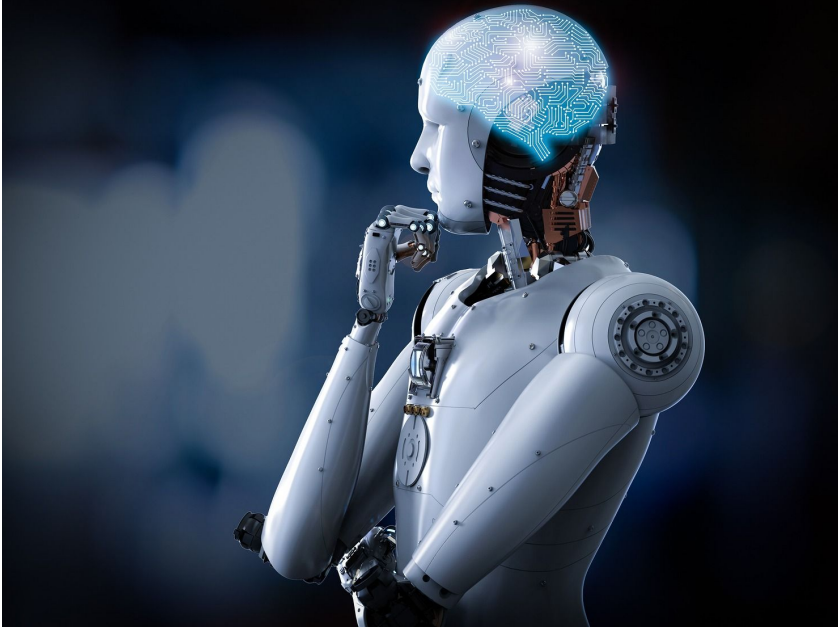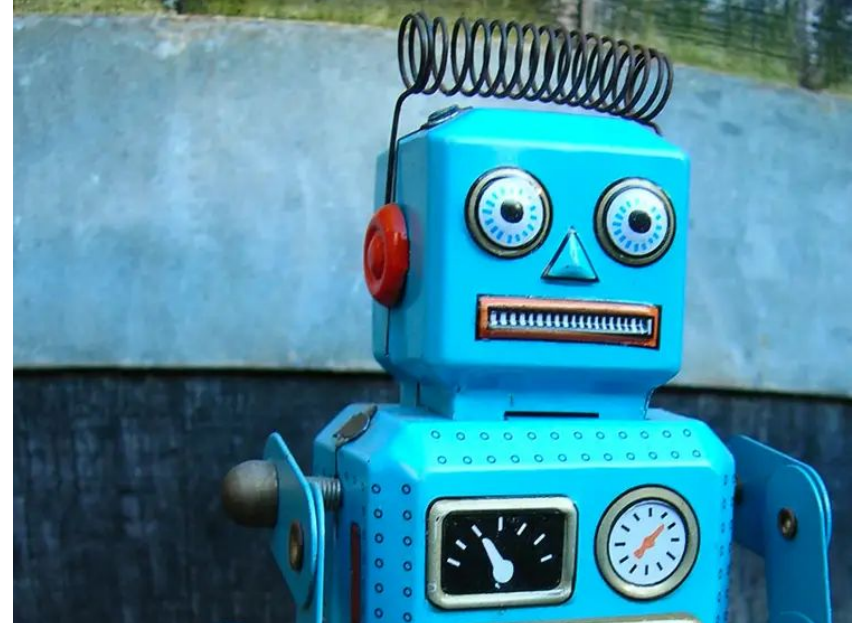# Introduction to Large Language Models

Eya Ben Charrada

# Too smart or too dumb?



ChatGPT has passed medical, law and business schools exams
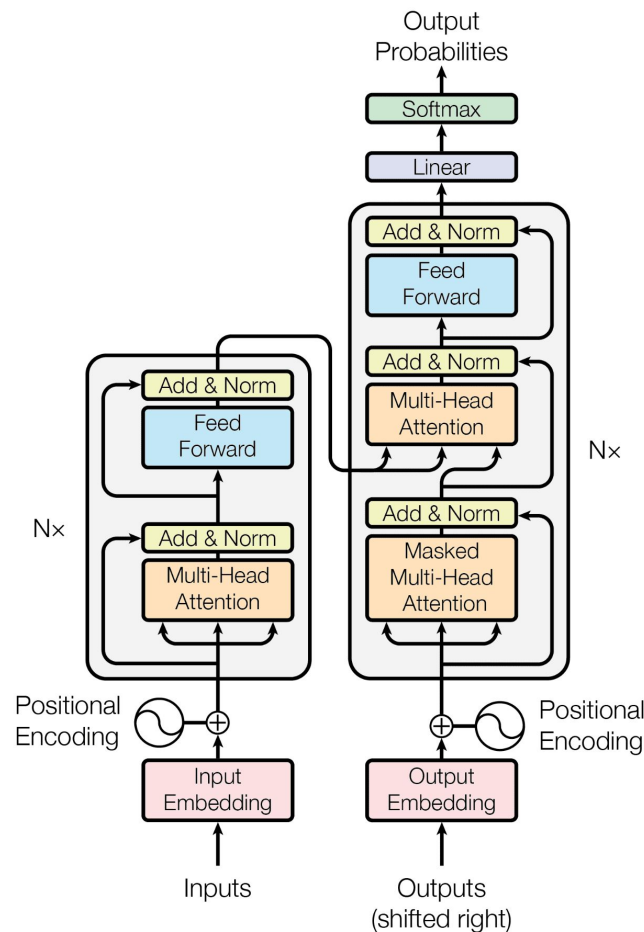


Fails to answer questions that a 5-years-old child can answer

2

# How does it work?
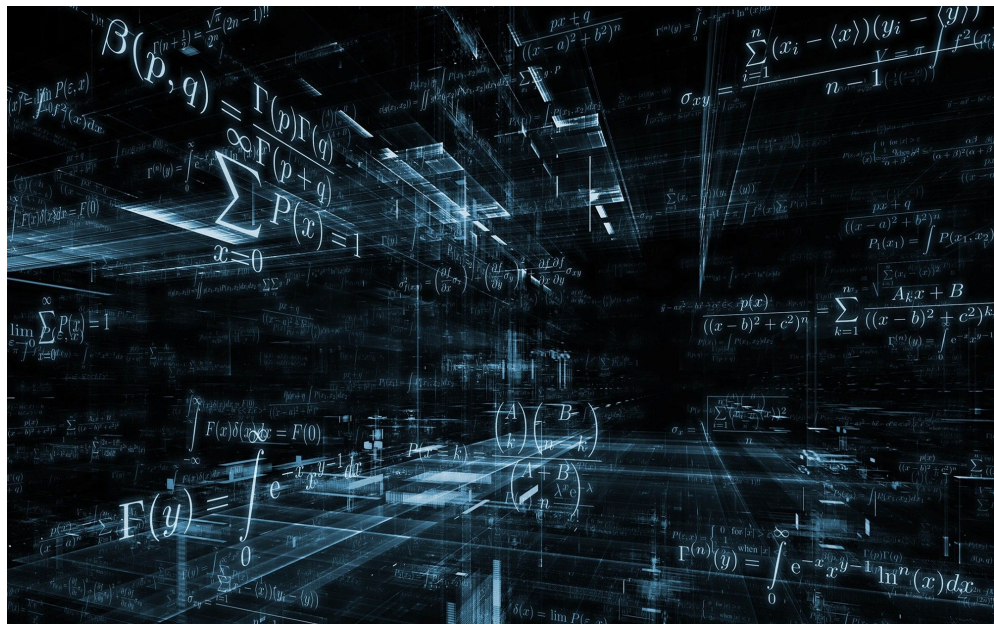
Famous google paper (2017)

**Attention Is All You Need**

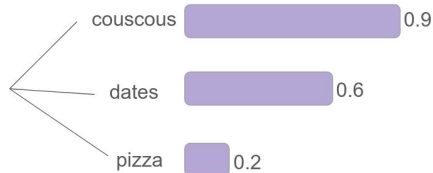$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$



3

# Goal

Understand how transformer models work :

Explore the hidden meaning behind mathematical equations

# Outline

## What is a language model

couscous 0.9

dates 0.6

pizza 0.2

## Before the transformer

I love Palestine <SOS>

Ich liebe Palästina

## Transformer architecture

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

Nx

Nx

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding

Inputs

Outputs (shifted right)

5

# Language Model

In Tunisia, I ate delicious $\longrightarrow$      $\longrightarrow$ ?

