

Deep Learning Transformer and Newer Architectures

Shrey Jain, Dareen Alharthi, Hao Chen

Fall 2025
Attendance:@1205

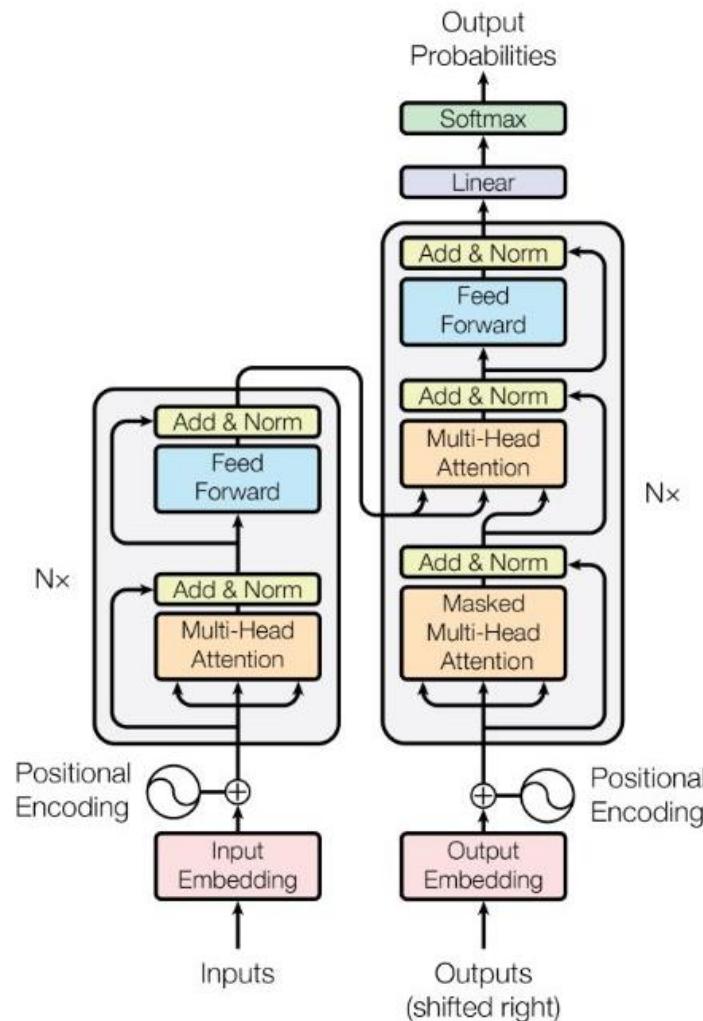
Content

- Transformer Architecture
- Improvements on Transformers
- Transformer for different modalities
- Scaling Laws
- Parameter Efficient Tuning
- Quantizing Transformers
- Interpreting Transformer – attention, logit lens

Content

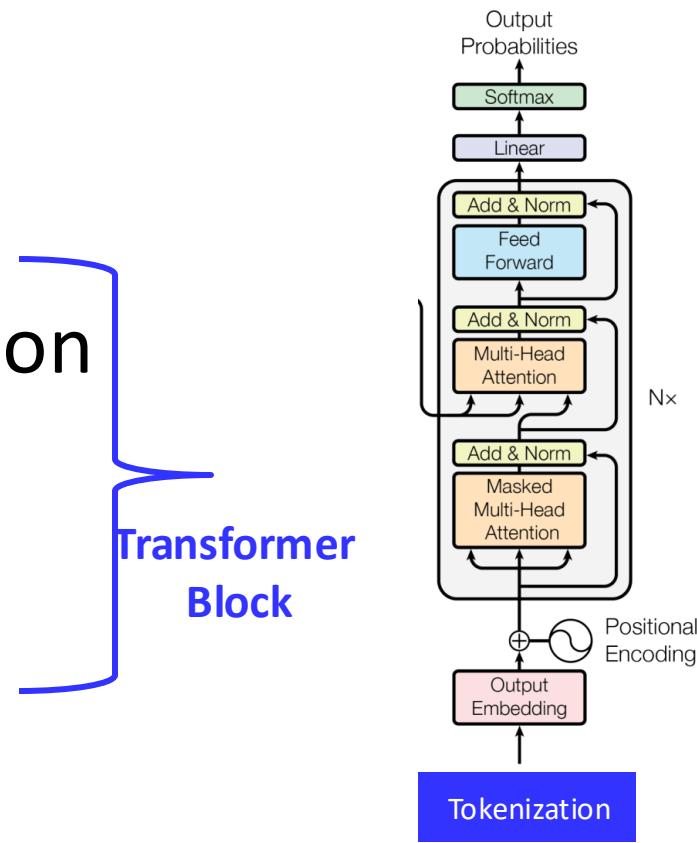
- Transformer Architecture
- Improvements on Transformers
- Transformer for different modalities
- Scaling Laws
- Parameter Efficient Tuning
- Quantizing Transformers
- Interpreting Transformer – attention, logit lens

Transformer Architecture



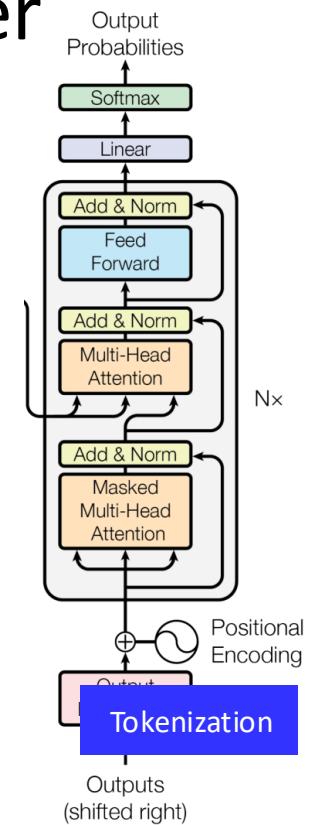
Transformer Architecture

- Word Tokenization
 - Word Embedding
 - (Masked) Multi-Head Attention
 - Position Encoding
 - Feed-Forward
 - Add & Norm
 - Output Projection Layer
- Transformer Block



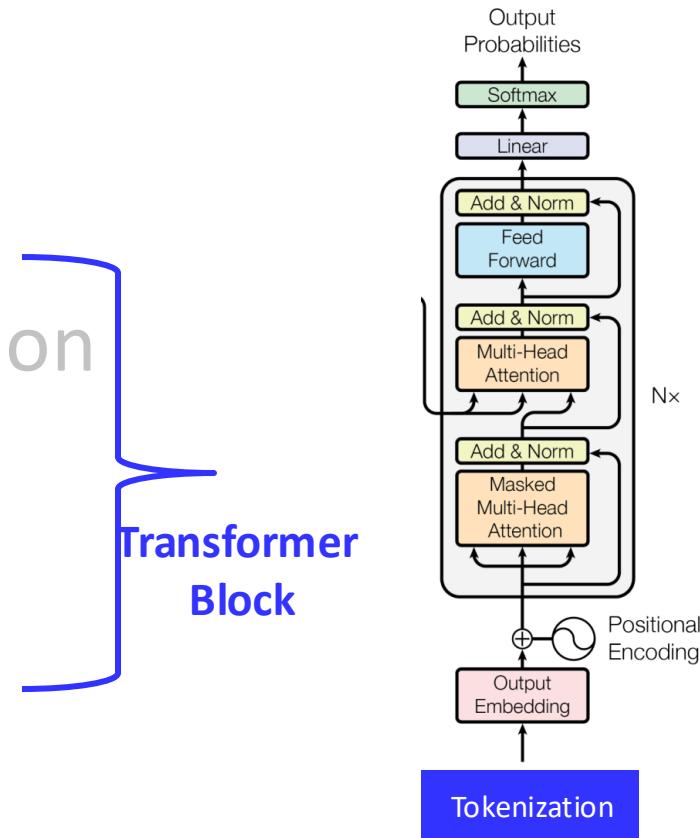
Transformer Architecture

- Transformer is a giant number cruncher
- But it takes in text, outputs text
- How do you go from text to numbers and back?
 - Old solution: one-hot vectors
 - But one-hot of what?



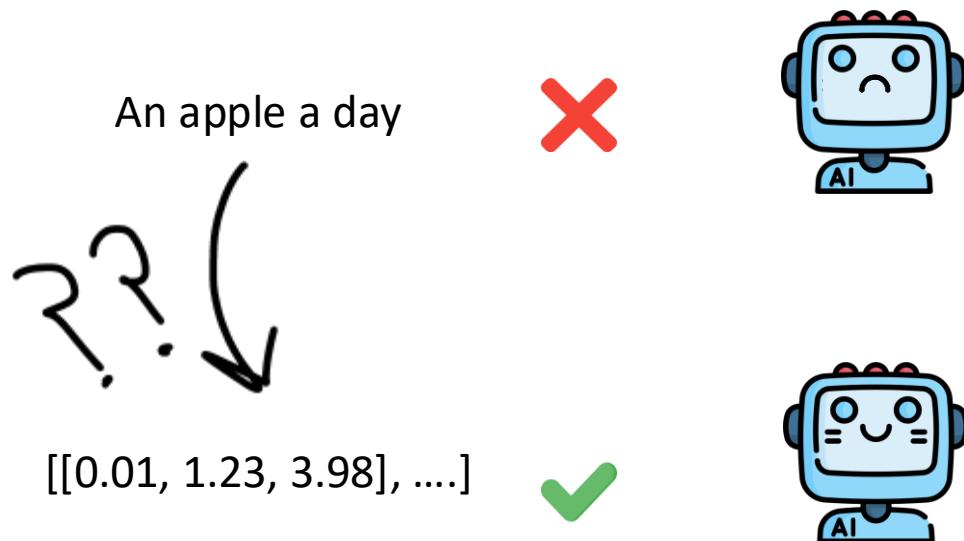
Transformer Architecture

- Word Tokenization
- Word Embedding
- (Masked) Multi-Head Attention
- Position Encoding
- Feed-Forward
- Add & Norm
- Output Projection Layer



Tokenization

- Computers don't understand text, they only get numbers
- How to represent text as numbers?



Tokenization

- Computers don't understand text, they only get numbers
- How to represent text as numbers?
- Words are basic semantic units, so, one-hot of a dictionary of all words?

$[[0,0,\dots,1,0,\dots,0],$
 $[0,0,\dots,0,1,\dots,0],$
 $[0,0,\dots,1,0,\dots,0],$
 $[0,\dots,1,0,0,\dots,0]]$



An apple a day



An apple a day

Tokenization

- Computers don't understand text, they only get numbers
- How to represent text as numbers?
- Words are basic semantic units, so, one-hot of a dictionary of all words?
 - What about new words? <UNK> helps no one!
 - 'Hello', 'hello' and 'HeLIO' almost same yet not same, a word for each permutation?

Tokenization

- Computers don't understand text, they only get numbers
- How to represent text as numbers?
- Words are basic semantic units, so, one-hot of a dictionary of all words?
 - Advantage would be never producing an illegible word
 - But it's too inflexible for input:
 - New words can't be handled
 - Typos can't be handled either
 - Permutation nightmare



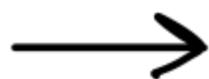
Tokenization

- Computers don't understand text, they only get numbers
- How to represent text as numbers?
- Index for each character? Representing Input would be simple!
 - Represent any input unambiguously
 - Tiny vocabulary!

Tokenization

- Computers don't understand text, they only get numbers
- How to represent text as numbers?
- Index for each character? Representing Input would be simple!

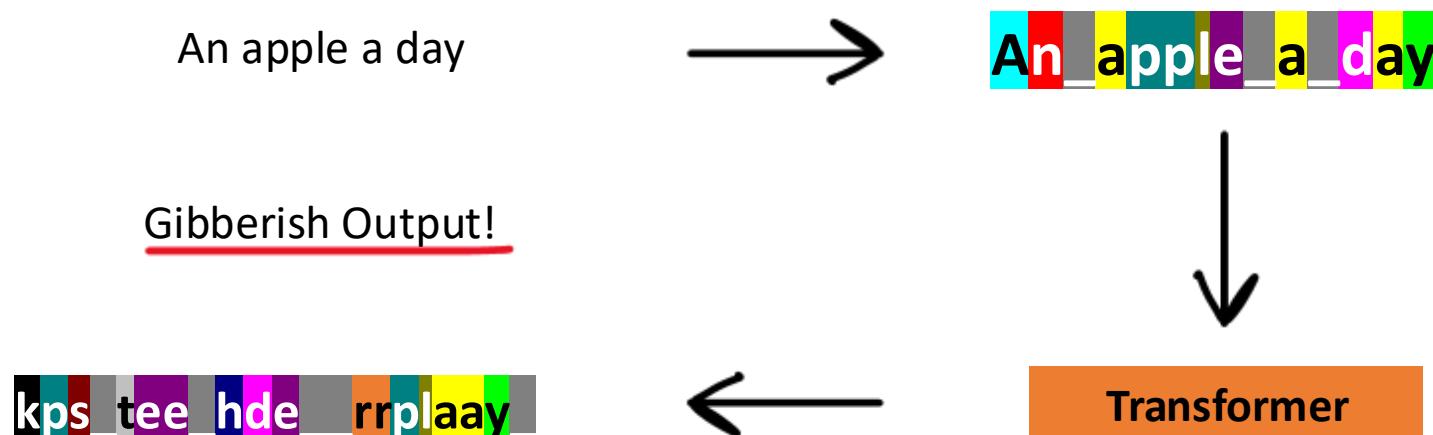
An apple a day



An_apple_a_day

Tokenization

- Computers don't understand text, they only get numbers
- How to represent text as numbers?
- Index for each character? Representing Input would be simple!



Tokenization

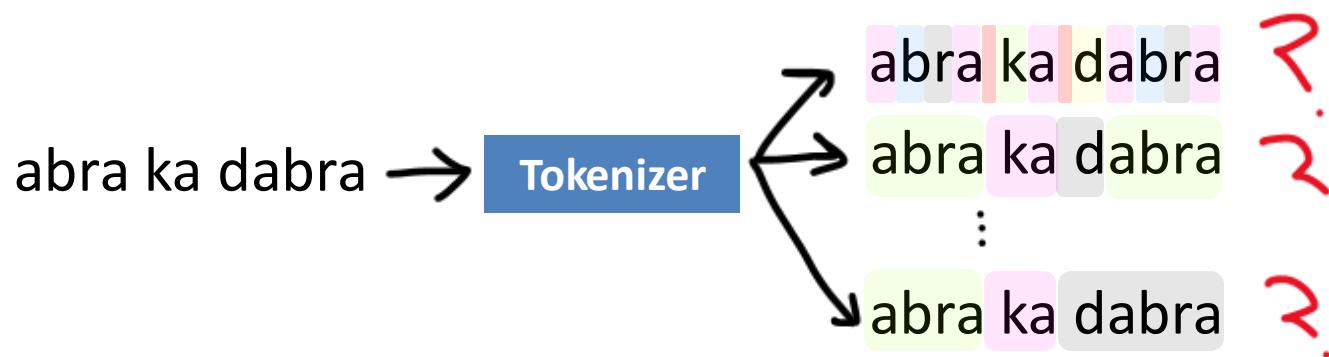
- Computers don't understand text, they only get numbers
- How to represent text as numbers?
 - Words are semantically perfect but inflexible
 - Characters too flexible and output is junk

Tokenization

- Computers don't understand text, they only get numbers
- How to represent text as numbers?
- Enter: Sub-word tokenization, a middle-ground

Tokenization

- Computers don't understand text, they only get numbers
- How to represent text as numbers?
- Enter: Sub-word tokenization, a middle-ground
 - Get meaningfulness of a word for the outputs
 - Get flexibility of smaller units for inputs
 - Units can even include punctuations!



Tokenization

- Computers don't understand text, they only get numbers
- How to represent text as numbers?
- Enter: Sub-word tokenization, a middle-ground
 - Get meaningfulness of a word for the outputs
 - Get flexibility of smaller units for inputs
 - Units can even include punctuations!
- Tokenizer, an external module, is learned to do this deterministically

abra ka dabra → **Tokenizer** → abra ka dabra

Tokenization

- Computers don't understand text, they only get numbers
- How to represent text as numbers?
- Enter: Sub-word tokenization, a middle-ground
 - Get meaningfulness of a word for the outputs
 - Get flexibility of smaller units for inputs
 - Units can even include punctuations!
- Tokenizer, an external module, is learned to do this deterministically
 - Methods: **Byte-pair encoding**, WordPiece

Tokenization - BPE

- Iteratively, replace most common token pair with a single token

Tokenization - BPE

- Iteratively, replace most common token pair with a single token
- Let's take an example:

Tokenize: aaabdaaabbac

Tokenization - BPE

- Iteratively, replace most common token pair with a single token
- Let's take an example:

Tokenize: aaabdaaaabac

Counts: {'aa': 4, 'ab': 2, 'bd': 1, 'da': 1, 'ba': 1, 'ac': 1}

Tokenization - BPE

- Iteratively, replace most common token pair with a single token
- Let's take an example:

Tokenize: aaabdaaaabac

Counts: {'aa': 4, 'ab': 2, 'bd': 1, 'da': 1, 'ba': 1, 'ac': 1}

Replace 'aa' with Z: aaabdaaaabac -> ZabdZabac

Tokenization - BPE

- Iteratively, replace most common token pair with a single token
- Let's take an example:

Tokenize: ZabdZabac

Counts: {'ab

Replace 'ab' with Y: ZabdZabac -> ZYdZYac

Tokenization - BPE

- Iteratively, replace most common token pair with a single token
- Let's take an example:

So we'd get:

aaabdaaaabac -> ZabdZabac -> ZYdZYac -> XdXac

Where $Z = aa$, $Y = ab$, $X = ZY = aaab$

Tokenization - BPE

- Iteratively, replace most common token pair with a single token
- Let's take an example:

So we'd get:

aaabdaaaabac -> ZabdZabac -> ZYdZYac -> XdXac

Where Z = aa, Y = ab, X = ZY = aaab

This compressed text to < ½ the original length

Tokenization - BPE

- Iteratively, replace most common token pair with a single token
- Let's take an example:

If we started with just lowercase vocab ('a': 1, 'b': 2...):

aaabdaaabac = [1, 1, 1, 2, 4, 1, 1, 1, 2, 1, 3]

Tokenization - BPE

- Iteratively, replace most common token pair with a single token
- Let's take an example:

If we started with just lowercase vocab ('a': 1, 'b': 2...):

aaabdaaaabac = [1, 1, 1, 2, 4, 1, 1, 1, 2, 1, 3]

We connect (1, 1) or ('a', 'a') -> 27 or 'aa'

Then (1, 2) -> 28 and finally (27, 28)-> 29

Tokenization - BPE

- Iteratively, replace most common token pair with a single token
- Let's take an example:

If we started with just lowercase vocab ('a': 1, 'b': 2...):
aabdaaabac = [1, 1, 1, 2, 4, 1, 1, 1, 2, 1, 3]

We get XdXac = [29, 4, 29, 1, 3]

Much shorter and simpler!

Tokenization - BPE

- Iteratively, replace most common token pair with a single token
- Let's take an example:

If we started with just lowercase vocab ('a': 1, 'b': 2...):
aabdaaabac = [1, 1, 1, 2, 4, 1, 1, 1, 2, 1, 3]

We get $X_d X_a c = [29, 4, 29, 1, 3]$

Much shorter and simpler!

Tokenization - BPE

- Iteratively, replace most common token pair with a single token
- Worked well for the small example,
how to handle for our multilingual complex world?

Tokenization - BPE

- Iteratively, replace most common token pair with a single token
- Worked well for the small example, how to handle for our multilingual complex world?
- Every known digital character part of Unicode (160K)
- Represent each char as a combination of bytes
 - Byte is 8-bits or 256 total possible values
 - e.g. G = [71], ü = [195, 188]

Tokenization - BPE

- Iteratively, replace most common token pair with a single token
- Worked well for the small example, how to handle for our multilingual complex world?
- Every known digital character part of Unicode (160K)
- Can represent each as combination of bytes
 - Byte is 8-bits or 256 total possible values
- Represent all data as bytes and then do BPE!

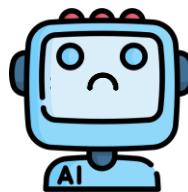
Only BPE isn't enough - Decoding

Code point \leftrightarrow UTF-8 conversion

First code point	Last code point	Byte 1	Byte 2	Byte 3	Byte 4
U+0000	U+007F	0yyz _{zzz}	0-127		
U+0080	U+07FF	110xx _{xyy}	10yyz _{zzz}	194-223, 128-191	224-239, 128-191
U+0800	U+FFFF	1110w _{www}	10xxxx _{yy}	10yyz _{zzz}	128-191, 128-191
U+010000	U+10FFFF	11110uvv	10vvwwww	10xxxx _{yy}	10yyz _{zzz}

240-247, 128-191, 128-191, 128-191

Decoding '\x80' or '128' fails



Only BPE isn't enough - Decoding

Code point ↔ UTF-8 conversion					
First code point	Last code point	Byte 1	Byte 2	Byte 3	Byte 4
U+0000	U+007F	0yyz _{zzzz}	0-127		
U+0080	U+07FF	110xx _{xyy}	10yyz _{zzzz}	194-223, 128-191	224-239, 128-191
U+0800	U+FFFF	1110w _{www}	10xxxx _{yy}	10yyz _{zzzz}	128-191, 128-191
U+010000	U+10FFFF	11110uvv	10vvwwww	10xxxx _{yy}	10yyz _{zzzz}
			240-247, 128-191,	128-191,	128-191

Decoding '\x80' or '128' fails

- handle gracefully: *errors='replace'* to use '?

Only BPE isn't enough - Decoding

Code point ↔ UTF-8 conversion					
First code point	Last code point	Byte 1	Byte 2	Byte 3	Byte 4
U+0000	U+007F	0yyz _{zzzz}	0-127		
U+0080	U+07FF	110xx _{xyy}	10yyz _{zzzz}	194-223, 128-191	224-239, 128-191
U+0800	U+FFFF	1110w _{www}	10xxxx _{yy}	10yyz _{zzzz}	128-191, 128-191
U+010000	U+10FFFF	11110uvv	10vvwwww	10xxxx _{yy}	10yyz _{zzzz}
			240-247, 128-191,	128-191,	128-191

Decoding '\x80' or '128' fails

- handle gracefully: *errors='replace'* to use '?'
- learn correct usage: model learns correct usage

Only BPE isn't enough

Code point \leftrightarrow UTF-8 conversion

First code point	Last code point	Byte 1	Byte 2	Byte 3	Byte 4
U+0000	U+007F	0yyz _{zzzz}	0-127		
U+0080	U+07FF	110xxxyy	10yyz _{zzzz}	194-223, 128-191	224-239, 128-191
U+0800	U+FFFF	1110wwww	10xxxxyy	10yyz _{zzzz}	128-191, 128-191
U+010000	U+10FFFF	11110uvv	10vvwwww	10xxxxyy	10yyz _{zzzz}

240-247, 128-191, 128-191, 128-191

Decoding '\x80' or '128' fails

- handle gracefully: *errors='replace'* to use '?'
- learn correct usage: model learns correct usage
- Constraints during decoding:

Only BPE isn't enough

Decoding '\x80' or '128' fails

- handle gracefully: *errors='replace'* to use '⁇'
 - learn correct usage: model learns correct usage
 - Constraints during decoding: Through HuggingFace's `transformers.LogitsProcessor` classes
-

Only BPE isn't enough

Decoding '\x80' or '128' fails

- handle gracefully: *errors='replace'* to use '⁇'
- learn correct usage: model learns correct usage
- Constraints during decoding: Through HuggingFace's `transformers.LogitsProcessor` classes

e.g. Decoding a json... output so far: ``{"age": ``

Next tokens: '{': 2.4 | '5': 2.91 | '}': 3.04 | ',': 4.1 | ...

Only BPE isn't enough

Decoding '\x80' or '128' fails

- handle gracefully: *errors='replace'* to use '█'
- learn correct usage: model learns correct usage
- Constraints during decoding: Through HuggingFace's `transformers.LogitsProcessor` classes

e.g. Decoding a json... output so far: ``{"age": ``

Next tokens: '{': 2.4 | '5': 2.91 | '}': 3.04 | ',': 4.1 | ...

Invalid format -> set logits to -inf before softmax

Only BPE isn't enough

Decoding '\x80' or '128' fails

- handle gracefully: *errors='replace'* to use '█'
- learn correct usage: model learns correct usage
- Constraints during decoding: Through HuggingFace's `transformers.LogitsProcessor` classes

SGLang LLM infer. framework leverages constrained decoding to get better format adherence while decoding e.g. JSON

Tokenization at a Large Scale

- How to avoid unintelligible strings from being decoded?
- Deliberate design choices pre-tokenization
 - Directly BPE on all docs?
'Hello world', 'Hello world.', 'Hello world!' ... other variations would be common
 - Do we learn all variants of world? ' world', ' world.', ' world!'
 - Unnecessarily redundant tokens

Tokenization at a Large Scale

- Case Study: GPT-2 [\[src\]](#)
- Regex Pattern:

```
's|'t|'re|'ve|'m|'ll|'d| ?\p{L}+| ?\p{N}+| ?[^s\p{L}\p{N}]+| \s+(?!\\S)| \s+
```

"Hello world I've learned to loveee! tokenization, I'M NOW THE #1Fan "

Tokenization at a Large Scale

- Case Study: GPT-2 [src]
- Regex Pattern:

```
's|'t|'re|'ve|'m|'ll|'d| ?\p{L}+| ?\p{N}+| ?[^s\p{L}\p{N}]+| \s+(?!S)|\s+
```

"Hello world I've learned to loveee! tokenization, I'M NOW THE #1Fan "

```
[ 'Hello', '_world', '_I', "'ve", '_learned', '_to', '_loveee',  
  '!', '_', '_tokeni', 'zation', ',', '_I', "'", 'M', '_NOW',  
  '_THE', '_#', '1', 'Fan', '_']
```

Tokenization at a Large Scale

- Case Study: GPT-2 [src]
- Regex Pattern:

```
's|'t|'re|'ve|'m|'ll|'d| ?\p{L}+| ?\p{N}+| ?[^s\p{L}\p{N}]+| \s+(?!S)|\s+
```

"Hello world I've learned to loveee! tokenization, I'M NOW THE #1Fan "

```
[ 'Hello', '_world', '_I', "'ve", '_learned', '_to', '_loveee',  
  '!', '_', '_tokeni', 'zation', ',', ',', '_I', "'", 'M', '_NOW',  
  '_THE', '_#', '1', 'Fan', '_']
```

Didn't set regex to ignore case

A bug in tokenization that was too costly to fix since they had already pre-trained the model!

Tokenization at a Large Scale

- Case Study: GPT-2 [src]
- Regex Pattern:

```
's|'t|'re|'ve|'m|'ll|'d| ?\p{L}+| ?\p{N}+| ?[^s\p{L}\p{N}]+|\s+(?!S)|\s+
```

"Hello world I've learned to loveee! tokenization, I'M NOW THE #1Fan "

```
[ 'Hello', '_world', '_I', "'ve", '_learned', '_to', '_loveee',  
  '!', '_', '_tokeni', 'zation', ',', '_I', "'", 'M', '_NOW',  
  '_THE', '_#', '1', 'Fan', '_']
```

- Convert all these chunks to bytes -> Do BPE till vocab size achieved

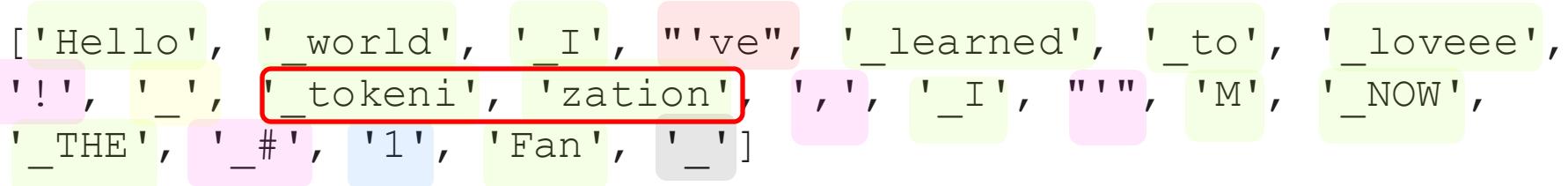
Tokenization at a Large Scale

- Case Study: GPT-2 [src]
- Regex Pattern:

```
's|'t|'re|'ve|'m|'ll|'d| ?\p{L}+| ?\p{N}+| ?[^s\p{L}\p{N}]+|\s+(?!S)|\s+
```

"Hello world I've learned to loveee! tokenization, I'M NOW THE #1Fan "

```
[ 'Hello', '_world', '_I', "'ve", '_learned', '_to', '_loveee',
  '!', '_', '_tokeni', 'zation', ',', '_', '_I', "'", 'M', '_NOW',
  '_THE', '_#', '1', 'Fan', '_' ]
```



The tokens are: 'Hello', '_world', '_I', "'ve", '_learned', '_to', '_loveee', '!', '_', '_tokeni', 'zation', ',', '_', '_I', "'", 'M', '_NOW', '_THE', '_#', '1', 'Fan', '_'. The tokens 'tokeni' and 'zation' are highlighted with a red rectangle.

