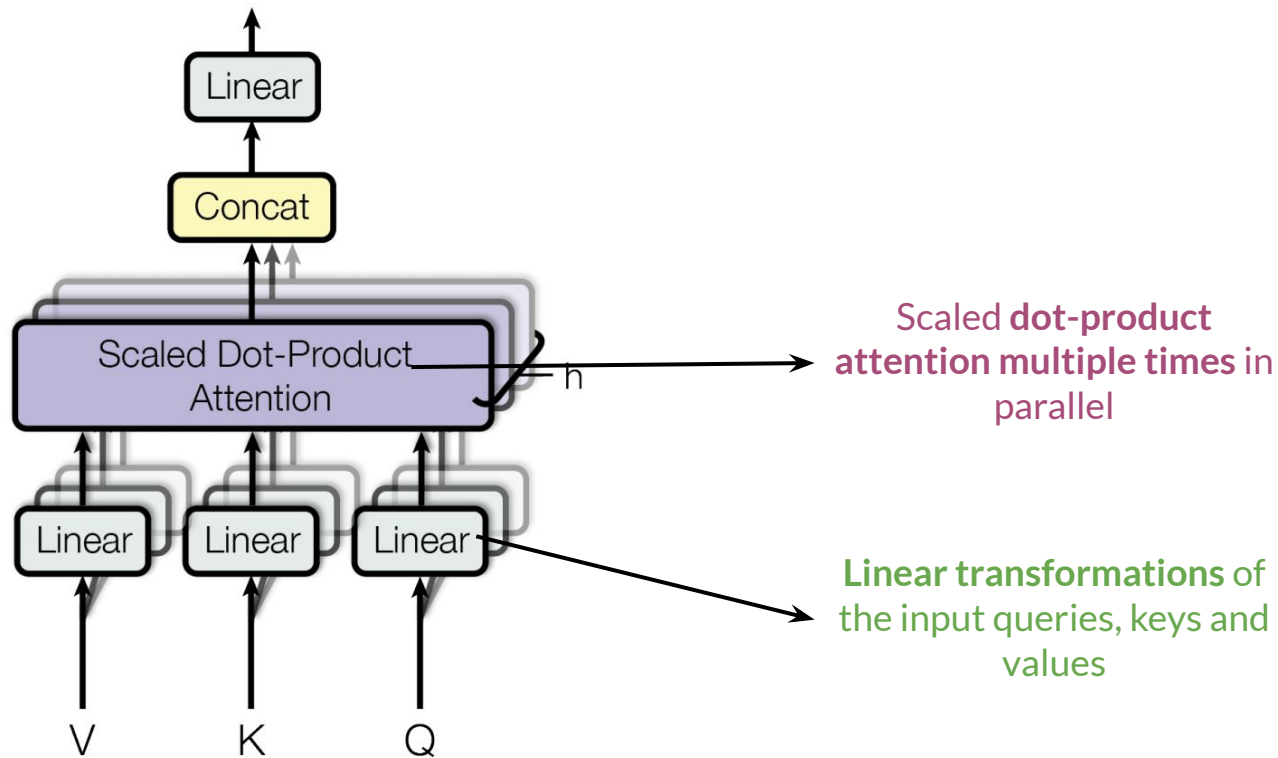**Illia Polosukhin**[*][‡]
illia.polosukhin@gmail.com

https://arxiv.org/abs/1706.03762

deeplearning.ai

# Scaled Dot-Product Attention



$$\text{softmax} \left( \frac{QK^\top}{\sqrt{d_k}} \right) V$$

Q  K  V  Values

Queries  Keys

(Vaswani et al., 2017)

deeplearning.ai

# Multi-Head Attention



Scaled **dot-product attention multiple times** in parallel

**Linear transformations** of the input queries, keys and values

# The Encoder



Nx

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Provides contextual representation of each item in the input sequence

