

# Transformers for NLP

Transformers and LLM

# Outline

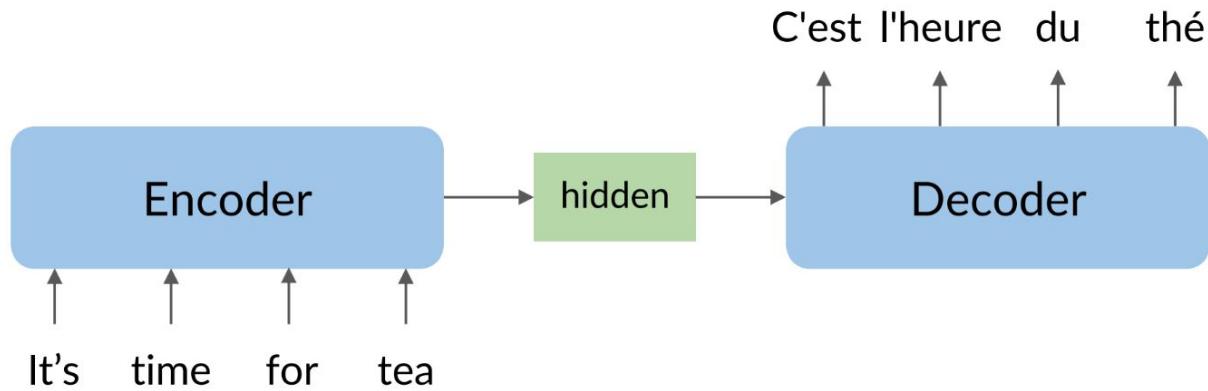
1. What is transformers
2. Explain how RNN's seq2seq works
3. Why transformers. Explain with RNN/LSTM reference. Explain problems of memory issue and vanishing gradients. Advantages of Transformers. Parallel processing, Alleviating Vanishing Gradient Problem, Long-Range Dependency Capture, Transfer Learning Capabilities.
4. Components of transformers. Everything in single image
5. Explain basic components one by one. 1. Encoder 2. Self-Attention 3. Multi-head self attention 4. Decoder 5. Masked self attention 6. Cross Self Attention [Encode-decoder attention]

# Seq2Seq model

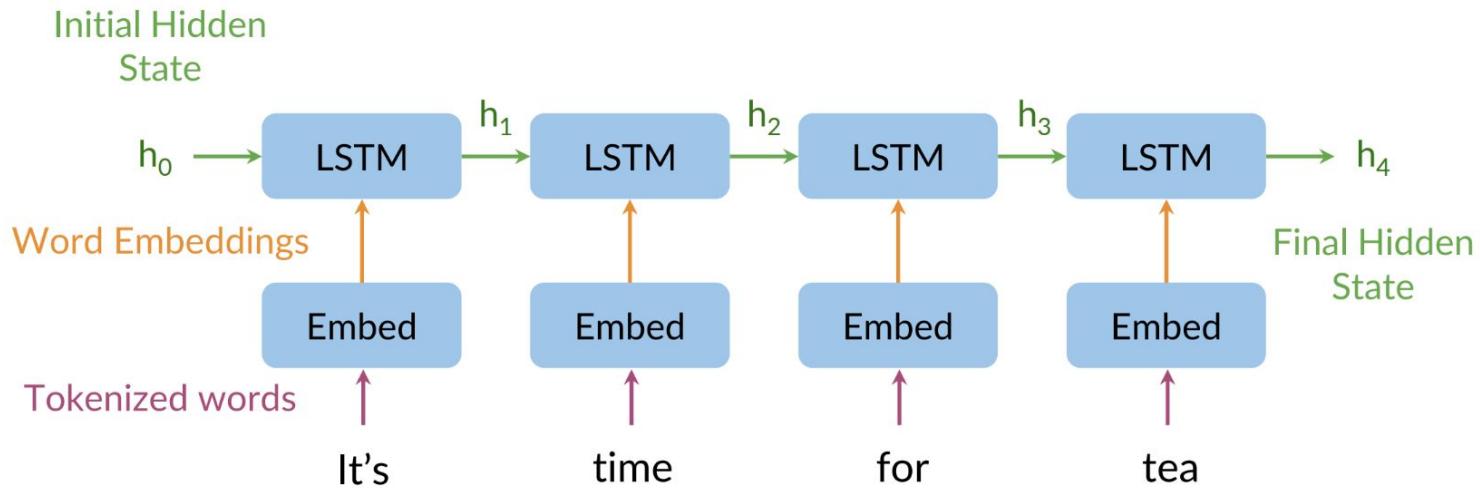
- Introduced by Google in 2014
- Maps variable-length sequences to fixed-length memory
- Inputs and outputs can have different lengths
- LSTMs and GRUs to avoid vanishing and exploding gradient problems



# Seq2Seq model

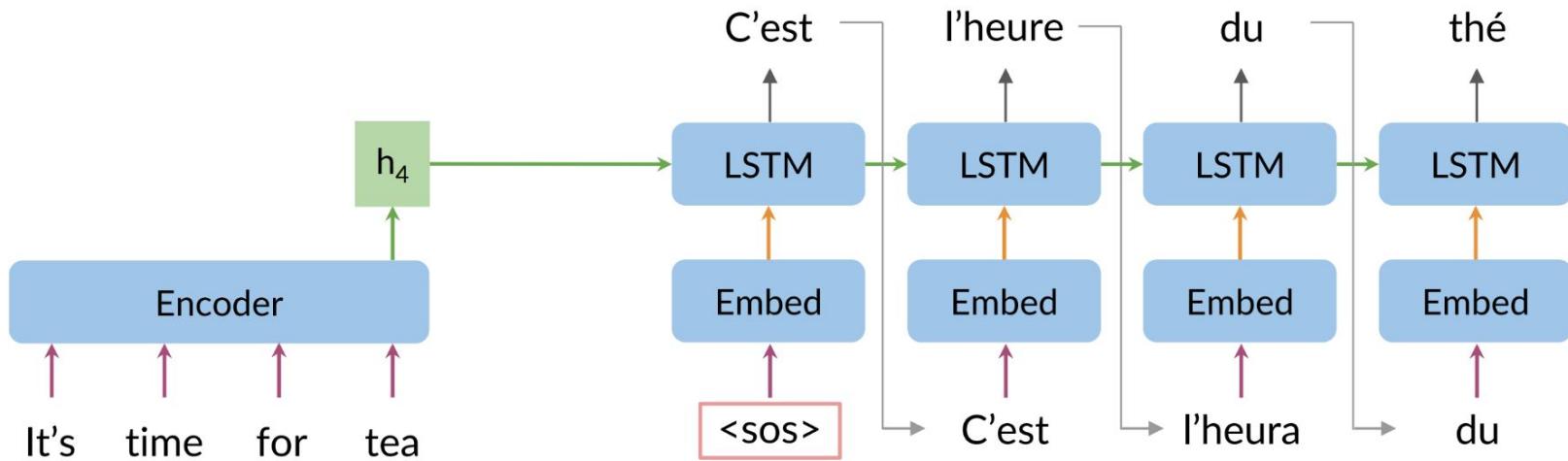


# Seq2Seq encoder



Encodes the overall meaning of the sentence

# Seq2Seq decoder



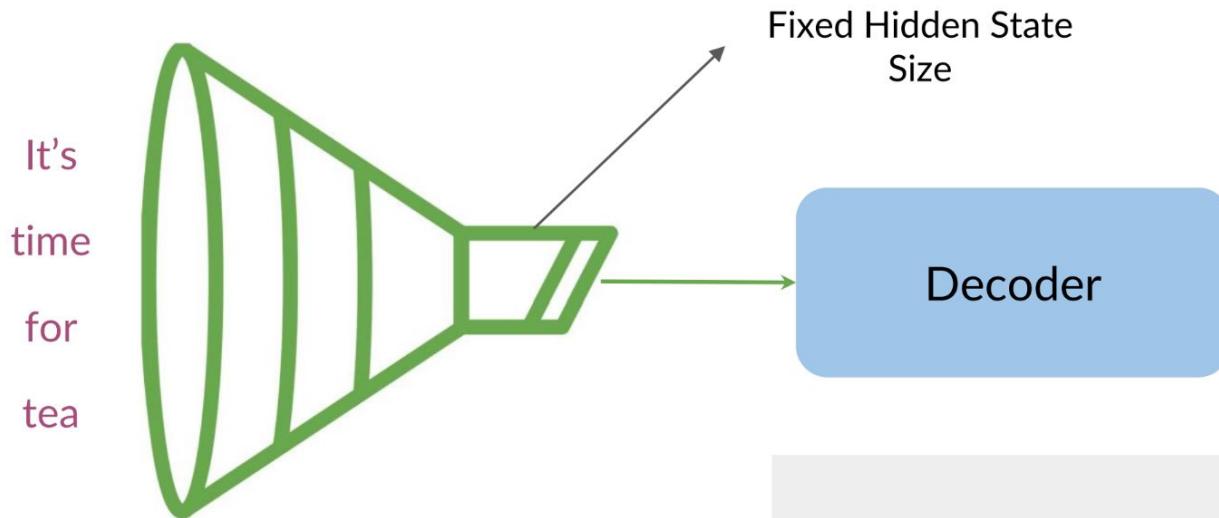
# Seq2Seq shortcomings

- Variable-length sentences + fixed-length memory =



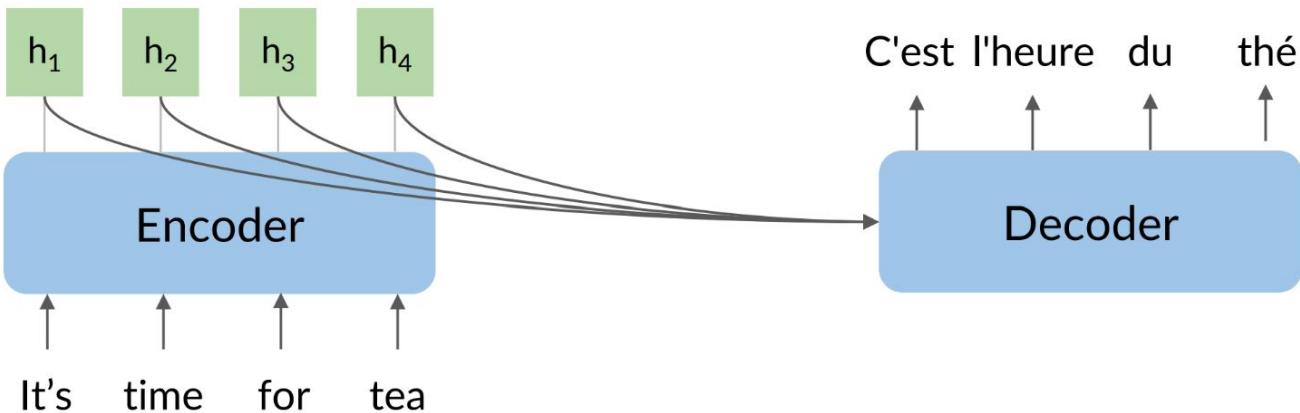
- As sequence size increases, model performance decreases