- Wouldn't we just learn words this way? No, it'll still be sub-words:
world is common, 1 token; tokenization is uncommon so 'token' + 'ization'

Tokenization - BPE

- Iteratively, replace most common token pair with a single token
- Worked well for the small example, but we live in a multilingual complex world -> how to handle for that?
- Represent all data as bytes and then do BPE!
- Tokenization can have a huge impact!

Tokenization – BPE Issues

- GPT-2 was great at natural English language, not at other languages or even code – why?

Tokenization – BPE Issues

- GPT-2 was great at natural English language, not at other languages or even code but GPT 4 is – why?
- Each space was one token, overwhelming context window

Token count
32

GPT 2

```
for i in range(10):  
    for j in range (10):  
        print("Hello!")
```

Token count
20

GPT 4

```
for i in range(10):  
    for j in range(10):  
        print("Hello!")
```

Tokenization – BPE Issues

- Mandarin had >2x tokens; recall issue with char-level tokenization?

Model	Token count	Text
gpt2	15	Hello everyone, welcome to lecture 19 of 11-X85 class!
gpt-4	15	Hello everyone, welcome to lecture 19 of 11-X85 class!
gpt2	34	大家好，欢迎来到11-X85班第19讲！
gpt-4	22	大家好，欢迎来到11-X85课程的第19讲！

Tokenization – BPE Issues

- GPT 2 was bad at math - it didn't limit length on numbers, GPT 4 limited to 3 digits
- Struggled to learn multiplication mechanics this way

The image shows two side-by-side tokenizers. On the left is a 'gpt2' tokenizer, which shows a 'Token count' of 1 and the text '1999' below it. On the right is a 'gpt-4' tokenizer, which shows a 'Token count' of 2 and the text '1999' where the digits are split into two tokens: '19' and '99'. This demonstrates that GPT-4 limits the token length for numbers to three digits.

Tokenization – BPE Issues

- What if a token appears a lot in tokenization data but not in training data?
 - SolidGoldMagikarp -> LLM would go berzerk!
- LLM can't spell words
 - The infamous 'how many r's in strawberry'

Tokenization - BPE

- Iteratively, replace most common token pair with a single token
- Worked well for the small example, but we live in a multilingual complex world -> how to handle for that?
- Represent all data as bytes and then do BPE!
 - All data = all text on the internet!
- GPT-4 handled a lot of edge cases, starting with 258 tokens (256+2 special) and did 100K merges to have a 100,258 vocab - fixed for problems with GPT-2

Poll @1206

What is the primary harm of allowing a BPE algorithm to merge tokens across pre-tokenization boundaries (e.g., merging the 'o' from 'Hello' with the ' w' from ' world')?

- It would violate UTF-8 encoding rules and cause the model to crash.
- It would fill the vocabulary with semantically meaningless but frequent 'junk' tokens (like 'o w' or 'd t') and lead to inefficient, broken tokenization.
- It is not harmful and is a standard part of all modern tokenizers.
- It would make the tokenizer too slow to run during inference.

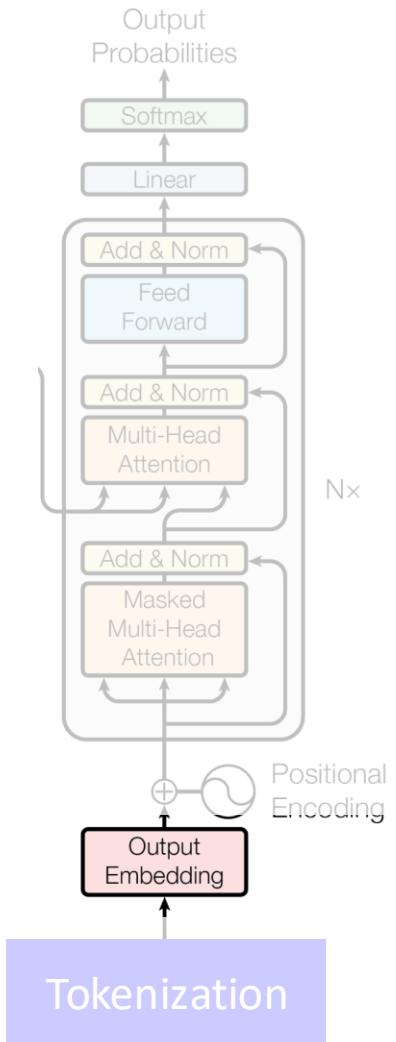
Poll @1206

What is the primary harm of allowing a BPE algorithm to merge tokens across pre-tokenization boundaries (e.g., merging the 'o' from 'Hello' with the ' w' from ' world')?

- It would violate UTF-8 encoding rules and cause the model to crash.
- It would fill the vocabulary with semantically meaningless but frequent 'junk' tokens (like 'o w' or 'd t') and lead to inefficient, broken tokenization.
- It is not harmful and is a standard part of all modern tokenizers.
- It would make the tokenizer too slow to run during inference.

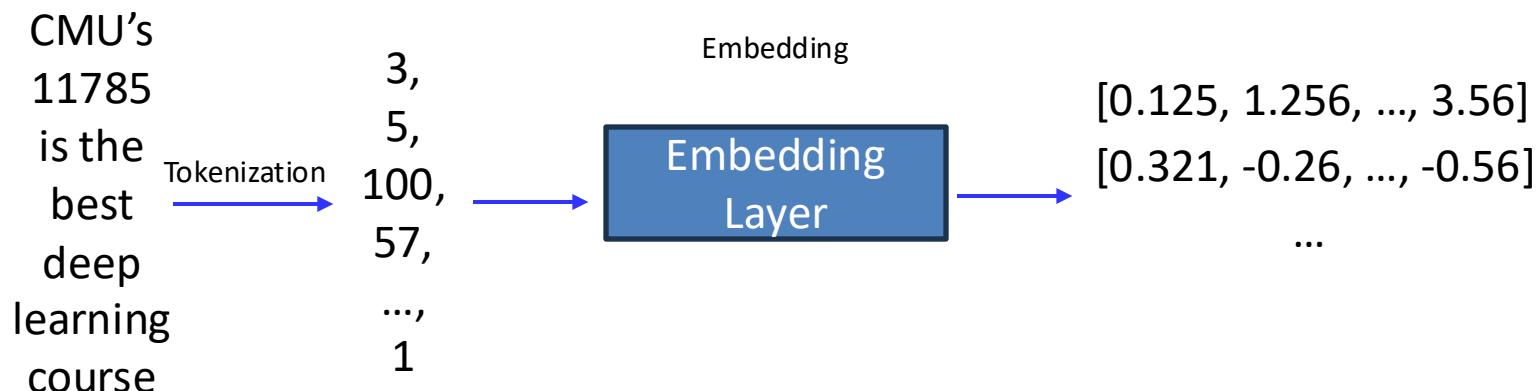
Transformer Architecture

- Word Tokenization
- **Word Embedding**
- (Masked) Multi-Head Attention
- Position Encoding
- Feed-Forward
- Add & Norm
- Output Projection Layer



Embedding

- Represents each discrete token index as continuous token embeddings



Embedding Layer is a Linear Layer

- $nn.Embedding$ is essentially a linear layer $Y = XW$

One-Hot Vector
Token Index $X \in \mathbb{R}^{L \times |V|}$

$$\begin{bmatrix} 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix}$$

Weight Matrix $W \in \mathbb{R}^{|V| \times D}$

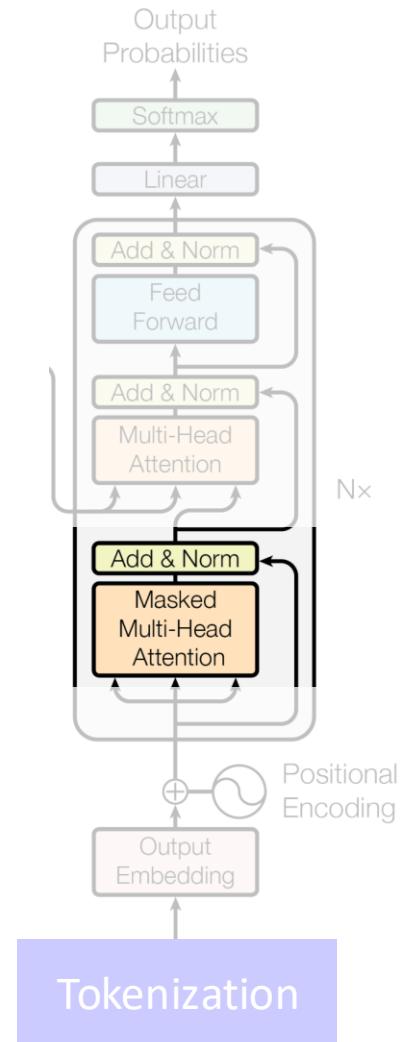
$$\begin{bmatrix} 0.235 & -1.256 & 3.513 & \dots & -0.187 \\ 1.291 & -2.012 & 0.624 & \dots & -1.291 \\ 0.535 & 0.012 & -0.024 & \dots & 2.345 \\ \dots & \dots & \dots & \dots & \dots \\ 0.131 & 2.102 & 0.935 & \dots & -0.125 \end{bmatrix}$$

$$\begin{bmatrix} 1.291 & -2.012 & 0.624 & \dots & -1.291 \\ 0.535 & 0.012 & -0.024 & \dots & 2.345 \\ 0.131 & 2.102 & 0.935 & \dots & -0.125 \end{bmatrix}$$

Token Embedding $Y \in \mathbb{R}^{L \times D}$

Transformer Architecture

- Word Tokenization
- Word Embedding
- **(Masked) Multi-Head Attention**
- Position Encoding
- Feed-Forward
- Add & Norm
- Output Projection Layer



Self-Attention

- Attention Operation

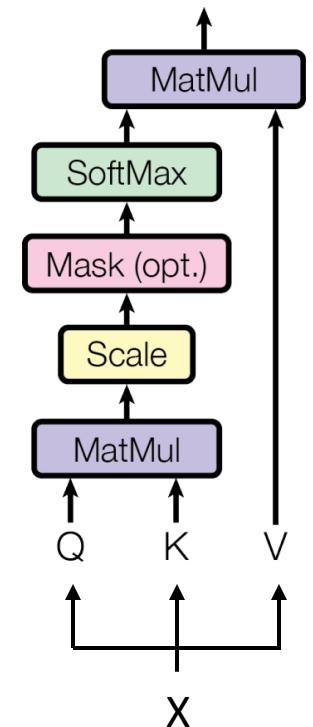
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Scaled Dot-Product Attention

- Query-Key-Value

- Linear affine from input X itself

- Weighted-sum of V based on similarity/correlation between Q and K
 - Each token's weights sum to one



Self-Attention

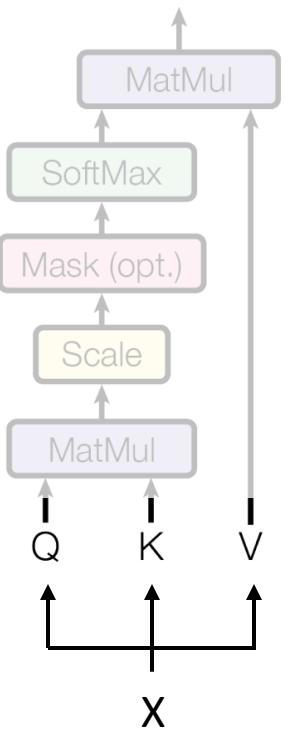
- Query-Key-Value from Three Linear Affine of X

$$\begin{array}{ccc} \mathbf{X} & & \mathbf{W}^Q \\ \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} & \times & \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} \\ & & = \\ & & \mathbf{Q} \end{array}$$

$$\begin{array}{ccc} \mathbf{X} & & \mathbf{W}^K \\ \mathbb{R}^{L \times D} & \times & \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} \\ & & = \\ & & \mathbf{K} \end{array}$$

$$\begin{array}{ccc} \mathbf{X} & & \mathbf{W}^V \\ \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} & \times & \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} \\ & & = \\ & & \mathbf{V} \end{array}$$

Scaled Dot-Product Attention



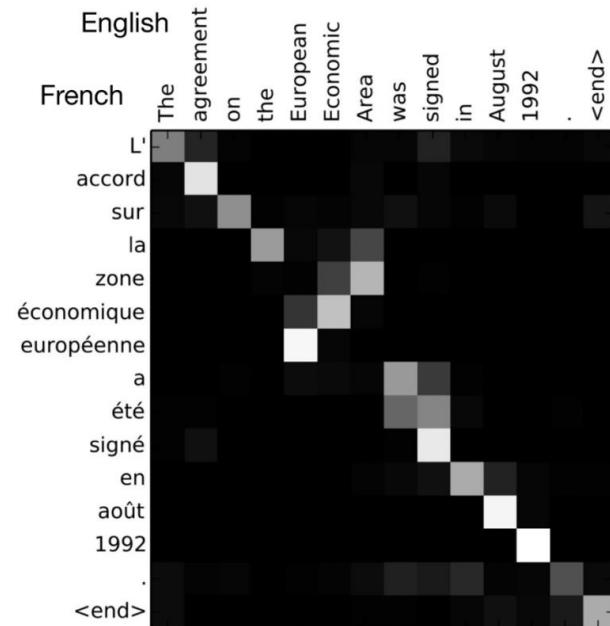
Self-Attention

- Attention weights

$$\text{softmax} \left(\frac{\begin{array}{c} \text{Q} \\ \times \\ \text{K}^T \end{array}}{\sqrt{d_k}} \right)$$

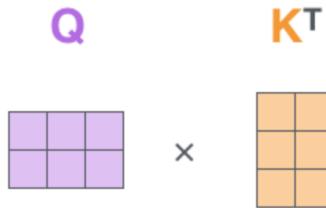
$\mathbb{R}^{L \times D}$

$\mathbb{R}^{L \times L}$



Self-Attention

- Why is the dot product not enough?

$$\begin{matrix} \mathbf{Q} & \mathbf{K^T} \\ \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} & \times \quad \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} \end{matrix}$$


The diagram illustrates the multiplication of two 3x3 matrices, Q and K^T. Matrix Q is represented by a 3x3 grid of purple squares, and matrix K^T is represented by a 3x3 grid of orange squares. A multiplication symbol 'x' is placed between the two matrices.

- Dot product could be arbitrary numbers
-

Self-Attention

- Why is the dot product not enough?

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

- Dot product could be arbitrary numbers
- We want a weighted average of all attention scores to get the appropriate $\mathbf{V} \Rightarrow \text{Softmax}$

Self-Attention

- That seems good enough, why do we need the $\text{sqrt}(d_k)$

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

The diagram illustrates the computation of attention weights. It shows two matrices, Q (purple) and K^T (orange), being multiplied together. The result is then passed through a softmax function to produce the attention scores.

Self-Attention

- That seems good enough, why do we need the $\text{sqrt}(d_k)$

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

The diagram shows the computation of attention weights. It consists of two parts: \mathbf{Q} and \mathbf{K}^T . \mathbf{Q} is represented by a 2x3 grid of purple squares, and \mathbf{K}^T is represented by a 3x2 grid of orange squares. A multiplication symbol (\times) is placed between them, indicating the dot product operation. Above the grids, the labels \mathbf{Q} and \mathbf{K}^T are written in purple and orange respectively. Below the multiplication symbol is a horizontal line.

- Softmax has a vanishing/exploding problem!

Self-Attention – why d_k

$$\begin{array}{|c|c|c|} \hline 2 & 1 & 3 \\ \hline 1 & -2 & -1 \\ \hline \end{array} \otimes \begin{array}{|c|c|} \hline 1 & 3 \\ \hline 2 & 1 \\ \hline 2 & -1 \\ \hline \end{array} = \begin{array}{|c|c|} \hline 10 & 4 \\ \hline -5 & 2 \\ \hline \end{array}$$

Self-Attention – why d_k

$$\text{Softmax} \left(\begin{array}{cc} 10 & 4 \\ -5 & 2 \end{array} \right) = \begin{array}{cc} 0.99 & 2e-3 \\ 9e-4 & 0.99 \end{array}$$

Self-Attention – why d_k

$$\text{Softmax} \left(\begin{array}{|c|c|} \hline 10 & 4 \\ \hline -5 & 2 \\ \hline \end{array} \right) = \begin{array}{|c|c|} \hline 0.99 & 2e-3 \\ \hline 9e-4 & 0.99 \\ \hline \end{array}$$

Softmax pushes disparity to extremes
=> vanishing gradients

Self-Attention – why d_k

$$\text{Softmax} \left(\begin{array}{cc} 10 & 4 \\ -5 & 2 \end{array} \right) / \sqrt{3}) \equiv \begin{array}{|c|c|} \hline 0.97 & 0.03 \\ \hline 0.02 & 0.98 \\ \hline \end{array}$$

Self-Attention – why d_k

$$\text{Softmax} \left(\begin{array}{cc} 10 & 4 \\ -5 & 2 \end{array} \right) / \sqrt{3}) \equiv \begin{array}{|c|c|} \hline 0.97 & 0.03 \\ \hline 0.02 & 0.98 \\ \hline \end{array}$$

Sqrt(d_k) is the std. dev. of the dot product

Scaling by that reduces the vanishing gradient!

Self-Attention

- Output

$$\text{softmax} \left(\frac{\begin{matrix} Q \\ \times \\ K^T \end{matrix}}{\sqrt{d_k}} \right) V = Z$$

The diagram illustrates the computation of self-attention. It shows the multiplication of two matrices, Q (purple) and K^T (orange), followed by division by $\sqrt{d_k}$, and finally the multiplication by matrix V (blue). The result is labeled Z (pink).

The diagram illustrates the computation of self-attention. It shows the multiplication of two matrices, Q (purple) and K^T (orange), followed by division by $\sqrt{d_k}$, and finally the multiplication by matrix V (blue). The result is labeled Z (pink).

Multi-Head Self-Attention

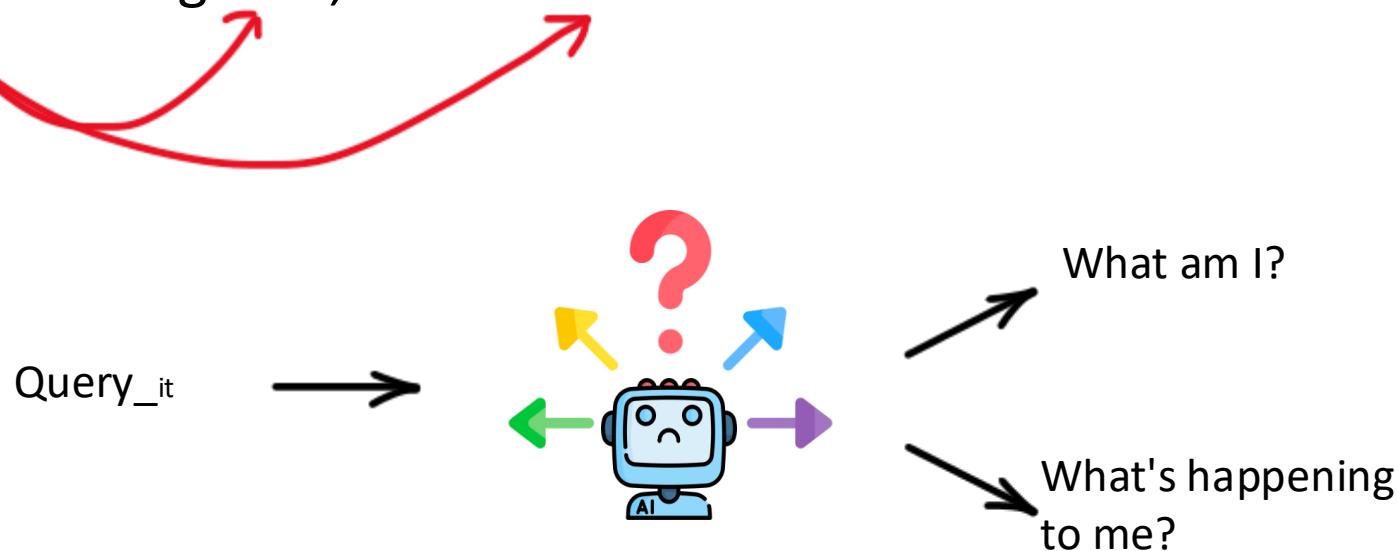
- Why do we need multiple self-attention operations?
Take the sentence:

- "I put the glass on the table, but **it** broke."
- "**it**" is the "glass", it is also "broke"



Multi-Head Self-Attention

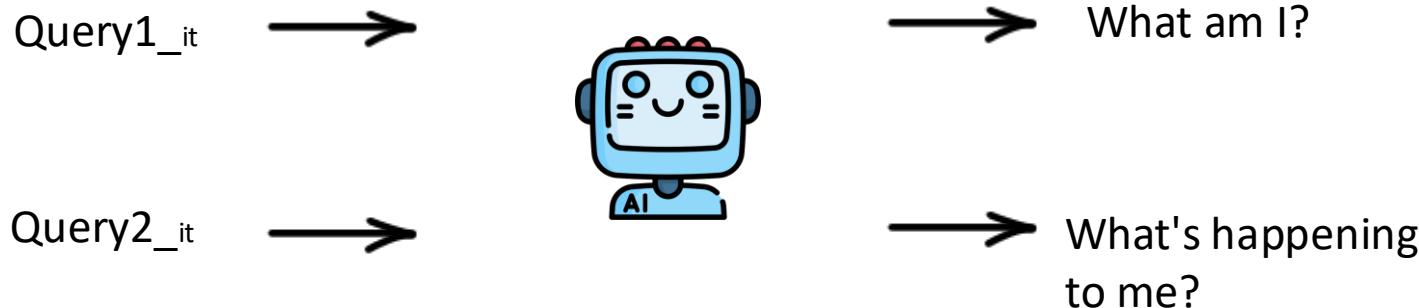
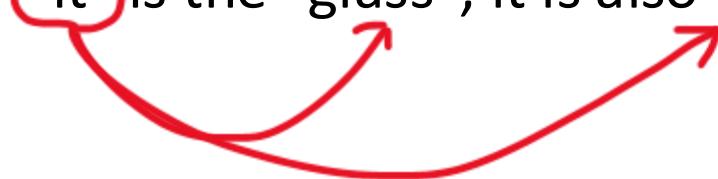
- Why do we need multiple self-attention operations?
Take the sentence:
- "I put the glass on the table, but **it** broke."
- "**it**" is the "glass", it is also "broke"



Multi-Head Self-Attention

- Why do we need multiple self-attention operations?
Take the sentence:

- "I put the glass on the table, but **it** broke."
- "**it**" is the "glass", it is also "broke"



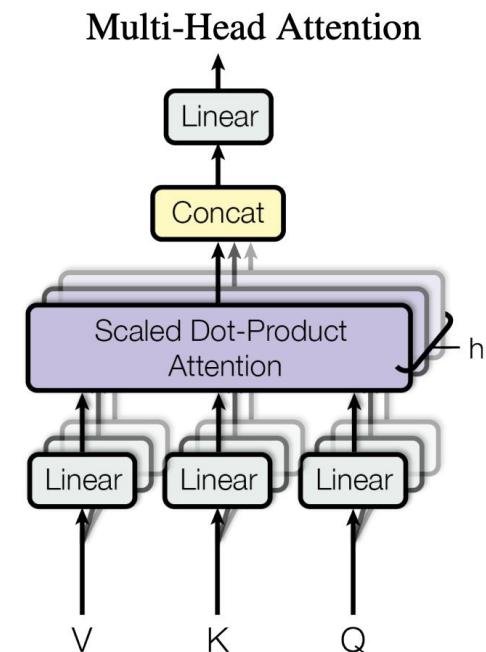
Multi-Head Self-Attention

- Multiple self-attention operations over the channel dimension

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

where $\text{head}_i = \text{Attention}\left(QW_i^Q, KW_i^K, VW_i^V\right)$

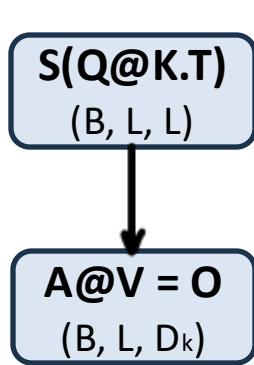
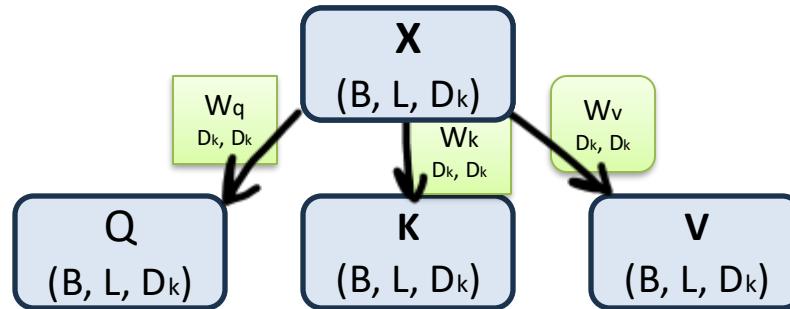
- Advantages:
 - Different attention maps capture different relationships
 - Parallel processing



(F.scaled_dot_product_attention)

SHA vs MHA

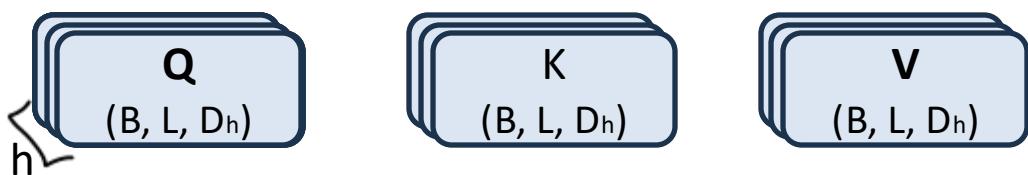
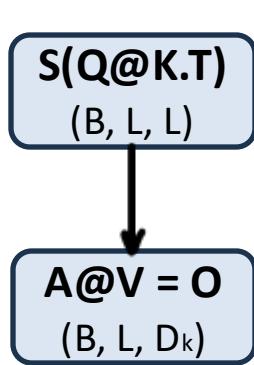
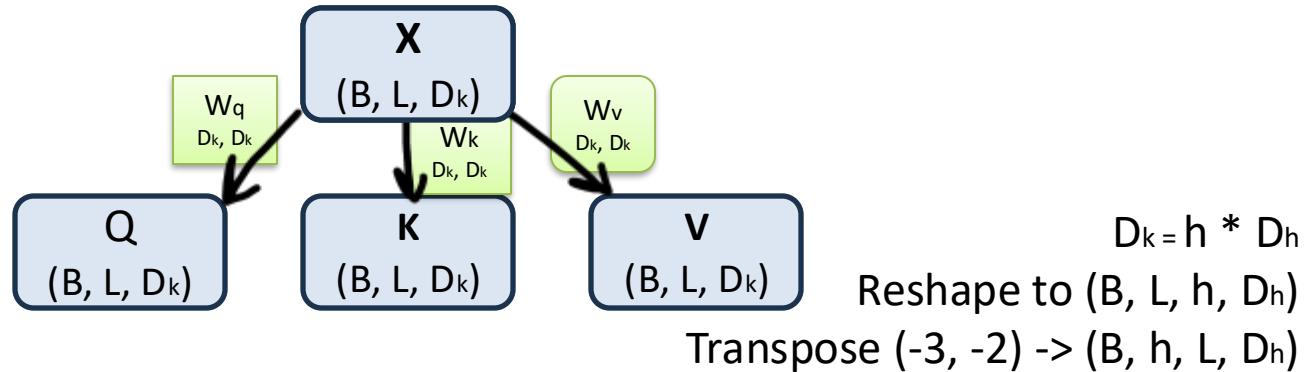
(torch.nn.MultiheadAttention)



