

Attention & Transformers

Ivo Verhoeven | Advanced Topics in Computational Semantics

About Me



- 2017 - 2020: BSc. Liberal Arts & Sciences
- 2020 – 2022: MSc. AI at University of Amsterdam
 - Thesis on meta-learning, morphology and translation
 - Took ATCS in 2021
- 2022 - ????: PhD at ILLC
 - Misinformation detection and generalisation with Katia Shutova

Vaswani et al.: Attention is All You Need

- Introduces the Transformer architecture in late 2017
 - Google Brain/Google Research collab

Vaswani et al. (2017). Attention is all you need. Advances in neural information processing systems, 30.

The image shows a screenshot of the research paper "Attention is All You Need" by Vaswani et al. The paper is presented in a clean, modern layout with a white background and black text. At the top center, the title "Attention Is All You Need" is displayed in a bold, sans-serif font. Below the title, there is a horizontal line. The authors' names and their respective institutions are listed in two rows. The first row includes Ashish Vaswani*, Noam Shazeer*, Niki Parmar*, and Jakob Uszkoreit*. The second row includes Llion Jones*, Aidan N. Gomez*, Lukasz Kaiser*, and Illia Polosukhin*. Each author's name is followed by a superscript asterisk or dagger symbol indicating their role. Below the author list, there is a section titled "Abstract" in bold. The abstract text discusses the limitations of traditional sequence transduction models and introduces the Transformer model, which uses attention mechanisms instead of recurrent or convolutional layers. It highlights the model's superior performance on machine translation tasks, achieving high BLEU scores and establishing new state-of-the-art results. The paper is dated December 6, 2017, and is part of the NeurIPS 2017 conference.

06.03762v5 [cs.CL] 6 Dec 2017

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

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

Lukasz Kaiser*
Google Brain
lukaszkaiser@google.com

Illia Polosukhin* ‡
illia.polesukhin@gmail.com

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-Frç our model establishes a new single-model state-of-the-art BLEU! **Attention & Transformers** | 3 of 59

Vaswani et al.: Attention is All You Need

- Introduces the Transformer architecture in late 2017
 - Google Brain/Google Research collab
- Paper currently has **169 248** citations
 - Or **~64 citations a day**

Vaswani et al. (2017). Attention is all you need. Advances in neural information processing systems, 30.

06.03762v5 [cs.CL] 6 Dec 2017

Attention Is All You Need

Ashish Vaswani^{*}
Google Brain
avaswani@google.com

Noam Shazeer^{*}
Google Brain
noam@google.com

Niki Parmar^{*}
Google Research
nikip@google.com

Jakob Uszkoreit^{*}
Google Research
usz@google.com

Llion Jones^{*}
Google Research
llion@google.com

Aidan N. Gomez^{* †}
University of Toronto
aidan@cs.toronto.edu

Łukasz Kaiser^{*}
Google Brain
lukasz.kaiser@google.com

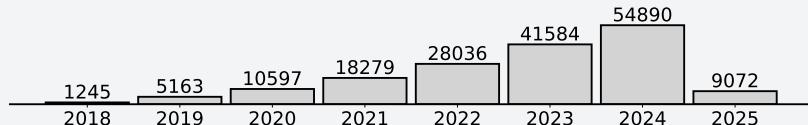
Ilia Polosukhin^{* †}
ilia.pолосухин@gmail.com

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-Frç our model establishes a new single-model state-of-the-art BLEU!  **Attention & Transformers** | 3 of 59

Vaswani et al.: Attention is All You Need

- Introduces the Transformer architecture in late 2017
 - Google Brain/Google Research collab
- Paper currently has **169 248** citations
 - Or **~64 citations a day**
- Number of citations is only accelerating



Vaswani et al. (2017). Attention is all you need. Advances in neural information processing systems, 30.

06.03762v5 [cs.CL] 6 Dec 2017

Attention Is All You Need

Ashish Vaswani^{*}
Google Brain
avaswani@google.com

Noam Shazeer^{*}
Google Brain
noam@google.com

Niki Parmar^{*}
Google Research
nikip@google.com

Jakob Uszkoreit^{*}
Google Research
usz@google.com

Llion Jones^{*}
Google Research
llion@google.com

Aidan N. Gomez^{* †}
University of Toronto
aidan@cs.toronto.edu

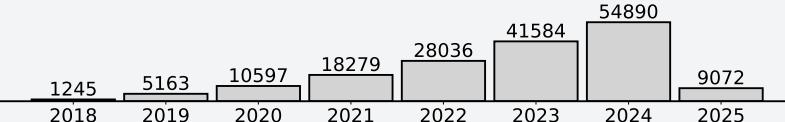
Łukasz Kaiser^{*}
Google Brain
lukasz.kaiser@google.com

Illia Polosukhin^{* †}
illia.polosukhin@gmail.com

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-Frç our model establishes a new single-model state-of-the-art BLEU! Attention & Transformers | 3 of 59

Vaswani et al.: Attention is All You Need

- Introduces the Transformer architecture in late 2017
 - Google Brain/Google Research collab
- Paper currently has **169 248** citations
 - Or **~64 citations a day**
- Number of citations is only accelerating
- Most cited paper ever has **233 829** citations

Lowry et al. (1951) Protein measurement with the folin phenol reagent.

Vaswani et al. (2017). Attention is all you need. Advances in neural information processing systems, 30.

06.03762v5 [cs.CL] 6 Dec 2017

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com Noam Shazeer*
Google Brain
noam@google.com Niki Parmar*
Google Research
nikip@google.com Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu Lukasz Kaiser*
Google Brain
lukaszkaiser@google.com

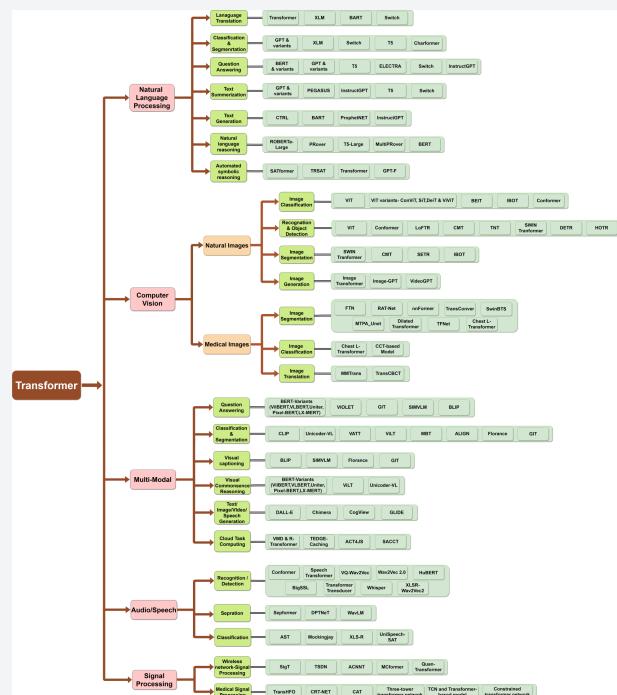
Illia Polosukhin* †
illia.polosukhin@gmail.com

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-Fr, our model establishes a new single-model state-of-the-art BLEU!  Attention & Transformers | 3 of 59

Vaswani et al.: Attention is All You Need

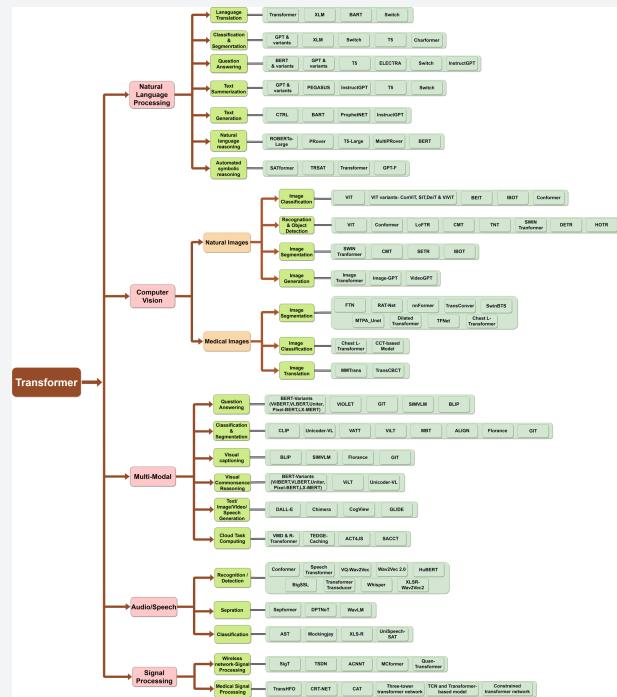
- It's hard to think of an AI area that hasn't been affected by the Transformer



Islam, et al. (2023). A Comprehensive Survey on Application Transformers for Deep Learning Tasks. arXiv:2306.07303.

Vaswani et al.: Attention is All You Need

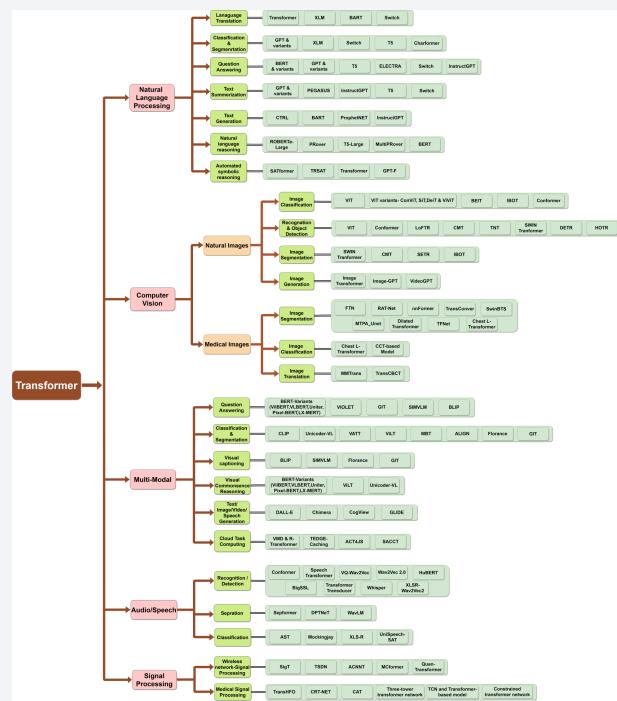
- It's hard to think of an AI area that hasn't been affected by the Transformer
 - **NLP:** Transformer > RNN
 - Seq-to-seq: what it was designed for
 - Classification: encoder-only transformers
 - Generation: decoder-only transformers



Islam, et al. (2023). A Comprehensive Survey on Application Transformers for Deep Learning Tasks. arXiv:2306.07303.

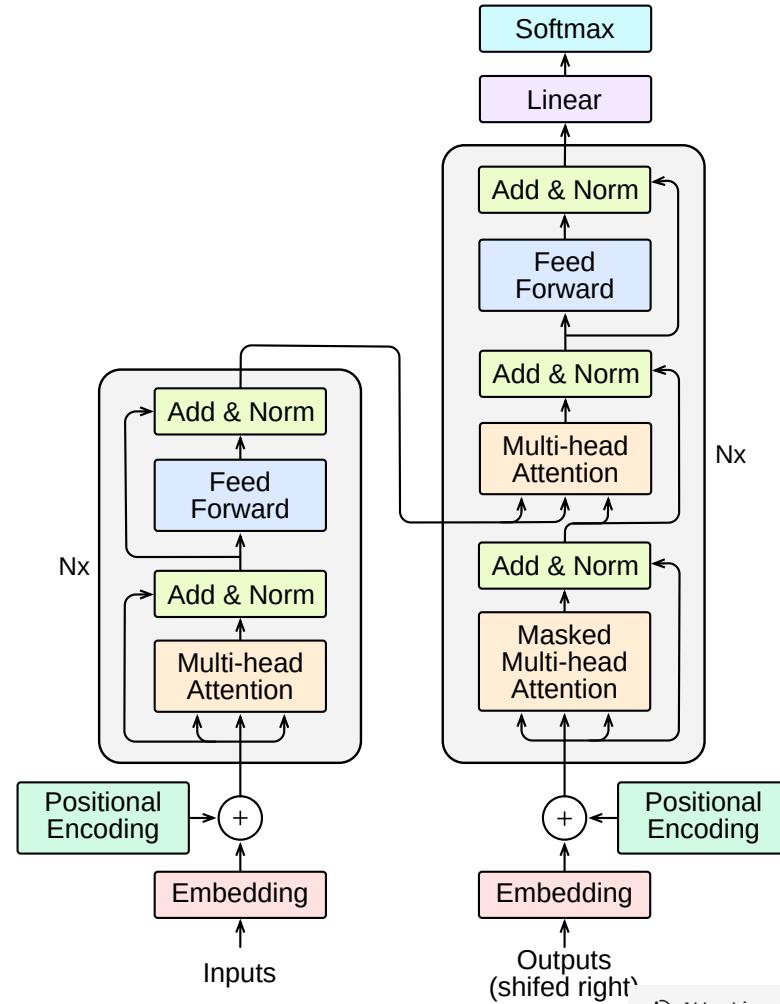
Vaswani et al.: Attention is All You Need

- It's hard to think of an AI area that hasn't been affected by the Transformer
 - **NLP:** Transformer > RNN
 - Seq-to-seq: what it was designed for
 - Classification: encoder-only transformers
 - Generation: decoder-only transformers
 - **CV:** ViT > CNN
 - **Multi-modal:** Transformer > different architectures
 - **Speech:** Transformer > CNN
 - **Graphs:** Transformer/Attention > GCN

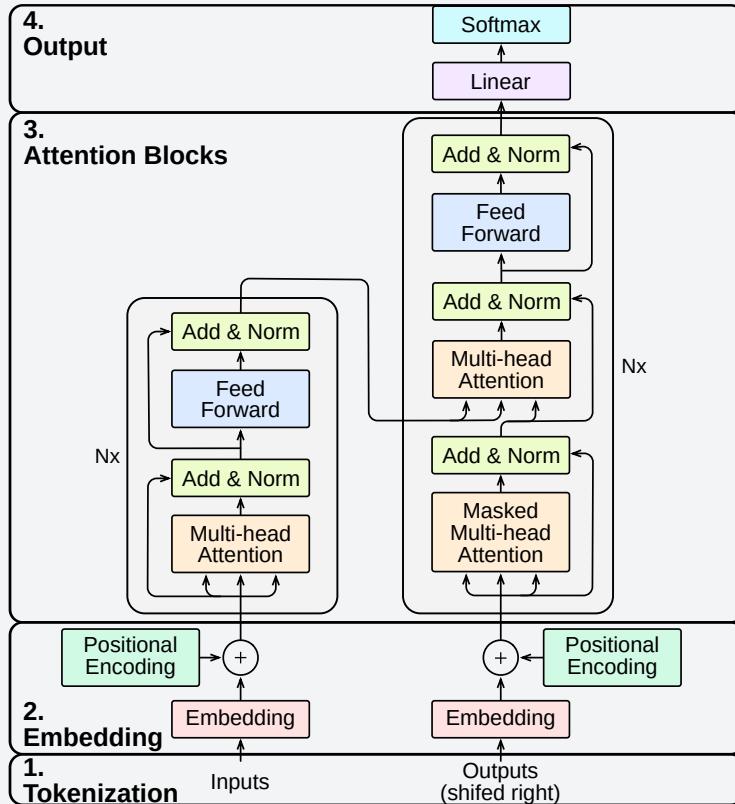


Islam, et al. (2023). A Comprehensive Survey on Application Transformers for Deep Learning Tasks. arXiv:2306.07303.

The Transformer



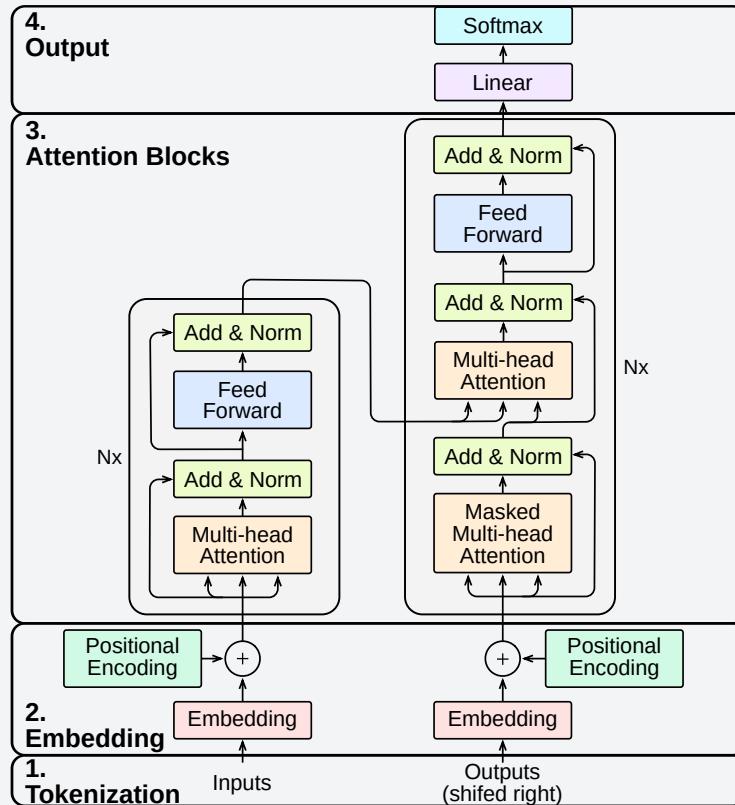
Breaking the Transformer into modules



Breaking the Transformer into modules

4. Output

- Softmax
- Linear



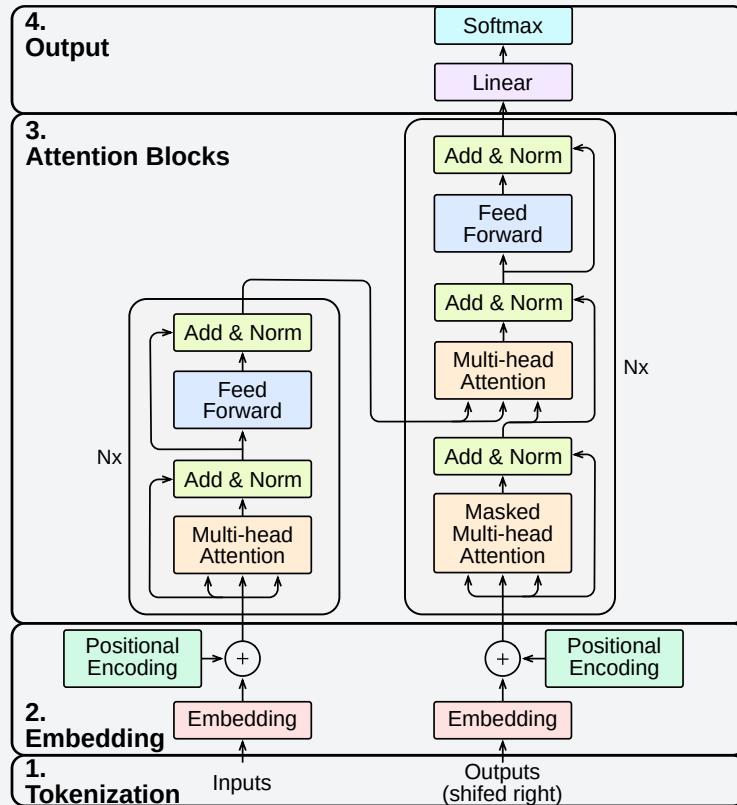
Breaking the Transformer into modules

4. Output

- Softmax
- Linear

3. Attention Blocks

- Multi-head Attention
- Add & Norm
- Feed Forward



Breaking the Transformer into modules

4. Output

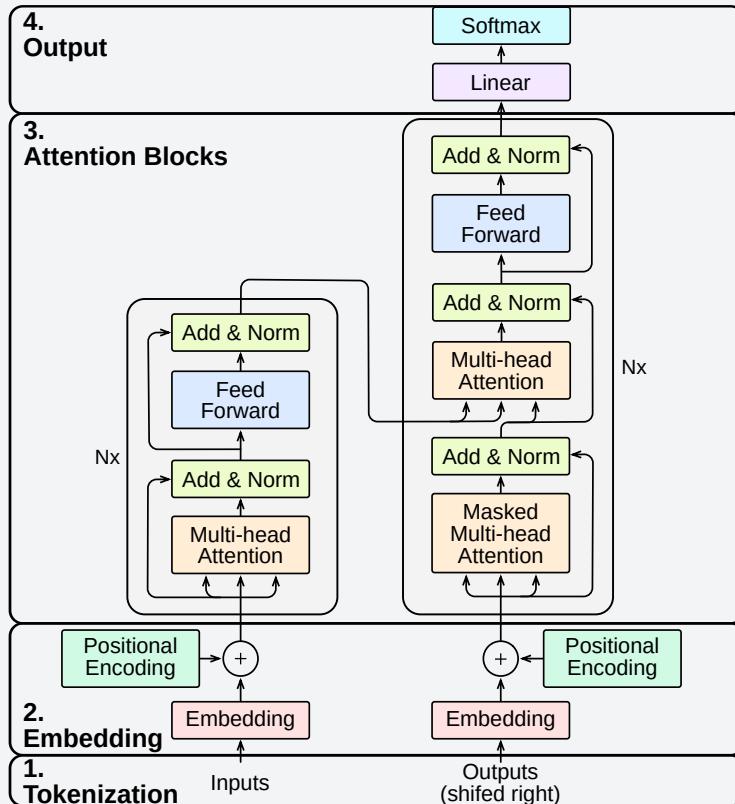
- Softmax
- Linear

3. Attention Blocks

- Multi-head Attention
- Add & Norm
- Feed Forward

2. Embedding

- Token Embedding
- Positional Encoding



Breaking the Transformer into modules

4. Output

- Softmax
- Linear

3. Attention Blocks

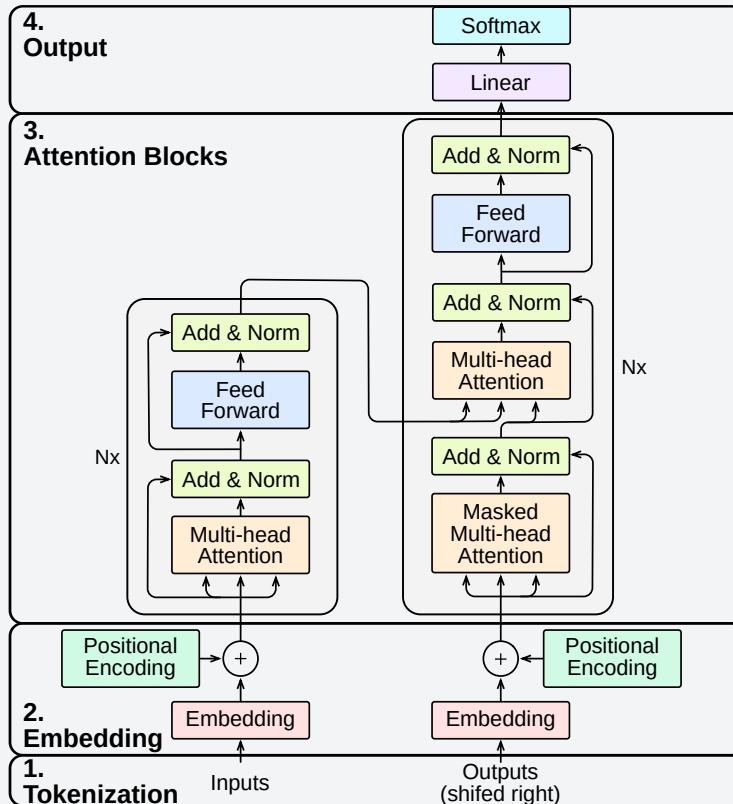
- Multi-head Attention
- Add & Norm
- Feed Forward

2. Embedding

- Token Embedding
- Positional Encoding

1. Tokenization

- (Not pictured)



Breaking the Transformer into modules

4. Output

- Softmax
- Linear

3. Attention Blocks

- Multi-head Attention
- Add & Norm
- Feed Forward

2. Embedding

- Token Embedding
- Positional Encoding

1. Tokenization

- (Not pictured)