# The information bottleneck

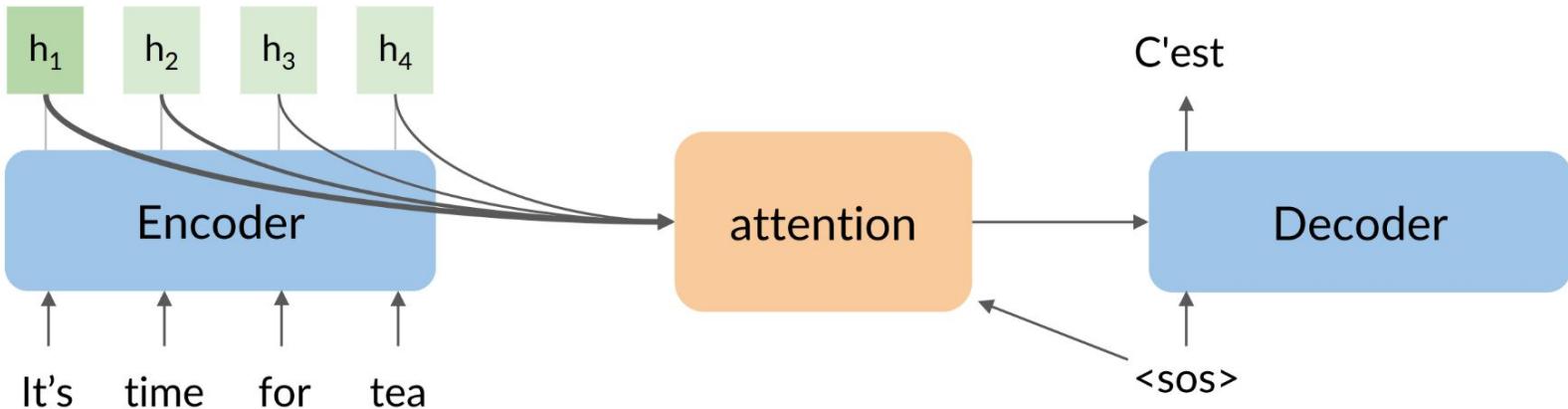


A fixed amount of informations goes from the encoder to the decoder

# Use all the encoder hidden states?



# Solution: focus attention in the right place



The model can focus on specific hidden states at every step

# Seq2Seq model with attention

---

# NEURAL MACHINE TRANSLATION BY JOINTLY LEARNING TO ALIGN AND TRANSLATE

**Dzmitry Bahdanau**

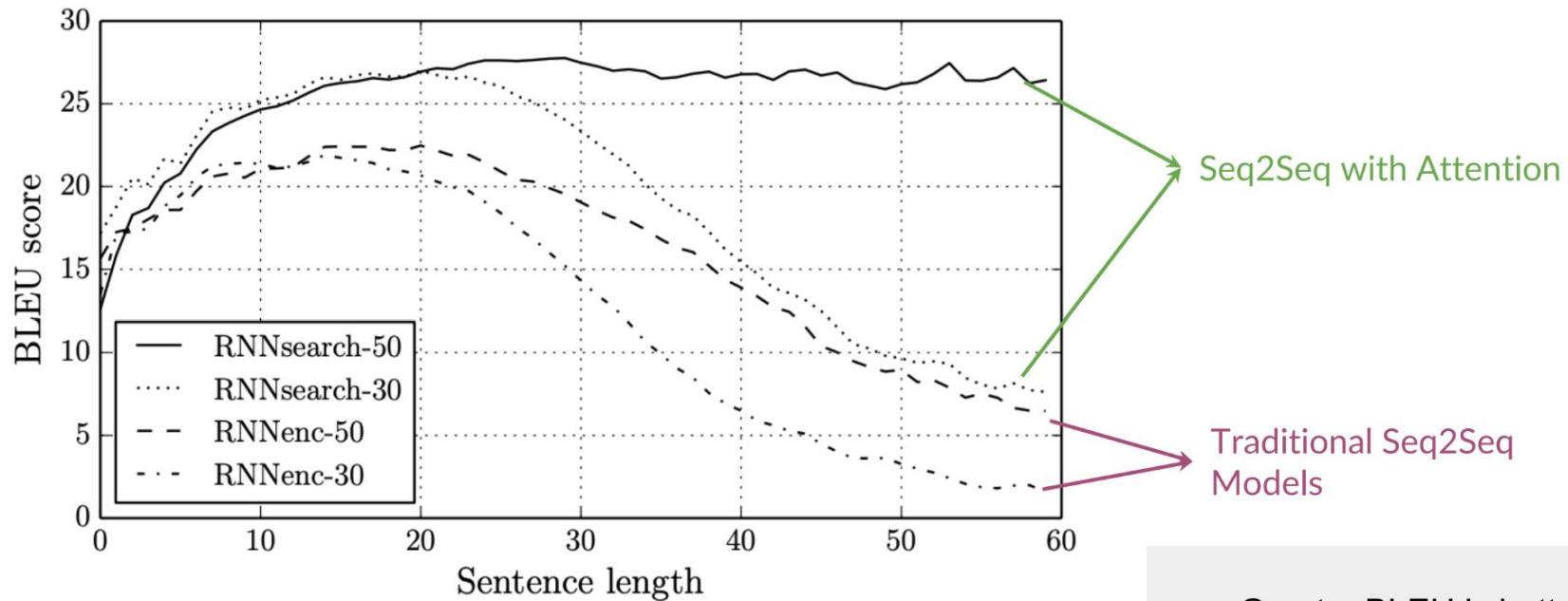
Jacobs University Bremen, Germany

**KyungHyun Cho    Yoshua Bengio\***

Université de Montréal

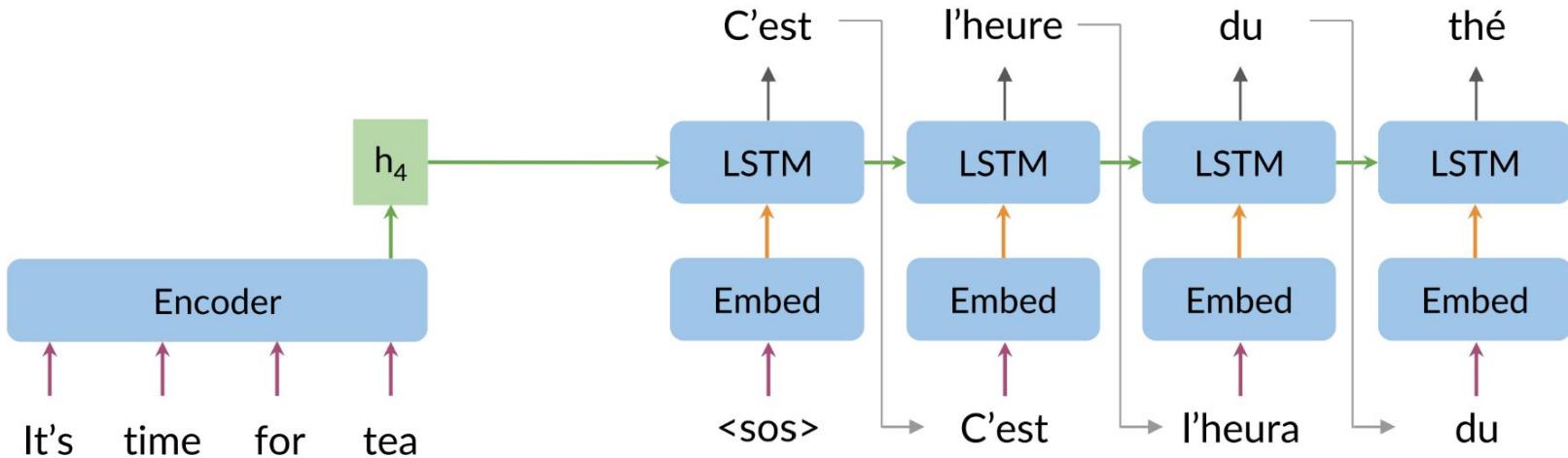
[1409.0473.pdf \(arxiv.org\)](https://arxiv.org/pdf/1409.0473.pdf)