**Self**-Attention

Every item in the input attends to every other item in the sequence

# The Decoder



**Encoder-Decoder** Attention

Every position from the decoder attents to the outputs from the encoder

**Masked Self**-Attention

Every position attends to **previous** positions

# RNNs vs Transformer: Positional Encoding

# The Transformer



Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

Encoder

Nx

Nx

Decoder

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding
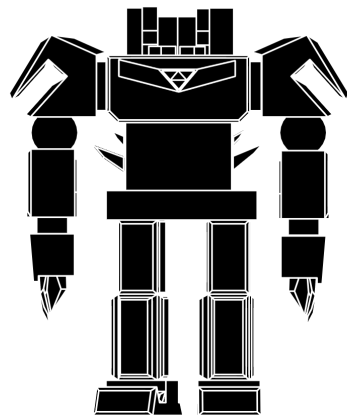
Inputs

Outputs (shifted right)

Easy to parallelize!

deeplearning.ai

# Summary

- In RNNs parallel computing is difficult to implement

- For long sequences in RNNs there is loss of information

- In RNNs there is the problem of vanishing gradient

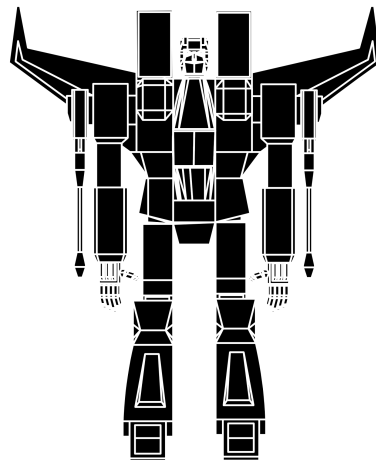- Transformers help with all of the above

deeplearning.ai
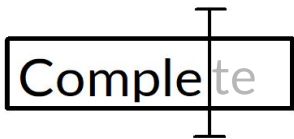
# Transformer Applications

# Outline

- Transformers applications in NLP

- Some Transformers

- Introduction to T5

# Transformer NLP applications


Text summarization


Auto-Complete



| The | wind | blows | hard |
|-----|------|-------|------|
| Article | Noun | Verb | Adjective |

Named entity recognition (NER)


Question answering (Q&A)

Translation 

Chat-bots 

Other NLP tasks

Sentiment Analysis
Market Intelligence
Text Classification
Character Recognition
Spell Checking

# State of the Art Transformers

**Radford, A., et al. (2018)**
**Open AI**

**GPT-2**: Generative Pre-training for Transformer

**Devlin, J., et al. (2018)**
**Google AI Language**

**BERT**:Bidirectional Encoder Representations from Transformers

**Colin, R., et al. (2019)**
**Google**

**T5:** Text-to-text transfer transformer

# T5: Text-To-Text Transfer Transformer



**Translate English into French:** "I am happy"

**Cola sentence:** "He bought fruits and."

***Cola** stands for "Corpus of Linguistic Acceptability"

**Cola sentence:** "He bought fruits and vegetables."

**Question:** Which volcano in Tanzania is the highest mountain in Africa?

**Translation**
**Classification**
**T5**
**Q&A**

"Je suis content"

Unacceptable

Acceptable

**Answer:** Mount Kilimanjaro

# T5: Text-To-Text Transfer Transformer

**Stsb sentence1:** "Cats and dogs are mammals." **Sentence2:** "There are four known forces in nature – gravity, electromagnetic, weak and strong."