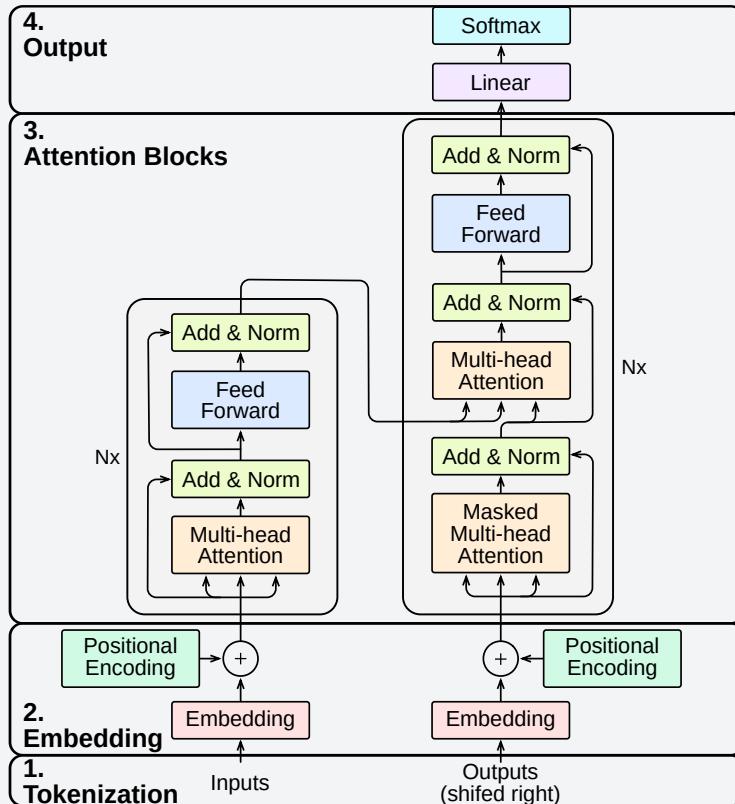
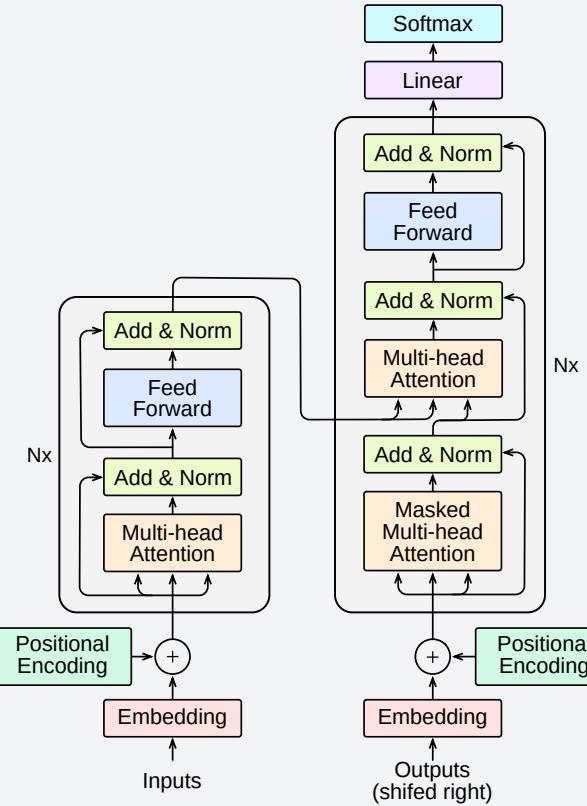
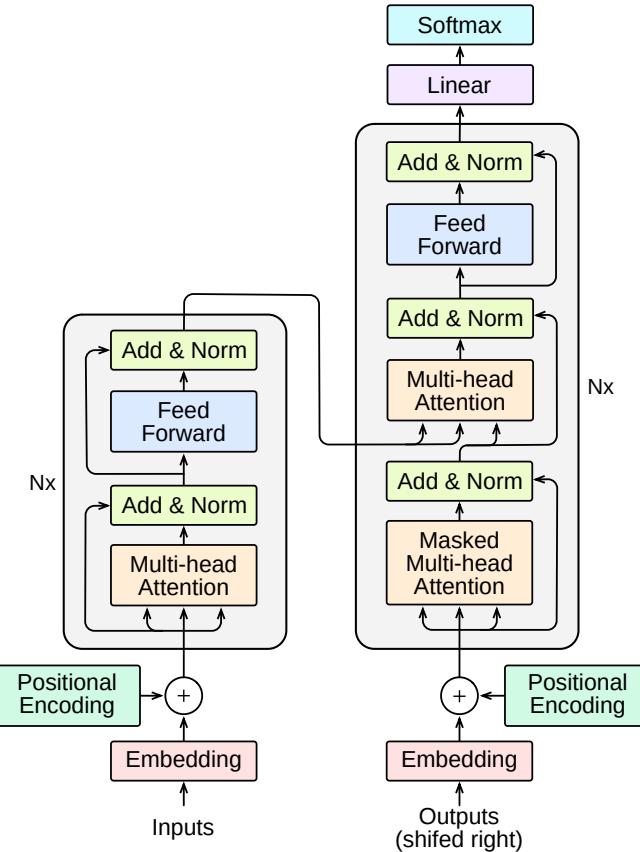


Table of Contents

1. Multi-head Attention
 1. Definition & Properties
 2. Non-Transformer Examples
 3. Attention in Transformers
 4. Multi-head Attention
 5. Why Attention?
2. Feed Forward
3. Add & Norm
 1. Residual Connections
 2. Layer Norm
4. Embeddings
 1. Token Embedding
 2. Positional Encoding
5. Tokenization
6. Training Transformers



Multi-head Attention



Definition & Properties

Multi-head Attention

Definition & Properties

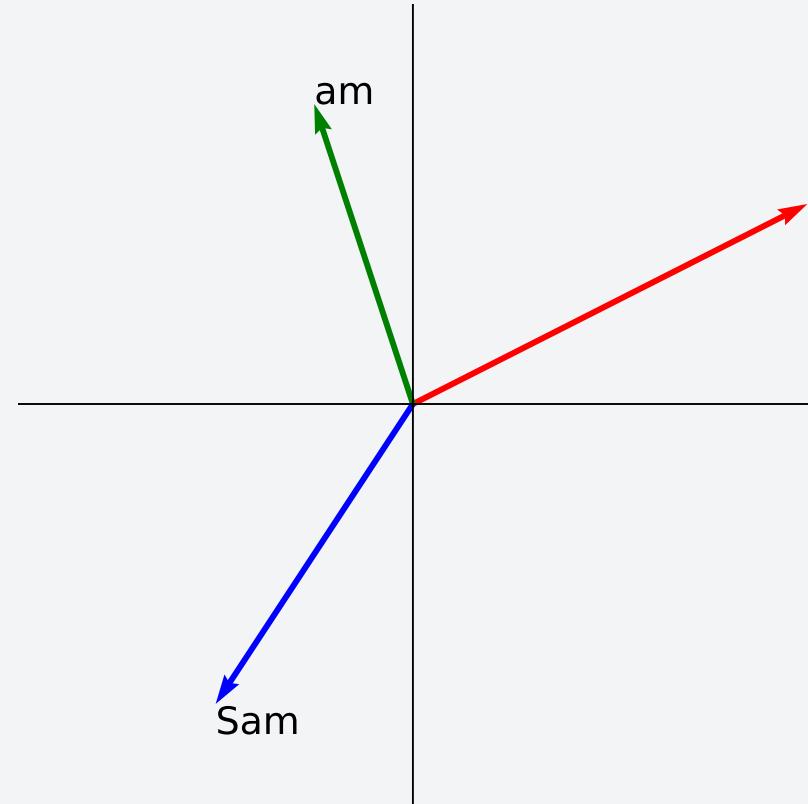
Multi-head Attention

Definition & Properties

Multi-head Attention

- The result of **Attention** is just a convex combination of \mathbf{V}

$$\begin{bmatrix} 0.6 & 0.1 & 0.3 \\ 0.3 & 0.5 & 0.2 \\ 0.2 & 0.1 & 0.7 \end{bmatrix} \begin{bmatrix} 2.0 & 1.0 \\ -0.5 & 2.0 \\ -1.0 & -0.5 \end{bmatrix} \begin{matrix} \mathbf{I} \\ \mathbf{am} \\ \mathbf{Sam} \end{matrix}$$



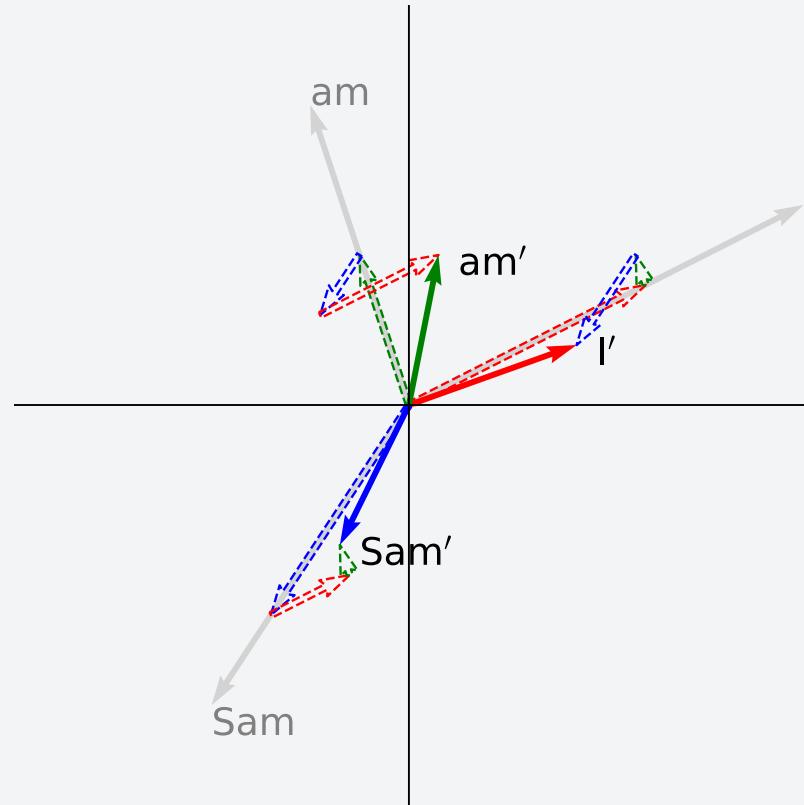
Definition & Properties

Multi-head Attention

- The result of **Attention** is just a convex combination of **V**

$$\begin{bmatrix} 0.6 & 0.1 & 0.3 \\ 0.3 & 0.5 & 0.2 \\ 0.2 & 0.1 & 0.7 \end{bmatrix} \begin{bmatrix} 2.0 & 1.0 \\ -0.5 & 2.0 \\ -1.0 & -0.5 \end{bmatrix} \begin{matrix} I \\ am \\ Sam \end{matrix}$$

$$= \begin{bmatrix} 0.6 * I + 0.1 * am + 0.3 * Sam \\ 0.3 * I + 0.5 * am + 0.2 * Sam \\ 0.2 * I + 0.1 * am + 0.7 * Sam \end{bmatrix}$$

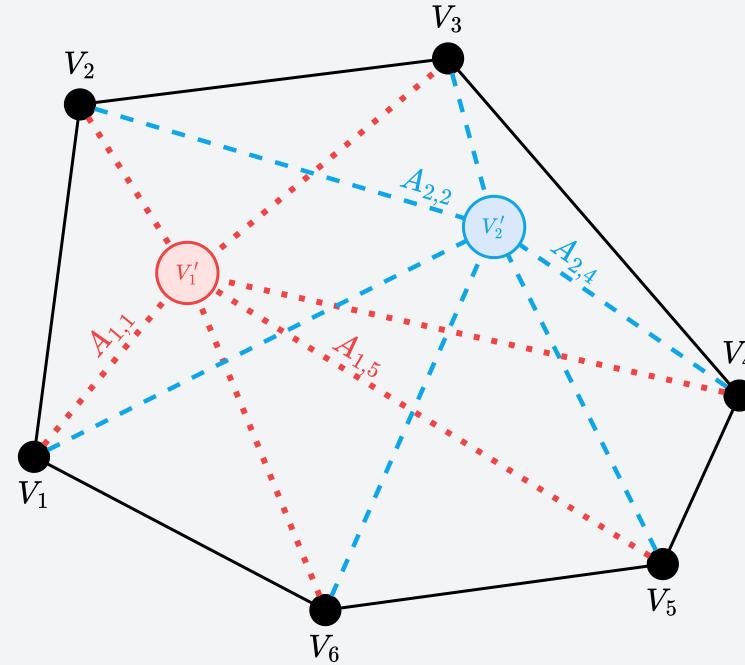


Definition & Properties

Multi-head Attention

Convex Combination

The elements of V' will lie inside the convex hull of all of the elements in V



Definition & Properties

Multi-head Attention

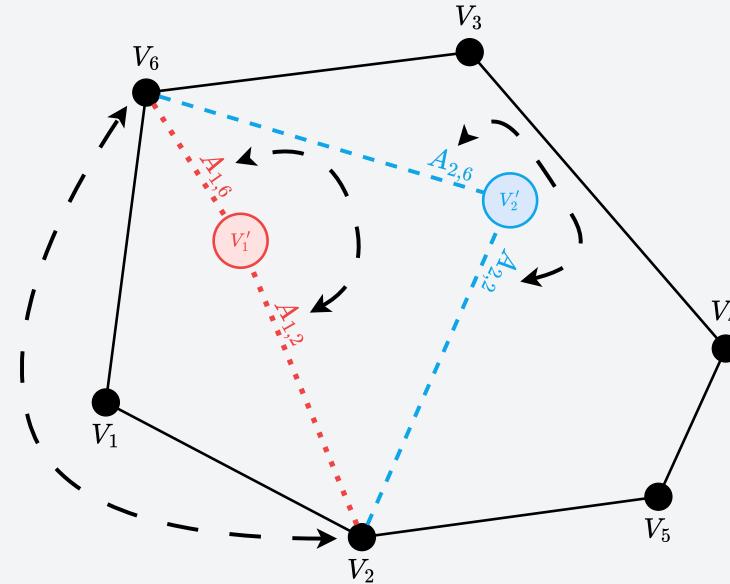
Convex Combination

The elements of V' will lie inside the convex hull of all of the elements in V

Permutation Equivariance

The elements of V' are *equivariant* to a change in the order of the columns of A and the rows of V

- Attention does not care about word order
 - 'I am Sam' ~ 'Sam I am'



Definition & Properties

Multi-head Attention

So is **Attention** just a linear map?

- Not quite

Linear maps are:

Definition & Properties

Multi-head Attention

So is **Attention** just a linear map?

- Not quite

Linear maps are:

- Inflexible in terms of sequence length

Definition & Properties

Multi-head Attention

So is **Attention** just a linear map?

- Not quite

Linear maps are:

- Inflexible in terms of sequence length
- Parameter inefficient

Definition & Properties

Multi-head Attention

So is **Attention** just a linear map?

- Not quite

Linear maps are:

- Inflexible in terms of sequence length
- Parameter inefficient
- Invariant to the input content

Definition & Properties

Multi-head Attention

- Let \mathbf{V} be a matrix of **value** vectors
 - It has a sequence length of t_V
 - It has a dimensionality of d_V
- Let \mathbf{K} be a matrix of **key** vectors
 - It has a sequence length of t_V
 - It has a dimensionality of d_K
- Let \mathbf{Q} be a matrix of **query** vectors
 - It has a sequence length of t_Q
 - It has a dimensionality of d_K

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \underbrace{\text{softmax}(f(\mathbf{Q}, \mathbf{K}))\mathbf{V}}_{\mathbf{A}}$$

$$\mathbf{A} \in (0, 1)^{[t_Q \times t_V]}$$

$$\mathbf{V} \in \mathbb{R}^{[t_V \times d_v]}$$

$$\mathbf{K} \in \mathbb{R}^{[t_V \times d_k]}$$

$$\mathbf{Q} \in \mathbb{R}^{[t_Q \times d_k]}$$

Definition & Properties

Multi-head Attention

- Let \mathbf{V} be a matrix of **value** vectors
 - It has a sequence length of t_V
 - It has a dimensionality of d_V
- Let \mathbf{K} be a matrix of **key** vectors
 - It has a sequence length of t_V
 - It has a dimensionality of d_K
- Let \mathbf{Q} be a matrix of **query** vectors
 - It has a sequence length of t_Q
 - It has a dimensionality of d_K
- Let $f(\mathbf{Q}, \mathbf{K})$ be some kernel function
 - Read: similarity function

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \underbrace{\text{softmax}(f(\mathbf{Q}, \mathbf{K}))\mathbf{V}}_{\mathbf{A}}$$

$$\mathbf{A} \in (0, 1)^{[t_Q \times t_V]}$$

$$\mathbf{V} \in \mathbb{R}^{[t_V \times d_v]}$$

$$\mathbf{K} \in \mathbb{R}^{[t_V \times d_k]}$$

$$\mathbf{Q} \in \mathbb{R}^{[t_Q \times d_k]}$$

Definition & Properties

Multi-head Attention

- Let \mathbf{V} be a matrix of **value** vectors
 - It has a sequence length of t_V
 - It has a dimensionality of d_V
- Let \mathbf{K} be a matrix of **key** vectors
 - It has a sequence length of t_V
 - It has a dimensionality of d_K
- Let \mathbf{Q} be a matrix of **query** vectors
 - It has a sequence length of t_Q
 - It has a dimensionality of d_K
- Let $f(\mathbf{Q}, \mathbf{K})$ be some kernel function
 - Read: similarity function

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \underbrace{\text{softmax}(f(\mathbf{Q}, \mathbf{K}))\mathbf{V}}_{\mathbf{A}}$$

$$\mathbf{A} \in (0, 1)^{[t_Q \times t_V]}$$

$$\mathbf{V} \in \mathbb{R}^{[t_V \times d_v]}$$

$$\mathbf{K} \in \mathbb{R}^{[t_V \times d_k]}$$

$$\mathbf{Q} \in \mathbb{R}^{[t_Q \times d_k]}$$

Non-Transformer Examples

Multi-head Attention

- **V** contains information
- **K** contains information about information (i.e., metadata)
- **Q** contains metadata about what we want from **V**
- $f(\mathbf{Q}, \mathbf{K})$ is high when **Q** is similar to **K**

Non-Transformer Examples

Multi-head Attention

- **V** contains information
- **K** contains information about information (i.e., metadata)
- **Q** contains metadata about what we want from **V**
- $f(\mathbf{Q}, \mathbf{K})$ is high when **Q** is similar to **K**

E Soft lookup

We want to find a textbook about NLP in the library (**V**). We search for titles (**K**) with "jurafsky" and "martin" as authors (**Q**). The computer returns books with similar titles (f)

The screenshot shows a library search interface. At the top, the University of Amsterdam logo and name are displayed, along with icons for search, filters, and help. The search bar contains the query "jurafsky martin". Below the search bar, a red header bar says "CataloguePlus". The main area shows a search result for "1-2 of 2 Results". The first result is a book titled "Speech and language processing : an introduction to natural language processing, computational linguistics, and speech recognition / Daniel Jurafsky, James H. Martin." by Jurafsky, Dan, 1962-; Martin, James H., 1959-. It includes a green "Available at Library Learning Centre" link and other location details.

Non-Transformer Examples

Multi-head Attention

- **Q** and **V** do not need to have the same sequence length
- The output of f is *always* a matrix of size $\mathbf{A} \in (0, 1)^{[t_Q \times t_V]}$

Non-Transformer Examples

Multi-head Attention

- \mathbf{Q} and \mathbf{V} do not need to have the same sequence length
- The output of f is *always* a matrix of size $\mathbf{A} \in (0, 1)^{[t_Q \times t_V]}$

E Nadaraya-Watson Kernel Regression

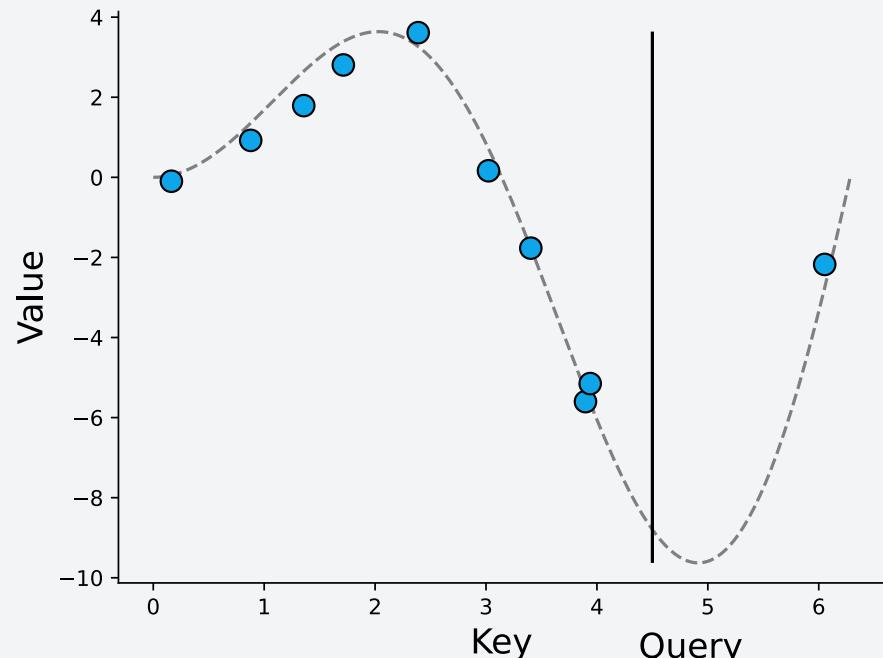
We have some sequence of values

$$\mathcal{D} = [(1.36, 1.79), (3.40, -1.77) \dots, (6.05, -2.17)]$$

We want to predict a new sample at $x = 4.25$

We compute the negative Euclidean distance of our new sample with all training samples (f). We normalize the outputs to lie between $(0, 1)$

We compute our predicted value as the mean of the seen values, weighted by the computed similarities



Non-Transformer Examples

Multi-head Attention

- \mathbf{Q} and \mathbf{V} do not need to have the same sequence length
- The output of f is *always* a matrix of size $\mathbf{A} \in (0, 1)^{[t_Q \times t_V]}$

E Nadaraya-Watson Kernel Regression

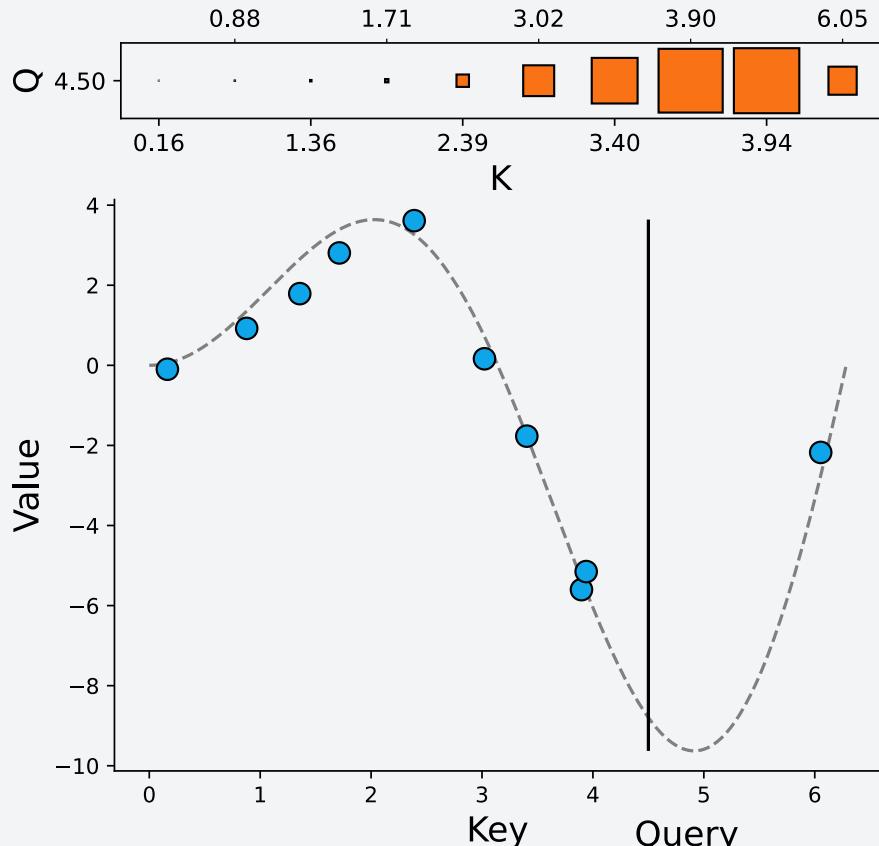
We have some sequence of values

$$\mathcal{D} = [(1.36, 1.79), (3.40, -1.77) \dots, (6.05, -2.17)]$$

We want to predict a new sample at $x = 4.25$

We compute the negative Euclidean distance of our new sample with all training samples (f). We normalize the outputs to lie between $(0, 1)$

