# Transformer Architecture

CSCI 601-471/671 (NLP: Self-Supervised Models)

# RNNs, Back to the Cons

- While RNNs in theory can represent long sequences, they quickly <span style="color:red">forget</span> portions of the input.

- Vanishing/exploding gradients

- Difficult to parallelize

- The alternative solution we will see: Transformers!

# Language Models: History Recap

- Probabilistic n-gram models of text generation [Jelinek+ 1980's, ...]
  - Applications: Speech Recognition, Machine Translation

- Statistical or shallow neural LMs (late 90's – mid 00's) [Bengio+ 2001, ...]

- Recurrent neural nets (2010s)

- Pre-training deep neural language models (2017's onward):
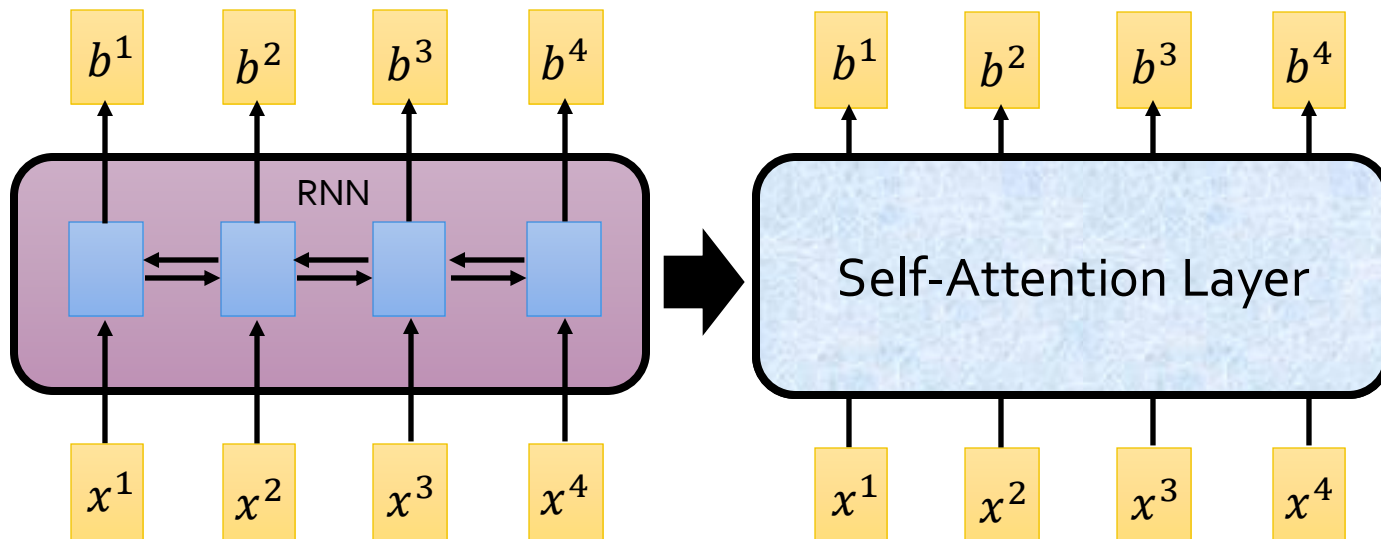  - Many models based on: Self-Attention

# Chapter Plan

1. Self-Attention module
2. Transformer architecture
3. Computation/space cost
4. Thinking about Transformer implementation

**Chapter goal** — getting very comfortable with nuances involved in Transformers.
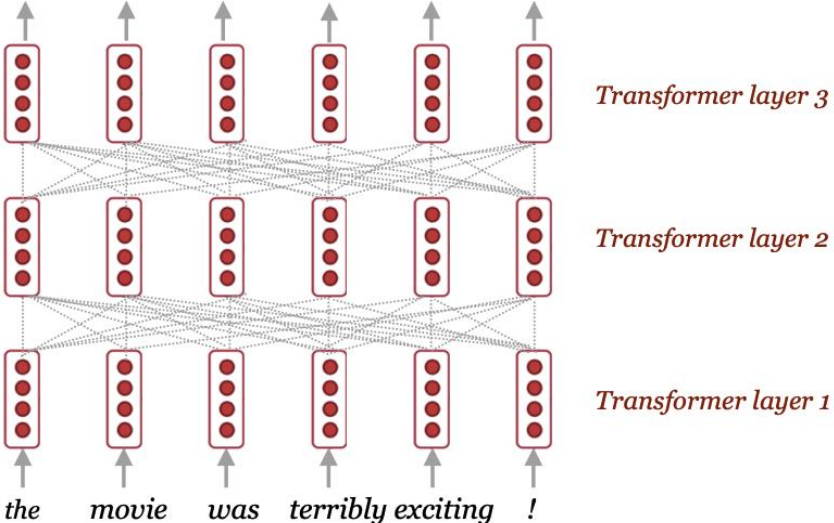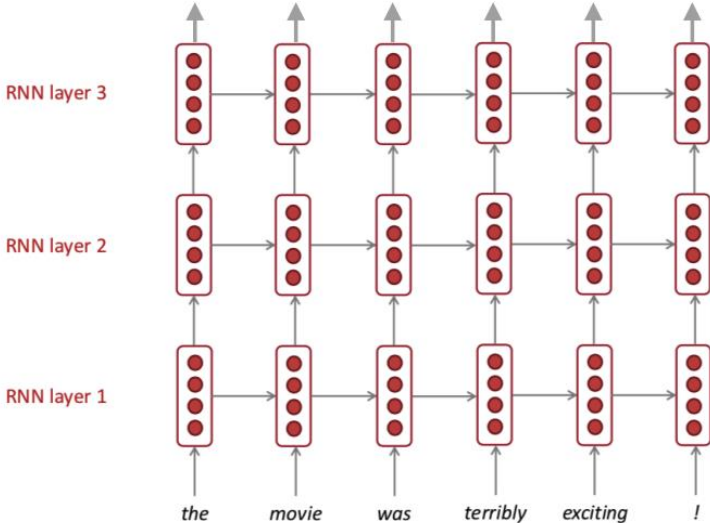
# Self-Attention Module

# Self-Attention

- $b^i$ is obtained based on the whole input sequence.
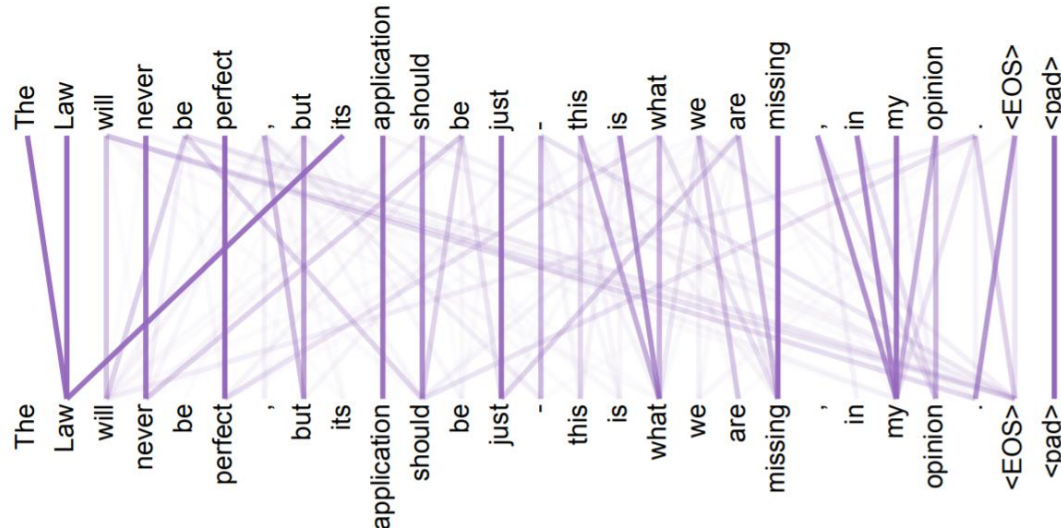- can be parallelly computed.



Idea: replace any thing done by RNN with self-attention.

"Neural machine translation by jointly learning to align and translate" Bahdanau etl. 2014;
"Attention is All You Need" Vaswani et al. 2017

# RNN vs Transformer



RNN layer 3

RNN layer 2

RNN layer 1

the    movie    was    terribly    exciting    !

Transformer layer 3

Transformer layer 2

Transformer layer 1

the    movie    was    terribly exciting    !

# Attention

- Core idea: build a mechanism to focus ("attend") on a particular part of the context.