# Performance

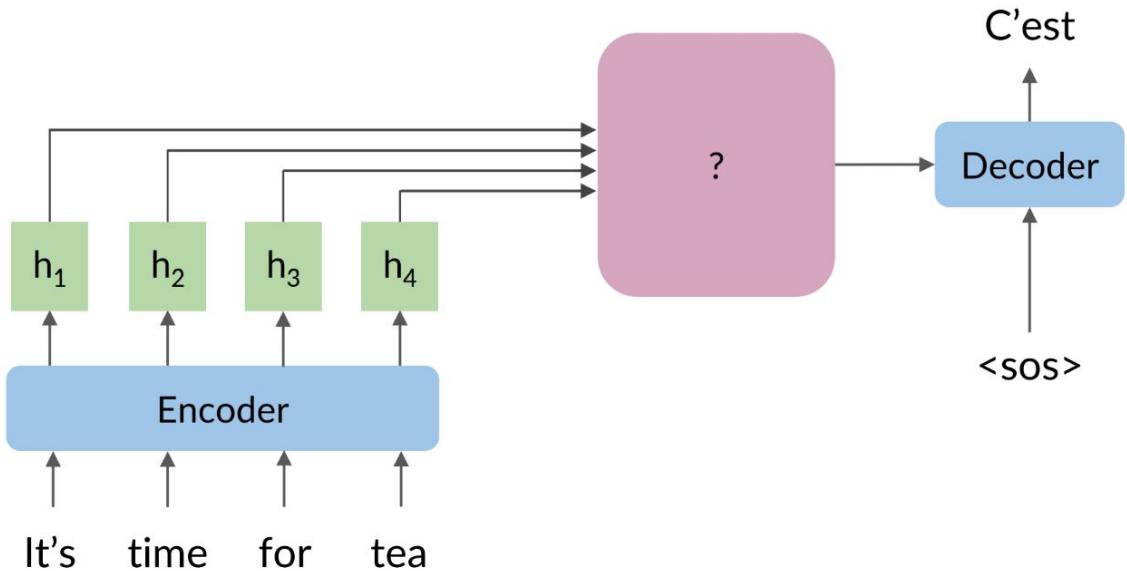


Greater BLEU is better

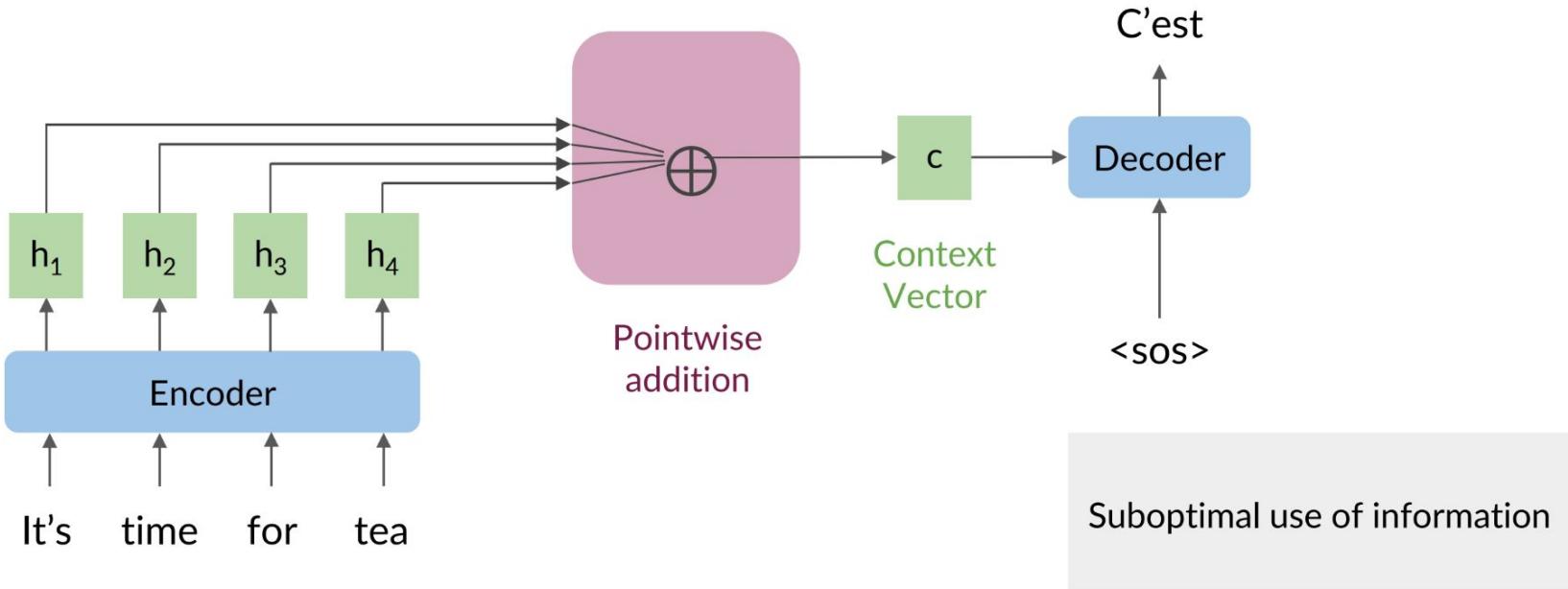
# Traditional seq2seq models



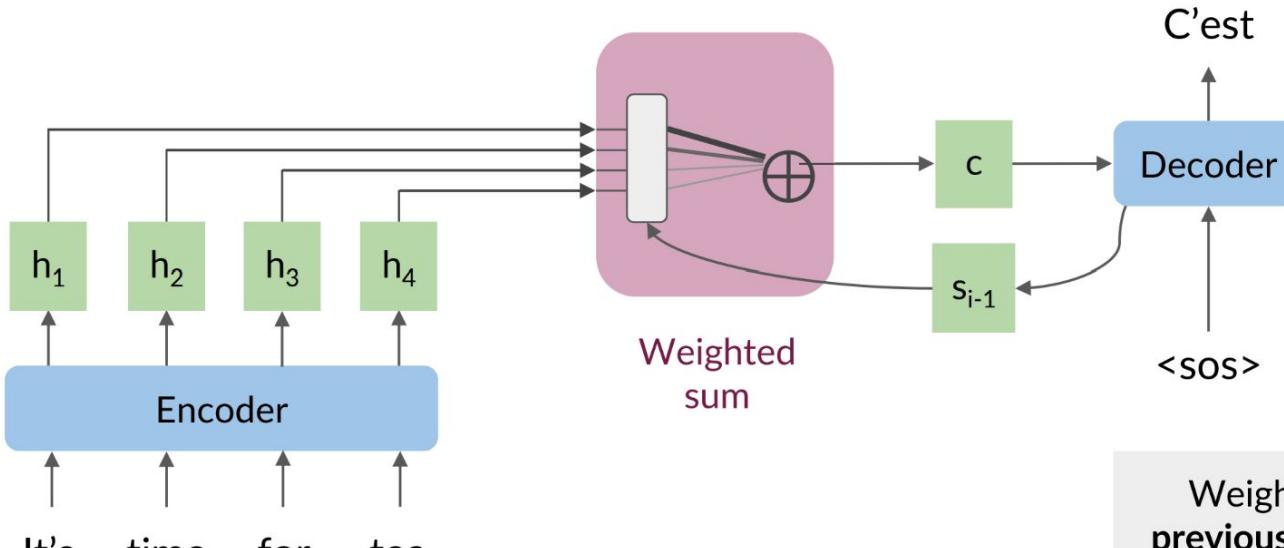
# How to use all the hidden states?



# How to use all the hidden states?



# How to use all the hidden states?



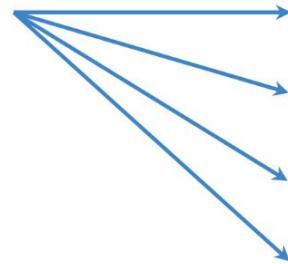
Weights depend on the  
**previous hidden state** in the  
decoder

# Queries, Keys, Values and Attention

---

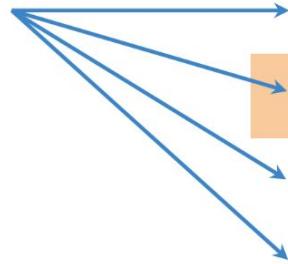
# Queries, Keys, Values

Query	Key	Value
l'heure	It's	[0.5, 0.2, -1.2, ..., ]
	time	[0.2, -0.7, 0.9, ..., ]
	for	[1.3, 0.3, 0.8, ..., ]
	tea	[-0.4, 0.6, -1.1, ..., ]



# Queries, Keys, Values

Query	Key	Value
l'heure	It's	[0.5, 0.2, -1.2, ..., ]
	time	[0.2, -0.7, 0.9, ..., ]
	for	[1.3, 0.3, 0.8, ..., ]
	tea	[-0.4, 0.6, -1.1, ..., ]

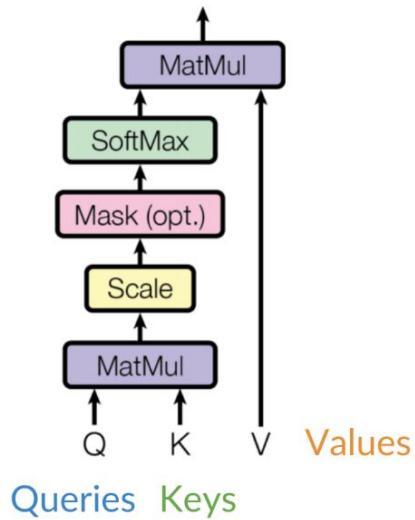


# Queries, Keys, Values

Query	Key	Value
l'heure	It's	[0.5, 0.2, -1.2, ..., ]
	time	[0.2, -0.7, 0.9, ..., ]
	for	[1.3, 0.3, 0.8, ..., ]
	tea	[-0.4, 0.6, -1.1, ..., ]

Similarity is used in for weighted sum

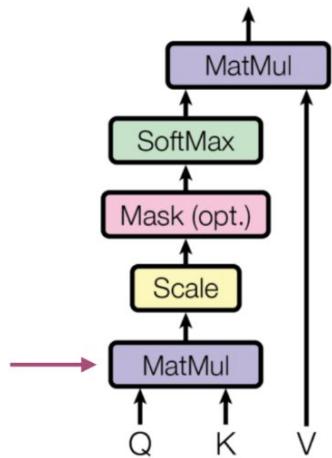
# Scaled dot-product attention



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

(Vaswani et al., 2017)

# Scaled dot-product attention

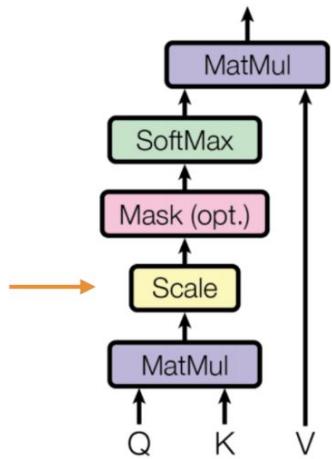


(Vaswani et al., 2017)

Similarity Between  
Q and K

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

# Scaled dot-product attention

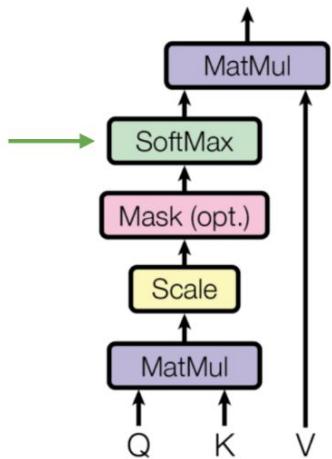


(Vaswani et al., 2017)

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

Scale using the root  
of the key vector  
size

# Scaled dot-product attention

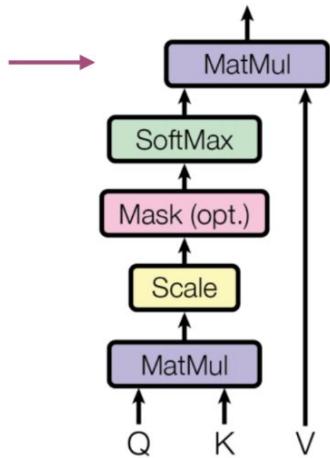


(Vaswani et al., 2017)