We compute our predicted value as the mean of the seen values, weighted by the computed similarities



Non-Transformer Examples

Multi-head Attention

- \mathbf{Q} and \mathbf{V} do not need to have the same sequence length
- The output of f is *always* a matrix of size $\mathbf{A} \in (0, 1)^{[t_Q \times t_V]}$

E Nadaraya-Watson Kernel Regression

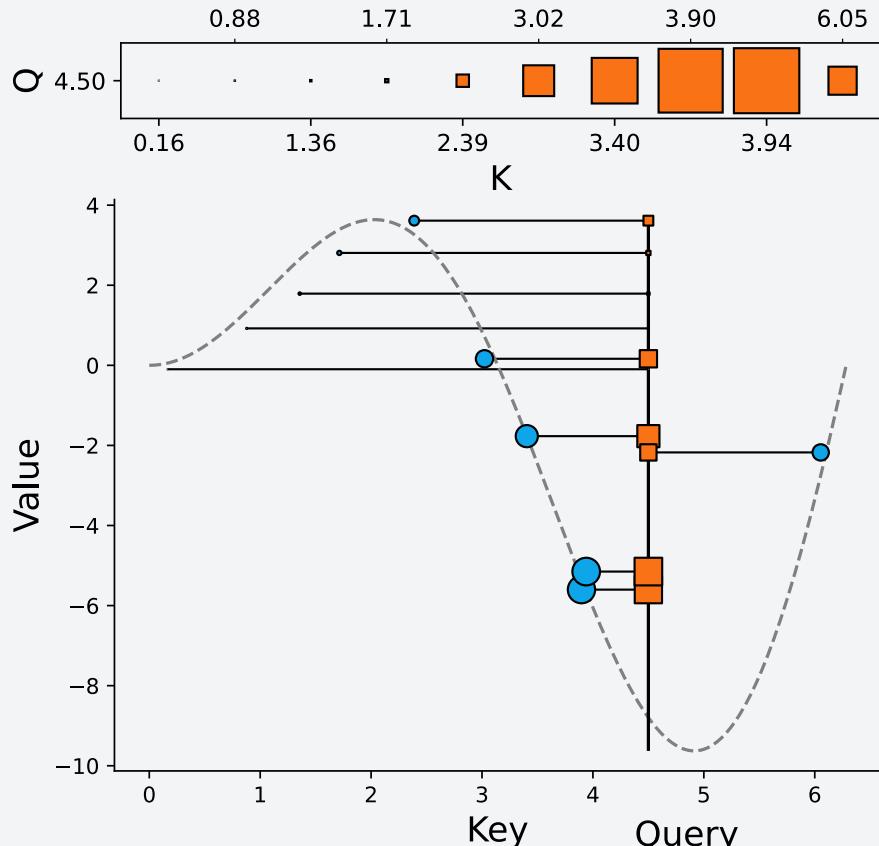
We have some sequence of values

$$\mathcal{D} = [(1.36, 1.79), (3.40, -1.77) \dots, (6.05, -2.17)]$$

We want to predict a new sample at $x = 4.25$

We compute the negative Euclidean distance of our new sample with all training samples (f). We normalize the outputs to lie between $(0, 1)$

We compute our predicted value as the mean of the seen values, weighted by the computed similarities



Non-Transformer Examples

Multi-head Attention

- \mathbf{Q} and \mathbf{V} do not need to have the same sequence length
- The output of f is *always* a matrix of size $\mathbf{A} \in (0, 1)^{[t_Q \times t_V]}$

E Nadaraya-Watson Kernel Regression

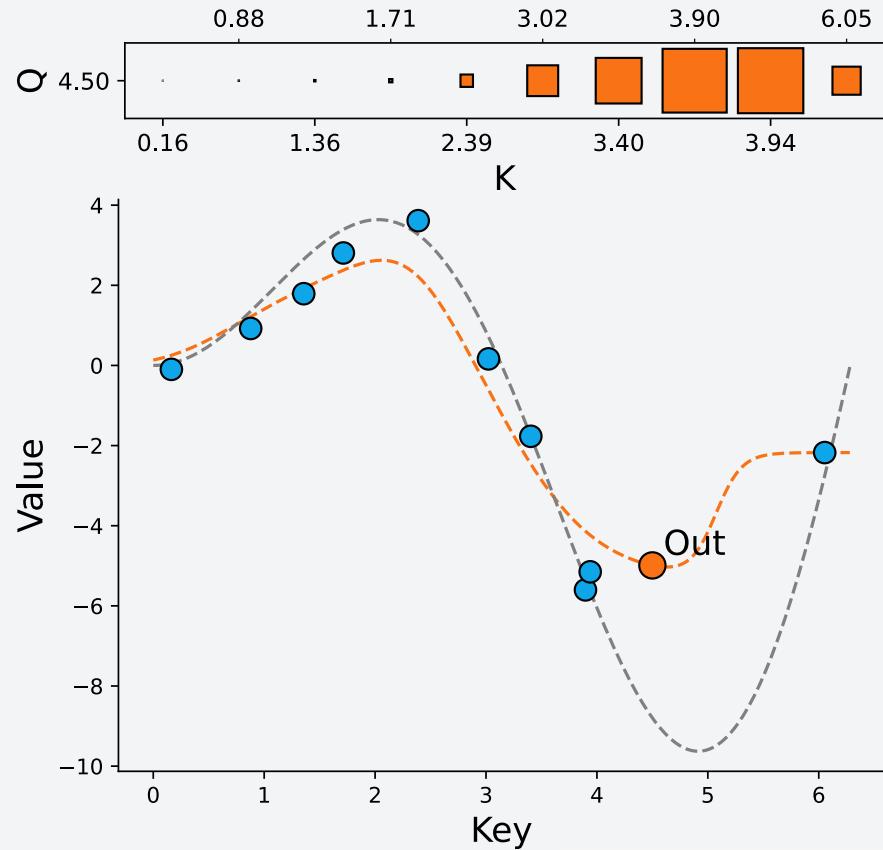
We have some sequence of values

$$\mathcal{D} = [(1.36, 1.79), (3.40, -1.77) \dots, (6.05, -2.17)]$$

We want to predict a new sample at $x = 4.25$

We compute the negative Euclidean distance of our new sample with all training samples (f). We normalize the outputs to lie between $(0, 1)$

We compute our predicted value as the mean of the seen values, weighted by the computed similarities



Non-Transformer Examples

Multi-head Attention

- $f(\mathbf{Q}, \mathbf{K})$ is high when q_i is similar to k_j
- Attention matrix tells us how important v_j is to q_i

Non-Transformer Examples

Multi-head Attention

- $f(\mathbf{Q}, \mathbf{K})$ is high when q_i is similar to k_j
- Attention matrix tells us how important v_j is to q_i

E Bahdanau et al. Alignment

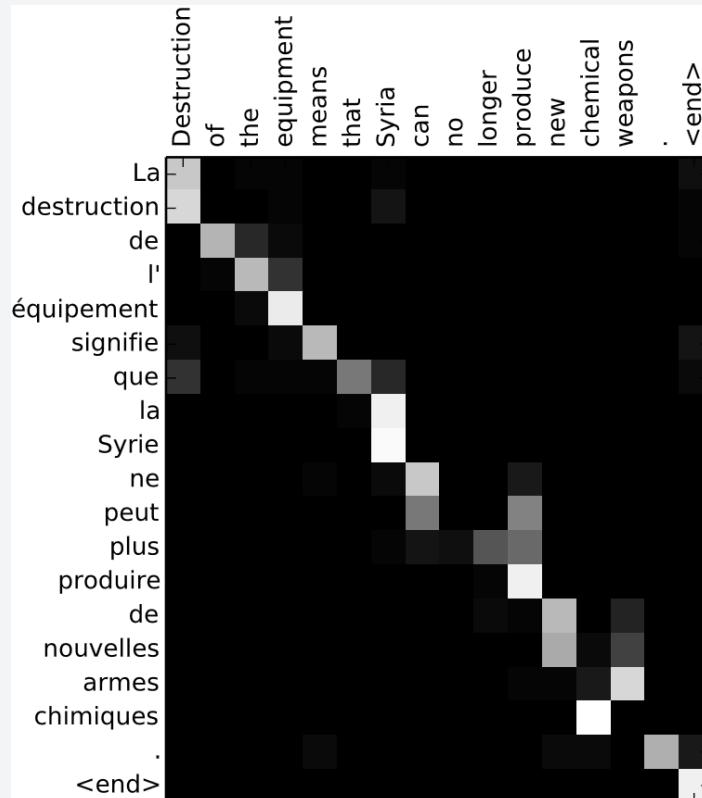
In Neural Machine Translation (NMT) the encoder generates a representation of the input language

The decoder needs to generate in a target language

Token in input language \neq token in output language

Solution: have each token in the target language (\mathbf{Q}) attend back to all input language tokens (\mathbf{K} , \mathbf{V})

Bahdanau, Cho & Bengio (2014). Neural machine translation by jointly learning to align and translate.
arXiv preprint arXiv:1409.0473.



Non-Transformer Examples

Multi-head Attention

$$\text{attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}(f(\mathbf{Q}, \mathbf{K})) \mathbf{V}$$

Non-Transformer Examples

Multi-head Attention

$$\text{attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}(f(\mathbf{Q}, \mathbf{K})) \mathbf{V}$$

Model	$f(\mathbf{Q}, \mathbf{K})$
Gaussian	$\log \exp \left(\frac{(\mathbf{q} - \mathbf{k})^2}{\sigma^2} \right), \quad \forall \mathbf{q}, \mathbf{k} \in \mathbf{Q}, \mathbf{K}$
Cosine	$\frac{\mathbf{Q}\mathbf{K}^\top}{\ \mathbf{Q}\ \ \mathbf{K}\ }$
Additive	$\mathbf{v}^\top \tanh(\mathbf{W} [\mathbf{q} \mathbf{k}]), \quad \forall \mathbf{q}, \mathbf{k} \in \mathbf{Q}, \mathbf{K}$
General	$\mathbf{Q}\mathbf{W}\mathbf{K}^\top$

Attention in Transformers

Multi-head Attention

- Transformer Attention uses a **masked scaled dot-product** kernel function

$$f(\mathbf{Q}, \mathbf{K}) = \frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}$$

- \mathbf{Q} is of size $t_Q \times d_K$
- \mathbf{K} is of size $t_V \times d_K$

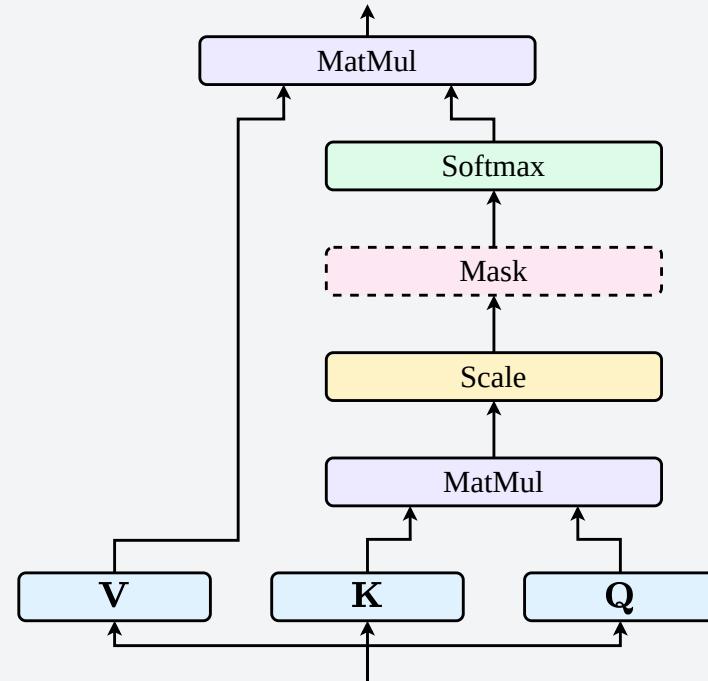
Attention in Transformers

Multi-head Attention

- Transformer Attention uses a **masked scaled dot-product** kernel function

$$f(\mathbf{Q}, \mathbf{K}) = \frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}$$

- \mathbf{Q} is of size $t_Q \times d_K$
- \mathbf{K} is of size $t_V \times d_K$



Attention in Transformers

Multi-head Attention

- Transformer Attention uses a **masked scaled dot-product** kernel function

$$f(\mathbf{Q}, \mathbf{K}) = \frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}$$

- Why scale?
 - Assume the elements in \mathbf{Q} and \mathbf{K} come from *independent* normal distributions:
 - The distribution of their dot-product is:

$$\mathbf{q}, \mathbf{k} \sim \mathcal{N}(0, 1)$$

$$\mathbf{q}^T \mathbf{k} \sim \mathcal{N}(0, \sqrt{d_k})$$

$$\begin{aligned}\text{var} [\mathbf{q}^T \mathbf{k}] &= \text{var} \left[\sum_i^{d_k} q_i k_i \right] \\ &= \sum_i^{d_k} \text{var} [q_i k_i] \\ &= \sum_i^{d_k} \text{var} [q_i] \text{var} [k_i] \\ &= \sum_i^{d_k} 1 \cdot 1 \\ &= d_k\end{aligned}$$

Attention in Transformers

Multi-head Attention

- Transformer Attention uses a **masked scaled dot-product** kernel function

$$f(\mathbf{Q}, \mathbf{K}) = \frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}$$

- Why mask?
 - Currently all tokens are treated equally

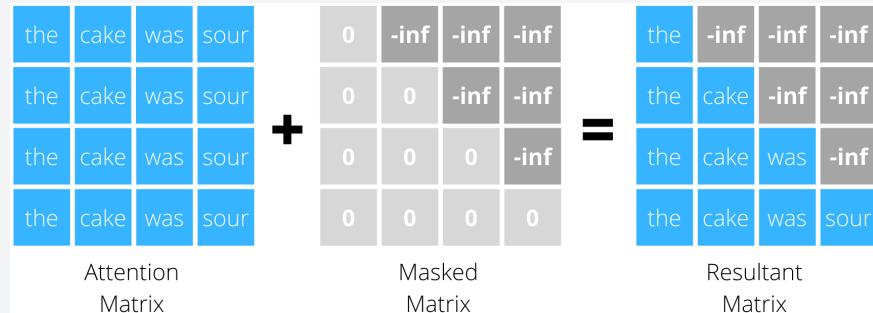
Attention in Transformers

Multi-head Attention

- Transformer Attention uses a **masked scaled dot-product** kernel function

$$f(\mathbf{Q}, \mathbf{K}) = \frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}$$

- Why mask?
 - Currently all tokens are treated equally
 - Causal masking:** decoder tokens should never attend to future tokens, only to the past



<https://krypticmouse.hashnode.dev/attention-is-all-you-need>

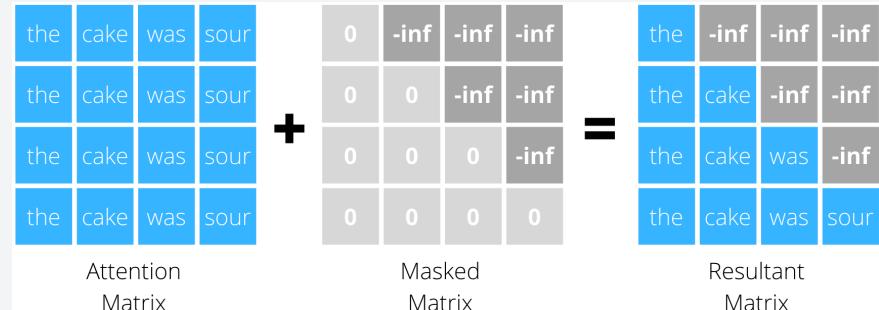
Attention in Transformers

Multi-head Attention

- Transformer Attention uses a **masked scaled dot-product** kernel function

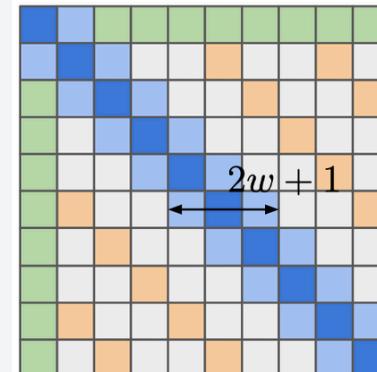
$$f(\mathbf{Q}, \mathbf{K}) = \frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}$$

- Why mask?
 - Currently all tokens are treated equally
 - Causal masking:** decoder tokens should never attend to future tokens, only to the past
 - Local/Global masking:** sometimes local attention is all you need

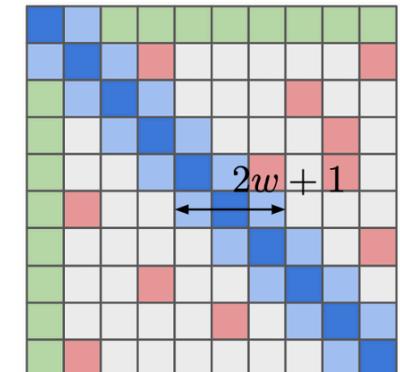


<https://krypticmouse.hashnode.dev/attention-is-all-you-need>

Longformer



Big Bird

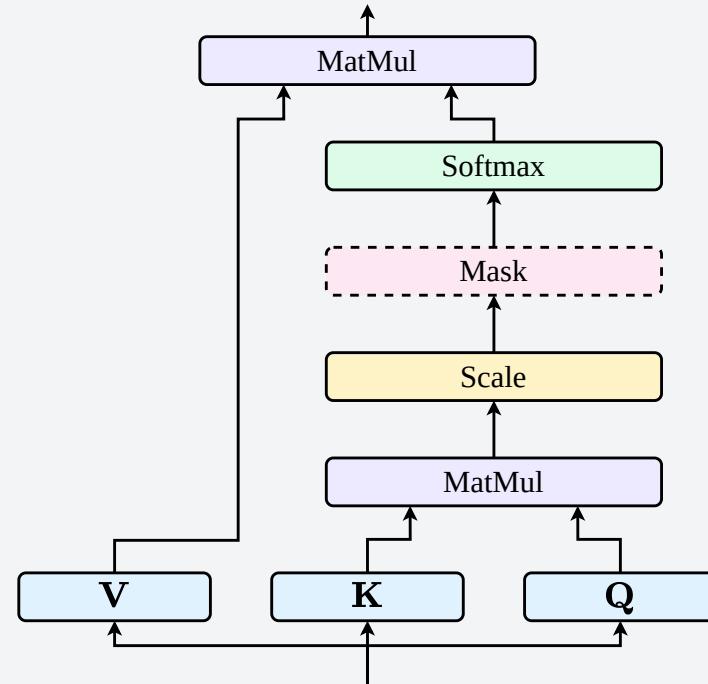


<https://lilianweng.github.io/posts/2023-01-27-the-transformer-family-v2/>

Attention in Transformers

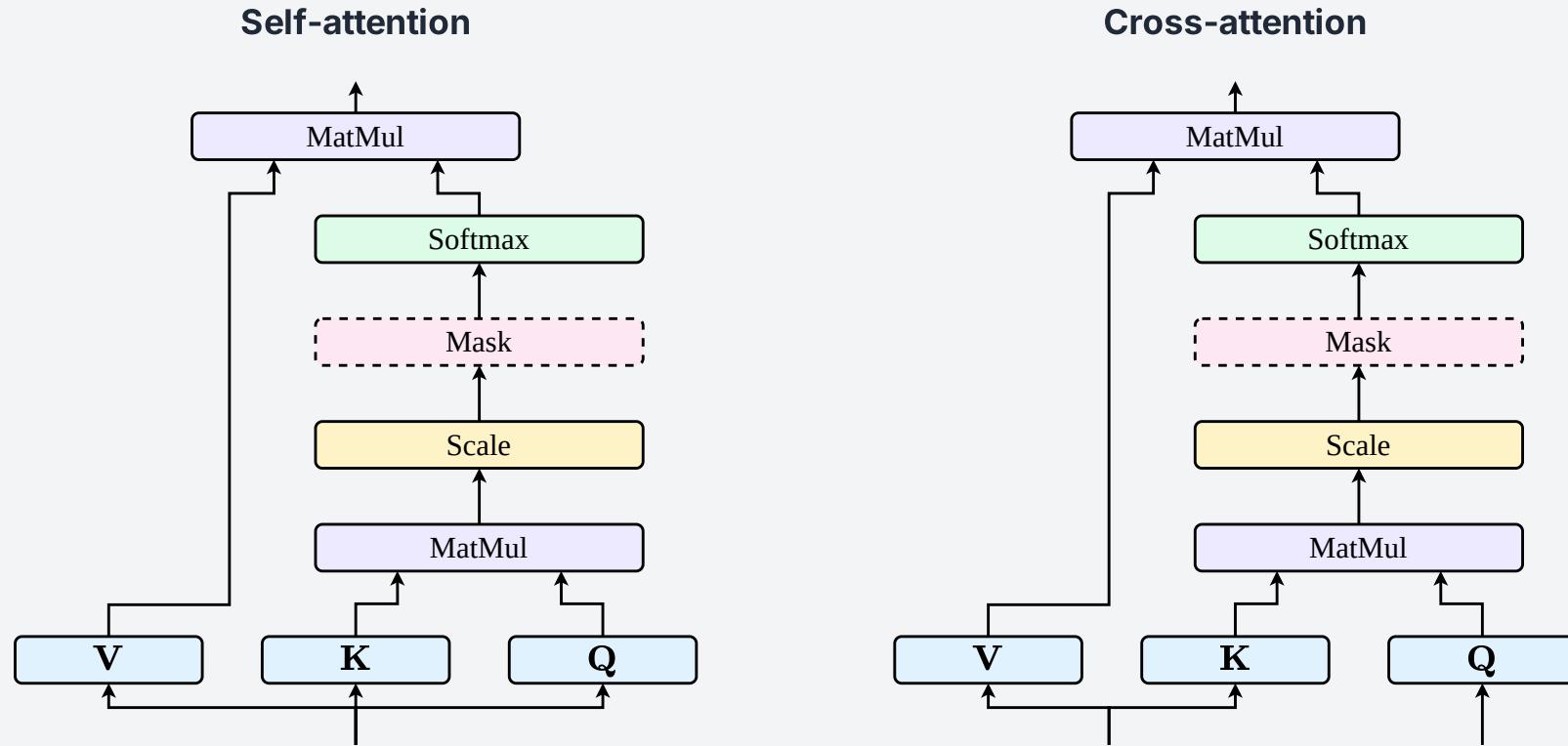
Multi-head Attention

- Where do \mathbf{V} , \mathbf{K} , \mathbf{Q} come from?



Attention in Transformers

Multi-head Attention

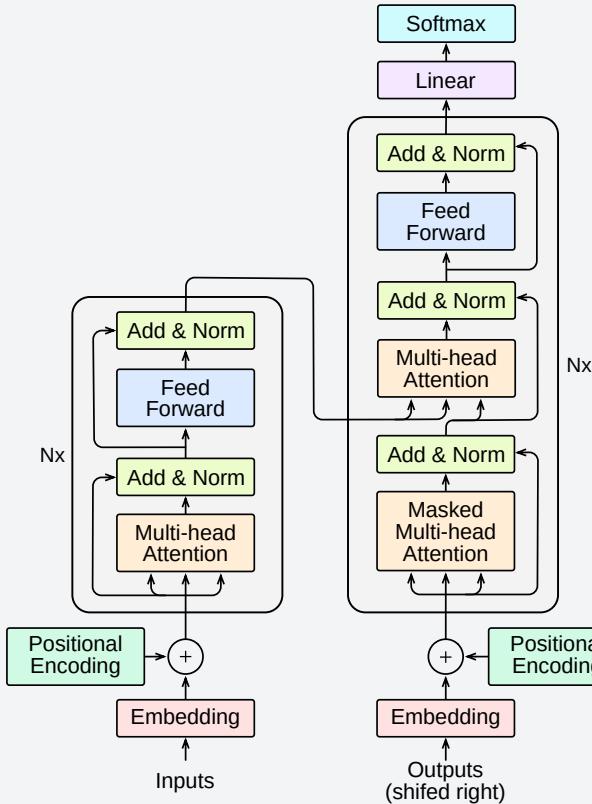


Attention in Transformers

Multi-head Attention

Where do \mathbf{V} , \mathbf{K} , \mathbf{Q} come from?

