

CME 295: Transformers & Large Language Models

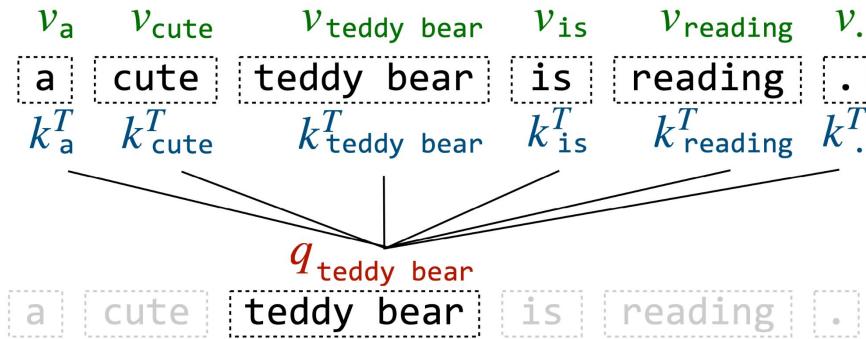


Afshine Amidi & Shervine Amidi

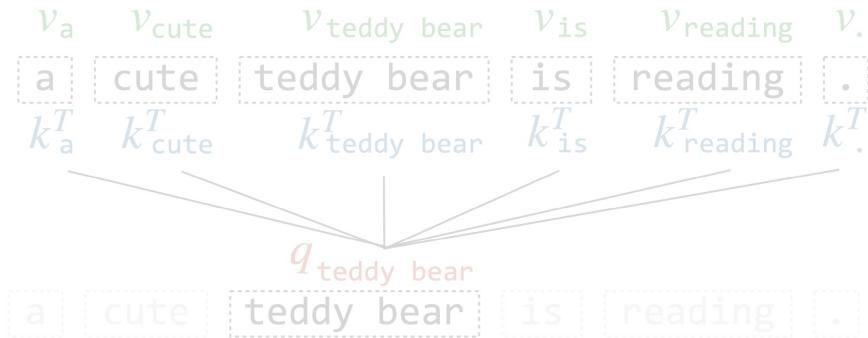


Recap of last episode...

Recap of last episode...



Recap of last episode...

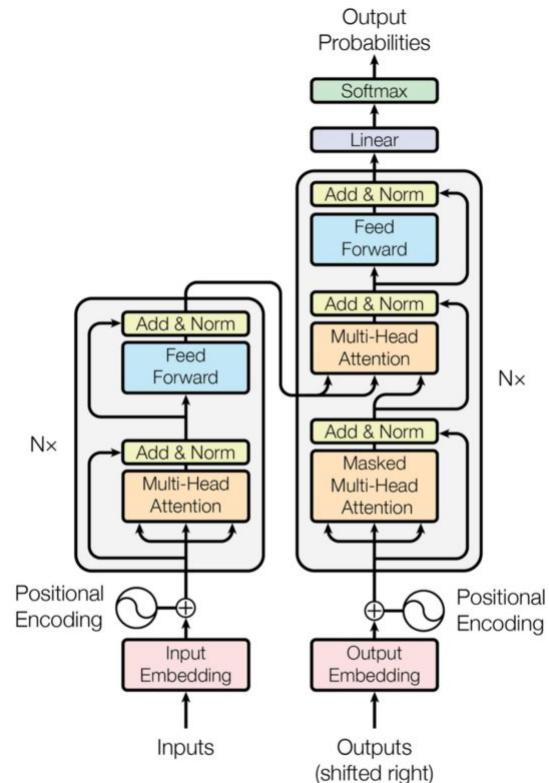


$$\text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$

Recap of last episode...



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

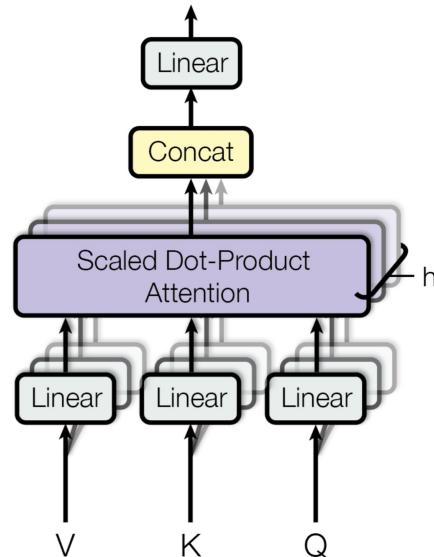


Recap of last episode...

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

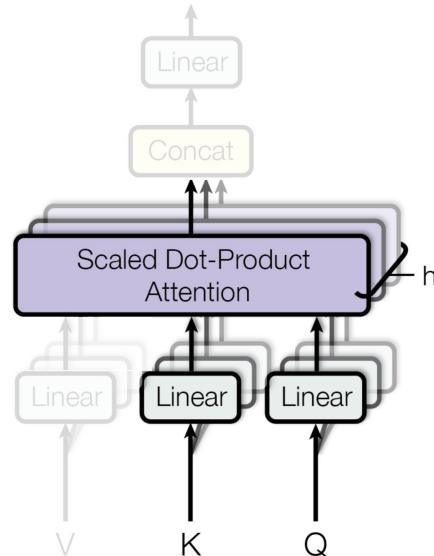
Recap of last episode...

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

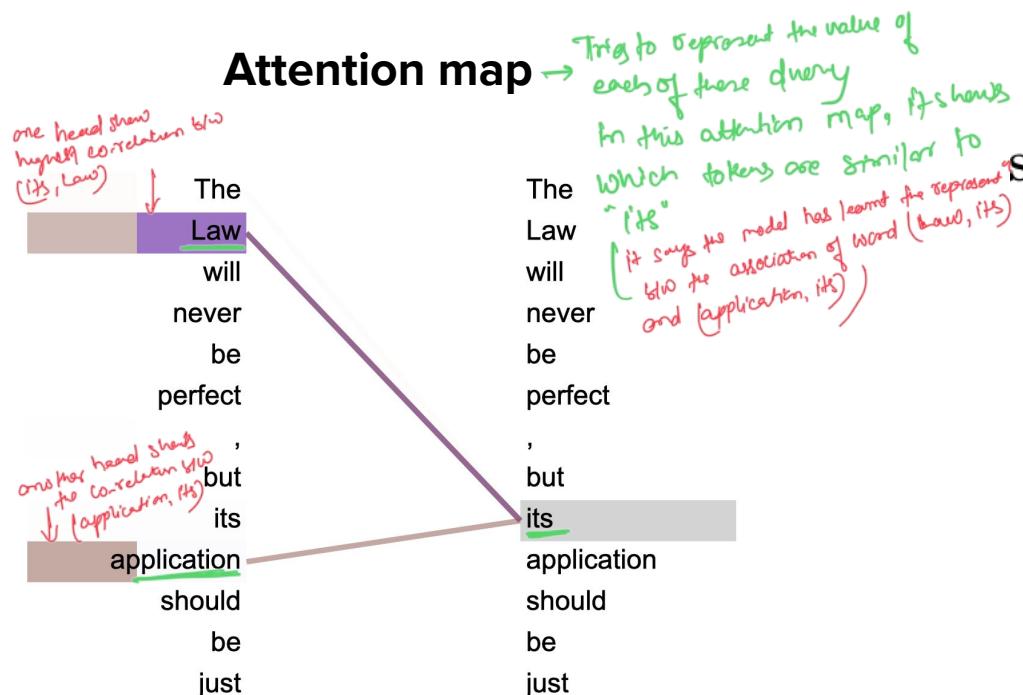


Recap of last episode...

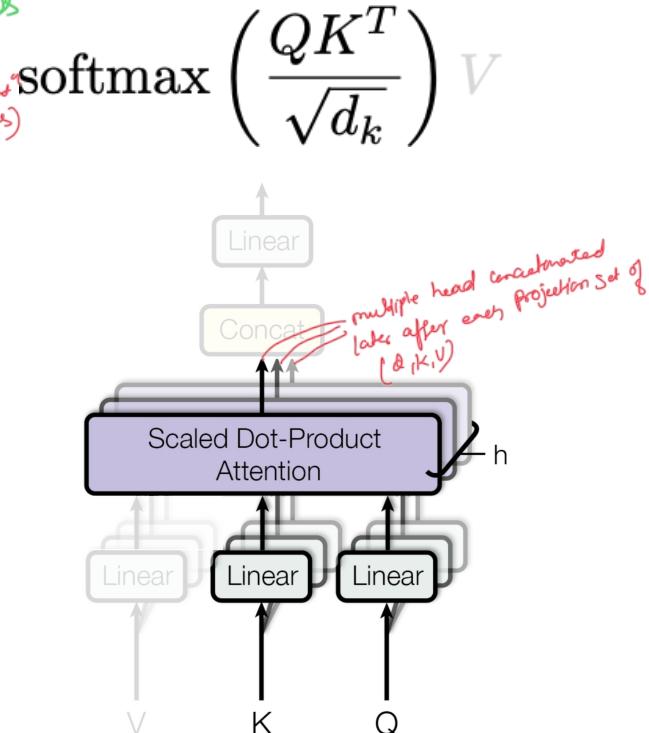
$$\text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$



Recap of last episode...



"Two attention heads, also in layer 5 of 6, apparently involved in anaphora resolution. [...] Isolated attentions from just the word 'its' for attention heads 5 and 6."



Suggested reading: original Transformer paper

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

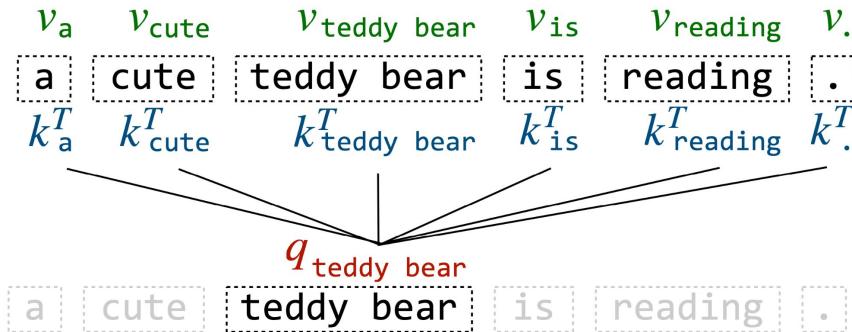


Transformers & Large Language Models

- Position embeddings
 - Layer normalization
 - Attention approximation
 - Transformer-based models
 - BERT deep dive
- In the case case of RNN LSTM where was processing each tokens to produce but having multiple temporal life it needs multiple possibility of combining but in transformer we process each token with all other tokens (sharing last order forget), to maintain the order position embedding played role in transformer

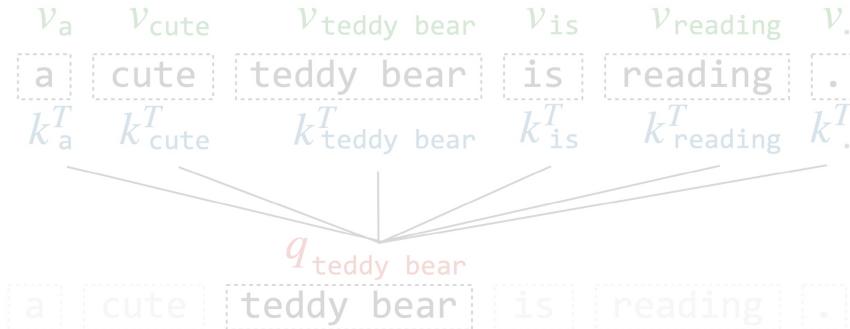
Need for position information

Motivation. Direct links "lose" position info



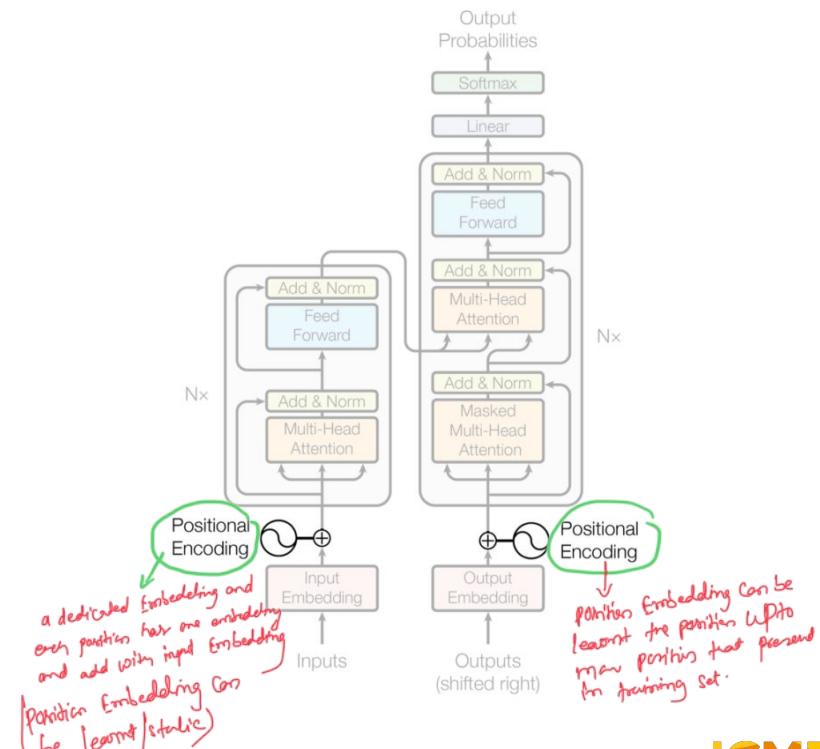
Need for position information

Motivation. Direct links "lose" position info



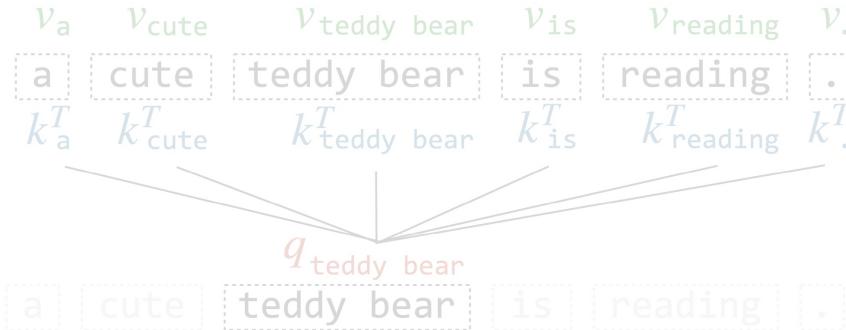
Idea. Add learned position-specific embedding

to token vector → "one embedding per position" = one way of
Position Embedding



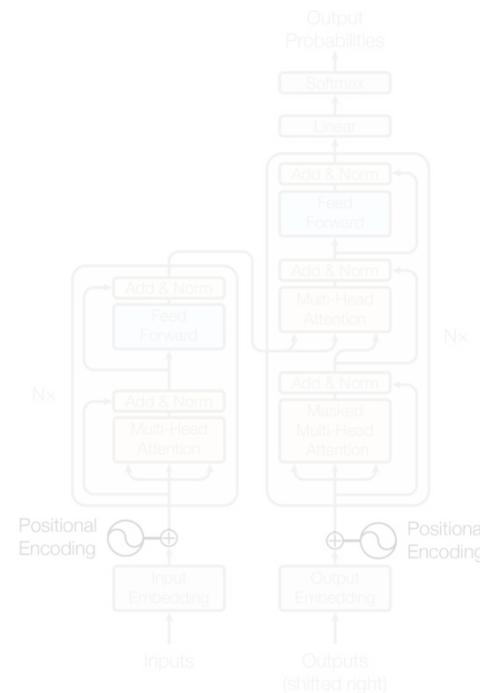
Need for position information

Motivation. Direct links "lose" position info



Idea. Add learned position-specific embedding to token vector

Limitations. Need to retrain for longer sequences



Hardcoded position embeddings

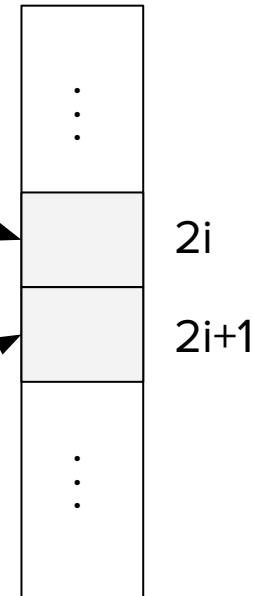
Idea. Hardcode values of position embedding