Translation, summarisation, question answering…

# Language Model

In Tunisia, I ate delicious

couscous      0.9

dates      0.6

pizza      0.2

# Large Language Model



llama-2-70b

parameters → 140GB

run.c → ~500 lines of C code

Image credit: Andrej Karpathy

# Training costs

10 TB of Text (scraped from the internet)

6000 GPUs for 12 days ~2M$*

* Data For Llama 2 - 70 B according to Karpathy

# Parameters

Calculated based on training

Parameters are not interpretable or understandable by human

No one knows how they collaborate

# How do language models work?

In Tunisia, I ate delicious → couscous

# Outline

## What is a language model

couscous     0.9

dates     0.6

pizza     0.2

## Before the transformer



Ich    liebe    Palästina

I    love    Palestine    <SOS>

## Transformer architecture



Output
Probabilities

Softmax

Linear

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

N×

Add & Norm

Masked
Multi-Head
Attention

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

N×

Positional
Encoding

Positional
Encoding

Input
Embedding

Output
Embedding

Inputs

Outputs
(shifted right)

12

# Neural network

# Sequence to sequence models - Translation



Takes a sequence as input and generated another sequence
Example: Recurrent neural network (RNN), Gated Recurrent Units (GRU),
Long-short-term memory (LSTM)

14

# Sequence to sequence models



Encoder

Decoder

# Sequence to sequence models



These models process the sequence word by word. This makes them very expensive to train.

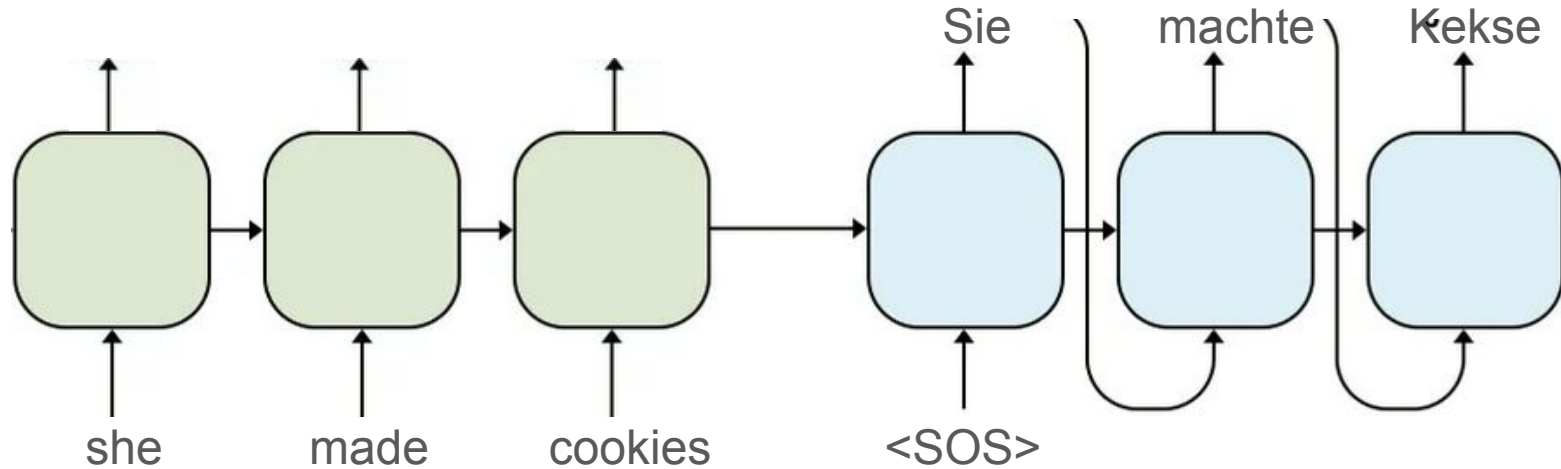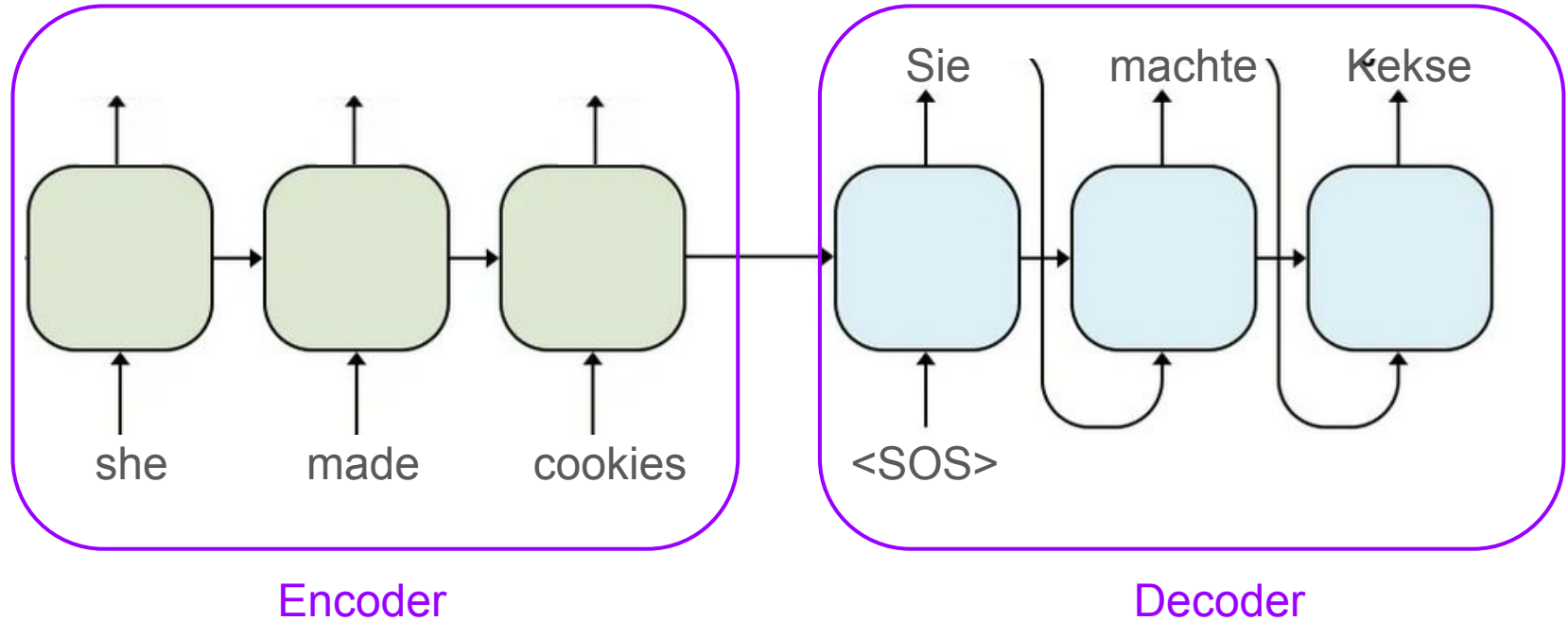# What is a language model

couscous — 0.9

dates — 0.6

pizza — 0.2

# Before the transformer

Ich   liebe   Palästina

I   love   Palestine   <SOS>

# Transformer architecture

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

N×

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

N×

Add & Norm

Masked Multi-Head Attention

Positional Encoding

Input Embedding

Inputs

Positional Encoding

Output Embedding

Outputs (shifted right)