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

Weights for the  
weighted sum

# Scaled dot-product attention



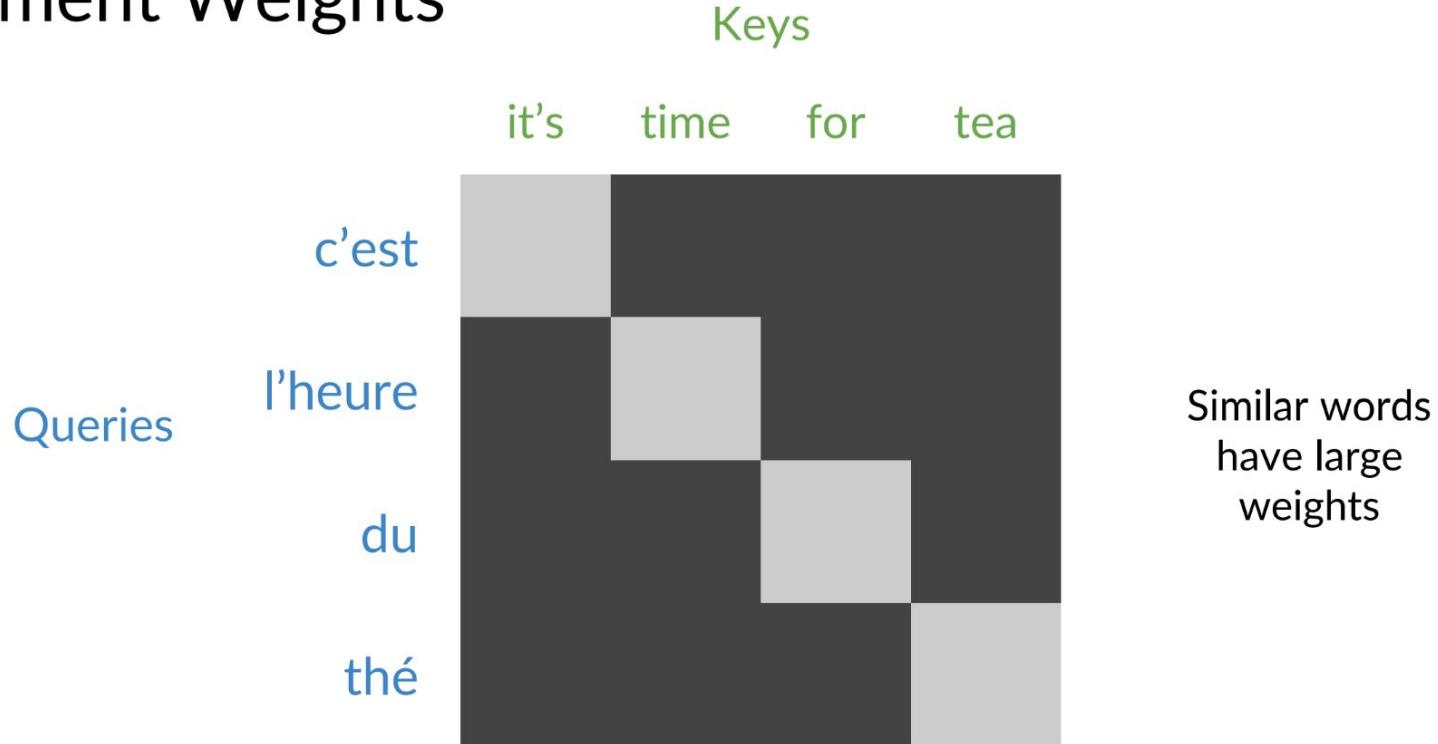
(Vaswani et al., 2017)

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

Weighted sum of values V

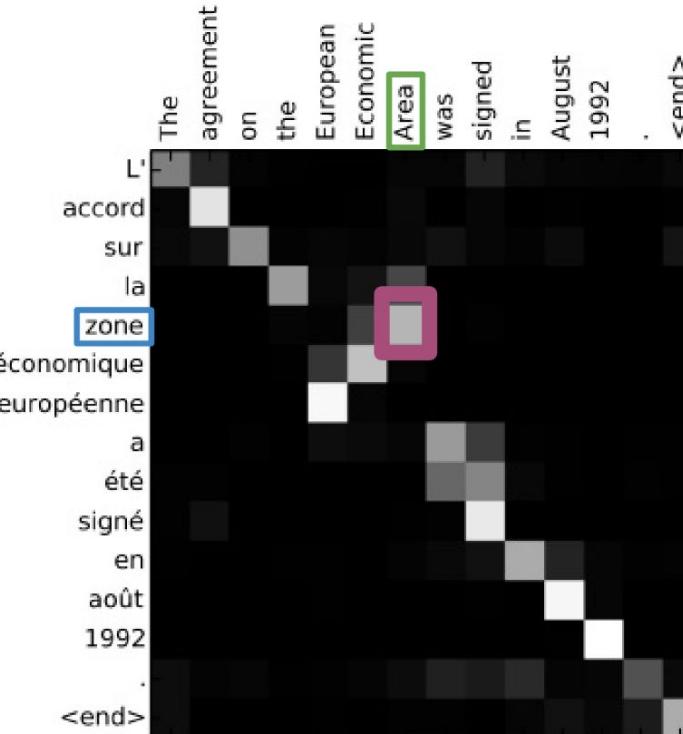
Just two matrix multiplications  
and a Softmax!

# Alignment Weights



# Flexible attention

Works for languages with different grammar structures!



Bahdanau et al.,  
2015

# Summary

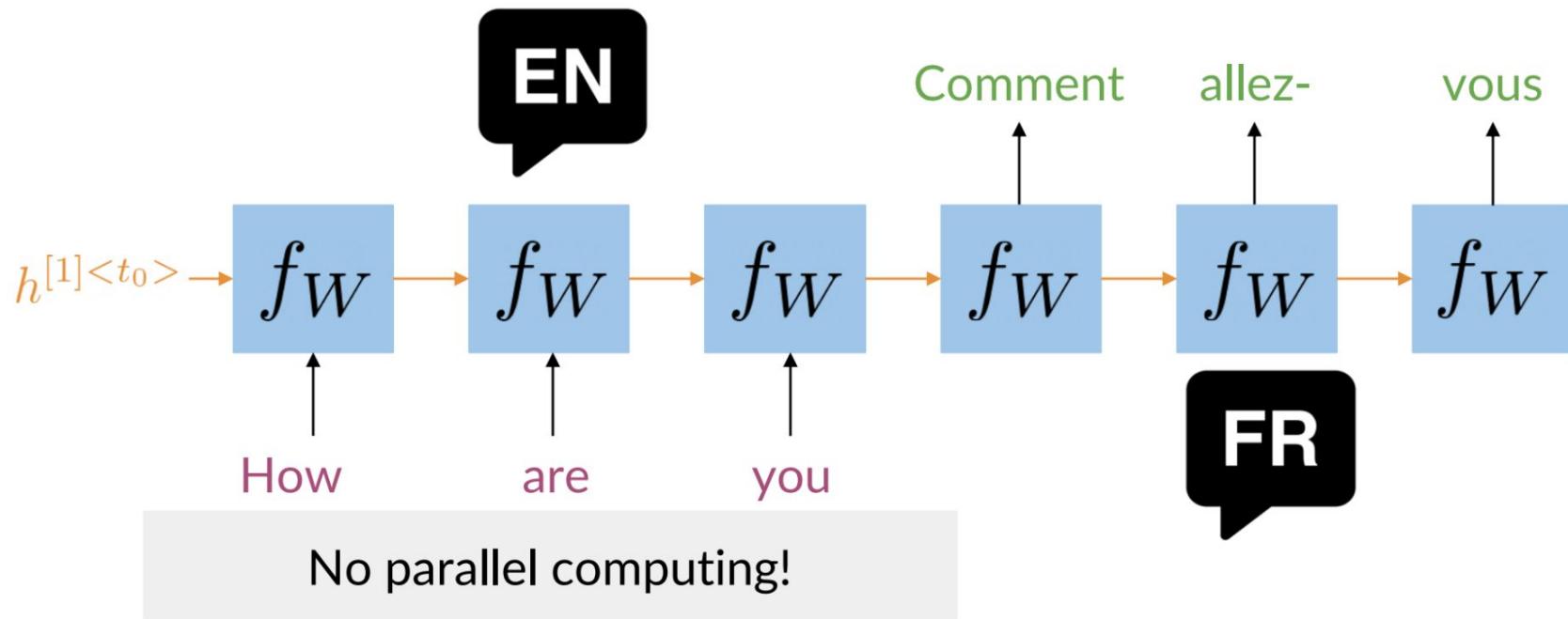
- Attention is a layer that lets a model focus on what's important
- Queries, Values, and Keys are used for information retrieval inside the Attention layer
- Works for languages with very different grammatical structures



