CS540 Summer 2023

# Attention and Transformers

Jiang, Yuye

PhD student in Computer Sciences

# Participation game (on TopHat)

What is the full name of "GPT"?

A.   Generic Pre-trained Transformers
B.   Generic Pre-trained Tensors
C.   Generative Pre-trained Transformers
D.   Generative Probabilistic Transformer

# Background

- Attention Mechanism
  - Bahdanau et. al. 2014. Neural Machine Translation by Jointly Learning to Align and Translate.
  - Originally developed as an enhancement of RNN applied to translation task
- Transformer Model
  - Vaswani et. al. 2017. Attention is All you Need.
  - First transduction model relying entirely on self-attention to compute representations of its input and output
  - Backbone of the modern large language models & CV models, etc
  - GP**T** = Generative Pre-trained **Transformer**

# Database Manager

The thing your customer is interested in

| Query |
|-------|
| 2 |

Database of (Key, Value) pairs

| Key | Value |
|-----|-------|
| 1 | 100 |
| 2 | 200 |
| 3 | 300 |

# Database Manager

The thing your customer is interested in

| Query |
|-------|
| 2 |

Database of (Key, Value) pairs

| Key | Value |     |
|-----|-------|-----|
| 1 | 100 | unrelated to 2 |
| 2 | 200 | equal to 2 |
| 3 | 300 | unrelated to 2 |

# Database Manager

The thing your customer is interested in

| Query |
|-------|
| 2 |

Database of (Key, Value) pairs

| Key | Value |
|-----|-------|
| 1 | 100 |
| 2 | 200 |
| 3 | 300 |

| | Attention weight |
|---|---|
| unrelated to 2 | 0 |
| equal to 2 | 1 |
| unrelated to 2 | 0 |

$0\% \times 100$

$100\% \times 200$

$0\% \times 300$

| Output |
|--------|
| 200 |

# Database Manager

The thing your customer is interested in

| Query |
|:---:|
| 2.5 |

Database of (Key, Value) pairs

| Key | Value |
|:---:|:---:|
| 1 | 100 |
| 2 | 200 |
| 3 | 300 |

| |
|:---:|
| far from 2.5 |
| close to 2.5 |
| close to 2.5 |

| Attention weight |
|:---:|
| 0 |
| 0.5 |
| 0.5 |

$0\% \times 100$

$50\% \times 200$

$50\% \times 300$

| Output |
|:---:|
| 250 |

# Attention Mechanism

$$Attention(q, D) = \sum_{i=1}^{m} \alpha(q, k_i) v_i$$

$q$

| Query |
|-------|
| 2 |

$D$

$\alpha(q, k)$

$$\sum_{i=1}^{m} \alpha(q, k_i) v_i$$

| Key | Value | | Attention weight | |
|-----|-------|---|------------------|---|
| 1 | 100 | unrelated to 2 | 0 | $0\% \times 100$ |
| 2 | 200 | equal to 2 | 1 | $100\% \times 200$ |
| 3 | 300 | unrelated to 2 | 0 | $0\% \times 300$ |

| Output |
|--------|
| 200 |

# Attention Mechanism

$$Attention(q, D) = \sum_{i=1}^{m} \alpha(q, k_i) v_i$$

# Attention Mechanism

$$Attention(q, D) = \sum_{i=1}^{m} \alpha(q, k_i)v_i$$

$$D$$

$$q \in \mathbb{R}^d$$

| Key | Value |
|-----|-------|

$$k \in \mathbb{R}^d \quad v \in \mathbb{R}^{d'}$$

Query

$$\alpha(q, k)$$

Attention Scoring function $a(q, k)$

Any function that captures similarity between q and k for your task

A common one: scaled dot-product attention

$$a(q, k) = \frac{q^T k}{\sqrt{d}}$$

vector direction more similar → dot-product higher
See: cosine similarity

- Angle θ close to 0
- Cos(θ) close to 1
- **Similar vectors**

- Angle θ close to 90
- Cos(θ) close to 0
- **Orthogonal vectors**

- Angle θ close to 180
- Cos(θ) close to -1
- **Opposite vectors**

# Attention Mechanism

$$Attention(q, D) = \sum_{i=1}^{m} \alpha(q, k_i) v_i$$

$$D$$

$$q \in \mathbb{R}^d \qquad k \in \mathbb{R}^d \quad v \in \mathbb{R}^{d'} \qquad a(q, k_i)$$

| | Key | Value | Attention score |
|---|---|---|---|
| Query | | | -0.9 |
| | | | 0.25 |
| | | | 0.9 |

Attention score does not follow probability distribution

convert this to a probability via softmax function

$$softmax(a(q, k_i)) = \frac{e^{a(q, k_i)}}{\sum_{j=1}^{m} e^{a(q, k_i)}}$$

# Attention Mechanism

$$Attention(q, D) = \sum_{i=1}^{m} \alpha(q, k_i) v_i$$



$$\alpha(q, k_i) = softmax(a(q, k_i))$$

# Attention Mechanism

$$Attention(q, D) = \sum_{i=1}^{m} \alpha(q, k_i) v_i$$

$D$

$q \in \mathbb{R}^d$

| Query |
|-------|

$k \in \mathbb{R}^d$    $v \in \mathbb{R}^{d'}$    $\alpha(q, k_i)$

| Key | Value | Attention weight |
|-----|-------|------------------|
|     |       | 0.098 |
|     |       | 0.309 |
|     |       | 0.593 |

$9.8\% \times v_1$

$30.9\% \times v_2$

$59.3\% \times v_3$

| Output |
|--------|

# Attention Mechanism



$$Attention(q, K, V) = softmax(\frac{qK^T}{\sqrt{d}})V$$

# Attention Mechanism



$$q \in \mathbb{R}^d \qquad K \in \mathbb{R}^{m \times d} \quad V \in \mathbb{R}^{m \times d'}$$

$$\text{softmax}\left( \frac{q \times K^T}{\sqrt{d}} \right) V = \text{output}$$

But how are these constructed?

$$Attention(q, K, V) = softmax\left(\frac{qK^T}{\sqrt{d}}\right)V$$

# Machine Translation

# Embedding

A cat is sleeping on a red sofa

A dog is sitting on a green chair

one-hot encoding:

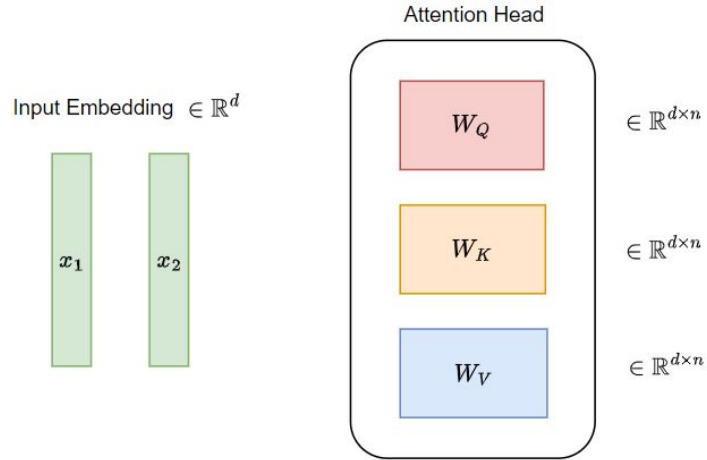| vocab | | cat = | | dog = | |
|---|---|---|---|---|---|
| | a | | 0 | | 0 |
| | cat | | 1 | | 0 |
| | is | | 0 | | 0 |
| | sleeping | | 0 | | 0 |
| | on | | 0 | | 0 |
| | red | | 0 | | 0 |
| | sofa | | 0 | | 0 |
| | dog | | 0 | | 1 |
| | sitting | | 0 | | 0 |
| | green | | 0 | | 0 |
| | chair | | 0 | | 0 |

Can be used for numerical computation

No similarity measurement

cannot tell "cat" & "dog" similar
and tend to appear in similar context &
position

# Embedding

2-d embedding

A cat is sleeping on a red sofa

A dog is sitting on a green chair

x2

chair
sofa

dog = [0.6, 0.6]
cat = [0.5, 0.5]

is
a on
x1
sleeping = [0.5, -0.5]
sitting = [0.6, -0.6]
red
green

Each word is converted to a vector

- can input to neural network
- can learn similarity between words

For the many words in languages, usually pick
a very large embedding dimension,
for example d = 768

# Embedding



Vocabulary

Word Embedding Layer

yes

I

...

today

hello

$|V|$

$\in \mathbb{R}^{|V| \times d}$

corpus

BPE
Sentence Piece
...

Learnable and updated with the model

d: embedding dimension, usually picked to be a high number like 768

# Embedding

# Embedding

# Embedding

# Embedding

# Embedding

# Self-attention

| When | people | go | camping | they | collect | log | to | build | a | bonfire |
|------|--------|-----|---------|------|---------|-----|-----|-------|---|---------|

"Every word needs to pay attention to each other"

"Every word should pay more attention to the other word thats is related to it"

# Reminder of our initial question



$q \in \mathbb{R}^d$  $K \in \mathbb{R}^{m \times d}$  $V \in \mathbb{R}^{m \times d'}$

But how are these constructed?

$$Attention(q, K, V) = softmax(\frac{qK^T}{\sqrt{d}})V$$

# Attention Head

Attention Head

Input Embedding $\in \mathbb{R}^d$

$x_1$ $x_2$

$W_Q$ $\in \mathbb{R}^{d \times n}$

$W_K$ $\in \mathbb{R}^{d \times n}$

$W_V$ $\in \mathbb{R}^{d \times n}$

# Attention Head

Multiply the first embedding vector with each matrix

Attention Head

Input Embedding $\in \mathbb{R}^d$

$x_1$

$W_Q$

$W_K$

$W_V$

$x_1^T W_Q = q_1$

$x_1^T W_K = k_1$

$x_1^T W_V = v_1$

$q_1 \in \mathbb{R}^n$

$k_1 \in \mathbb{R}^n$

$v_1 \in \mathbb{R}^n$

# Attention Head

Multiply the second embedding vector with each matrix

# Attention Head



The query, key and value for each word is calculated from all the m words in the same sentence, using shared learnable matrices

# Attention Head



Attention Head

Input Embedding $\in \mathbb{R}^d$

$x_1$  $x_2$

sequence length = m

$W_Q$  $W_K$  $W_V$

$q_1$ $q_2$ $Q$
$k_1$ $k_2$ $K$
$v_1$ $v_2$ $V$

$$Attention(q, K, V) = softmax\left(\frac{qK^T}{\sqrt{d}}\right)V$$

$q_1$  $K^T$  $V$  output

$$softmax\left(\frac{q_1 \times K^T}{\sqrt{d}}\right) V = $$

$q_2$  $K^T$  $V$  output

$$softmax\left(\frac{q_2 \times K^T}{\sqrt{d}}\right) V = $$

In Matrix form:

$Q$  $K^T$  $V$  $O$

$$softmax\left(\frac{Q \times K^T}{\sqrt{d}}\right) V = $$

# Attention Head

# Attention Head



$$\text{softmax} \left( \frac{Q \times K^T}{\sqrt{d}} \right) V = O$$

$x_1 \rightarrow x_1$   $x_1 \rightarrow x_2$

| 0.88 | 0.12 |
|------|------|
| 0.25 | 0.75 |

Attention Weight matrix $\in \mathbb{R}^{m \times m}$

$x_2 \rightarrow x_1$   $x_2 \rightarrow x_2$

# Attention Head



Attention Head

Input Embedding $\in \mathbb{R}^d$

$x_1$ $x_2$

sequence length = m

$W_Q$

$W_K$

$W_V$

learnable parameters

Output Embedding $\in \mathbb{R}^n$

$x_1$ $x_2$

sequence length = m

# Multi-headed Attention

# Multi-headed Attention

# Multi-headed Attention

# Batch of input sentences

Good morning
How are you
Hello

$\longrightarrow$

Good morning [PAD]
How   are        you
Hello [PAD]    [PAD]

$\longrightarrow$

Batch = 3

$\longrightarrow$

| Good | morning | [PAD] |
|------|---------|-------|
| How | are | you |
| Hello | [PAD] | [PAD] |

$\in \mathbb{R}^{b \times maxlen \times d}$

Embedding layer

I,
am,
a,
student

[PAD]

$\longrightarrow$

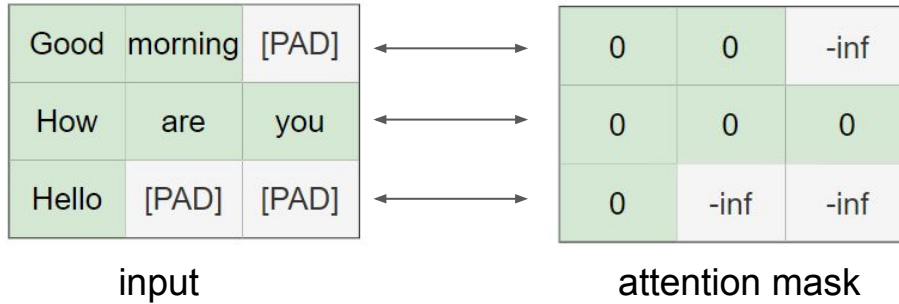| I | $\in \mathbb{R}^d$ |
| am | $\in \mathbb{R}^d$ |
| a | $\in \mathbb{R}^d$ |
| student | $\in \mathbb{R}^d$ |

Insert padding

maxlen = maximum sequence length
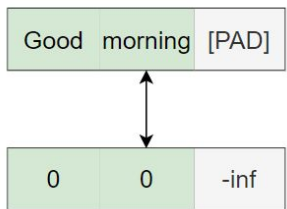in this batch