17

# Transformer architecture

"Attention is all you need"
(2017)

# Transformer architecture

Encoder



Input transformed matrix to a Neural Network

Transform input sequence using attention

Add information about word position

Transform input words to vectors

19

# Transformer architecture

Encoder

Input transformed matrix to a Neural Network

Transform input sequence using attention

Add information about word position

Transform input words to vectors

Word Embedding

Positional Encoding

N×

**Add & Norm**

**Feed Forward**

**Add & Norm**

**Multi-Head Attention**

**Input Embedding**

Inputs

20

# Sequence as Matrix

apple

| 3.4 |
|-----|
| 5.2 |
| -0.3 |
| 0.4 |

She made apple cake

| 0.1 | 2.2 | 3.4 | 6.4 |
|-----|-----|------|-----|
| 5.2 | 1.0 | 5.2 | 2.2 |
| 0.7 | 2.3 | -0.3 | 5.3 |
| 2.4 | -1.0 | 0.4 | 3.1 |

# Word embedding

Words are represented as a vectors in a space of n dimensions

Similar words are closer to each other

*Distance is usually calculated using cosine similarity:* cos(θ)

# Word embedding

I made an <u>apple</u> <span style="color:red">juice</span>

I made <u>orange</u> <span style="color:red">juice</span>

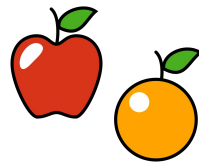I made a <span style="color:purple">fruit</span> salad with <u>apples</u> and <u>oranges</u>

I go to <span style="color:blue">school</span> by <u>car</u>

The <span style="color:blue">school</span> <u>bus</u> arrived on time

# Word embedding

The distance between words reflects the relation between the words

Example: Tokyo to Japan is like Cairo to Egypt

Cairo
Tokyo
Egypt
Tunisia
Japan
Tunis
Man
King
Women
Queen

# Transformer architecture

Encoder

Input transformed matrix to a Neural Network

Attention

Transform input sequence using attention

Add information about word position

Word Embedding

Transform input words to vectors

# Attention mechanism

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

The math is too simple. 😄

But wait…

What's the meaning of multiplying the input 3 times? 🤔

Ideas?

MatMul

SoftMax

Mask (opt.)

Scale

MatMul

Q    K    V

26

# Attention mechanism - Explained

The value of the dot product depends on the cosine distance between the vectors



$$a . b = |a||b| \cos \theta$$

So similar words will have a higher dot product

Sequence: she made apple cake

|  | Food | Tech | Other |
|---|---|---|---|
| cake | 3 | 0 | 0 |
| apple | 2 | 2 | 0 |
| made | 0 | 0 | 2 |
| she | 0 | 0 | 3 |

Q

**Food**

●

| cake | apple | made | she |
|---|---|---|---|
| 3 | 2 | 0 | 0 |
| 0 | 2 | 0 | 0 |
| 0 | 0 | 2 | 3 |

$Q^T$

**Tech**

Example inspired by Serrano.Academy

# First dot product

*Cookies* and *dates* have a higher dot product than *dates* and *she*

The multiplication gives higher values to similar words in the input sequence

|  | cookies | apple | made | she |
|---|---|---|---|---|
| cookies | 9 | 6 | 0 | 0 |
| dates | 6 | 8 | 0 | 0 |
| made | 0 | 0 | 4 | 6 |
| she | 0 | 0 | 6 | 9 |

# Scaling & Softmax

Scaling:

For large values of $d_k$, the dot products grow large in magnitude. To counteract this effect, the dot product is scaled by $1/\sqrt{d_k}$

Softmax: $\dfrac{e^{z_i}}{\sum_{j=1}^{K} e^{z_j}}$

Returns values between 0 and 1 while preserving the order of the elements



$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$
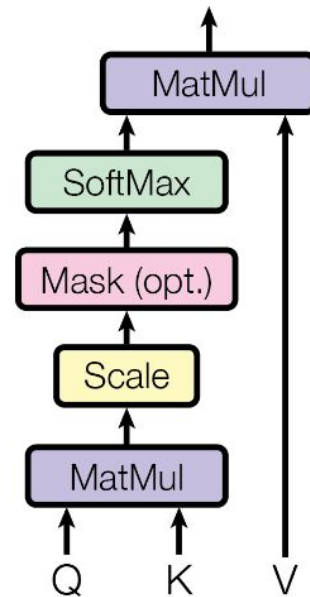
# First dot product

*Orange* and *apple* have a higher dot product than *apple* and *an*

The multiplication gives higher values to similar words in the input sequence

|  | cookies | apple | made | she |
|---|---|---|---|---|
| cookies | 0.8 | 0.2 | 0 | 0 |
| apple | 0.2 | 0.8 | 0 | 0 |
| made | 0 | 0 | 0.6 | 0.3 |
| she | 0 | 0 | 0.4 | 0.5 |

# Second dot product

|  | cake | apple | made | she |
|---|---|---|---|---|
| cake | 0.8 | 0.2 | 0 | 0 |
| apple | 0.2 | 0.8 | 0 | 0 |
| made | 0 | 0 | 0.6 | 0.3 |
| she | 0 | 0 | 0.4 | 0.5 |

●

|  | Food | Tech | Other |
|---|---|---|---|
| cake | 3 | 0 | 0 |
| apple | 2 | 2 | 0 |
| made | 0 | 0 | 2 |
| she | 0 | 0 | 3 |

# Second dot product

|  | cake | apple | made | she |
|------|------|-------|------|-----|
| cake | 0.8 | 0.2 | 0 | 0 |
| apple | 0.2 | 0.8 | 0 | 0 |
| made | 0 | 0 | 0.6 | 0.3 |
| she | 0 | 0 | 0.4 | 0.5 |

●

|  | Food | Tech | Other |
|------|------|------|-------|
| cake | 3 | 0 | 0 |
| apple | 2 | 2 | 0 |
| made | 0 | 0 | 2 |
| she | 0 | 0 | 3 |

=

|  | | |
|------|------|------|
| cake | 2.8 | 0.4 |
| apple | | |
| made | | |
| she | | |

33

# Second dot product

| | cake | apple | made | she |
|---|---|---|---|---|
| cake | 0.8 | 0.2 | 0 | 0 |
| apple | 0.2 | 0.8 | 0 | 0 |
| made | 0 | 0 | 0.6 | 0.3 |
| she | 0 | 0 | 0.4 | 0.5 |

●

| | Food | Tech | Other |
|---|---|---|---|
| cake | 3 | 0 | 0 |
| apple | 2 | 2 | 0 |
| made | 0 | 0 | 2 |
| she | 0 | 0 | 3 |

=

| | | | |
|---|---|---|---|
| cake | 2.8 | 0.4 | 0 |
| apple | 2.2 | 1.6 | 0 |
| made | 0 | 0 | 2.1 |
| she | 0 | 0 | 2.3 |