- **Self-attention:** everything comes from the same sequence
- **Cross-attention:** \mathbf{V} , \mathbf{K} come from source sequence, \mathbf{Q} comes from target sequence



Attention in Transformers

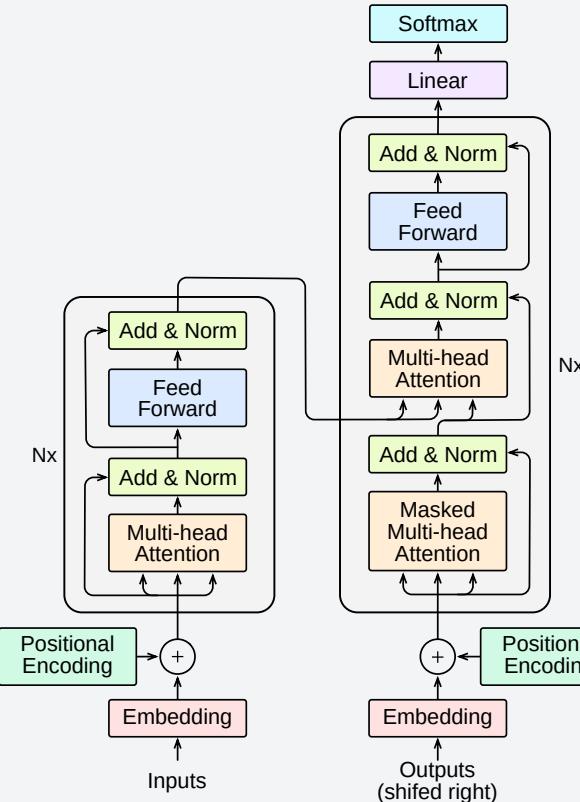
Multi-head Attention

Where do \mathbf{V} , \mathbf{K} , \mathbf{Q} come from?

- **Self-attention:** everything comes from the same sequence
- **Cross-attention:** \mathbf{V} , \mathbf{K} come from source sequence, \mathbf{Q} comes from target sequence

All components constructed from a projection of the token embeddings

1. $\mathbf{V} = \mathbf{XW}_V$
2. $\mathbf{K} = \mathbf{XW}_K$
3. $\mathbf{Q} = \underbrace{\mathbf{XW}_Q}_{\text{Self-attention}} \text{ or } \underbrace{\mathbf{YW}_Q}_{\text{Cross-attention}}$



Multi-head Attention

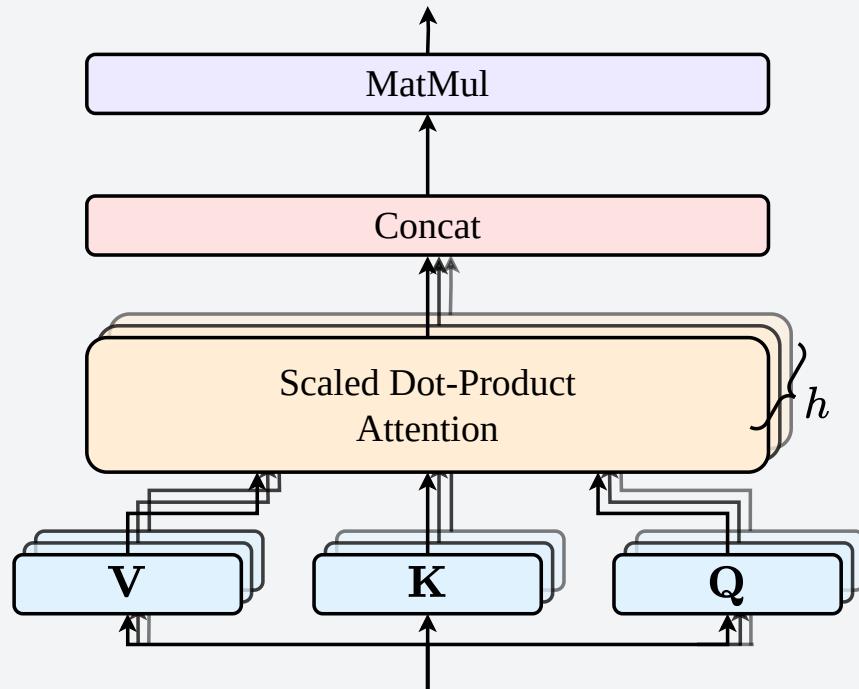
Multi-head Attention

- Currently we use 1 set of attention weights
 - Can only process 1 query type

Multi-head Attention

Multi-head Attention

- Currently we use 1 set of attention weights
 - Can only process 1 query type
- With h attention heads, we learn h concepts
 - To reduce cost, reduce dimensionality $d_{K,V}/h$

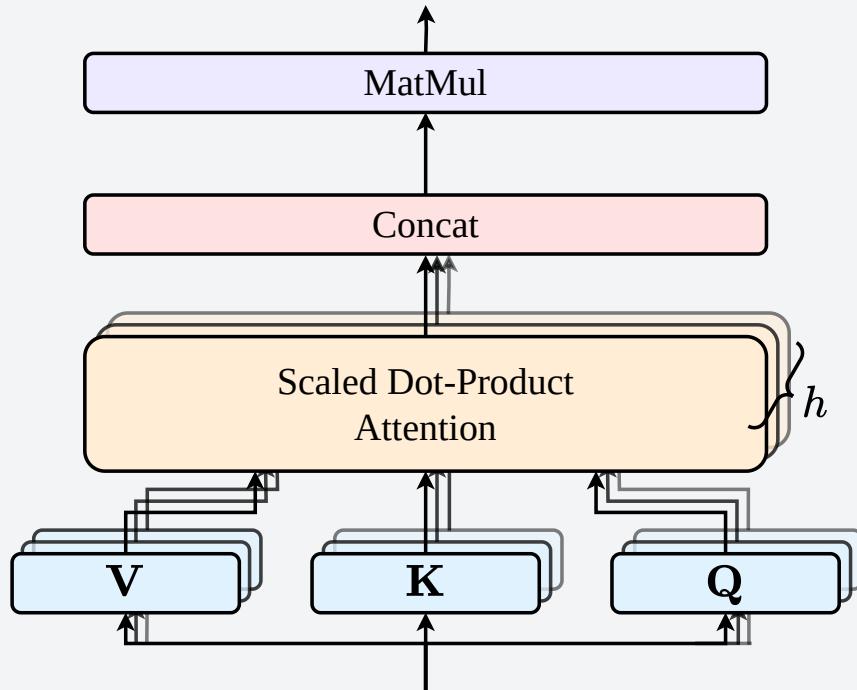


Multi-head Attention

Multi-head Attention

- Currently we use 1 set of attention weights
 - Can only process 1 query type
- With h attention heads, we learn h concepts
 - To reduce cost, reduce dimensionality
 $d_{K,V}/h$

```
self.attention_heads = [  
    AttentionHead(d=self.d // self.h) for i in range(self.h)  
]  
  
self.mha_proj = nn.Linear(self.d, self.d)  
  
mha = torch.concat([  
    attention_heads[i](x) for i in range(self.h)  
])  
  
out = self.mha_proj(mha)
```

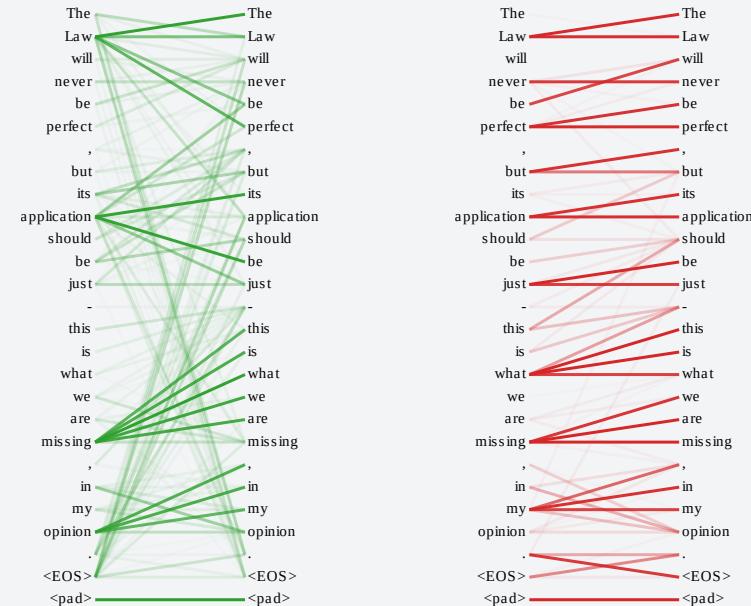


Multi-head Attention

Multi-head Attention

Multi-head attention allows the model to jointly attend to information from different representation subspaces at different positions. [One] attention head, averaging inhibits this.

Vaswani et al. (2017). Attention is all you need. Advances in neural information processing systems, 30. (p. 5 & 1



Multi-head Attention

Multi-head Attention

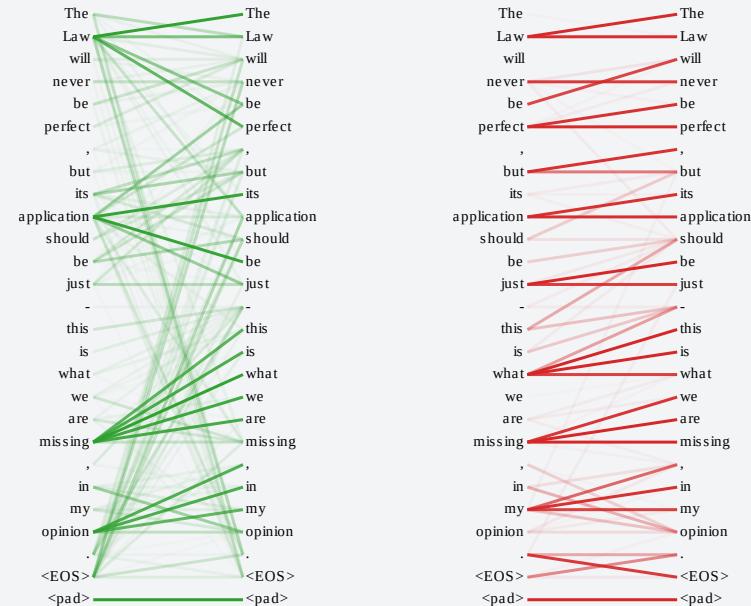
Multi-head attention allows the model to jointly attend to information from different representation subspaces at different positions. [One] attention head, averaging inhibits this.

Vaswani et al. (2017). Attention is all you need. Advances in neural information processing systems, 30. (p. 5 & 1)

Multiple heads, multiple different queries
processed in parallel

- Positional heads
- Syntactic heads
- Rare words?

Voita et al. (2019). Analyzing Multi-Head Self-Attention: Specialized Heads Do the Heavy Lifting, the Rest Can Be Pruned. Association for Computational Linguistics.



Multi-head Attention

Multi-head Attention

Do different heads attend to different concepts?

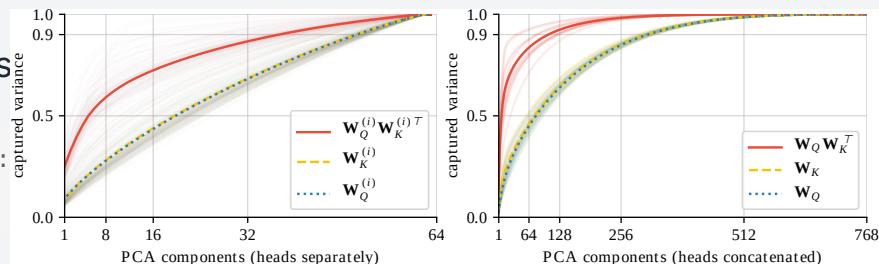
Multi-head Attention

Multi-head Attention

Do different heads attend to different concepts?

- Individual heads = high rank, concatenated heads = low rank

Cordonnier, Loukas & Jaggi (2020). Multi-head attention:
Collaborate instead of concatenate. arXiv:2006.16362.



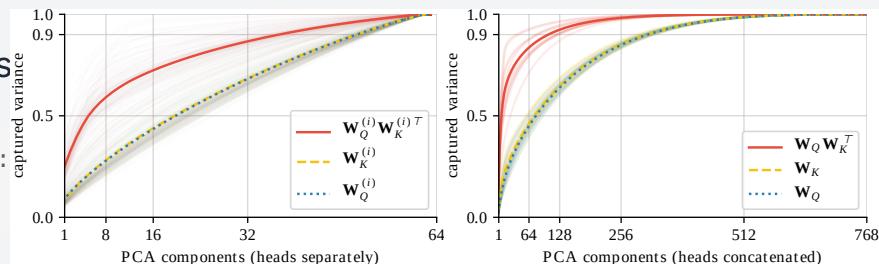
Multi-head Attention

Multi-head Attention

Do different heads attend to different concepts?

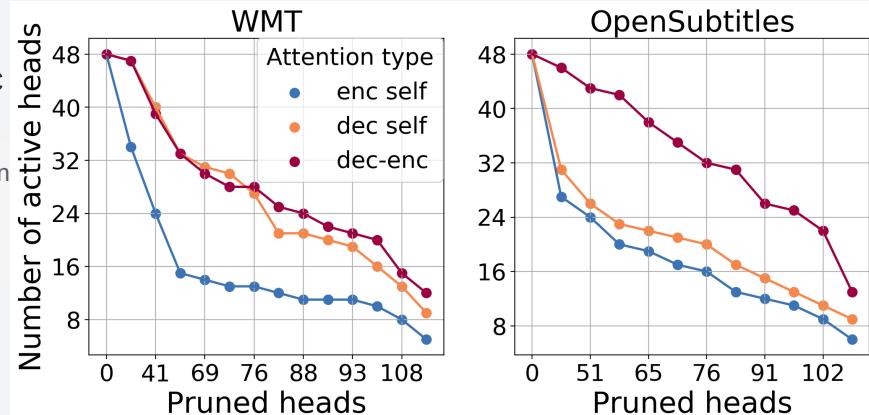
- Individual heads = high rank, concatenated heads = low rank

Cordonnier, Loukas & Jaggi (2020). Multi-head attention: Collaborate instead of concatenate. arXiv:2006.16362.



- Most heads can be pruned away
- Enc-Dec heads are more important than Enc-Enc heads

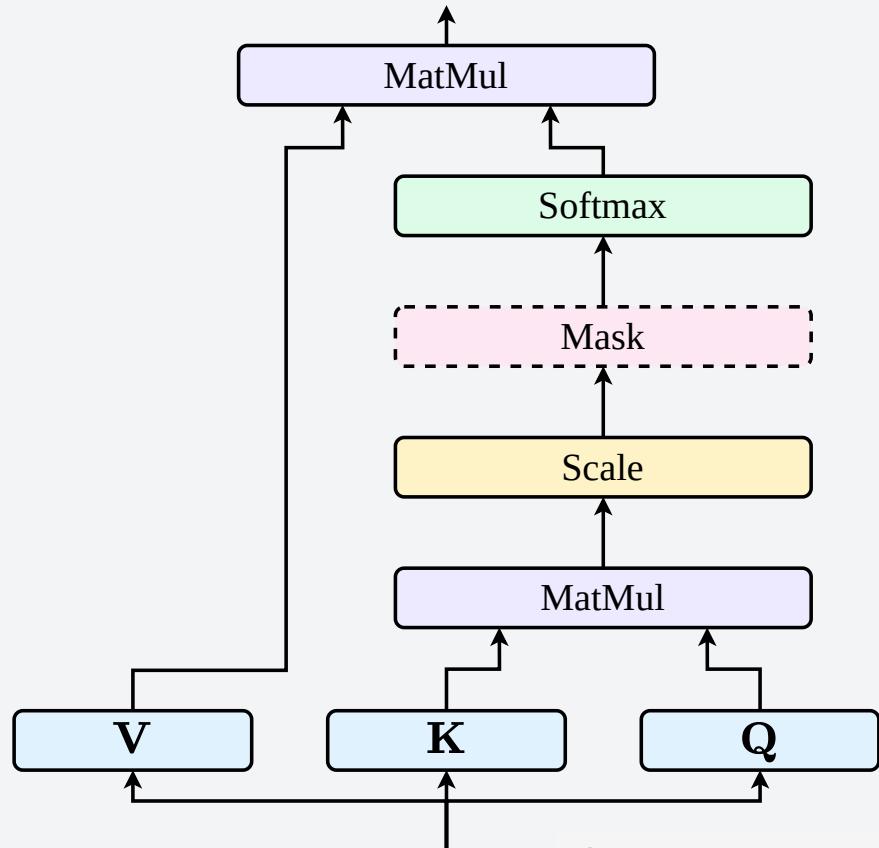
Voita et al. (2019). Analyzing Multi-Head Self-Attention Specialized Heads Do the Heavy Lifting, the Rest Can Be Pruned. Association for Computational Linguistics.



Why Attention?

Multi-head Attention

Summary

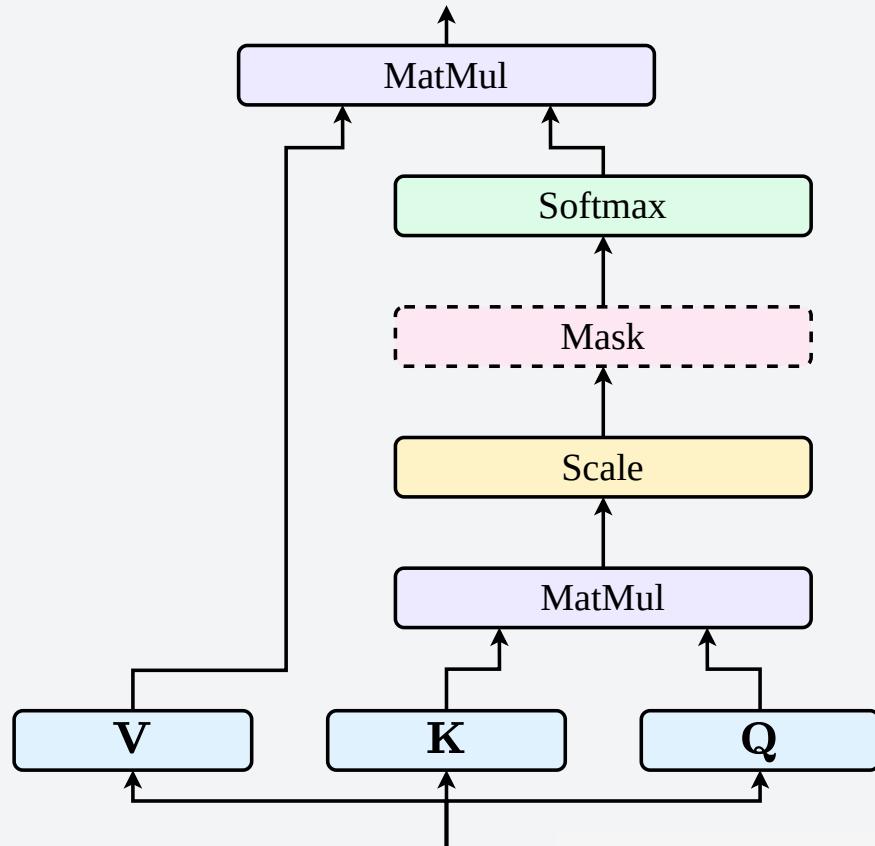


Why Attention?

Multi-head Attention

Summary

1. Attention is a **linear map \mathbf{AV}** where \mathbf{A} is dynamically constructed from $f(\mathbf{Q}, \mathbf{K})$

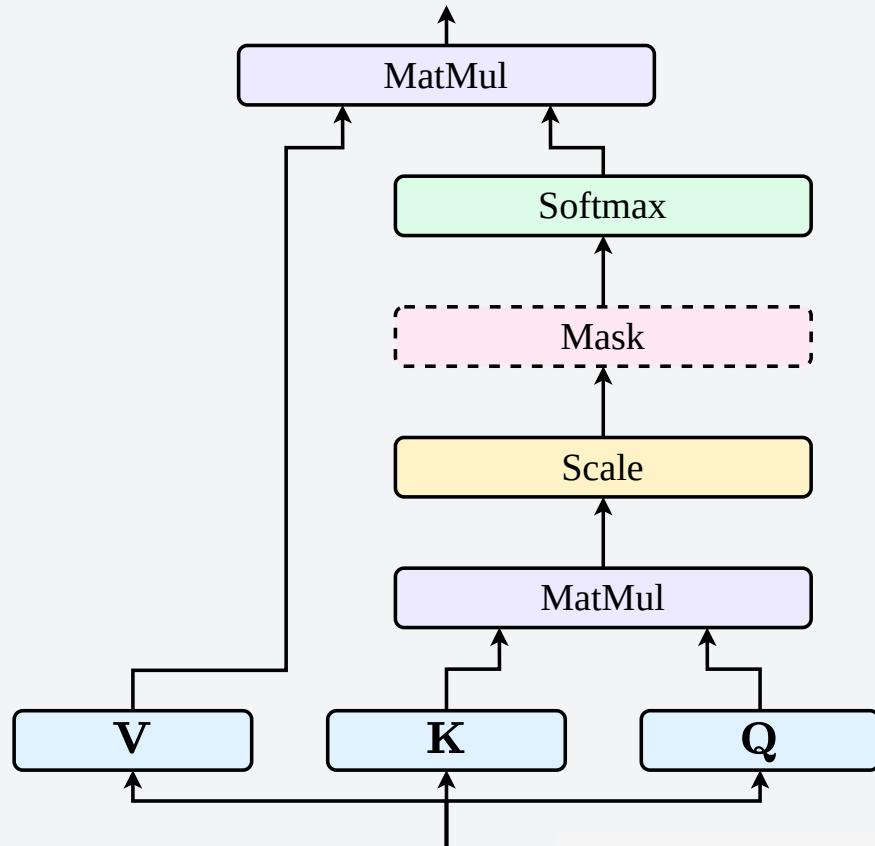


Why Attention?

Multi-head Attention

Summary

1. Attention is a **linear map \mathbf{AV}** where \mathbf{A} is dynamically constructed from $f(\mathbf{Q}, \mathbf{K})$
2. The values of \mathbf{A} are all in $(0, 1)$, making \mathbf{AV} a convex combination/weighted mean of \mathbf{V}

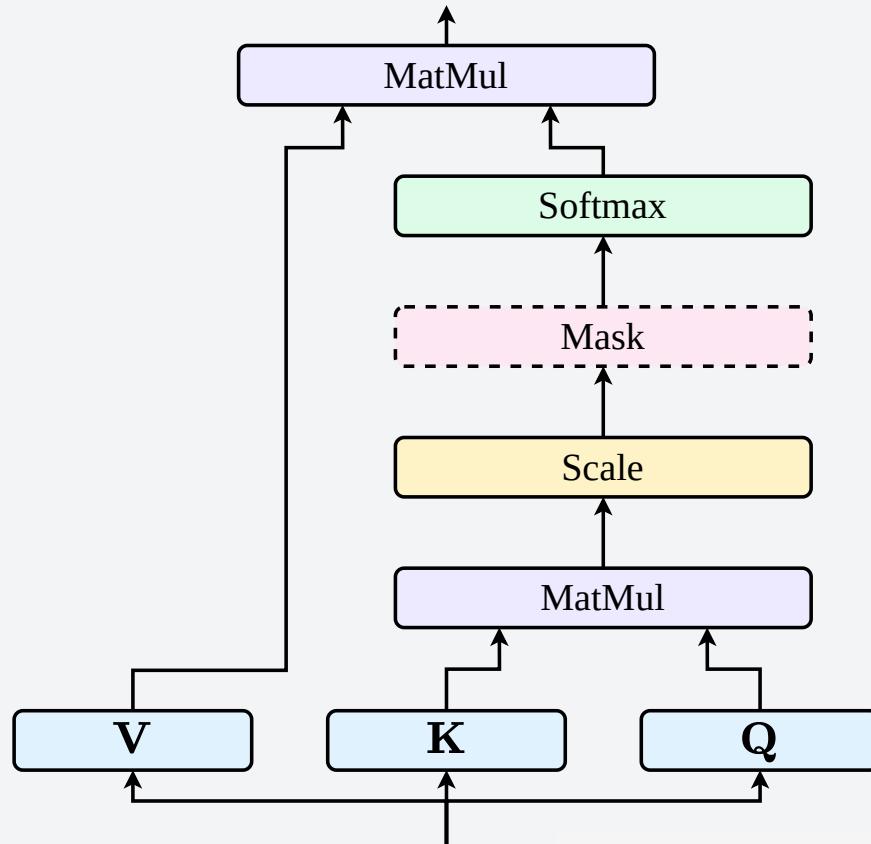


Why Attention?