Don't want padding to affect training

# Attention Mask

| Good | morning | [PAD] |
|------|---------|-------|
| How | are | you |
| Hello | [PAD] | [PAD] |

input

| 0 | 0 | -inf |
|---|---|------|
| 0 | 0 | 0 |
| 0 | -inf | -inf |

attention mask

-inf represents a very small negative number

# Attention Mask

# Attention Mask

Good | morning | [PAD]

0 | 0 | -inf

Before softmax

$$\text{softmax} \left( \frac{Q \times K^T}{\sqrt{d}} \right) V = O$$

$x_1 \to x_1 \quad x_1 \to x_2 \quad x_1 \to x_3$

| 1.2 | 4.8 | -0.1 |
| 3.4 | -9.3 | 2.2 |
| 2.3 | -9.5 | 3.0 |

Attention **Score** matrix $\in \mathbb{R}^{m \times m}$

| 0 | 0 | -inf |
| 0 | 0 | -inf |
| 0 | 0 | -inf |

Attention Mask $\in \mathbb{R}^{m \times m}$

+

$$\text{softmax} \left( \begin{array}{ccc} 1.2 & 4.8 & -inf \\ 3.4 & -9.3 & -inf \\ 2.3 & -9.5 & -inf \end{array} \right) = \begin{array}{ccc} 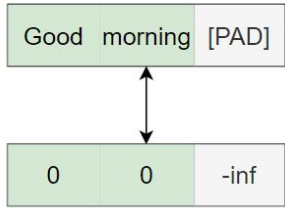0.26 & 0.97 & 9e\text{-}4346 \\ 0.99 & 3e\text{-}6 & 3e\text{-}4345 \\ 0.99 & 7e\text{-}6 & 1e\text{-}4344 \end{array}$$

Attention **Weight** matrix $\in \mathbb{R}^{m \times m}$

what about pad's attention to other tokens?

# Batch of input sentences
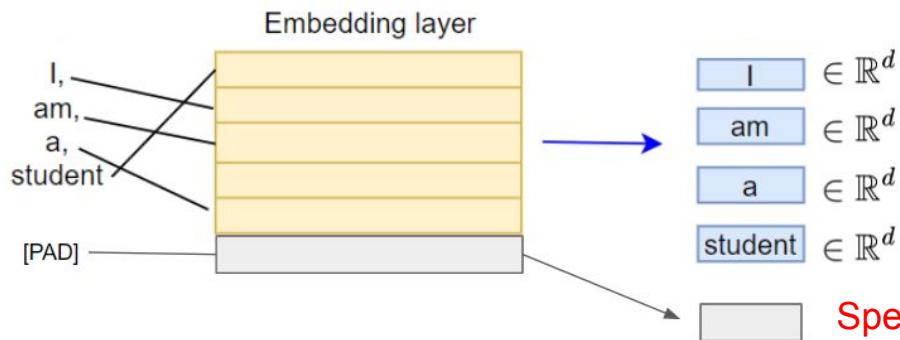
Good morning
How are you
Hello

→

Good morning [PAD]
How   are        you
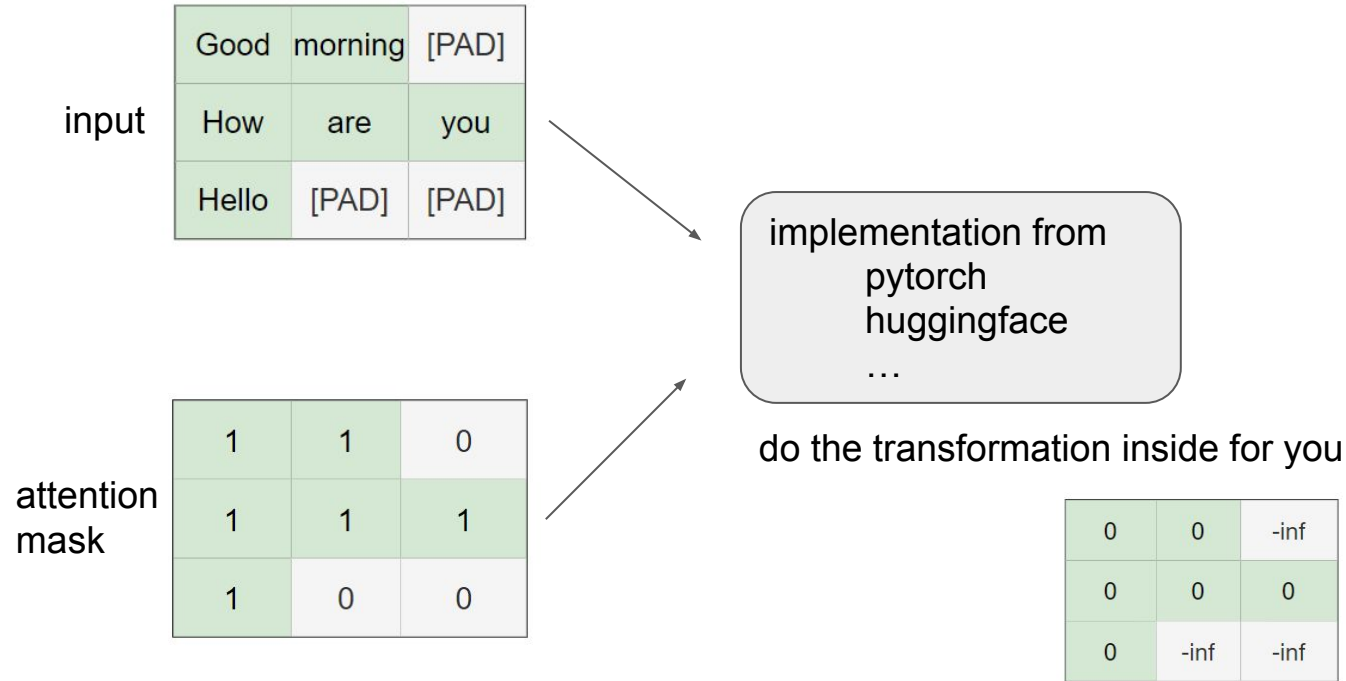Hello [PAD]    [PAD]

→ Batch = 3 →

| Good | morning | [PAD] |
|------|---------|-------|
| How  | are     | you   |
| Hello | [PAD]  | [PAD] |

$\in \mathbb{R}^{b \times maxlen \times d}$

Embedding layer

I,
am,
a,
student

[PAD]

→

| I | $\in \mathbb{R}^d$ |
| am | $\in \mathbb{R}^d$ |
| a | $\in \mathbb{R}^d$ |
| student | $\in \mathbb{R}^d$ |