(F.scaled_dot_product_attention)

SHA vs MHA

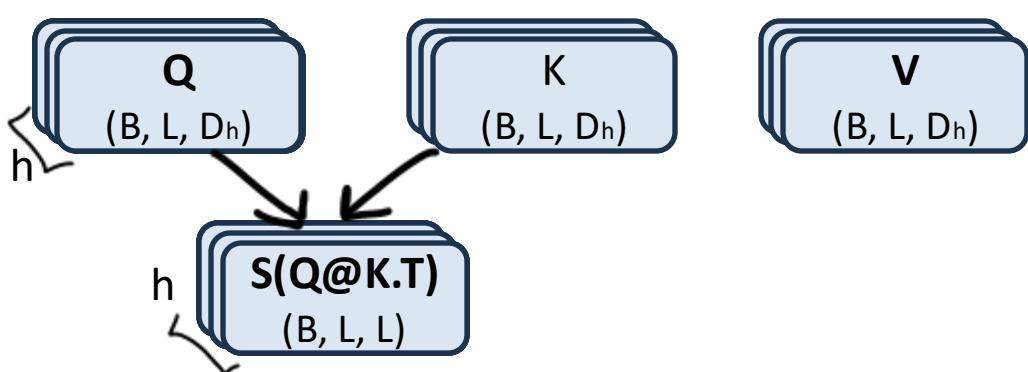
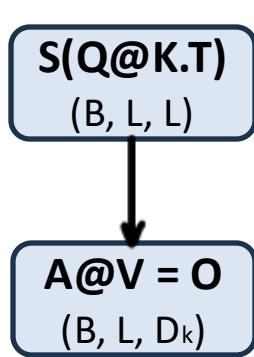
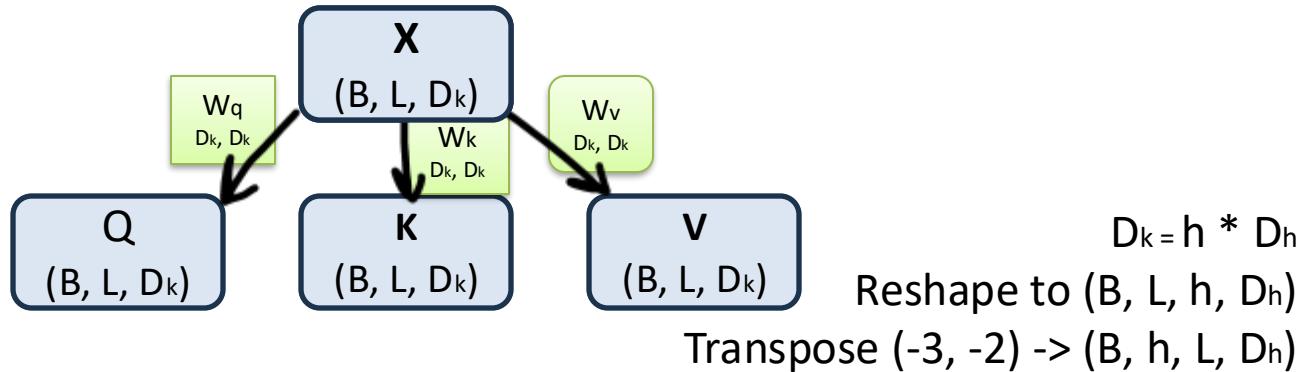
(torch.nn.MultiheadAttention)



(F.scaled_dot_product_attention)

SHA vs MHA

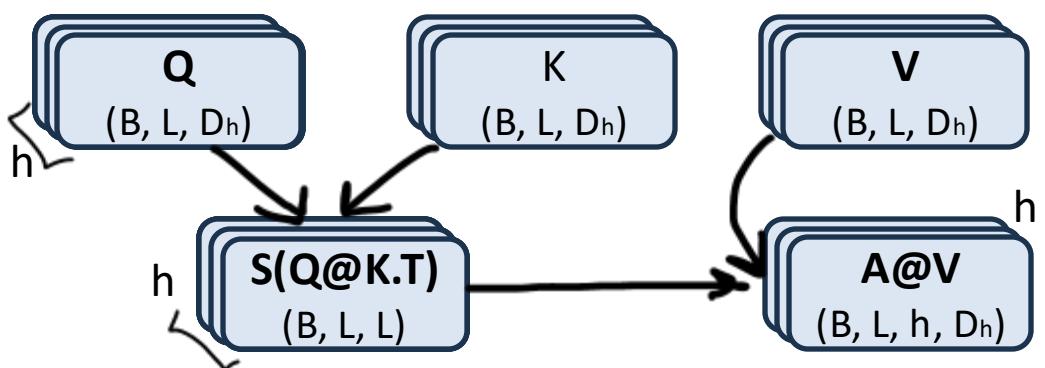
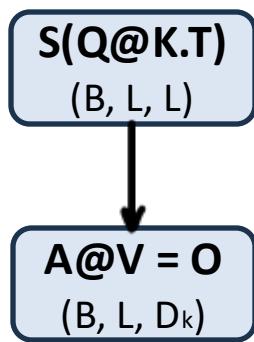
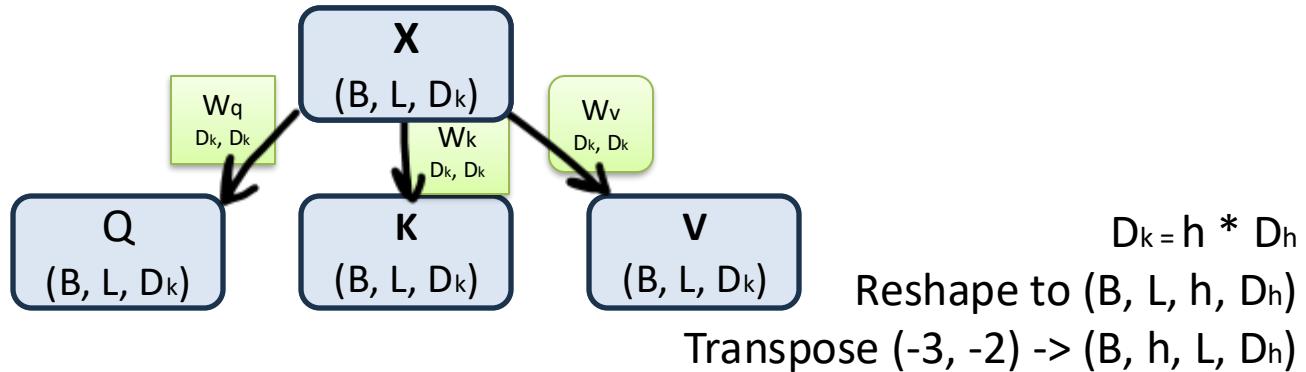
(torch.nn.MultiheadAttention)



(F.scaled_dot_product_attention)

SHA vs MHA

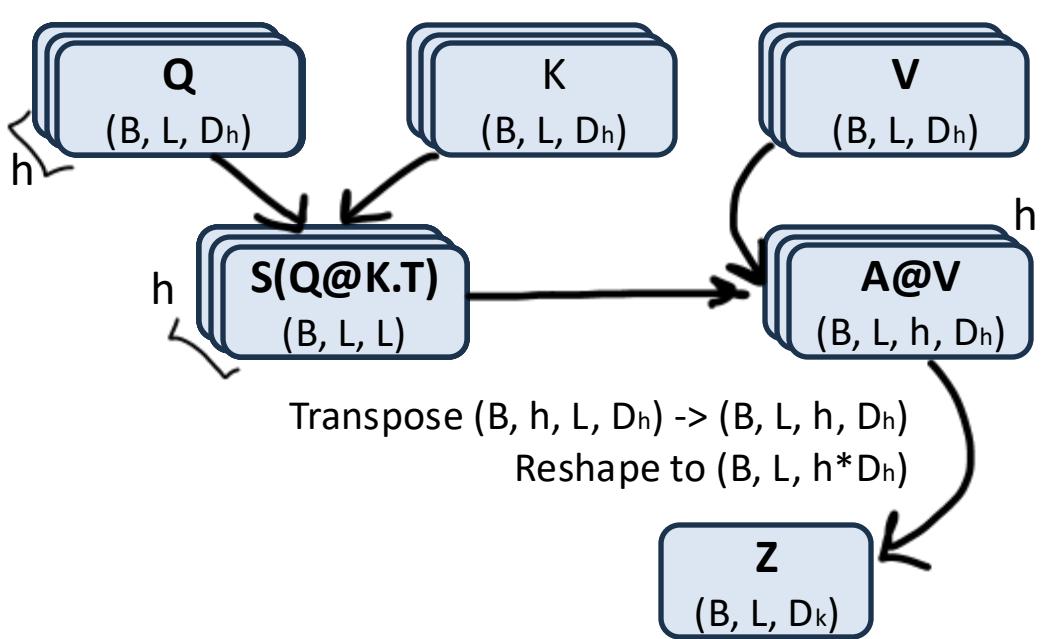
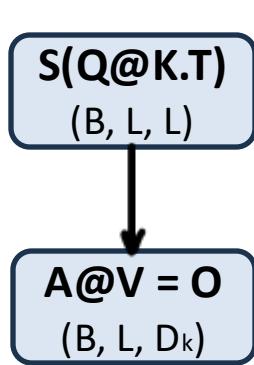
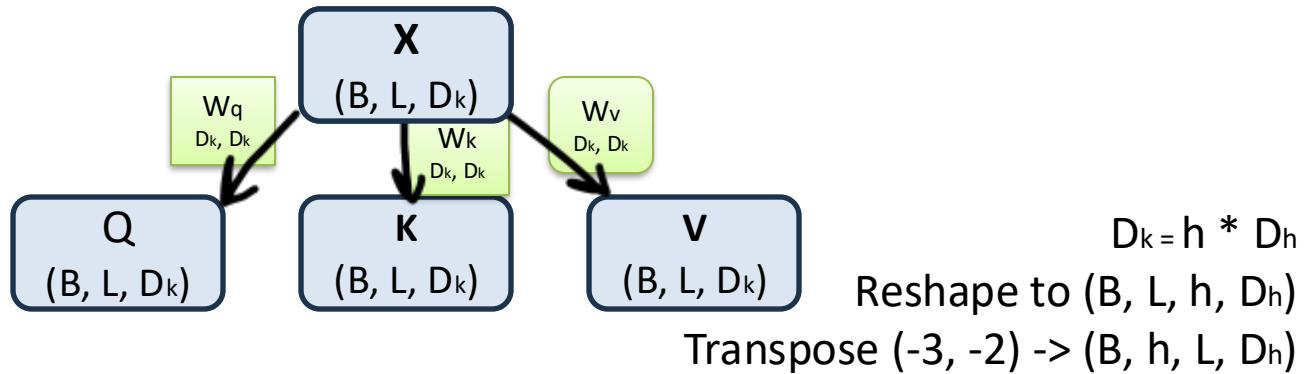
(torch.nn.MultiheadAttention)



(F.scaled_dot_product_attention)

SHA vs MHA

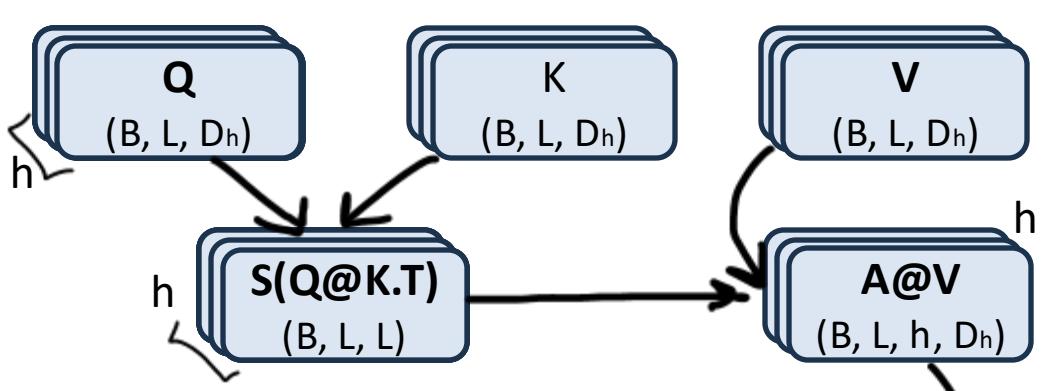
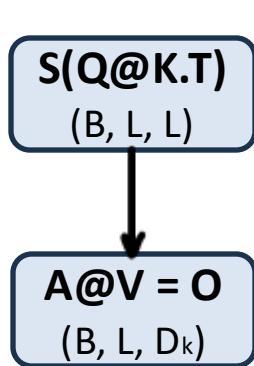
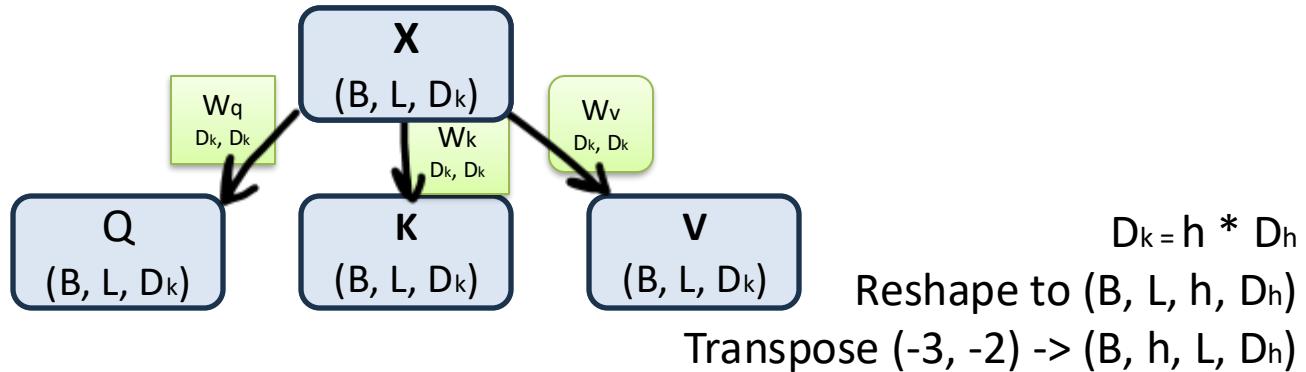
(torch.nn.MultiheadAttention)



(F.scaled_dot_product_attention)

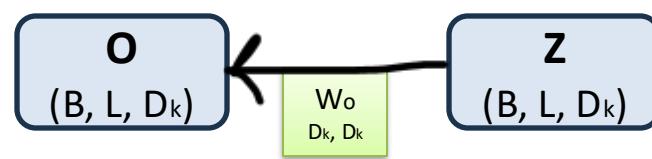
SHA vs MHA

(torch.nn.MultiheadAttention)

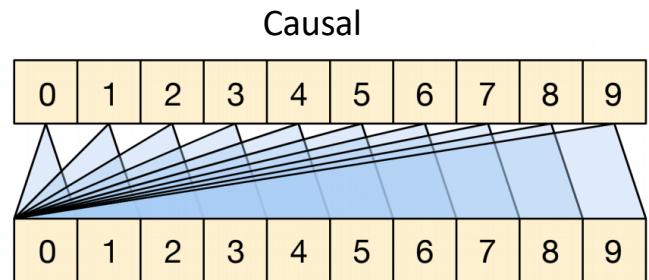
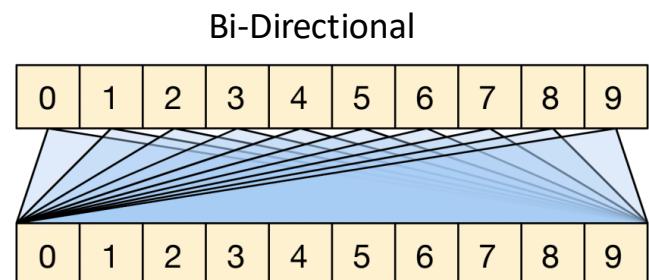
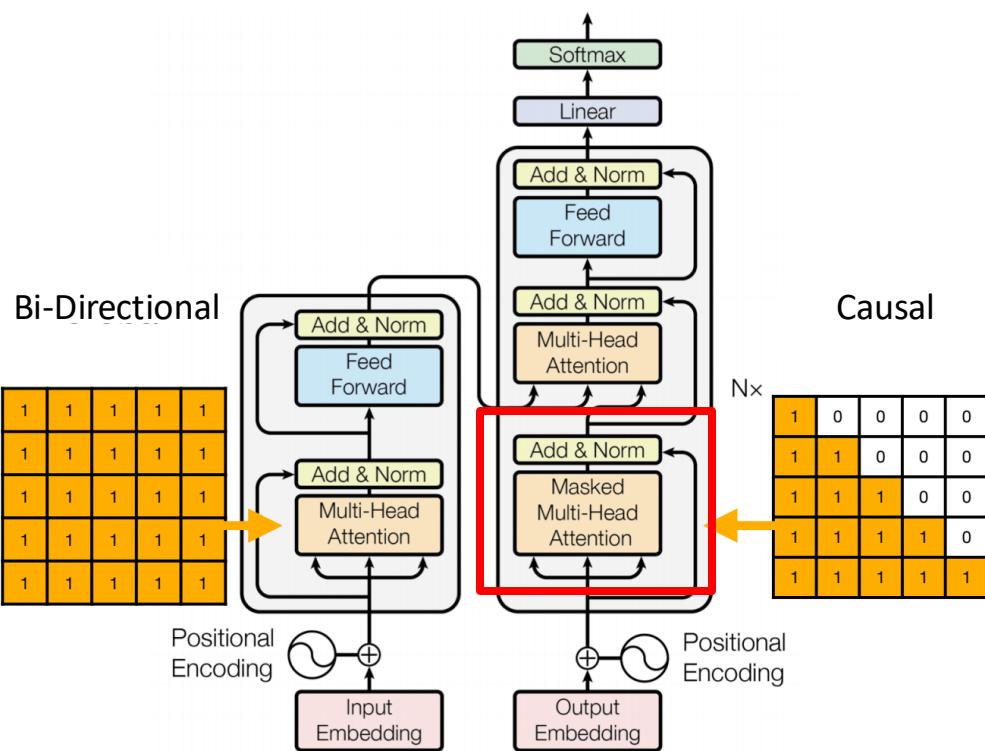


Transpose $(B, h, L, D_h) \rightarrow (B, L, h, D_h)$

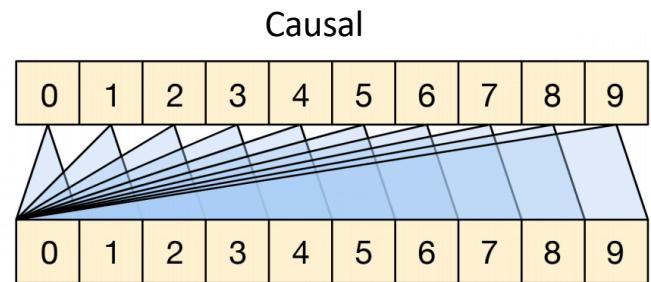
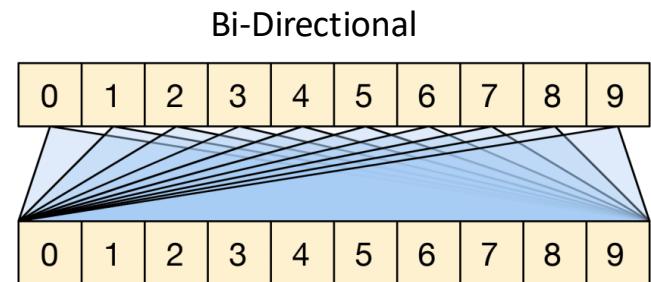
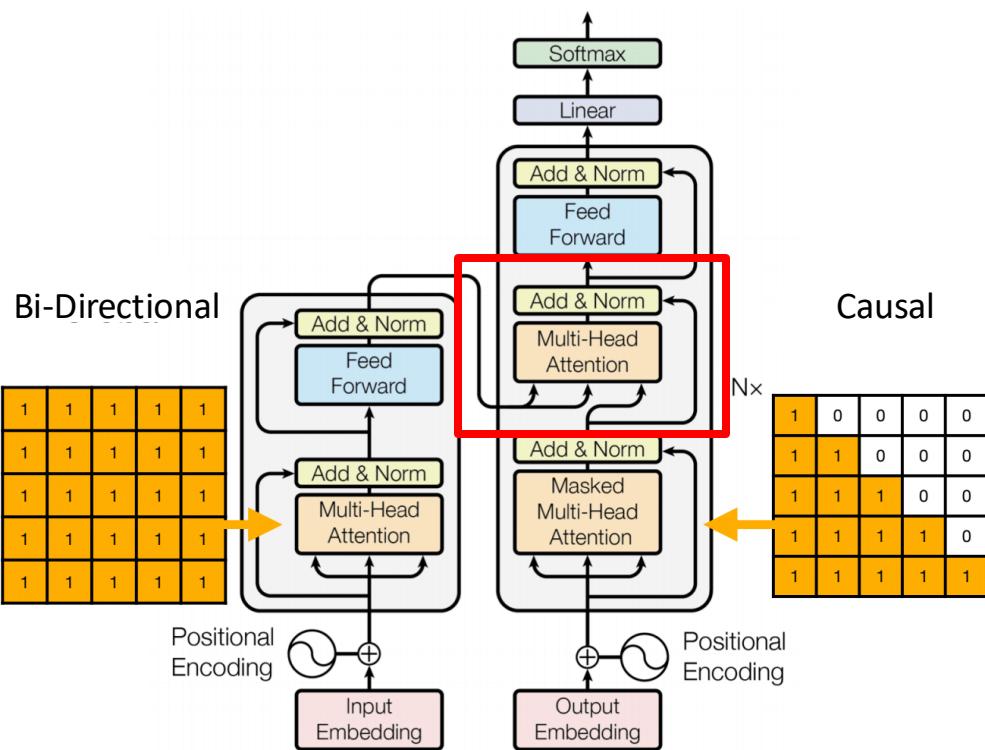
Reshape to $(B, L, h * D_h)$



Attention Masking

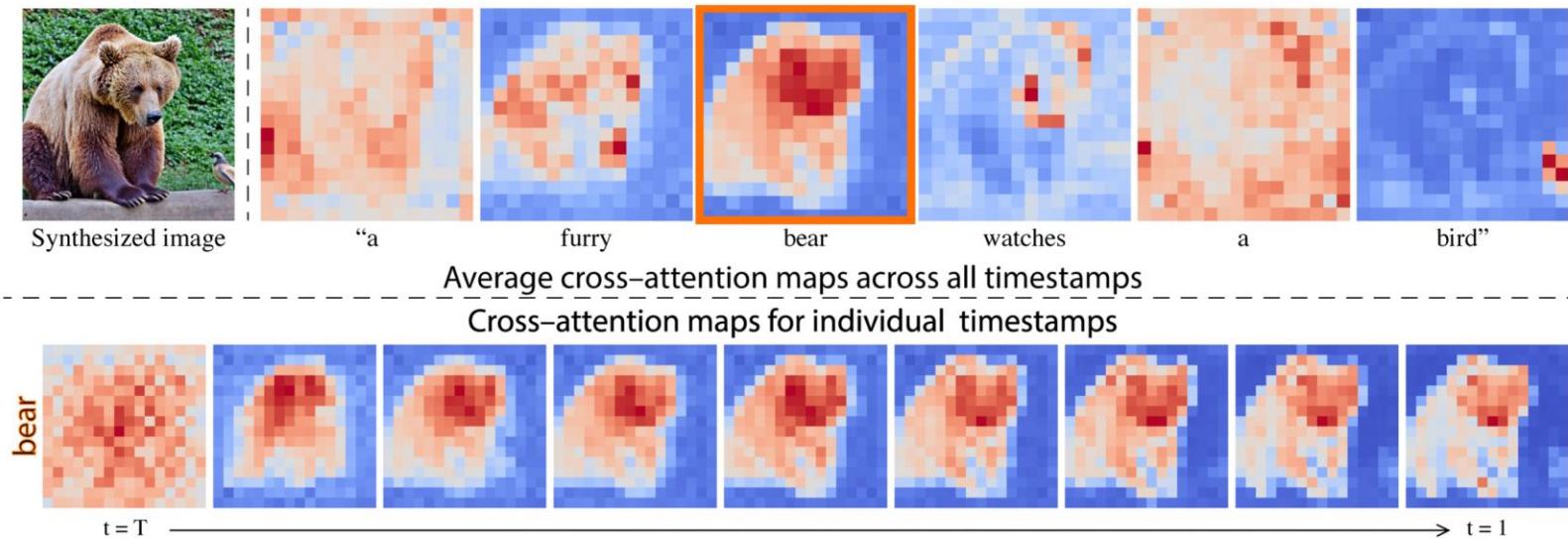


Attention Masking



Cross-Attention

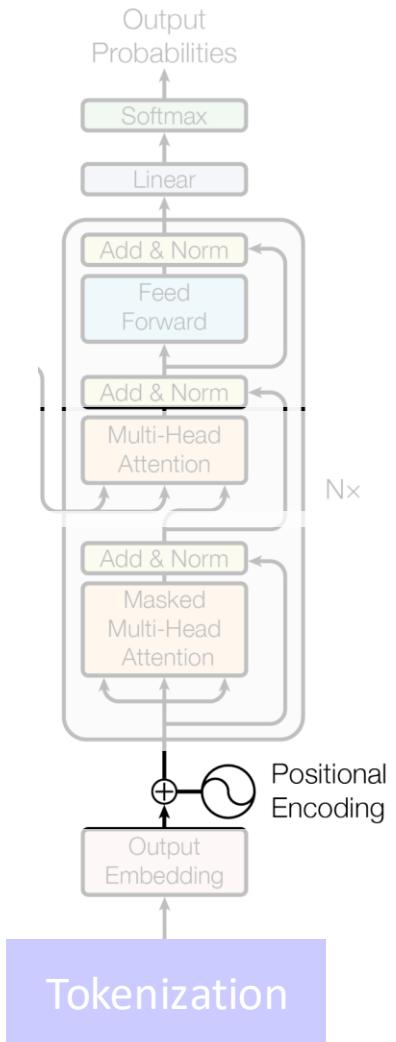
query, "a furry bear watches a bird."