Multi-head Attention

Summary

1. Attention is a **linear map \mathbf{AV}** where \mathbf{A} is dynamically constructed from $f(\mathbf{Q}, \mathbf{K})$
2. The values of \mathbf{A} are all in $(0, 1)$, making \mathbf{AV} a convex combination/weighted mean of \mathbf{V}
3. Attention does not understand word order

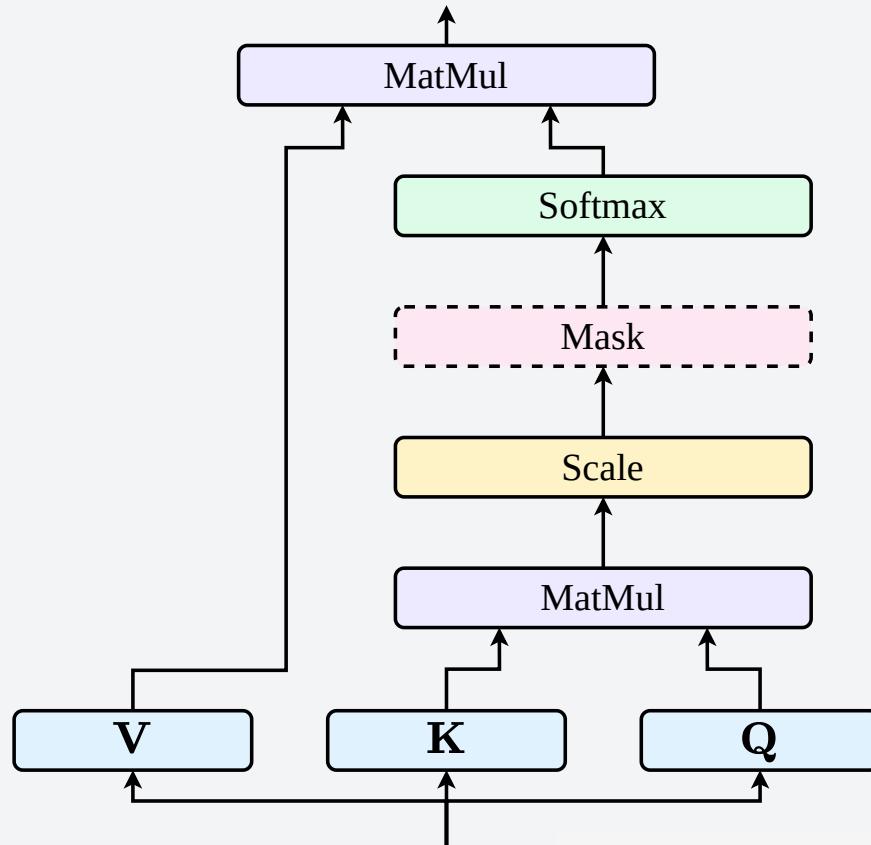


Why Attention?

Multi-head Attention

Summary

1. Attention is a **linear map \mathbf{AV}** where \mathbf{A} is dynamically constructed from $f(\mathbf{Q}, \mathbf{K})$
2. The values of \mathbf{A} are all in $(0, 1)$, making \mathbf{AV} a convex combination/weighted mean of \mathbf{V}
3. Attention does not understand word order
4. In Transformers, Attention is used to :
 - add context from self (self-attention)
 - add context from others (cross-attention)

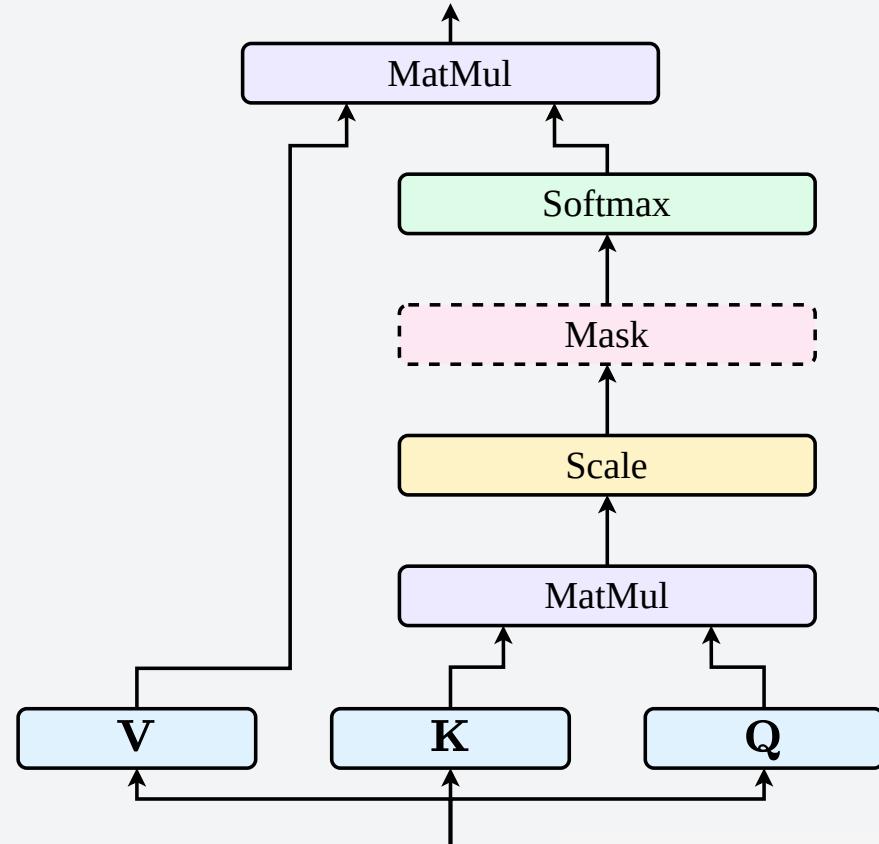


Why Attention?

Multi-head Attention

Summary

1. Attention is a **linear map \mathbf{AV}** where \mathbf{A} is dynamically constructed from $f(\mathbf{Q}, \mathbf{K})$
2. The values of \mathbf{A} are all in $(0, 1)$, making \mathbf{AV} a convex combination/weighted mean of \mathbf{V}
3. Attention does not understand word order
4. In Transformers, Attention is used to :
 - add context from self (self-attention)
 - add context from others (cross-attention)
5. Attention cost scales quadratically with sequence length

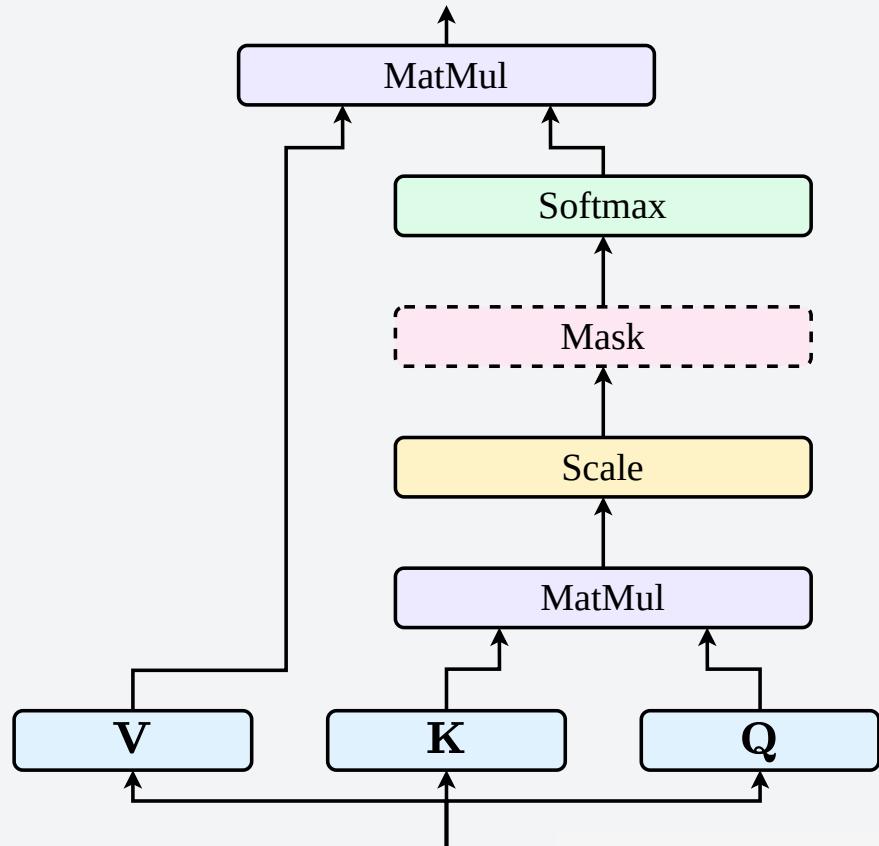


Why Attention?

Multi-head Attention

- Transformer attention between two sequences, \mathbf{X} and \mathbf{Y} has a computational cost of (excluding projections):

$$\mathcal{O} \left(\underbrace{t_x \cdot t_y \cdot d_k}_{\text{MatMul 1}} + \underbrace{t_x \cdot t_y \cdot d_v}_{\text{MatMul 2}} \right)$$



Why Attention?

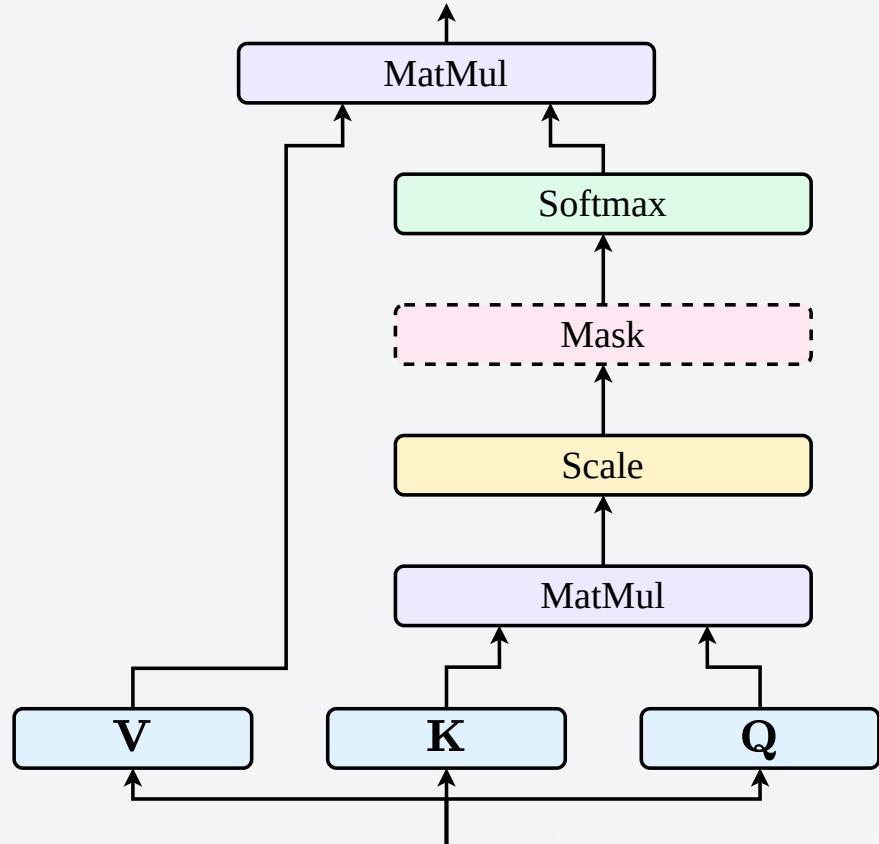
Multi-head Attention

- Transformer attention between two sequences, \mathbf{X} and \mathbf{Y} has a computational cost of (excluding projections):

$$\mathcal{O} \left(\underbrace{t_x \cdot t_y \cdot d_k}_{\text{MatMul 1}} + \underbrace{t_x \cdot t_y \cdot d_v}_{\text{MatMul 2}} \right)$$

- But RNNs have linear time complexity...

$$\mathcal{O} \left(\underbrace{t_x \cdot d_x \cdot d_h}_{x_t} + \underbrace{t_x \cdot d_h^2}_{h_{t-1}} \right)$$

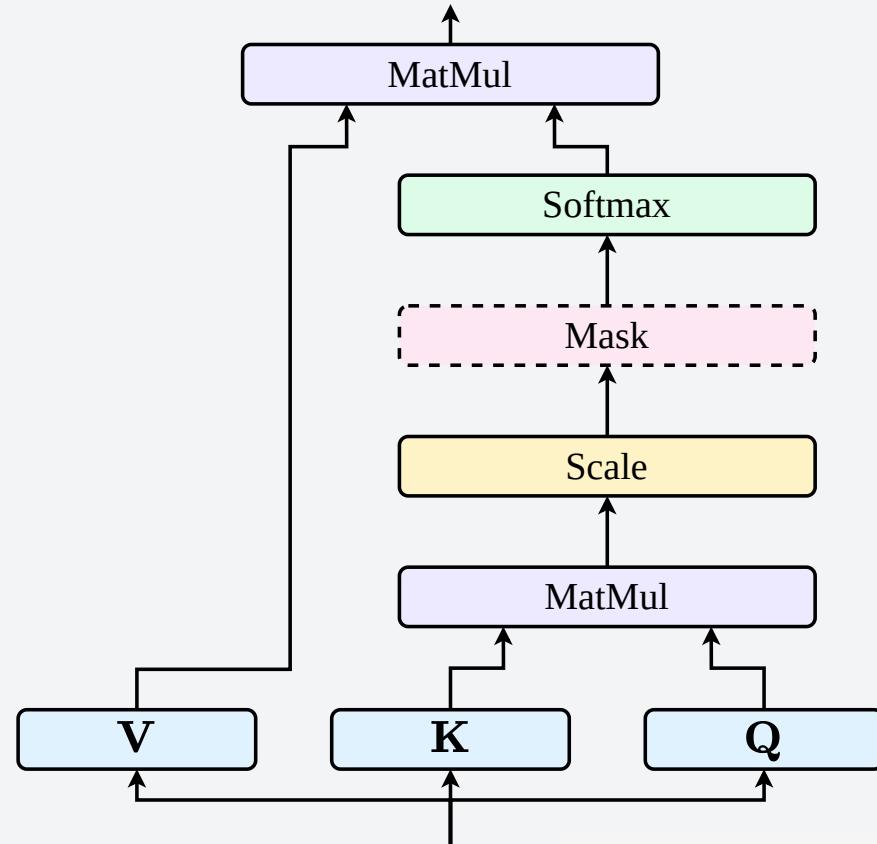


Why Attention?

Multi-head Attention

Summary

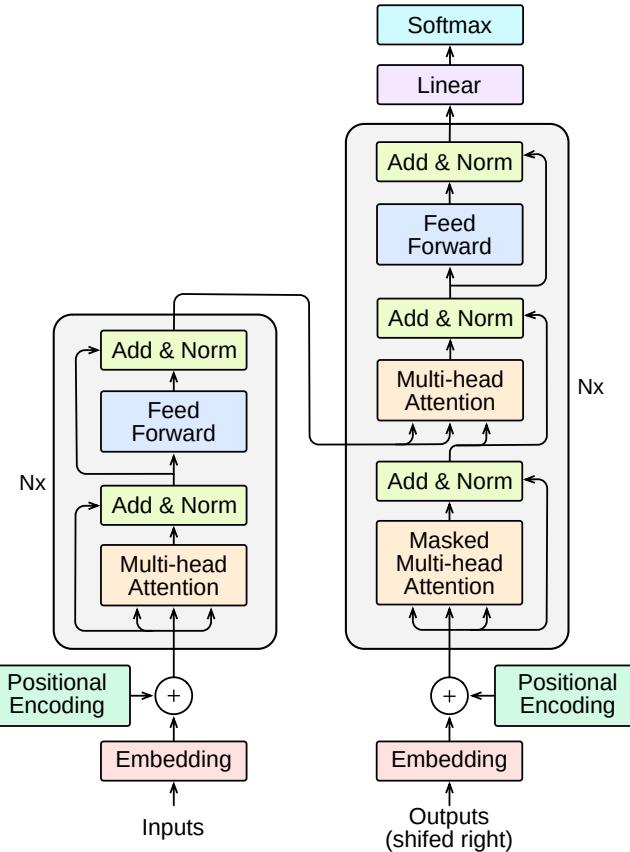
1. Attention is a **linear map \mathbf{AV}** where \mathbf{A} is dynamically constructed from $f(\mathbf{Q}, \mathbf{K})$
2. The values of \mathbf{A} are all in $(0, 1)$, making \mathbf{AV} a convex combination/weighted mean of \mathbf{V}
3. Attention does not understand word order
4. In Transformers, Attention is used to:
 - add context from self (self-attention)
 - add context from others (cross-attention)
5. Attention cost scales quadratically with sequence length
6. **Attention is parallelizable**





Jakob Uszkoreit (August 31, 2017). Transformer: A Novel Neural Network Architecture for Language Understanding.
<https://research.google/blog/transformer-a-novel-neural-network-architecture-for-language-understanding/>

Feed Forward



Non-linearities

Feed Forward

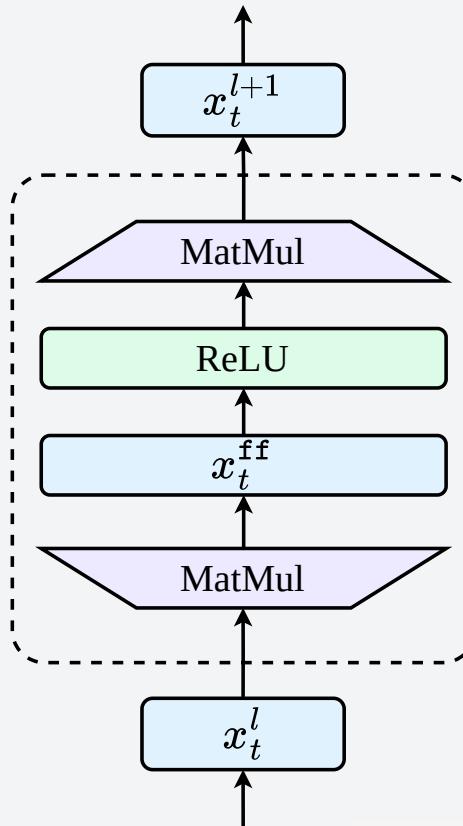
- Attention is just a fancy linear map*
 - Neural nets need non-linear operations

* More or less

Non-linearities

Feed Forward

- Attention is just a fancy linear map*
 - Neural nets need non-linear operations
- Add pointwise feed-forward nets for non-linear expressiveness
 - Just an MLP applied to each time step

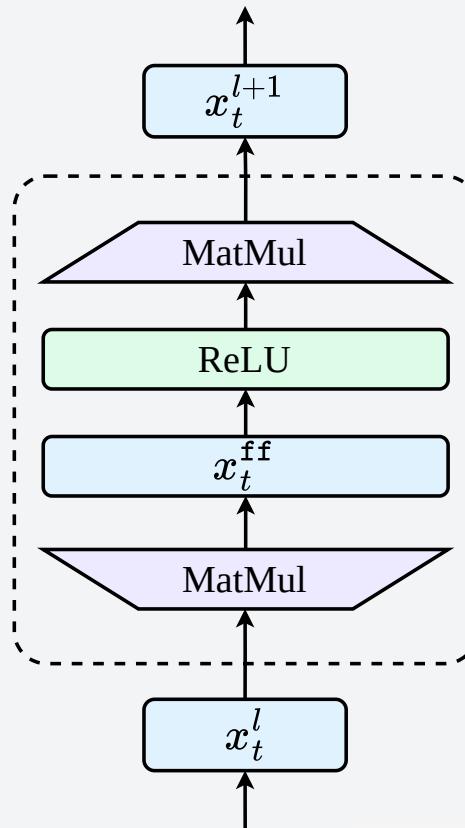


* More or less

Non-linearities

Feed Forward

- Attention is just a fancy linear map*
 - Neural nets need non-linear operations
- Add pointwise feed-forward nets for non-linear expressiveness
 - Just an MLP applied to each time step
- Assume $d_{\text{ff}} \gg d_V$
 - Depth is serial
 - Width is parallel



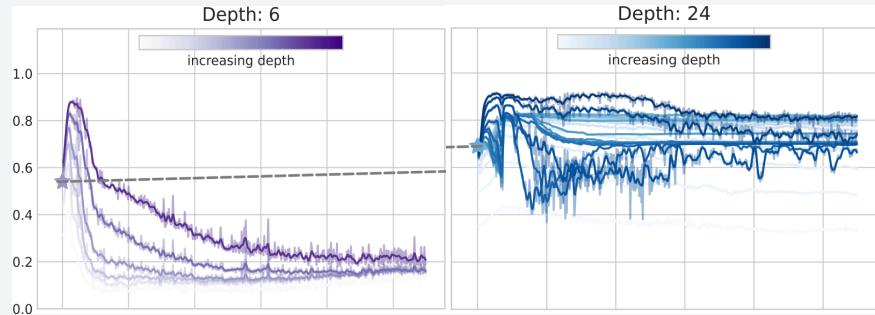
* More or less

Necessity of Feed Forward Nets

Feed Forward

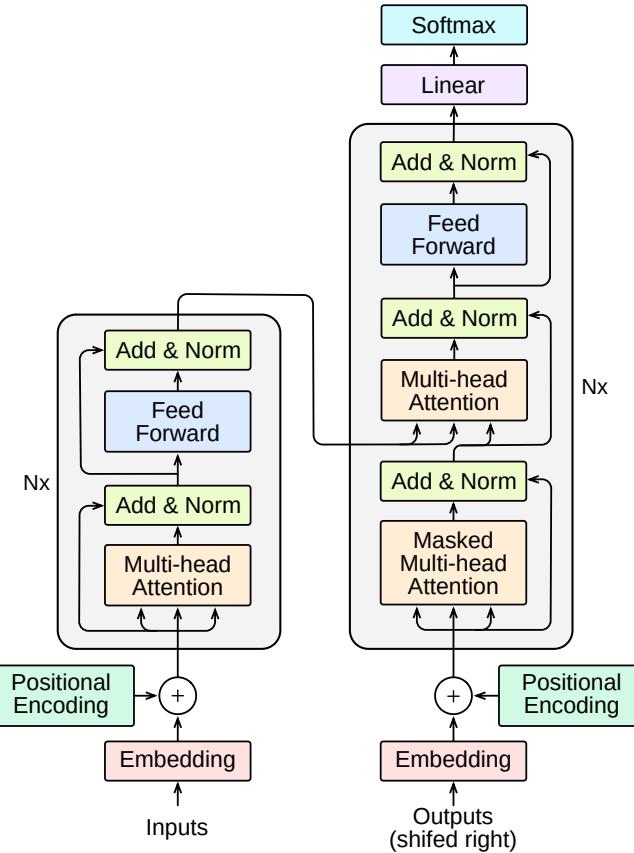
- Is attention all you need?
- Without feed forward nets, transformer token representations collapse
 - Occurs *doubly exponential* with depth (very, very fast)

Dong, Cordonnier, & Loukas (2021). Attention is not all you need: Pure attention loses rank doubly exponentially with depth. PMLR



Noci et al. (2022). Signal Propagation in Transformers: Theoretical Perspectives and the Role of Rank Collapse. arXiv:2206.03126.