Special vector
The gradient is not calculated

# Attention Mask in libraries
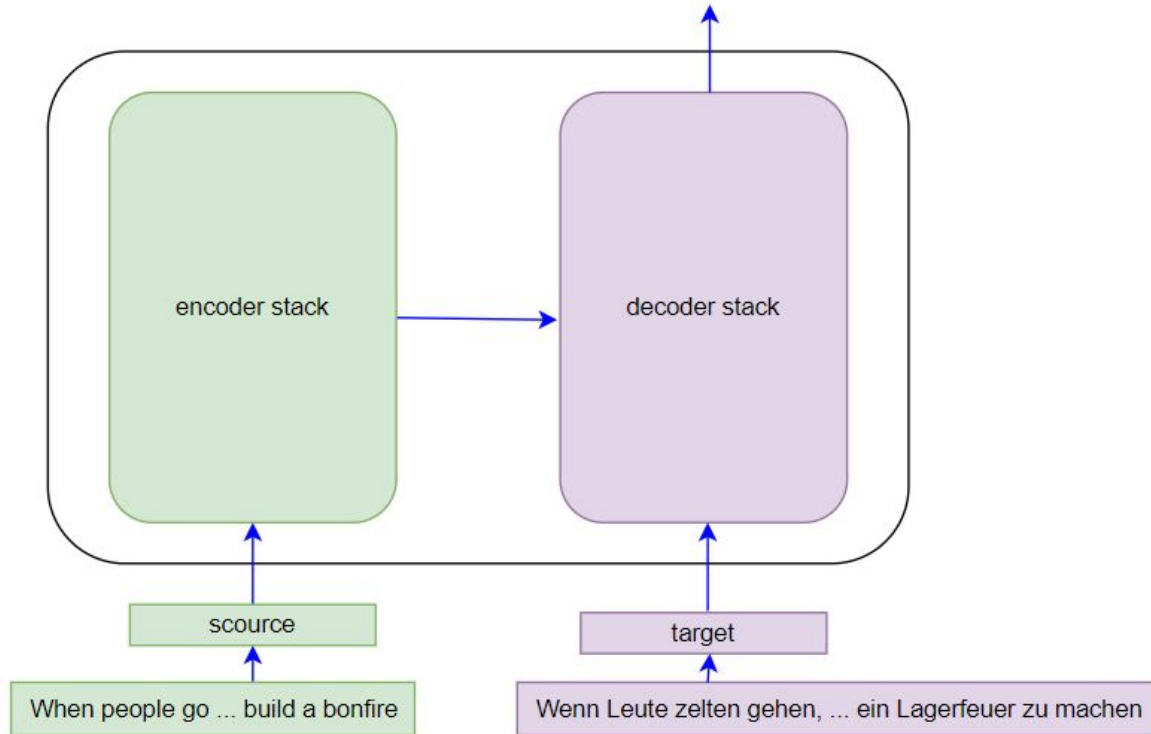
# Machine Translation

# The Transformer Model

At Testing Time:

# The Transformer Model

At Training Time:
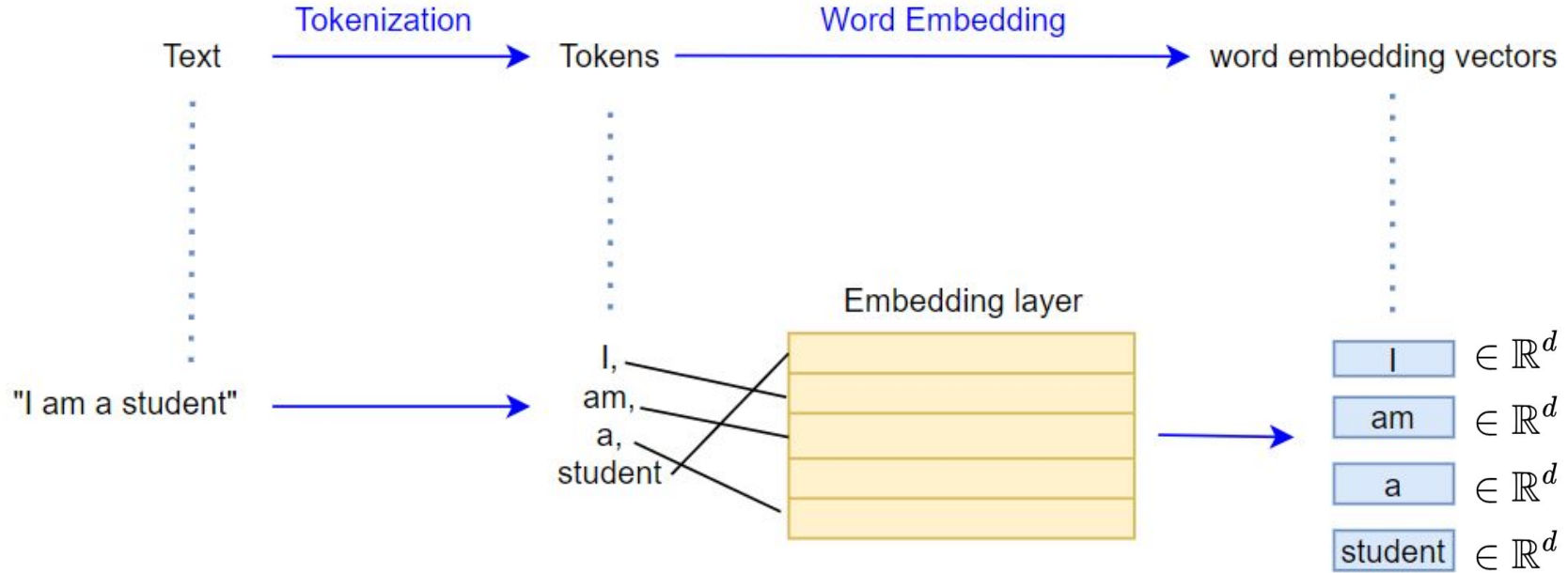
# Input Embedding



1. word embedding

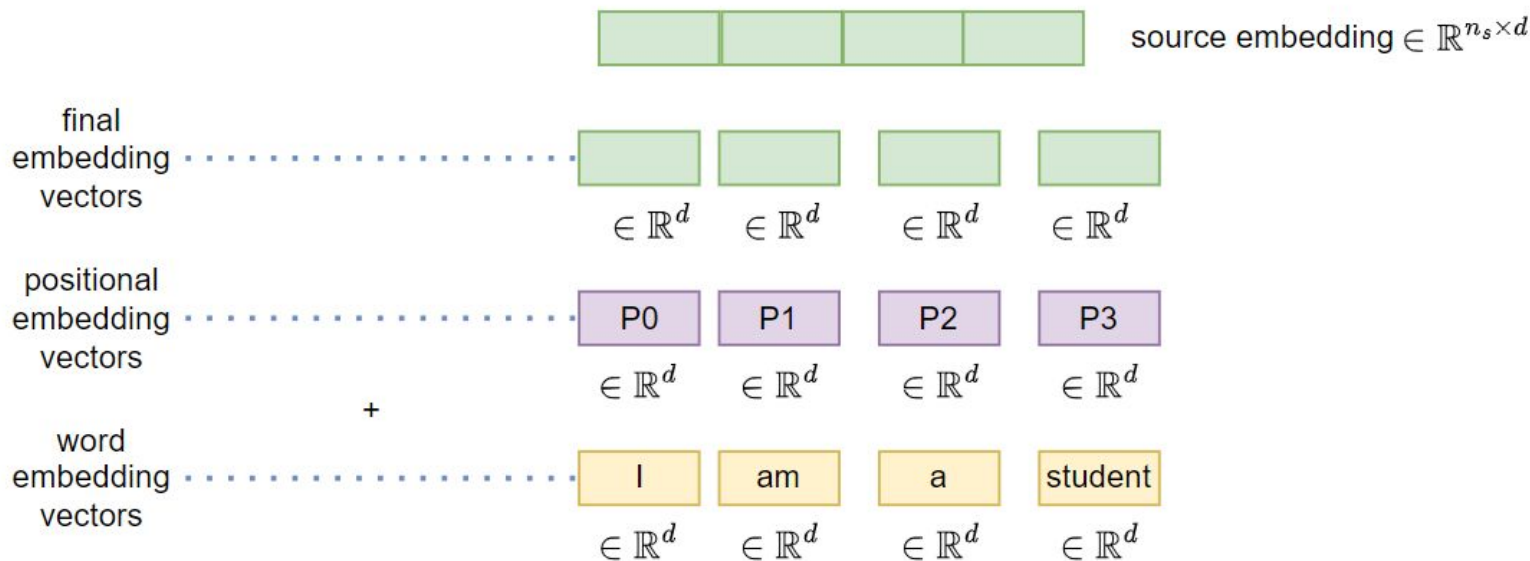2. positional embedding

# Word Embedding

# Positional Embedding

order of words matters in some language

A attacked B.

B attacked A.

# Positional Embedding

Way 1: learned



**Vocabulary**

| yes |
| I |
| ... |
| today |
| hello |

$|V|$

**Word Embedding Layer**

$\in \mathbb{R}^{|V| \times d}$

**Indices**

| 0 |
| 1 |
| ... |
| P-2 |
| P-1 |

$|P|$
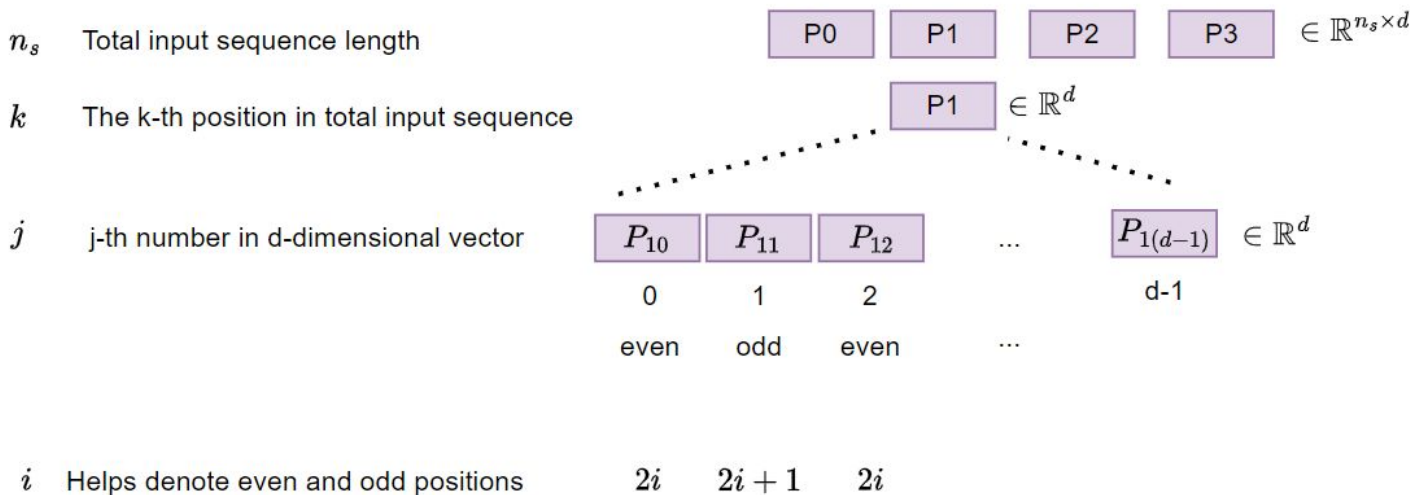
**Positional Embedding Layer**

$\in \mathbb{R}^{|P| \times d}$

P is a pre-decoded maximum length
The model cannot accept length > P, unless trimmed

# Positional Embedding

Way 2: calculated



$n_s$    Total input sequence length

$k$    The k-th position in total input sequence

$j$    j-th number in d-dimensional vector

$i$    Helps denote even and odd positions

Even elements    $P(k, 2i) = sin(\frac{k}{n^{2i/d}})$

Odd elements    $P(k, 2i+1) = cos(\frac{k}{n^{2i/d}})$

$n$   user defined large scalar, = 10,000 in paper

# Positional Embedding

Even elements
$$P(k, 2i) = sin(\frac{k}{n^{2i/d}})$$

$n$ user defined large scalar, = 10,000 in paper

Odd elements
$$P(k, 2i + 1) = cos(\frac{k}{n^{2i/d}})$$

| P0 | $\in \mathbb{R}^4$ |
|---|---|

$k = 0$
$n = 100$

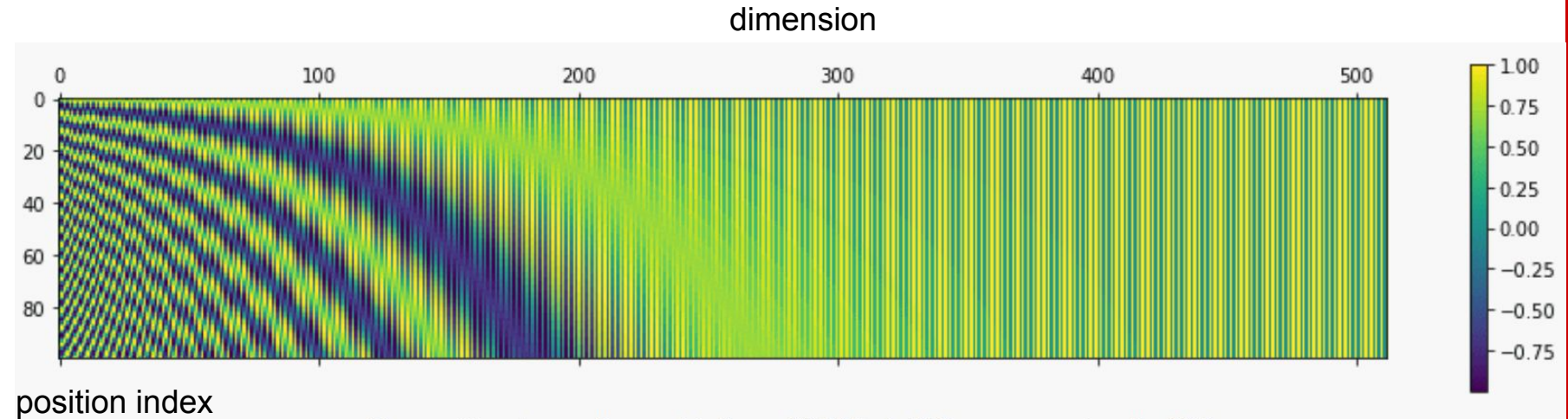| $P_{00}$ | $j = 0, i = 0$ | $P(k, 0) = sin(0) = 0$ |
|---|---|---|
| $P_{01}$ | $j = 1, i = 0$ | $P(k, 1) = cos(0) = 1$ |
| $P_{02}$ | $j = 2, i = 1$ | $P(k, 2) = sin(0) = 0$ |
| $P_{03}$ | $j = 3, i = 1$ | $P(k, 3) = cos(0) = 1$ |