The model iteratively denoise the noise vector based on the given text query to generate an Image

Transformer Architecture

- Word Tokenization
- Word Embedding
- (Masked) Multi-Head Attention
- **Position Encoding**
- Feed-Forward
- Add & Norm
- Output Projection Layer



Position Encoding

- Why do we need them?
 - Self-attention is permutation-invariant!
 - "James Bond signed a Bond with the bank"
 - Both 'Bond' will have same embeddings and then equal attention weight -> that can't be right, no?
 - We need something that gives idea of where the token is positioned in the context

Position Encoding

- Position encoding = P_t
 - Vaswani et al. 2017 decided to add it to X
 - So, $Q_m \text{pe} = Q_m + P_m$; same for $K_m \text{pe}$

$$(Q_m + P_m) \cdot (K_n + P_n) = Q_m \cdot K_n + \boxed{Q_m P_n + K_n P_m}$$

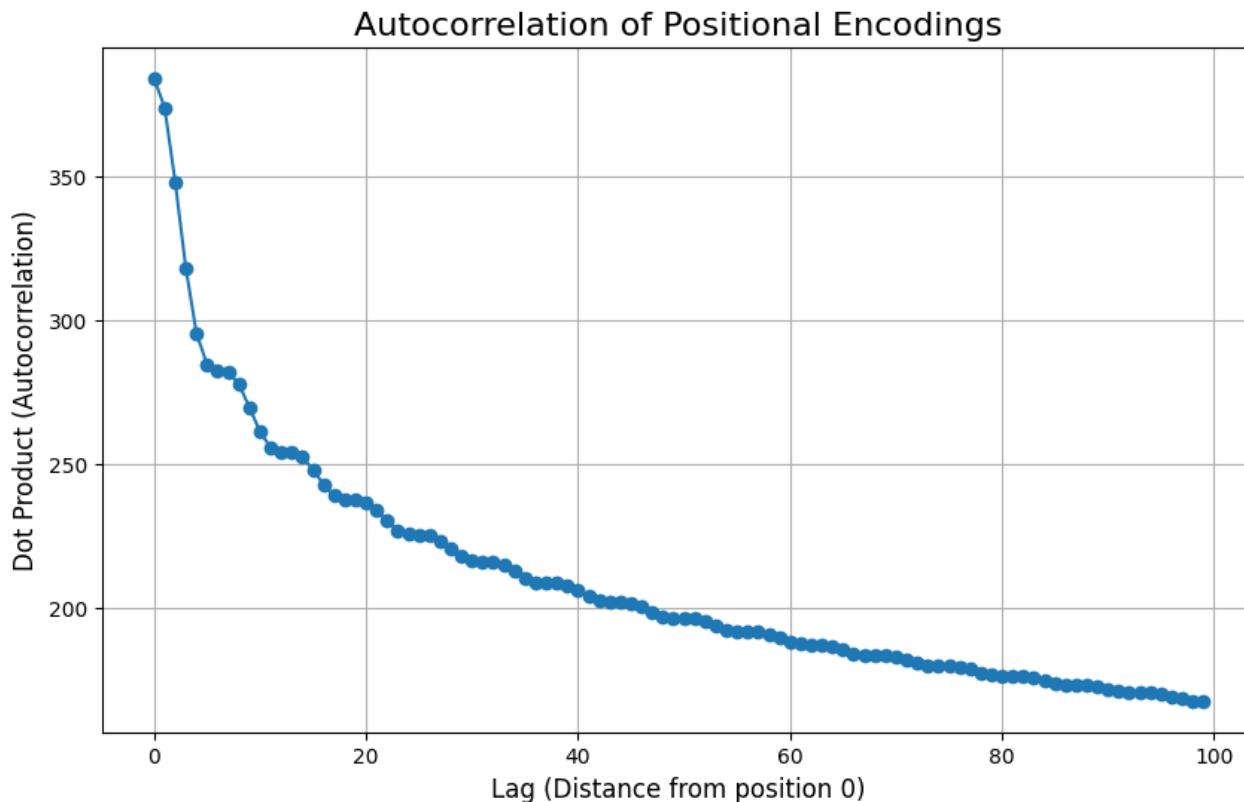
helpful signal noise

Position Encoding

- Position encoding= P_t
 - The $P_m @ P_n$ term is used to give positional insight

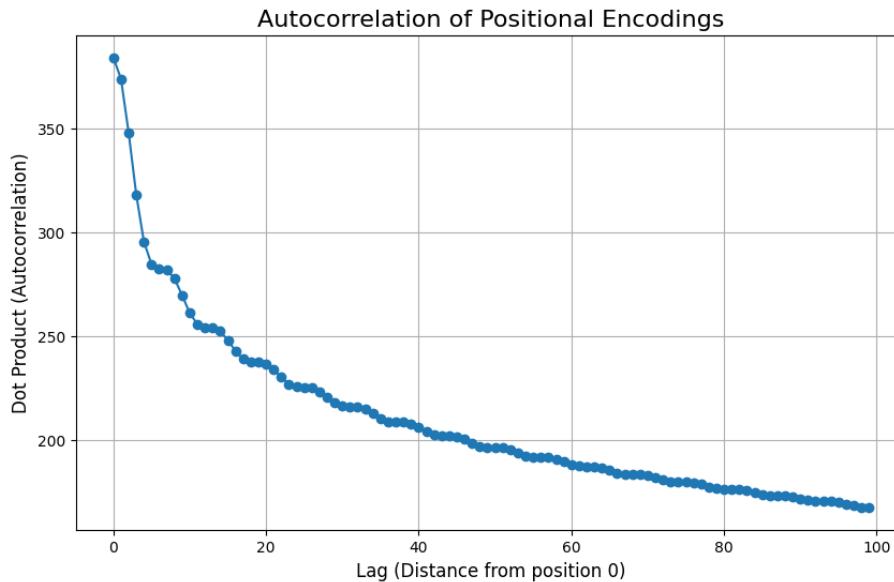
Position Encoding

- Position embedding = P_t
 - Actual autocorrelation curve for Vaswani et al 2017.



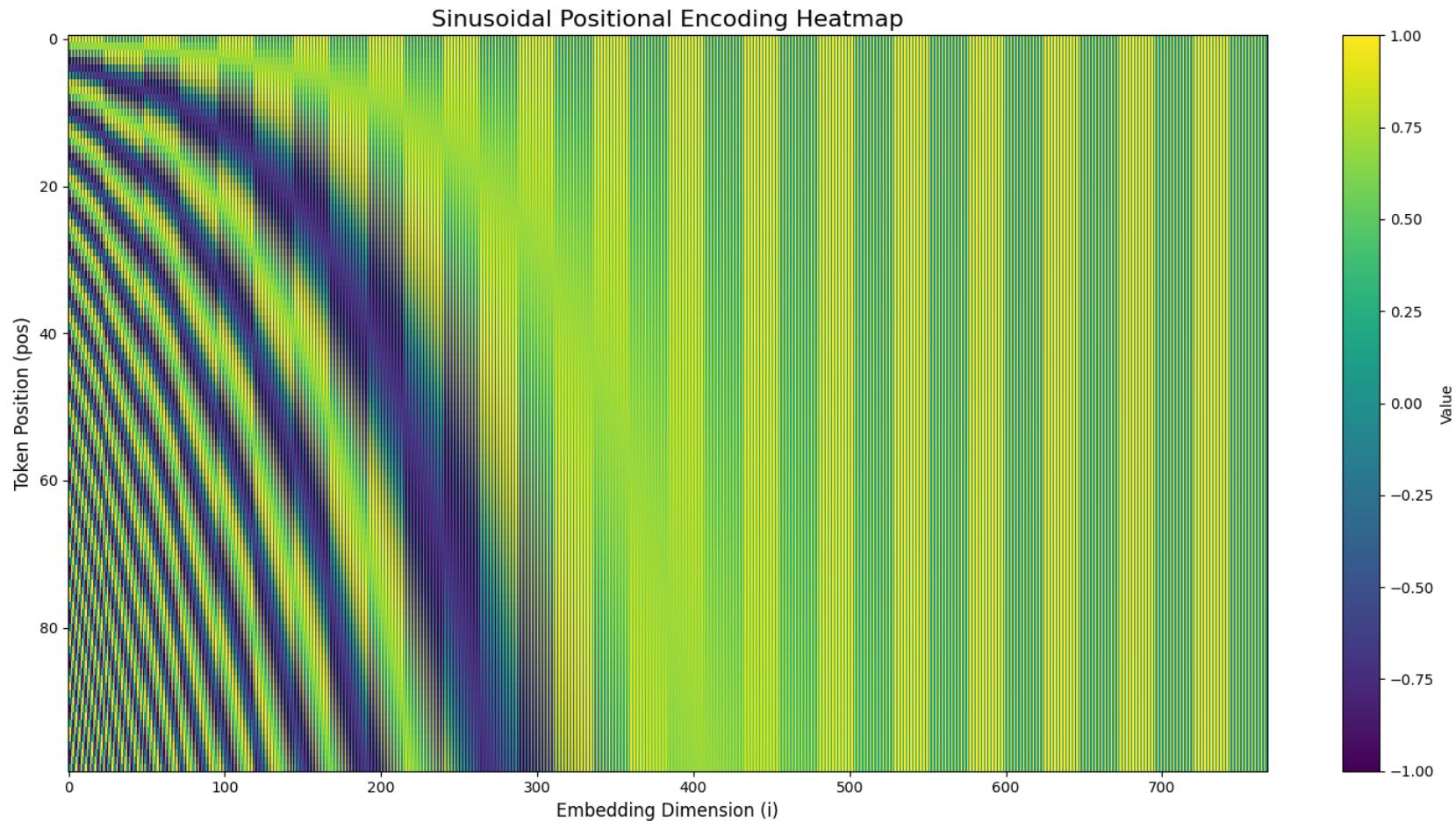
Position Encoding

- Position encoding= P_t
 - Attention gets relative positional insight from $P_m @ P_n$
 - M and N terms will be some distance apart, this dot product essentially an autocorrelation term.
 - If it falls smoothly with distance, we embed bias that closer tokens are more important and farther ones aren't



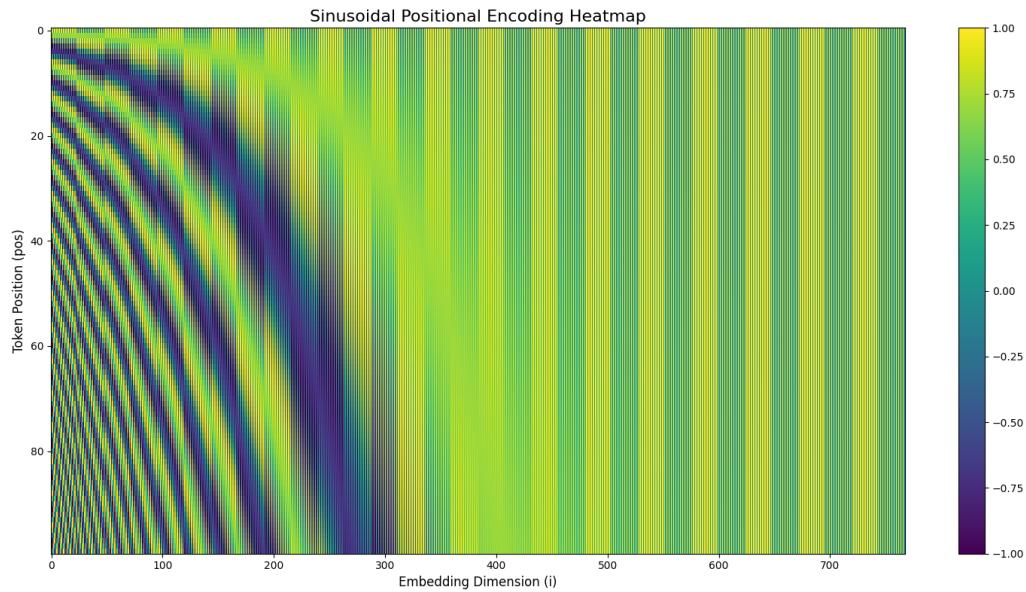
Positional Encoding

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



Positional Encoding

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



- The vector is made of sines and cosines of a harmonic series of frequencies
- All positions get unique values, it never repeats!

Poll @1207

A set of positional encodings (let's call them \mathbf{P}_t for position t) must effectively encode both absolute and relative position. Which of the following statements incorrectly describes a necessary property or mechanism?

- The dot product $\mathbf{P}_t \cdot \mathbf{P}_{t'}$ must be a mathematical component of the final $\mathbf{Q}_t \cdot \mathbf{K}_{t'}$ attention score, allowing the model to factor in position when calculating relevance.
- The embeddings must be designed so that the relationship between any two positions (e.g., $\mathbf{P}_t \cdot \mathbf{P}_{t+k}$) depends only on the offset k , not on the absolute position t .
- To ensure all positions are unique, the embeddings for any two different positions (\mathbf{P}_t and $\mathbf{P}_{t'}$) must be perfectly orthogonal (i.e., $\mathbf{P}_t \cdot \mathbf{P}_{t'} = 0$).

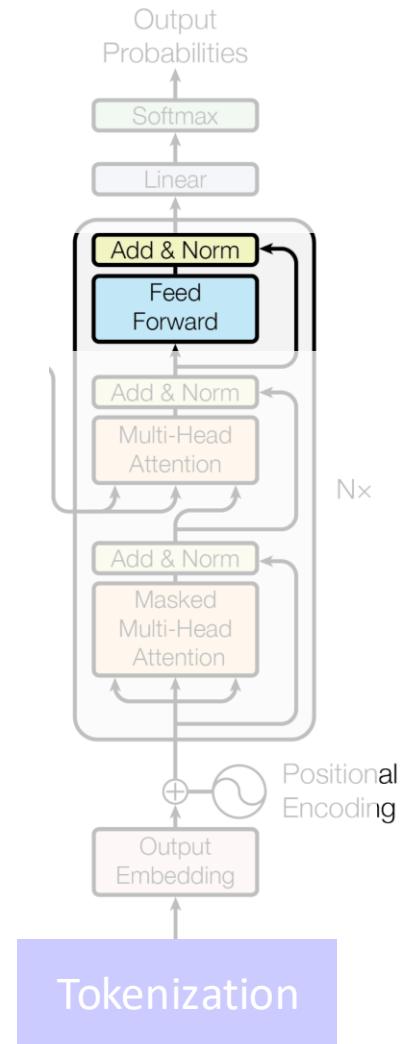
Poll @1207

A set of positional encodings (let's call them P_t for position t) must effectively encode both absolute and relative position. Which of the following statements correctly describes a necessary property or mechanism?

- The dot product $P_t \cdot P_{t'}$ must be a mathematical component of the final $Q_t \cdot K_{t'}$ attention score, allowing the model to factor in position when calculating relevance.
- The embeddings must be designed so that the relationship between any two positions (e.g., $P_t \cdot P_{t+k}$) depends only on the offset k , not on the absolute position t .
- To ensure all positions are unique, the embeddings for any two different positions (P_t and $P_{t'}$) must be perfectly orthogonal (i.e., $P_t \cdot P_{t'} = 0$).

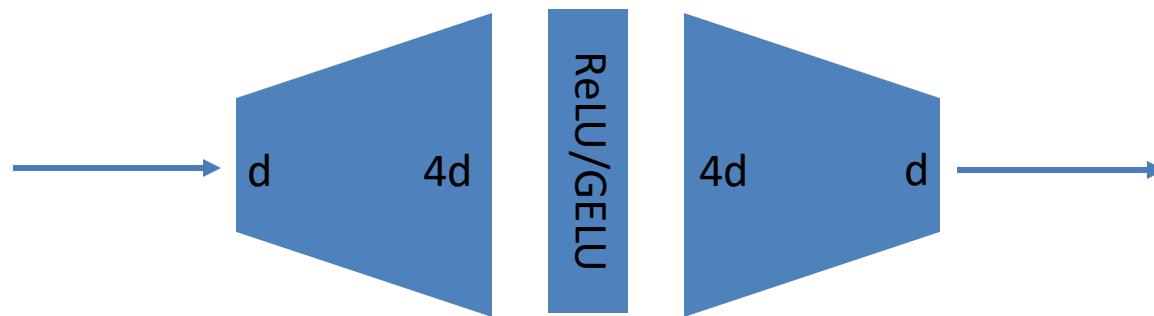
Transformer Architecture

- Word Tokenization
- Word Embedding
- (Masked) Multi-Head Attention
- Position Encoding
- **Feed-Forward**
- Add & Norm
- Output Projection Layer



Feed-Forward Block

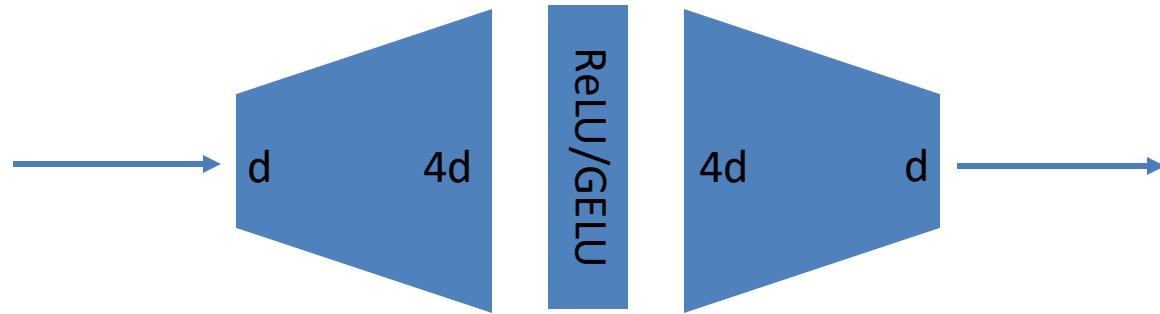
- Just a MLP!



- This is where the main memorization happens!

Feed-Forward Block

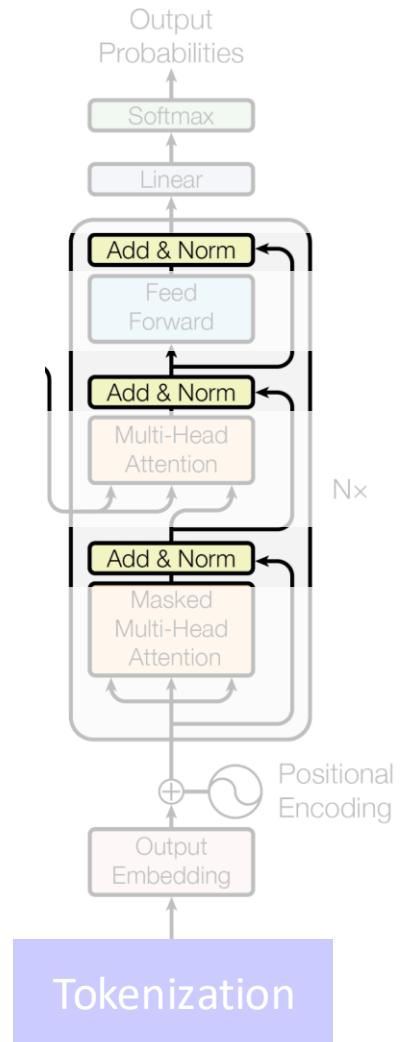
- Just a MLP!



- Why expand: to let it understand context better ->4x more non-linearities
- More connections created -> more 'memory'

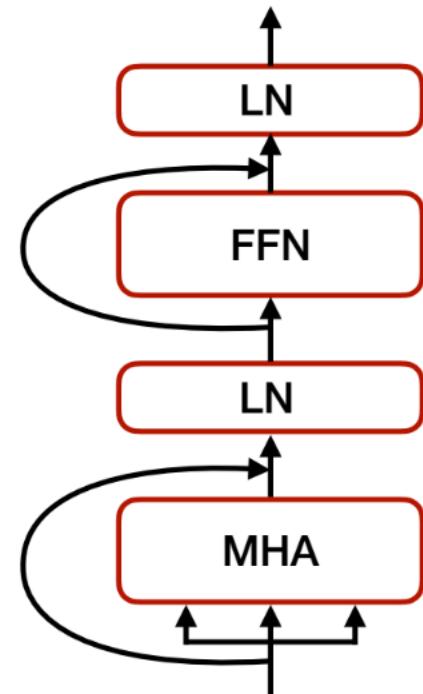
Transformer Architecture

- Word Tokenization
- Word Embedding
- (Masked) Multi-Head Attention
- Position Encoding
- Feed-Forward
- **Add & Norm**
- Output Projection Layer



Residual and Normalization

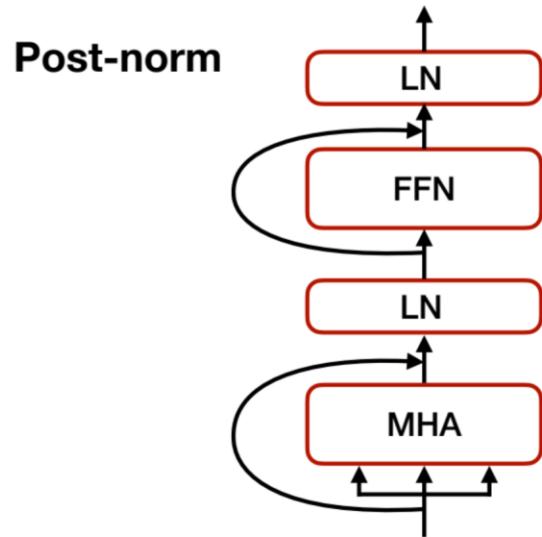
- Each layer in Transformer has:
 - A residual connection
 - A normalization layer
- Layer Norm. normalize each token by its embedding size dimension
 - For more stable training



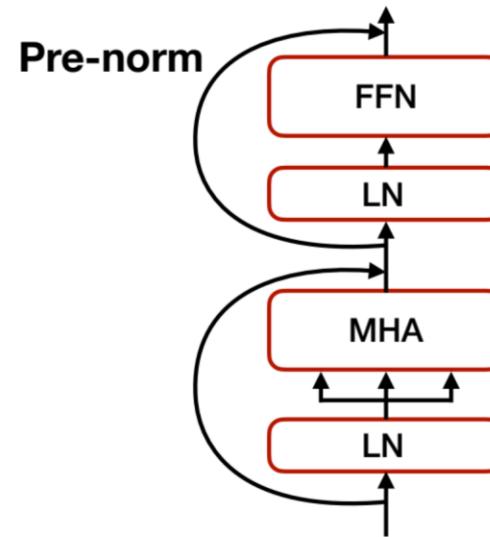
Position of Normalization

- Post-Norm vs Pre-Norm

$$\mathbf{x}_{t+1} = \text{Norm}(\mathbf{x}_t + F_t(\mathbf{x}_t))$$



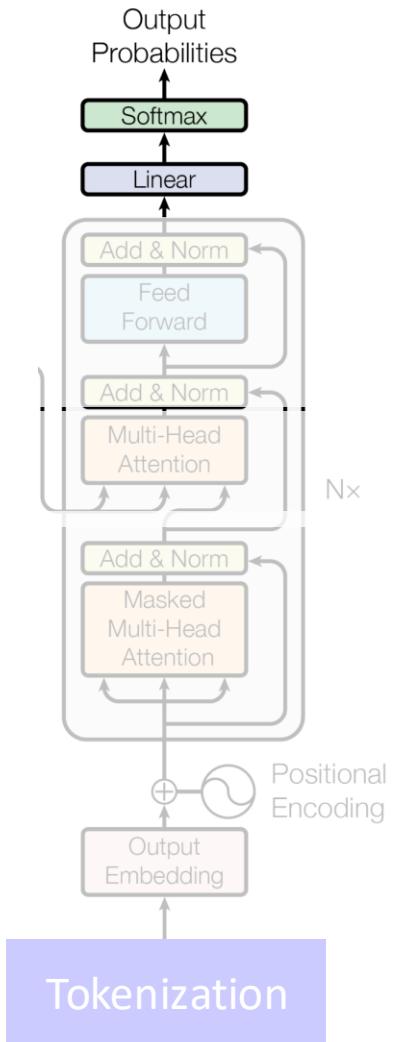
$$\mathbf{x}_{t+1} = \mathbf{x}_t + F_t(\text{Norm}(\mathbf{x}_t))$$



- Pre-Norm is easier and more stable to train
- Post-Norm tends to present better performance if properly trained

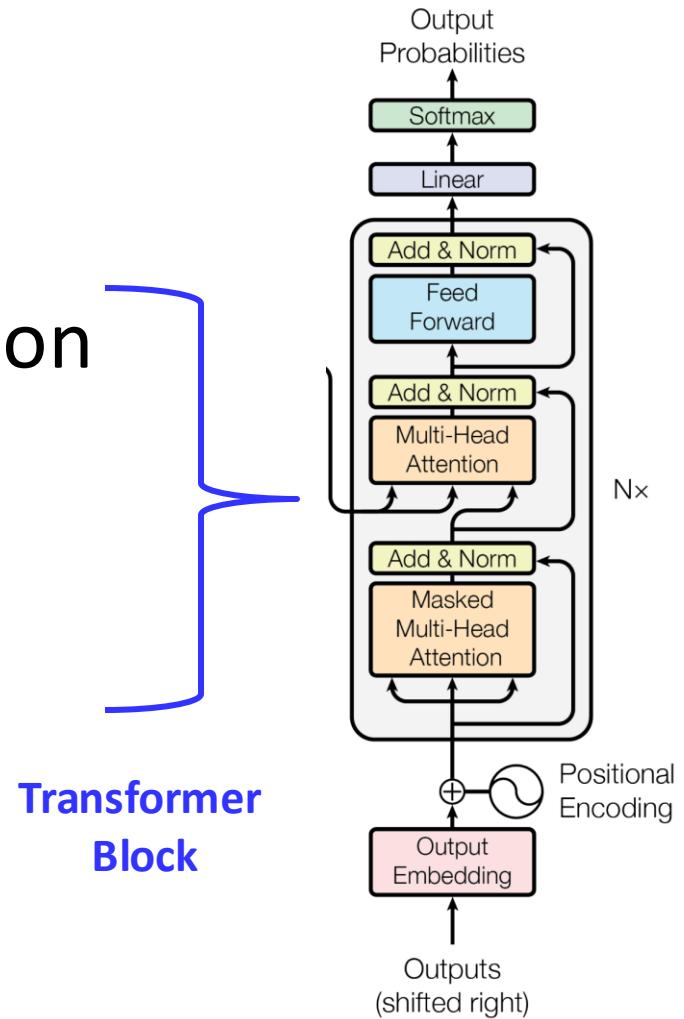
Transformer Architecture

- Word Tokenization
- Word Embedding
- (Masked) Multi-Head Attention
- Feed-Forward
- Add & Norm
- Position Encoding
- **Output Projection Layer**
 - Just a linear layer



Putting Them Together - Transformer

- Word Tokenization
- Word Embedding
- (Masked) Multi-Head Attention
- Position Encoding
- Feed-Forward
- Add & Norm
- Output Projection Layer



Poll @1208

Which of the following are true about self-attention?

- Self-attention is permutation invariant without position information
- The attention weights are scaled by the dimension d before computing softmax
- The attention weights are scaled by \sqrt{d} before computing softmax
- Q , K , V are affine transformations of input X

Poll @1208

Which of the following are true about self-attention?