Add & Norm



Residual Connections

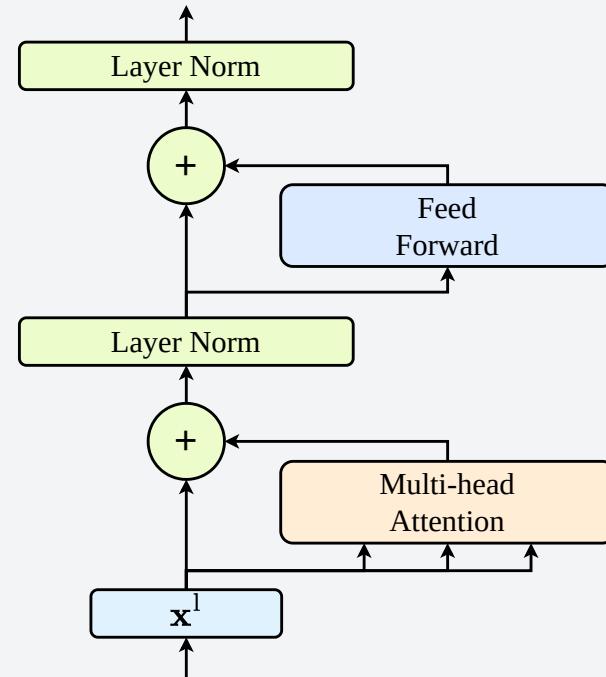
Add & Norm

- After each operation, the input is added back in

$$x^{l+1} = \text{SubLayer}(x^l) + x^l$$

- Makes deep architectures optimisable

He et al. (2015). Deep residual learning for image recognition. arXiv:1512.03385



Layer Norm

Add & Norm

- After each operation and residual connection, normalize

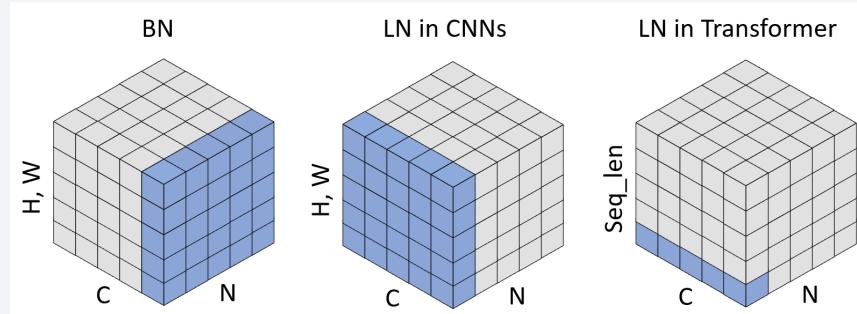
$$x^{l+1} = \text{LayerNorm}(\text{SubLayer}(x^l) + x^l)$$

- LayerNorm in Transformers is per token
- Makes neural nets converge faster

Ba, Kiros & Hinton (2016). Layer normalization.
arXiv:1607.06450

- Modern transformers prefer **LayerNorm both** before and after operations

Shleifer, Weston, & Ott (2021). Normformer: Improved transformer pretraining with extra normalization.
arXiv:2110.09456.



Yao et al. (2021). Leveraging batch normalization for vision transformers. In Proceedings of the IEEE/CVF International Conference on Computer Vision (pp. 413-422).

Embeddings

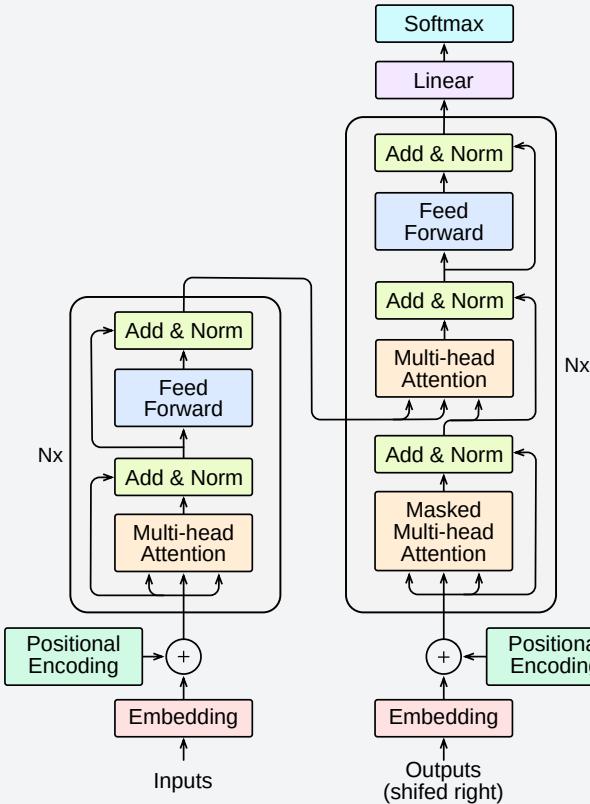
From words to vectors and back

Embedding Types

Embeddings

Transformers apply two embeddings:

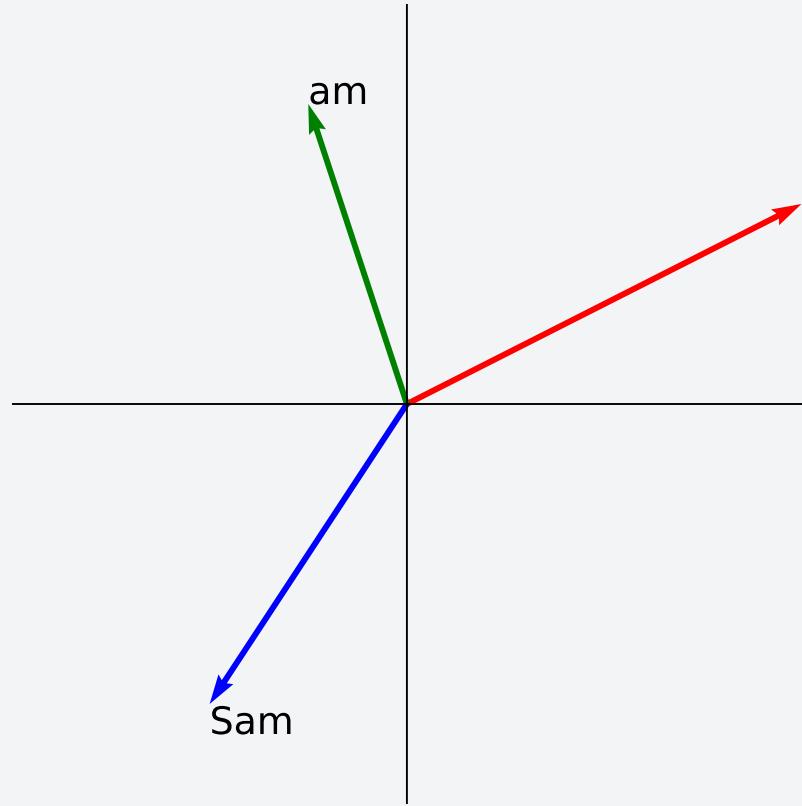
1. **Token Embedding** transform token IDs into vectors
2. **Position Encodings** add information about token location



Token to vector

Token Embedding

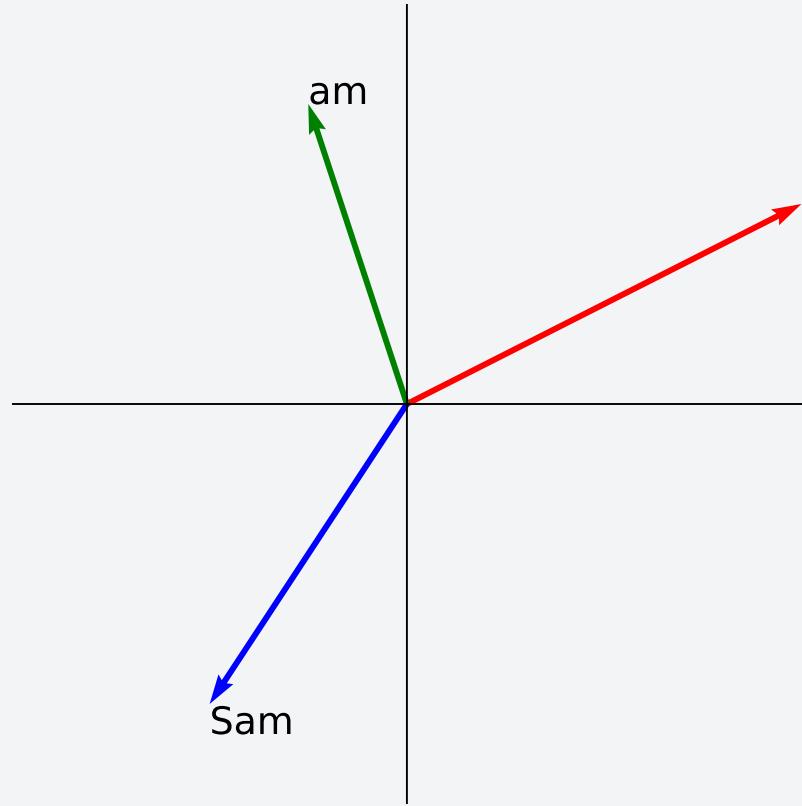
- Convert strings to dense representations
 - Typically done using dictionary lookup



Token to vector

Token Embedding

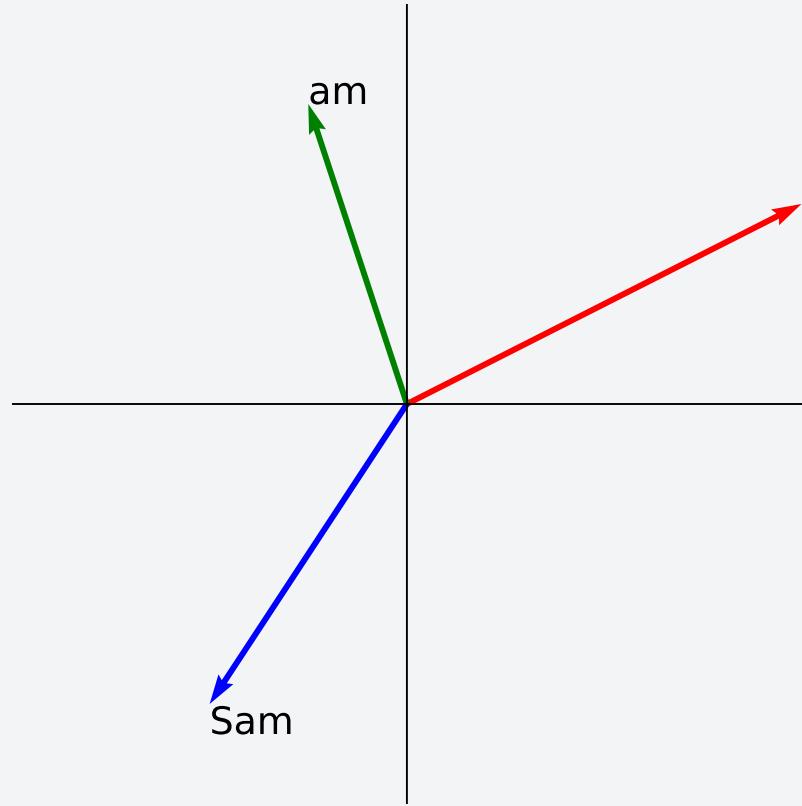
- Convert strings to dense representations
 - Typically done using dictionary lookup
- Memory cost: $\mathcal{O}(|\mathcal{V}| \cdot d_v)$
 - Usually most expensive operation
 - Can save cost by have embedding matrices share weights



Token to vector

Token Embedding

- Convert strings to dense representations
 - Typically done using dictionary lookup
- Memory cost: $\mathcal{O}(|\mathcal{V}| \cdot d_v)$
 - Usually most expensive operation
 - Can save cost by have embedding matrices share weights
- Typically *not* pre-trained



Permutation Equivariance Revisited

Positional Encoding

Permutation Equivariance

Attention is permutation *equivariant* to changes in the word order of its input

- From NLP1: word order is pretty important for modelling language...

'Only I love you'

=

'I only love you'

=

'I love only you'

=

'I love only you, only...'

Permutation Equivariance Revisited

Positional Encoding

Permutation Equivariance

Attention is permutation *equivariant* to changes in the word order of its input

- From NLP1: word order is pretty important for modelling language...
- Solution: add token position information to token embedding

`embed(string[t]) + position(t)`

'Only I love you'

=

'I only love you'

=

'I love only you'

=

'I love only you, only...'

Learning Position Embeddings

Positional Encoding

Proposal 1: learn per-position offsets manually

Learning Position Embeddings

Positional Encoding

Proposal 1: learn per-position offsets manually

1. Deterministic [✓]
2. Distinct at all time-steps [?]

Learning Position Embeddings

Positional Encoding

Proposal 1: learn per-position offsets manually

1. Deterministic [✓]
2. Distinct at all time-steps [?]
3. Extends to different sequence lengths [✗]
4. Encodes relative positions [?]

Learning Position Embeddings

Positional Encoding

Proposal 1: learn per-position offsets manually

1. Deterministic [✓]
2. Distinct at all time-steps [?]
3. Extends to different sequence lengths [✗]
4. Encodes relative positions [?]

Not a terrible idea

Model	PPL (dev)	BLEU (test)	Params $\times 10^6$
Base	4.92	25.8	65
Learned Positional Encoding	4.92	25.7	65

Bits as Position Encoding

Positional Encoding

Proposal 2: add the offsets in bit representations to the embeddings

0	0	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1
1	0	0	1	1	0	0	1	1	0	0	1	1	0	0	1	1	0	0
2	0	0	0	0	1	1	1	1	0	0	0	0	1	1	1	1	1	1
3	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1

<S> The best grad i ent desc ent method is grad stud ent desc ent </S>

Bits as Position Encoding

Positional Encoding

Proposal 2: add the offsets in bit representations to the embeddings

1. Deterministic [✓]
2. Distinct at all time-steps [✓]
3. Extends to different sequence lengths [✓]*
4. Encodes relative positions [✗]

0	0	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1
1	0	0	1	1	0	0	1	1	0	0	1	1	0	0	1	1	0	0
2	0	0	0	0	1	1	1	1	0	0	0	0	1	1	1	1	1	1
3	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1

</s> The best grad i ent desc ent method is grad stud ent desc ent </s>

* Maximum sequence length is now $2^{d_{\text{model}}}$

Bits as Position Encoding

Positional Encoding

Proposal 2: add the offsets in bit representations to the embeddings

1. Deterministic [✓]
2. Distinct at all time-steps [✓]
3. Extends to different sequence lengths [✓]*
4. Encodes relative positions [✗]
5. Elegant [✗]

0	0	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1
1	0	0	1	1	0	0	1	1	0	0	1	1	0	0	1	1	0	0
2	0	0	0	0	1	1	1	1	0	0	0	0	1	1	1	1	1	1
3	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1

</s> The best grad i ent desc ent method is grad stud ent desc ent </s>

* Maximum sequence length is now $2^{d_{\text{model}}}$

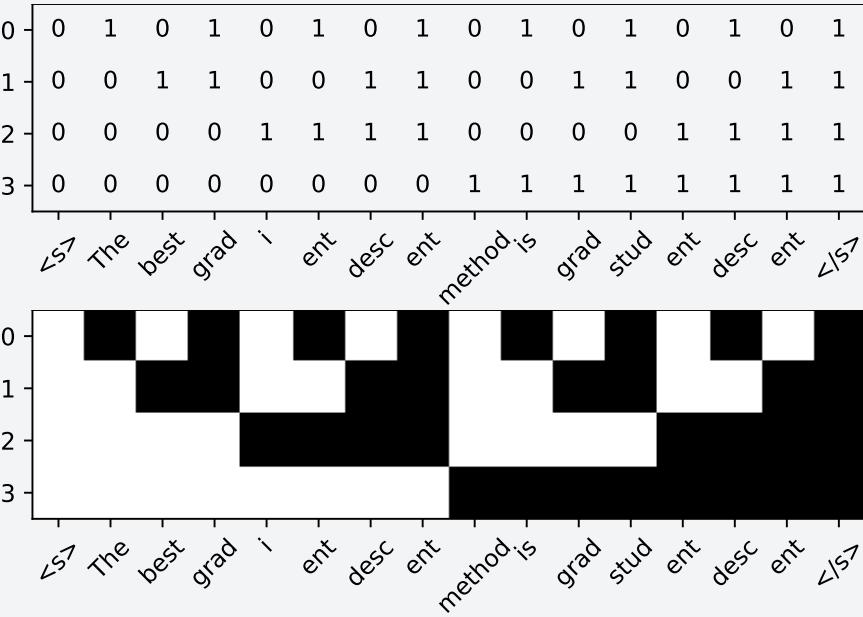
Bits as Position Encoding

Positional Encoding

Proposal 2: add the offsets in bit representations to the embeddings

1. Deterministic [✓]
2. Distinct at all time-steps [✓]
3. Extends to different sequence lengths [✓]*
4. Encodes relative positions [✗]
5. Elegant [✗]

* Maximum sequence length is now $2^{d_{\text{model}}}$



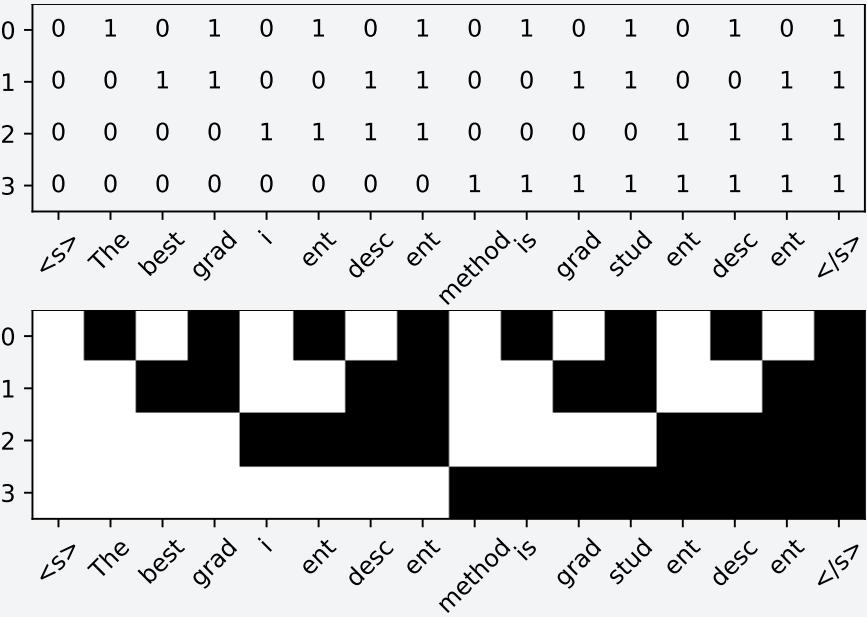
Bits as Position Encoding

Positional Encoding

Proposal 2: add the offsets in bit representations to the embeddings

1. Deterministic [✓]
2. Distinct at all time-steps [✓]
3. Extends to different sequence lengths [✓]*
4. Encodes relative positions [✗]
5. Elegant [✗]

* Maximum sequence length is now $2^{d_{\text{model}}}$



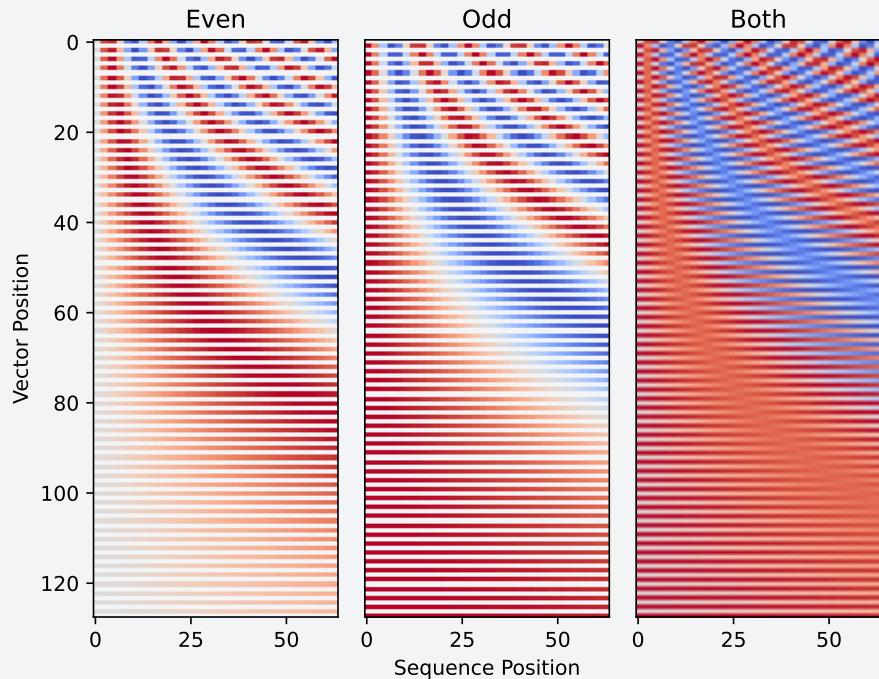
Sinusoidal Position Encoding

Positional Encoding

Proposal 3: use sinusoidal position encoding

$$\text{Position}(t) = \begin{cases} \sin(\omega_i \cdot t), & \text{mod } (i, 2) = 0 \\ \cos(\omega_i \cdot t), & \text{mod } (i, 2) \neq 0 \end{cases}$$
$$\omega_i = k^{-2i/d_{\text{model}}}$$

- First dimensions get very quickly oscillating sinusoid
- Last dimensions get very slow oscillating sinusoid



Sinusoidal Position Encoding

Positional Encoding

Relative positions can now be represented as a linear mapping, dependent on the offset

Neural nets can learn this offset

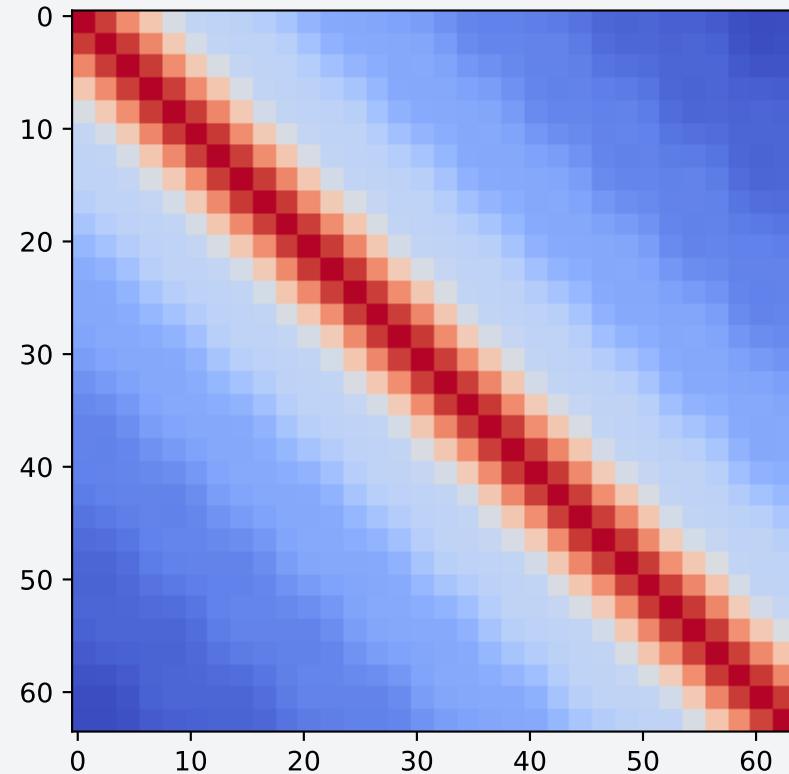
$$\begin{aligned}\text{Position}(t + \Delta t) &= \begin{bmatrix} \sin(\omega_i \cdot (t + \Delta t)) \\ \cos(\omega_i \cdot (t + \Delta t)) \end{bmatrix} \\ &= \begin{bmatrix} \cos \Delta t & \sin \Delta t \\ -\sin \Delta t & \cos \Delta t \end{bmatrix} \begin{bmatrix} \sin(\omega_i \cdot t) \\ \cos(\omega_i \cdot t) \end{bmatrix} \\ &= \begin{bmatrix} \cos \Delta t & \sin \Delta t \\ -\sin \Delta t & \cos \Delta t \end{bmatrix} \text{Position}(t) \\ &= \mathbf{A}_{\Delta t} \text{Position}(t)\end{aligned}$$