[Attention Is All You Need, Vaswani et al. 2017]

8

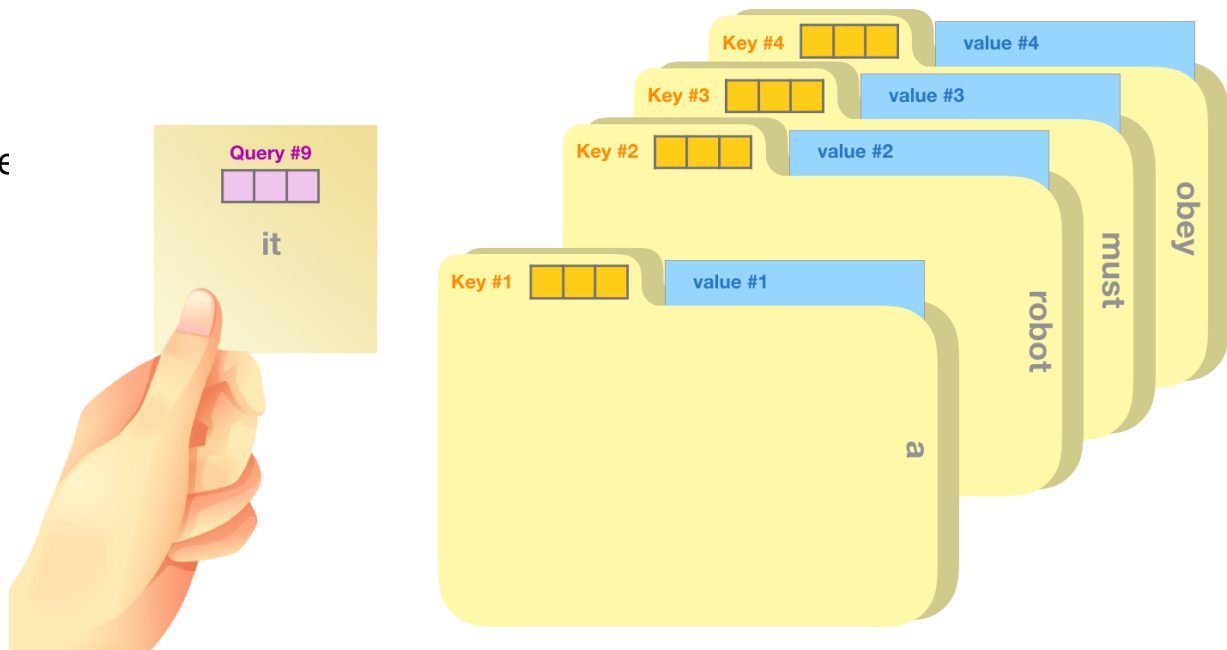# Defining Self-Attention

- Terminology:
  - Query: to match others
  - Key: to be matched
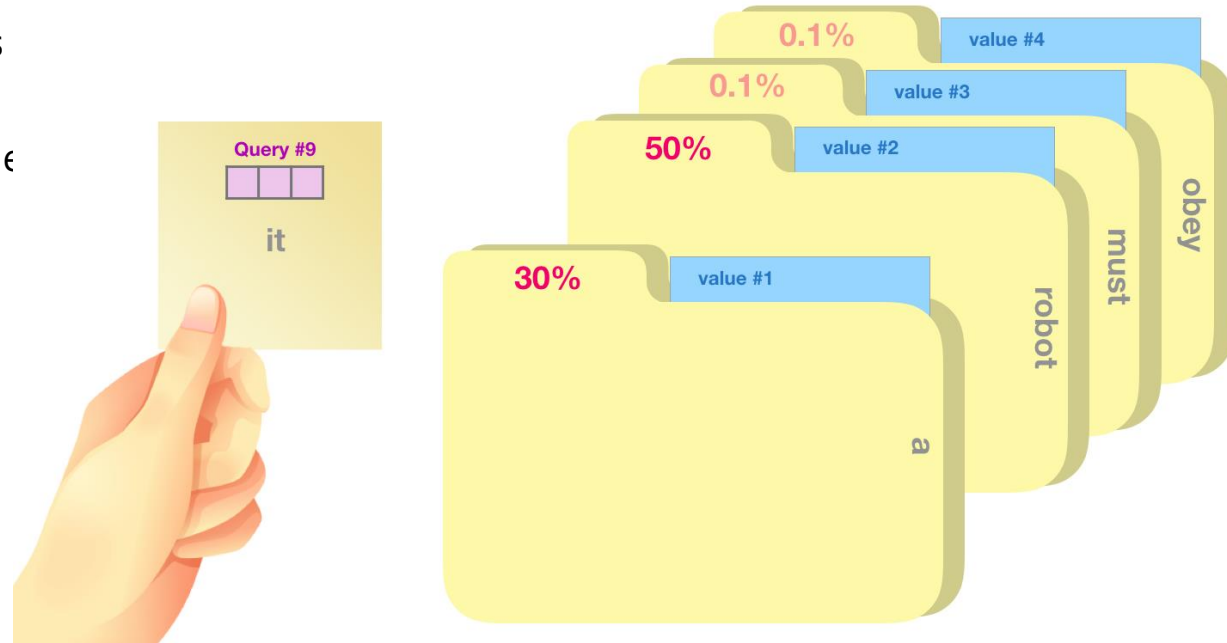  - Value: information to be extracted

[Vaswani et al. 2017: https://arxiv.org/abs/1706.03762]

# Defining Self-Attention

An analogy ....

- Terminology:
  - o Query: to match others
  - o Key: to be matched
  - o Value: information to be

Query #9

it

Key #4   value #4

Key #3   value #3

Key #2   value #2

Key #1   value #1

obey

must

robot

a

[Vaswani et al. 2017: https://arxiv.org/abs/1706.03762]

# Defining Self-Attention

- Terminology:
  - Query: to match others
  - Key: to be matched
  - Value: information to be



[Vaswani et al. 2017: https://arxiv.org/abs/1706.03762]
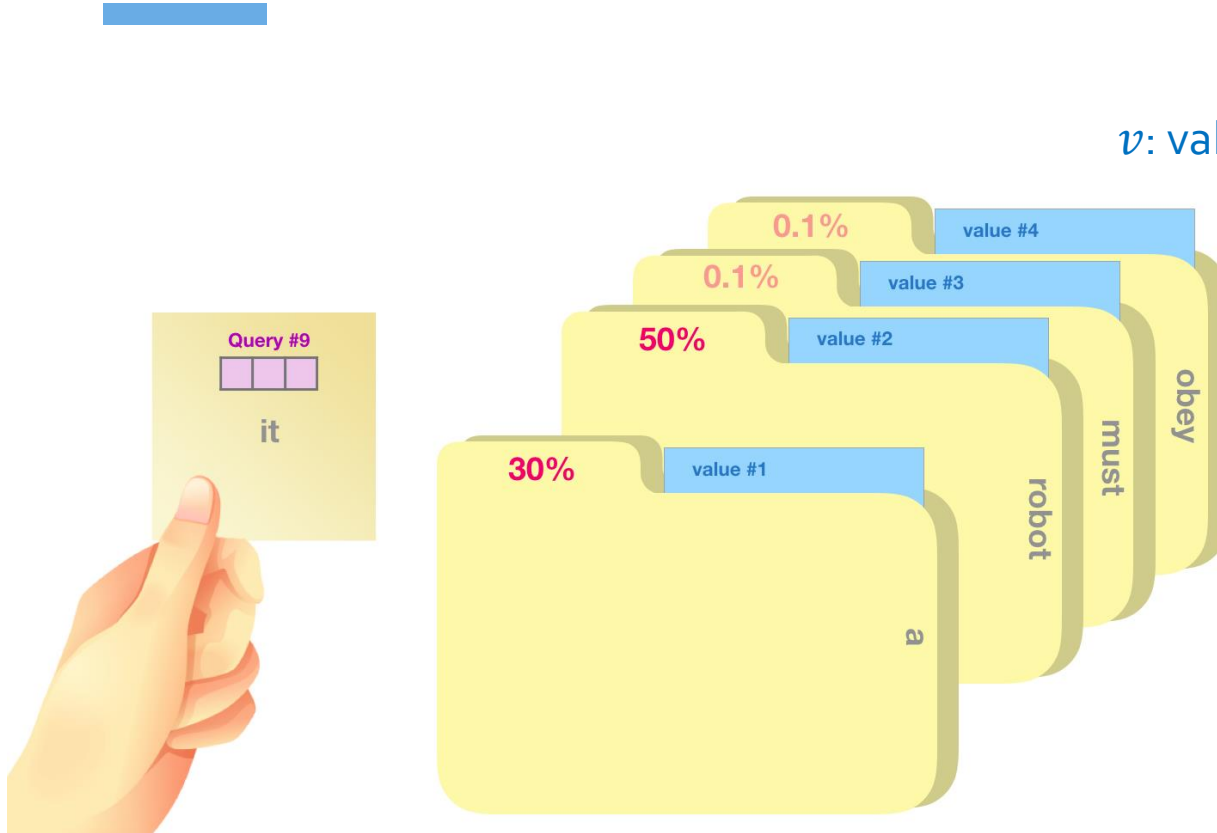
$q$: query (to match others)
$$q_i = W^q x_i$$