- Self-attention is permutation invariant without position information
- The attention weights are scaled by the dimension d before computing softmax
- The attention weights are scaled by \sqrt{d} the dimension d before computing softmax
- Q, K, V are affine transformations of input X

Content

- Transformer Architecture
- **Improvements on Transformers**
- Transformer for different modalities
- Parameter Efficient Tuning
- Scaling Laws
- Quantizing Transformers
- Interpreting Transformer – attention, logit lens

Overview

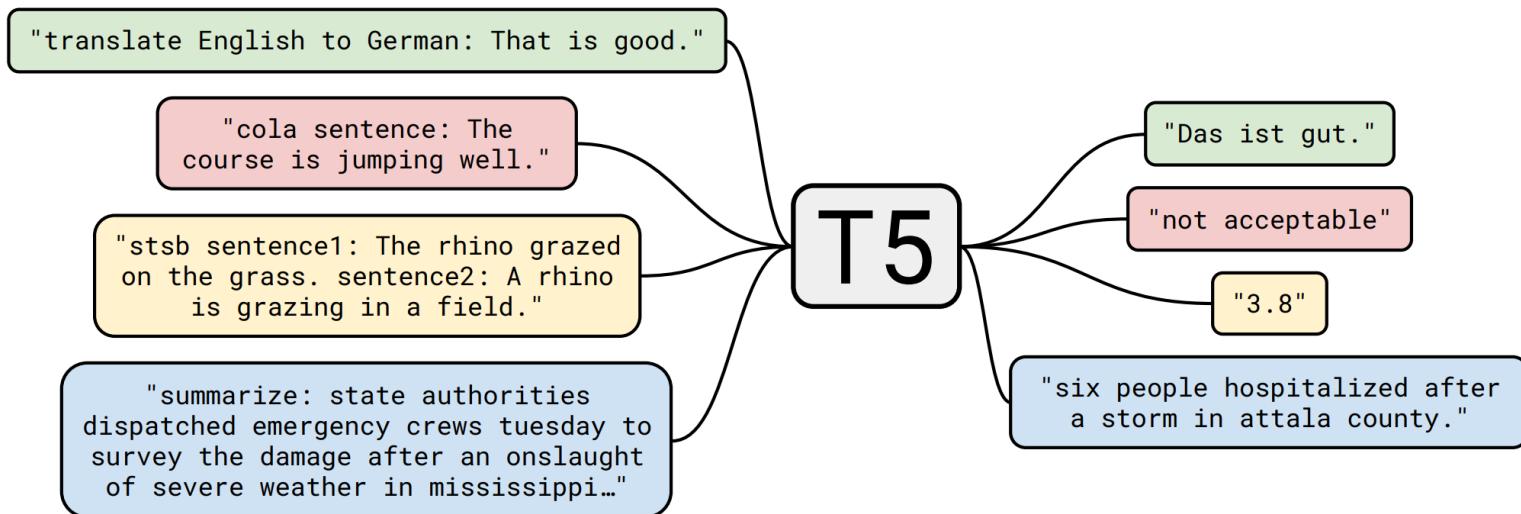
- Architecture
 - Encoder-Decoder
 - Encoder-Only
 - Decoder-Only
- Position Encoding
 - Relative Position Encoding
 - Rotary Position Encoding
- Efficient Attention Mechanism
 - Grouped Query Attention
 - Multi Query Attention
 - Flash Attention
 - Multi-head Latent Attention

Overview

- Architecture
 - Encoder-Decoder
 - Encoder-Only
 - Decoder-Only
- Position Encoding
 - Relative Position Encoding
 - Rotary Position Encoding
- Efficient Attention Mechanism
 - Grouped Query Attention
 - Multi Query Attention
 - Flash Attention
 - Multi-head Latent Attention

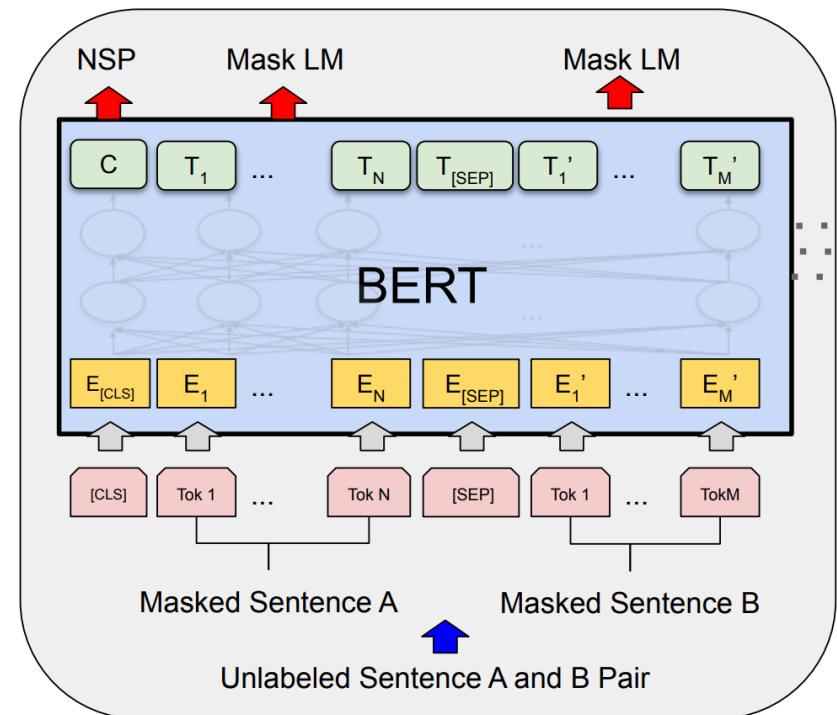
Encoder-Decoder - T5

- Encoder-Decoder architecture as in the original transformer paper
- A text-to-text model on various NLP tasks



Encoder-Only - BERT

- Bidirectional Encoder Representations from Transformers (BERT)
 - Encoder-only arch.
- Trained with
 - Mask token prediction
 - Next sentence prediction



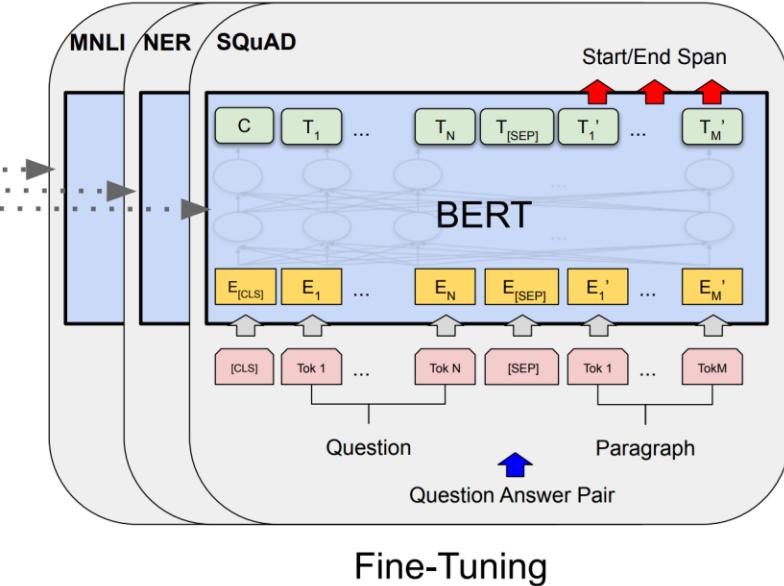
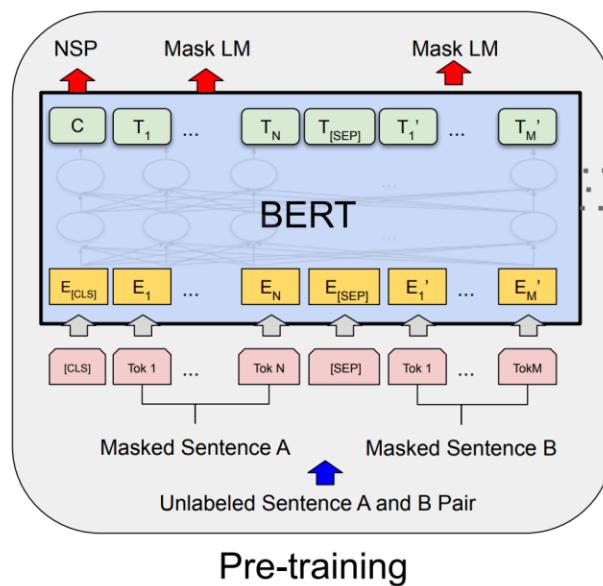
Pre-training and then Fine-Tuning

Pre-training on a proxy task

- Masked token prediction
- Next sentence prediction

Fine-tuning on specific downstream tasks

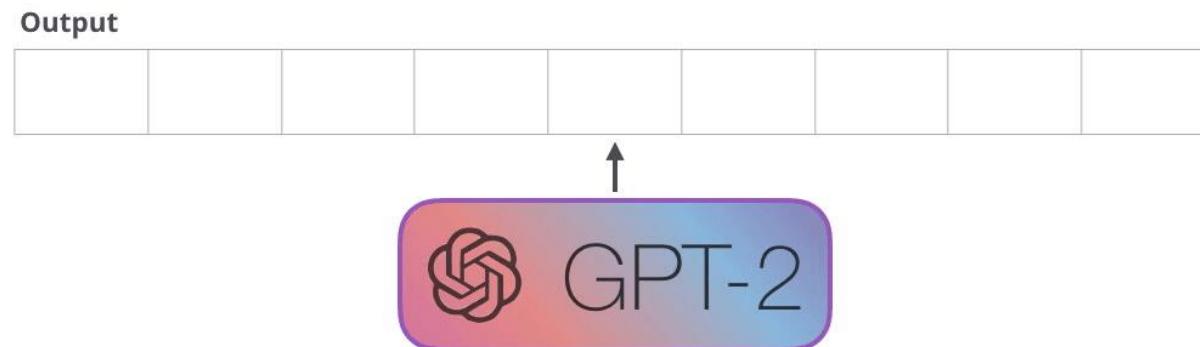
- Machine translation
- Question answering



Decoder-Only - GPT

- Generative Pre-training (GPT)
 - Decoder-only
- Trained with next token prediction
 - A language model!

$$L_1(\mathcal{U}) = \sum_i \log P(u_i | u_{i-k}, \dots, u_{i-1}; \Theta)$$



Large Language Model

- GPT-2
 - Pre-training and fine-tuning on specific tasks
- GPT-3
 - zero-shot capability
 - in-context learning
 - ChatGPT!
- GPT-4

Overview

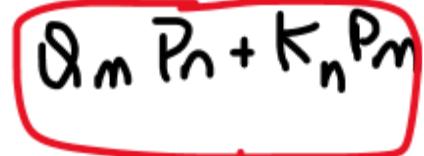
- Architecture
 - Encoder-Decoder
 - Encoder-Only
 - Decoder-Only
- Position Encoding
 - Relative Position Encoding
 - Rotary Position Encoding
- Efficient Attention Mechanism
 - Grouped Query Attention
 - Multi Query Attention
 - Flash Attention
 - Multi-head Latent Attention

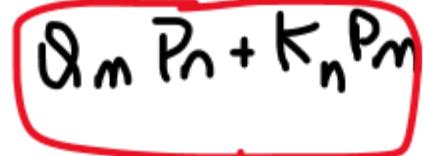
Absolute Position Encoding

- Absolute position embedding fuses the position information into input embeddings

$$(Q_m + P_m) \cdot (K_n + P_n) = Q_m \cdot K_n + P_m P_n$$

helpful signal \leftarrow 

$Q_m \cdot K_n$ + 
 $P_m P_n$

 \downarrow noise

- But, it's for a fixed length! Unable to generalize to longer input sequence and then there are noise terms as well.

Absolute Position Encoding

- Absolute position embedding fuses the position information into input embeddings

$$(Q_m + P_m) \cdot (K_n + P_n) = Q_m \cdot K_n + P_m \cdot P_n + Q_m P_n + K_n P_m$$

helpful signal \leftarrow $Q_m \cdot K_n + P_m \cdot P_n$

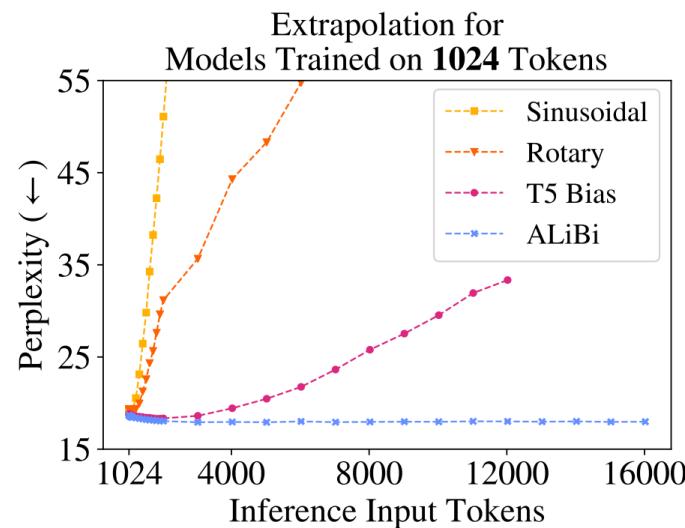
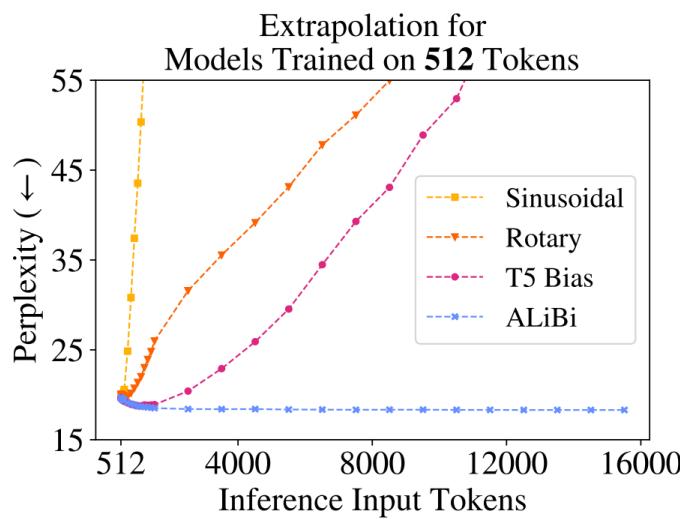
noise \downarrow $Q_m P_n + K_n P_m$

- But, it's for a fixed length! Unable to generalize to longer input sequence and then there's noise.
- We only need the relative position insight in attention, why not add just that?

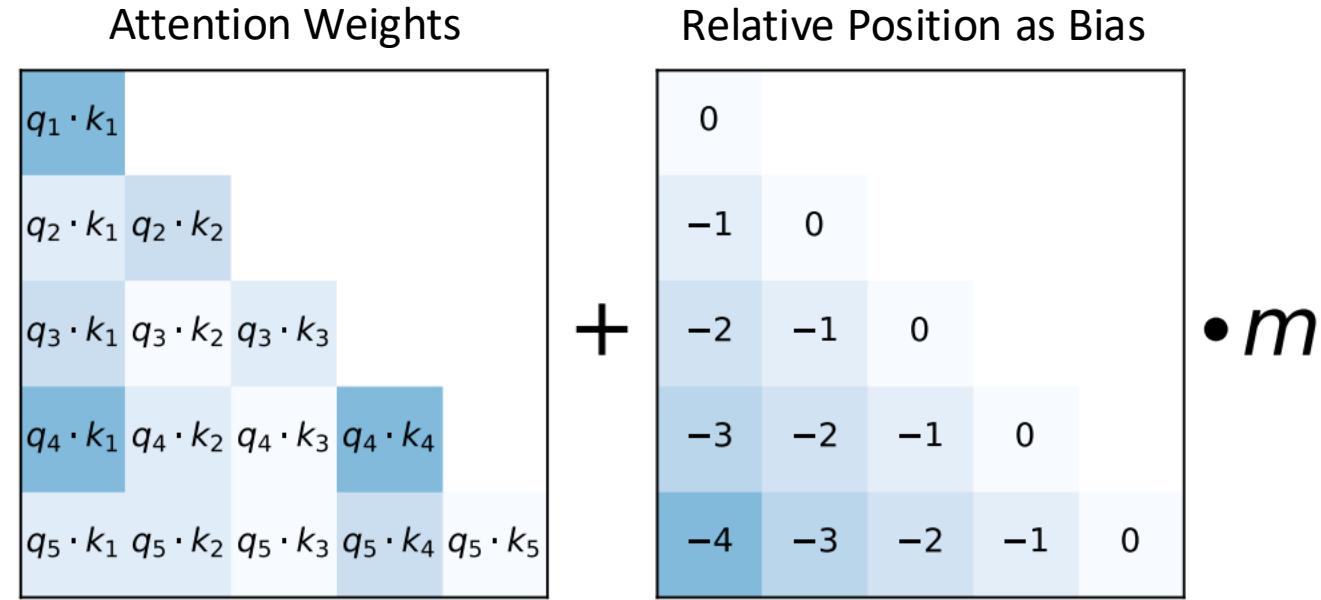
$$\text{Att}_{m,n} = Q_m \cdot K_n + r[m-n] | P_m \cdot P_n$$

Relative Position Encoding

- Relative position embedding fuses position information into attention matrices
- Attention with linear bias
 - Input length extrapolation!



Relative Position Encoding



- Relative distance as offset added to attention matrix
- Absolute position embedding not needed

Rotary Position Encoding

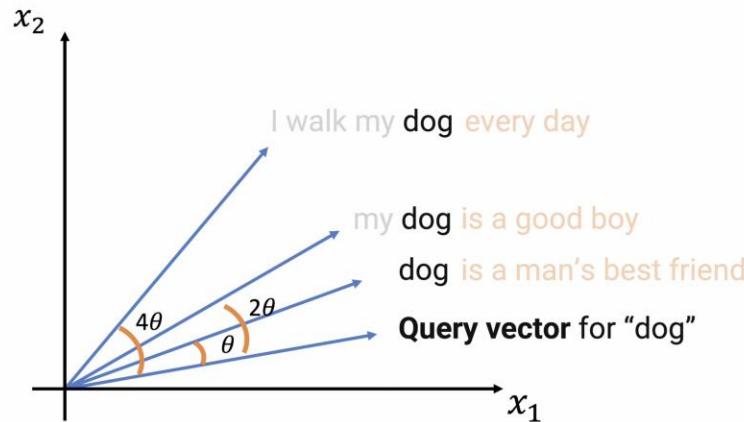
- Used in Large Language Models such as LLAMA
- Rotate the embedding in 2D space

$$\langle f_q(x_m, m), f_k(x_n, n) \rangle = g(x_m, x_n, m - n)$$

$$f_q(x_m, m) = (W_q x_m) e^{im\theta}$$

$$f_k(x_n, n) = (W_k x_n) e^{in\theta}$$

$$g(x_m, x_n, m - n) = \operatorname{Re} \left[(W_q x_m)(W_k x_n)^* e^{i(m-n)\theta} \right]$$



Rotary Position Encoding

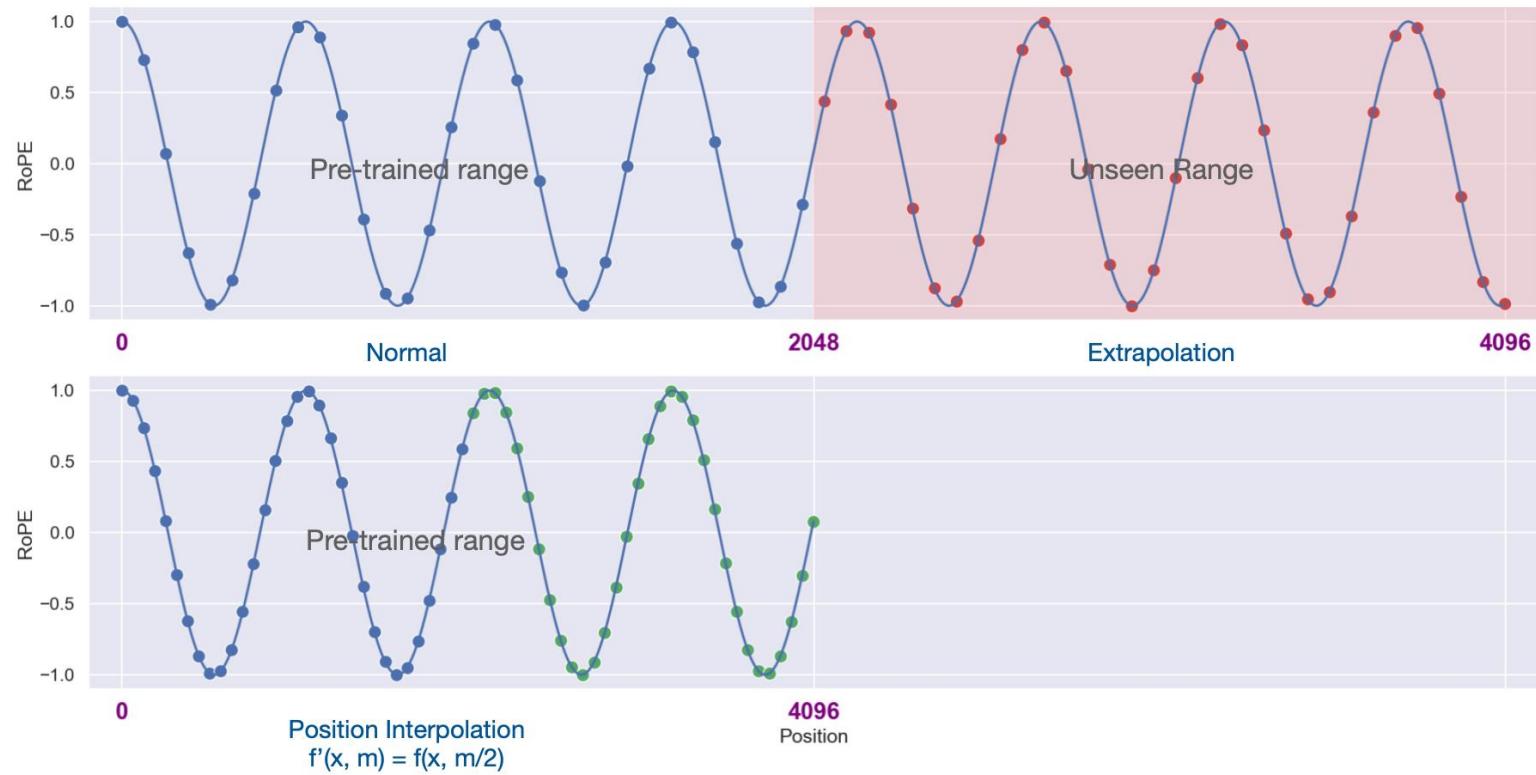
- General form

$$\begin{aligned}f_q(x_m, m) &= (W_q x_m) e^{im\theta} \\f_k(x_n, n) &= (W_k x_n) e^{in\theta} \\g(x_m, x_n, m - n) &= \operatorname{Re} \left[(W_q x_m)(W_k x_n)^* e^{i(m-n)\theta} \right]\end{aligned}$$

$$f_{\{q,k\}}(\mathbf{x}_m, m) = \mathbf{R}_{\Theta, m}^d \mathbf{W}_{\{q,k\}} \mathbf{x}_m$$
$$\mathbf{R}_{\Theta, m}^d = \begin{pmatrix} \cos m\theta_1 & -\sin m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ \sin m\theta_1 & \cos m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cos m\theta_2 & -\sin m\theta_2 & \cdots & 0 & 0 \\ 0 & 0 & \sin m\theta_2 & \cos m\theta_2 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & \cos m\theta_{d/2} & -\sin m\theta_{d/2} \\ 0 & 0 & 0 & 0 & \cdots & \sin m\theta_{d/2} & \cos m\theta_{d/2} \end{pmatrix}$$

Rotary Position Encoding

- Allows extension of the context window



Overview

- Architecture
 - Encoder-Decoder
 - Encoder-Only
 - Decoder-Only
- Position Encoding
 - Relative Position Encoding
 - Rotary Position Encoding
- Efficient Attention Mechanism
 - Linear Attention
 - Flash Attention
 - Grouped Query Attention
 - Multi Query Attention
 - Multi-head Latent Attention

Quadratic Complexity

- Self-attention has quadratic complexity to input length

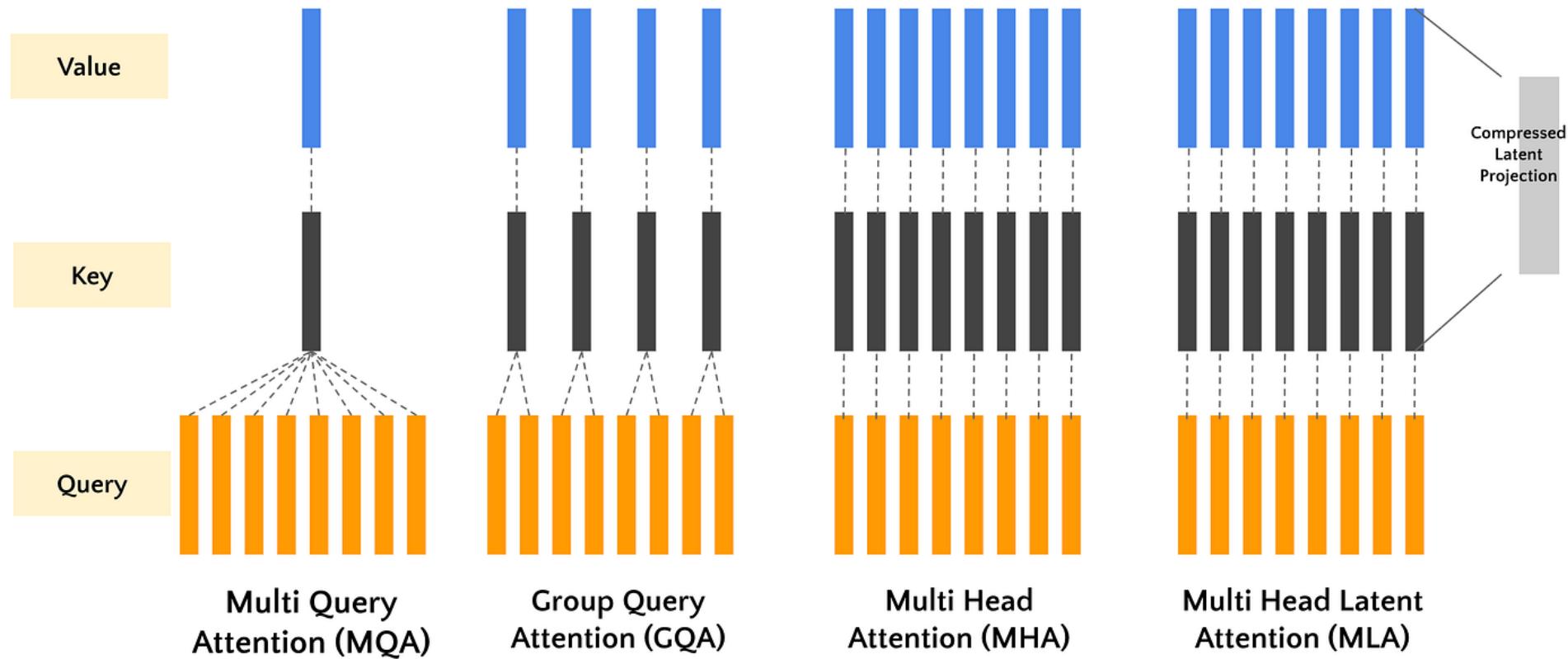
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

- $O(L^2d)$ FLOPS

- Many attempts for reducing the quadratic complexity to linear
 - Linear Attention
 - Flash Attention
 - Grouped Query Attention
 - Multi Query Attention
 - Multi-head Latent Attention

Attention Types

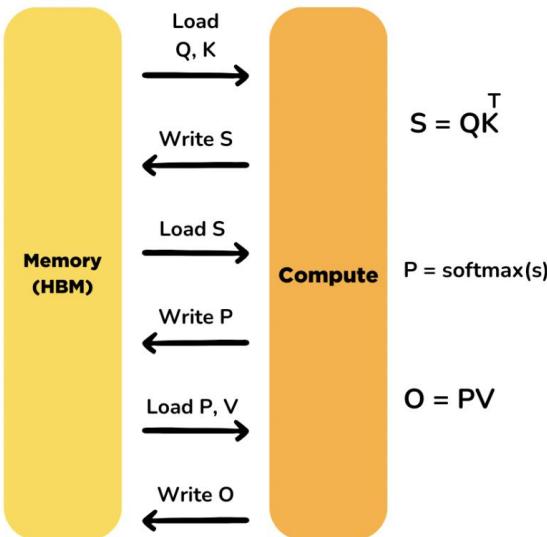
Attention – MQA | GQA | MHA | MLA



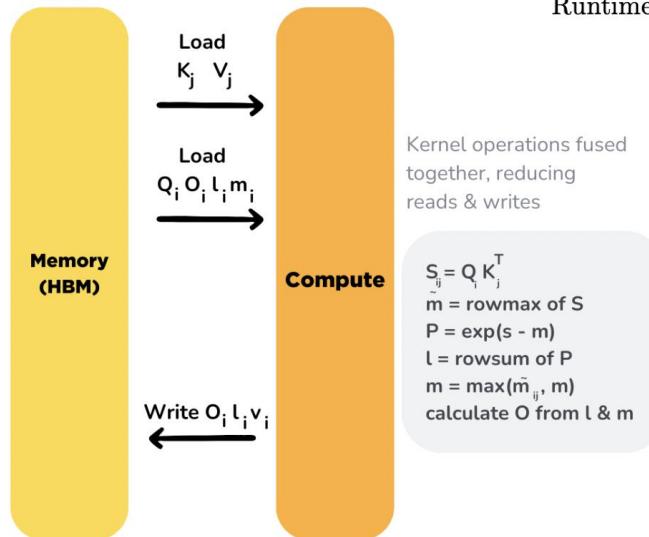
Flash Attention

- Reduce GPU <-> CPU memory trips

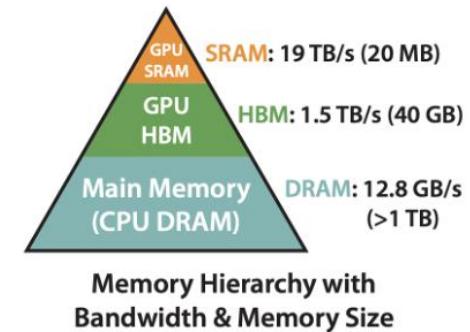
Standard Attention Implementation



Flash Attention



Attention	Standard	FLASHATTENTION
GFLOPs	66.6	75.2
HBM R/W (GB)	40.3	4.4
Runtime (ms)	41.7	7.3



Initialize O , L and m matrices with zeroes. m and L are used to calculate cumulative softmax. Divide Q , K , V into blocks (due to SRAM's memory limits) and iterate over them, for i is row & j is column.