Words are attracted to each other. Similar words will have more attraction effect on each other than non similar words.
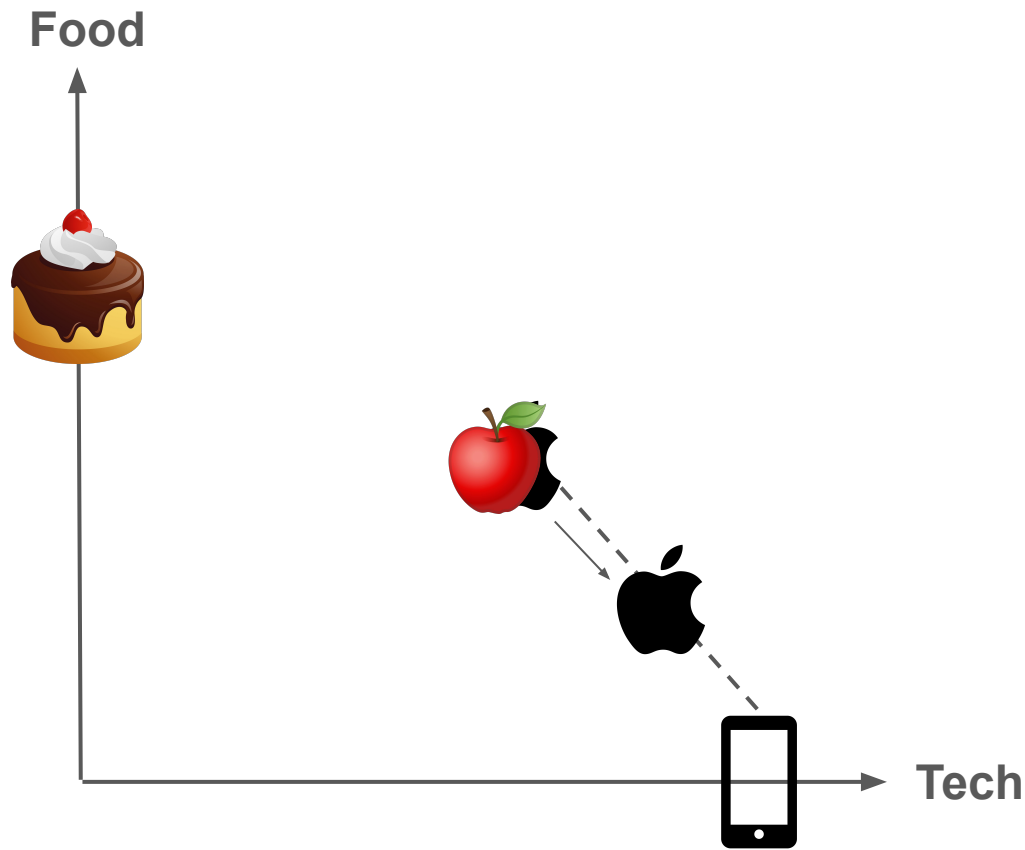
34

# Dot product for context awareness

*She made apple cake*

Using the matrix dot product will move "**apple**" toward "**cake**"



**Food**

**Tech**

Example inspired by Serrano.Academy

# Dot product for context awareness

**Food**

*new phone from apple*

**Tech**

# Dot product provides context awareness

*She made a fruit cake with apples, pears and almond.*
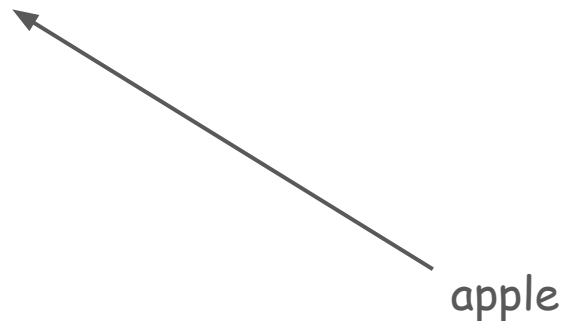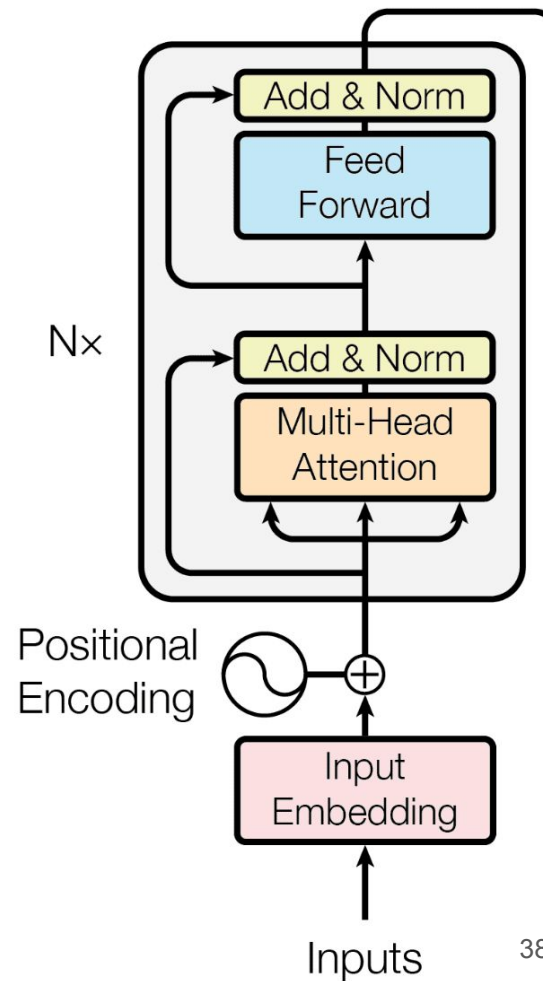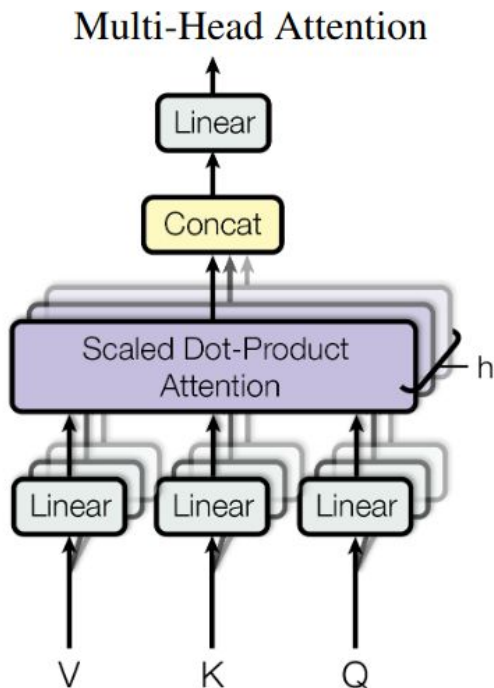
Gravity effect