→ on another way of position Embedding.
using sin/cos

for every Embedding

$$PE_{m,2i} = \sin\left(\frac{m}{10000^{\frac{2i}{d_{model}}}}\right)$$

$$PE_{m,2i+1} = \cos\left(\frac{m}{10000^{\frac{2i}{d_{model}}}}\right)$$

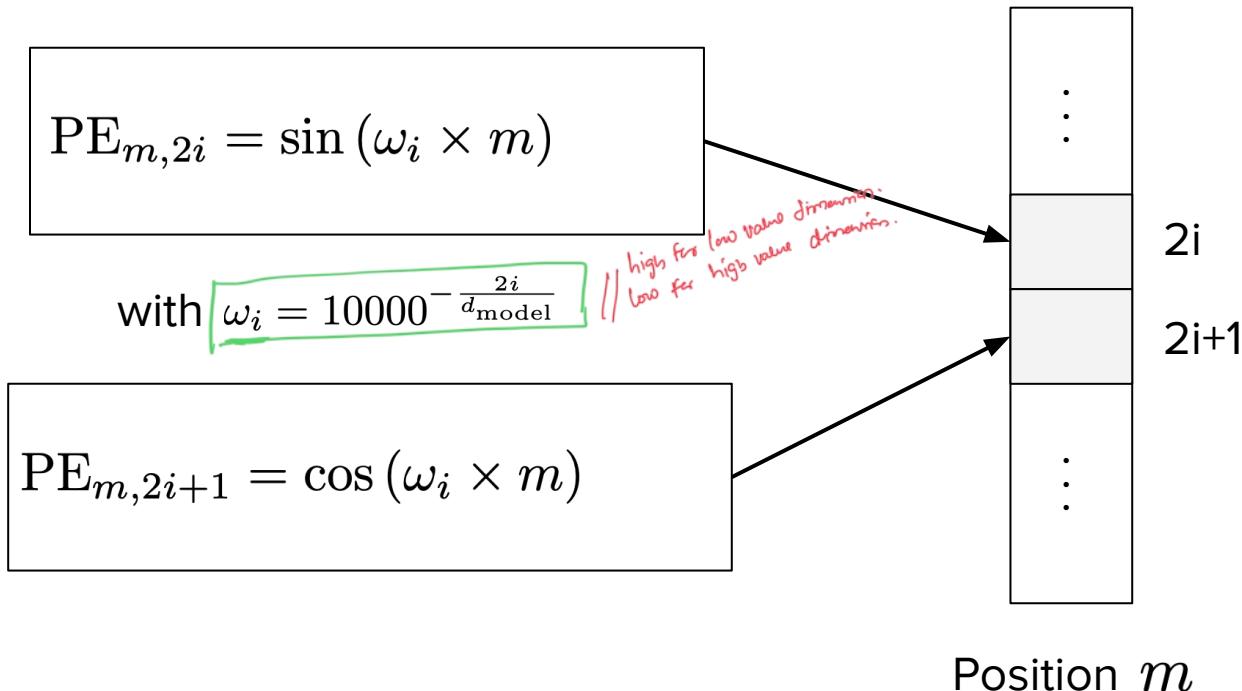


d_{model} : dimensions of token Embeddings

Position m

Hardcoded position embeddings

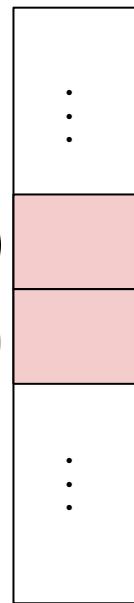
Idea. Hardcode values of position embedding



Hardcoded position embeddings

Idea. Hardcode values of position embedding

So writing cosine/Sine; It tries to find the relationship b/w Embedding and closer Embedding will have close angle
Cosine - relation b/w [mth position, nth position] for embedding i



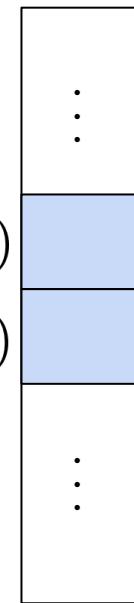
$$\text{PE}_{m,2i} = \sin(\omega_i \times m)$$

$$\text{PE}_{m,2i+1} = \cos(\omega_i \times m)$$

Position *m*

$$\text{PE}_{n,2i} = \sin(\omega_i \times n)$$

$$\text{PE}_{n,2i+1} = \cos(\omega_i \times n)$$



Position *n*

Hardcoded position embeddings

Idea. Hardcode values of position embedding

$$\cos(a - b) = \cos(a) \cos(b) + \sin(a) \sin(b)$$

Hardcoded position embeddings

Idea. Hardcode values of position embedding

$$\cos(a - b) = \cos(a) \cos(b) + \sin(a) \sin(b)$$

$$\cos(\omega_i(m - n)) = \cos(\omega_i \times m) \cos(\omega_i \times n) + \sin(\omega_i \times m) \sin(\omega_i \times n)$$

Hardcoded position embeddings

Idea. Hardcode values of position embedding

$$\cos(a - b) = \cos(a) \cos(b) + \sin(a) \sin(b)$$

$$\cos(\omega_i(m - n)) = \cos(\omega_i \times m) \cos(\omega_i \times n) + \sin(\omega_i \times m) \sin(\omega_i \times n)$$

$$\langle \text{PE}_{\textcolor{red}{m}}, \text{PE}_{\textcolor{blue}{n}} \rangle = \dots + \cos(\omega_i(\textcolor{red}{m} - \textcolor{blue}{n})) + \dots$$

If calculates cosine similarity
b/w (m^{th} position, n^{th} position) embedding i

Hardcoded position embeddings

Idea. Hardcode values of position embedding

$$\cos(a - b) = \cos(a) \cos(b) + \sin(a) \sin(b)$$

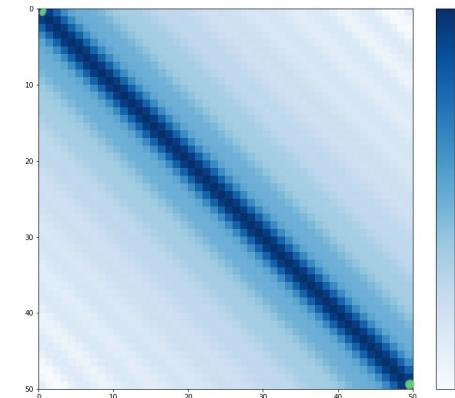
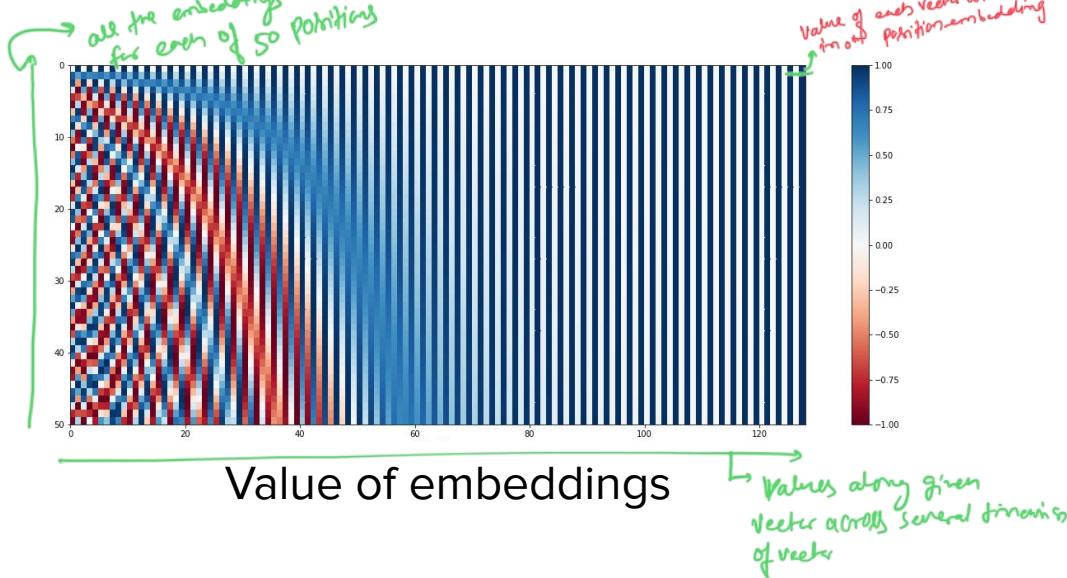
$$\cos(\omega_i(m - n)) = \cos(\omega_i \times m) \cos(\omega_i \times n) + \sin(\omega_i \times m) \sin(\omega_i \times n)$$

$$\langle \text{PE}_{\textcolor{brown}{m}}, \text{PE}_{\textcolor{brown}{n}} \rangle = \dots + \cos(\omega_i(m - n)) + \dots$$

$$\langle \text{PE}_{\textcolor{brown}{m}}, \text{PE}_{\textcolor{blue}{n}} \rangle = f(\textcolor{brown}{m} - \textcolor{blue}{n})$$

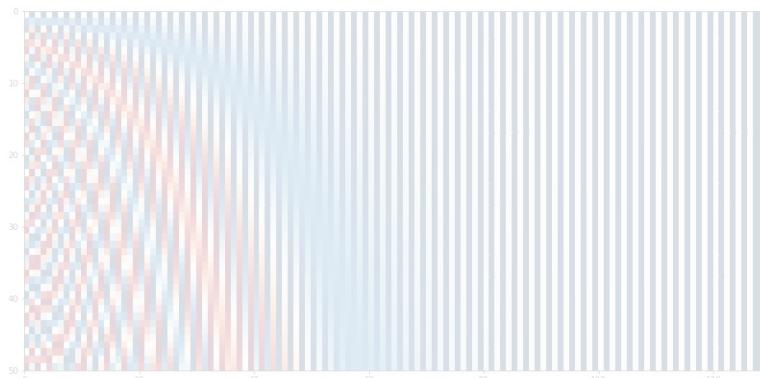
Hardcoded position embeddings

Idea. Hardcode values of position embedding

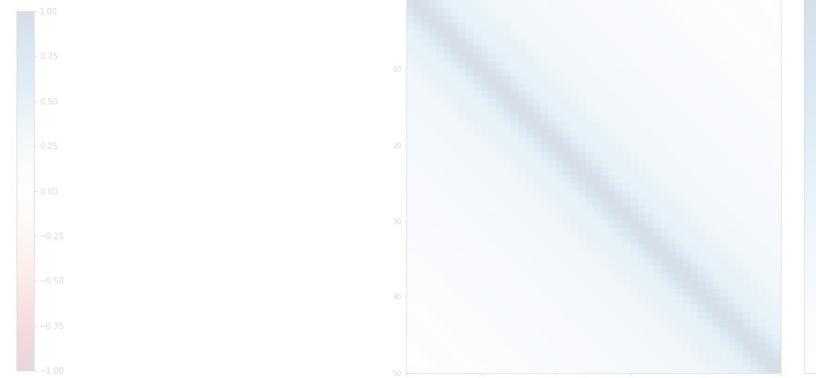


Hardcoded position embeddings

Idea. Hardcode values of position embedding



Value of embeddings

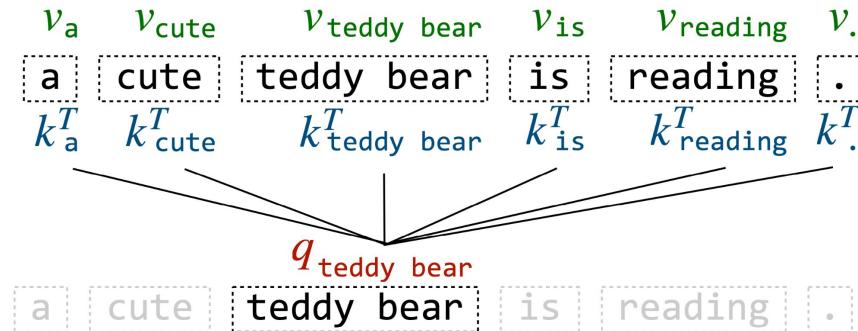


Similarity between positions

Benefits. Can extend to any sequence length

From absolute to relative position info

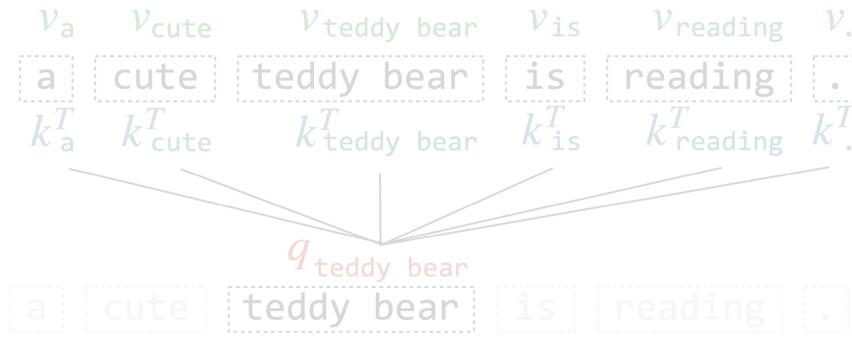
Motivation. Actually, we really care about relative position **with attention layer**



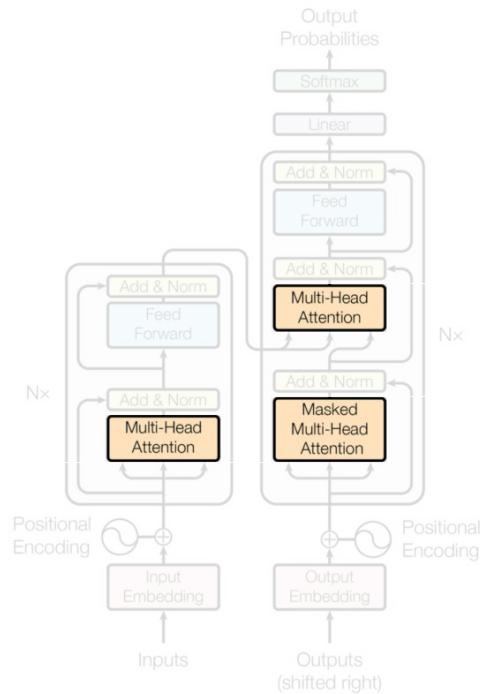
Position Embedding does train/learn for positions
but need to add with attention layer
but we care about finding the similarity b/w tokens to
self-attention component by why not integrating
self-attention layer with position embeddings.