$k$: key (to be matched)
$$k_i = W^k x_i$$

$v$: value (information to be extracted)
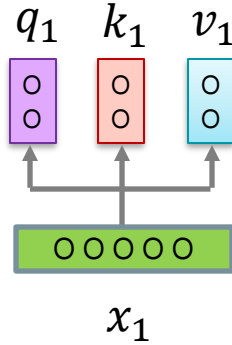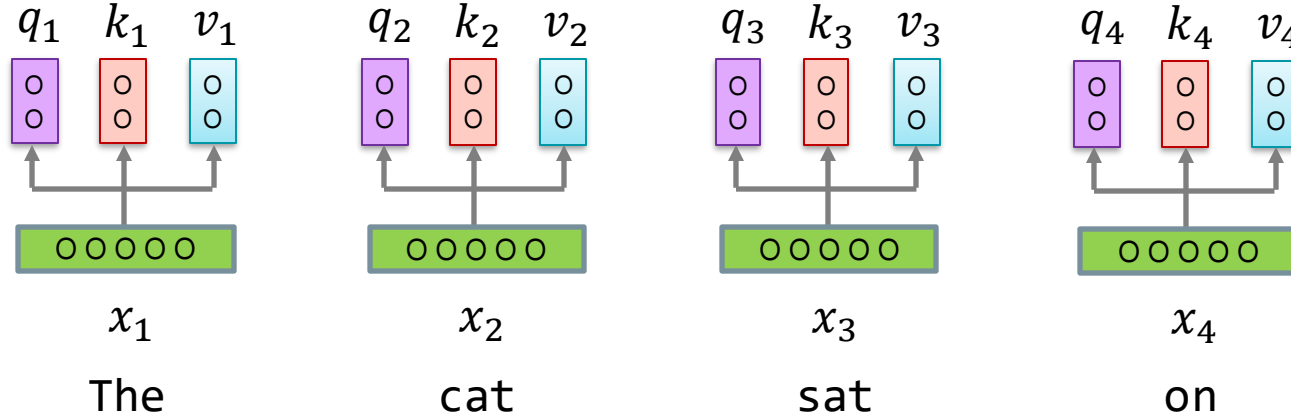$$v_i = W^v x_i$$

$q$: query (to match others)
$$q_i = W^q x_i$$

$k$: key (to be matched)
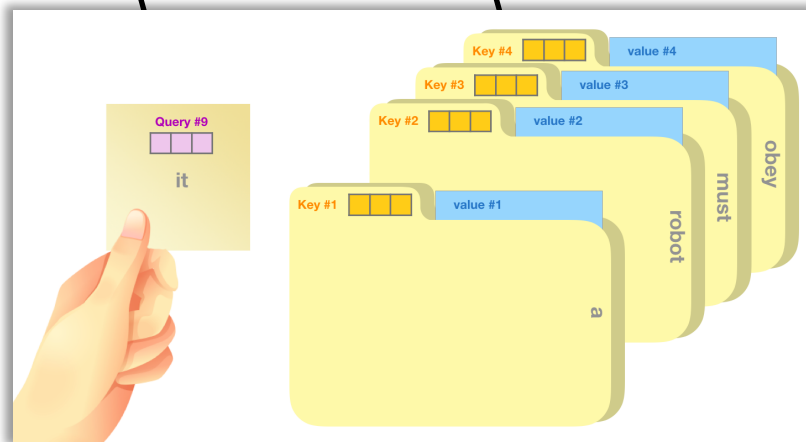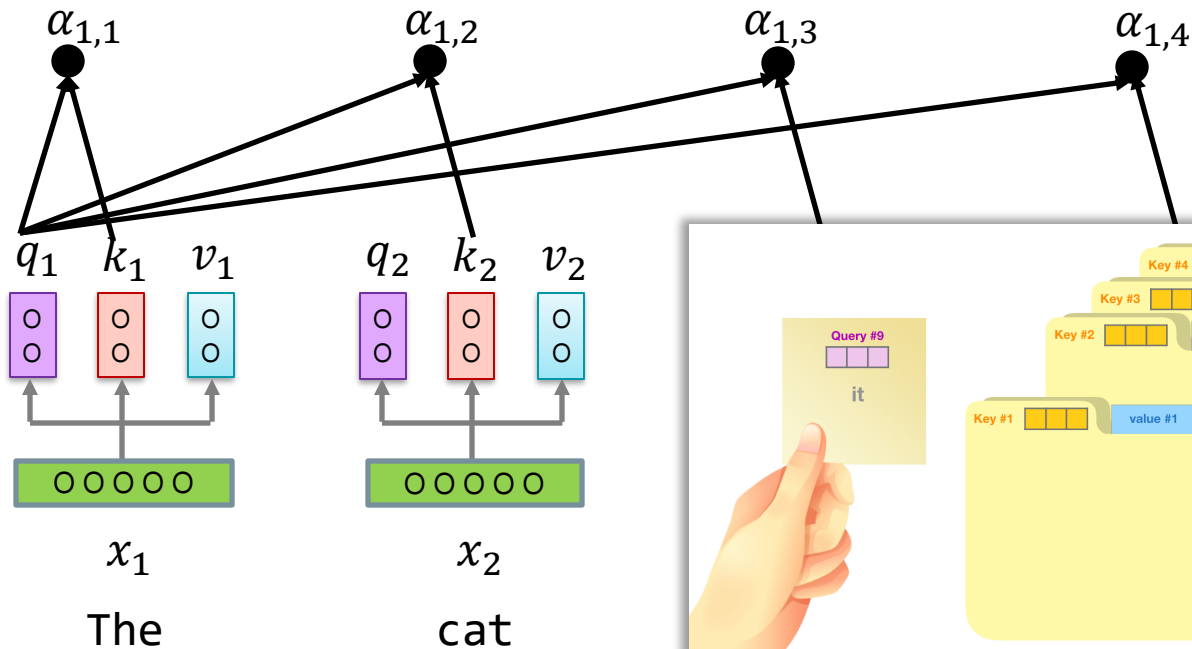$$k_i = W^k x_i$$

$v$: value (information to be extracted)
$$v_i = W^v x_i$$

$q_1 \quad k_1 \quad v_1$

$x_1$

The

$q$: query (to match others)
$$q_i = W^q x_i$$

$k$: key (to be matched)
$$k_i = W^k x_i$$

$v$: value (information to be extracted)
$$v_i = W^v x_i$$

$q_1 \quad k_1 \quad v_1 \qquad q_2 \quad k_2 \quad v_2 \qquad q_3 \quad k_3 \quad v_3 \qquad q_4 \quad k_4 \quad v_4$

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

The             cat             sat             on

15

$$\sigma(z)_i = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$

$\hat{\alpha}_{1,1}$  $\hat{\alpha}_{1,2}$  $\hat{\alpha}_{1,3}$  $\hat{\alpha}_{1,4}$

Softmax

How much should "The" attend to other positions?