**Stsb sentence1:** "Cats and dogs are mammals." **Sentence2:** "Cats, dogs, and cows are domesticated."

**Summarize:** "State authorities dispatched emergency crews Tuesday to survey the damage after an onslaught of severe weather in mississippi…"

**Regression**

**T5**

**Summarization**

0.0

2.6

"Six people hospitalized after a storm in Attala county"

# T5: Demo

# Summary

- Transformers are suitable for a wide range of NLP applications

- Some transformers include GPT, BERT and T5

- T5 is a powerful multi-task transformer

# Scaled Dot-Product Attention

deeplearning.ai
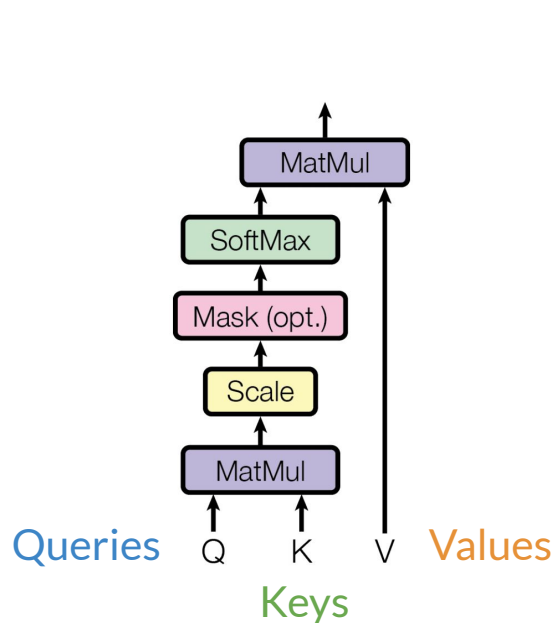
# Outline

- Revisit scaled dot product attention

- Mathematics behind Attention

# Scaled dot-product attention



Queries  Q  K  V  Values

Keys

(Vaswani et al., 2017)

Weights add up to 1

Improves performance

$$\text{softmax}\left(\frac{QK^\top}{\sqrt{d_k}}\right)V$$

Weighted sum of values V

Just two matrix multiplications and a Softmax!

# Queries, Keys and Values

# Attention Math



$$\text{softmax}\left(\frac{QK^{\top}}{\sqrt{d_k}}\right)V$$

softmax $\left(\dfrac{\begin{array}{c}Q \quad K^{\top}\end{array}}{\sqrt{d_k}}\right)$ $V$ = Context vectors for each query

Number of queries

Size of the value vector

Weight assigned to the **third key** for the **second query**

deeplearning.ai

# Summary

- Scaled Dot-product Attention is essential for Transformer

- The input to Attention are queries, keys, and values

- GPUs and TPUs

Masked
Self-Attention

# Outline

- Ways of Attention

- Overview of masked Self-Attention

# Encoder-Decoder Attention

Queries from one sentence, keys and values from another



Weight matrix

# Self-Attention

Queries, keys and values come from the **same sentence**



Weight matrix

Meaning of each word **within** the sentence

deeplearning.ai

# Masked Self-Attention

Queries, keys and values come from the **same sentence**. Queries don't attend to future positions.



Weight matrix

# Masked self-attention math



$$\mathrm{softmax}\left(\frac{Q\,K^\top}{\sqrt{d_k}} + \begin{matrix} 0 & \blacksquare & \blacksquare \\ 0 & 0 & \blacksquare \\ 0 & 0 & 0 \end{matrix}\right)V$$

■ → Minus infinity

Weights assigned to future positions are equal to 0

deeplearning.ai

# Summary

- There are three main ways of Attention: Encoder/Decoder, self-attention and masked self-attention.