Flash Attention

- Classic GPU trick – tiled multiplication

With Tiling 16 access

$$\begin{array}{c}
 \text{A} \\
 \begin{array}{|cc|cc|} \hline & & & \\ 1 & 2 & 0 & 5 \\ & & & \\ \hline & & & \\ 4 & 8 & 2 & -1 \\ & & & \\ \hline & & & \\ 3 & 1 & 6 & 3 \\ & & & \\ \hline & & & \\ -7 & 5 & 0 & 8 \\ & & & \\ \hline \end{array}
 \end{array}
 \times
 \begin{array}{c}
 \text{B} \\
 \begin{array}{|cc|cc|} \hline & & & \\ 3 & 1 & 6 & 3 \\ -7 & 5 & 0 & 8 \\ \hline & & & \\ 1 & 2 & 0 & 5 \\ & & & \\ \hline & & & \\ 0 & 3 & 5 & 1 \\ & & & \\ \hline \end{array}
 \end{array}
 =
 \begin{array}{c}
 \text{C} \\
 \begin{array}{|c|c|} \hline & \\ \text{C1,1} & \text{C1,2} \\ & \\ \hline & \\ \text{C2,1} & \text{C2,2} \\ & \\ \hline \end{array}
 \end{array}$$

$C_{11} = A_{11} \times B_{11} + A_{12} \times B_{21}$
 $C_{12} = A_{11} \times B_{12} + A_{12} \times B_{22}$
 $C_{21} = A_{21} \times B_{11} + A_{22} \times B_{21}$
 $C_{22} = A_{21} \times B_{12} + A_{22} \times B_{22}$

Dao, Tri, et al. "Flashattention: Fast and memory-efficient exact attention with io-awareness." Advances in Neural Information Processing Systems 35 (2022): 16344-16359.

Flash Attention

- online softmax - revitalizes safe softmax

Softmax

Computation requires two loops: one to calculate the normalizing factor (the sum of exponentials) and another to compute the attention weights by dividing each exponentiated value by this factor.

$$\frac{e^{x_i}}{\sum_{j=1}^N e^{x_j}}$$

Safe Softmax

Requires three loops: one to find the maximum value (for numerical stability), one to compute the normalizing factor, and one to obtain the attention weights.

```
for i ← 1, N do
     $m_i \leftarrow \max(m_{i-1}, x_i)$ 
for i ← 1, N do
     $d_i \leftarrow d_{i-1} + e^{x_i - m_N}$ 
for i ← 1, N do
     $a_i \leftarrow \frac{e^{x_i - m_N}}{d_N}$ 
```

Online Softmax

Requires two loops: one to find the maximum value (and to compute the normalizing factor, and one to obtain the attention weights.

```
for i ← 1, N do
     $m_i \leftarrow \max(m_{i-1}, x_i)$ 
     $d'_i \leftarrow d'_{i-1} e^{m_{i-1} - m_i} + e^{x_i - m_i}$ 
for i ← 1, N do
     $a_i \leftarrow \frac{e^{x_i - m_N}}{d'_N}$ 
```

Flash Attention

Models	ListOps	Text	Retrieval	Image	Pathfinder	Avg	Speedup
Transformer	36.0	63.6	81.6	42.3	72.7	59.3	-
FLASHATTENTION	37.6	63.9	81.4	43.5	72.7	59.8	2.4×
Block-sparse FLASHATTENTION	37.0	63.0	81.3	43.6	73.3	59.6	2.8×
Linformer [84]	35.6	55.9	77.7	37.8	67.6	54.9	2.5×
Linear Attention [50]	38.8	63.2	80.7	42.6	72.5	59.6	2.3×
Performer [12]	36.8	63.6	82.2	42.1	69.9	58.9	1.8×
Local Attention [80]	36.1	60.2	76.7	40.6	66.6	56.0	1.7×
Reformer [51]	36.5	63.8	78.5	39.6	69.4	57.6	1.3×
Smyrf [19]	36.1	64.1	79.0	39.6	70.5	57.9	1.7×

IO complexity:

$$\text{Flash Attention: } \mathbf{O}\left(\frac{N^2 d^2}{M}\right)$$

$$\text{Standard Attention: } \Omega(Nd + N^2)$$

Where **N** is sequence length, **d** head dimensions and **M** the size of SRAM.

KV- Caching

Step 1

Without Cache

$$\begin{array}{ccc} Q & K^T & QK^T \\ \text{Query Token 1} & \text{Key Token 1} & Q_i K_i \\ x & & = \\ (1, \text{emb_size}) & (\text{emb_size}, 1) & (1, 1) \end{array} \quad \left| \quad \begin{array}{ccc} V & \text{Attention} \\ \text{Value Token 1} & \text{Token 1} \\ x & = \\ (1, \text{emb_size}) & (1, \text{emb_size}) \end{array} \right.$$

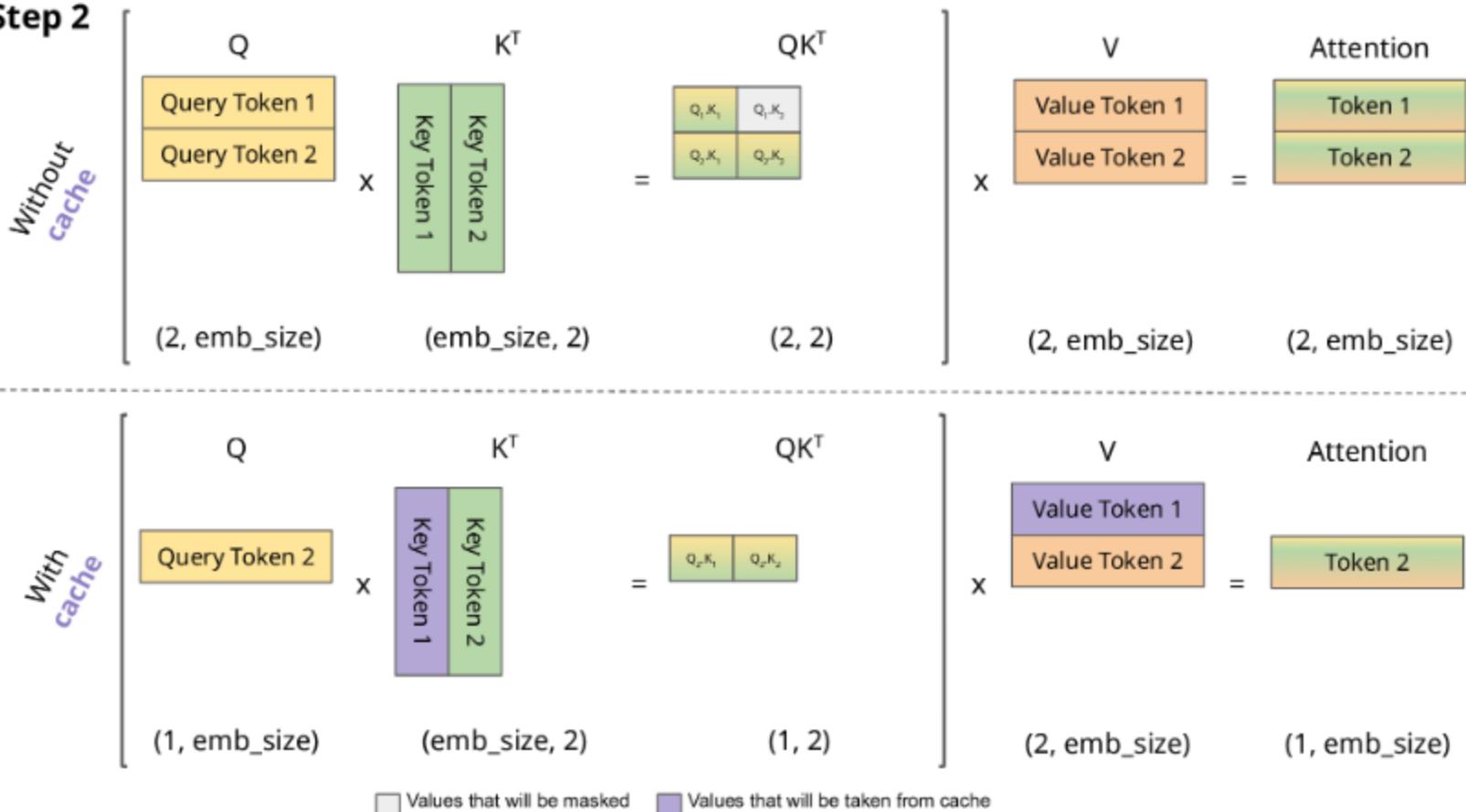
With cache

$$\begin{array}{ccc} Q & K^T & QK^T \\ \text{Query Token 1} & \text{Key Token 1} & Q_i K_i \\ x & & = \\ (1, \text{emb_size}) & (\text{emb_size}, 1) & (1, 1) \end{array} \quad \left| \quad \begin{array}{ccc} V & \text{Attention} \\ \text{Value Token 1} & \text{Token 1} \\ x & = \\ (1, \text{emb_size}) & (1, \text{emb_size}) \end{array} \right.$$

 Values that will be masked  Values that will be taken from cache

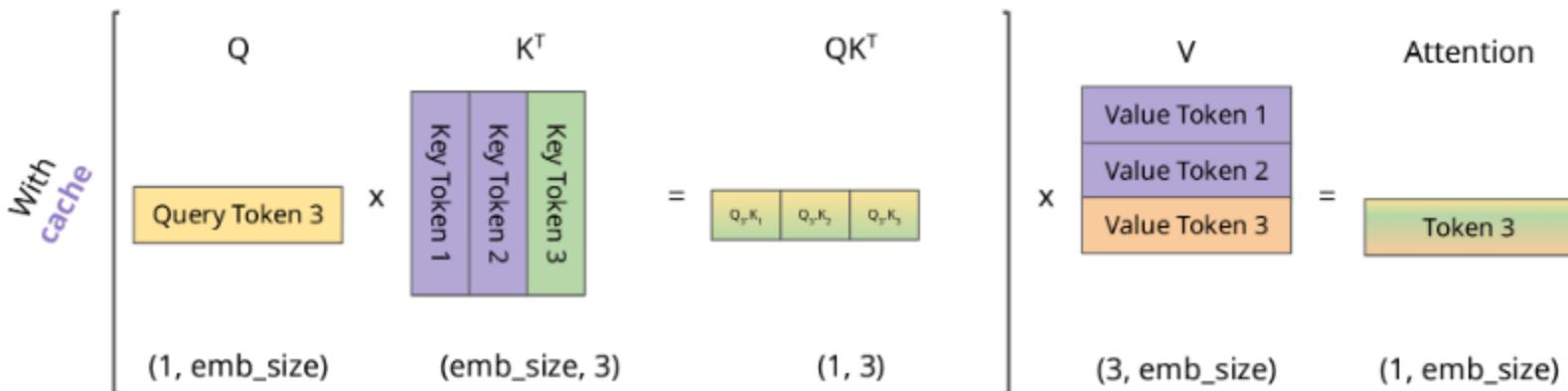
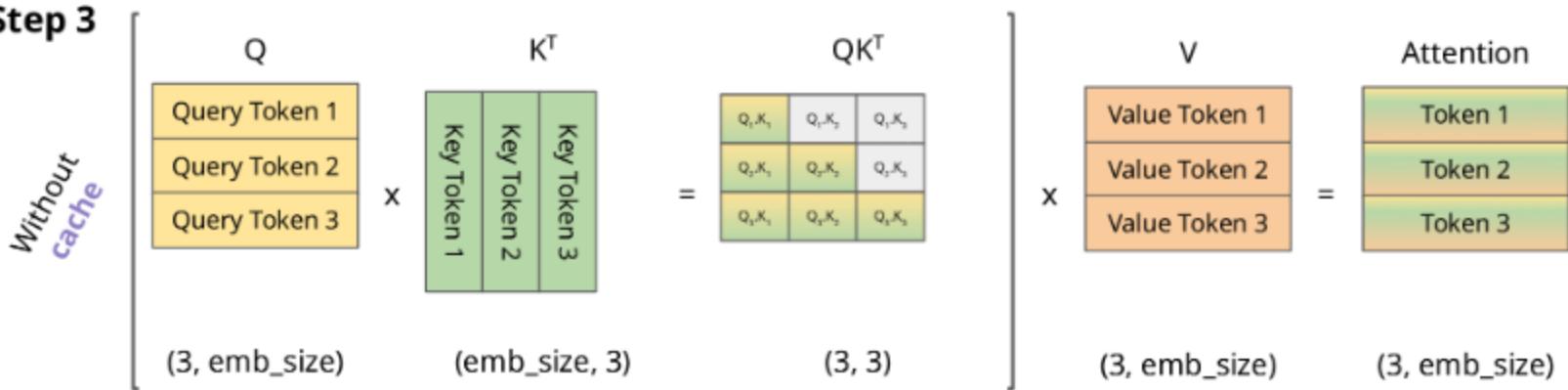
KV- Caching

Step 2



KV- Caching

Step 3

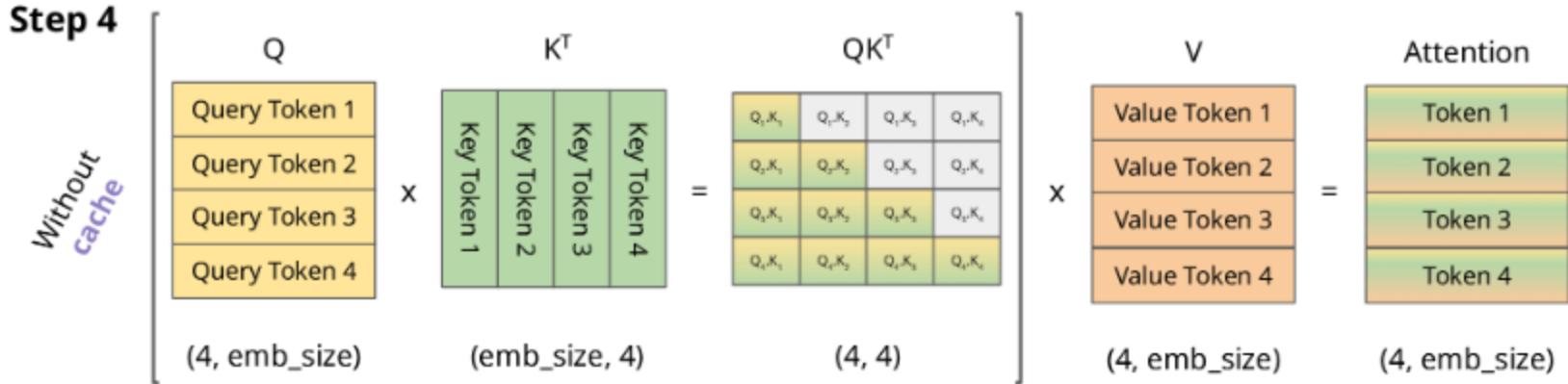


 Values that will be masked

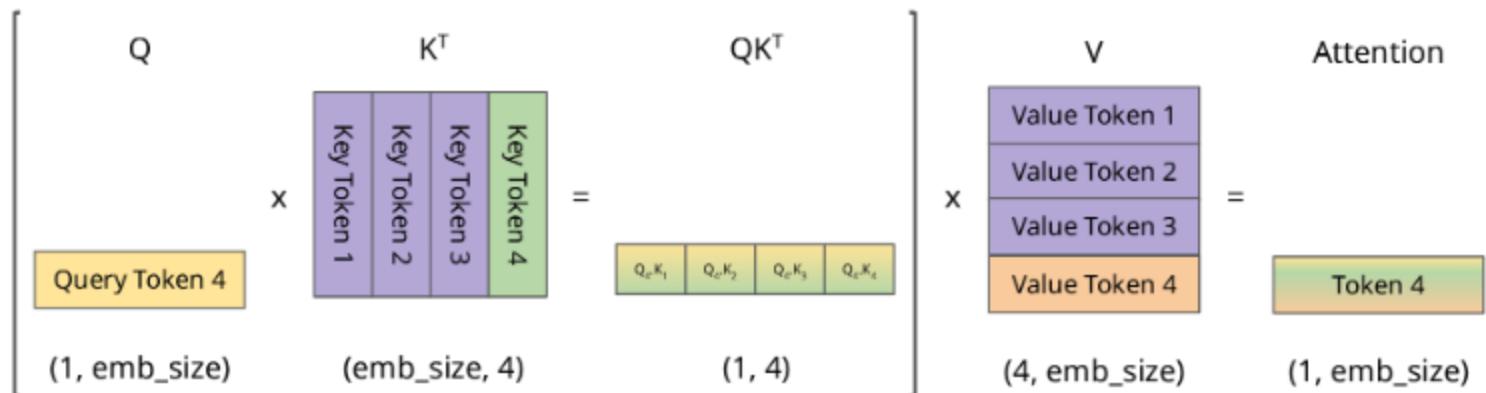
 Values that will be taken from cache

KV- Caching

Step 4



With
cache



Values that will be masked Values that will be taken from cache

Poll @1209

Which of the following statements is true?

- FlashAttention is particularly effective for long sequences, as it stores the full attention matrix in memory, which would otherwise grow quadratically with sequence length due to a higher number of memory accesses.
- FlashAttention improves efficiency by splitting computations into blocks that fit in fast SRAM, reducing memory access overhead while maintaining mathematical equivalence to standard attention.
- FlashAttention performs worse than standard attention implementations because the block-wise computation approach introduces additional computational overhead that outweighs any memory benefits.
- FlashAttention is primarily designed for CPU optimization and shows minimal performance improvements when implemented on GPU hardware.

Poll @1209

Which of the following statements is true?

- FlashAttention is particularly effective for long sequences, as it stores the full attention matrix in memory, which would otherwise grow quadratically with sequence length due to a higher number of memory accesses.
- FlashAttention improves efficiency by splitting computations into blocks that fit in fast SRAM, reducing memory access overhead while maintaining mathematical equivalence to standard attention.
- FlashAttention performs worse than standard attention implementations because the block-wise computation approach introduces additional computational overhead that outweighs any memory benefits.
- FlashAttention is primarily designed for CPU optimization and shows minimal performance improvements when implemented on GPU hardware.

Content

- Transformer Architecture
- Improvements on Transformers
- **Transformer for different modalities**
- Parameter Efficient Tuning
- Scaling Laws
- Quantizing Transformers
- Interpreting Transformer – attention, logitlens

Transformer in Vision and Audio

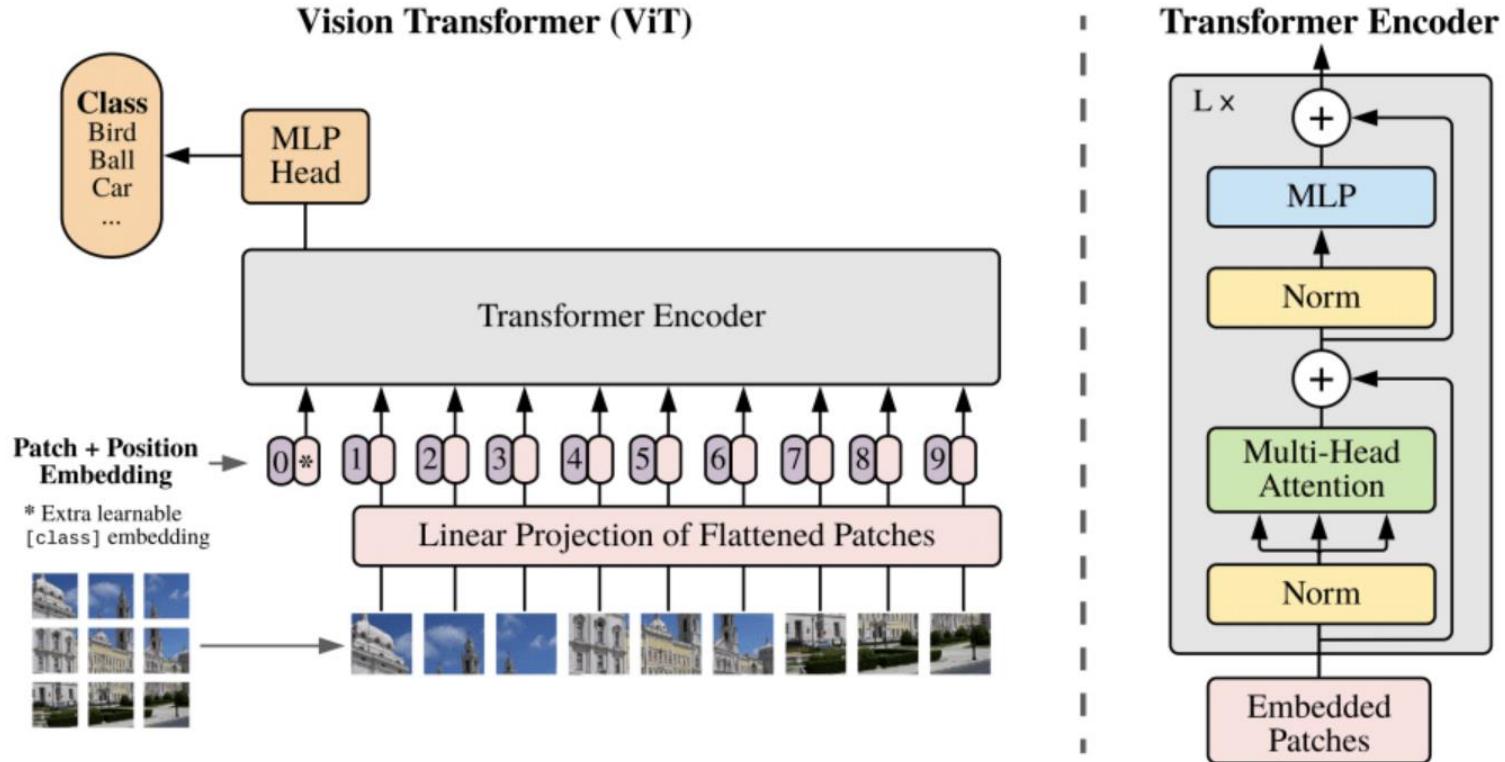
Overview

- Vision Transformer Architecture
- Transformer in Audio
- Tokenizer

Overview

- Vision Transformer Architecture
- Transformer in Audio
- Tokenizers

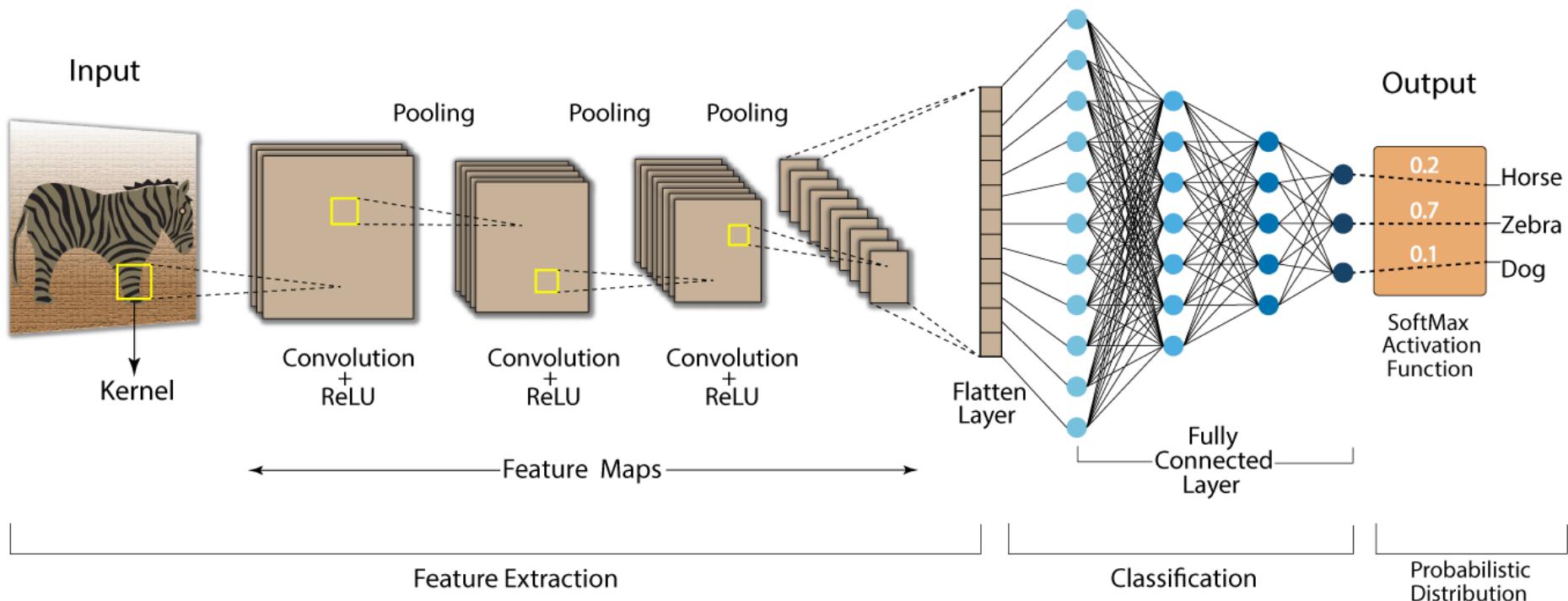
Vision Transformer (ViT)



- Transformer architecture can also be used for images
- How do we process an image into tokens?