$\alpha_{1,1}$  $\alpha_{1,2}$  $\alpha_{1,3}$  $\alpha_{1,4}$

$q_1$  $k_1$  $v_1$    $q_2$  $k_2$  $v_2$    $q_3$  $k_3$  $v_3$    $q_4$  $k_4$  $v_4$

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

The    cat    sat    on

$$b^1 = \sum_i \hat{\alpha}_{1,i} v^i$$

Representation of "The" given the attention weights

$\hat{\alpha}_{1,1}$  $\hat{\alpha}_{1,2}$  $\hat{\alpha}_{1,3}$  $\hat{\alpha}_{1,4}$

Softmax

$\alpha_{1,1}$  $\alpha_{1,2}$  $\alpha_{1,3}$  $\alpha_{1,4}$

$q_1$  $k_1$  $v_1$    $q_2$  $k_2$  $v_2$    $q_3$  $k_3$  $v_3$    $q_4$  $k_4$  $v_4$

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

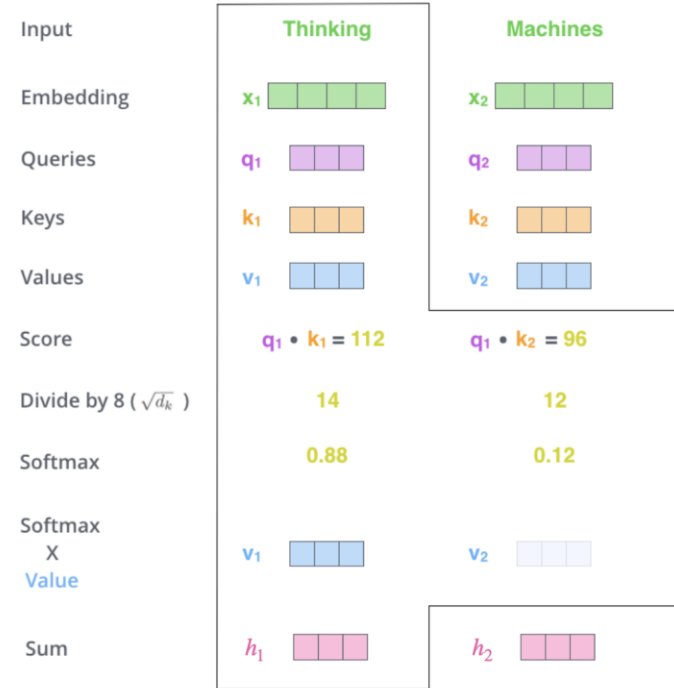The    cat    sat    on

# Question

- What would be the output vector for the word "Thinking"?

    (a)  $0.5\mathbf{v}_1 + 0.5\mathbf{v}_2$

    (b)  $0.54\mathbf{v}_1 + 0.46\mathbf{v}_2$

    (c)  $0.88\mathbf{v}_1 + 0.12\mathbf{v}_2$

    (d)  $0.12\mathbf{v}_1 + 0.88\mathbf{v}_2$

| | Thinking | Machines |
|---|---|---|
| Input | | |
| Embedding | $x_1$ | $x_2$ |
| Queries | $q_1$ | $q_2$ |
| Keys | $k_1$ | $k_2$ |
| Values | $v_1$ | $v_2$ |
| Score | $q_1 \cdot k_1 = 112$ | $q_1 \cdot k_2 = 96$ |
| Divide by 8 ( $\sqrt{d_k}$ ) | 14 | 12 |
| Softmax | 0.88 | 0.12 |
| Softmax X Value | $v_1$ | $v_2$ |
| Sum | $h_1$ | $h_2$ |

JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

# Self-Attention: Matrix Notation

$X \in \mathbb{R}^{n \times d_1}$      (n = input length)

$Q = XW^Q$     $K = XW^K$     $V = XW^V$

$W^Q \in \mathbb{R}^{d_1 \times d_q}, W^K \in \mathbb{R}^{d_1 \times d_k}, W^V \in \mathbb{R}^{d_1 \times d_v}$

$n \times d_q$      $d_k \times n$

$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

$n \times$

Q: What is this softmax operation?

$$\text{softmax}\left( \frac{Q \times K^T}{\sqrt{d_k}} \right) V$$

$= H$

JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

# Self-Attention

- Can write it in matrix form:

- Given input **x**:

$$Q = \mathbf{W}^q \mathbf{x}$$
$$K = \mathbf{W}^k \mathbf{x}$$
$$V = \mathbf{W}^v \mathbf{x}$$

$$\text{Attention}(\mathbf{x}) = \text{softmax}\left(\frac{QK^{\mathrm{T}}}{\sqrt{d}}\right) V$$

**hardmaru**
@hardmaru

The most important formula in deep learning after 2018

**Self-Attention**

**What is self-attention?** Self-attention calculates a weighted average of feature representations with the weight proportional to a similarity score between pairs of representations. Formally, an input sequence of $n$ tokens of dimensions $d$, $X \in \mathbf{R}^{n \times d}$, is projected using three matrices $W_Q \in \mathbf{R}^{d \times d_q}$, $W_K \in \mathbf{R}^{d \times d_k}$, and $W_V \in \mathbf{R}^{d \times d_v}$ to extract feature representations $Q$, $K$, and $V$, referred to as query, key, and value respectively with $d_k = d_q$. The outputs $Q$, $K$, $V$ are computed as

$$Q = XW_Q, \quad K = XW_K, \quad V = XW_V. \quad (1)$$

So, self-attention can be written as,

$$S = D(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_q}}\right)V, \quad (2)$$
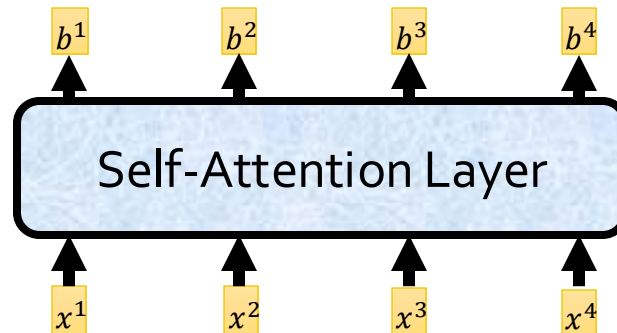
where softmax denotes a *row-wise* softmax normalization function. Thus, each element in $S$ depends on all other elements in the same row.

9:08 PM · Feb 9, 2021 · Twitter Web App

**553** Retweets      **42** Quote Tweets      **3,338** Likes

# Self-Attention: Back to Big Picture

- Attention is a powerful mechanism to create context-aware representations
- A way to focus on select parts of the input

$$b^1 \quad b^2 \quad b^3 \quad b^4$$

Self-Attention Layer

$$x^1 \quad x^2 \quad x^3 \quad x^4$$

- Better at maintaining long-distance dependencies in the context.

[Attention Is All You Need, Vaswani et al. 2017]

# Computational and Space Complexity

- The attention function:

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

- $\dim(QK^T) = N^2 \rightarrow O(N^2 d_k)$ time complexity to calculate $QK$.

- Attention matrix $\dim\left(\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)\right) = N \times N$

  o Storing the attention matrix for each head $\rightarrow O(N^2 h)$.

- If $N \gg d_k, h$, the time and space complexity is $O(N^2)$.
  o Scalability, resource consumption, adoption, etc.

[Slide credit: CS886 at UWaterloo]

# Computational and Space Complexity (2)

| Layer Type | Complexity per Layer | Sequential Operations |
|---|---|---|
| Self-Attention | $O(n^2 \cdot d)$ | $O(1)$ |
| Recurrent | $O(n \cdot d^2)$ | $O(n)$ |

- $n$ = sequence length, $d$ = hidden dimension
- Quadratic complexity, but:
  - O(1) sequential operations (not linear like in RNN)

- Can be efficiently parallelized

[Attention Is All You Need, Vaswani et al. 2017]

# **Multi-Headed** Self-Attention

- Multiple parallel attention layers.
  - Each attention layer has its own parameters.
  - Concatenate the results and run them through a linear projection.

- Main idea: Allows model to jointly attend to information from different representation subspaces (like ensembling)



Self-Attention Layer

$x^1$   $x^2$   $x^3$   $x^4$

[Attention Is All You Need, Vaswani et al. 2017]

# Multi-Headed Self-Attention



- Just concatenate all the heads and apply an output projection matrix.

$$\text{head}_i = \text{Attention}\left(\mathbf{W}_i^q \mathbf{x}, \mathbf{W}_i^k \mathbf{x}, \mathbf{W}_i^v \mathbf{x}\right)$$

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)\boldsymbol{W}^O$$

- In practice, we use a reduced dimension for each head.
  - Denote: $d = \text{hidden dimension}$, $\text{m} = \text{number of heads}$

  $$\mathbf{W}_i^q \in \mathbb{R}^{d \times \frac{d}{m}}, \qquad \mathbf{W}_i^k \in \mathbb{R}^{d \times \frac{d}{m}}, \qquad \mathbf{W}_i^v \in \mathbb{R}^{d \times \frac{d}{m}}, \qquad \boldsymbol{W}^O \in \mathbb{R}^{d \times d}$$

- The total computational cost is similar to that of single-hear attention with full dimensionality.

# Combine with FFN

- Add a feed-forward network on top it to add more expressivity.
  - This allows the model to apply another transformation to the contextual representations (or "post-process" them).
  - Usually, the dimensionality of the hidden feedforward layer is 2-8 times larger than the input dimension.

$$\text{FFN}(\mathbf{x}) = f(cW_1 + b_1)W_2 + b_2$$



Feedforward Net: Refresher

A fully-connected network of nodes and weights.