From absolute to relative position info

Motivation. Actually, we really care about relative position **with attention layer**

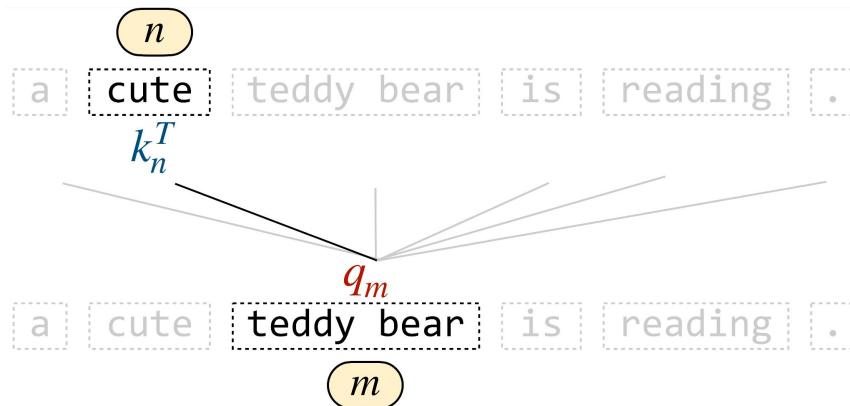


Idea. Let's change the attention layer instead



Linear bias in attention layer

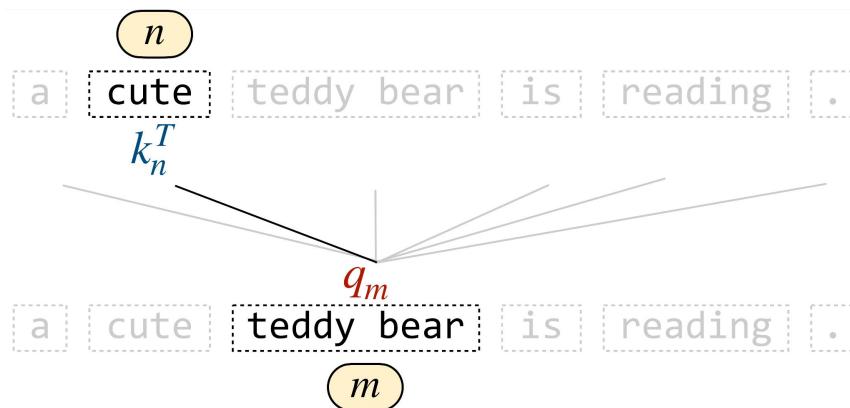
Idea. Add bias to query-key scores



$$\text{softmax} \left(\frac{\langle q_m, k_n \rangle}{\sqrt{d_k}} + \text{bias}(m, n) \right)$$

Linear bias in attention layer

Idea. Add bias to query-key scores



T5 bias. Bias is learned per head

$$\text{softmax} \left(\frac{\langle q_m, k_n \rangle}{\sqrt{d_k}} + \text{bias}(m, n) \right)$$

it basically says that some tokens are more similar compared to other tokens.

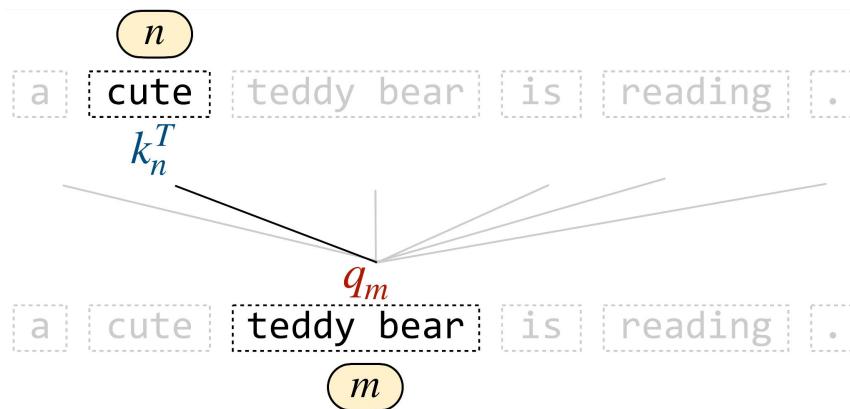
① approach

$$\text{bias}(m, n) = \beta_{\text{bucket}(m-n)}$$

\downarrow
bias(m, n) is getting learnt using the concept of Bucketization of all (m, n) into some buckets and learn which bucket will be best fit for injecting into softmax.

Linear bias in attention layer

Idea. Add bias to query-key scores



$$\text{softmax} \left(\frac{\langle q_m, k_n \rangle}{\sqrt{d_k}} + \text{bias}(m, n) \right)$$

T5 bias. Bias is learned per head

ALiBi. Bias is **linear, deterministic**, not bounded

↓
attention
with linear biases

$$\text{bias}(m, n) = \beta_{\text{bucket}(m-n)}$$

② approach

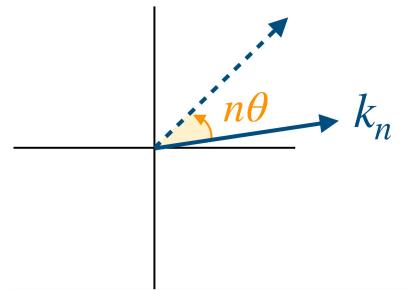
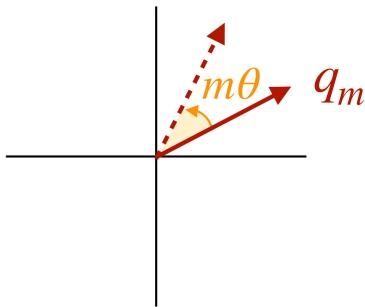
$$\text{bias}(m, n) = \mu \times (n - m)$$

↳ instead of learning
more biases, let's have deterministic formula.

Default choice nowadays: rotations in attention layer

↓
important | default choice for attention layer

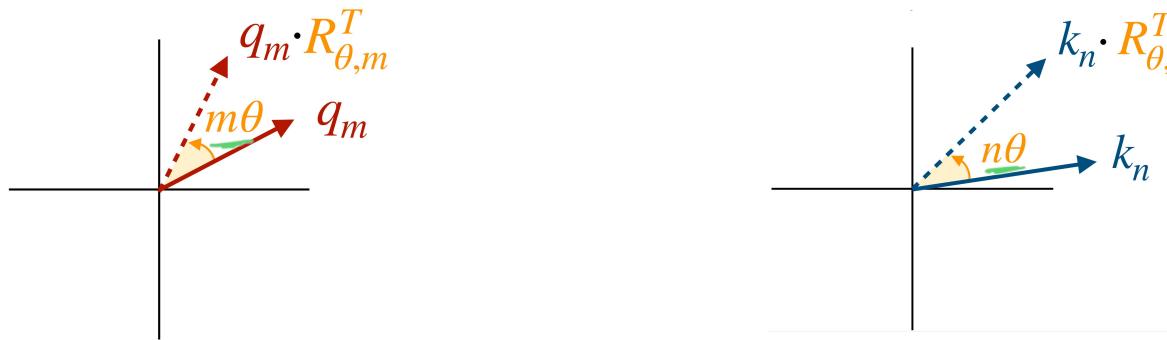
Idea. Rotate query and key vectors with rotation matrix



Default choice nowadays: rotations in attention layer

↓
Rotary Position Embeddings.

Idea. Rotate query and key vectors with rotation matrix

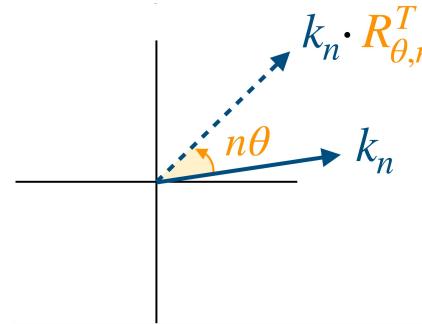
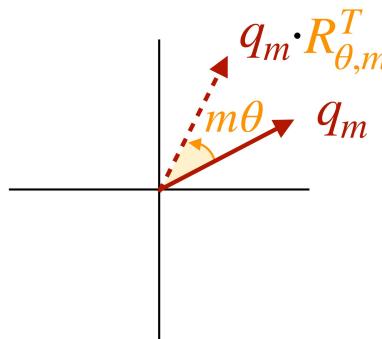


Formula. Use rotation matrix: $R_{\theta,m} = \begin{pmatrix} \cos(m\theta) & -\sin(m\theta) \\ \sin(m\theta) & \cos(m\theta) \end{pmatrix}$

Default choice nowadays: rotations in attention layer

on ideal position embeddings ← [RoPE = Rotary Position Embeddings]

Idea. Rotate query and key vectors with rotation matrix → *End up with duality that will be function w/ relative distance b/w two.*



Formula. Use rotation matrix: $R_{\theta,m} = \begin{pmatrix} \cos(m\theta) & -\sin(m\theta) \\ \sin(m\theta) & \cos(m\theta) \end{pmatrix}$

Default choice nowadays: rotations in attention layer

Extension to dimension $d > 2$. Rotate every block of dimension 2

$$R_{\theta,m} = \begin{pmatrix} \text{block}_1 & 0 & \cdots & 0 \\ 0 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & \text{block}_{\frac{d_k}{2}} \end{pmatrix} \quad \text{block}_i = \begin{pmatrix} \cos(m\theta_i) & -\sin(m\theta_i) \\ \sin(m\theta_i) & \cos(m\theta_i) \end{pmatrix}$$

Default choice nowadays: rotations in attention layer

Extension to dimension $d > 2$. Rotate every block of dimension 2

$$R_{\theta,m} = \begin{pmatrix} \text{block}_1 & 0 & \cdots & 0 \\ 0 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & \text{block}_{\frac{d_k}{2}} \end{pmatrix} \quad \text{block}_i = \begin{pmatrix} \cos(m\theta_i) & -\sin(m\theta_i) \\ \sin(m\theta_i) & \cos(m\theta_i) \end{pmatrix}$$

Benefits. Relative distance nicely captured:

$$q_m k_n^T = x_m W_q R_{\theta,n-m} W_k^T x_n^T$$

Annotations:

- green arrow from q_m to $R_{\theta,n-m}$: rotating key by angle $m\theta$
- green arrow from k_n^T to $R_{\theta,n-m}$: rotation matrix
- green arrow from q_m to x_m : rotating query by angle $m\theta$
- green arrow from k_n^T to x_n^T : rotating query by angle $n\theta$

Default choice nowadays: rotations in attention layer

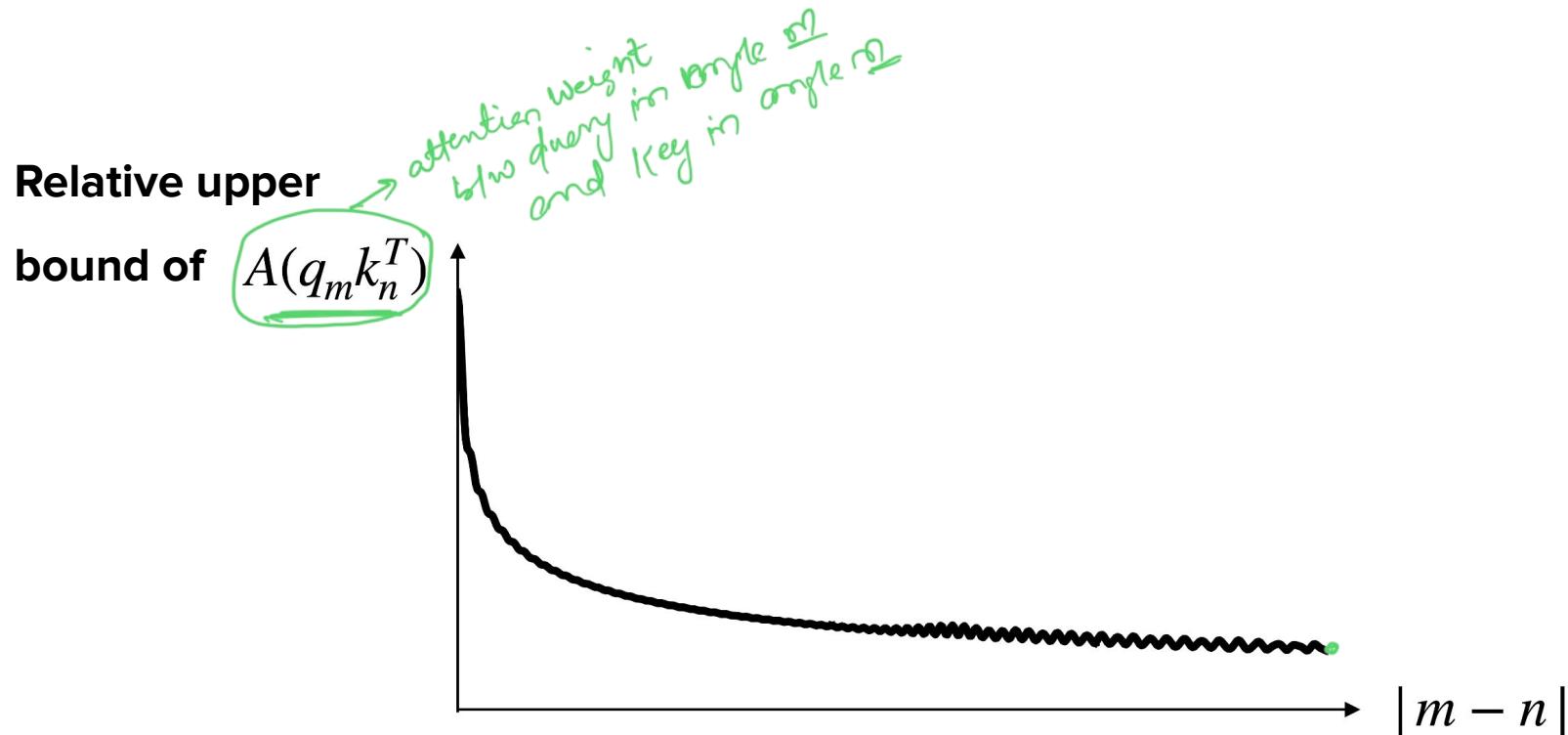
Extension to dimension $d > 2$. Rotate every block of dimension 2

$$R_{\theta,m} = \begin{pmatrix} \text{block}_1 & 0 & \cdots & 0 \\ 0 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & \text{block}_{\frac{d_k}{2}} \end{pmatrix} \quad \text{block}_i = \begin{pmatrix} \cos(m\theta_i) & -\sin(m\theta_i) \\ \sin(m\theta_i) & \cos(m\theta_i) \end{pmatrix}$$

Benefits. Relative distance nicely captured:

$$q_{\mathbf{m}} k_{\mathbf{n}}^T = x_{\mathbf{m}} W_q R_{\theta,\mathbf{n}-\mathbf{m}} W_k^T x_{\mathbf{n}}^T$$

Attention weight has a long-term decay





Transformers & Large Language Models

Position embeddings

Layer normalization

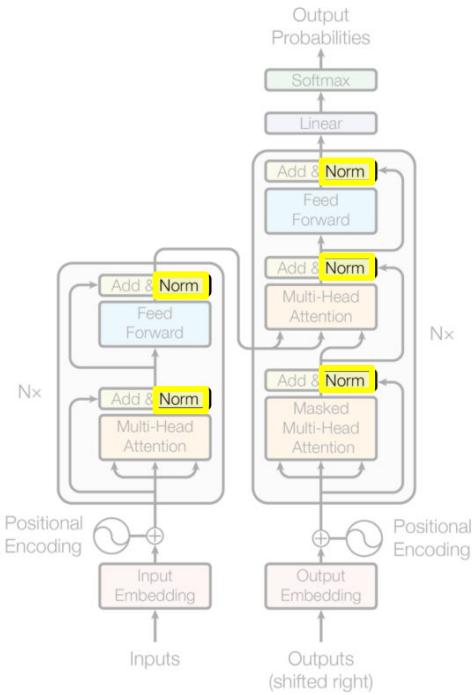
Attention approximation

Transformer-based models

BERT deep dive

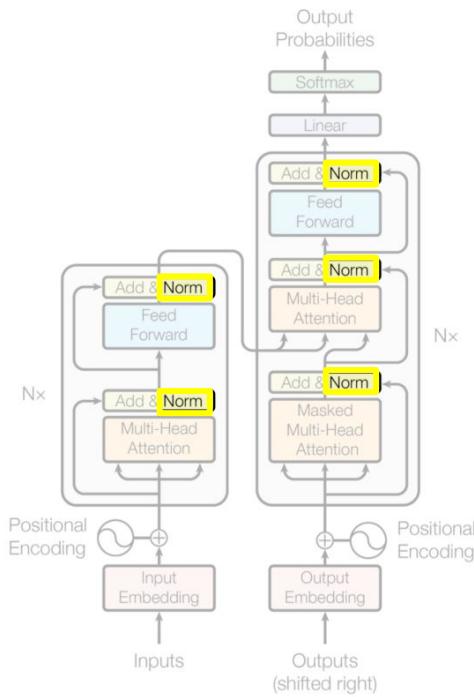
↑ used for normalizing
components of vector.