- In self-attention, queries and keys come from the same sentence

- In masked self-attention queries cannot attend to the future

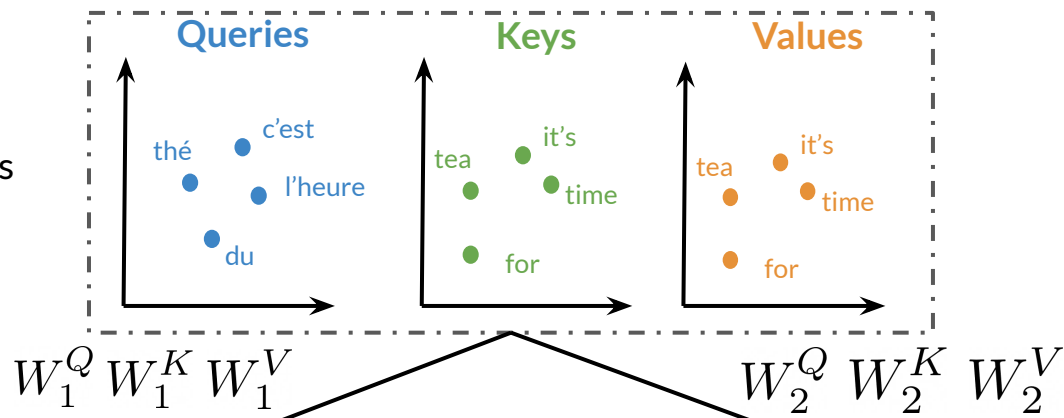Multi-head Attention

deeplearning.ai

# Outline

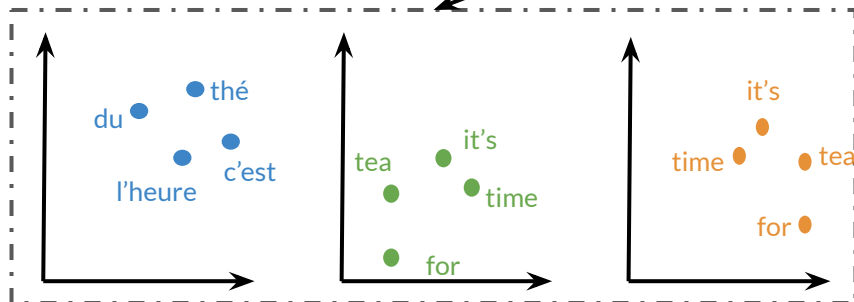- Intuition Multi-Head Attention

- Math of Multi-Head Attention

# Multi-Head Attention - Overview



Original Embeddings
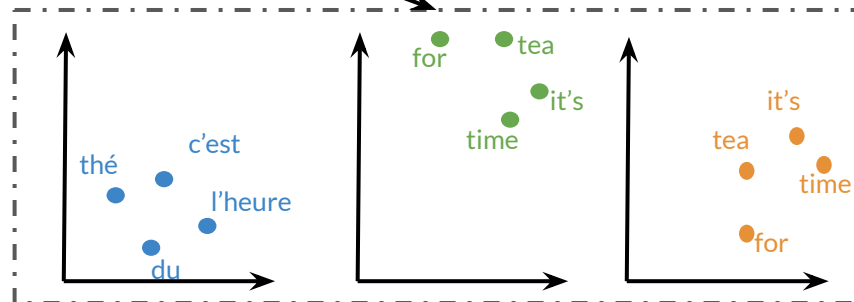