Similar words attract each other more

Cluster similar words have a stronger attraction effect

**Food**

cake

Almond

pear

fruit

apple

Example inspired by Serrano.Academy
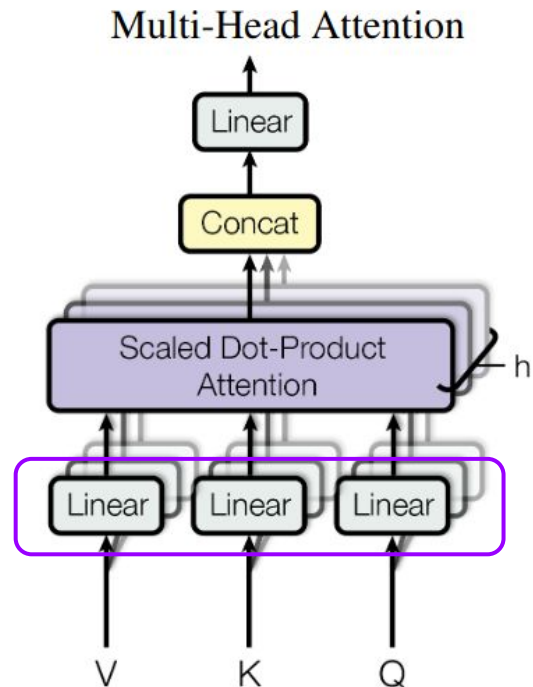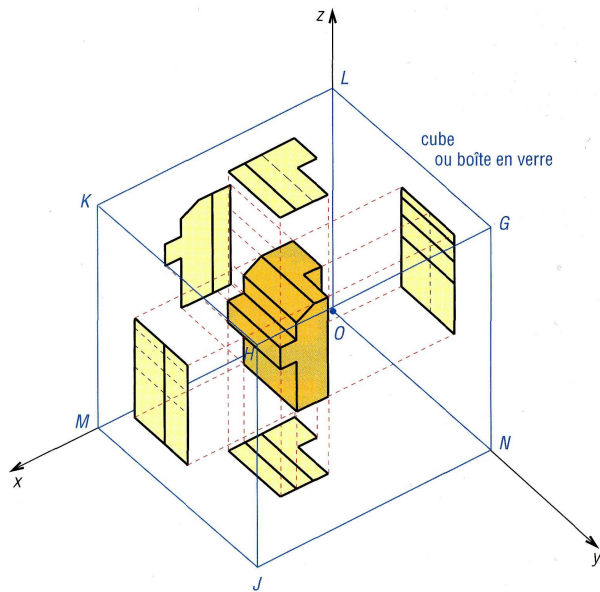
# Multi-Head Attention



Multi-Head Attention

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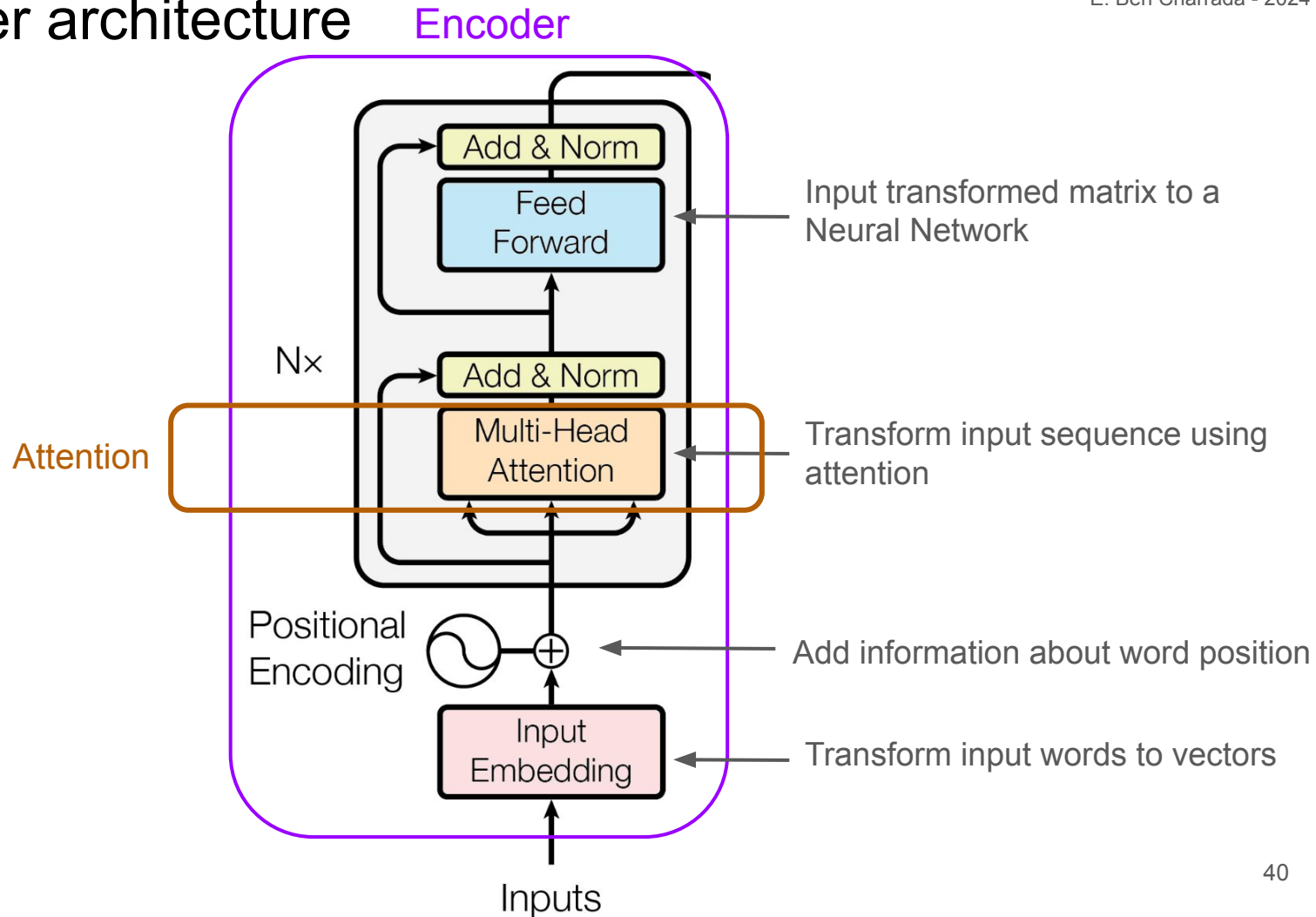
# Projections

V, K and Q are the result of multiplying the sequence S with parameter matrices

The multiplication is a linear projection of the sequence

# Transformer architecture

Encoder

Attention

Input transformed matrix to a Neural Network

Transform input sequence using attention

Add information about word position

Transform input words to vectors

40

# Transformer architecture

Encoder



Input transformed matrix to a Neural Network

Transform input sequence using attention

Add information about word position

Transform input words to vectors

N×

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Positional Encoding

Input Embedding

Inputs

**She made date cookies**

41

# What is a language model

| | |
|---|---|
| couscous | 0.9 |
| dates | 0.6 |
| pizza | 0.2 |

# Before the transformer



# Transformer architecture



Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

Positional Encoding

Output Embedding

Outputs (shifted right)