Layer normalization



Layer normalization

LN = Layer Normalization



Layer normalization

LN = Layer Normalization

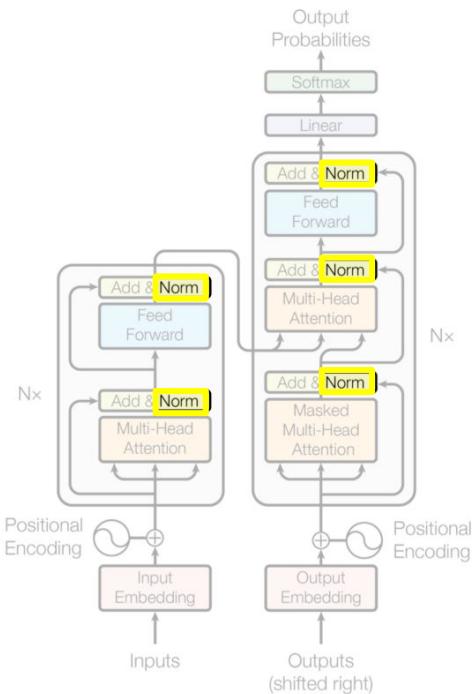
↳ helps training stability
and convergence time

Learnable parameters:

$$\text{LN}(x) = \gamma \hat{x} + \beta$$

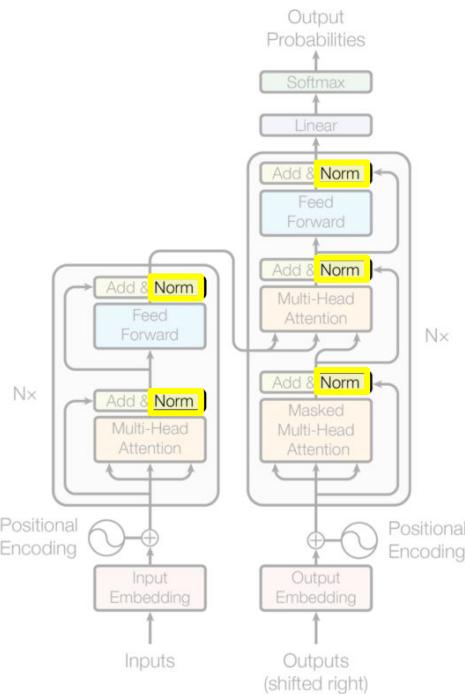
with $\hat{x} = \frac{x - \mu}{\sqrt{\sigma^2 + \epsilon}}$

$$\text{where } \mu = \frac{1}{d} \sum_{i=1}^d x_i \quad \text{and} \quad \sigma^2 = \frac{1}{d} \sum_{i=1}^d (x_i - \mu)^2$$



Layer normalization

LN = Layer Normalization



$$\text{LN}(x) = \gamma \hat{x} + \beta$$

$$\text{with } \hat{x} = \frac{x - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

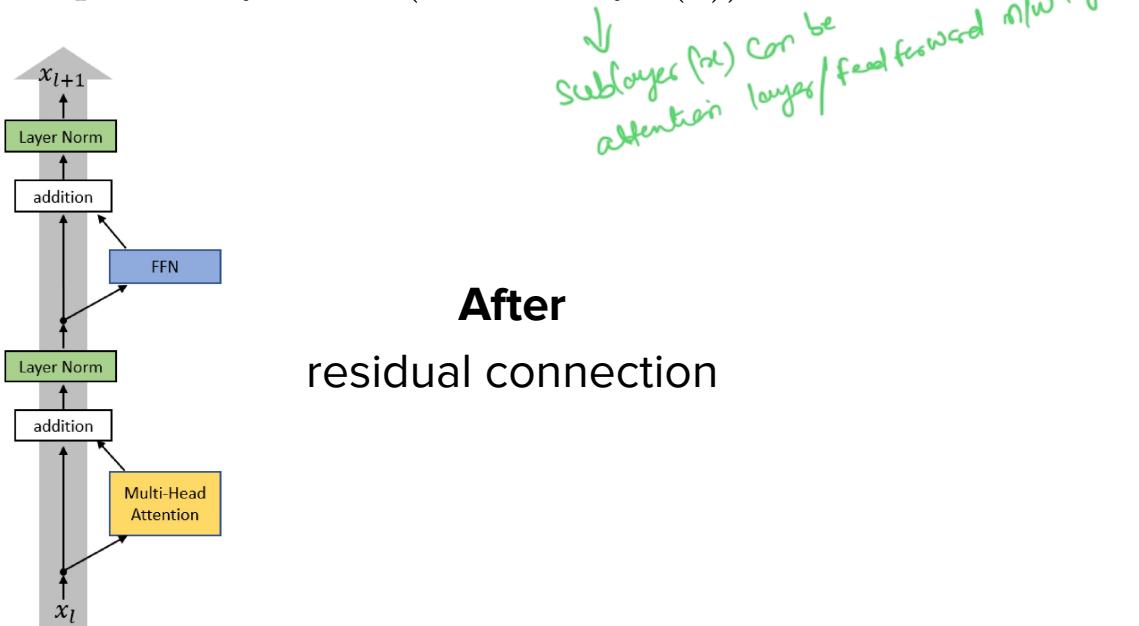
$$\text{where } \mu = \frac{1}{d} \sum_{i=1}^d x_i \quad \text{and} \quad \sigma^2 = \frac{1}{d} \sum_{i=1}^d (x_i - \mu)^2$$

Benefits. Helps with training stability and convergence.

Different kinds of layer normalization

Post-Norm

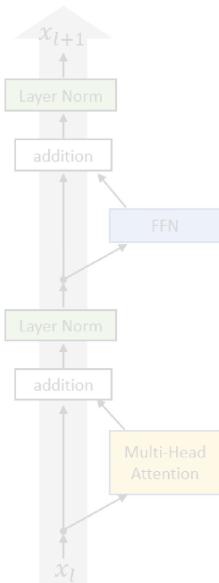
$$\text{Output} = \text{LayerNorm}(x + \text{SubLayer}(x))$$



Different kinds of layer normalization

Post-Norm

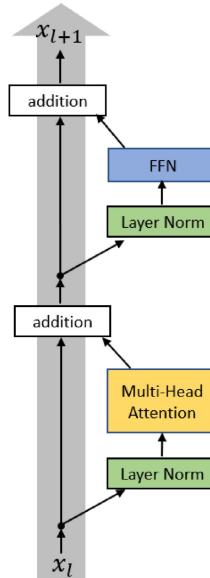
$$\text{Output} = \text{LayerNorm}(x + \text{SubLayer}(x))$$



After
residual connection

Pre-Norm

$$\text{Output} = x + \text{SubLayer}(\text{LayerNorm}(x))$$

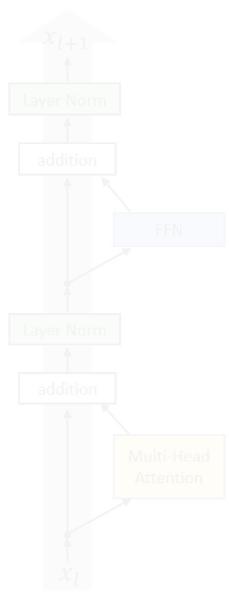


Before
residual connection

Different kinds of layer normalization

Post-LN

$$\text{Output} = \text{LayerNorm}(x + \text{SubLayer}(x))$$



Nowadays:

Pre-Norm

$$\text{Output} = x + \text{SubLayer}(\text{LayerNorm}(x))$$

After
residual connection



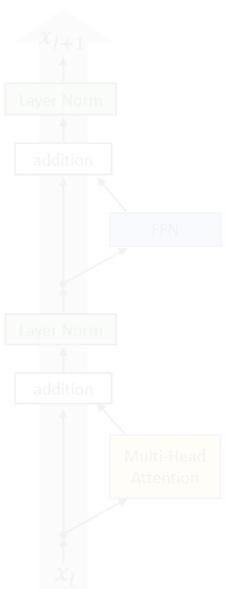
Pre-LN

Before
residual connection

Different kinds of layer normalization

Post-LN

$$\text{Output} = \text{LayerNorm}(x + \text{SubLayer}(x))$$



Nowadays:

Pre-LN

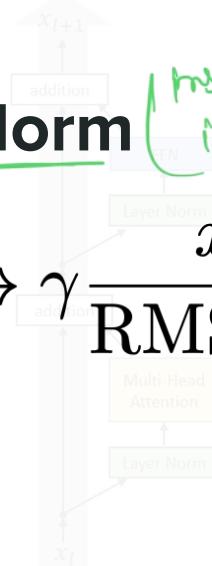
$$\text{Output} = x + \text{SubLayer}(\text{LayerNorm}(x))$$

Pre-Norm +

RMSNorm

After
residual connection

$$\gamma \hat{x} + \beta \longrightarrow \gamma \frac{x}{\text{RMS}(x)}$$



Instead of writing layernorm;
it is also useful to use for
benefits of fewer parameters to
learn; In layer norm,
we learn γ, β but here
we have γ to learn



Transformers & Large Language Models

Position embeddings

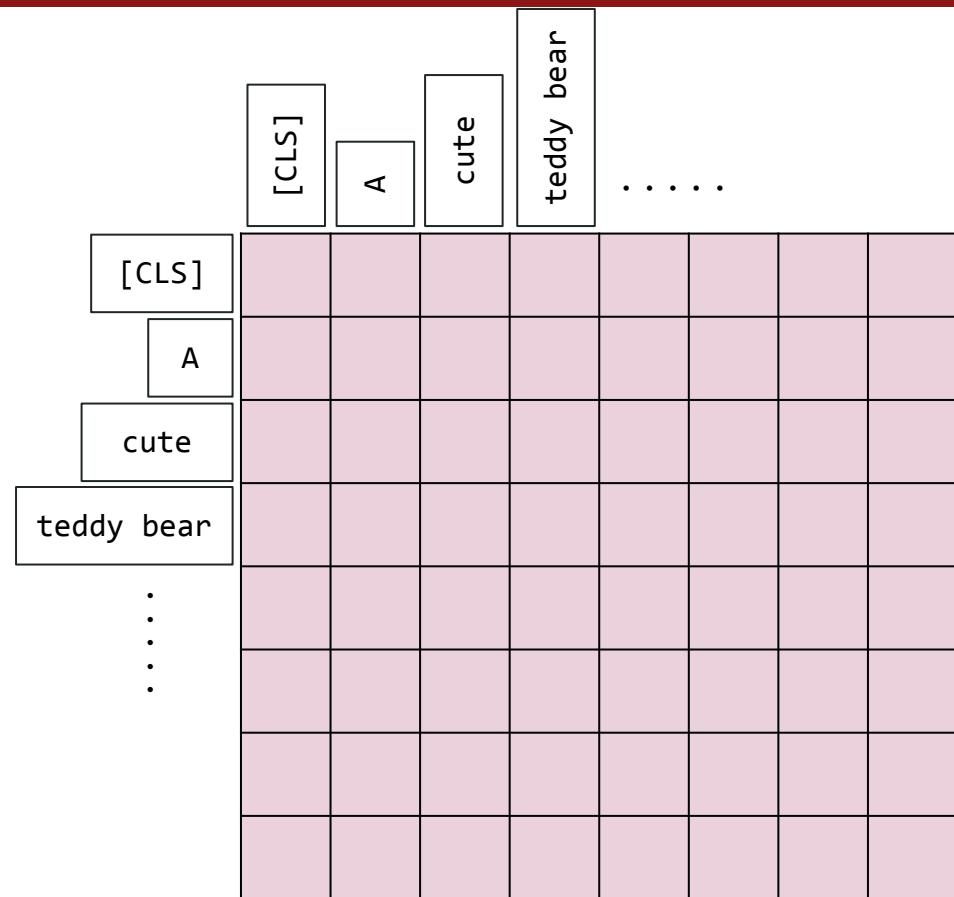
Layer normalization

Attention approximation

Transformer-based models

BERT deep dive

Sparse attention: Longformer



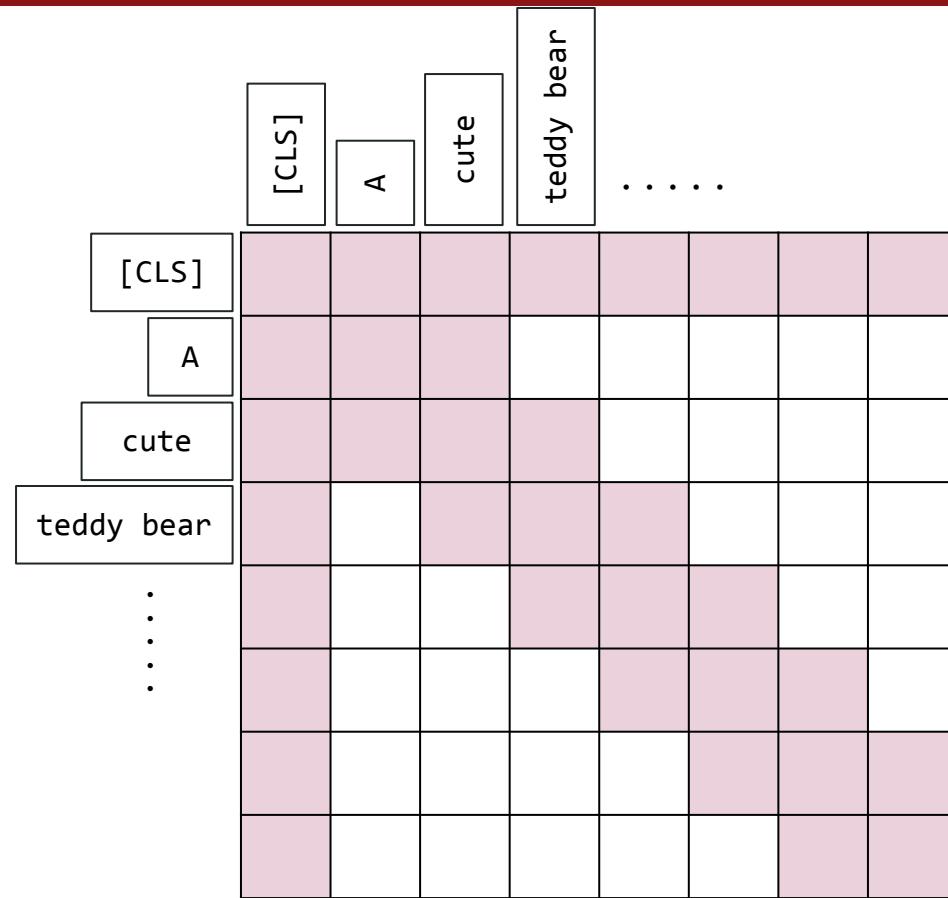
if when each embedding will attend each embedding for attention mechanism calculation then basically it becomes n^2 calculations
So, to get rid of this much of calculation,
the attention can be approximated.
using Longformer

it restricts the window of interact
it only allow neighbour embedding
to include into attention mechanism
and remaining approximated
using sliding window Attention
mechanism.