| P0 | = [0,1,0,1] |
|---|---|

| P1 | $\in \mathbb{R}^4$ |
|---|---|

$k = 1$
$n = 100$

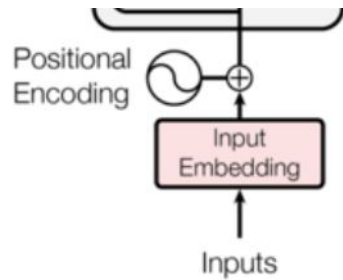| $P_{10}$ | $j = 0, i = 0$ | $P(k, 0) = sin(1/n^{\frac{0}{4}}) = sin(1) \approx 0.84$ |
|---|---|---|
| $P_{11}$ | $j = 1, i = 0$ | $P(k, 1) = cos(1/n^{\frac{0}{4}}) = cos(1) \approx 0.54$ |
| $P_{12}$ | $j = 2, i = 1$ | $P(k, 2) = sin(1/100^{\frac{2}{4}}) = sin(1/10) \approx 0.10$ |
| $P_{13}$ | $j = 3, i = 1$ | $P(k, 3) = cos(1/100^{\frac{2}{4}}) = cos(1/10) \approx 1.0$ |

| P1 | = [0.84,0.54,0.1,1] |
|---|---|

# Positional Embedding

dimension



The positional encoding matrix for n=10,000, d=512, sequence length=100

position index

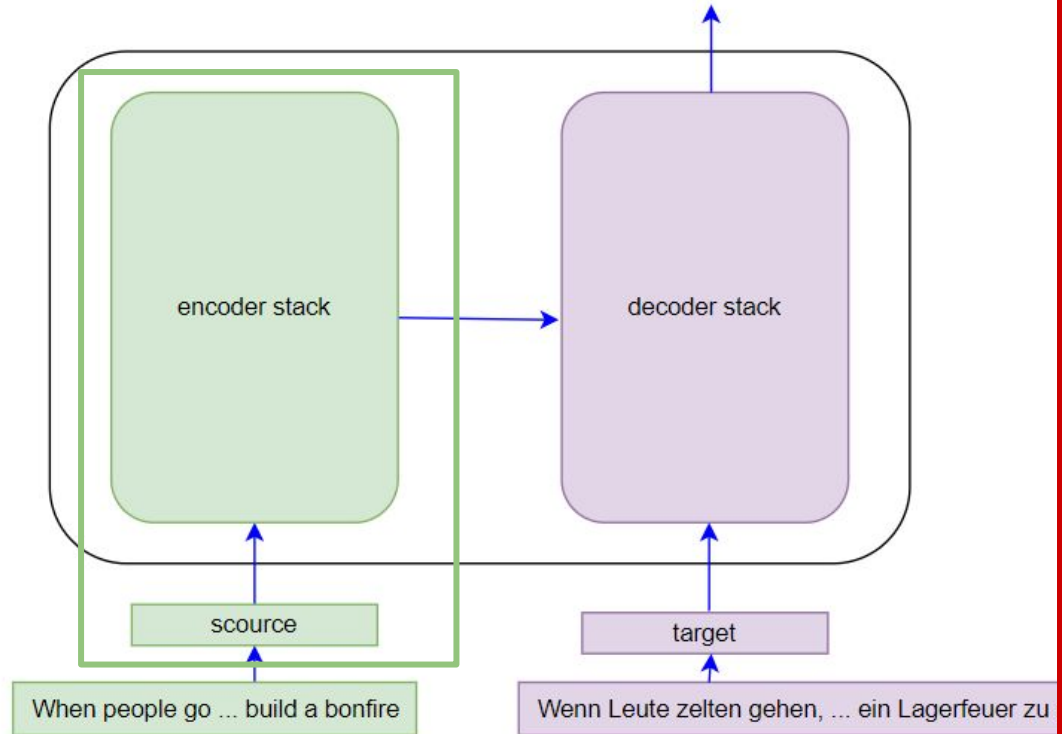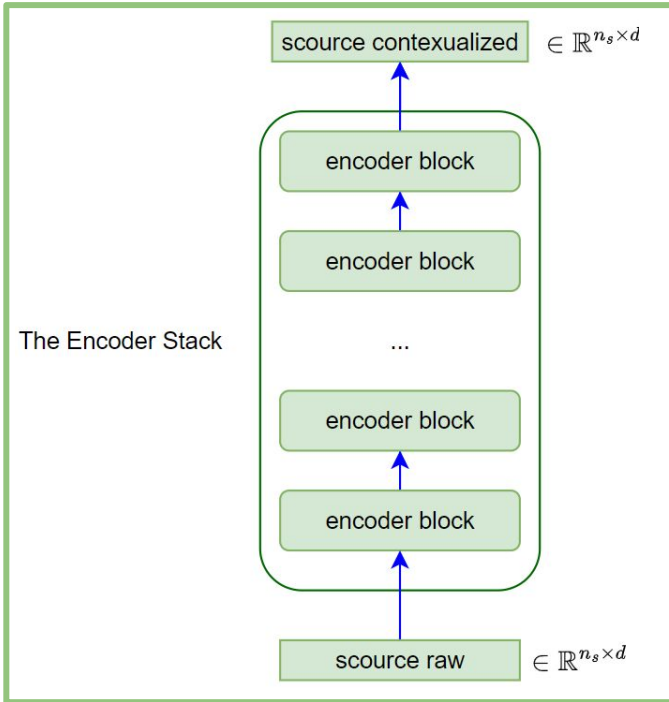# Input Embedding



Positional
Encoding

Input
Embedding

Inputs

In original paper

1. word embedding

2. positional embedding

encoder stack

decoder stack

scource

When people go ... build a bonfire

target

Wenn Leute zelten gehen, ... ein Lagerfeuer zu machen

# The Encoder Stack



The Encoder Stack

scource contextualized $\in \mathbb{R}^{n_s \times d}$

encoder block

encoder block

...