Nx Decoder

Encoder

Nx

Attention

Positional Encoding

Word Embedding

Input Embedding

Inputs

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

42

# Transformer architecture



Decoder

Masked attention on output sequence

Masking

Add information about word position

Transform input words to vectors

<SOS> Sie machte Apfelkuchen

43

# Masked attention

Mask out dependency with following words by setting it to -∞
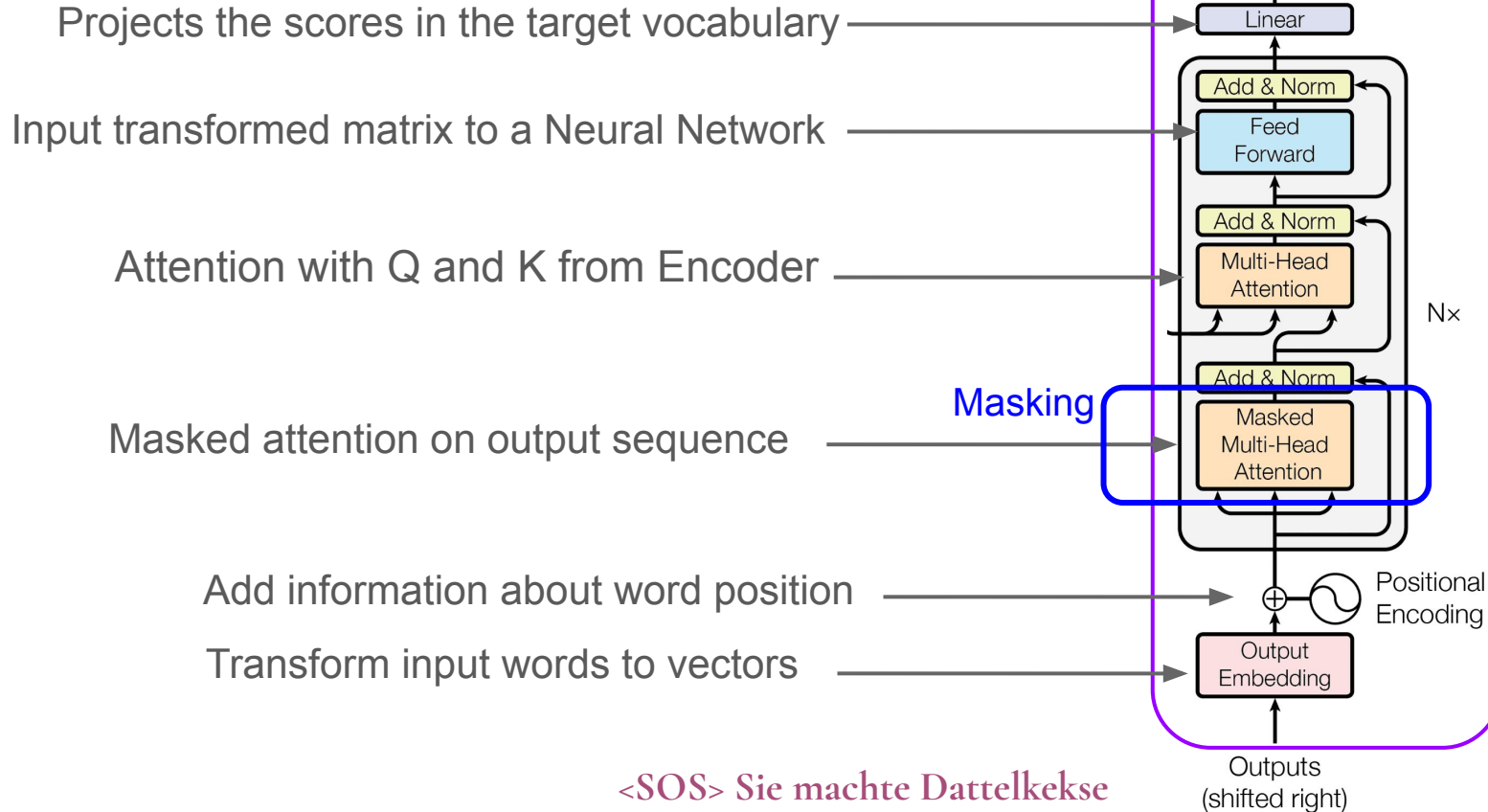
Only previous words are taken into consideration

The softmax will transform the -∞ to 0

Allows processing the whole output sequence at once during learning

|  | <SOS> | Sie | machte | Apfelkuchen |
|---|---|---|---|---|
| <SOS> | 0.8 | -∞ | -∞ | -∞ |
| Sie | 0.1 | 0.7 | -∞ | -∞ |
| machte | 0 | 0.4 | 0.6 | -∞ |
| Apfelkuchen | 0.2 | 0.1 | 0.4 | 0.5 |

# Transformer architecture

Decoder

Projects the scores in the target vocabulary —————————— Linear

Input transformed matrix to a Neural Network ——————— Feed Forward

Attention with Q and K from Encoder ————————— Multi-Head Attention

Masking

Masked attention on output sequence ————————— Masked Multi-Head Attention

Add information about word position ————————→ Positional Encoding

Transform input words to vectors ——————————→ Output Embedding

Output Probabilities

Softmax

Add & Norm

Add & Norm

Add & Norm

Nx

<SOS> Sie machte Dattelkekse

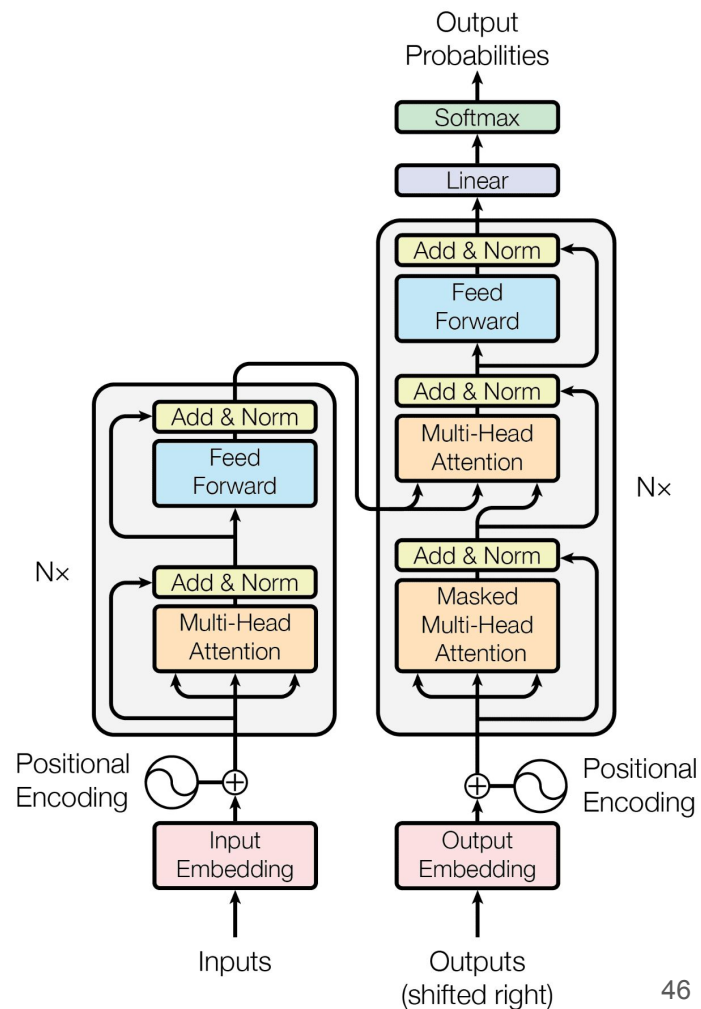Outputs (shifted right)

45

# Costs of attention

Learning:

    Input/output sequences processed in one step
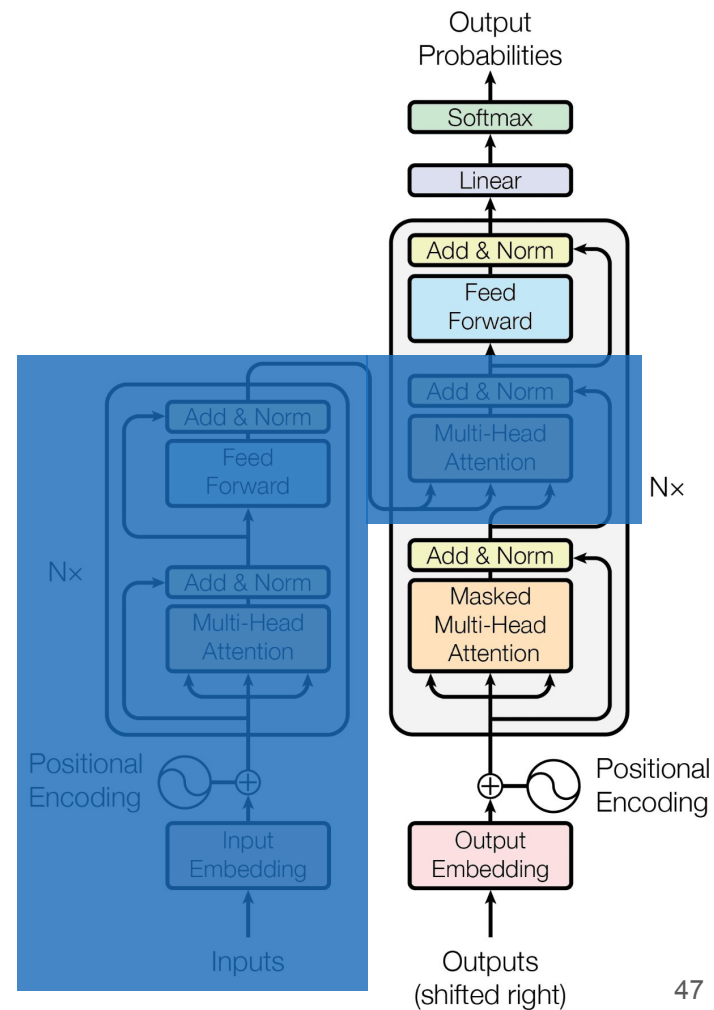
Inference:

    Input sequence processed in one step

    Output inferred token by token
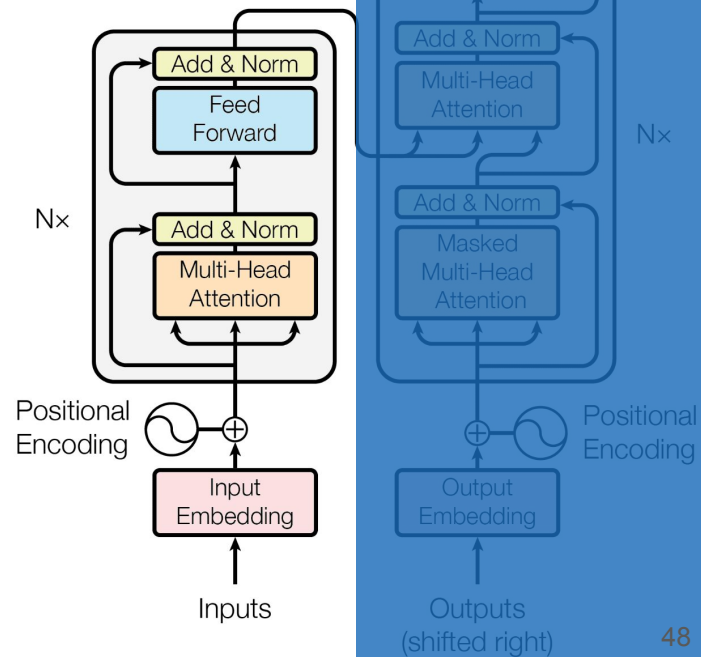
# GPT

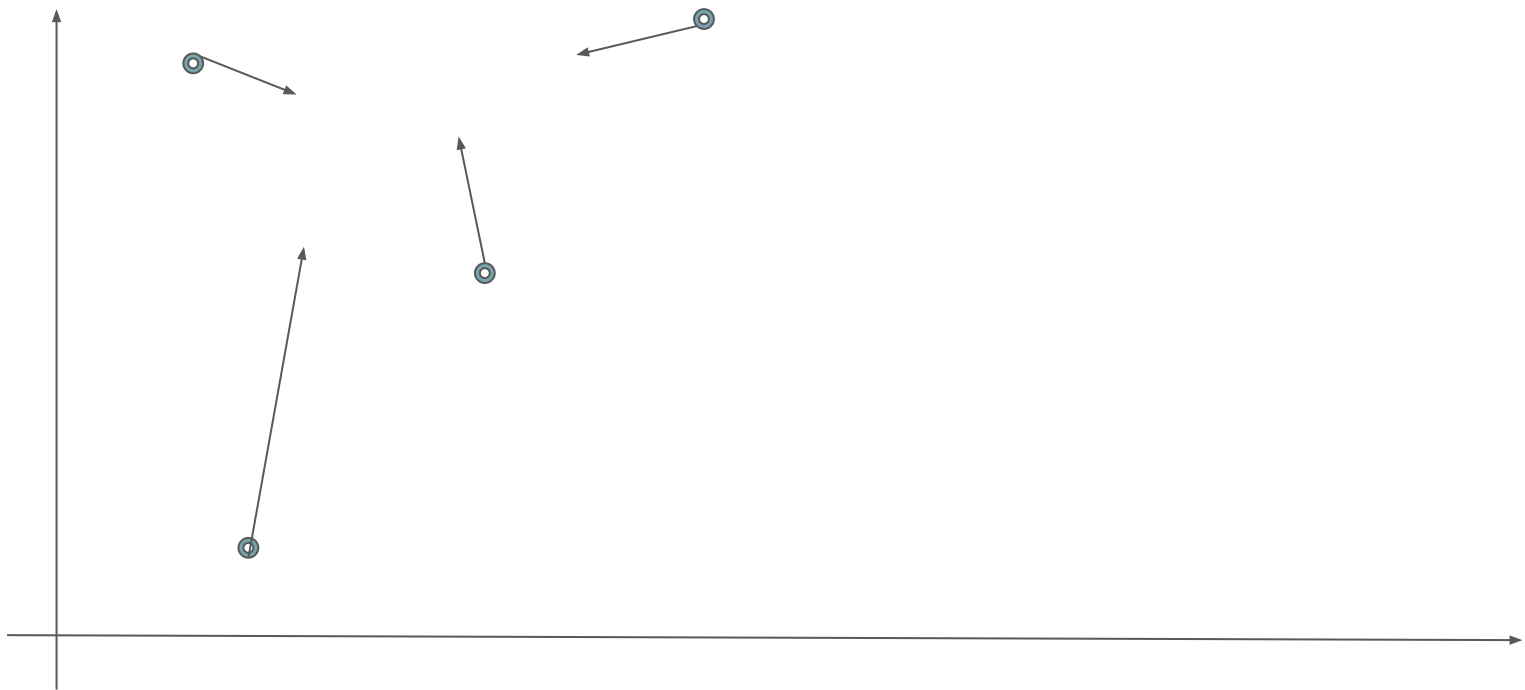Generative Pre-trained Transformer

Decoder only



Output
Probabilities

Softmax

Linear

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

N×

Add & Norm

Masked
Multi-Head
Attention

Positional
Encoding

Output
Embedding

Positional
Encoding

Outputs
(shifted right)