# How Do We Prevent Vanishing Gradients?

- Residual connections let the model "skip" layers
  - These connections are particularly useful for training deep networks

- Use layer normalization to stabilize the network and allow for proper gradient flow

# Putting it Together: Self-Attention Block

Given input **x**:

$$\text{out} = LN(\tilde{\boldsymbol{c}} + \boldsymbol{c}')$$
$$\tilde{\boldsymbol{c}} = \text{FFN}(\boldsymbol{c}') = f(\boldsymbol{c}'W_1 + b_1)W_2 + b_2$$

$$\boldsymbol{c}' = LN(\boldsymbol{c} + \boldsymbol{x})$$
$$\boldsymbol{c} = \text{MultiHeadedAttention}(\boldsymbol{x}; \mathbf{W}^q, \mathbf{W}^k, \mathbf{W}^v)$$

out

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

$x$: input sequence

[Attention Is All You Need, Vaswani et al. 2017]

# Summary: Self-Attention Block

- **Self-Attention:** A critical building block of modern language models.
    - The idea is to compose meanings of words weighted according some similarity notion.

- **Next:** We will combine self-attention blocks to build various architectures known as Transformer.

# Transformer

# How Do We Make it **Deep?**

- Stack more layers!

# From Representations to Prediction

- To perform prediction, add a classification head on top of the final layer of the transformer.

- This can be per token (Language modeling)

- Or can be for the entire sequence (only one token)

$$\text{out} \in \mathbb{R}^{S \times d} \quad (S: \text{Sequence length})$$

$$\text{logits} = \text{Linear}_{(d, v)}(out) = f\big(out \cdot W_V\big) \in \mathbb{R}^{S \times V}$$

$$\text{probabilies} = \text{softmax}(\text{logits}) \in \mathbb{R}^{S \times V}$$

Token Embeddings

output token probabilities (logits)

Decoder #12, Position #1 output vector

X

=

| 0.19850038 | aardvark |
| 0.7089803 | aarhus |
| 0.46333563 | aaron |
| | ... |
| | ... |
| | ... |
| | ... |
| | ... |

Pick an output token based on its probability (sample)

The

books

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

N×

Add & Norm

Multi-Head Attention

Positional Encoding

Input Embedding

Inputs

# One last wrinkle though …

$b^1 = \sum_i \hat{\alpha}_{1,i} v^i$

One issue: the model doesn't know word positions/ordering.

$\hat{\alpha}_{1,1}$  $\hat{\alpha}_{1,2}$  $\hat{\alpha}_{1,3}$  $\hat{\alpha}_{1,4}$

Softmax

$\alpha_{1,1}$  $\alpha_{1,2}$  $\alpha_{1,3}$  $\alpha_{1,4}$

$q_1$  $k_1$  $v_1$    $q_2$  $k_2$  $v_2$    $q_3$  $k_3$  $v_3$    $q_4$  $k_4$  $v_4$

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

The    cat    sat    on

# Absolute Positional Embeddings

- Why "add"? Why not, say, "concatenate and then project"?
  - "concatenate and then project" would be a more general approach with more trainable parameters.
  - In practice, "sum" works fine that
  - The intuition here is that "summing" forms point clouds of word embedding information around position embeddings unique to each position.

$p_i$ are positional embeddings

Allows model to learn relative positioning

$p_1$ $x_1$ $p_2$ $x_2$ $p_3$ $x_3$ $p_4$ $x_4$

# **Absolute Positional Embeddings**

- The idea is to create vectors that uniquely encoder each position.
- For example, consider vectors of binary values.
  - Example below shows 4-dimensional position encodings for 16 positions.

| 0 : | 0 0 0 0 | 8 :  | 1 0 0 0 |
|-----|---------|------|---------|
| 1 : | 0 0 0 1 | 9 :  | 1 0 0 1 |
| 2 : | 0 0 1 0 | 10 : | 1 0 1 0 |
| 3 : | 0 0 1 1 | 11 : | 1 0 1 1 |
| 4 : | 0 1 0 0 | 12 : | 1 1 0 0 |
| 5 : | 0 1 0 1 | 13 : | 1 1 0 1 |
| 6 : | 0 1 1 0 | 14 : | 1 1 1 0 |
| 7 : | 0 1 1 1 | 15 : | 1 1 1 1 |

The issue with binary encoding is that the positional information is localized around a few bits.

https://kazemnejad.com/blog/transformer_architecture_positional_encoding/

JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

# Math Recap: Sine and Cosine Functions

# Absolute Positional Embeddings

- Let $t$ be a desired position. Then the $i$-th element of the positional vector is:

$$\overrightarrow{p_t}^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k.t), & \text{if } i = 2k \\ \cos(\omega_k.t), & \text{if } i = 2k + 1 \end{cases} \qquad \omega_k = \frac{1}{10000^{2k/d}}$$

- Here $d$ is the maximum dimension.
- This provides unique vectors for each position.

https://kazemnejad.com/blog/transformer_architecture_positional_encoding/

# Quiz

- Let $t$ be a desired position:

$$\overrightarrow{p_t}^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k \cdot t), & \text{if } i = 2k \\ \cos(\omega_k \cdot t), & \text{if } i = 2k+1 \end{cases} \qquad \omega_k = \frac{1}{10000^{2k/d}}$$

- **Q:** Are the frequencies increasing with dimension i ?
- **Answer:** The frequencies are decreasing along the vector dimension.

# Visualizing Absolute Positional Embeddings

- Here positions range from 0-50, for an embedding dimension of 130.

An approach:
Sine/Cosine encoding

$$p_i = \begin{pmatrix} \sin(i/10000^{2*1/d}) \\ \cos(i/10000^{2*1/d}) \\ \vdots \\ \sin(i/10000^{2*\frac{d}{2}/d}) \\ \cos(i/10000^{2*\frac{d}{2}/d}) \end{pmatrix}$$



$p_i$ are positional embeddings

Allows model to learn relative positioning

$p_1$    $x_1$

$p_2$    $x_2$

$p_3$    $x_3$

$p_4$    $x_4$

# Transformer-based Language Modeling

Image by http://jalammar.github.io/illustrated-gpt2/

# Training a Transformer Language Model

- **Goal:** Train a Transformer for language modeling (i.e., predicting the next word).
- **Approach:** Train it so that each position is predictor of the next (right) token.
  - We just shift the input to right by one, and use as labels

EOS special token

(gold output) $Y =$    cat    sat    on    the    mat    </s>

```
X = text[:, :-1]
Y = text[:, 1:]
```

## TRANSFORMER

$X =$    the    cat    sat    on    the    mat

[Slide credit: Arman Cohan]

JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

44

# Training a Transformer Language Model

- For each position, compute their corresponding **distribution** over the whole vocab.

(gold output) $Y = $ cat   sat   on   the   mat   </s>



$X = $ the cat sat on the mat

# Training a Transformer Language Model

- For each position, compute the **loss** between the distribution and the gold output label.

(gold output) $Y$ = cat   sat   on   the   mat   \</s\>



TRANSFORMER

$X$ = the  cat  sat  on  the  mat

# Training a Transformer Language Model

- Sum the position-wise loss values to a obtain a **global loss**.

(gold output) $Y =$ cat   sat   on   the   mat   &lt;/s&gt;

$\mathcal{L} = \updownarrow + \updownarrow + \updownarrow + \updownarrow + \updownarrow + \updownarrow$

TRANSFORMER

$X =$ the  cat  sat  on  the  mat

# Training a Transformer Language Model

- Using this loss, do **Backprop** and **update** the Transformer parameters.

(gold output) $Y$ = cat    sat    on    the    mat    </s>

$\mathcal{L}$ = $\updownarrow$ + $\updownarrow$ + $\updownarrow$ + $\updownarrow$ + $\updownarrow$ + $\updownarrow$

$\nabla\mathcal{L}$

TRANSFORMER

$X$ = the cat sat on the mat

Well, this is not quite right 🐼
…
what is the problem with this?

# Training a Transformer Language Model

- The model would solve the task by copying the next token to output (data leakage).
  - Does not learn anything useful



(gold output) $Y$ = cat    sat    on    the    mat    </s>

$\mathcal{L}$ =

$\nabla\mathcal{L}$

TRANSFORMER

$X$ =    the    cat    sat    on    the    mat

# Training a Transformer Language Model

- We need to **prevent information leakage** from future tokens! How?