Queries    Keys    Values
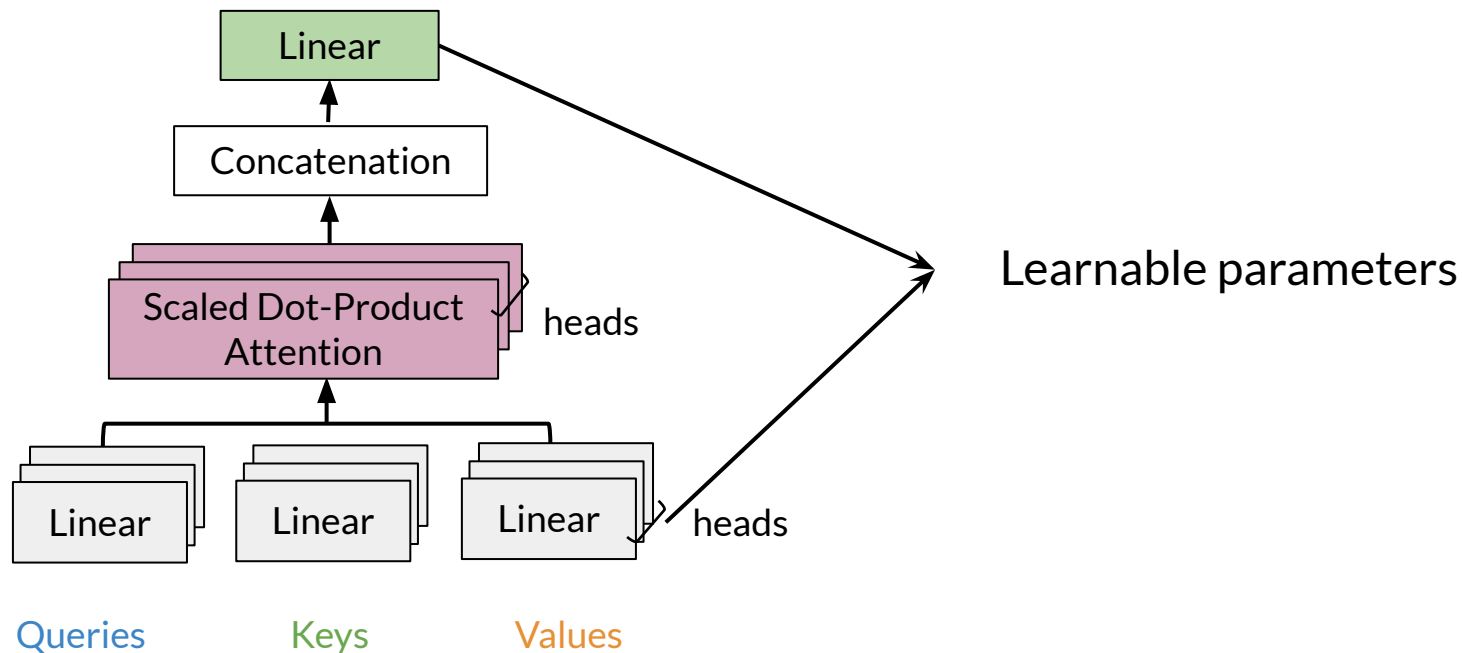
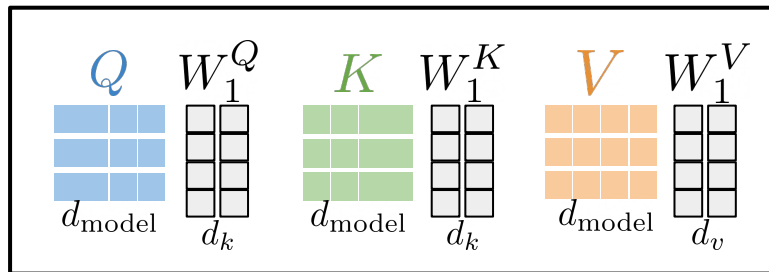$W_1^Q \ W_1^K \ W_1^V$    $W_2^Q \ W_2^K \ W_2^V$

Head 1    Head 2

# Multi-Head Attention - Overview

# Multi-Head Attention

Head 1



Head 2

Attention

Attention

Concat

$W^O$

Context vectors
for each query

$h \cdot d_v$

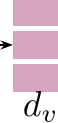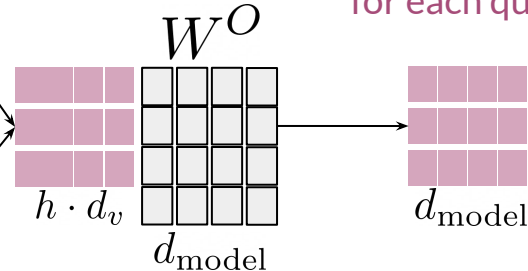$d_{\text{model}}$