encoder block

encoder block

scource raw $\in \mathbb{R}^{n_s \times d}$

encoder stack

decoder stack

scource

target

When people go ... build a bonfire

Wenn Leute zelten gehen, ... ein Lagerfeuer zu

# The Encoder Block



scource out $\in \mathbb{R}^{n_s \times d}$

encoder block

Add & Norm

feed-forward layer

Add & Norm

self-attention layer

encoder block
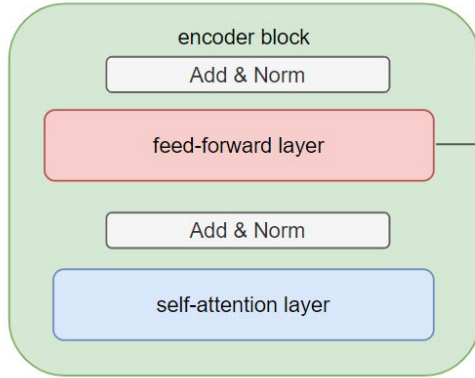
scource raw $\in \mathbb{R}^{n_s \times d}$
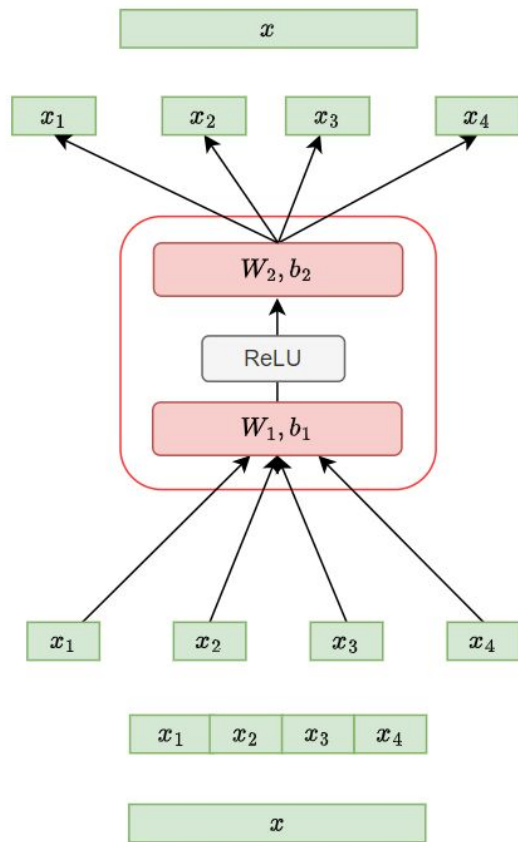
# Multi-headed Attention

# Position-wise Feed-forward Layer



$$x_{out} = (ReLU(x_{in}^T W_1 + b_1))W_2 + b_2$$
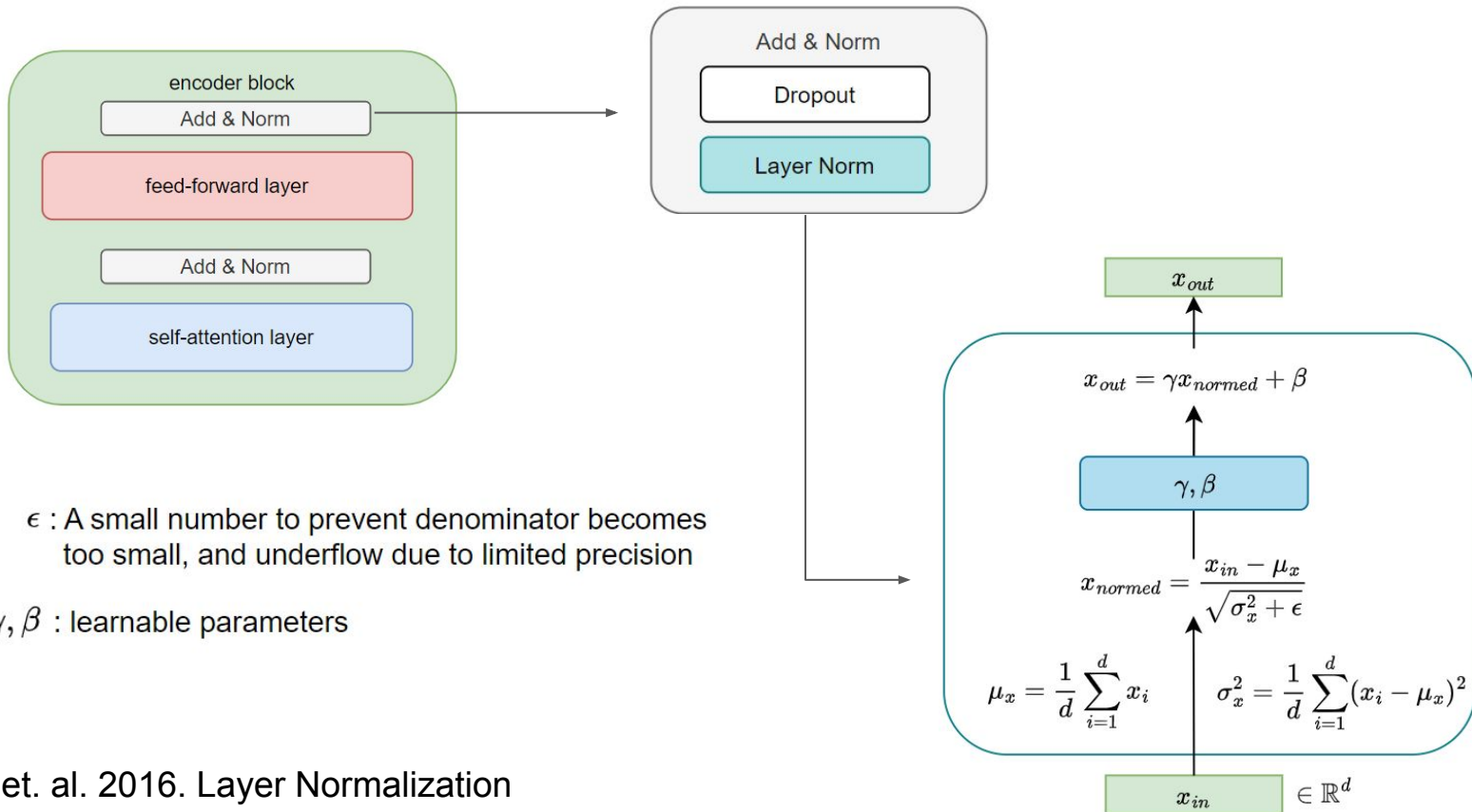
# Position-wise Feed-forward Layer



