

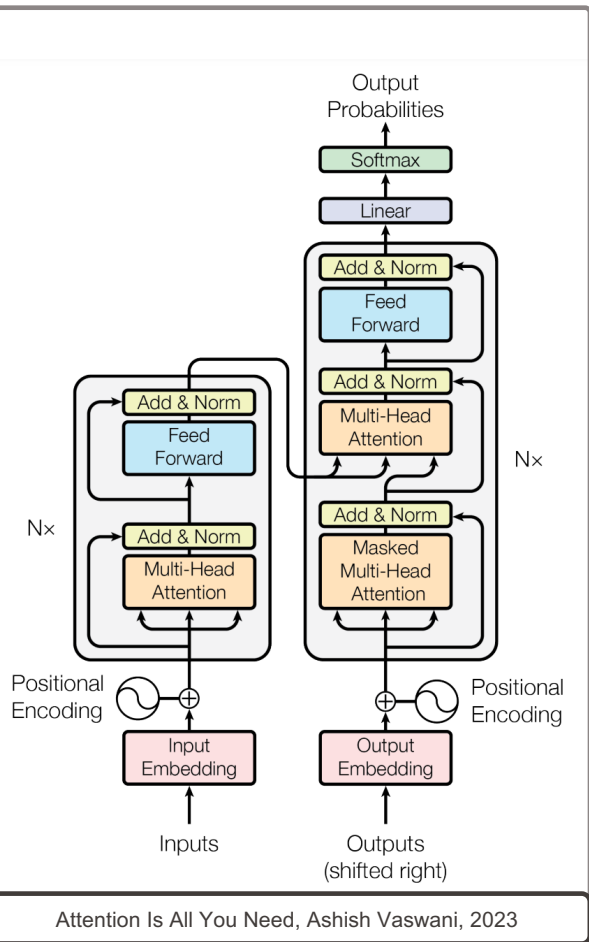


Exercise 5

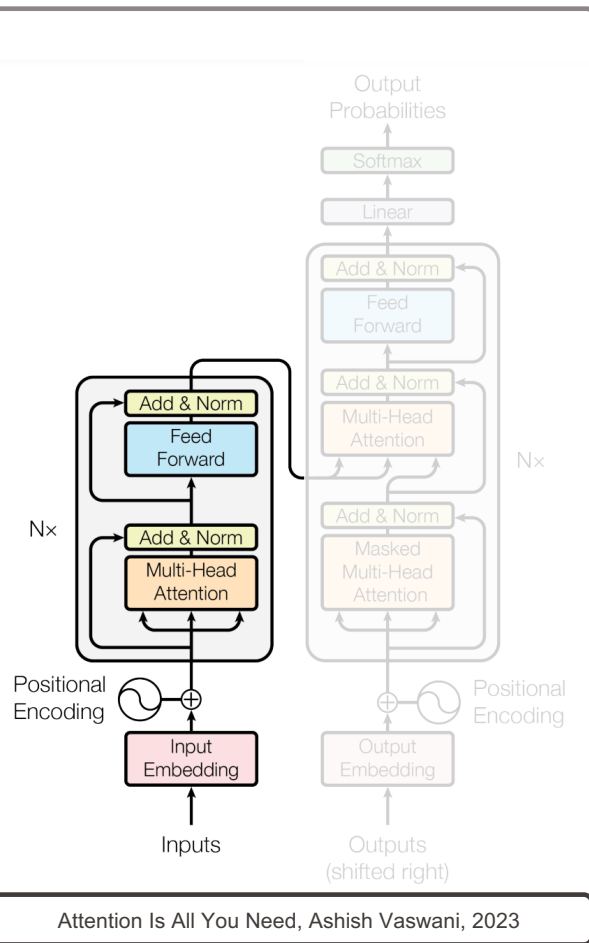
Transformers

EPFL – 20.10.2023

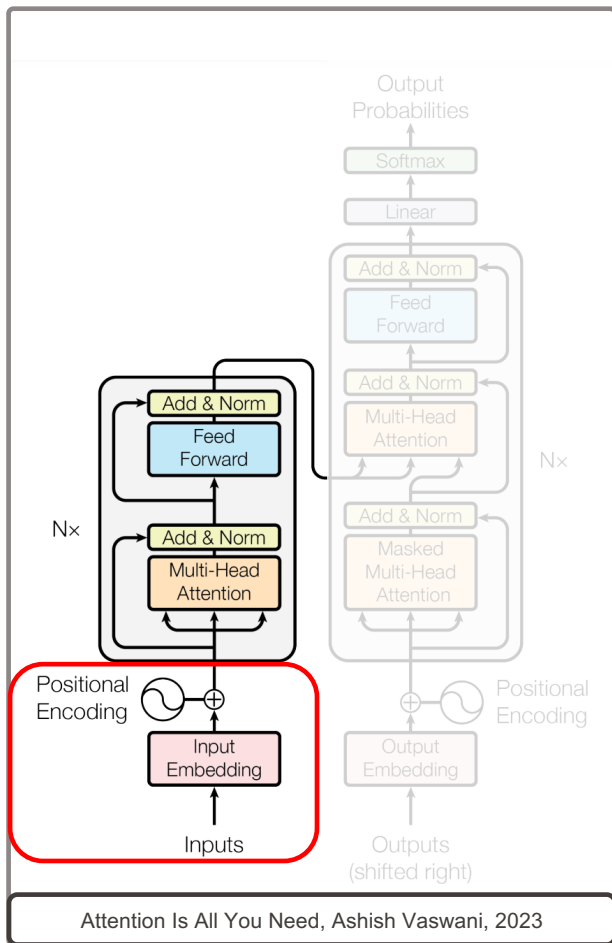
Transformer - Recap



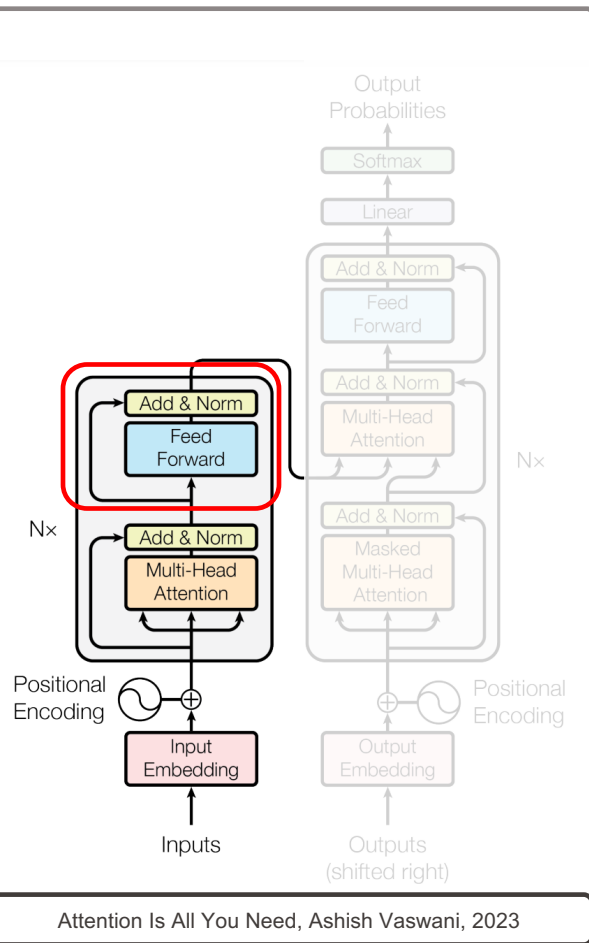
Transformer - Recap



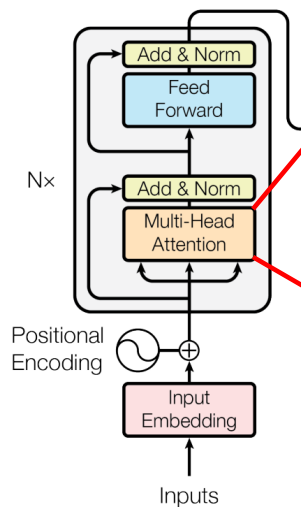
Transformer - Recap



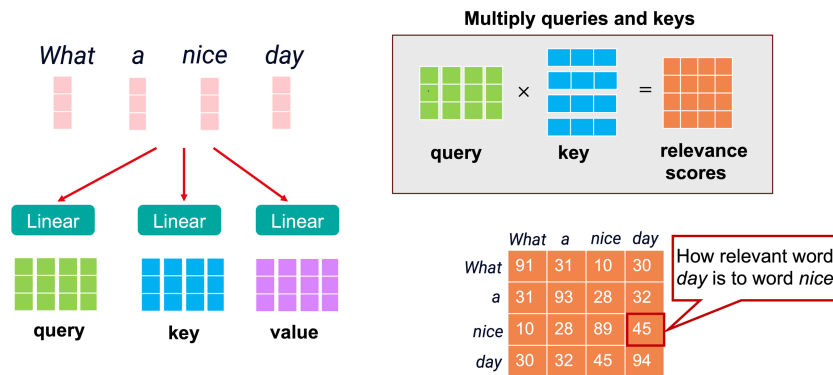
Transformer - Recap

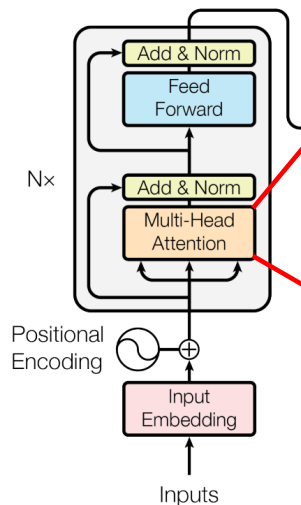


Transformer - Recap

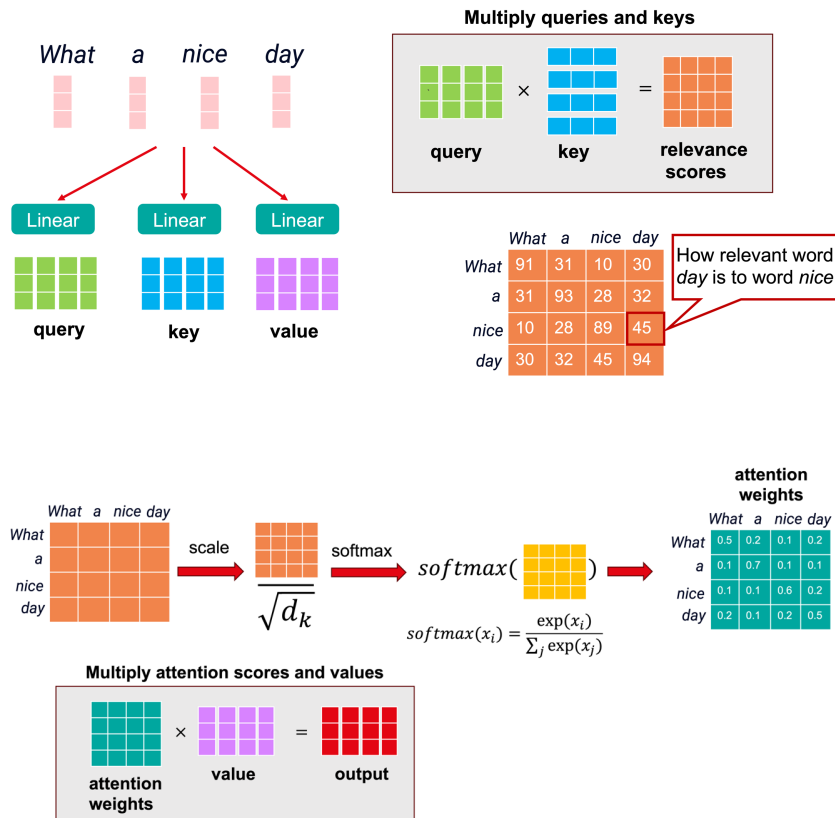


Attention Is All You Need, Ashish Vaswani, 2023



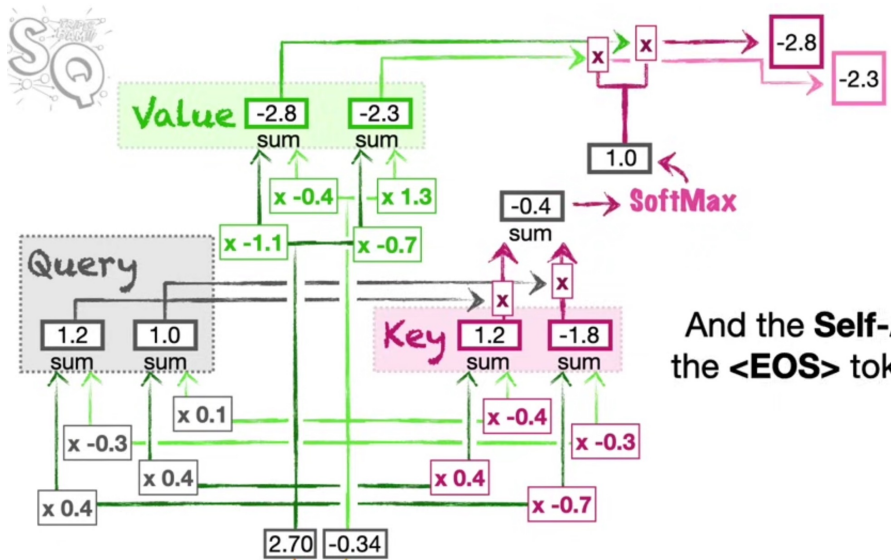


Attention Is All You Need, Ashish Vaswani, 2023



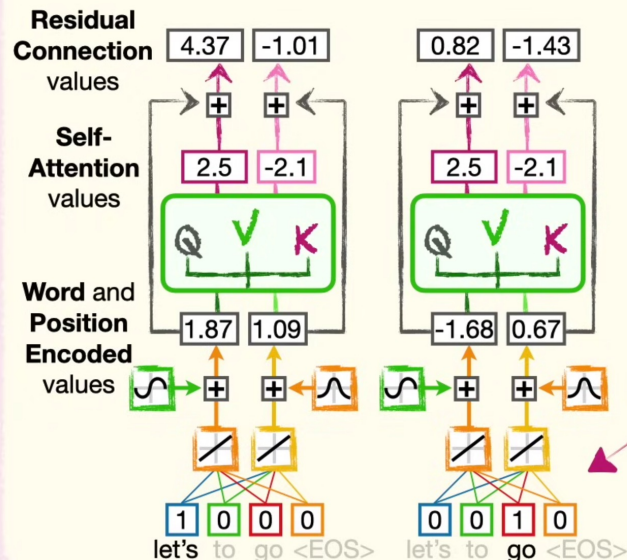
Tutorial Suggestion:

<https://www.youtube.com/watch?v=zxQyTK8quyY&t=1087s>



And the **Self-**
the **<EOS>** tok

Encoder



Quadratic dependency on Sequence Length for
Memory and Computation

	What	a	nice	day
What				
a				
nice				
day				

Quadratic dependency on Sequence Length for Memory and Computation

	What	a	nice	day
What				
a				
nice				
day				

Sparse Attention

- use fixed pattern of attention to reduce number of connections

Flash Attention

- Split score matrix into blocks
- Recompute attention matrix

Longformer

- local attention for most tokens
- global attention for a few selected tokens

Quadratic dependency on Sequence Length for Memory and Computation

	What	a	nice	day
What				
a				
nice				
day				

Sparse Attention

- use fixed pattern of attention to reduce number of connections

Flash Attention

- Split score matrix into blocks
- Recompute attention matrix

Longformer

- local attention for most tokens
- global attention for a few selected tokens

Quadratic dependency on Sequence Length for Memory and Computation

	What	a	nice	day
What				
a				
nice				
day				

Sparse Attention

- use fixed pattern of attention to reduce number of connections

Flash Attention

- Split score matrix into blocks
- Recompute attention matrix

Longformer

- local attention for most tokens
- global attention for a few selected tokens

Quadratic dependency on Sequence Length for Memory and Computation

	What	a	nice	day
What				
a				
nice				
day				

Sparse Attention

- use fixed pattern of attention to reduce number of connections

Flash Attention

- Split score matrix into blocks
- Recompute attention matrix

Longformer

- local attention for most tokens
- global attention for a few selected tokens



Debugging Tips - 1

Memory Usage

Which parameters impact the model's memory usage?

Dimensions Mismatch

How are input and output dimensions defined in a transformer?

Training and Evaluation Metrics

Are the metrics behaving as expected?

Memory Usage

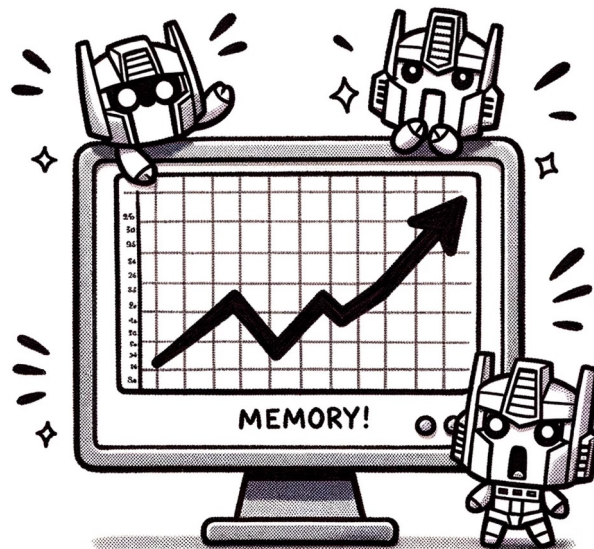
Which parameters impact the model's memory usage?