Sinusoidal Position Encoding

Positional Encoding

Cosine distance between time-steps decreases with offset magnitude

For proofs & more, see here

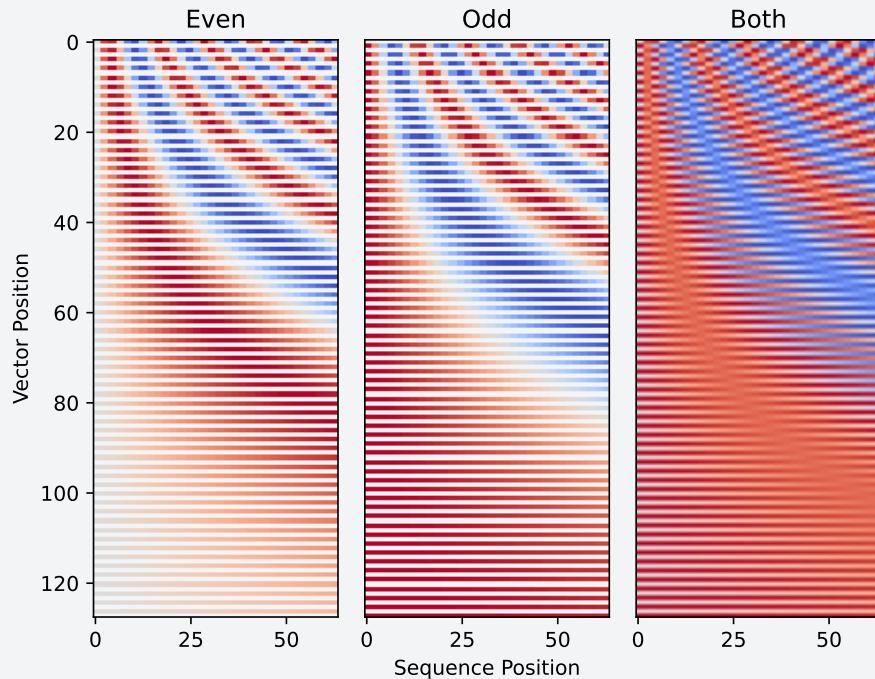


Sinusoidal Position Encoding

Positional Encoding

Proposal 3: use sinusoidal position encoding

1. Deterministic [✓]
2. Distinct at all time-steps [✓]
3. Extends to different sequence lengths [✓]
4. Encodes relative positions [✓]
5. Elegant [✓]

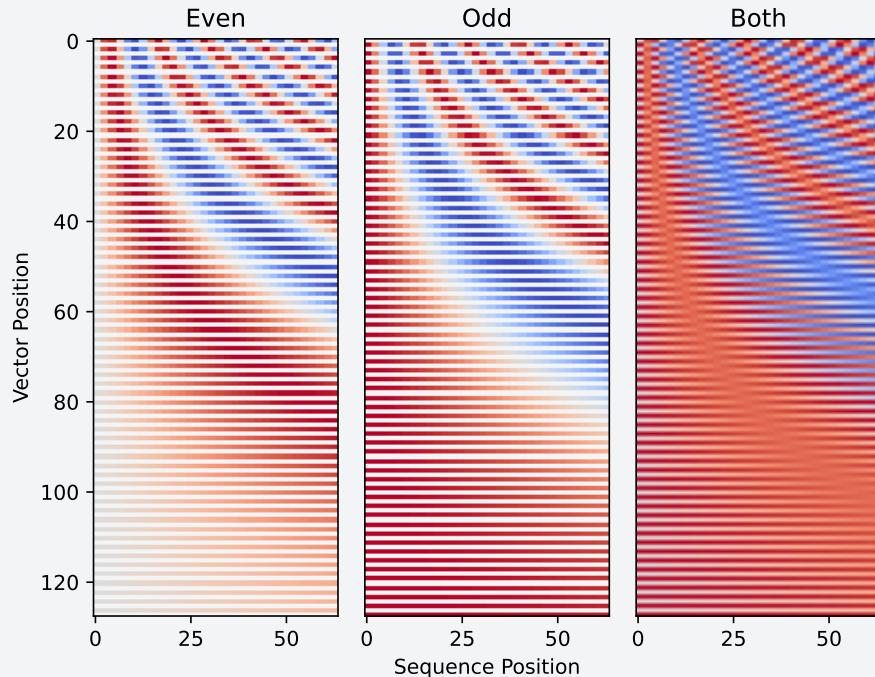


Sinusoidal Position Encoding

Positional Encoding

Proposal 3: use sinusoidal position encoding

1. Deterministic [✓]
2. Distinct at all time-steps [✓]
3. Extends to different sequence lengths [✓]
4. Encodes relative positions [✓]
5. Elegant [✓]
6. Optimal [???
7. Necessary [???



See this HF blog for more good ideas

Tokenization

Vocabulary Size

Tokenization

- Base transformer achieved SoTA translation with only 65M parameters

Name	N	d_{model}	d_{ff}	h	d_k	d_v	Params $\times 10^6$
base	6	512	2048	8	64	64	65
big	6	1024	4096	16	64	64	213

Vocabulary Size

Tokenization

- Base transformer achieved SoTA translation with only 65M parameters
- English has $|\mathcal{V}_{\text{src}}| \approx 1M$ words

Name	N	d_{model}	d_{ff}	h	d_k	d_v	Params $\times 10^6$
base	6	512	2048	8	64	64	65
big	6	1024	4096	16	64	64	213

Vocabulary Size

Tokenization

- Base transformer achieved SoTA translation with only 65M parameters
- English has $|\mathcal{V}_{\text{src}}| \approx 1M$ words
- Assume target language has $|\mathcal{V}_{\text{tgt}}| \approx 1M^*$

Name	N	d_{model}	d_{ff}	h	d_k	d_v	Params $\times 10^6$
base	6	512	2048	8	64	64	65
big	6	1024	4096	16	64	64	213

* Not realistic

Vocabulary Size

Tokenization

- Base transformer achieved SoTA translation with only 65M parameters
- English has $|\mathcal{V}_{\text{src}}| \approx 1M$ words
- Assume target language has $|\mathcal{V}_{\text{tgt}}| \approx 1M^*$
- Cost of word embedding and unembedding **alone** is

$$3 \times 1 \cdot 10^6 \times 512 \approx 1.5 \cdot 10^9$$

parameters...

* Not realistic

Name	N	d_{model}	d_{ff}	h	d_k	d_v	Params $\times 10^6$
base	6	512	2048	8	64	64	65
big	6	1024	4096	16	64	64	213

Vocabulary Size

Tokenization

- **Problem:** embedding cost scales with $|\mathcal{V}|$
 - Most words occur infrequently
 - Some unknown words are likely to exist

Vocabulary Size

Tokenization

- **Problem:** embedding cost scales with $|\mathcal{V}|$
 - Most words occur infrequently
 - Some unknown words are likely to exist
- **Solution:** sub-word tokenization

Vocabulary Size

Tokenization

- **Problem:** embedding cost scales with $|\mathcal{V}|$
 - Most words occur infrequently
 - Some unknown words are likely to exist
- **Solution:** sub-word tokenization

Unit	Vocab Size	Semantic Density	Sequence Length
Words	Large	High	Low
Characters	Small	Low	High
Subwords	Medium	Medium	Medium

Vocabulary Size

Tokenization

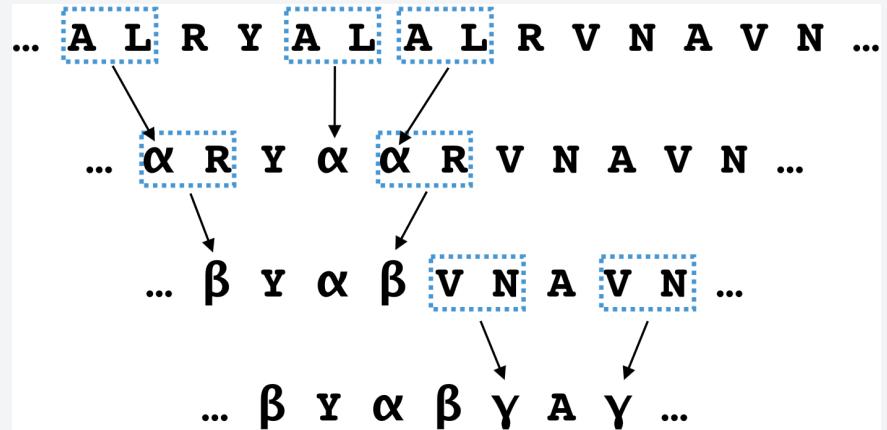
- **Problem:** embedding cost scales with $|\mathcal{V}|$
 - Most words occur infrequently
 - Some unknown words are likely to exist
- **Solution:** sub-word tokenization
- Which subwords?

Unit	Vocab Size	Semantic Density	Sequence Length
Words	Large	High	Low
Characters	Small	Low	High
Subwords	Medium	Medium	Medium

Subwords Selection

Tokenization

Vaswani et al. used Byte Pair Encodings (BPE) to automatically select 37k subword units



Kawano Koide & Imamura (2019). Seq2seq fingerprint with byte-pair encoding for predicting changes in protein stability upon single point mutation. IEEE/ACM transactions on computational biology and bioinformatics, 17(5), 1762-1772.

For other subword tokenizers, see [here](#) or [here](#)

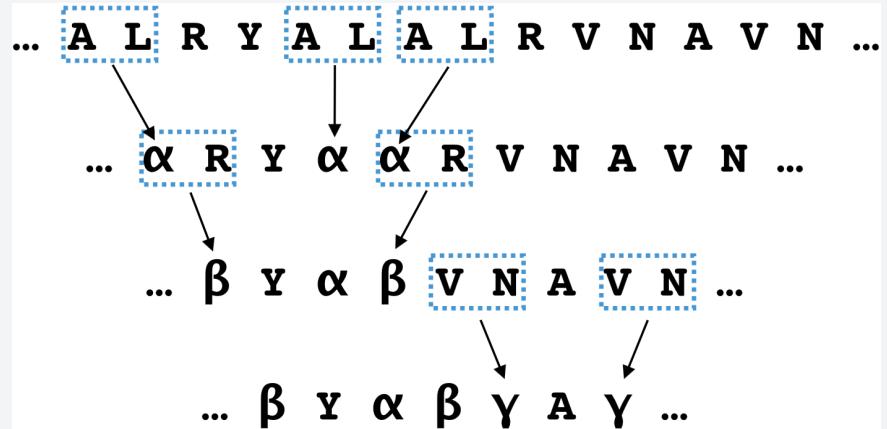
Subwords Selection

Tokenization

Vaswani et al. used Byte Pair Encodings (BPE) to automatically select 37k subword units

"Hello, y'all! How are you 😊 ?"

```
<_Hello> <,> <_y'> <all> <!> <_How>  
<are> <you> <[UNK]> <?>
```



Kawano Koide & Imamura (2019). Seq2seq fingerprint with byte-pair encoding for predicting changes in protein stability upon single point mutation. IEEE/ACM transactions on computational biology and bioinformatics, 17(5), 1762-1772.

For other subword tokenizers, see [here](#) or [here](#)

Subwords Selection

Tokenization

Vaswani et al. used Byte Pair Encodings (BPE) to automatically select 37k subword units

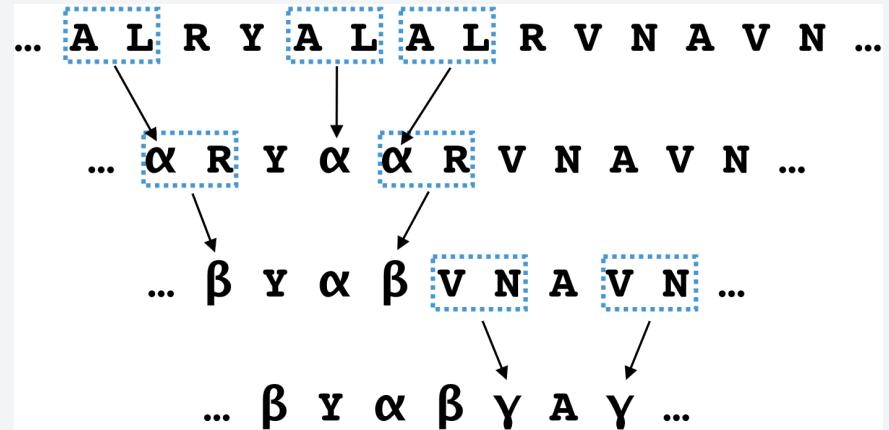
"Hello, y'all! How are you 😊 ?"

```
<_Hello> <,> <_y'> <all> <!> <_How>  
<are> <you> <[UNK]> <?>
```

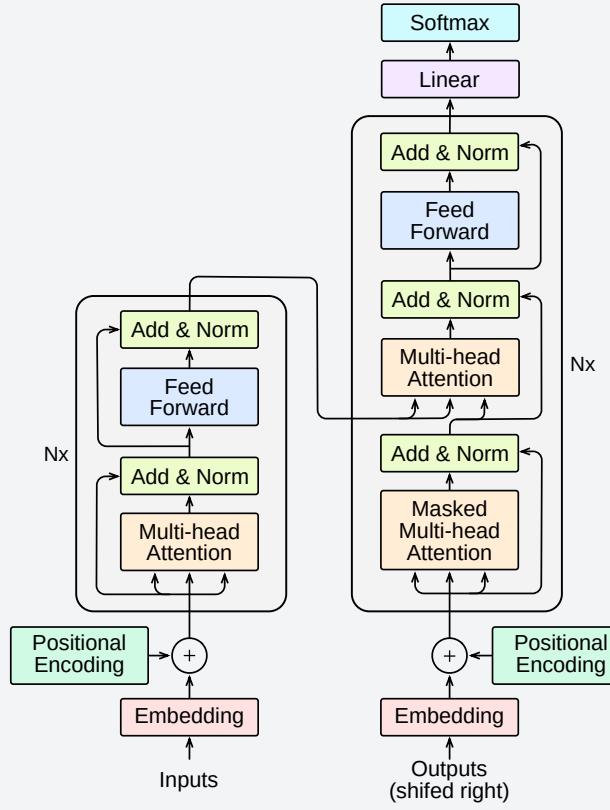
"The best gradient descent method is grad student descent"

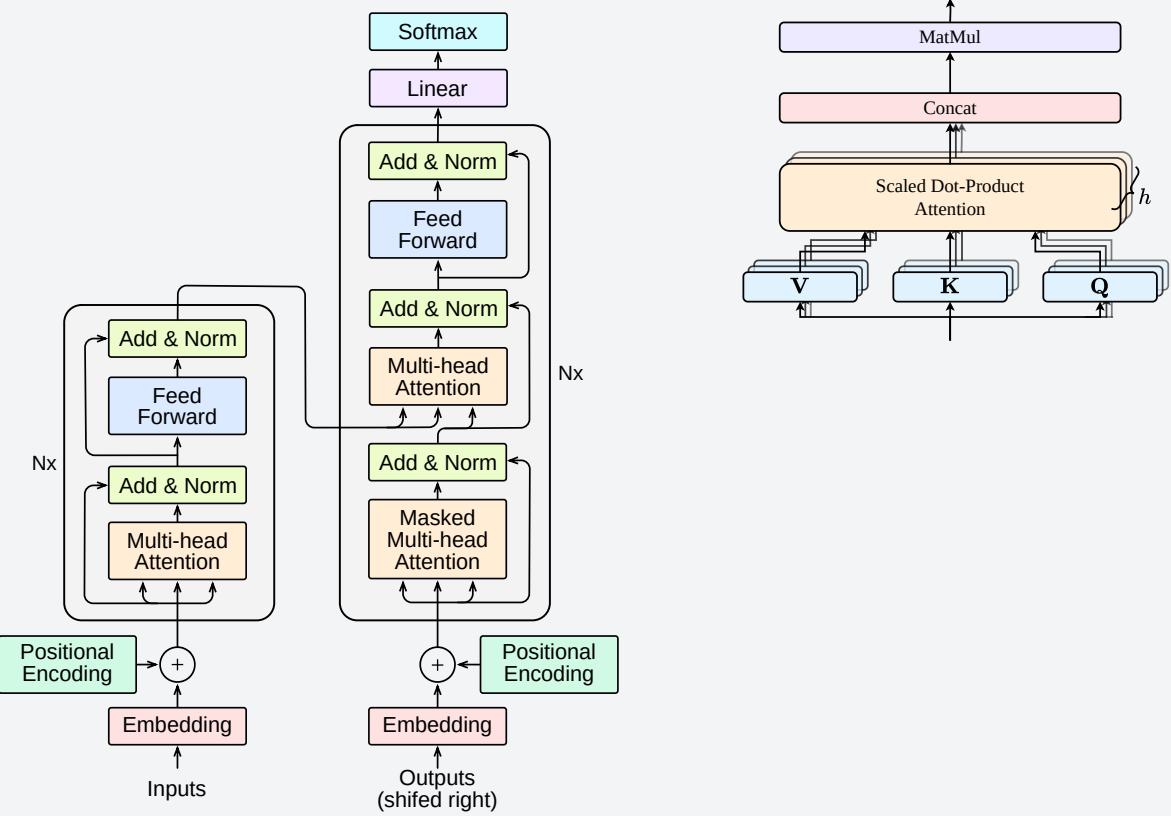
```
<_The> <_best> <_grad> <i> <ent> <_desc>  
<ent> <_method> <_is> <_grad> <_stud>  
<ent> <_desc> <ent>
```

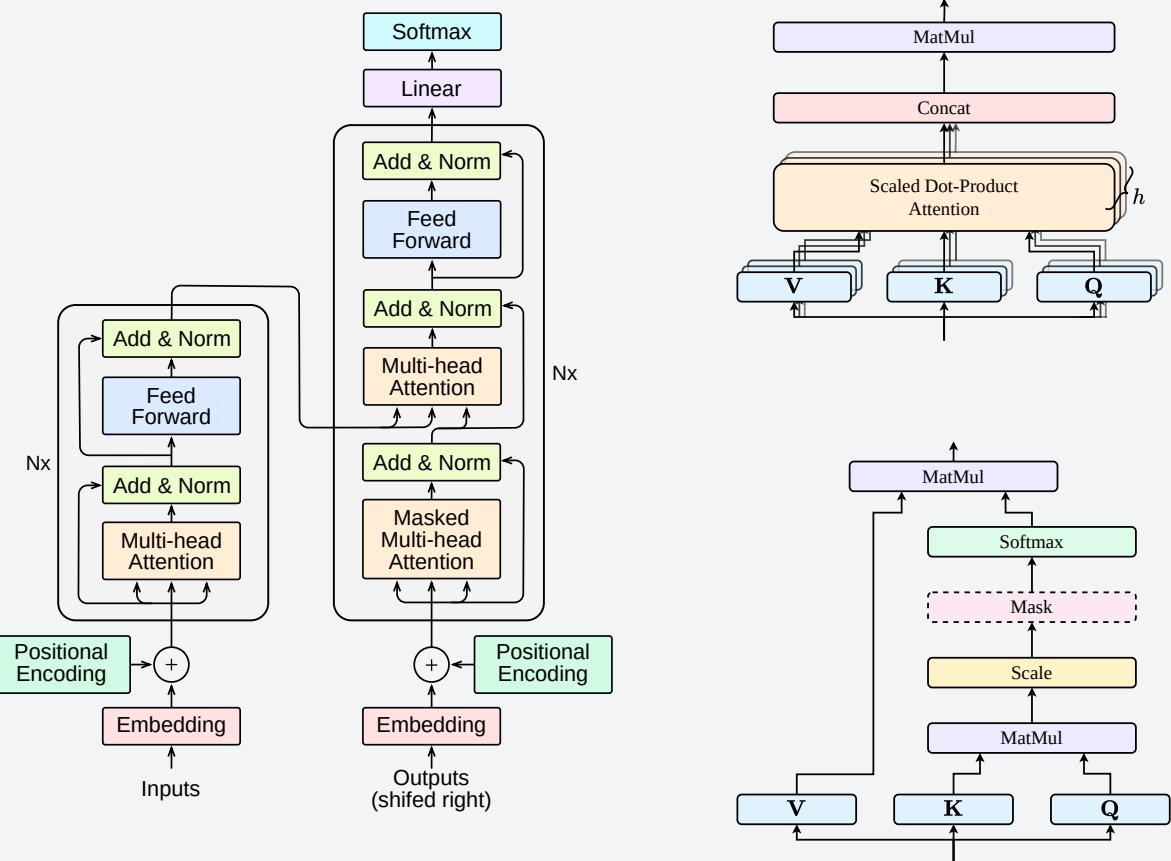
For other subword tokenizers, see [here](#) or [here](#)