CNN

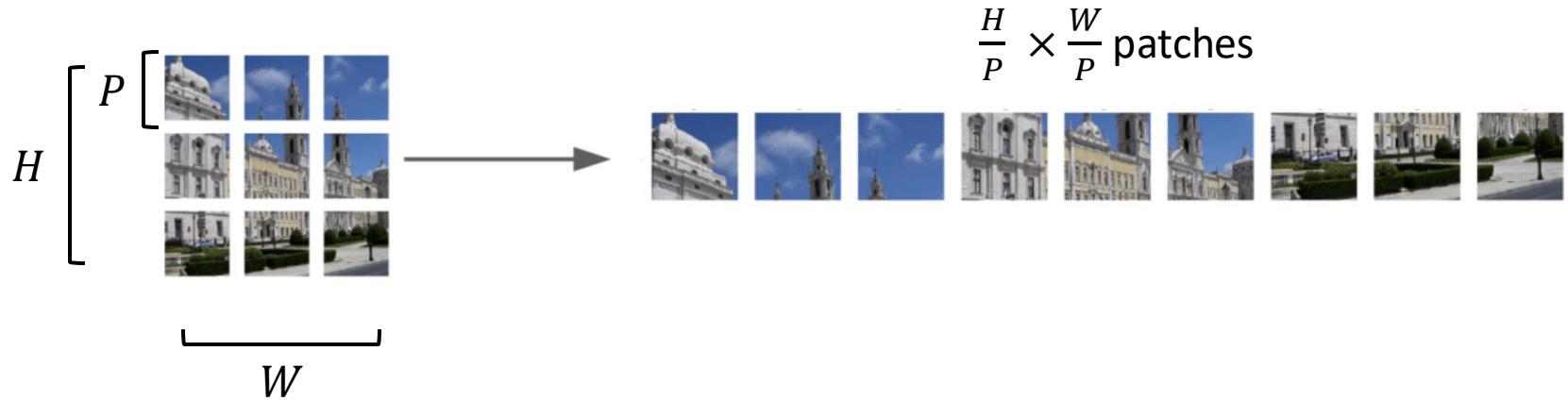
Convolution Neural Network (CNN)



- Naturally fits to 2D images

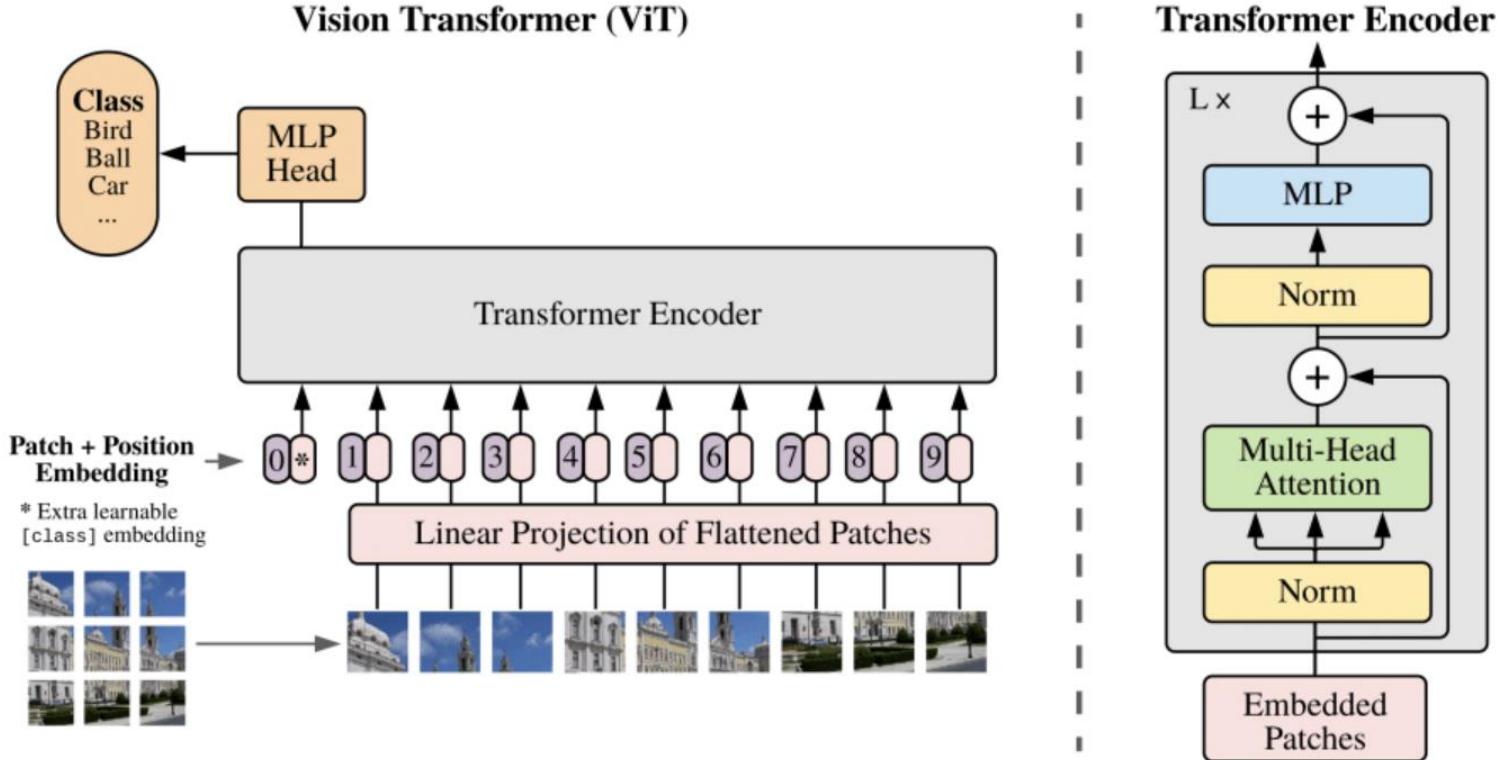
ViT

- Split images into a sequence of **patches**



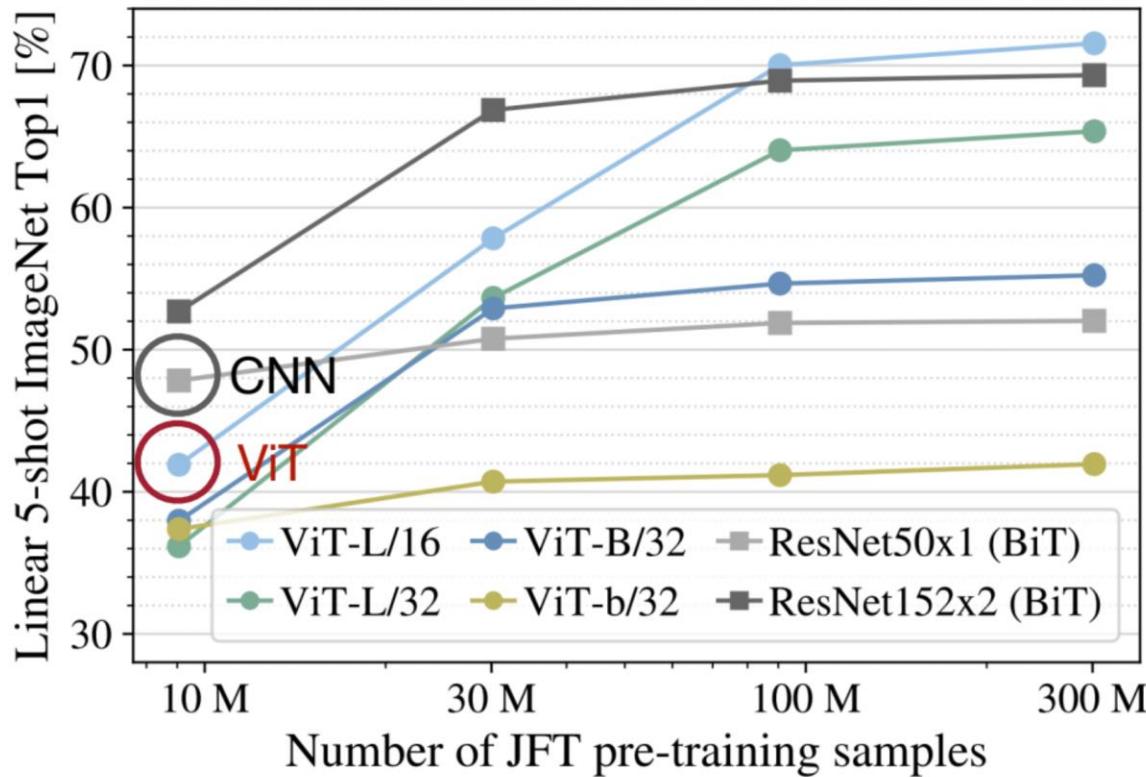
- Each patch is treated as one token as input to ViT
 - A convolution layer with kernel P and stride P !
 - Or a linear layer on the flatten pixels

ViT



- The remaining is same as Transformer
 - As an encoder-only model

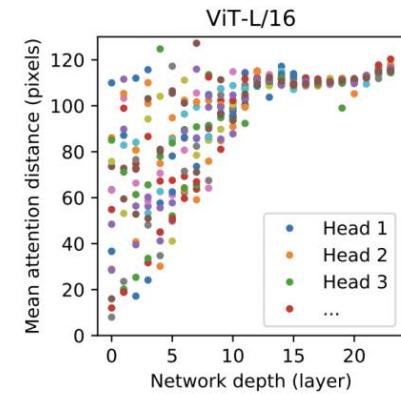
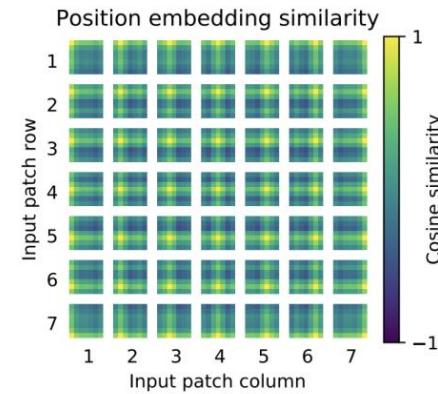
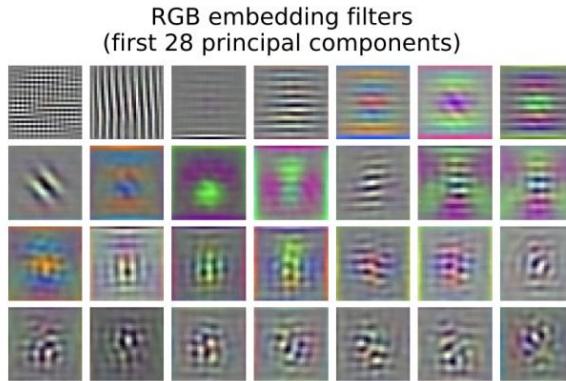
Image Classification



- Inferior performance compared to CNN when dataset size is limited – Why?

Inductive Bias

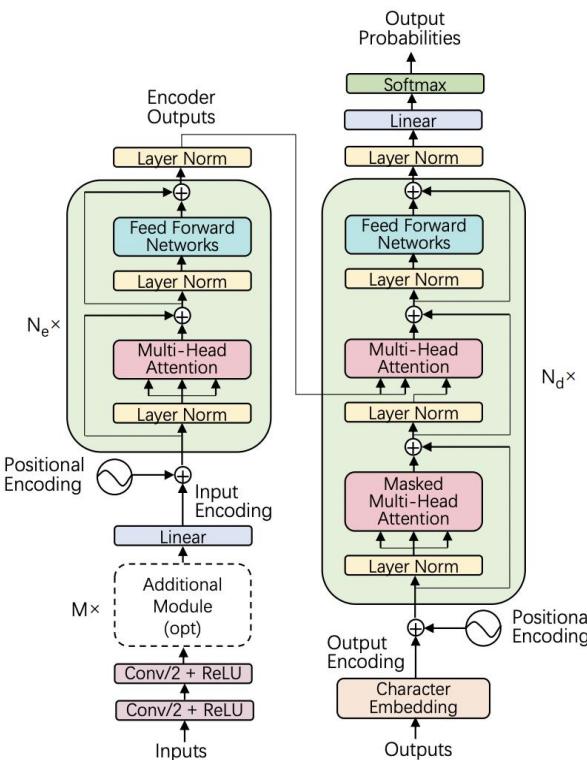
- Convolutional Neural Networks
 - Locality
 - Sharing weights
- Vision Transformer
 - None!
 - Has to learn locality and dependency from data!
 - A lot lot lot lot lot lot lot lot of data!



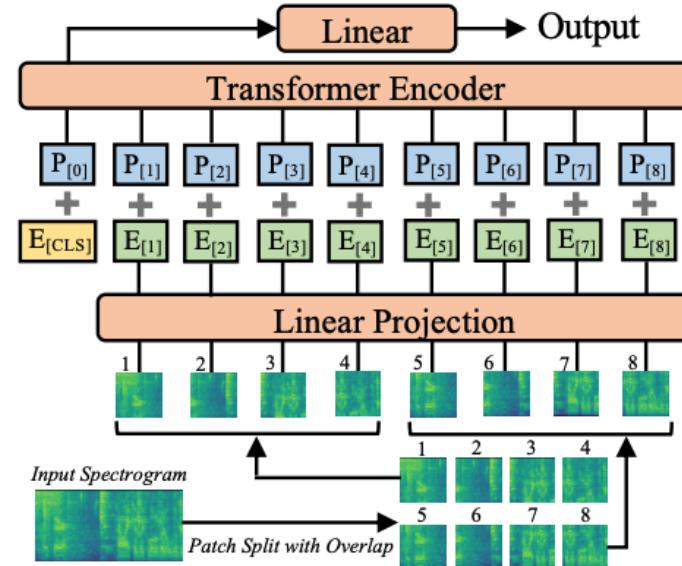
Overview

- Vision Transformer Architecture
- Transformer in Audio
- Tokenizer

Transformer in Audio



Speech Transformer for ASR



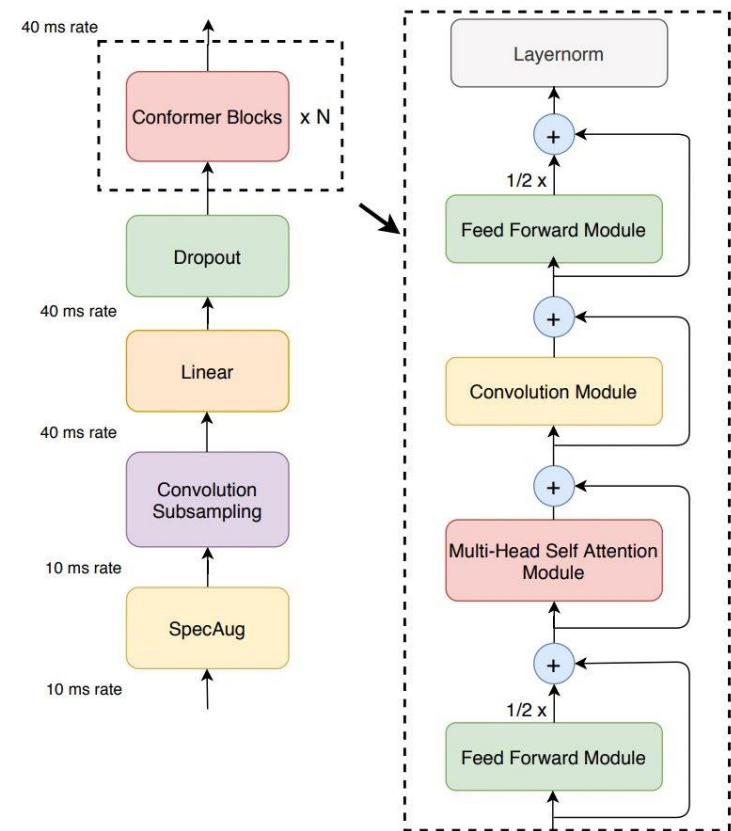
Audio Spectrogram Transformer

[1] Dong, Linhao, Shuang Xu, and Bo Xu. "Speech-transformer: a no-recurrence sequence-to-sequence model for speech recognition." *2018 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2018.

[2] Gong, Yuan, Yu-An Chung, and James Glass. "Ast: Audio spectrogram transformer." *arXiv preprint arXiv:2104.01778* (2021).

Conformer

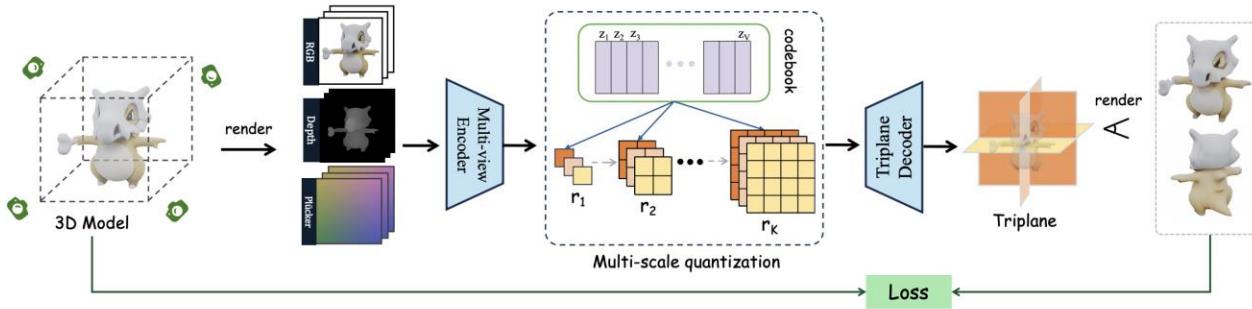
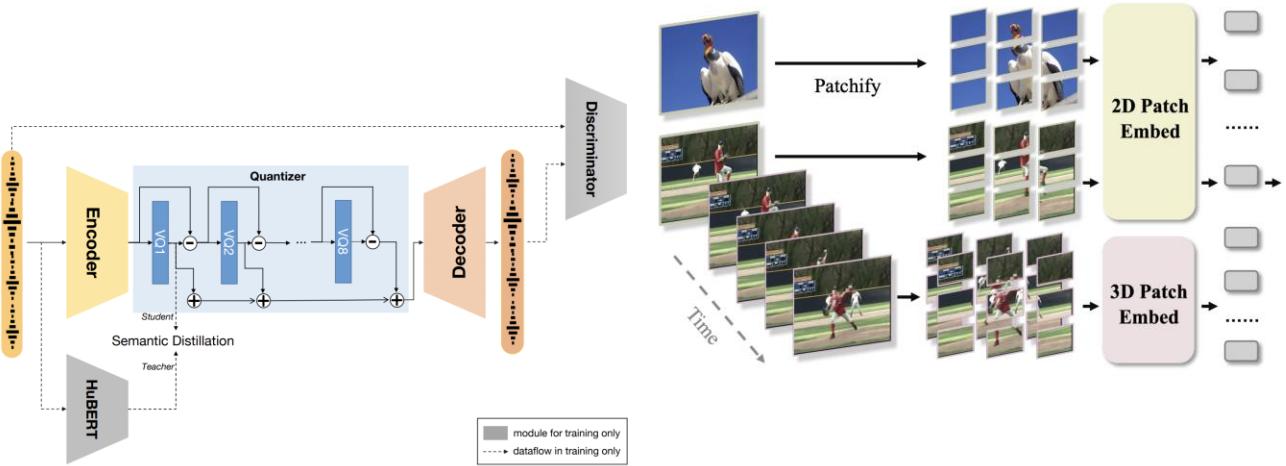
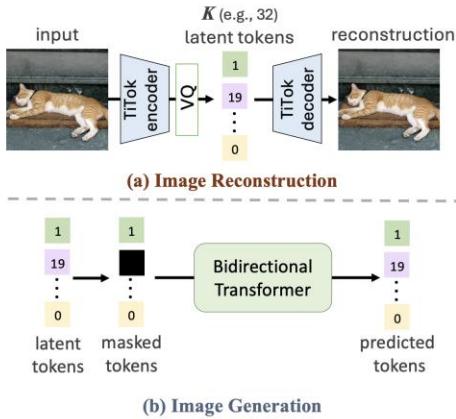
- The Conformer architecture augments a transformer by embedding convolution layers within the transformer blocks.
- Transformers capture global dependencies, CNNs capture local features efficiently.



Overview

- Vision Transformer Architecture
- Transformer in Audio: Conformer
- Tokenizer

Tokenizers



Zhang, Xin, et al. "Speechtokenizer: Unified speech tokenizer for speech language models." The Twelfth International Conference on Learning Representations. 2024.

Chen, Yongwei, et al. "SAR3D: Autoregressive 3D object generation and understanding via multi-scale 3D VQVAE." arXiv preprint arXiv:2411.16856 (2024).

Yu, Qihang, et al. "An Image is Worth 32 Tokens for Reconstruction and Generation." arXiv preprint arXiv:2406.07550 (2024).

Wang, Junke, et al. "OmniTokenizer: A Joint Image-Video Tokenizer for Visual Generation." arXiv preprint arXiv:2406.09399 (2024).

Poll @ 1210,1211

Which ones of the following are properties of ViT, compared to CNN?

- Weight sharing
- Dynamic weights from data
- Locality
- Global dependency from data

Which of the following statements about the Conformer architecture is correct?

- The Conformer uses convolution layers to replace self-attention entirely
- Conformer blocks have convolutional modules placed after the self-attention module
- The Conformer architecture eliminates the need for Feed Forward modules
- Conformer was primarily designed for computer vision tasks rather than speech recognition

Poll @ 1210,1211

Which ones of the following are properties of ViT, compared to CNN?

- Weight sharing
- Dynamic weights from data
- Locality
- Global dependency from data

Which of the following statements about the Conformer architecture is correct?

- The Conformer uses convolution layers to replace self-attention entirely
- Conformer blocks have convolutional modules placed after the self-attention module
- The Conformer architecture eliminates the need for Feed Forward modules
- Conformer was primarily designed for computer vision tasks rather than speech recognition

Content

- Transformer Architecture
- Improvements on Transformers
- Transformer for different modalities
- **Parameter Efficient Tuning**
- Scaling Laws
- Quantizing Transformers
- Interpreting Transformer – attention, logitlens

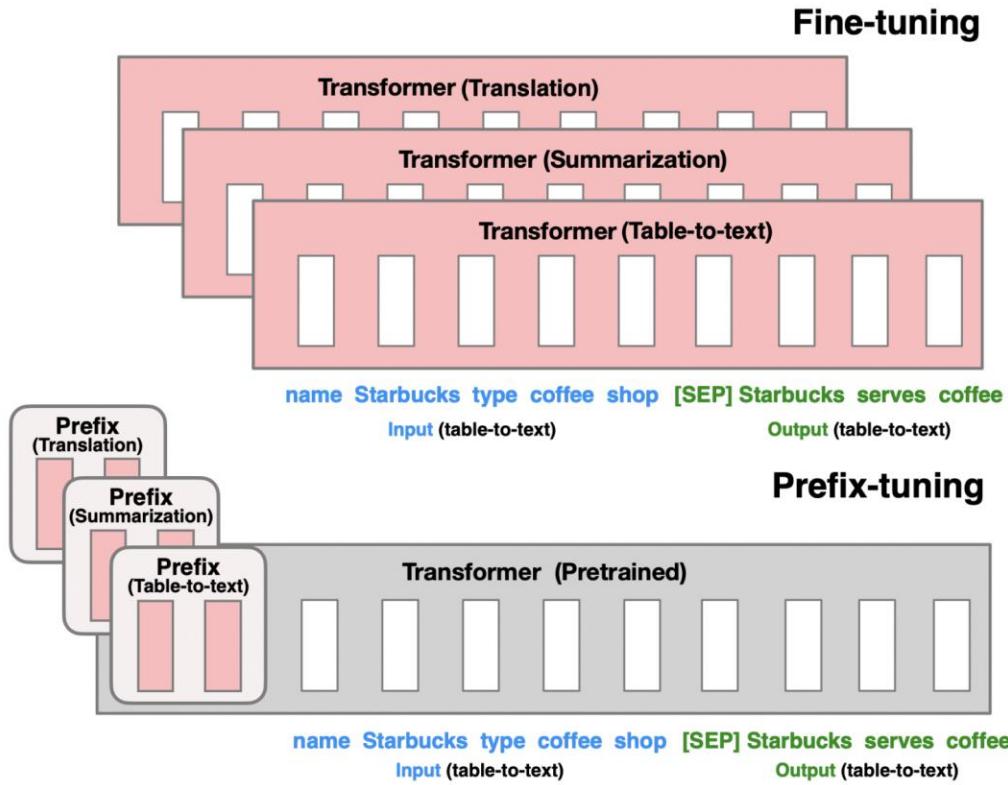
Parameter Efficient Tuning

Parameter Efficient Tuning

- Traditionally, you need to fine-tune entire network on specific downstream tasks
- Parameter Efficient Tuning – Only tune a small proportion of parameters of the pre-trained transformer
 - Prefix Tuning
 - Prompt tuning
 - Adapter
 - LoRA

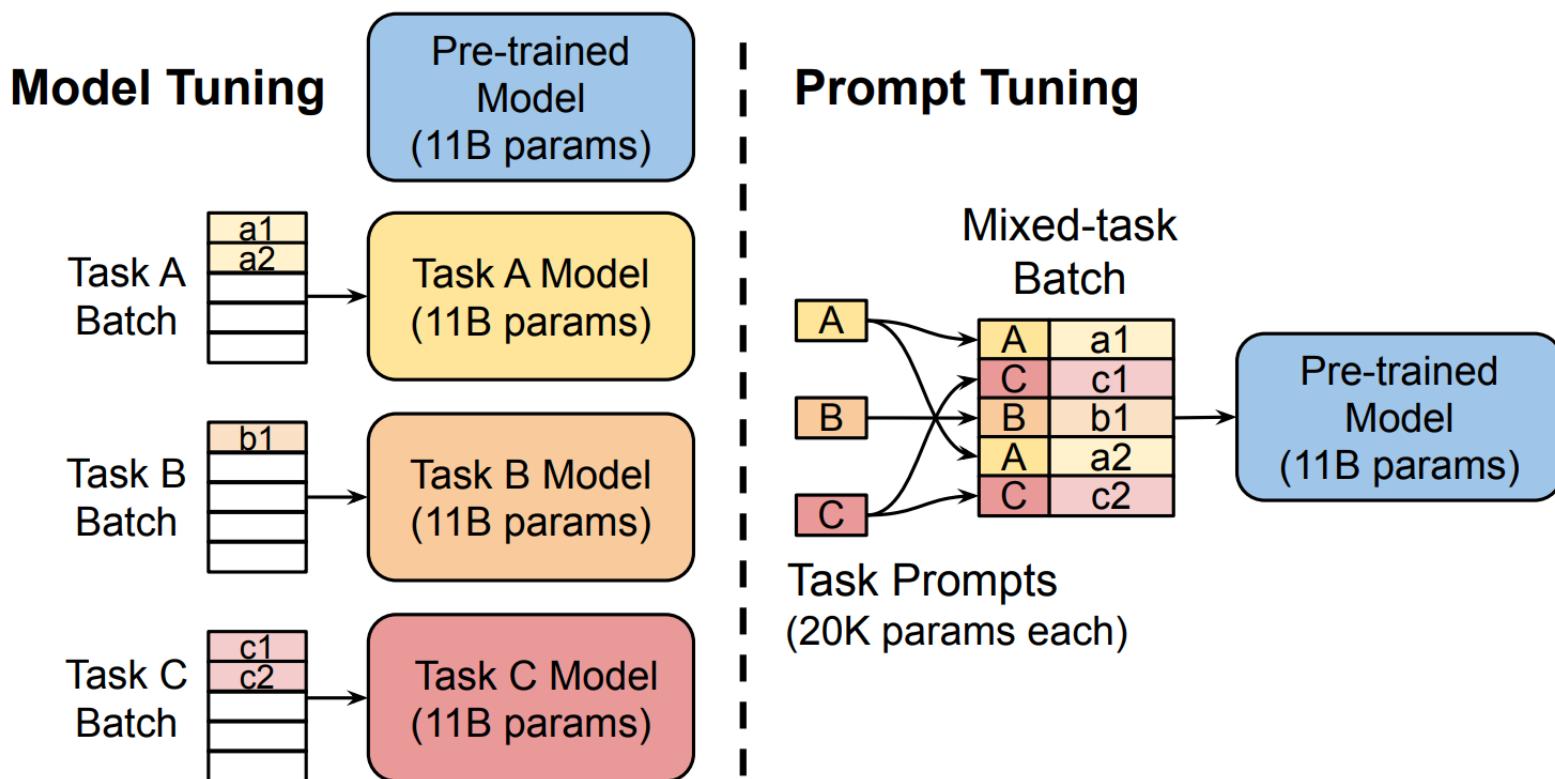
Prefix Tuning

- Only learns a set continuous prefixes)added to the input and transformer layers for each task.



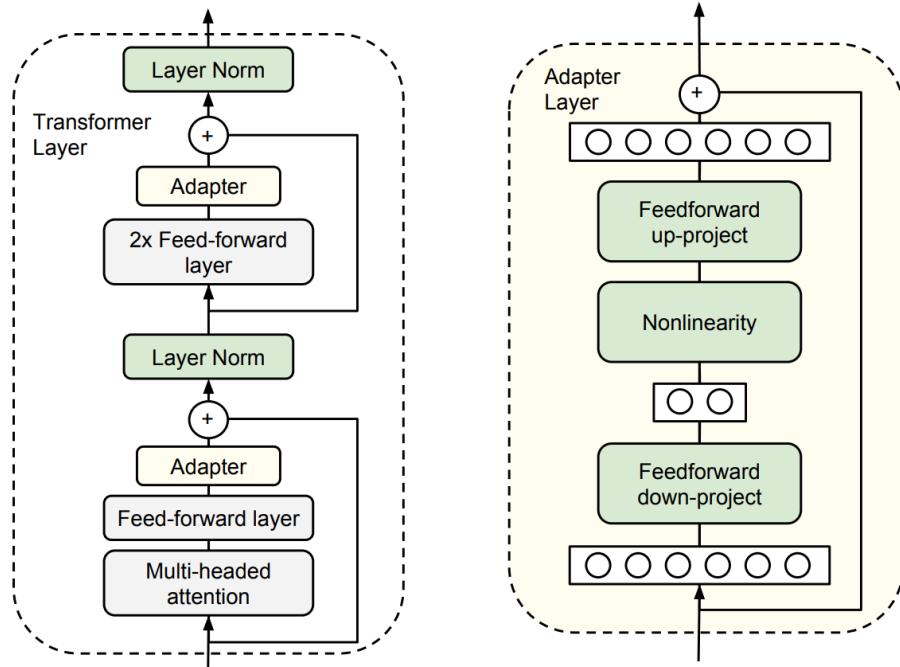
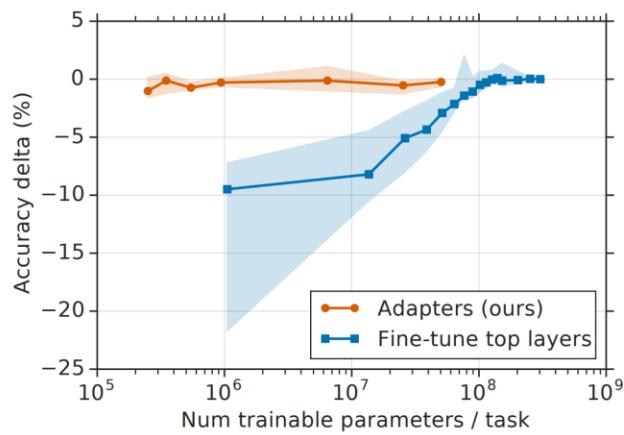
Prompt Tuning

- Only learns a set of ‘prompt’ or ‘token’ for each task



Adapter

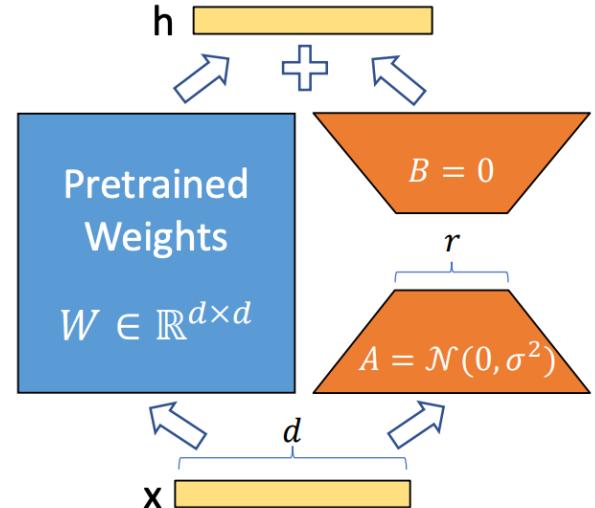
- Insert MLP at Feed-forward layers



LoRA

- Low-rank Adaptation (LoRA)
- No activation in-between
- A and B can be fused into W

$$h = W_0x + \Delta Wx = W_0x + BAx$$

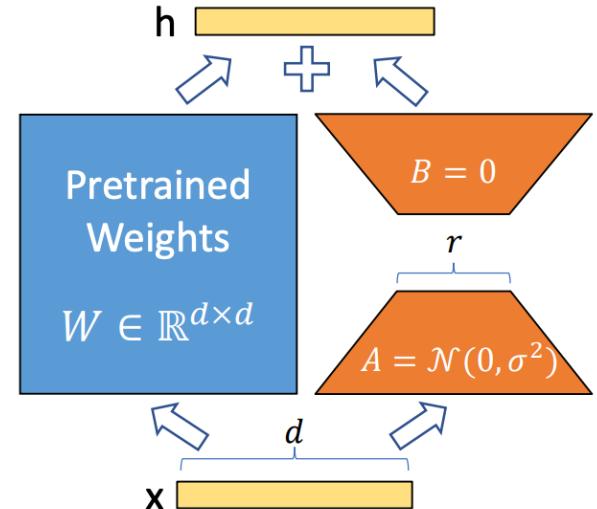


- By default, only W_Q and W_V have the AB pairs
 - Just changing which part to focus on works!