$d_{\text{model}}$

$d_{\text{model}}$ : Embedding size

Usual choice of dimensions
$$d_k = d_v = d_{\text{model}}/h$$

# Summary

- Multi-Headed models attend to information from different representations

- Parallel computations
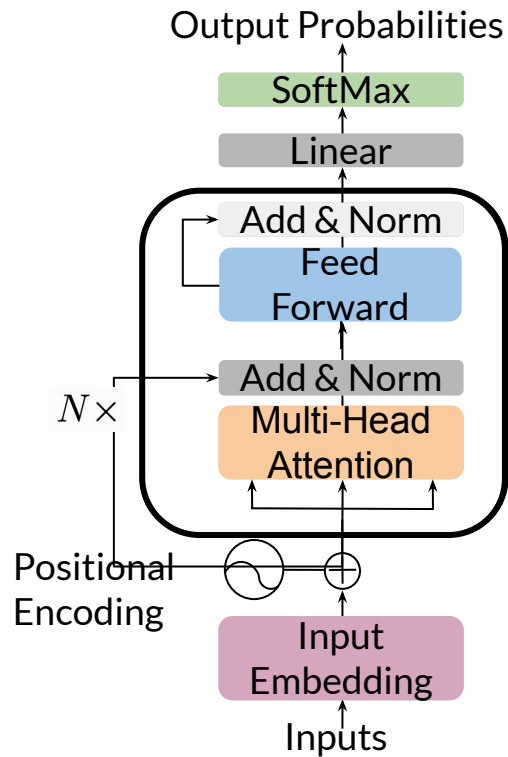
- Similar computational cost to single-head attention

# Outline

- Overview of Transformer decoder

- Implementation (decoder and feed-forward block)

# Transformer decoder

Output Probabilities

SoftMax

Linear

Add & Norm

Feed Forward

Add & Norm

$N\times$

Multi-Head Attention

Positional Encoding

Input Embedding

Inputs
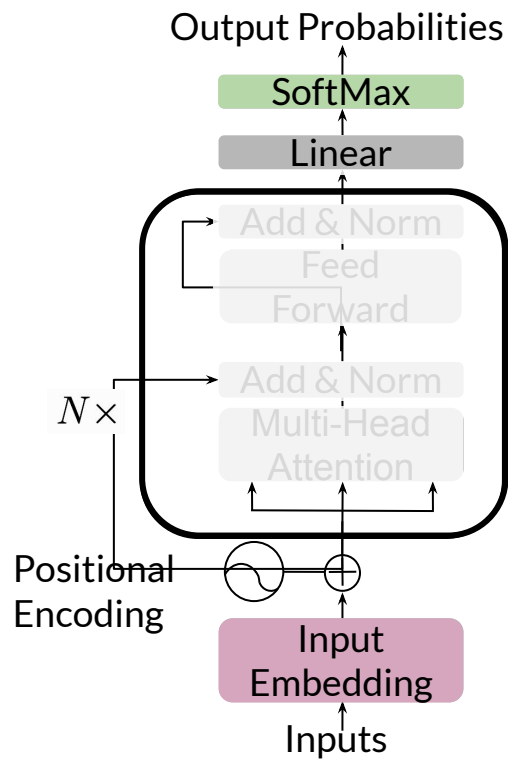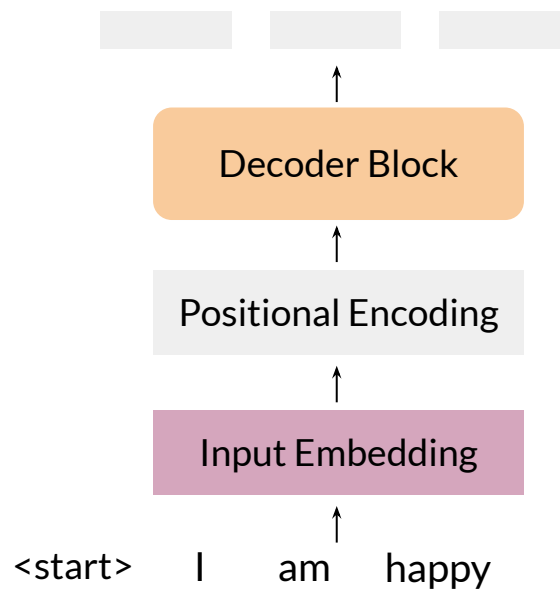
## Overview