(gold output) $Y = $ cat sat on the mat </s>

$\mathcal{L} = $

$\nabla\mathcal{L}$

TRANSFORMER

$X = $ the cat sat on the mat

# Attention mask



Attention raw scores

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | -0.08 | 1.24 | 0.69 | -0.98 | 1.43 | -0.6 | 0.7 | 0.16 | 0.93 | 1.28 | -1.61 | -1.1 |
| 1 | -0.09 | -0.0 | -0.7 | 0.06 | 0.25 | 0.23 | 0.26 | 0.18 | 0.78 | -0.21 | -1.01 | 1.01 |
| 2 | 0.86 | 1.19 | 1.59 | 0.86 | -0.13 | -0.15 | -2.13 | -0.98 | -0.87 | -1.72 | 1.87 | -0.72 |
| 3 | 0.12 | -0.03 | -0.02 | 0.88 | -0.46 | -0.7 | 0.54 | -0.42 | -1.89 | -0.38 | 0.04 | -0.84 |
| 4 | 0.51 | 0.17 | 0.13 | -1.64 | 0.24 | -0.02 | 1.68 | -0.36 | 0.64 | 0.36 | 0.27 | 0.66 |
| 5 | 0.24 | -1.44 | 0.43 | 0.74 | 0.96 | -1.21 | -0.31 | 1.54 | 1.66 | 1.14 | 0.58 | -1.44 |
| 6 | 0.26 | -0.1 | 0.93 | 0.72 | -0.38 | 1.65 | 0.47 | -0.96 | -0.17 | -0.9 | -1.57 | 0.22 |
| 7 | -0.55 | 0.81 | 0.71 | 1.7 | -0.8 | -1.14 | -0.32 | 1.78 | -0.7 | -0.04 | 1.54 | 0.81 |
| 8 | 0.74 | -0.76 | -0.44 | -0.08 | -1.38 | -0.13 | 1.25 | -1.37 | 1.84 | 0.3 | 0.57 | 0.74 |
| 9 | -0.97 | -0.91 | 0.15 | 0.35 | -0.81 | 0.11 | 1.14 | -1.52 | 1.06 | 1.87 | 0.5 | -0.3 |
| 10 | 1.56 | 0.9 | 0.39 | 1.46 | 1.44 | -1.05 | 0.9 | -0.73 | 0.36 | -0.67 | -0.62 | -0.43 |
| 11 | 0.32 | 0.74 | 0.44 | -0.1 | 1.19 | 0.83 | 0.29 | 2.06 | 0.51 | -0.26 | 1.51 | 0.11 |

What we want          What we have

# Attention mask



Attention raw scores

Attention mask

large negative numbers,
which leads to softmax($-\infty$) ≈ 0

Slide credit: Arman Cohan

# Attention mask



Attention raw scores



X

Attention mask



Note matrix multiplication is quite fast in GPUs.

Slide credit: Arman Cohan

53

# Attention mask

# **Attention mask**

The effect is more than just pruning out some of the wirings in self-attention block.



X

softmax →

Attention probabilities

# Training a Transformer Language Model

- We need to **prevent information leakage** from future tokens! How?

(gold output) $Y = $ cat    sat    on    the    mat    </s>

$\mathcal{L} = $ $\updownarrow$ + $\updownarrow$ + $\updownarrow$ + $\updownarrow$ + $\updownarrow$ + $\updownarrow$

$\nabla\mathcal{L}$

+ masking

TRANSFORMER

$X = $ the   cat   sat   on   the   mat

# How to use the model to generate text?

- Use the output of previous step as input to the next step repeatedly

# How to use the model to generate text?

- Use the output of previous step as input to the next step repeatedly

The probabilities get revised upon adding a new token to the input.

sample → on

TRANSFORMER

the cat sat

# How to use the model to generate text?

- Use the output of previous step as input to the next step repeatedly

# How to use the model to generate text?

- Use the output of previous step as input to the next step repeatedly

The probabilities get revised upon adding a new token to the input.



sample → mat

TRANSFORMER

the cat sat on the

# How to use the model to generate text?

- Use the output of previous step as input to the next step repeatedly

The probabilities get revised upon adding a new token to the input.

sample → </s>

TRANSFORMER

the cat sat on the mat

# Summary

- This is a very generic Transformer!
- We will implement this in HW5 to build a simple Transformer Language Model!!

- **Next:**
  - Architectural variants

# Transformer Architectural Variants

# Encoder-Decoder Architectures

- It is useful to think of generative models as two sub-models.

"The cat sat on the [MASK]" → Some model →

# Encoder-Decoder Architectures

- It is useful to think of generative models as two sub-models

Representation (compression) of the context

"The cat sat on the [MASK]"

Encoder

Decoder

Processes the context and compiles it into a vector.

Produces the output sequence item by item using the representation of the context.

# Encoder-decoder models

- Transformer is two blocks

- Encoder = read or encode the input,
  - Architecture is as we've seen

- Decoder = generate or decode the output
  - Architecture is identical to the encoder but we give it the ability to also attend to the input

Le   chat   est   mignon

Encoder

Decoder

The   cat   is   cute

&lt;s&gt;   Le   chat   est

# Transformer [Vaswani et al. 2017]

JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

# Transformer [Vaswani et al. 2017]

JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

# **Transformer** [Vaswani et al. 2017]

JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

# Transformer [Vaswani et al. 2017]

JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

# **Transformer** [Vaswani et al. 2017]

- Computation of **encoder** attends to both sides.



Encoder Self-Attention



[Attention Is All You Need, Vaswani et al. 2017]

# **Transformer** [Vaswani et al. 2017]

- At any step of **decoder**, it attends to previous computation of **encoder**



Encoder-Decoder Attention

[Attention Is All You Need, Vaswani et al. 2017]

# **Transformer** [Vaswani et al. 2017]

- At any step of **decoder**, it attends to previous computation of **encoder** as well as **decoder's** own generations



MaskedDecoder Self-Attention



[Attention Is All You Need, Vaswani et al. 2017]

# Transformer [Vaswani et al. 2017]

- At any step of **decoder**, it attends to previous computation of **encoder** as well as **decoder's** own generations

- At any step of decoder, **re-use** previous computation of encoder.

- Computation of decoder is **linear**, instead of quadratic.



[Attention Is All You Need, Vaswani et al. 2017]

74

# Recap: Transformer

- Yaaay we know Transformers now! 🥳
- An encoder-decoder architecture
- 3 forms of attention



Encoder-Decoder Attention

Encoder Self-Attention

MaskedDecoder Self-Attention

[Attention Is All You Need, Vaswani et al. 2017]

# Quiz: Enc-Dec Cost

- Source data (large!):
  - The references for a Wikipedia article.
  - Web search using article section titles, ~ 10 web pages per query.

- For a passage of length N and a summary of length M, the complexity of the attention is:
  - $O(N) + O(M)$
  - $O(N) + O(M) + O(NM)$
  - $O(N^2) + O(M^2) + O(NM)$
  - $O(N^2) + O(M^2)$

No, self attention is all-to-all and so quadratic.

JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

# Quiz: Enc-Dec Cost

- Source data (large!):
  - The references for a Wikipedia article.
  - Web search using article section titles, ~ 10 web pages per query.

- For a passage of length N and a summary of length M, the complexity of the attention is:
  - $O(N) + O(M)$
  - $O(N) + O(M) + O(NM)$
  - $O(N^2) + O(M^2) + O(NM)$
  - $O(N^2) + O(M^2)$

No, self attention is all-to-all and so quadratic in $M$ and $N$.

JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

# Quiz: Enc-Dec Cost

- Source data (large!):
  - The references for a Wikipedia article.
  - Web search using article section titles, ~ 10 web pages per query.

- For a passage of length N and a summary of length M, the complexity of the attention is:
  - $O(N) + O(M)$
  - $O(N) + O(M) + O(NM)$
  - $O(N^2) + O(M^2) + O(NM)$
  - $O(N^2) + O(M^2)$

No, self attention is all-to-all and so quadratic in $M$ and $N$.

JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

# Quiz: Enc-Dec Cost

- Source data (large!):
  - The references for a Wikipedia article.
  - Web search using article section titles, ~ 10 web pages per query.

- For a passage of length N and a summary of length M, the complexity of the attention is:
  - $O(N) + O(M)$
  - $O(N) + O(M) + O(NM)$
  - $O(N^2) + O(M^2) + O(NM)$
  - $O(N^2) + O(M^2)$

No, cross attention is missing.

[Slide: John Canny]

# Quiz: Enc-Dec Cost

- Source data (large!):
    - The references for a Wikipedia article.
    - Web search using article section titles, ~ 10 web pages per query.