Dimensions Mismatch

How are input and output dimensions defined in a transformer?

Training and Evaluation Metrics

Are the metrics behaving as expected?



Memory Usage

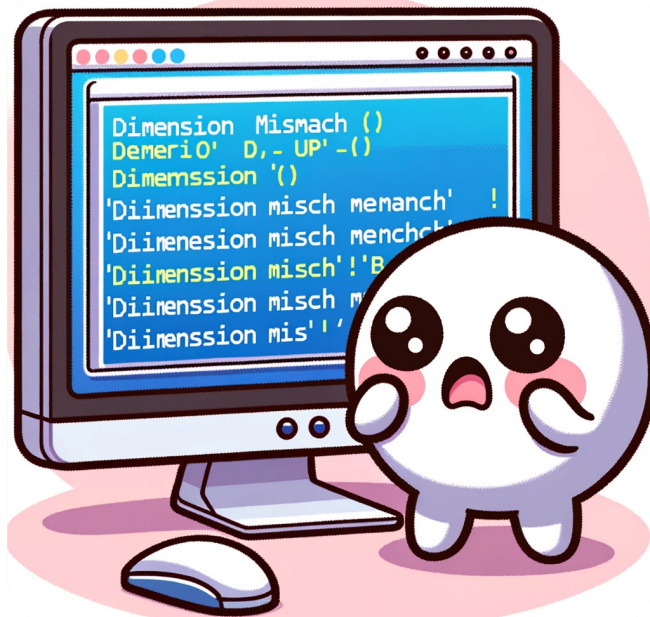
Which parameters impact the model's memory usage?

Dimensions Mismatch

How are input and output dimensions defined in a transformer?

Training and Evaluation Metrics

Are the metrics behaving as expected?



Memory Usage

Which parameters impact the model's memory usage?

Dimensions Mismatch

How are input and output dimensions defined in a transformer?

Training and Evaluation Metrics

Are the metrics behaving as expected?



Start Small

We provide you a Debug and Run mode

Performance Indicator

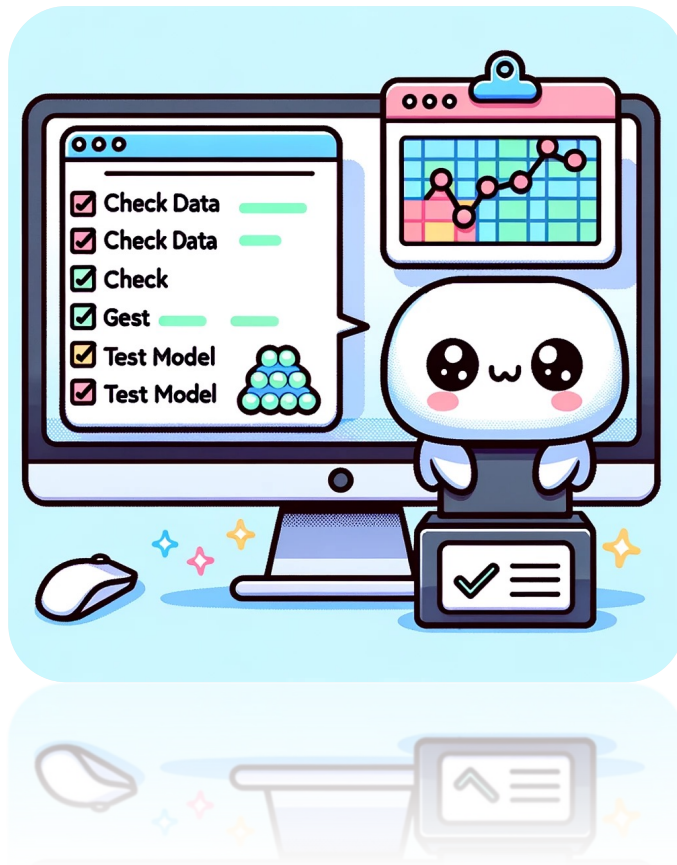
Expect a training perplexity around 20 after 50 epochs

Start Small

We provide you a Debug and Run mode

Performance Indicator

Expect a training perplexity around 20 after 50 epochs



Start Small

We provide you a Debug and Run mode

Performance Indicator

Expect a training perplexity around 20 after 50 epochs