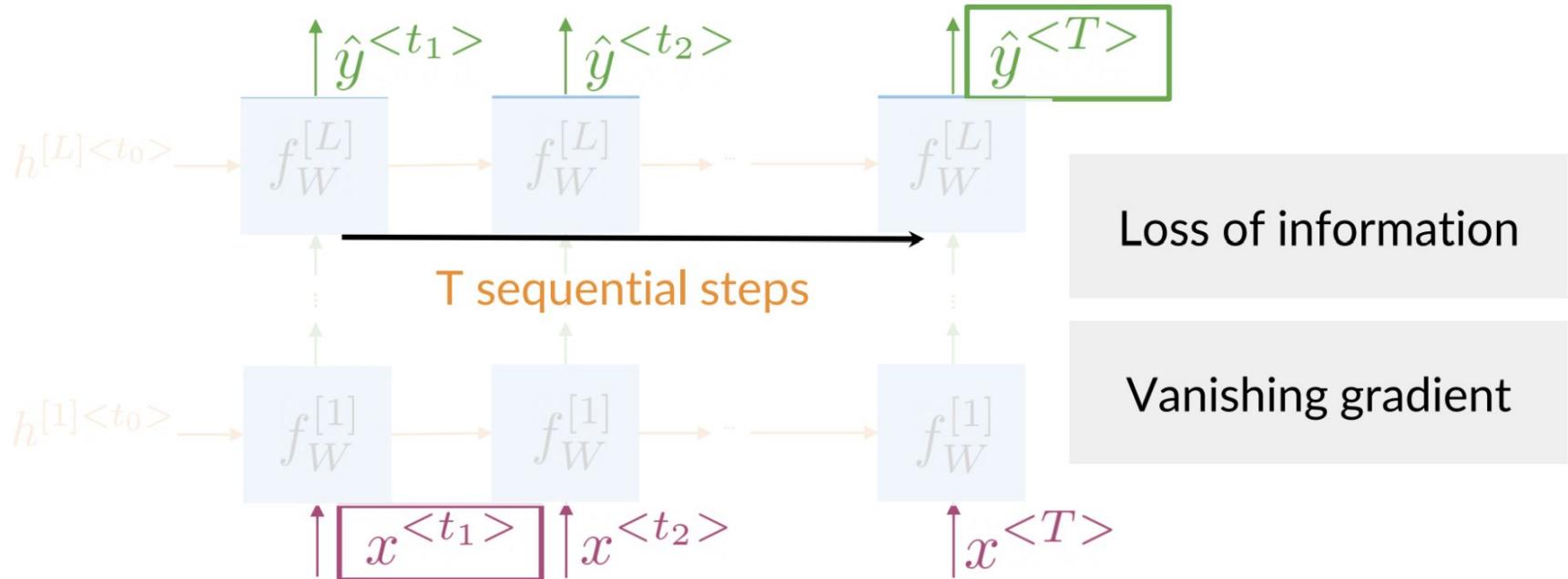
# Transformers vs RNNs

---

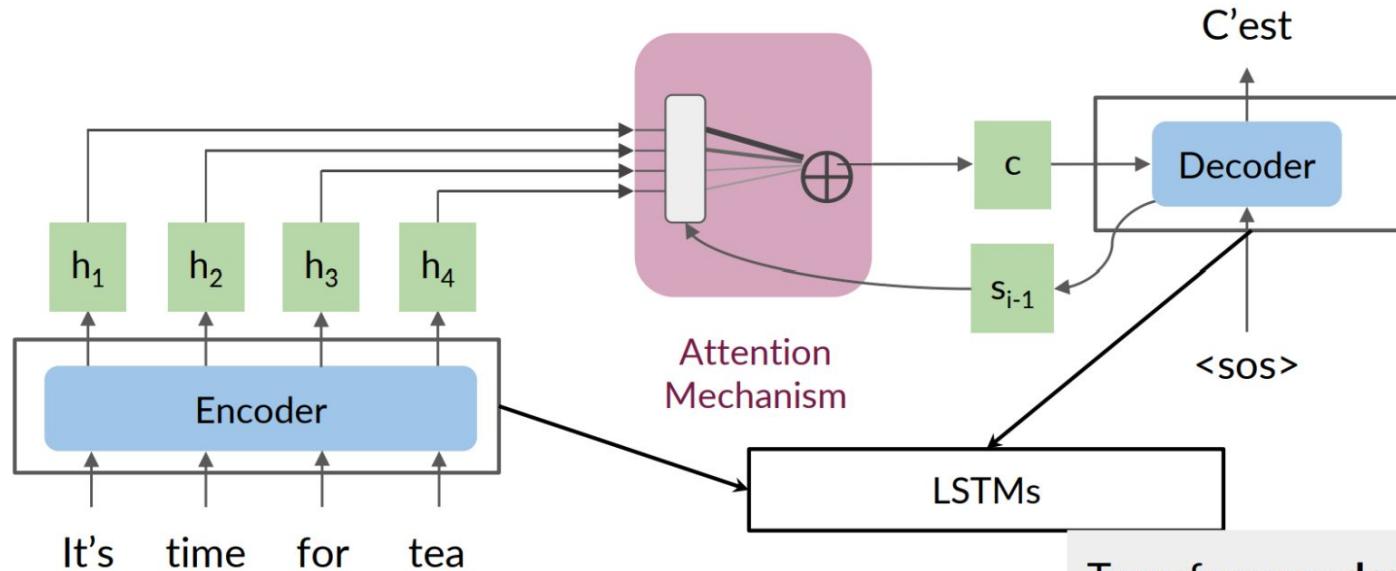
# Neural Machine Translation



# Seq2Seq Architectures



# RNNs vs Transformer: Encoder-Decoder



Transformers **don't** use RNNs, such as LSTMs or GRUs

# Transformers

# Overview

---

---

# 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\***

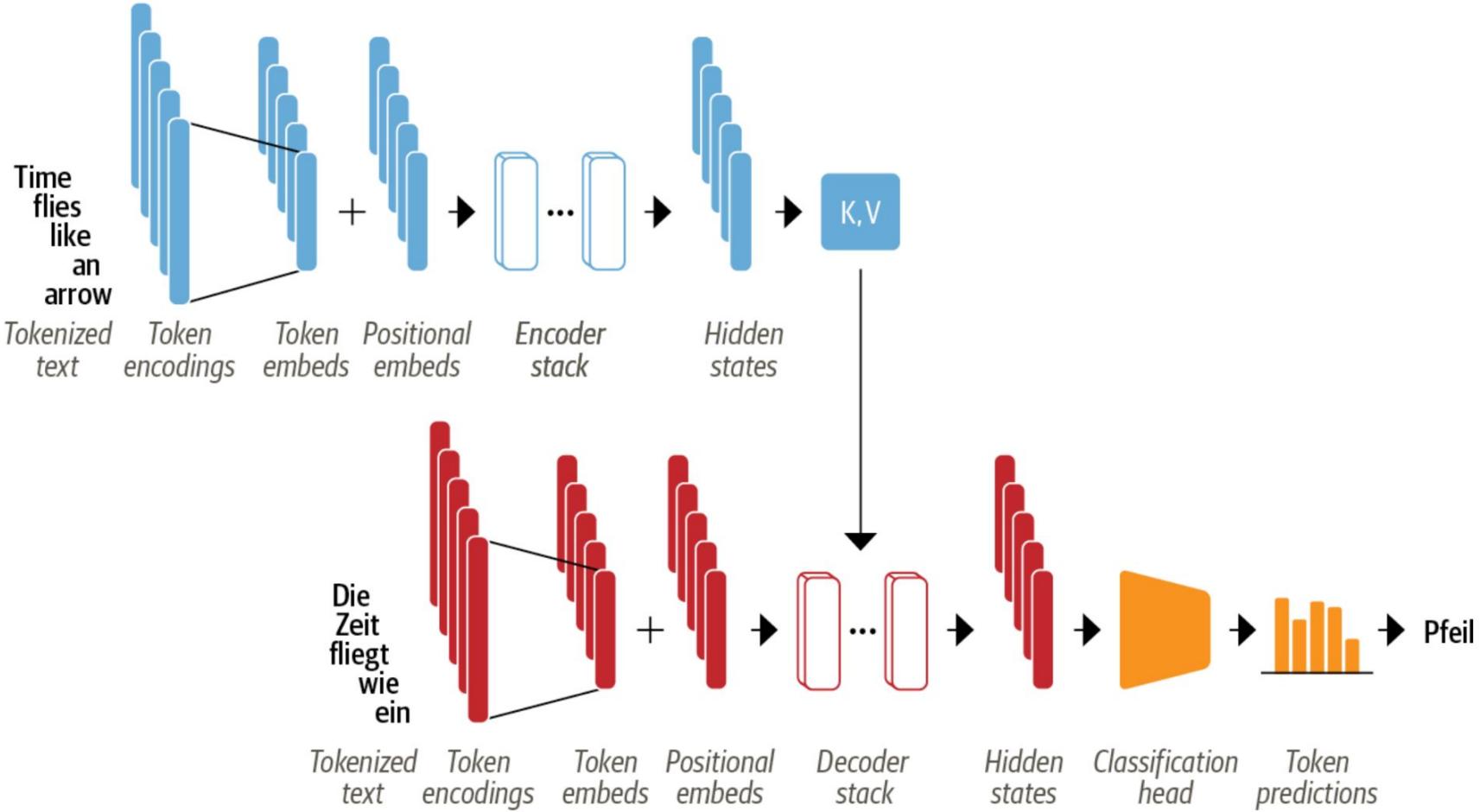
Google Brain

lukaszkaiser@google.com

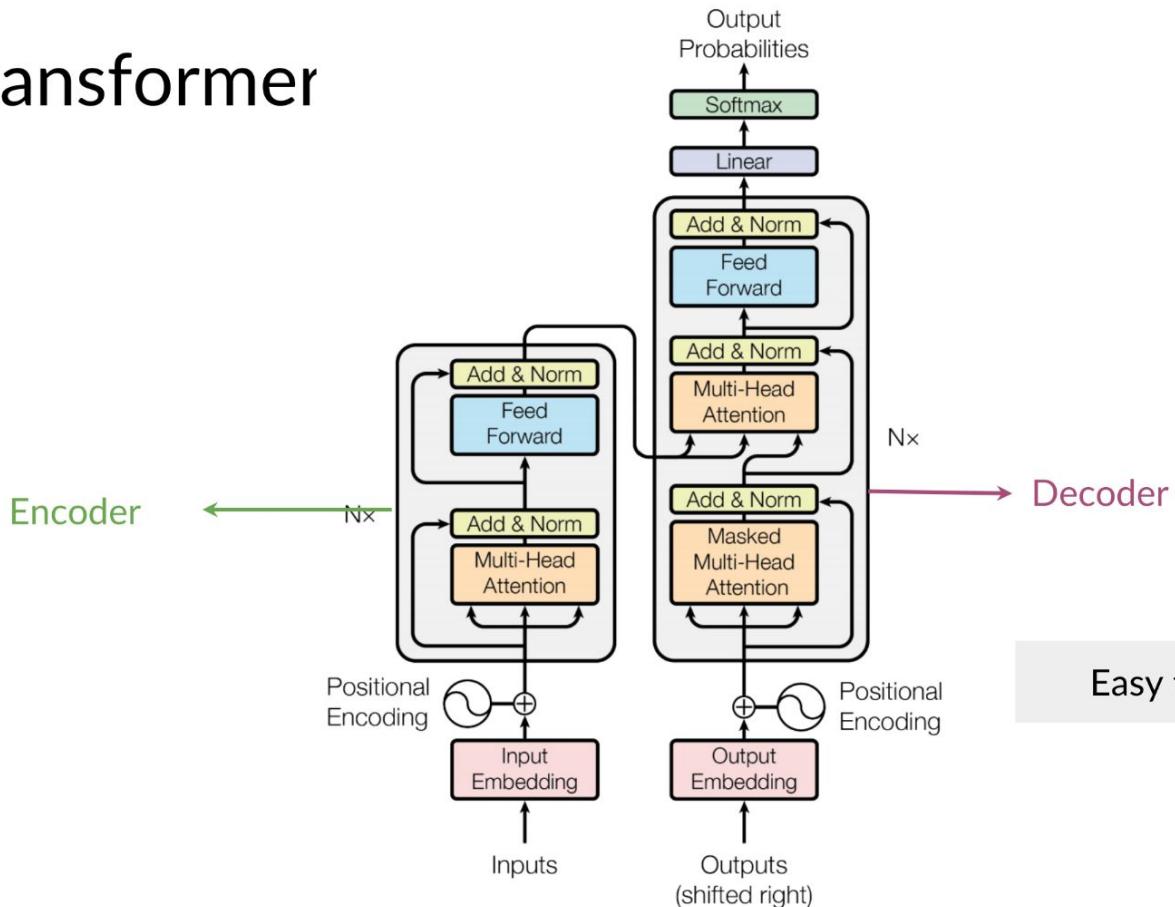
**Illia Polosukhin\*** ‡

illia.polosukhin@gmail.com

[1706.03762.pdf \(arxiv.org\)](https://arxiv.org/abs/1706.03762)