- input: sentence or paragraph
  - we predict the next word
- sentence gets embedded, add positional encoding
  - (vectors representing $\{0, 1, 2, \ldots, K\}$)
- multi-head attention looks at previous words
- feed-forward layer with ReLU
  - that's where most parameters are!
- residual connection with layer normalization
- repeat N times
- dense layer and softmax for output

# Transformer decoder

## Output Probabilities

| SoftMax |
| Linear |

Add & Norm
Feed Forward

Add & Norm
Multi-Head Attention

$N\times$

Positional Encoding

Input Embedding

Inputs

## Explanation

Decoder Block

Positional Encoding

Input Embedding

<start>    I    am    happy

deeplearning.ai

# The Transformer decoder



Decoder Block

Output Vector

Positional input embedding

Decoder Block

Add & Norm

Feed Forward | Feed Forward | Feed Forward

LayerNorm ( [    ] + [    ] )

Output Vector

Multi-Head Attention

Positional input embedding

deeplearning.ai

# The Transformer decoder



Output Probabilities

SoftMax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

$N\times$

Positional Encoding

Input Embedding

Inputs

## Feed forward layer



Feed Forward (ReLu)

Feed Forward (ReLu)

Self Attention

# The Transformer decoder

Output Probabilities

SoftMax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

$N\times$

Positional Encoding

Input Embedding

Inputs

## Feed forward layer

Feed Forward (ReLu)

Feed Forward (ReLu)

Self Attention

# Summary

- Transformer decoder mainly consists of three layers

- Decoder and feed-forward blocks are the core of this model code

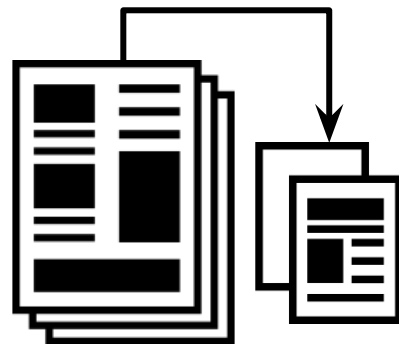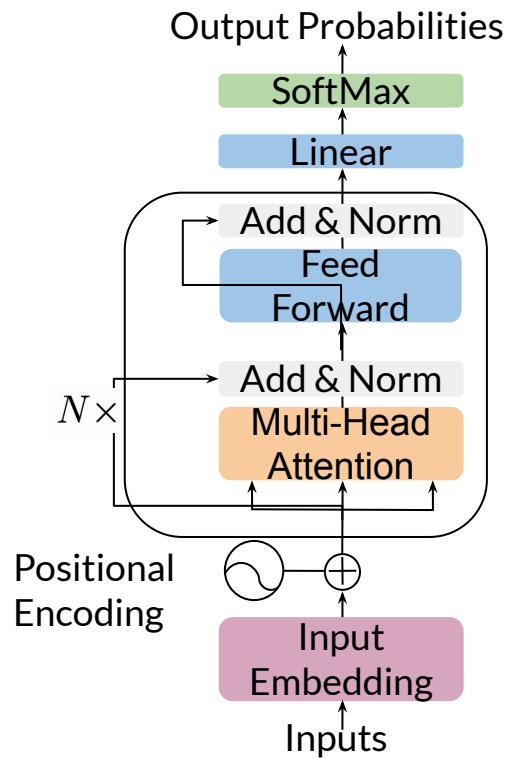- It also includes a module to calculate the cross-entropy loss
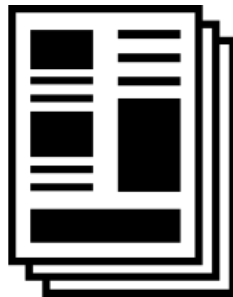
# Outline