Add & Norm

Feed
Forward

N×

Add & Norm

Multi-Head
Attention

Positional
Encoding

Input
Embedding

Inputs

47

# BERT

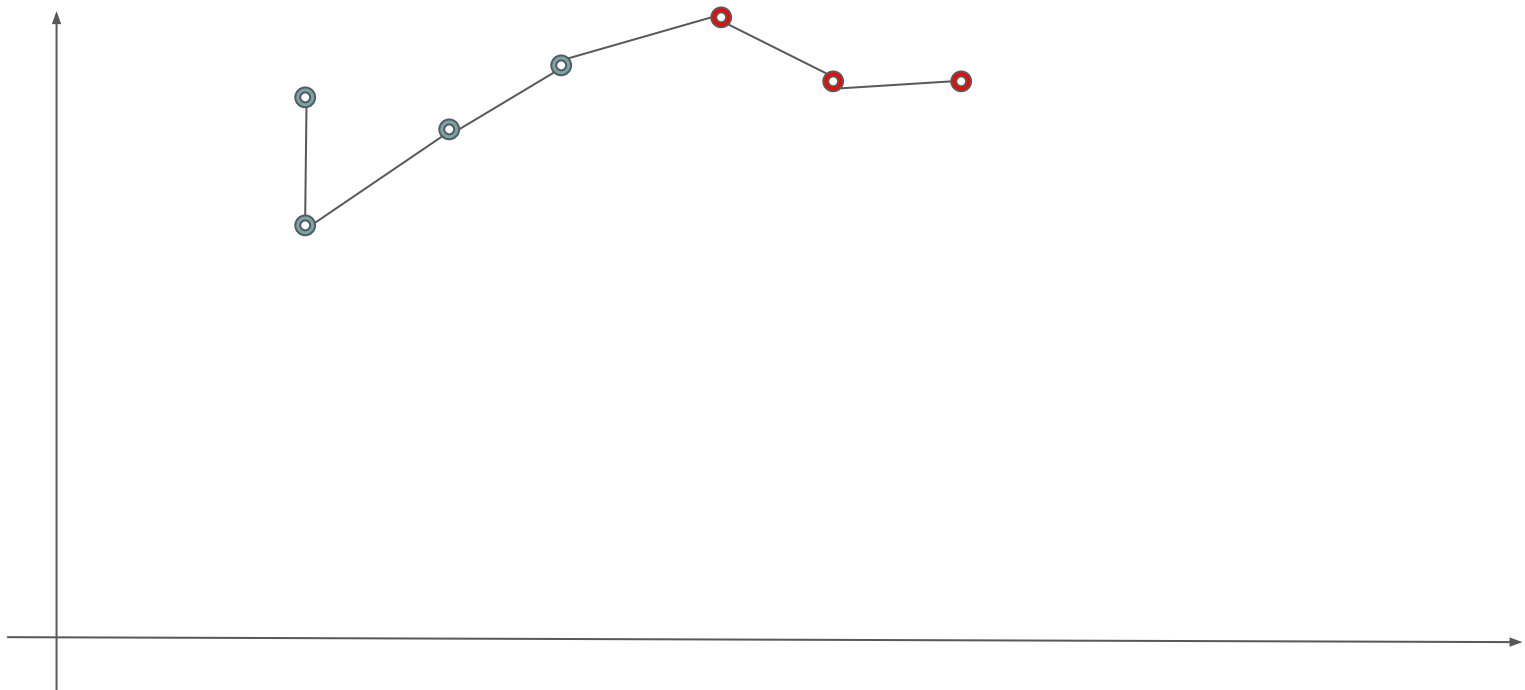Bidirectional Encoder Representations from Transformers

Encoder only

# Smart or Stupid

# Smart or Stupid

# Summary

## What is a language model
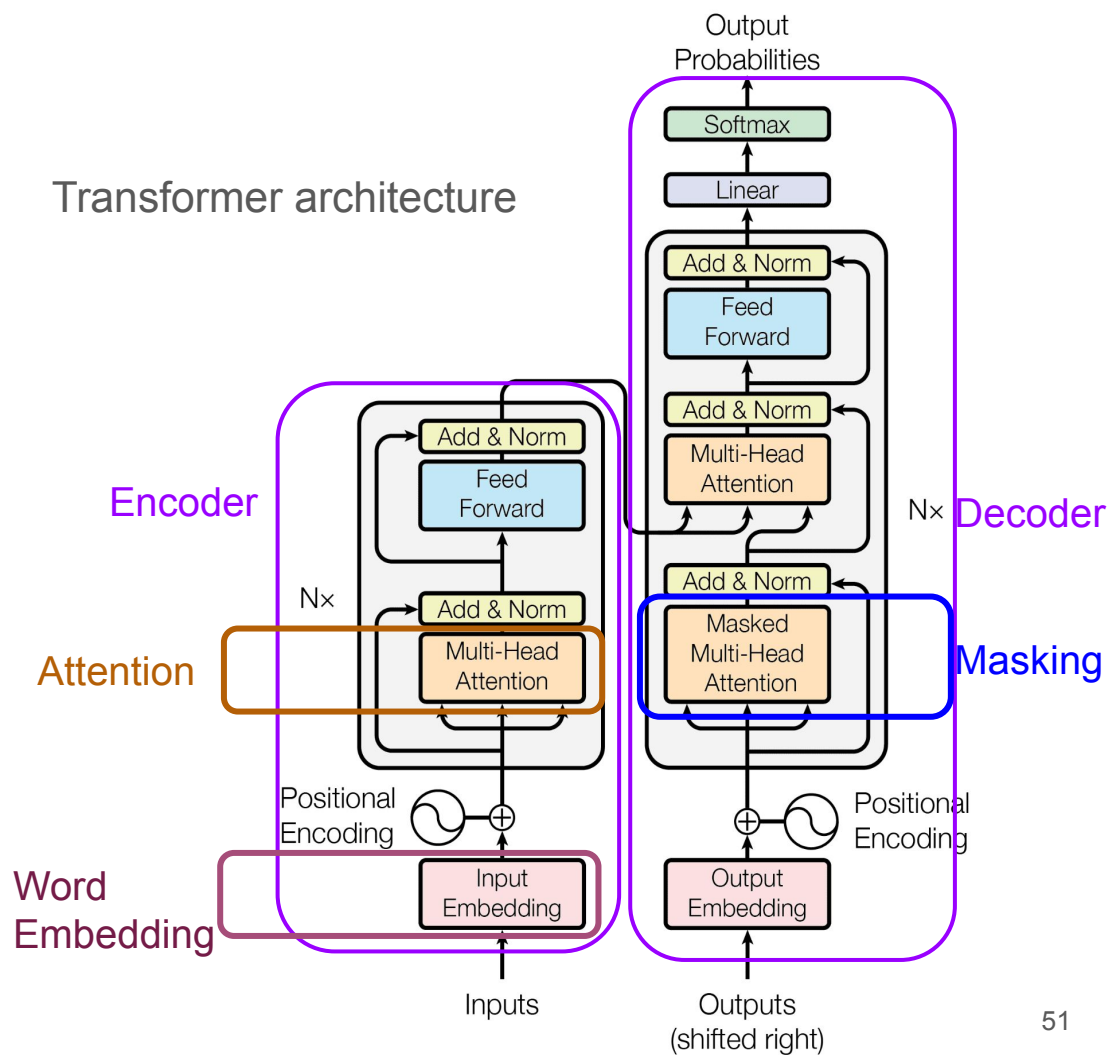


couscous — 0.9
dates — 0.6
pizza — 0.2

## Before the transformer



I  love  Palestine  <SOS>  Ich  liebe  Palästina

## Transformer architecture



Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

Encoder

Decoder

Nx

Attention

Masking

Word Embedding

Positional Encoding

Input Embedding

Output Embedding

Positional Encoding

Inputs

Outputs (shifted right)

# References

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

[1hr Talk] Intro to Large Language Models. Andrej Karpathy
https://youtu.be/zjkBMFhNj_g?si=wpJVQf6ah18LM30z

The Attention Mechanism in Large Language Models. Serrano.Academy
https://youtu.be/OxCpWwDCDFQ?si=qpKl2hgWtgAoEH3n