LoRA

- Low-rank Adaptation (LoRA)
- No activation in-between
- A and B can be fused into W

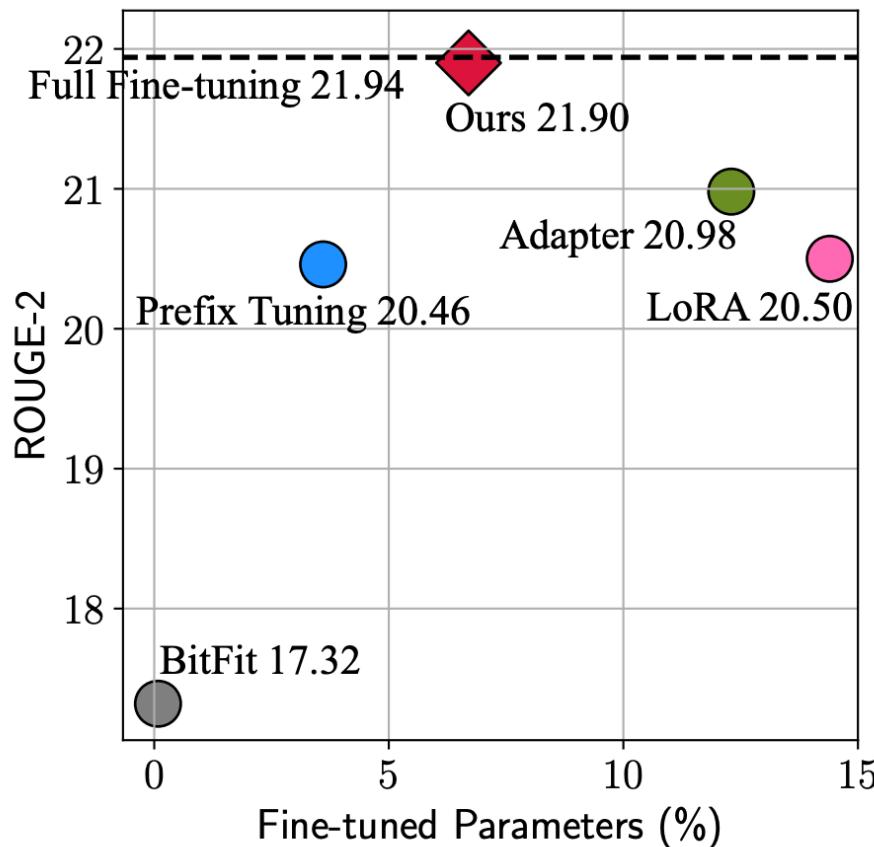
$$h = W_0x + \Delta Wx = W_0x + BAx$$



- By default, only W_Q and W_V have the AB pairs
 - Just changing which part to focus on works!
- For MLP when we want to add more knowledge

Parameter-Efficient Tuning

- Performance close to full fine-tuning while just train less than 15% of original parameters

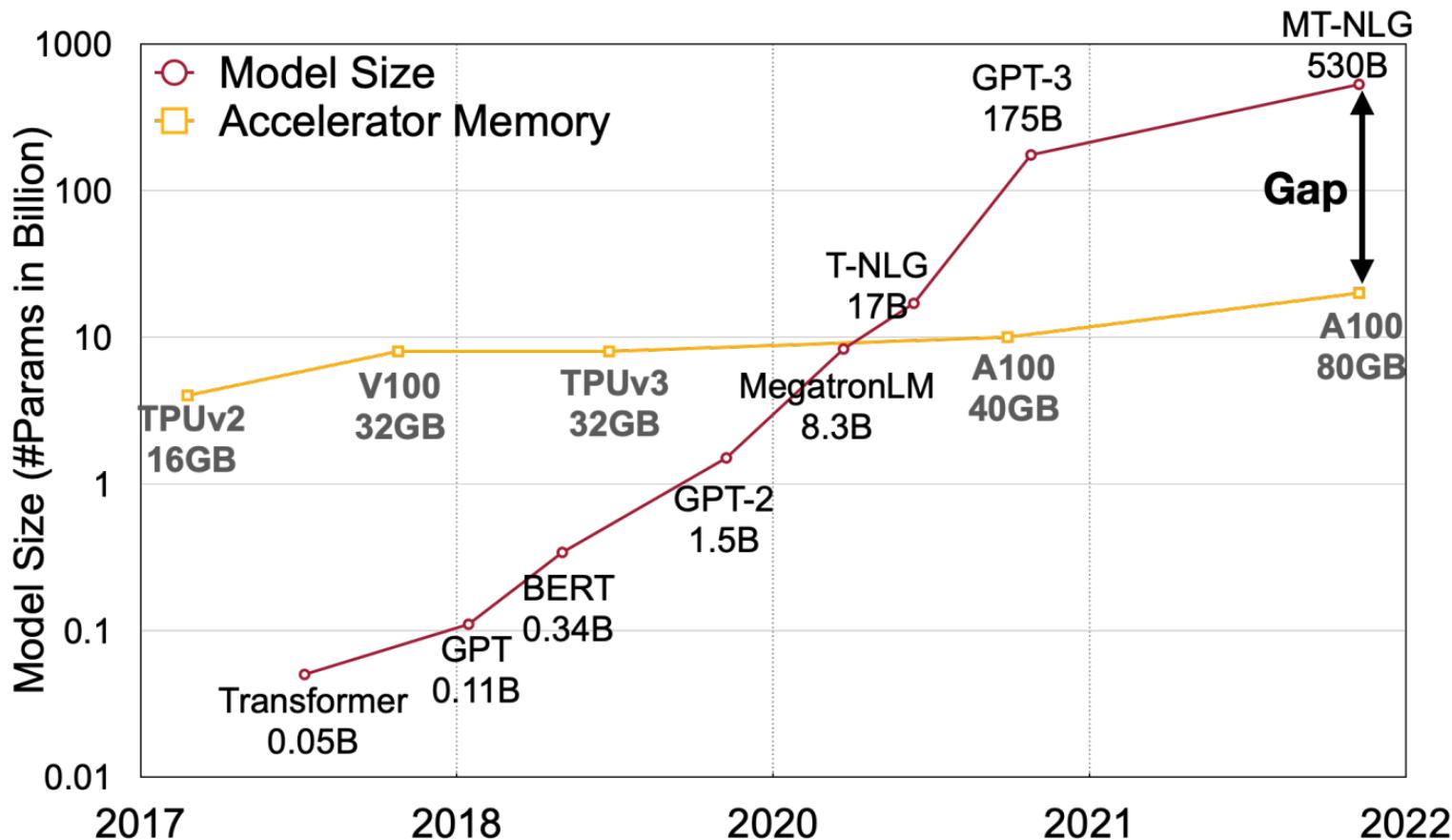


Content

- Transformer Architecture
- Improvements on Transformers
- Transformer for different modalities
- Parameter Efficient Tuning
- **Scaling Laws**
- Quantizing Transformers
- Interpreting Transformer – attention, logitlens

Scaling Laws

“Magic” of Transformer - Scaling

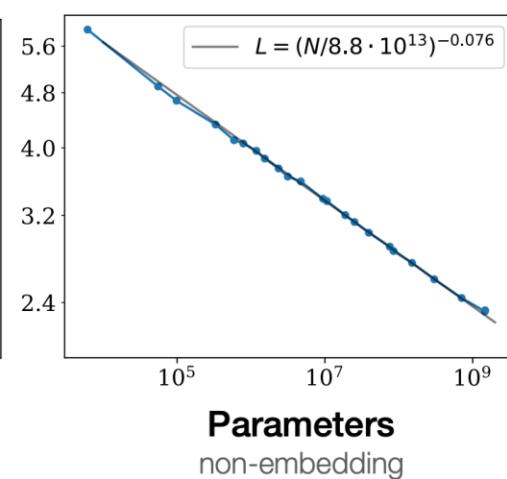
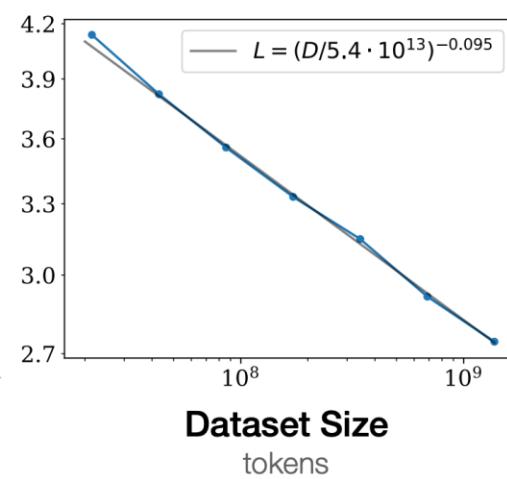
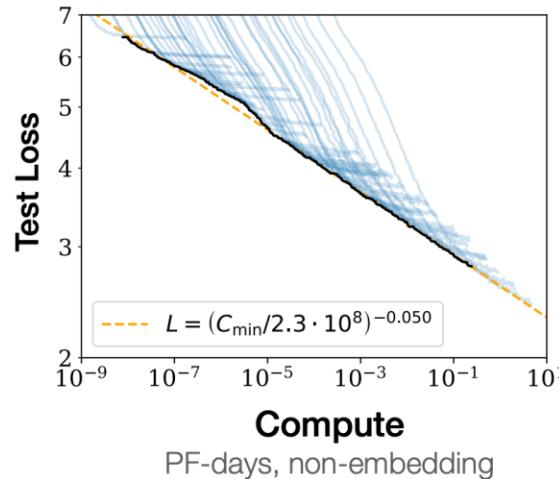


- Performance gets better as transformer scales up

Scaling Law

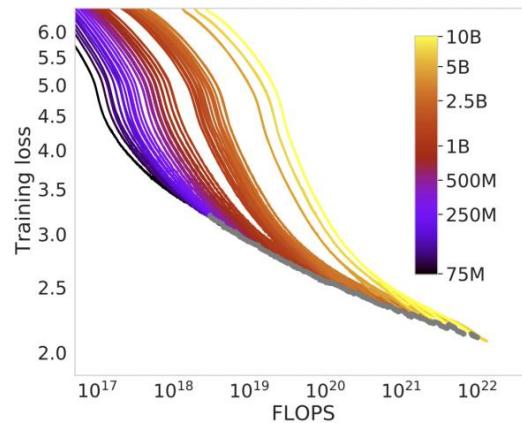
- For decoder-only models, the final performance is only related to **Compute**, **Data Size**, and **Parameter Size**
 - power law relationship for each factor
 - w/o constraints by the others

$$C = \alpha DP$$

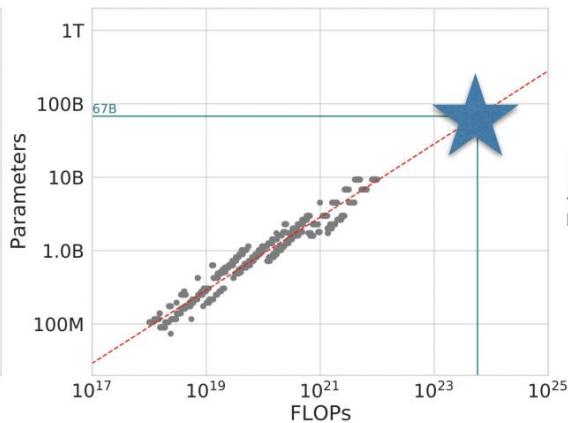


$\text{PF-day} = 10^{15} \times 24 \times 3600 = 8.64 \times 10^{19}$ floating point operations.

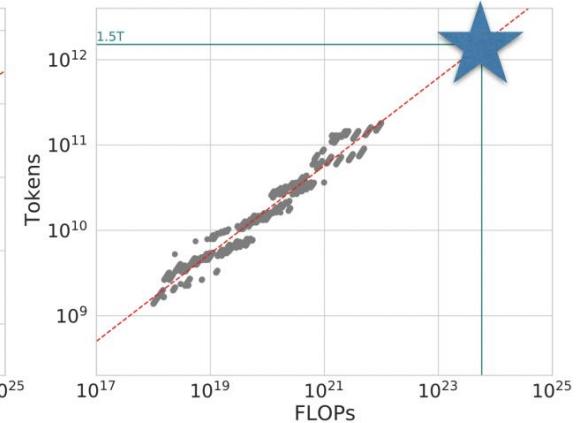
Scaling Law



Run experiments

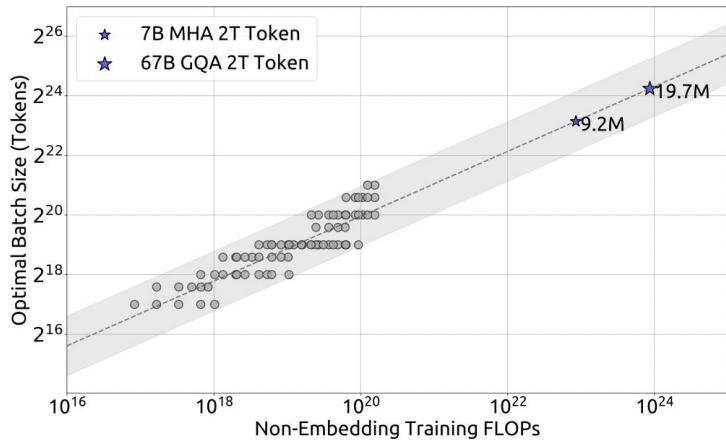


Fit a line and
predict optimal
model size

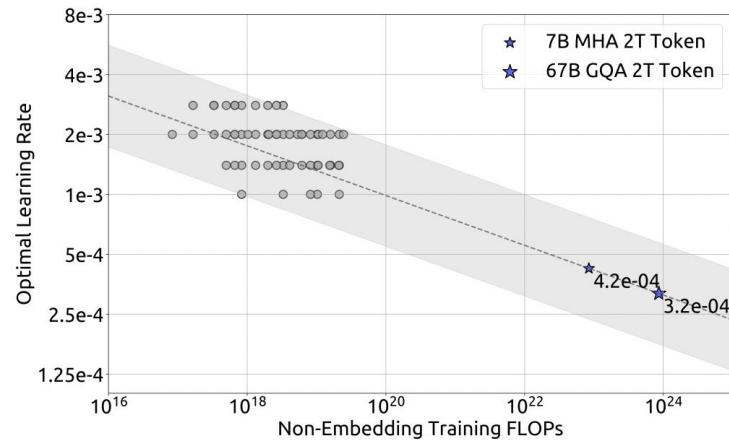


Fit a line and
predict optimal
of tokens

Scaling Law

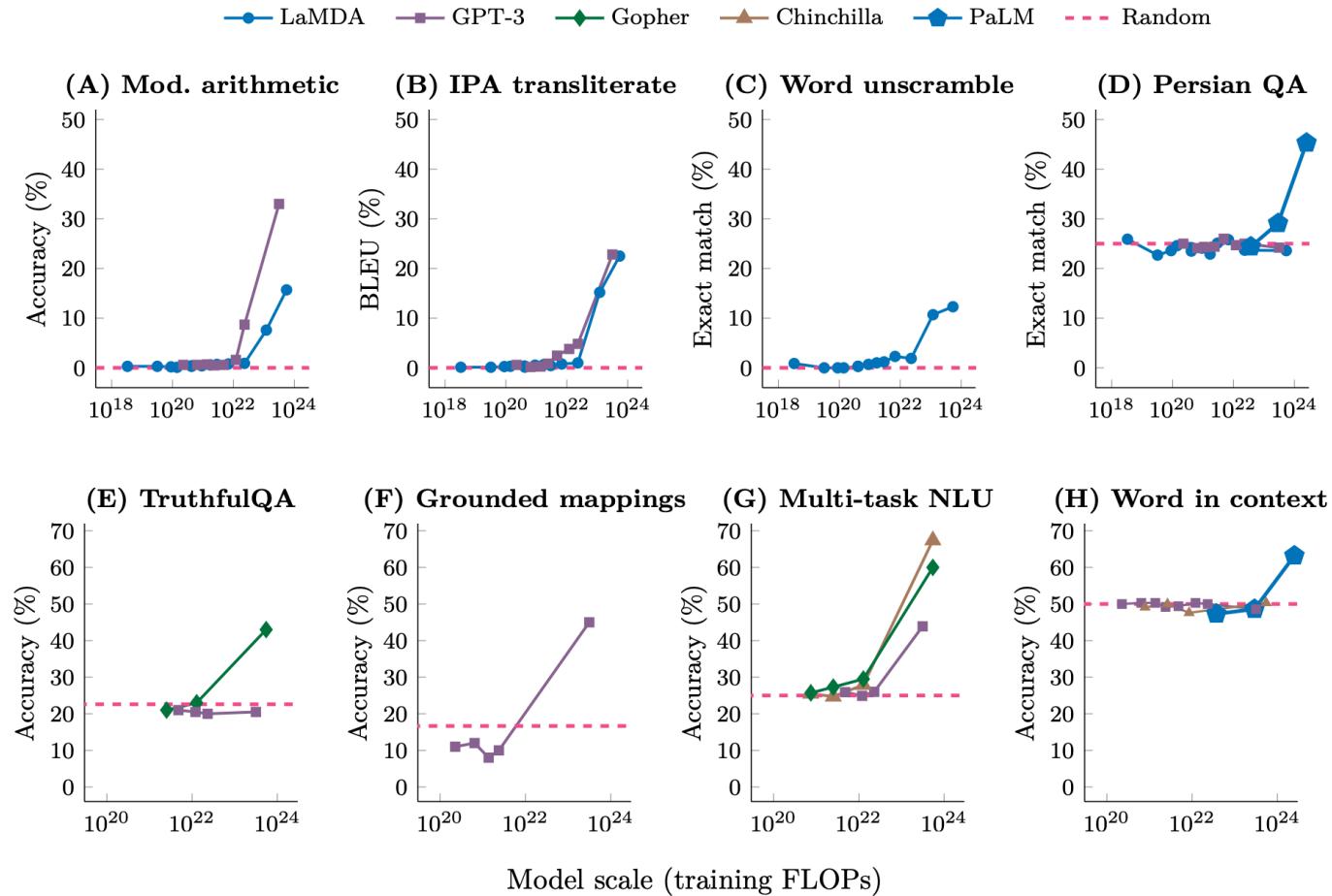


(a) Batch size scaling curve



(b) Learning rate scaling curve

“Emergent” Capability



In-Context Learning

- Scaled models can generalize to new tasks without fine-tuning!
 - Zero-shot
 - Few-shot

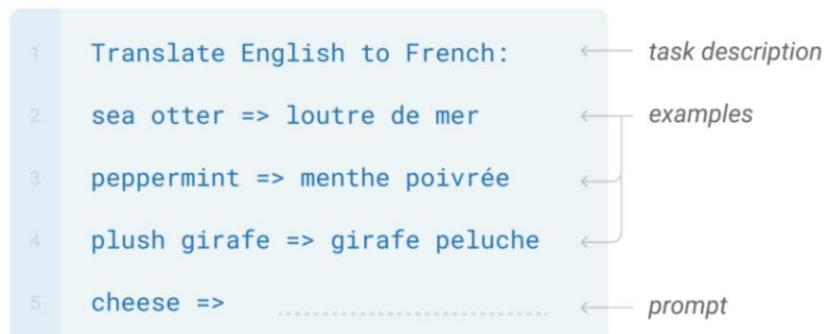
Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.



Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.



Content

- Transformer Architecture
- Improvements on Transformers
- Transformer for different modalities
- Parameter Efficient Tuning
- Scaling Laws
- Quantizing Transformers
- Interpreting Transformer – attention, logitlens

Quantization

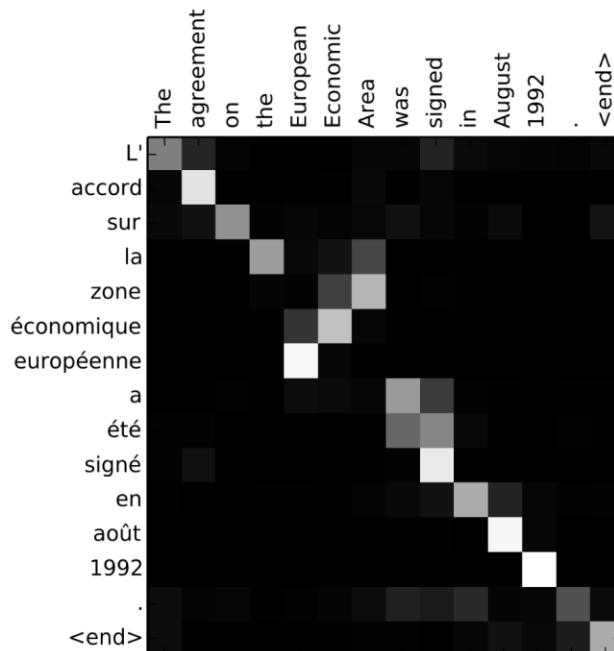
- Human neurons can represent roughly 24 states
- Transformers are in 32 bits, only needs ~4.5
- Basic idea: Weights are floats in a range, represent range with fewer bits
- GPTQ – cuts down to ≤ 4 bits weight, <5% perplexity difference for 1B+ models

Content

- Transformer Architecture
- Improvements on Transformers
- Transformer for different modalities
- Parameter Efficient Tuning
- Scaling Laws
- Quantizing Transformers
- Interpreting Transformer – attention, logitlens

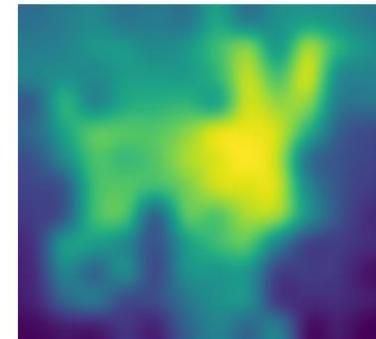
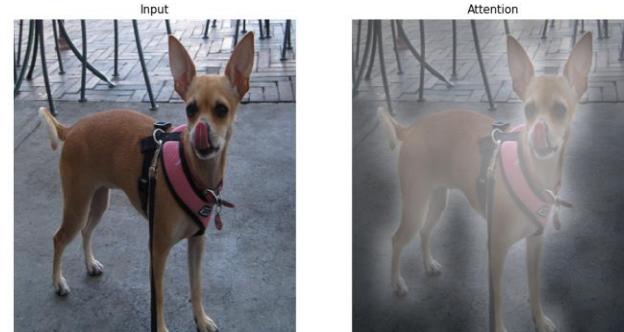
Attention Maps

- Many methods exist, common ones are:
 - Attention maps: aggregate attention weights across all heads in last layer to get avg weights for different tokens



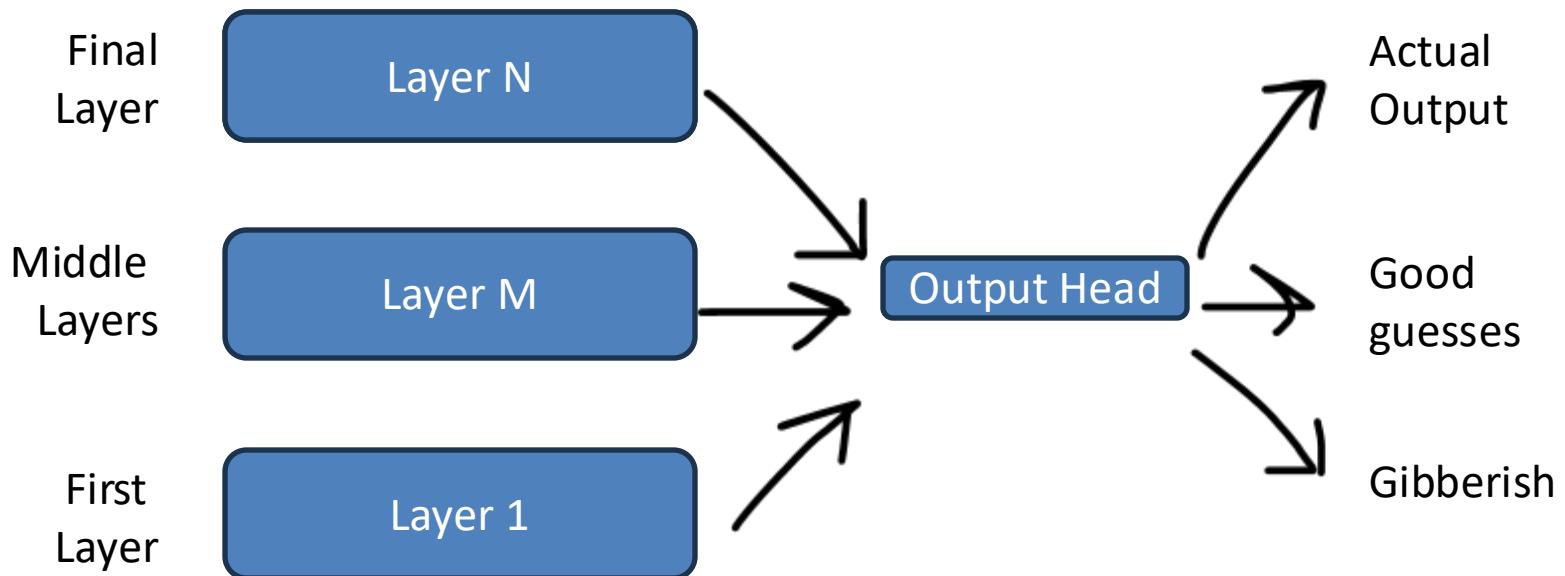
Attention-matrix heatmap

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



Logit Lens

- Many methods exist, common ones are:
 - Logit Lens: output head applied to intermediate layers



We learned...

- Transformer Architecture
- Improvements on Transformers
- Transformer for different modalities
- Parameter Efficient Tuning
- Scaling Laws
- Quantizing Transformers
- Interpreting Transformer – attention, logitlens