- For a passage of length N and a summary of length M, the complexity of the attention is:
    - $O(N) + O(M)$
    - $O(N) + O(M) + O(NM)$
    - $O(N^2) + O(M^2) + O(NM)$  Yes. The three terms are respectively the Encoder self-attention, Decoder self-attention, and Cross attention.
    - $O(N^2) + O(M^2)$

JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

# Quiz: Enc-Dec Connections

- Which best represents encoder-decoder connections?



Incorrect

Correct

[Slide credit: CS886 at UWaterloo]

# Considerations about computational cost in Transformers

# Making decoding more efficient



$$Q = \mathbf{W}^q \mathbf{x}$$
$$K = \mathbf{W}^k \mathbf{x}$$
$$V = \mathbf{W}^v \mathbf{x}$$

$$\text{Attention}(\mathbf{x}) = \text{softmax}\left(\frac{QK^{\mathrm{T}}}{\sqrt{d}}\right)V$$

[Slide credit: Arman Cohan]

# Making decoding more efficient

$Q = \mathbf{W}^q \mathbf{x}$
$K = \mathbf{W}^k \mathbf{x}$
$V = \mathbf{W}^v \mathbf{x}$

$\text{Attention}(\mathbf{x}) = \text{softmax}\left(\dfrac{QK^{\mathrm{T}}}{\sqrt{d}}\right) V$

x

q

q: the next token

K

V

previous context

# Making decoding more efficient



$$Q = \mathbf{W}^q \mathbf{x}$$
$$K = \mathbf{W}^k \mathbf{x}$$
$$V = \mathbf{W}^v \mathbf{x}$$

$$\text{Attention}(\mathbf{x}) = \text{softmax}\left(\frac{QK^{\mathrm{T}}}{\sqrt{d}}\right) V$$

q

q: the next token

$K = W_k x$

$V = W_v x$

previous context

The cat sat on the

[Slide credit: Arman Cohan]

# Making decoding more efficient

$$Q = \mathbf{W}^q \mathbf{x}$$
$$K = \mathbf{W}^k \mathbf{x}$$
$$V = \mathbf{W}^v \mathbf{x}$$

$$\text{Attention}(\mathbf{x}) = \text{softmax}\left(\frac{QK^{\mathrm{T}}}{\sqrt{d}}\right) V$$

q

q: the next token

$$K = W_k x$$

$$V = W_v x$$

previous context

The cat sat on the

# Making decoding more efficient

$$Q = \mathbf{W}^q \mathbf{x}$$
$$K = \mathbf{W}^k \mathbf{x}$$
$$V = \mathbf{W}^v \mathbf{x}$$

$$\text{Attention}(\mathbf{x}) = \text{softmax}\left(\frac{QK^{\mathrm{T}}}{\sqrt{d}}\right) V$$

q

q: the next token

$$K = W_k x$$

$$V = W_v x$$

previous context

The cat sat on the

[Slide credit: Arman Cohan]

# Making decoding more efficient

$Q = \mathbf{W}^q \mathbf{x}$
$K = \mathbf{W}^k \mathbf{x}$
$V = \mathbf{W}^v \mathbf{x}$

$\text{Attention}(\mathbf{x}) = \text{softmax}\left(\dfrac{QK^{\text{T}}}{\sqrt{d}}\right) V$



q

q: the next token

$K = W_k x$

$V = W_v x$

previous context

The cat sat on the

[Slide credit: Arman Cohan]

90

# Making decoding more efficient

$Q = \mathbf{W}^q \mathbf{x}$
$K = \mathbf{W}^k \mathbf{x}$
$V = \mathbf{W}^v \mathbf{x}$

$$\text{Attention}(\mathbf{x}) = \text{softmax}\left(\frac{QK^{\mathrm{T}}}{\sqrt{d}}\right)V$$

q

q: the next token

$K = W_k x$

$V = W_v x$

previous context

The cat sat on the

[Slide credit: Arman Cohan]

# Making decoding more efficient

$$Q = \mathbf{W}^q \mathbf{x}$$
$$K = \mathbf{W}^k \mathbf{x}$$
$$V = \mathbf{W}^v \mathbf{x}$$

$$\text{Attention}(\mathbf{x}) = \text{softmax}\left(\frac{QK^{\mathrm{T}}}{\sqrt{d}}\right)V$$

- We are computing the Keys and Values many times!
  - Let's reduce redundancy! 🥵

q

q: the next token

$$K = W_k x$$

$$V = W_v x$$

previous context

The cat sat on the

[Slide credit: Arman Cohan]

# Making decoding more efficient

$$Q = \mathbf{W}^q \mathbf{x}$$
$$K = \mathbf{W}^k \mathbf{x}$$
$$V = \mathbf{W}^v \mathbf{x}$$

- We are computing the Keys and Values many times!
  - Let's reduce redundancy! 🥵

$$\text{Attention}(\mathbf{x}) = \text{softmax}\left(\frac{QK^{\mathrm{T}}}{\sqrt{d}}\right)V$$

$$k_{new} = W_k \mathbf{x}[:, :-1]$$

q

q: the next token

K Cached

V Cached

previous context

$$v_{new} = W_v \mathbf{x}[:, :-1]$$

The cat sat on the

[Slide credit: Arman Cohan]

# Making decoding more efficient

$Q = \mathbf{W}^q \mathbf{x}$
$K = \mathbf{W}^k \mathbf{x}$
$V = \mathbf{W}^v \mathbf{x}$

- **Question:** How much memory does this K, V cache require?

$$\text{Attention}(\mathbf{x}) = \text{softmax}\left(\frac{QK^{\mathrm{T}}}{\sqrt{d}}\right)V$$

$$k_{new} = W_k \mathbf{x}[:, :-1]$$

q

q: the next token

K Cached

V Cached

previous context

$$v_{new} = W_v \mathbf{x}[:, :-1]$$

The cat sat on the

[Slide credit: Arman Cohan]

# Writing our own Transformer

# Clone Helper Function

- Create N copies of pytorch nn.Module
- The Transformer's structure contains a lot of design repetition (like VGG)
- Remember these clones shouldn't share parameters (for the most part)

```python
def clones(module, N):
    "Produce N identical layers."
    return nn.ModuleList([copy.deepcopy(module) for _ in range(N)])
```

[Slide credit: CS886 at UWaterloo]

# Create Embedding

- Create vector representation of sequence vocabulary
- nn.Embedding creates a lookup table to map sequence vocabulary to unique vectors

```python
class Embeddings(nn.Module):
    def __init__(self, d_model, vocab):
        super(Embeddings, self).__init__()
        self.lut = nn.Embedding(vocab, d_model)
        self.d_model = d_model

    def forward(self, x):
        return self.lut(x) * math.sqrt(self.d_model)
```

[Slide credit: CS886 at UWaterloo]

# Positional Encoding

- Add information about an element's position in a sequence to its representation
- Element wise addition of sinusoidal encoding

Positional Encoding

```python
class PositionalEncoding(nn.Module):
    "Implement the PE function."

    def __init__(self, d_model, dropout, max_len=5000):
        super(PositionalEncoding, self).__init__()
        self.dropout = nn.Dropout(p=dropout)

        # Compute the positional encodings once in log space.
        pe = torch.zeros(max_len, d_model)
        position = torch.arange(0, max_len).unsqueeze(1)
        div_term = torch.exp(
            torch.arange(0, d_model, 2) * -(math.log(10000.0) / d_model)
        )
        pe[:, 0::2] = torch.sin(position * div_term)
        pe[:, 1::2] = torch.cos(position * div_term)
        pe = pe.unsqueeze(0)
        self.register_buffer("pe", pe)

    def forward(self, x):
        x = x + self.pe[:, : x.size(1)].requires_grad_(False)
        return self.dropout(x)
```

[Slide credit: CS886 at UWaterloo]

# Attention block



```python
def attention(query, key, value, mask=None, dropout=None):
    "Compute 'Scaled Dot Product Attention'"
    d_k = query.size(-1)
    scores = torch.matmul(query, key.transpose(-2, -1)) / math.sqrt(d_k)
    if mask is not None:
        scores = scores.masked_fill(mask == 0, -1e9)
    p_attn = scores.softmax(dim=-1)
    if dropout is not None:
        p_attn = dropout(p_attn)
    return torch.matmul(p_attn, value), p_attn
```

$-1e9$ is a large negative number, which leads to softmax(-1e9) $\approx 0$

[Slide credit: CS886 at UWaterloo]