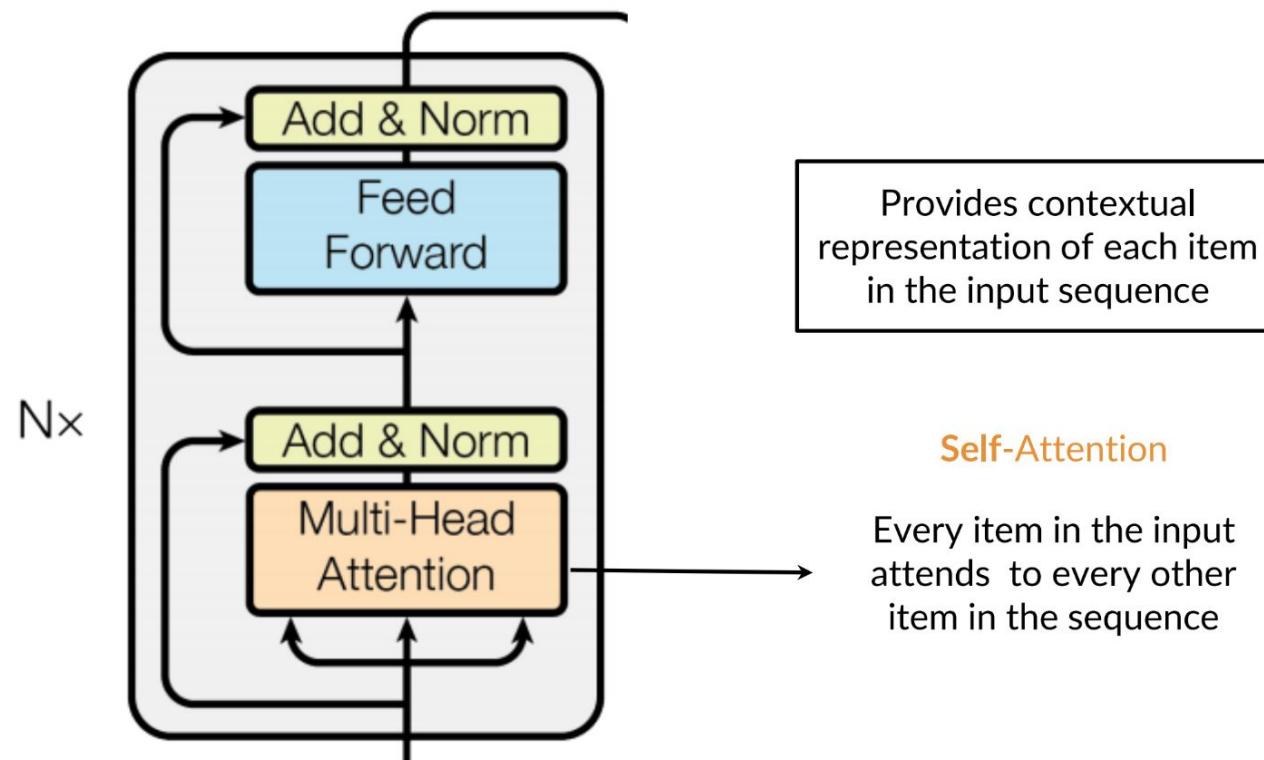


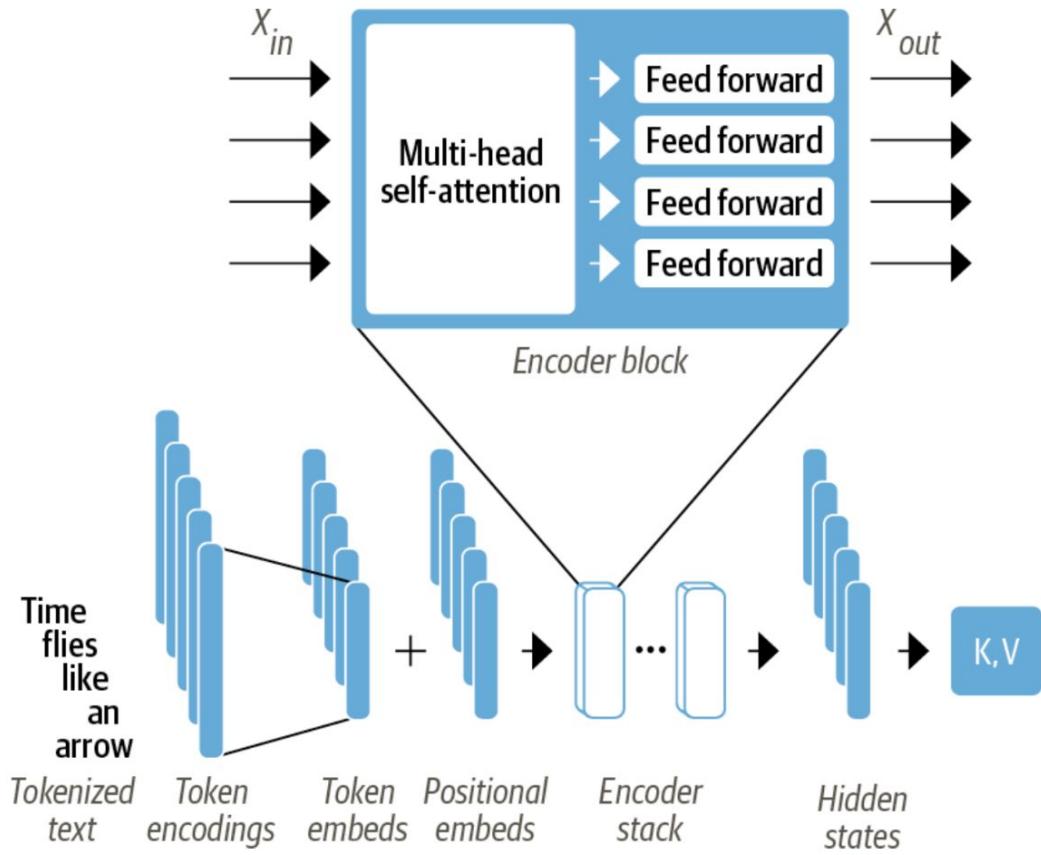
# The Transformer



Easy to parallelize!

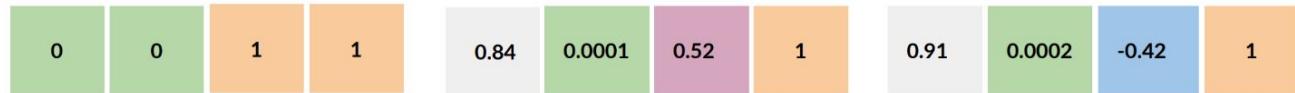
# The Encoder





# RNNs vs Transformer: Positional Encoding

POSITIONAL  
ENCODING

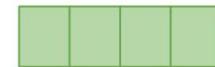
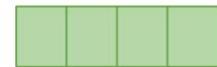
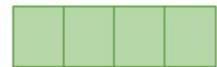


+

+

+

EMBEDDINGS



INPUT

Je

suis

content

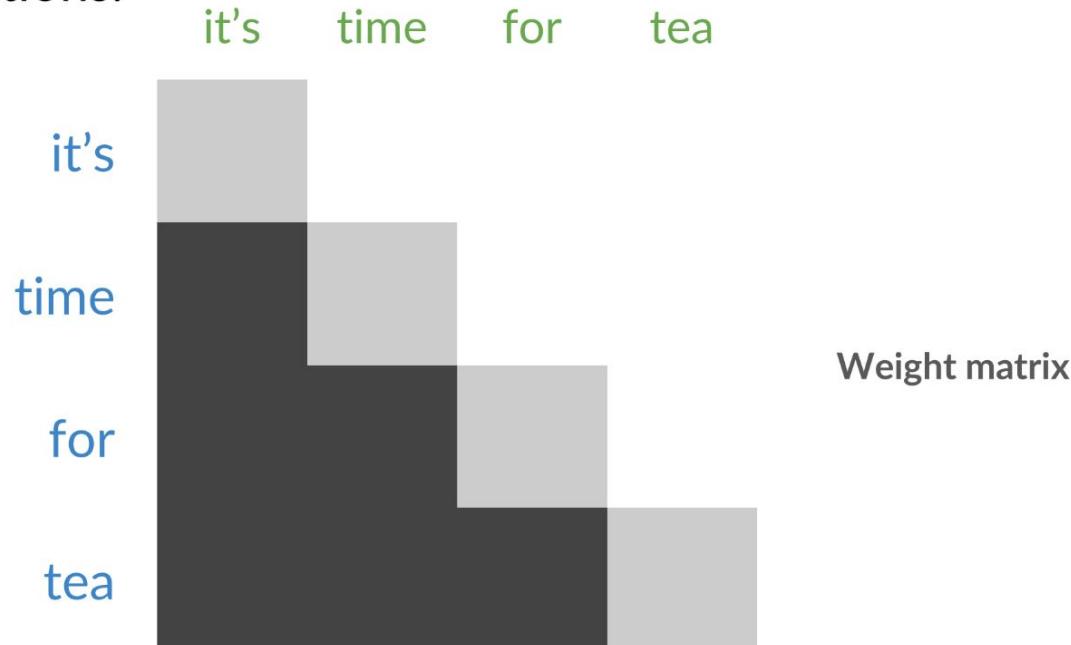
# Self-Attention

Queries, **keys** and **values** come from the **same sentence**



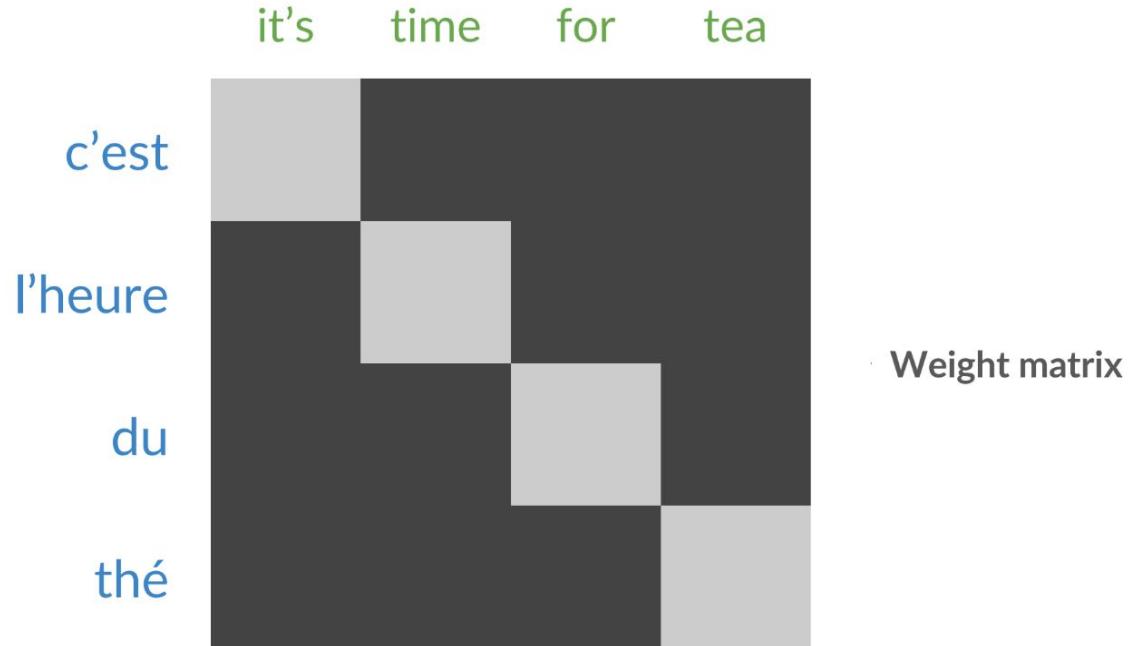
# Masked Self-Attention

Queries, keys and values come from the same sentence. Queries don't attend to future positions.