Sparse attention: Longformer

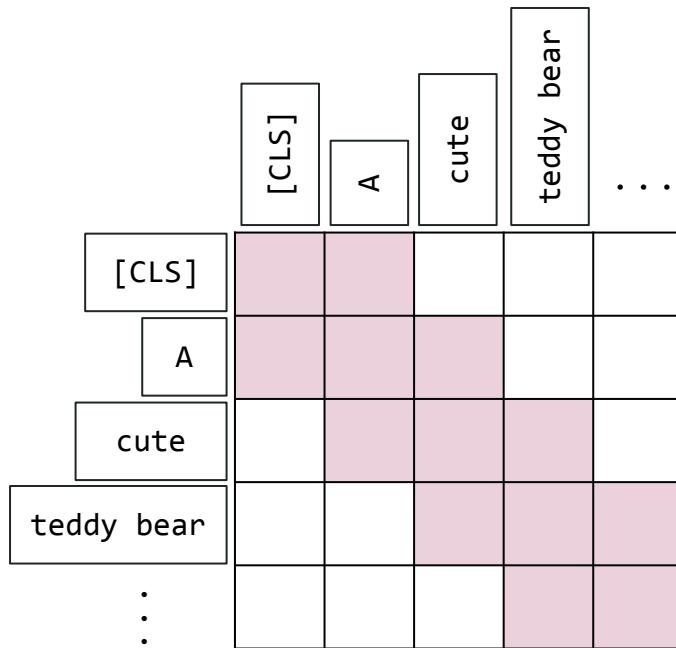
so basically ^{!}} sliding window

Attention mechanism gives the window size of embeddings to be attended. [Like receptive field in computer vision]



Leveraging local and global attention

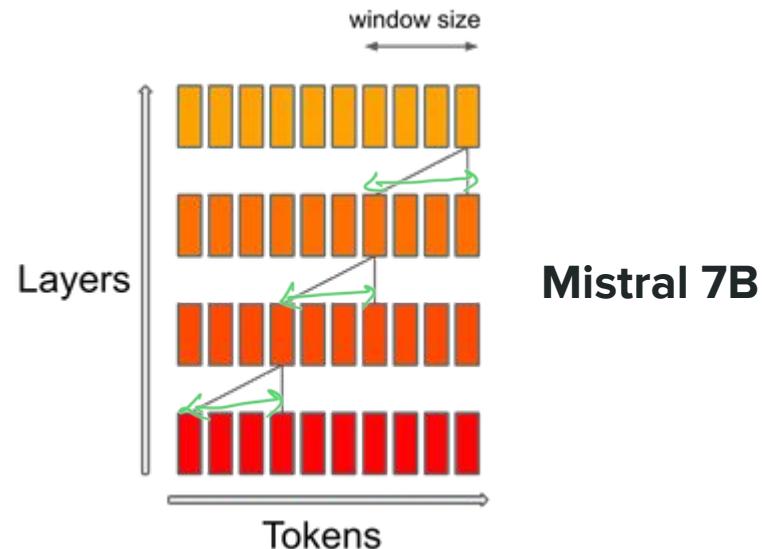
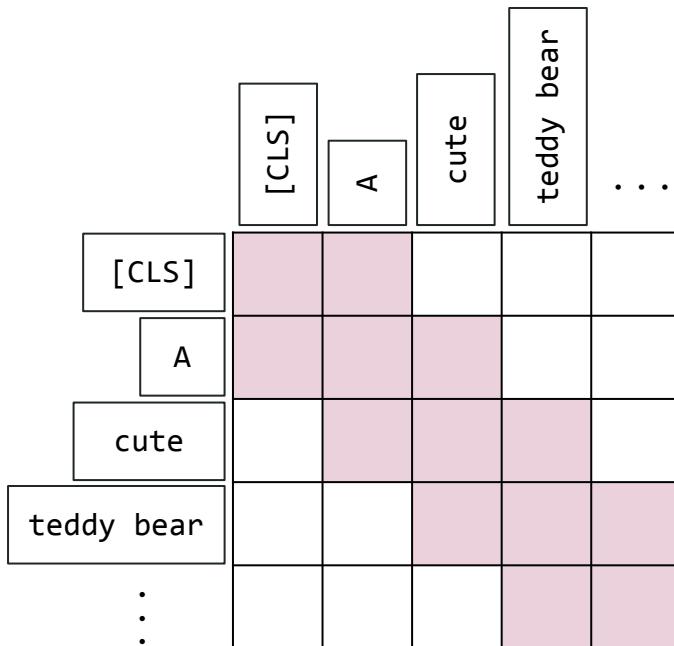
SWA = Sliding Window Attention



Variations include interleaving local and global attention layers

Leveraging local and global attention

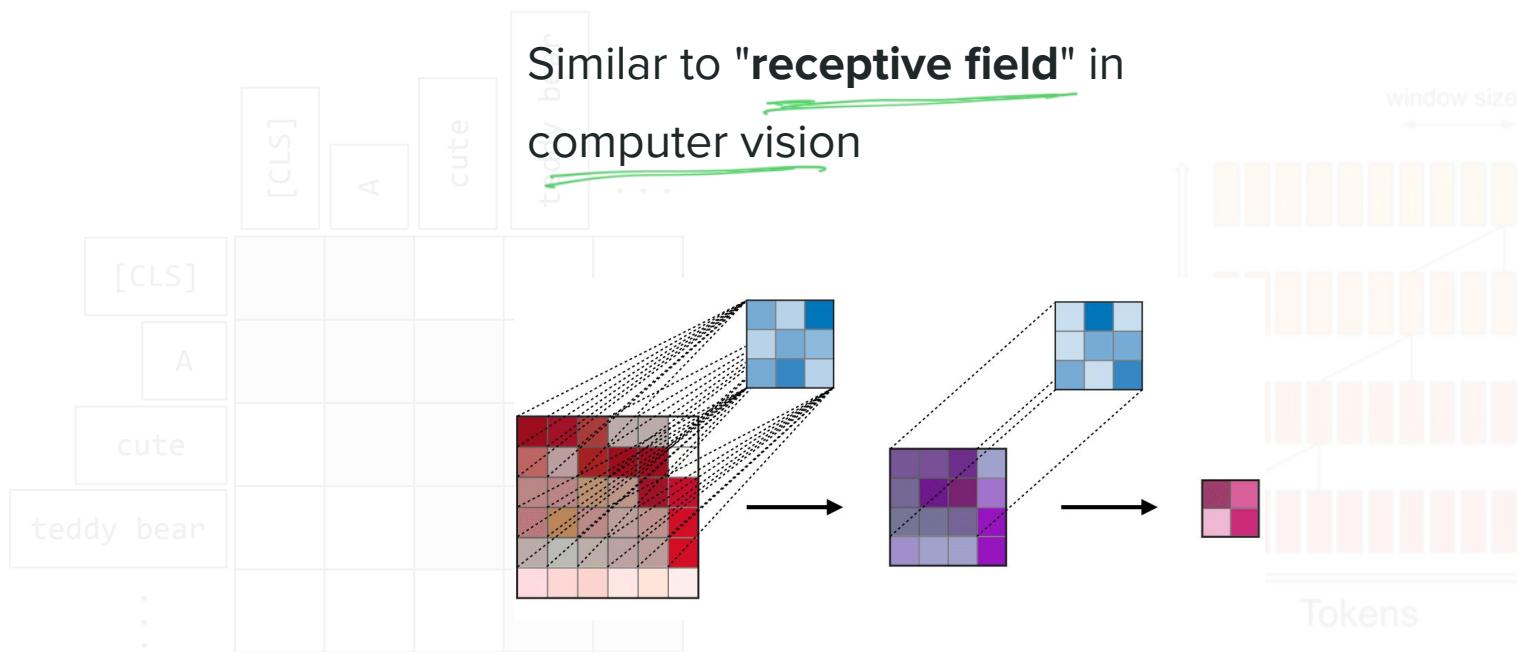
SWA = Sliding Window Attention



Mistral 7B

Leveraging local and global attention

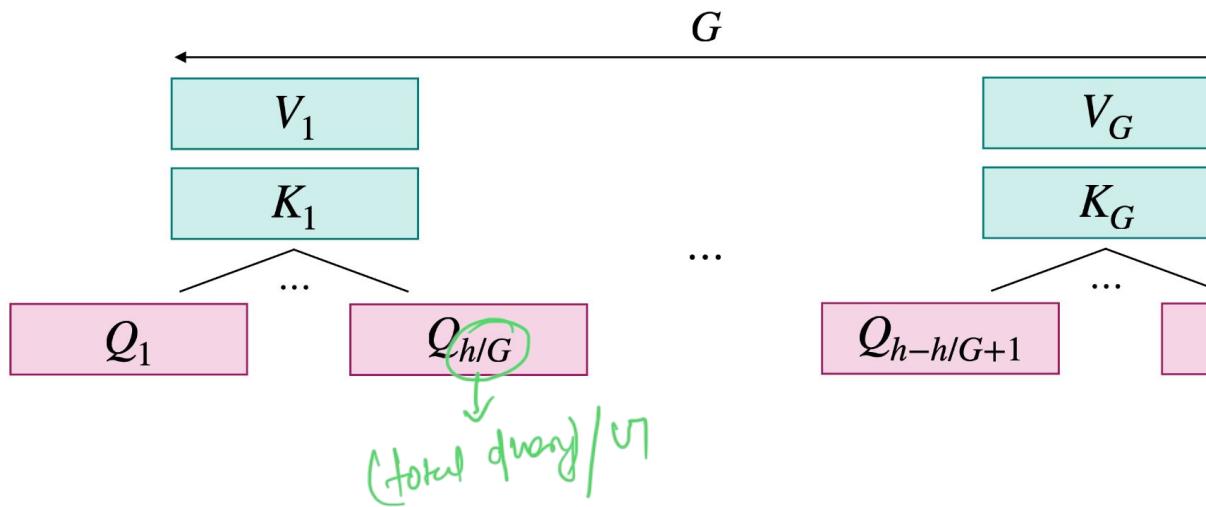
SWA = Sliding Window Attention



Mistral 7B

Sharing attention heads

Idea. Share key/value attention heads within groups of queries



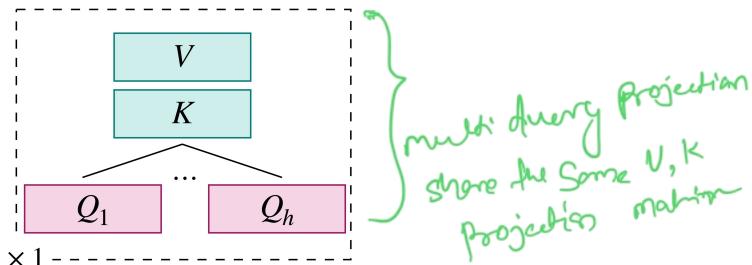
Sharing attention heads

of query.

$$G = 1$$

Multi-Query Attention (MQA)

All h head share all the
some projection matrix for V, K



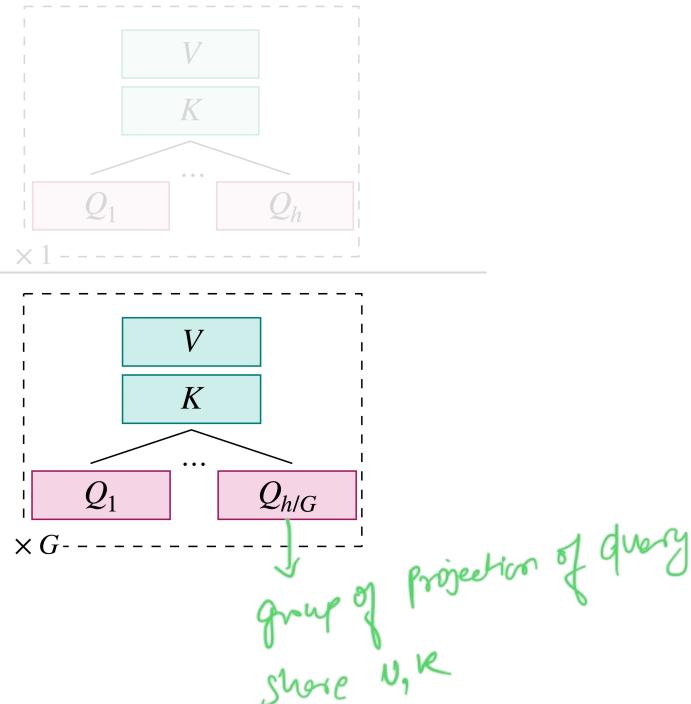
Sharing attention heads

$G = 1$

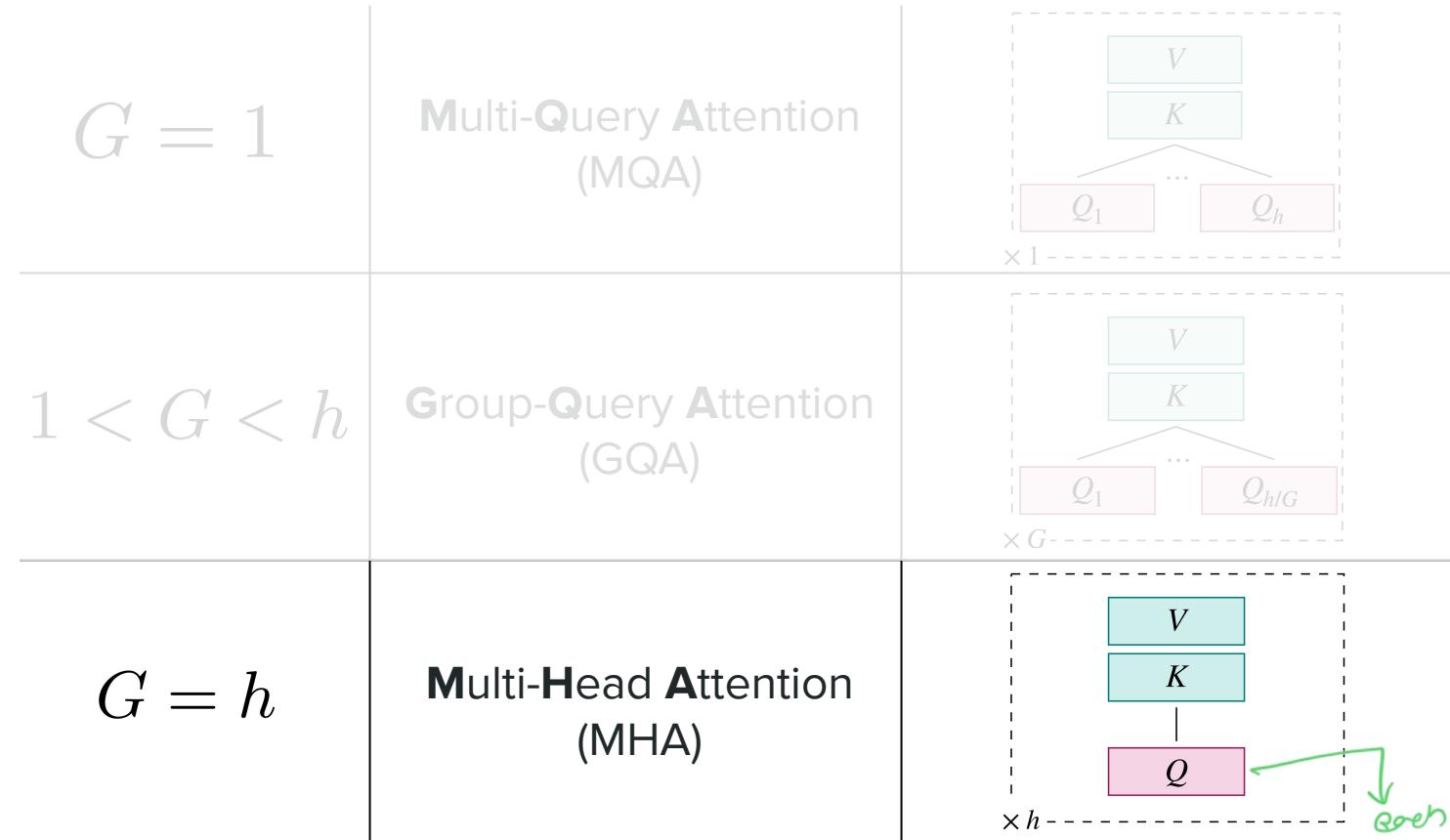
**Multi-Query Attention
(MQA)**

$1 < G < h$

**Group-Query Attention
(GQA)**



Sharing attention heads



each query projection
share all V, K



Transformers & Large Language Models

Attention approximation

Position embeddings

Layer normalization

Transformer-based models

BERT deep dive

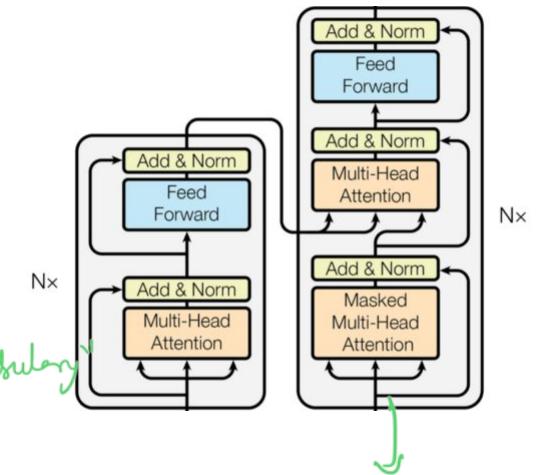
Overview of Transformer-based models

3 categories of Transformer-based models:

Encoder-decoder	Text to text	T5, mT5, ByT5
-----------------	--------------	---------------

in short, it processes the encoder with input text with blank spaces where decoder tries to predict the missing tokens.

vanilla paper
introduced "Span Corruption task"
tokenizer free model but byte level
"less size of vocabulary"
Instead of putting all the sentences if put some blanks which misses multiple tokens.
(Known as Control tokens)



it finds each of Span Corruption token (which was blank in the input text)

Overview of Transformer-based models

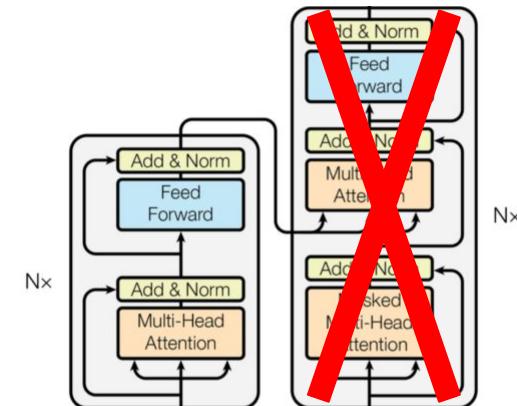
it decodes from many heads

3 categories of Transformer-based models:

Encoder-decoder	Text to text	T5, mT5, ByT5
<u>Encoder-only</u>	Projection of embedding for class prediction (e.g. sentiment extraction)	BERT, DistilBERT, RoBERTa

Classification task
i.e. it uses only Encoder
parts for transformer

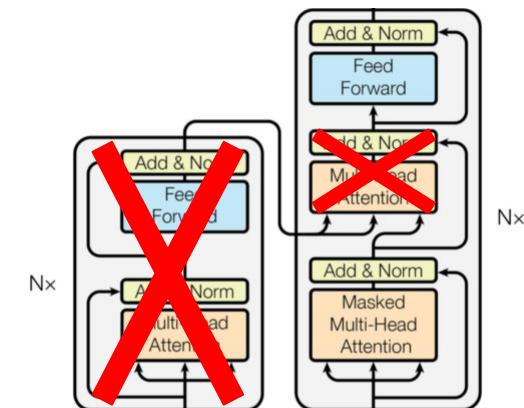
only Encoder parts being
used for classification task



Overview of Transformer-based models

3 categories of Transformer-based models:

Encoder-decoder	Text to text	T5, mT5, ByT5
Encoder-only	Projection of embedding for class prediction (e.g. sentiment extraction)	BERT, DistilBERT, RoBERTa
<u>Decoder-only</u>	Text to text	GPT series



Overview of Transformer-based models

3 categories of Transformer-based models:

Encoder-decoder	Text to text	T5, mT5, ByT5
Encoder-only	Projection of embedding for class prediction (e.g. sentiment extraction)	BERT, DistilBERT, RoBERTa
Decoder-only	Text to text	GPT series

**Popular in
~2018-2022**

```
graph LR; A[Encoder-decoder] --> C[Popular in ~2018-2022]; B[Encoder-only] --> C;
```

Overview of Transformer-based models

3 categories of Transformer-based models:

Encoder-decoder	Text to text	T5, mT5, ByT5
Encoder-only	Projection of embedding for class prediction (e.g. sentiment extraction)	BERT, DistilBERT, RoBERTa
Decoder-only	Text to text	GPT series

← **Popular now!**



Transformers & Large Language Models

Attention approximation

Position embeddings

Layer normalization

Transformer-based models

BERT deep dive

Why "BERT"?

BERT = Bidirectional Encoder Representations from Transformers

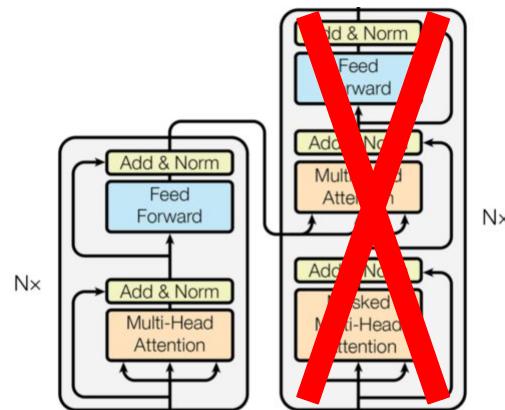
Encoder only architecture
used for classification tasks

Why "BERT"?

BERT = Bidirectional Encoder Representations from Transformers

Why "BERT"?

BERT = Bidirectional Encoder Representations from Transformers



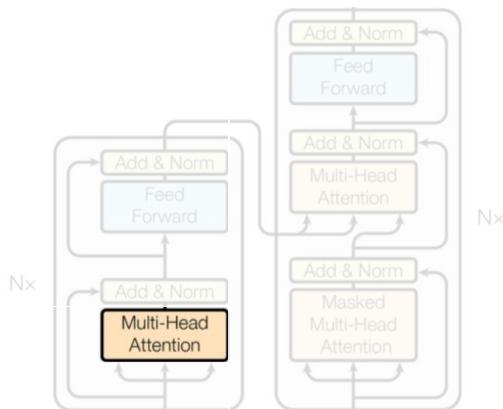
Why "BERT"?

BERT = Bidirectional Encoder Representations from Transformers

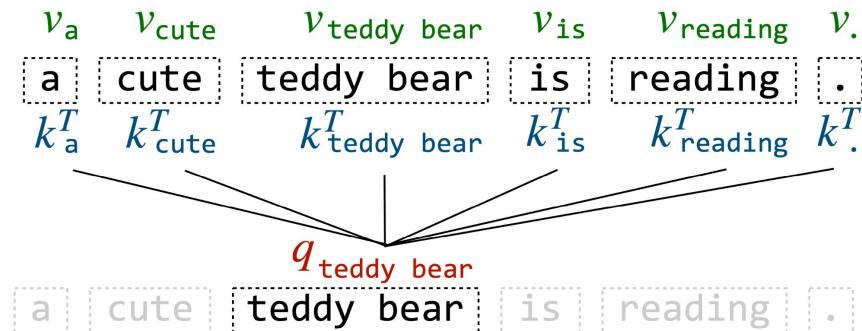
Why "BERT"?

BERT = Bidirectional Encoder Representations from Transformers

\downarrow
Bidirectional processing in BERT
architecture is because all tokens
attend all tokens for the representations

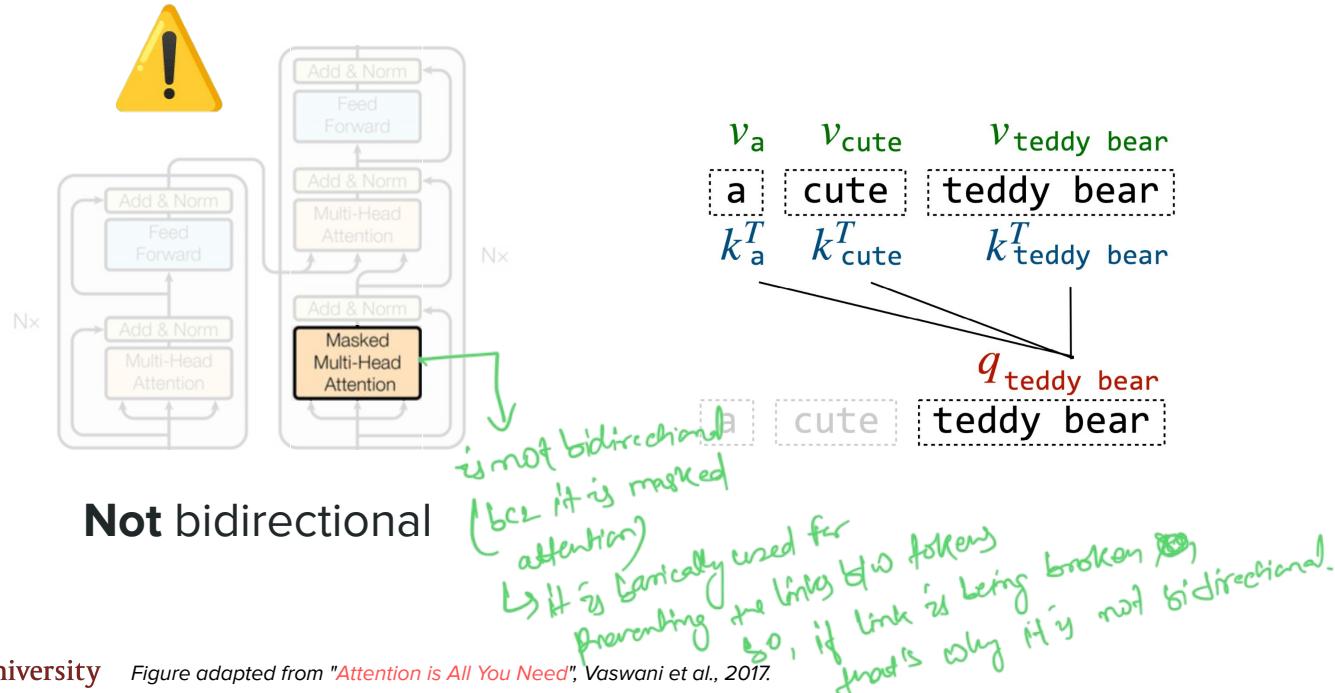


Bidirectional



Why "BERT"?

BERT = Bidirectional Encoder Representations from Transformers



Why "BERT"?

Deep contextualized word representations

Matthew E. Peters[†], Mark Neumann[†], Mohit Iyyer[†], Matt Gardner[†],
`{matthewp, markn, mohiti, mattg}@allenai.org`

Christopher Clark*, Kenton Lee*, Luke Zettlemoyer^{†*}
`{csquared, kentonl, lsz}@cs.washington.edu`

[†]Allen Institute for Artificial Intelligence

*Paul G. Allen School of Computer Science & Engineering, University of Washington

Abstract

We introduce a new type of *deep contextualized* word representation that models both (1) complex characteristics of word use (e.g., syn-

guage model (LM) objective on a large text corpus. For this reason, we call them ELMo (Embeddings from Language Models) representations. Unlike previous approaches for learning contextu-

ELMo paper
[submitted in Feb 2018]

BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

Jacob Devlin Ming-Wei Chang Kenton Lee Kristina Toutanova
Google AI Language
`{jacobdevlin, mingweichang, kentonl, kristout}@google.com`

Abstract

We introduce a new language representation model called BERT, which stands for Bidirectional Encoder Representations from Transformers. Unlike recent language representation models (Peters et al. 2018a; Rad-

There are two existing strategies for applying pre-trained language representations to downstream tasks: *feature-based* and *fine-tuning*. The feature-based approach, such as ELMo (Peters et al., 2018a), uses task-specific architectures that include the pre-trained representations as addi-

BERT paper
[submitted in Oct 2018]

Why "BERT"?



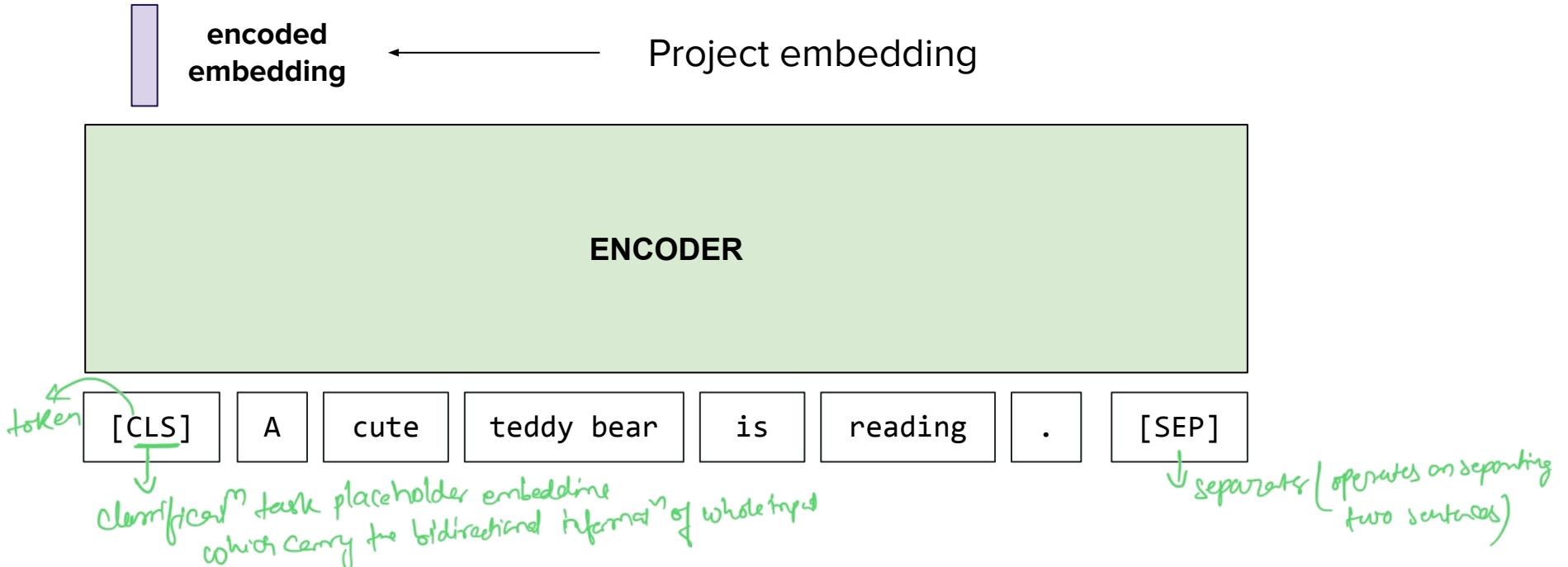
Elmo



Bert

BERT and encoder-only models

Idea. Keep only **encoders** + leverage **encoded embeddings** to make predictions.



BERT and encoder-only models

Strategy.

- **Step 1:** Pretraining with proxy tasks (MLM and NSP)
- **Step 2:** Finetuning for given end task

↑
Training from scratch

→ marked language model
(to learn the internal structure of texts)

↓
using pretrained model for
target task.

→ next Sentence prediction
(seen a way to make ordering
of words/sentences make sens!)

BERT and encoder-only models

Strategy.

- **Step 1:** Pretraining with proxy tasks (MLM and NSP)
- **Step 2:** Finetuning for given end task

Discussion.

Pros	Cons
<ul style="list-style-type: none">● Finetuning does not need a lot of data● Good performance	<ul style="list-style-type: none">● Not suited for a range of tasks (e.g. text generation)● Finetuning is a required step

BERT and encoder-only models

Strategy.