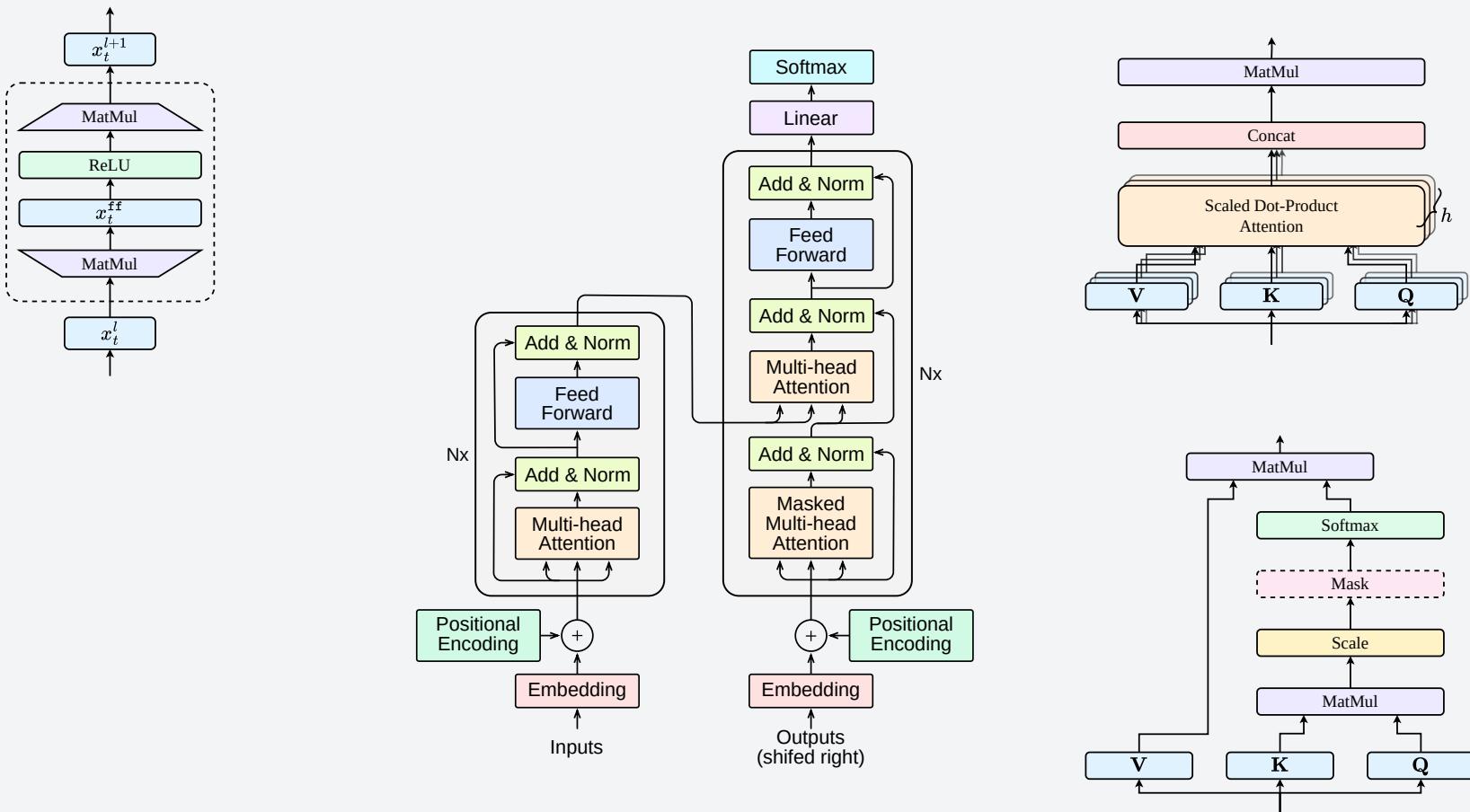


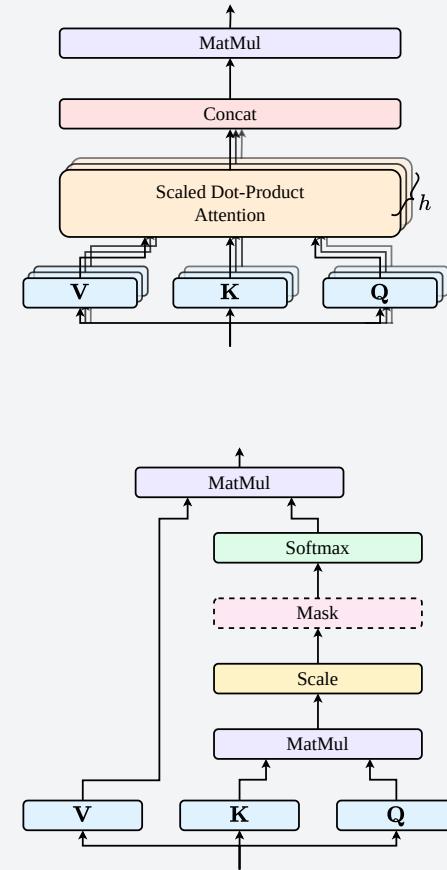
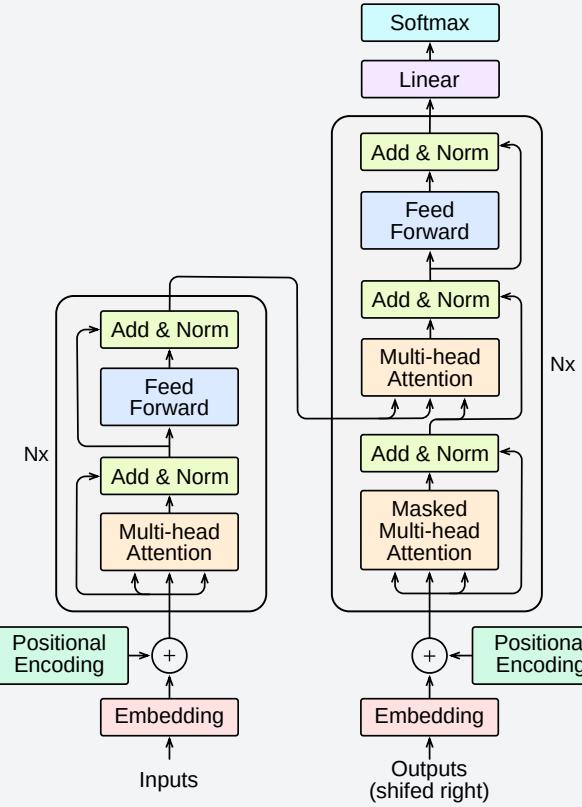
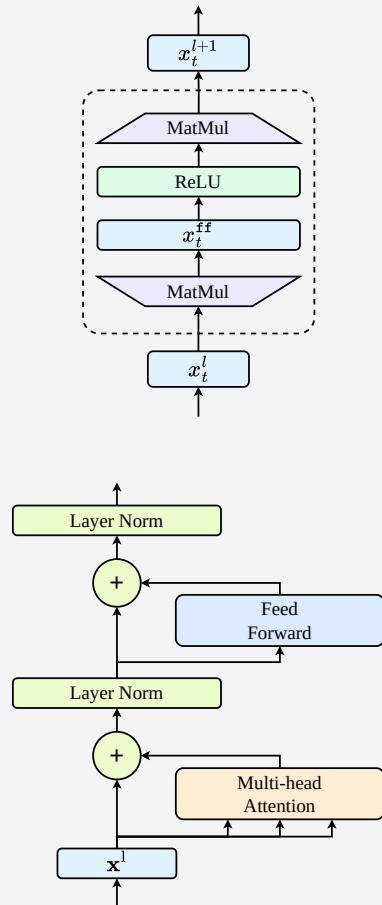
Kawano Koide & Imamura (2019). Seq2seq fingerprint with byte-pair encoding for predicting changes in protein stability upon single point mutation. IEEE/ACM transactions on computational biology and bioinformatics, 17(5), 1762-1772.











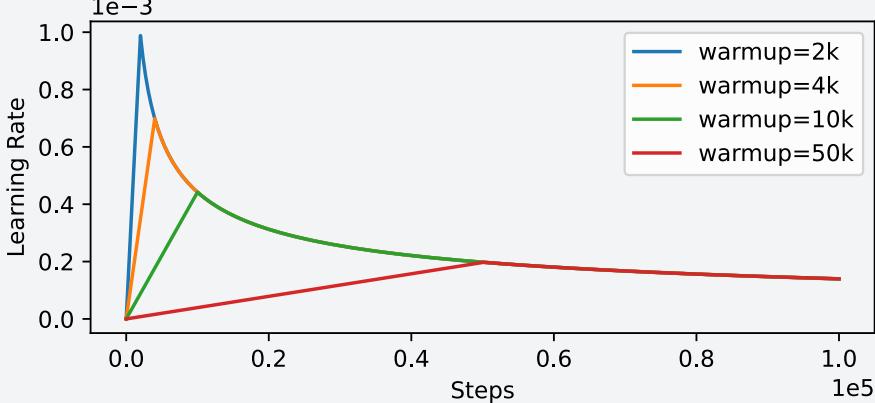
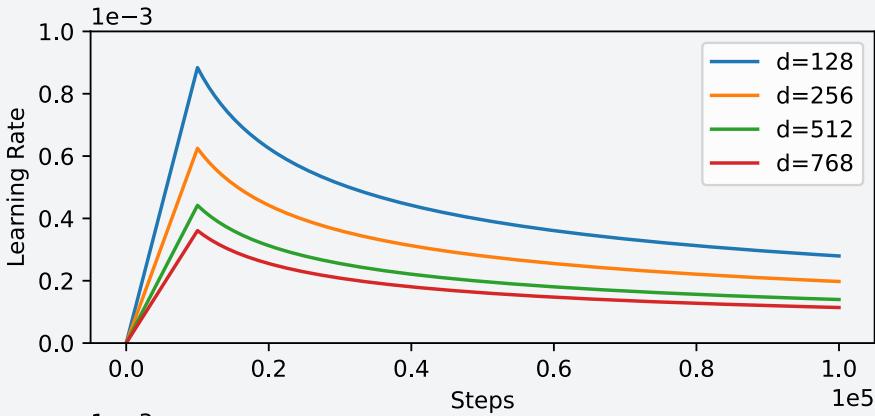
Training Transformers

Learning Rate Scheduling

Training Transformers

- Pre-training Transformers is unstable
- Vaswani et al. use a combination of **linear warmup** and **inverse square root decay**

$$\text{lrate} = d_{\text{model}}^{-0.5} \cdot \min(t^{-0.5}, t \cdot t_{\text{warmup}}^{-1.5})$$



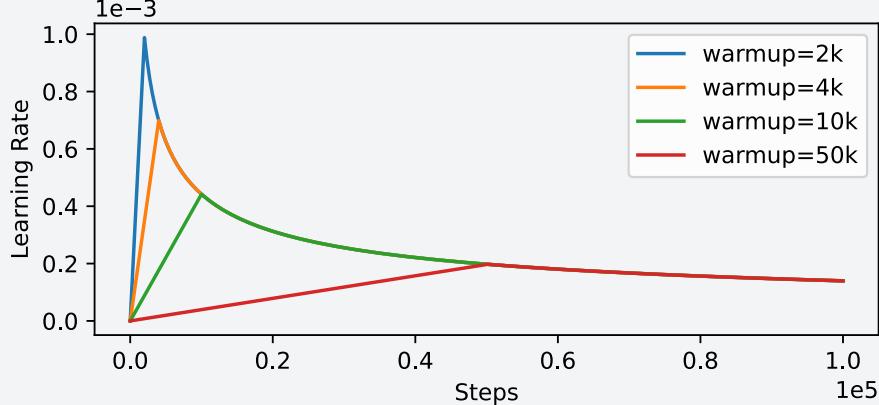
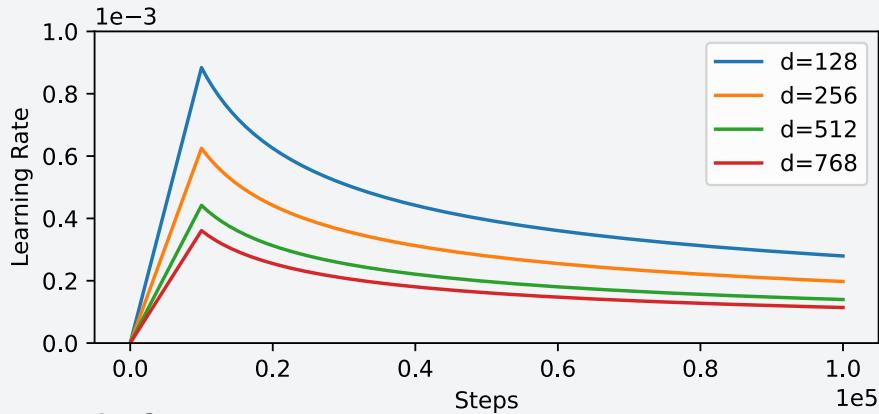
Learning Rate Scheduling

Training Transformers

- Pre-training Transformers is unstable
- Vaswani et al. use a combination of **linear warmup** and **inverse square root decay**

$$\text{lrate} = d_{\text{model}}^{-0.5} \cdot \min(t^{-0.5}, t \cdot t_{\text{warmup}}^{-1.5})$$

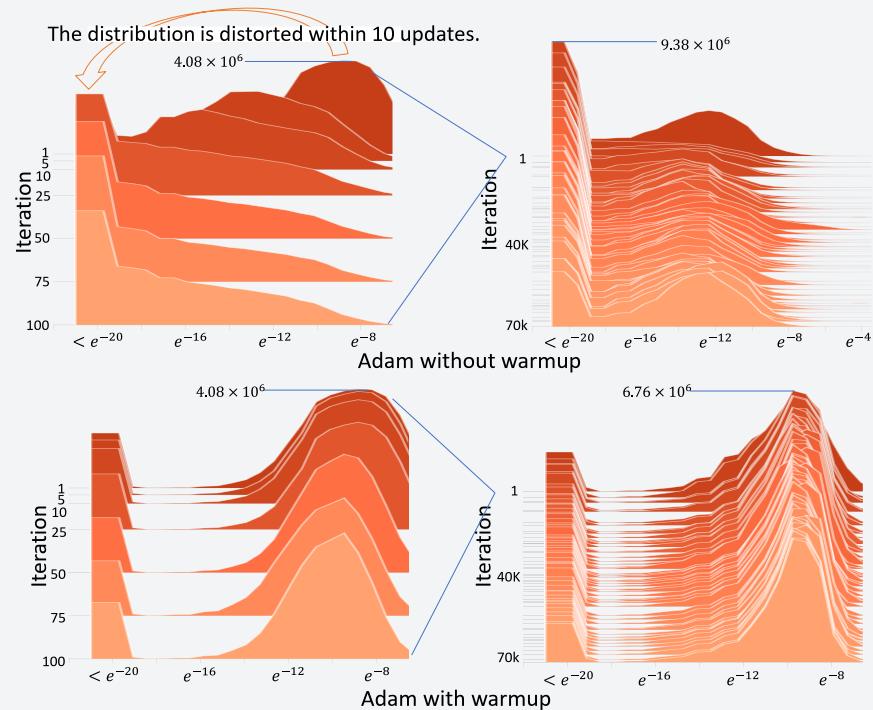
- Decreases impact of earliest steps (noise)
- Avoids plateaus towards end of training



Learning Rate Scheduling

Training Transformers

- Without warmup Adam gradient distributions collapse

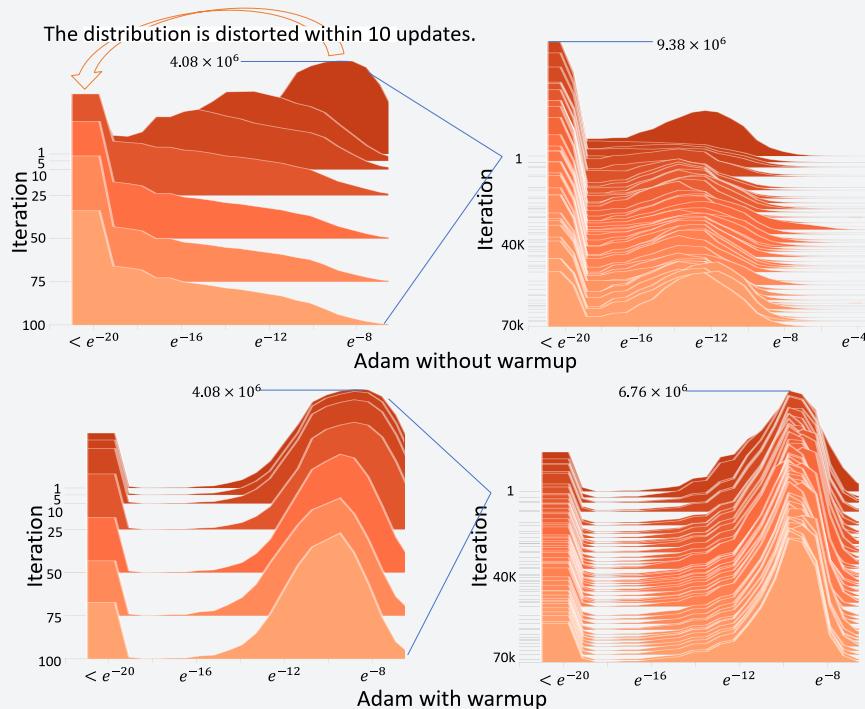


Liu et al. (2019). On the variance of the adaptive learning rate and beyond. arXiv:1908.03265.

Learning Rate Scheduling

Training Transformers

- Without warmup Adam gradient distributions collapse
- Other warmup strategies that stabilize training:
 - Batch size
 - Sequence length

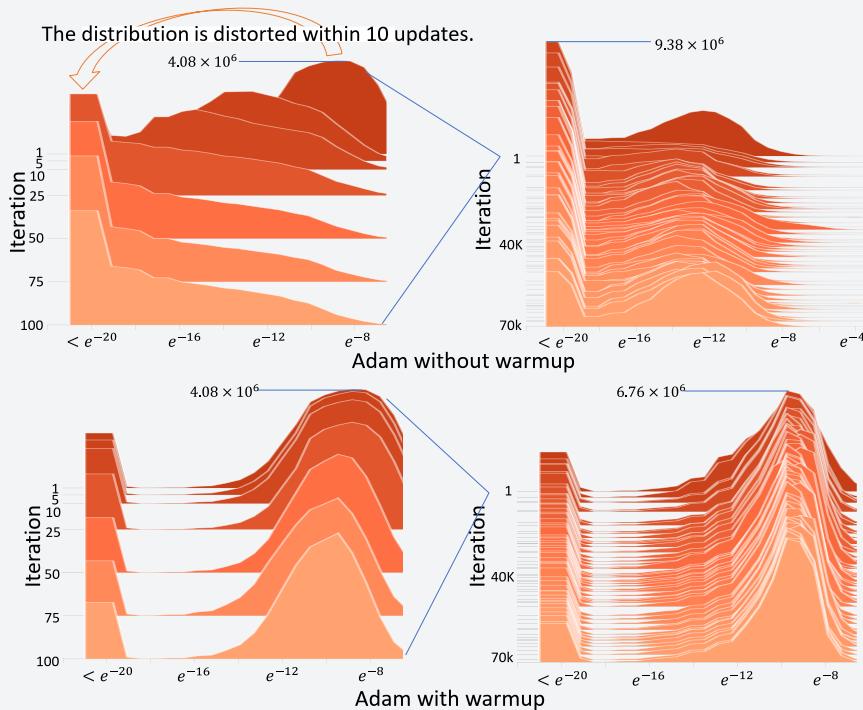


Liu et al. (2019). On the variance of the adaptive learning rate and beyond. arXiv:1908.03265.

Learning Rate Scheduling

Training Transformers

- Without warmup Adam gradient distributions collapse
- Other warmup strategies that stabilize training:
 - Batch size
 - Sequence length
- Adam variants have been proposed
 - For fine-tuning, AdamW is still king



Liu et al. (2019). On the variance of the adaptive learning rate and beyond. arXiv:1908.03265.

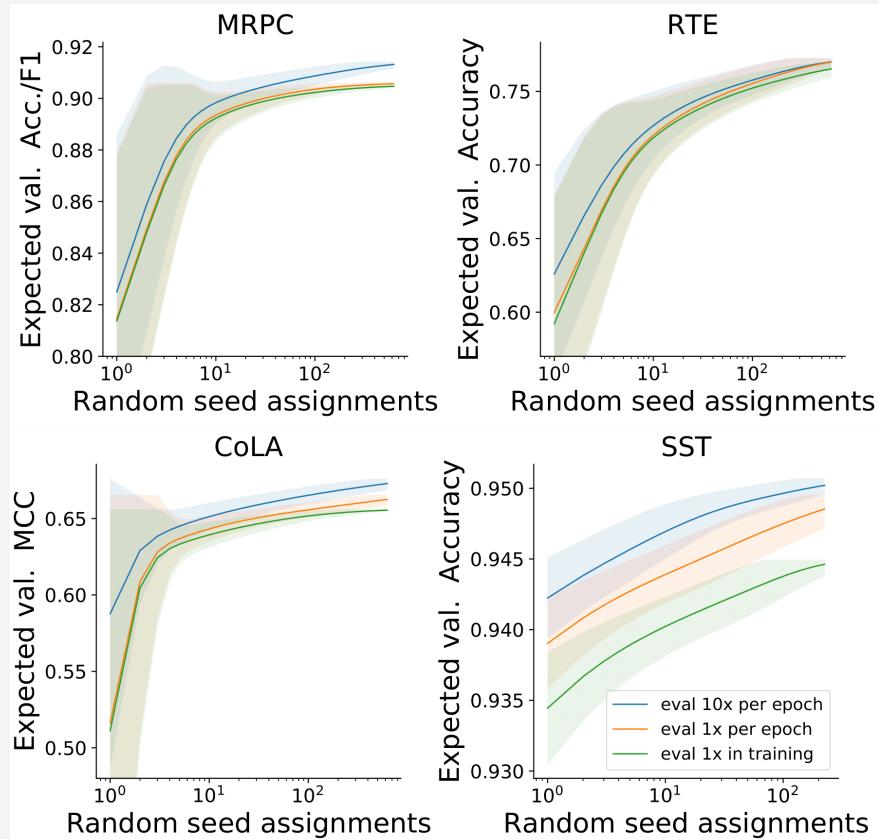
Fine-tuning Approaches

Training Transformers

- High variance in random seeds

Start many, stop early, continue some

Dodge et al. (2020) Fine-tuning pretrained language models: Weight initializations, data orders, and early stopping. arXiv:2002.06305



Fine-tuning Approaches

Training Transformers

- Assume $\mathcal{D}^{(\text{pre-train})} \gg \mathcal{D}^{(\text{fine-tune})}$

Fine-tuning Approaches

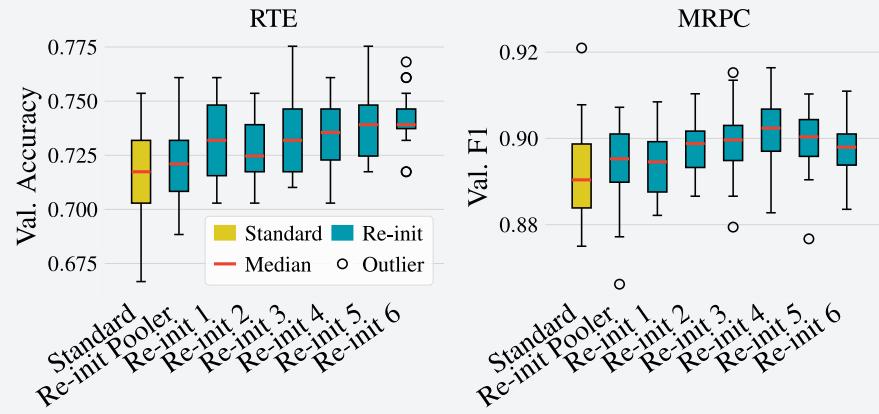
Training Transformers

- Assume $\mathcal{D}^{(\text{pre-train})} \gg \mathcal{D}^{(\text{fine-tune})}$
- Some tips:
 - **Intermediate task transfer:** first fine-tune on large related task, then fine-tun on small task

Fine-tuning Approaches

Training Transformers

- Assume $\mathcal{D}^{(\text{pre-train})} \gg \mathcal{D}^{(\text{fine-tune})}$
- Some tips:
 - **Intermediate task transfer:** first fine-tune on large related task, then fine-tun on small task
 - **Re-init:** completely re-train last layers

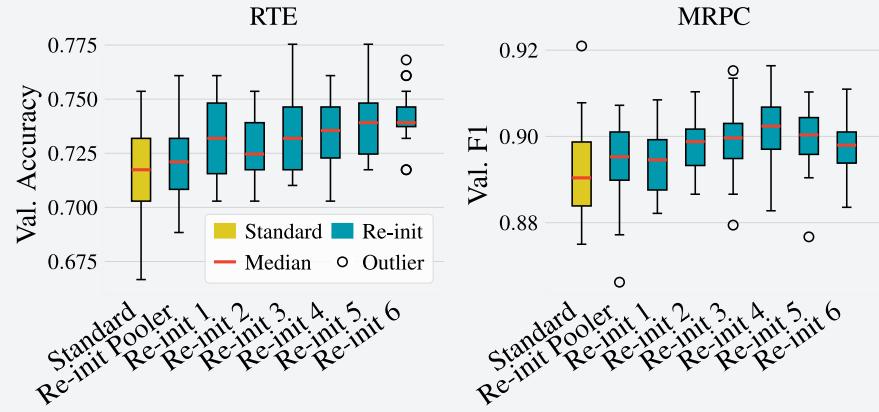


Zhang (2020). Revisiting few-sample BERT fine-tuning.
arXiv:2006.05987.

Fine-tuning Approaches

Training Transformers

- Assume $\mathcal{D}^{(\text{pre-train})} \gg \mathcal{D}^{(\text{fine-tune})}$
- Some tips:
 - **Intermediate task transfer:** first fine-tune on large related task, then fine-tun on small task
 - **Re-init:** completely re-train last layers
 - **Layer-wise Learning Rate Decay:**
 - low layers are general → low lr
 - high layers are task-specific → high lr

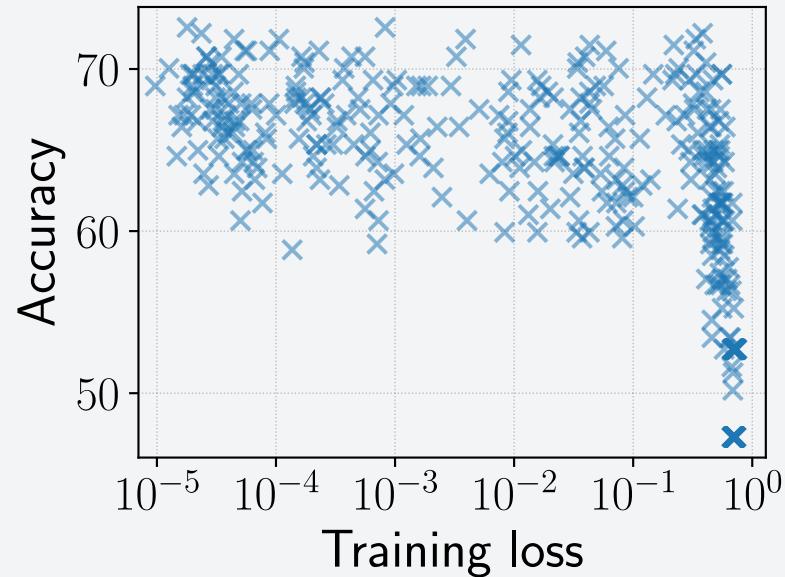


Zhang (2020). Revisiting few-sample BERT fine-tuning.
arXiv:2006.05987.

Fine-tuning Approaches

Training Transformers

- Biggest constraints are:
 1. Amount of data
 2. Compute budget
- If you have little data, long fine-tuning runs might still help



Mosbach, Andriushchenko, Klakow (2020). On the stability of fine-tuning bert: Misconceptions, explanations, and strong baselines. arXiv:2006.04884.

The End