# Encoder-Decoder Attention

Queries from one sentence, keys and values from another



# Summary

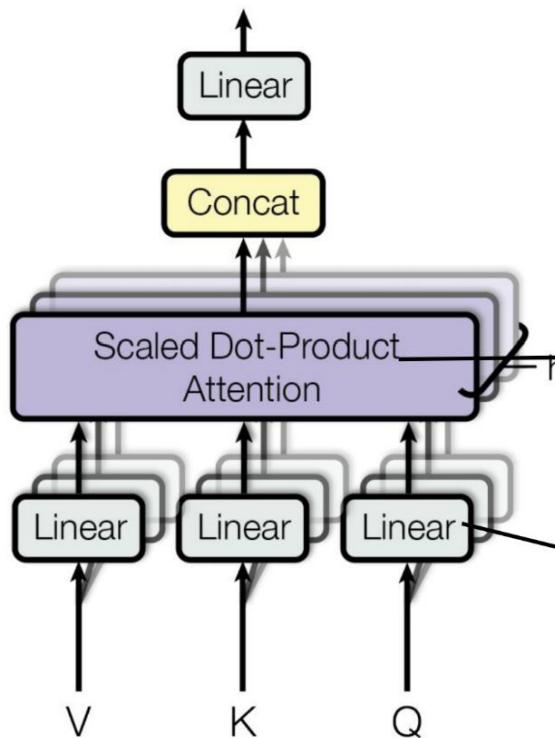
- There are three main ways of Attention: Encoder/Decoder, self-attention and masked self-attention.
- In self-attention, queries and keys come from the same sentence
- In masked self-attention queries cannot attend to the future



# Multi-head Attention

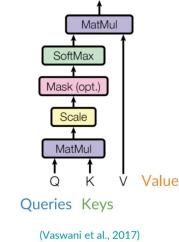
---

# Multi-Head Attention



Scaled dot-product  
attention multiple times in  
parallel

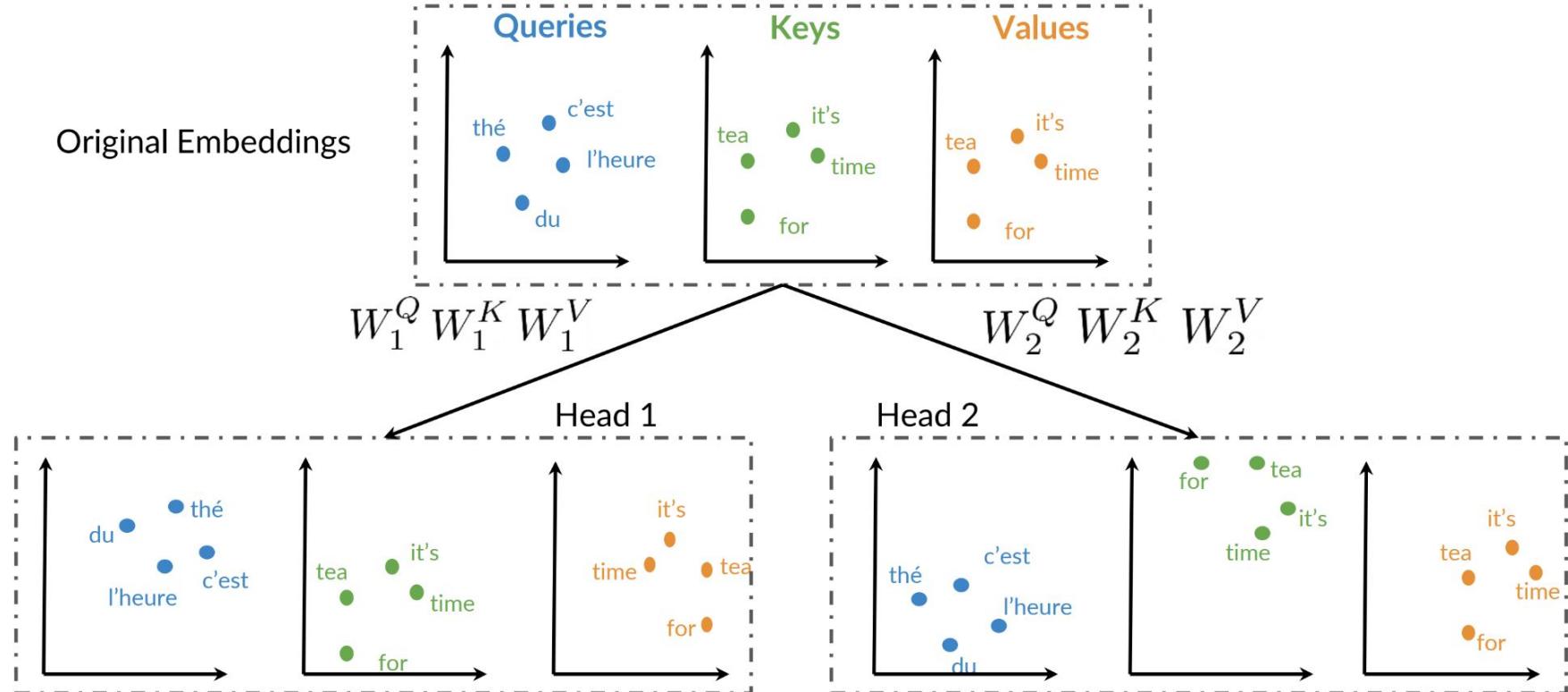
Linear transformations of  
the input queries, keys and  
values



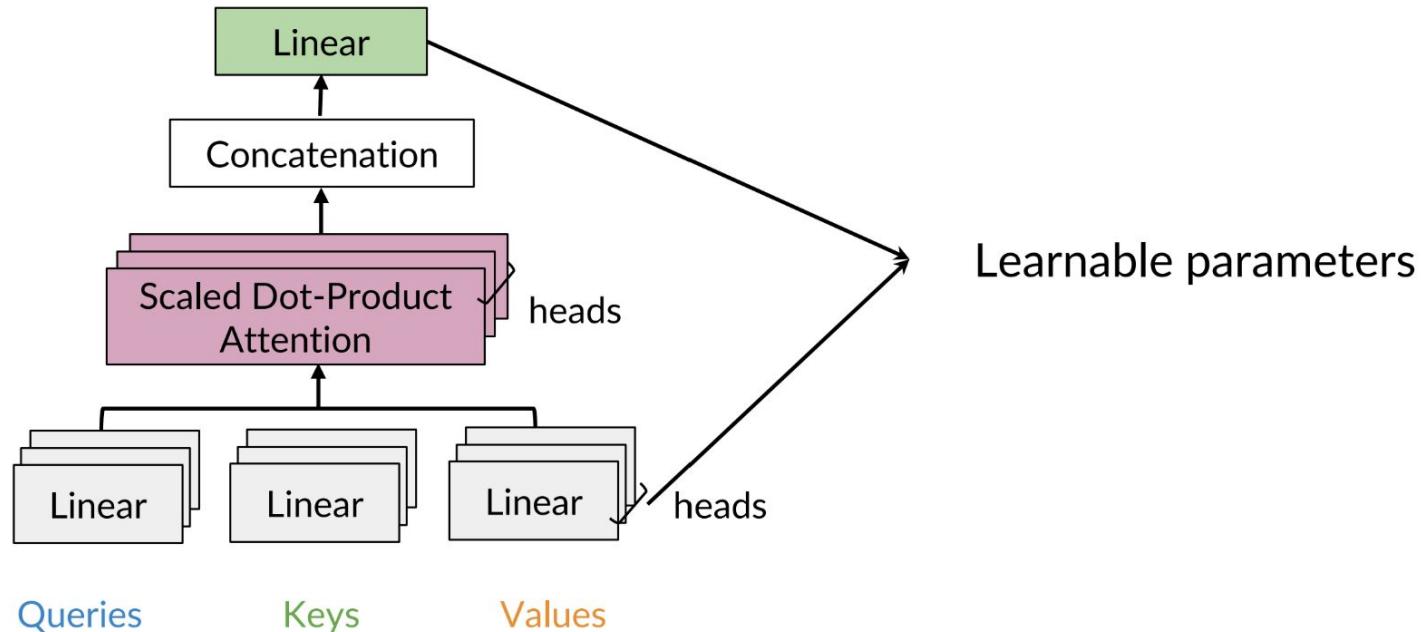
$$\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

(Vaswani et al., 2017)