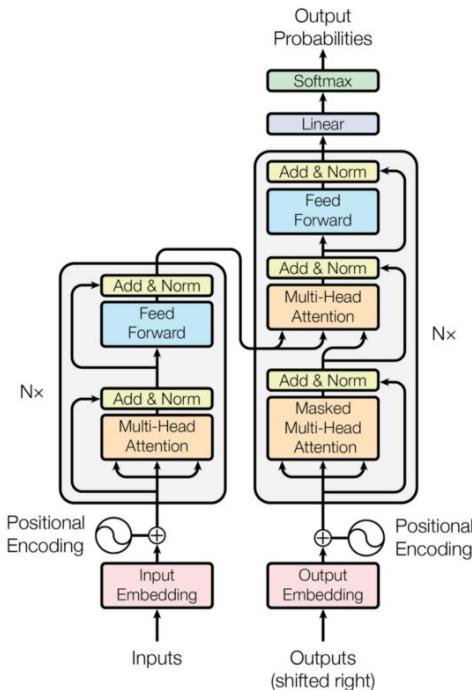
- **Step 1:** Pretraining with proxy tasks (MLM and NSP)
- **Step 2:** Finetuning for given end task

Discussion.

Pros	Cons
<ul style="list-style-type: none">● Finetuning does not need a lot of data● Good performance	<ul style="list-style-type: none">● Not suited for a range of tasks (e.g. text generation)● Finetuning is a required step

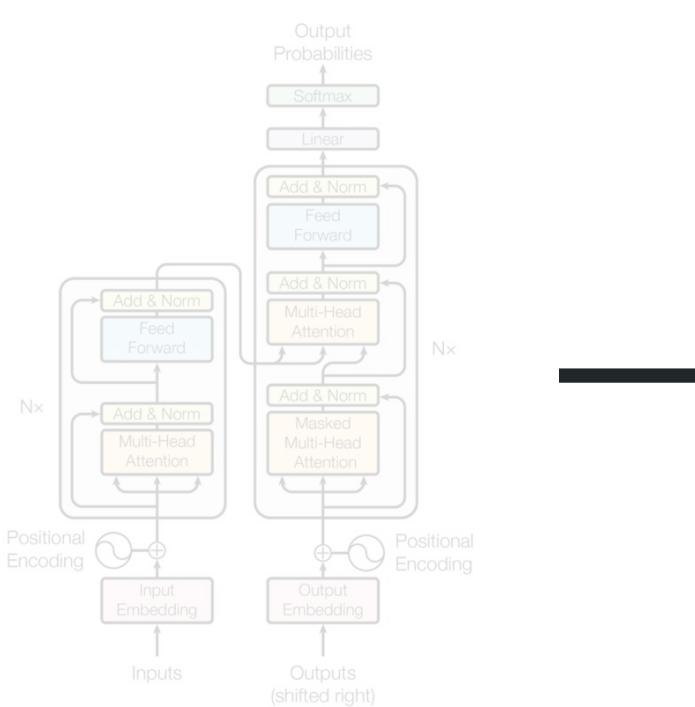
Variants. RoBERTa, DistilBERT, ALBERT

BERT architecture

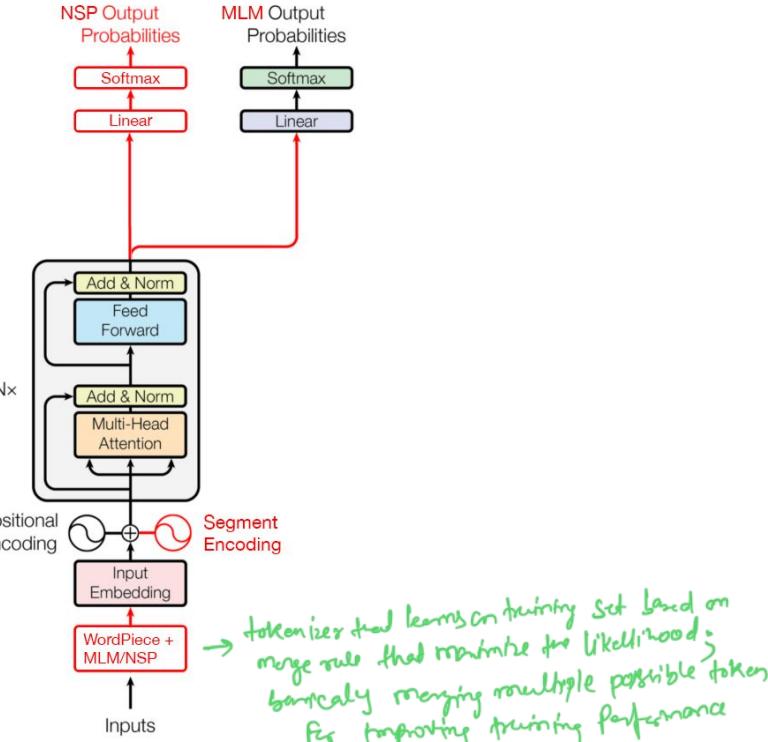


Original transformer (2017)

BERT architecture



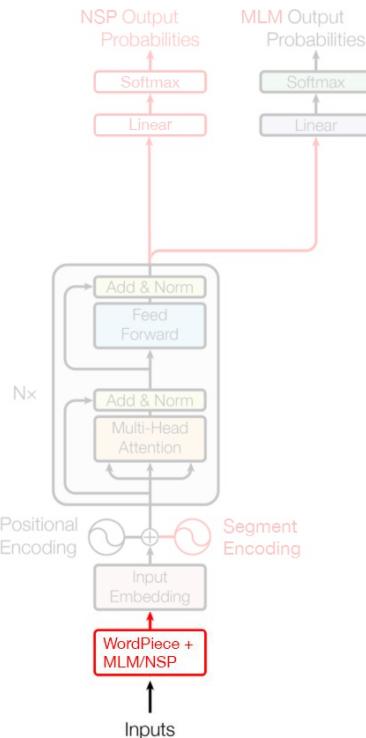
Original transformer (2017)



BERT (2018)

→ tokenizer that learns on training set based on merge rule that maintains the likelihood; basically merging multiple possible tokens for improving training performance

Input processing



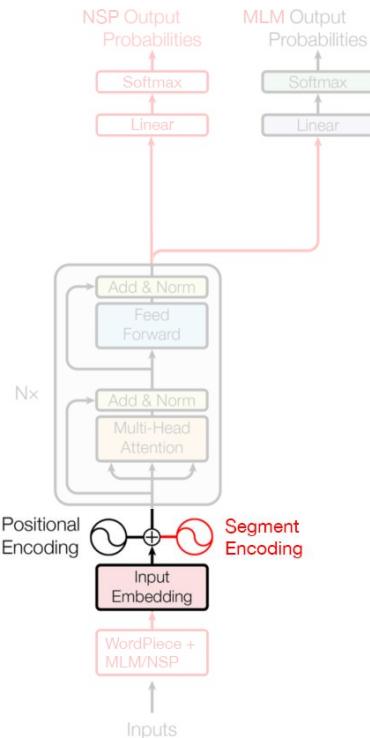
WordPiece algorithm

- Tokenizer trained on a training set beforehand
- Vocabulary size: ~30,000
- Great at detecting common particles

NSP / MLM task processing

- Add [CLS] token at the beginning of the input
- Separate consecutive segments with the [SEP] token and put another one at the end
- Use [MASK] to mask inputs
 - ↳ masking for token for predicting something meaningful new

Input embedding



Input embeddings

- Gigantic lookup table
- Learns an embedding for each word of the vocabulary

Positional encoding

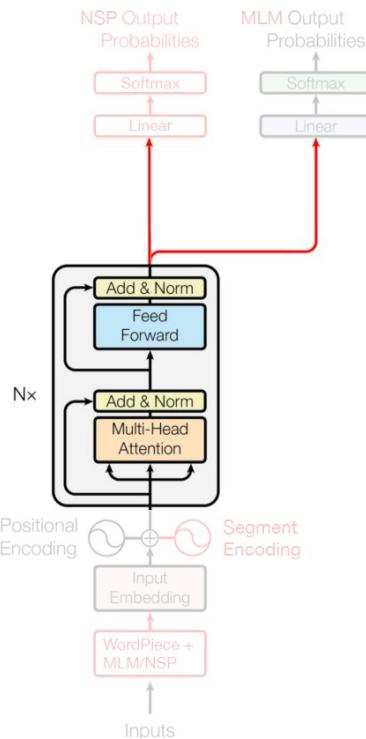
- Helps the network associate tokens with a position
- Encoding either learnt or fixed with cosines and sines

(new!) Segment encoding

- Shared embedding for a segment

*— used for helping NSP task
(so basically learned embedding from Segment A added with learned embedding from Segment B) for learning combined better representation*

Encoder-only model



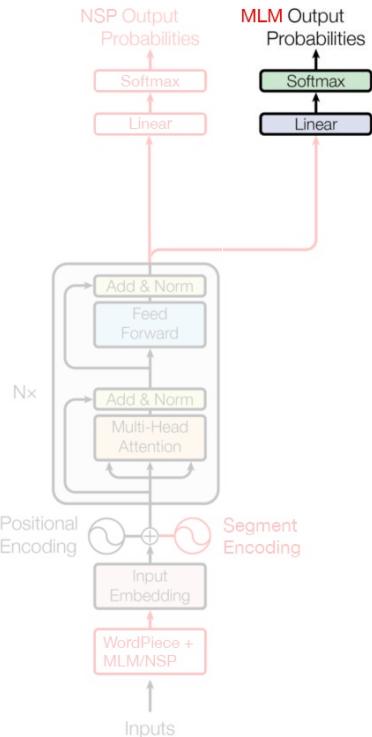
Model

Encoder part of the original transformers paper

Goal

- Represent input data with features (hopefully) needed for NLP tasks
- Leverage the Transformer's self-attention mechanism
- Use learned embedding towards classification-oriented tasks

Proxy task: Masked Language Modeling - *forcing model to learn the context.*



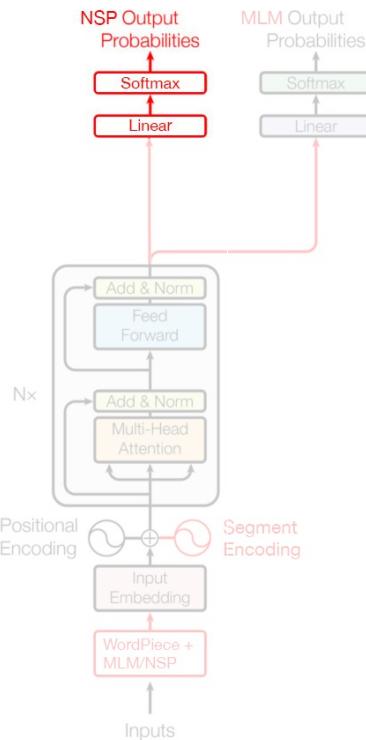
Idea. 15% of input tokens are set up for prediction where:

- 80% are masked *special token replaced with certain token/word in a sentence.*
- 10% are changed to a random word
- 10% are unchanged

Benefits

- Network learns language modeling based on contextual information
- Regularization reflects probabilistic nature of language

Proxy task: Next Sentence Prediction



Idea. Pick two sentences from the corpus, where:

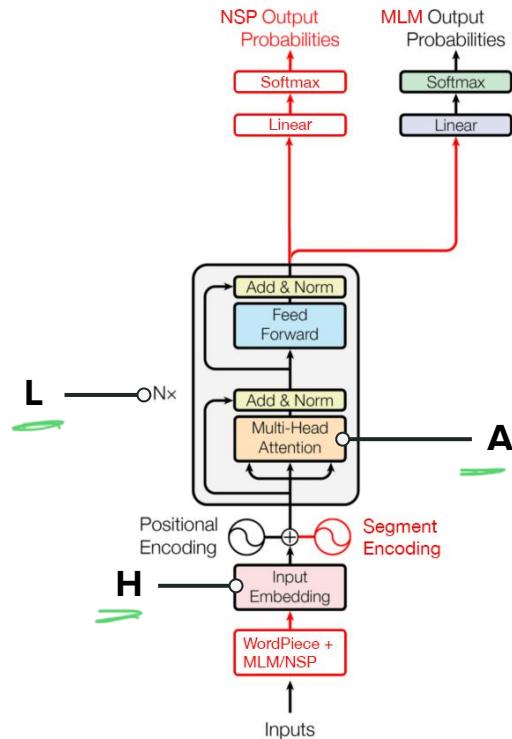
- 50% of the time, they follow each other
- 50% of the time, they **do not** follow each other

Task. Predict if they actually follow each other.

Benefits

- Network implicitly learns to detect useful contextual information
- Easy classification task that does not require any labels

Hyperparameters



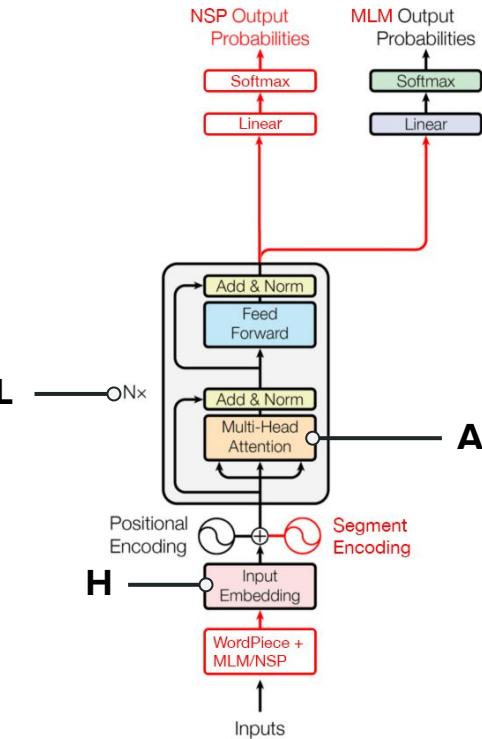
Model.

- **L** Layers.
- **H** Hidden layer size aka dimension of embeddings.
- **A** Attention heads operating in parallel.

Data.

- Language-specific / multilingual. Languages on which the model has been trained on.
- Cased / uncased. Whether inputs are converted to lowercase or not.

Some numbers



BERT

	L	H	A	Parameters
BERT-Tiny	2	128	2	4M
BERT-Mini	4	256	4	11M
BERT-Small	4	512	8	30M
BERT-Medium	8	512	8	42M
BERT-Base	12	768	12	110M
BERT-Large	24	1024	16	340M

BERT

Overview of BERT finetuning

Goal. Leverage embeddings learned by BERT for a “sister” task

Overview of BERT finetuning

Goal. Leverage embeddings learned by BERT for a “sister” task

Tricks

- Use weights from **already massively pretrained** model
- **Freezing** early layers: sometimes better trade-off complexity/performance
- Great results possible with **minimal labeled data** (depending on complexity/proximity to pretrained data distribution + objectives)

Overview of BERT finetuning

Goal. Leverage embeddings learned by BERT for a “sister” task

Tricks

- Use weights from **already massively pretrained** model
- **Freezing** early layers: sometimes better trade-off complexity/performance
- Great results possible with **minimal labeled data** (depending on complexity/proximity to pretrained data distribution + objectives)

Use cases

- Sequence classification: e.g. sentiment extraction
- Token classification: e.g. question answering

Finetuning: sentiment extraction

This teddy bear is SO CUTE!

Finetuning: sentiment extraction

↓ put everything in lower case

this teddy bear is so cute!

Finetuning: sentiment extraction

↓ tokenization mechanism

this teddy bear is so cute !

Finetuning: sentiment extraction

Add [CLS] token as a placeholder for sentiment.

↓ adding [CLS] token for classifier?

[CLS] this teddy bear is so cute !

Finetuning: sentiment extraction



Finetuning: sentiment extraction



Finetuning: sentiment extraction



position embedding



embedding

[CLS]

0

this

1

teddy

2

bear

3

is

4

so

5

cute

6

!