# Multi-Head Attention

```python
class MultiHeadedAttention(nn.Module):
    def __init__(self, h, d_model, dropout=0.1):
        "Take in model size and number of heads."
        super(MultiHeadedAttention, self).__init__()
        assert d_model % h == 0
        # We assume d_v always equals d_k
        self.d_k = d_model // h
        self.h = h
        self.linears = clones(nn.Linear(d_model, d_model), 4)
        self.attn = None
        self.dropout = nn.Dropout(p=dropout)
```

```python
def forward(self, query, key, value, mask=None):
    "Implements Figure 2"
    if mask is not None:
        # Same mask applied to all h heads.
        mask = mask.unsqueeze(1)
    nbatches = query.size(0)

    # 1) Do all the linear projections in batch from d_model => h x d_k
    query, key, value = [
        lin(x).view(nbatches, -1, self.h, self.d_k).transpose(1, 2)
        for lin, x in zip(self.linears, (query, key, value))
    ]

    # 2) Apply attention on all the projected vectors in batch.
    x, self.attn = attention(
        query, key, value, mask=mask, dropout=self.dropout
    )

    # 3) "Concat" using a view and apply a final linear.
    x = (
        x.transpose(1, 2)
        .contiguous()
        .view(nbatches, -1, self.h * self.d_k)
    )
    del query
    del key
    del value
    return self.linears[-1](x)
```

JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

[Slide credit: CS886 at UWaterloo]

# FeedForward Layer

```python
class PositionwiseFeedForward(nn.Module):
    "Implements FFN equation."

    def __init__(self, d_model, d_ff, dropout=0.1):
        super(PositionwiseFeedForward, self).__init__()
        self.w_1 = nn.Linear(d_model, d_ff)
        self.w_2 = nn.Linear(d_ff, d_model)
        self.dropout = nn.Dropout(dropout)

    def forward(self, x):
        return self.w_2(self.dropout(self.w_1(x).relu()))
```

[Slide credit: CS886 at UWaterloo]

# Sublayer Connections

```python
class SublayerConnection(nn.Module):
    """
    A residual connection followed by a layer norm.
    Note for code simplicity the norm is first as opposed to last.
    """

    def __init__(self, size, dropout):
        super(SublayerConnection, self).__init__()
        self.norm = LayerNorm(size)
        self.dropout = nn.Dropout(dropout)

    def forward(self, x, sublayer):
        "Apply residual connection to any sublayer with the same size."
        return x + self.dropout(sublayer(self.norm(x)))
```

JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

# Encoder Layer

```python
class EncoderLayer(nn.Module):
    "Encoder is made up of self-attn and feed forward (defined below)"

    def __init__(self, size, self_attn, feed_forward, dropout):
        super(EncoderLayer, self).__init__()
        self.self_attn = self_attn
        self.feed_forward = feed_forward
        self.sublayer = clones(SublayerConnection(size, dropout), 2)
        self.size = size

    def forward(self, x, mask):
        "Follow Figure 1 (left) for connections."
        x = self.sublayer[0](x, lambda x: self.self_attn(x, x, x, mask))
        return self.sublayer[1](x, self.feed_forward)
```

[Slide credit: CS886 at UWaterloo]

# Decoder Layer

- Same as encoder layers other than:
  - the additional multi-head attention block to preform cross-attention with the output representation from the encoder

```python
class DecoderLayer(nn.Module):
    "Decoder is made of self-attn, src-attn, and feed forward (defined below)"

    def __init__(self, size, self_attn, src_attn, feed_forward, dropout):
        super(DecoderLayer, self).__init__()
        self.size = size
        self.self_attn = self_attn
        self.src_attn = src_attn
        self.feed_forward = feed_forward
        self.sublayer = clones(SublayerConnection(size, dropout), 3)

    def forward(self, x, memory, src_mask, tgt_mask):
        "Follow Figure 1 (right) for connections."
        m = memory
        x = self.sublayer[0](x, lambda x: self.self_attn(x, x, x, tgt_mask))
        x = self.sublayer[1](x, lambda x: self.src_attn(x, m, m, src_mask))
        return self.sublayer[2](x, self.feed_forward)
```
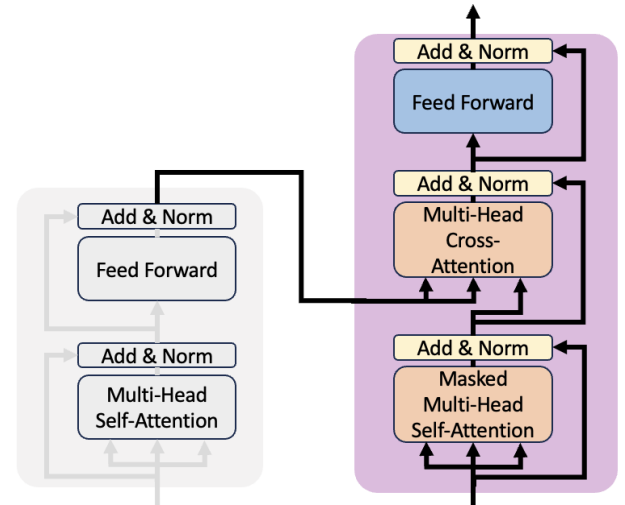


[Slide credit: CS886 at UWaterloo]

# The Prediction Head

- A final linear mapping
- Apply softmax to convert logits to probabilities

```python
class Generator(nn.Module):
    "Define standard linear + softmax generation step."

    def __init__(self, d_model, vocab):
        super(Generator, self).__init__()
        self.proj = nn.Linear(d_model, vocab)

    def forward(self, x):
        return log_softmax(self.proj(x), dim=-1)
```



Softmax

Linear

Output
Probabilities
For next area

# Build each block

```python
class Encoder(nn.Module):
    "Core encoder is a stack of N layers"

    def __init__(self, layer, N):
        super(Encoder, self).__init__()
        self.layers = clones(layer, N)
        self.norm = LayerNorm(layer.size)

    def forward(self, x, mask):
        "Pass the input (and mask) through each layer in turn."
        for layer in self.layers:
            x = layer(x, mask)
        return self.norm(x)
```

```python
class Decoder(nn.Module):
    "Generic N layer decoder with masking."

    def __init__(self, layer, N):
        super(Decoder, self).__init__()
        self.layers = clones(layer, N)
        self.norm = LayerNorm(layer.size)

    def forward(self, x, memory, src_mask, tgt_mask):
        for layer in self.layers:
            x = layer(x, memory, src_mask, tgt_mask)
        return self.norm(x)
```

[Slide credit: CS886 at UWaterloo]

# Putting it Together

```python
class EncoderDecoder(nn.Module):
    """
    A standard Encoder-Decoder architecture. Base for this and many
    other models.
    """

    def __init__(self, encoder, decoder, src_embed, tgt_embed, generator):
        super(EncoderDecoder, self).__init__()
        self.encoder = encoder
        self.decoder = decoder
        self.src_embed = src_embed
        self.tgt_embed = tgt_embed
        self.generator = generator

    def forward(self, src, tgt, src_mask, tgt_mask):
        "Take in and process masked src and target sequences."
        return self.decode(self.encode(src, src_mask), src_mask, tgt, tgt_mask)

    def encode(self, src, src_mask):
        return self.encoder(self.src_embed(src), src_mask)

    def decode(self, memory, src_mask, tgt, tgt_mask):
        return self.decoder(self.tgt_embed(tgt), memory, src_mask, tgt_mask)
```

[Slide credit: CS886 at UWaterloo]

JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

# Initialize the model

```python
def make_model(
    src_vocab, tgt_vocab, N=6, d_model=512, d_ff=2048, h=8, dropout=0.1
):
    "Helper: Construct a model from hyperparameters."
    c = copy.deepcopy
    attn = MultiHeadedAttention(h, d_model)
    ff = PositionwiseFeedForward(d_model, d_ff, dropout)
    position = PositionalEncoding(d_model, dropout)
    model = EncoderDecoder(
        Encoder(EncoderLayer(d_model, c(attn), c(ff), dropout), N),
        Decoder(DecoderLayer(d_model, c(attn), c(attn), c(ff), dropout), N),
        nn.Sequential(Embeddings(d_model, src_vocab), c(position)),
        nn.Sequential(Embeddings(d_model, tgt_vocab), c(position)),
        Generator(d_model, tgt_vocab),
    )

    # This was important from their code.
    # Initialize parameters with Glorot / fan_avg.
    for p in model.parameters():
        if p.dim() > 1:
            nn.init.xavier_uniform_(p)
    return model
```

[Slide credit: CS886 at UWaterloo]