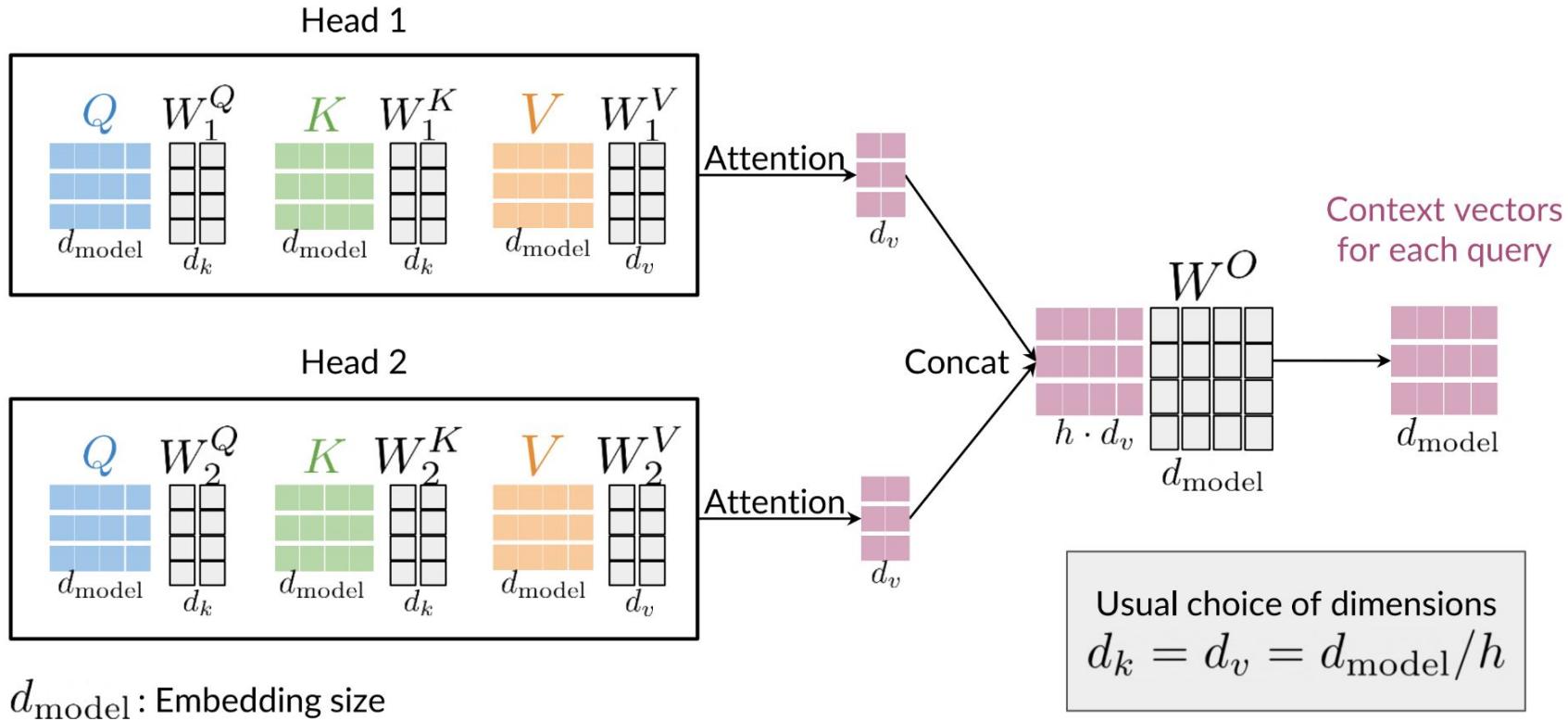
# Multi-Head Attention - Overview



# Multi-Head Attention - Overview



# Multi-Head Attention

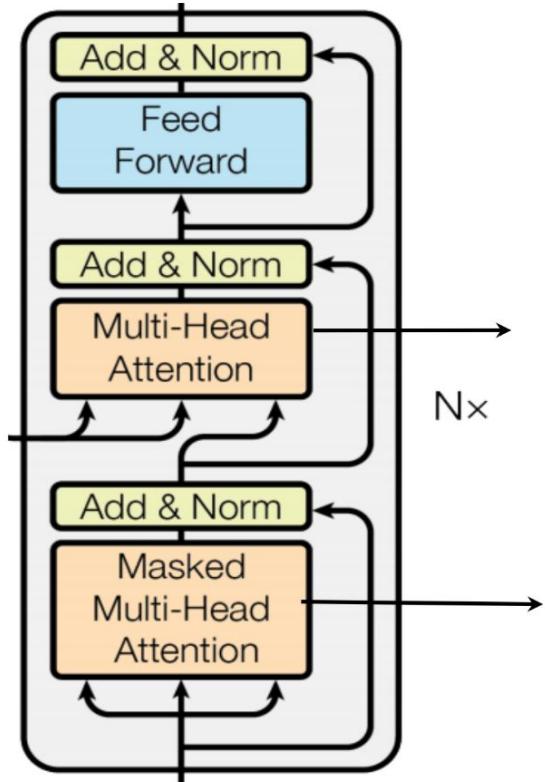


# Summary

- Multi-Headed models attend to information from different representations
- Parallel computations
- Similar computational cost to single-head attention



# The Decoder

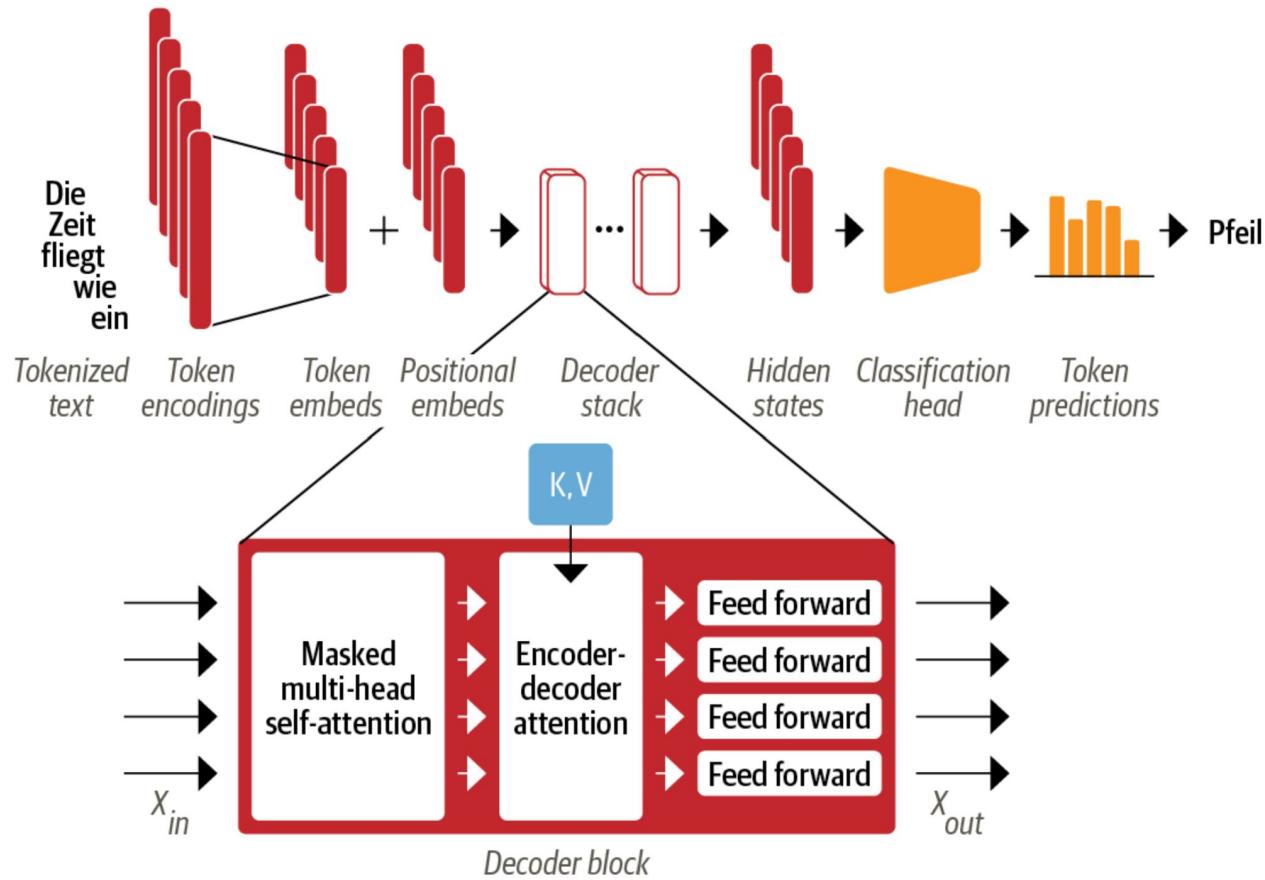


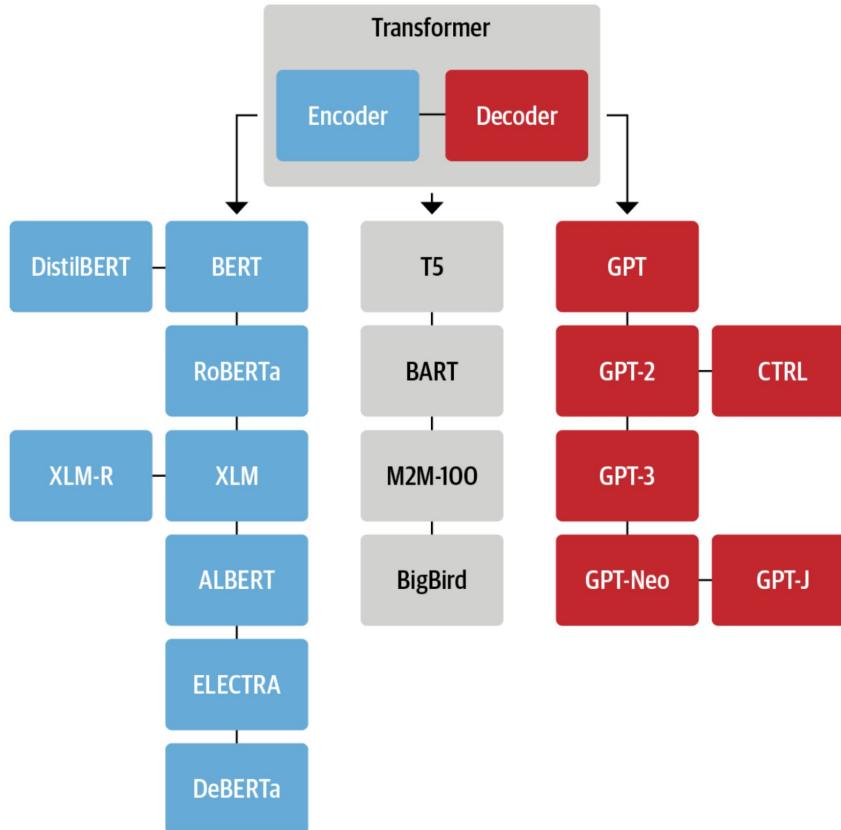
**Encoder-Decoder  
Attention**

Every position from the decoder attends to the outputs from the encoder

**Masked Self-Attention**

Every position attends to **previous** positions

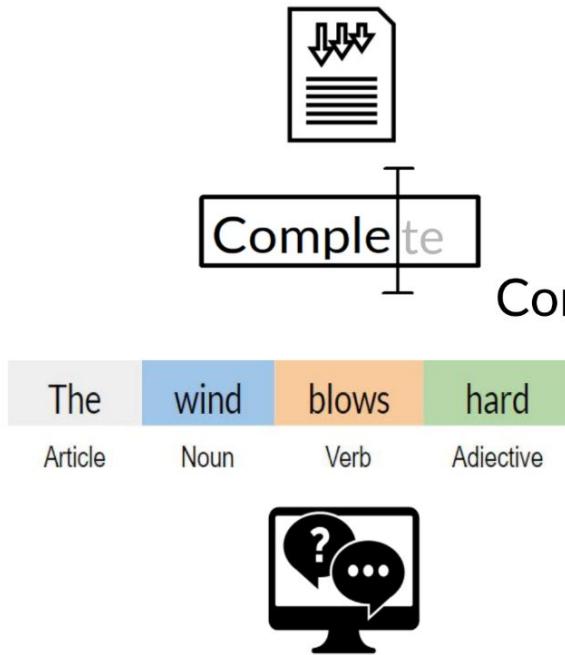




# Transformer Applications

---

# Transformer NLP applications



Text  
summarization

Auto-  
Complete

Named entity  
recognition (NER)

Question  
answering (Q&A)

Translation



Chat-bots



Other NLP tasks

[Sentiment Analysis](#)  
[Market Intelligence](#)  
[Text Classification](#)  
[Character Recognition](#)  
[Spell Checking](#)

# State of the Art Transformers

Radford, A., et al. (2018)  
Open AI

Devlin, J., et al. (2018)  
Google AI Language

Colin, R., et al. (2019)  
Google

GPT-2: Generative Pre-training for  
Transformer

BERT: Bidirectional Encoder  
Representations from Transformers

T5: Text-to-text transfer transformer