7

[SEP]

8

[PAD]

9

[PAD]

10

[PAD]

11

[PAD]

12

[PAD]

13

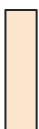
[PAD]

14

Finetuning: sentiment extraction



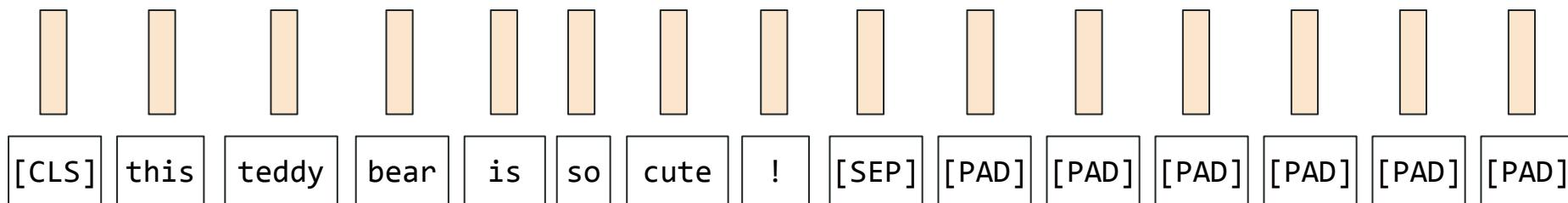
Finetuning: sentiment extraction



position- and segment-aware embedding



Finetuning: sentiment extraction

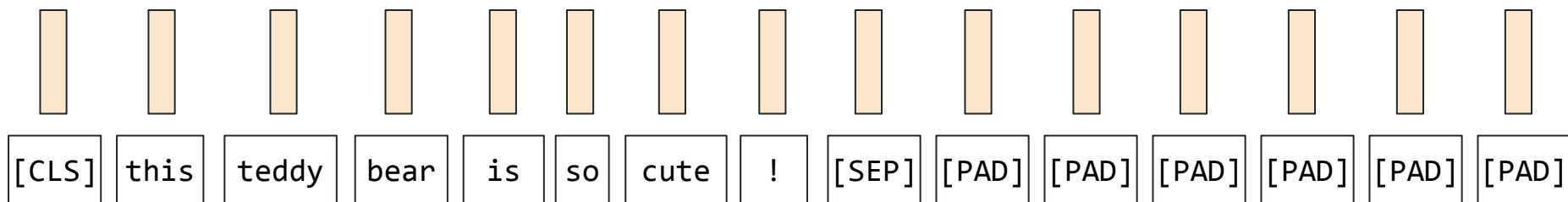


Finetuning: sentiment extraction

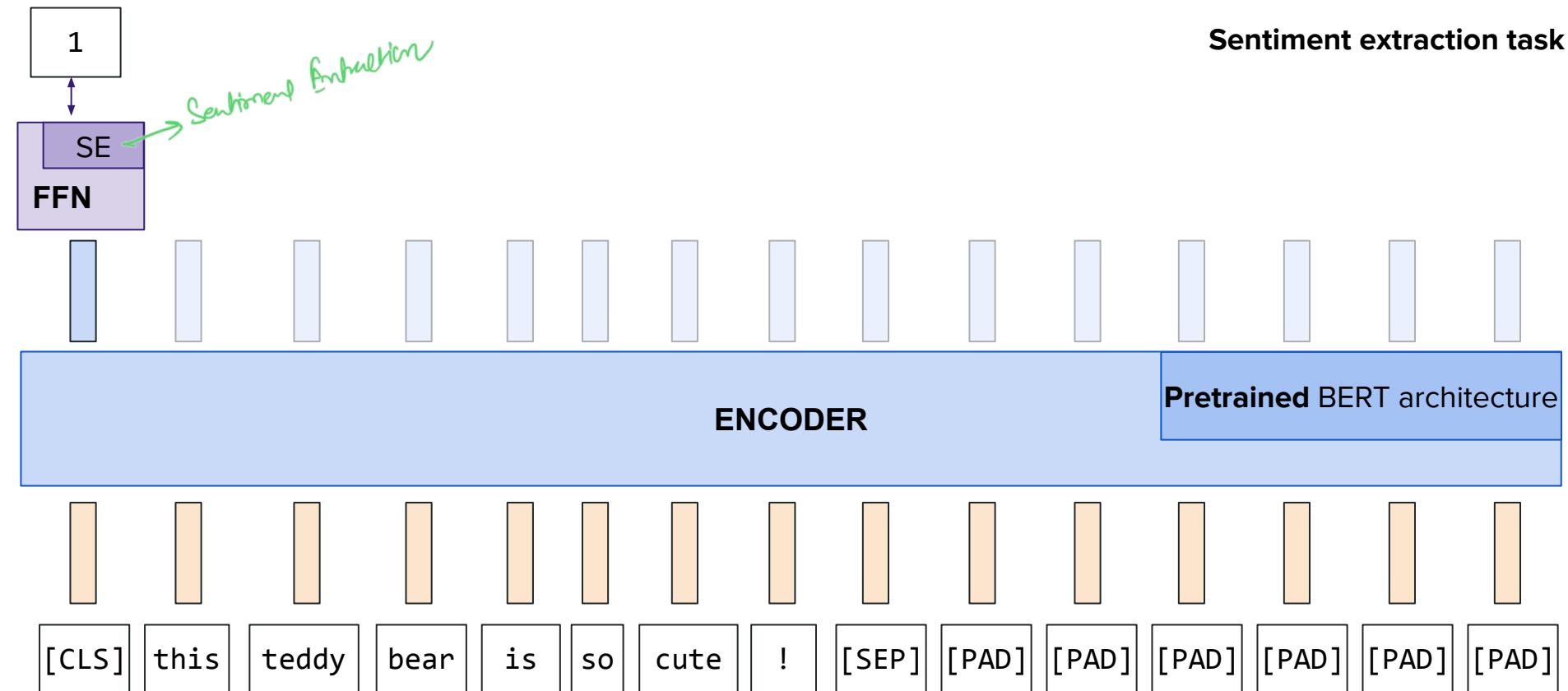
taking the pretrained
BERT architecture

ENCODER

Pretrained BERT architecture



Finetuning: sentiment extraction



Takeaways and shortcomings

Benefits.

- State of the art results
- ~True contextual representation of words
- Adaptable to many classification tasks

Takeaways and shortcomings

Benefits.

- State of the art results
- ~True contextual representation of words
- Adaptable to many classification tasks

Applications. Widely used in the industry for anything related to encoding

Takeaways and shortcomings

Benefits.

- State of the art results
- ~True contextual representation of words
- Adaptable to many classification tasks

Applications. Widely used in the industry for anything related to encoding

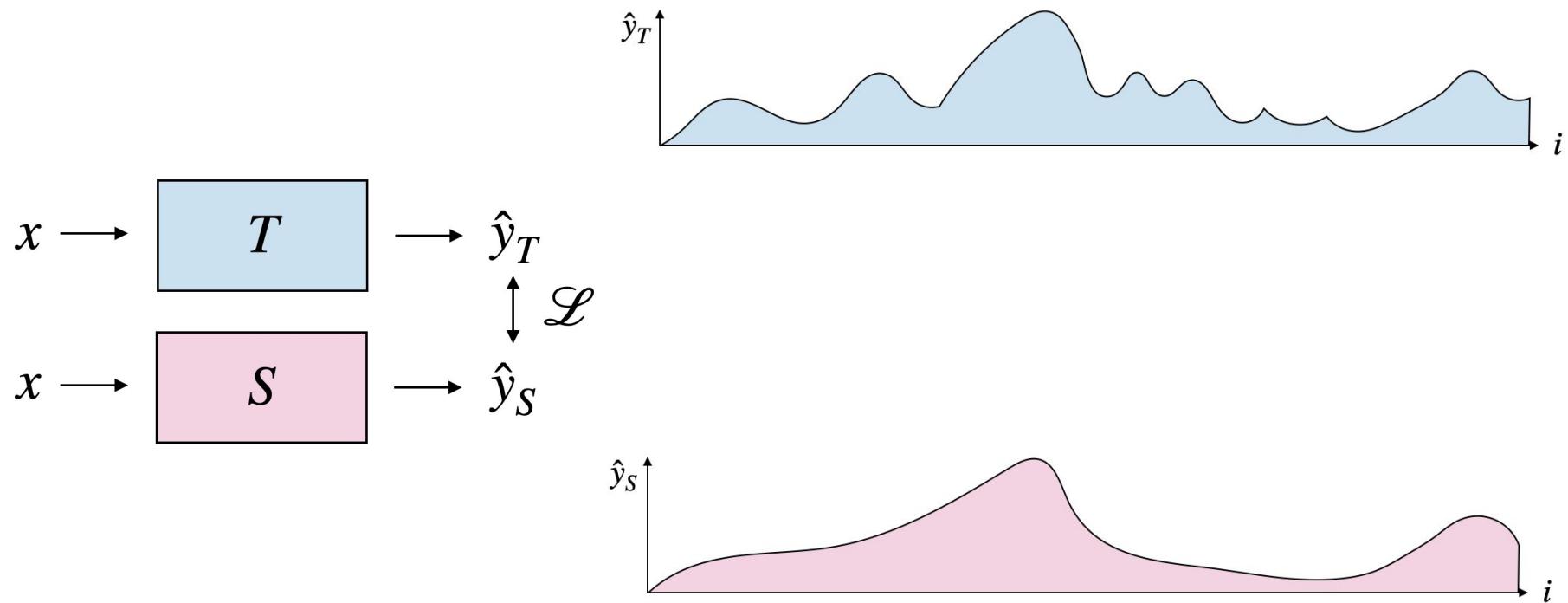
Limitations.

- Context window size is limited
- Computationally expensive: hard sell for low-latency/cost-sensitive applications
- Training paradigm is complex: MLM/NSP + finetuning

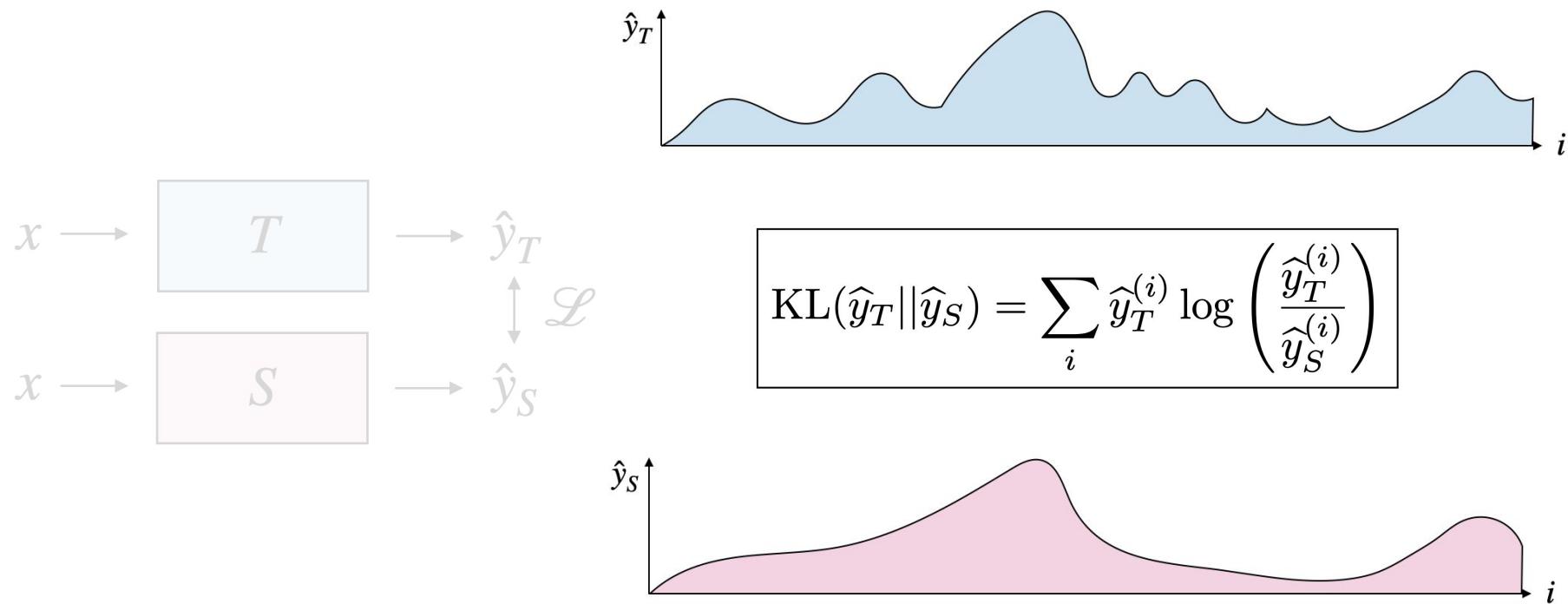
Trick: distillation

"The **soft targets** contain almost all the **knowledge**."

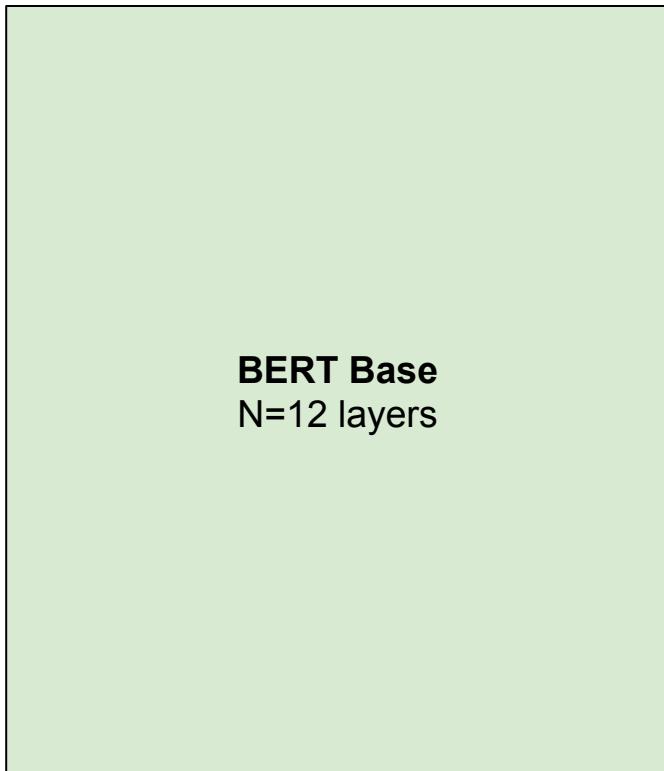
Trick: distillation



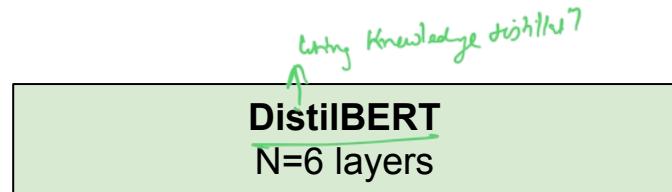
Trick: distillation



Variant for efficiency: DistilBERT



~1.6x faster



~97% performance

Variant for performance: RoBERTa

Goal. Thorough analysis of directions of optimization

Variant for performance: RoBERTa

Goal. Thorough analysis of directions of optimization

Modeling

- Removing NSP/segment encodings: ~no effect!*
- Static → dynamic masking across epochs *Dynamic masking dynamically*

*DistilBERT had also removed it

Variant for performance: RoBERTa

Goal. Thorough analysis of directions of optimization

Modeling

- Removing NSP/segment encodings: ~no effect!
- Static → dynamic masking across epochs

Data

- Richer: 16 GB → 160 GB of pretraining corpus size
- For longer: 1M steps of batch size 256 vs. 500k of batch size 8k

Variant for performance: RoBERTa

Goal. Thorough analysis of directions of optimization

Modeling

- Removing NSP/segment encodings: ~no effect!
- Static → dynamic masking across epochs

Data

- Richer: 16 GB → 160 GB of pretraining corpus size
- For longer: 1M steps of batch size 256 vs. 500k of batch size 8k

Result. +4% across benchmarks with same architecture

Thank you for your attention!