- Overview of Transformer summarizer

- Technical details for data processing

- Inference with a Language Model
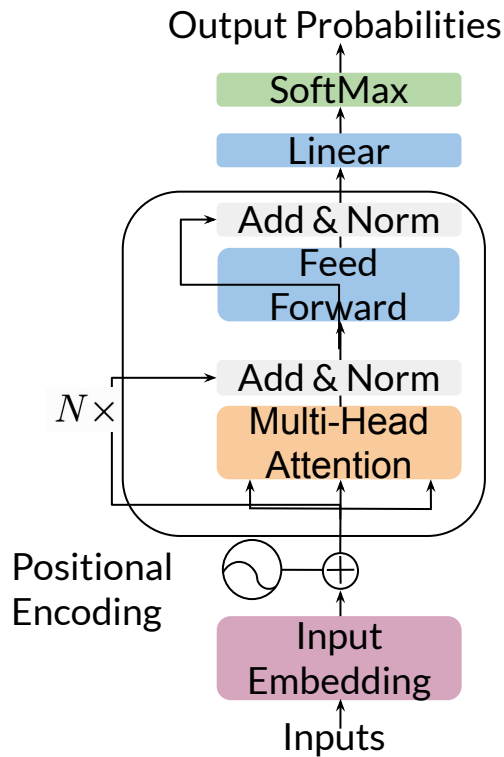
# Transformer for summarization

Output Probabilities

SoftMax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

$N\times$

Positional Encoding

Input Embedding

Inputs

Input

Output: Summary

# Technical details for data processing



Output Probabilities

SoftMax

Linear

Add & Norm

Feed Forward

Add & Norm

$N\times$

Multi-Head Attention

Positional Encoding

Input Embedding

Inputs

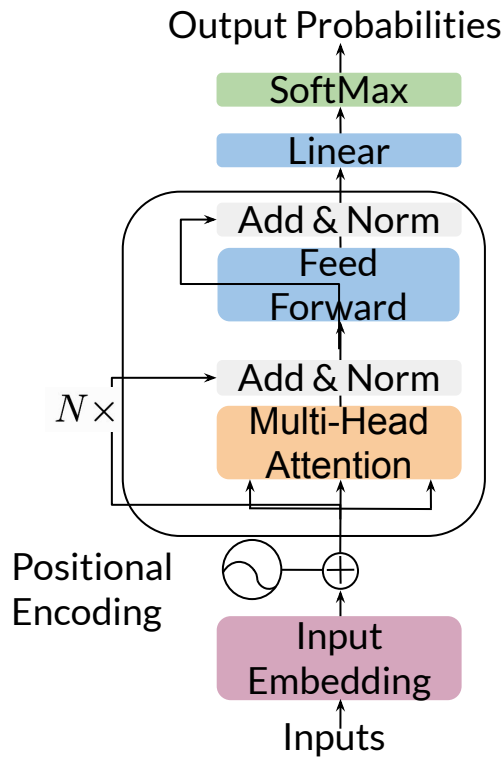**Model Input:**

```
ARTICLE TEXT <EOS> SUMMARY <EOS> <pad> …
```

**Tokenized version:**

```
[2,3,5,2,1,3,4,7,8,2,5,1,2,3,6,2,1,0,0]
```

Loss weights: 0s until the first `<EOS>` and then 1 on the start of the summary.

# Cost function

Output Probabilities

SoftMax

Linear

Add & Norm

Feed Forward

Add & Norm

$N\times$

Multi-Head Attention

Positional Encoding

Input Embedding

Inputs

## Cross entropy loss

$$J = -\frac{1}{m}\sum_{j}^{m}\sum_{i}^{K} y_j^i \log \hat{y}_j^i$$

$j$ : over summary

$i$ : bach elements

# Inference with a Language Model

## Model input:

`[Article] <EOS>  [Summary] <EOS>`

## Inference:

- Provide: [Article] <EOS>

- Generate summary word-by-word
  - until the final <EOS>

- Pick the next word by random sampling
  - each time you get a different summary!

# Summary

- For summarization, a  weighted loss function is optimized

- Transformer Decoder summarizes predicting the next word using

- The transformer uses tokenized versions of the input