$$x_{out} = (ReLU(x_{in}^T W_1 + b_1))W_2 + b_2$$

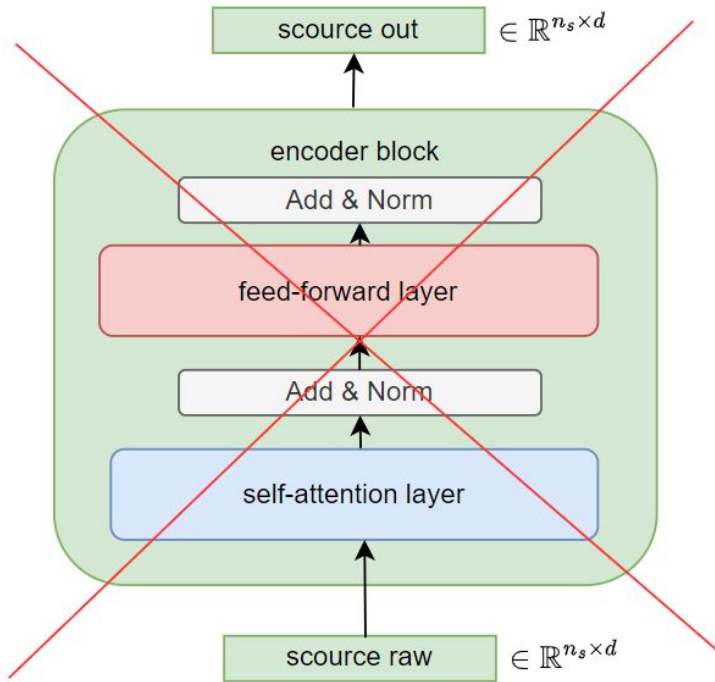"x_in" is each individual word
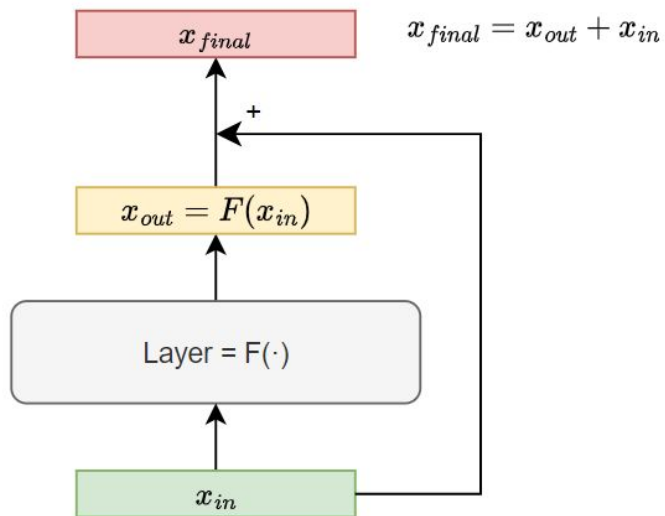
not a whole (n_s) * d sequence

# Add & Norm



$\epsilon$ : A small number to prevent denominator becomes too small, and underflow due to limited precision

$\gamma, \beta$ : learnable parameters

Ba et. al. 2016. Layer Normalization

Diagram labels:

encoder block
- Add & Norm
- feed-forward layer
- Add & Norm
- self-attention layer

Add & Norm
- Dropout
- Layer Norm

$x_{out}$

$x_{out} = \gamma x_{normed} + \beta$

$\gamma, \beta$

$x_{normed} = \dfrac{x_{in} - \mu_x}{\sqrt{\sigma_x^2 + \epsilon}}$

$\mu_x = \dfrac{1}{d}\sum_{i=1}^{d} x_i \qquad \sigma_x^2 = \dfrac{1}{d}\sum_{i=1}^{d}(x_i - \mu_x)^2$

$x_{in} \in \mathbb{R}^d$

# The Encoder Block



NOT simply pass through one layer after another

# Residual Connection



$$x_{final} = x_{out} + x_{in}$$

$$x_{out} = F(x_{in})$$

Layer = F(·)

$$x_{in}$$

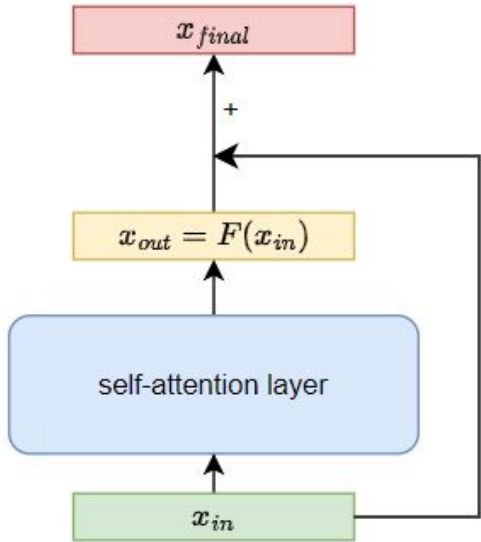$$x_{final}$$

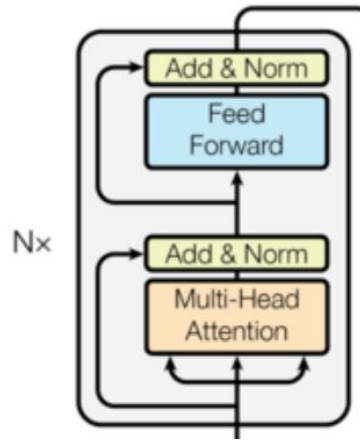He et. al. 2015. Deep Residual Learning for Image Recognition

# Residual Connection

# Full encoder block

In original paper

# Full encoder stack
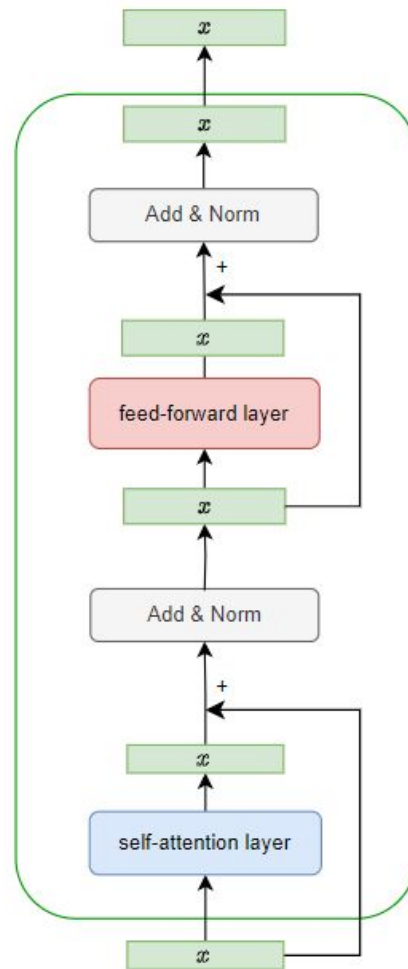
scource contexualized $\in \mathbb{R}^{n_s \times d}$

The Encoder Stack

encoder block

encoder block

...

encoder block

encoder block

scource raw $\in \mathbb{R}^{n_s \times d}$

$x$

$x$

Add & Norm

+

$x$

feed-forward layer

$x$

Encoder Block

Add & Norm

+

$x$

self-attention layer

$x$

# Full encoder stack

# The Transformer Model

At Training Time:

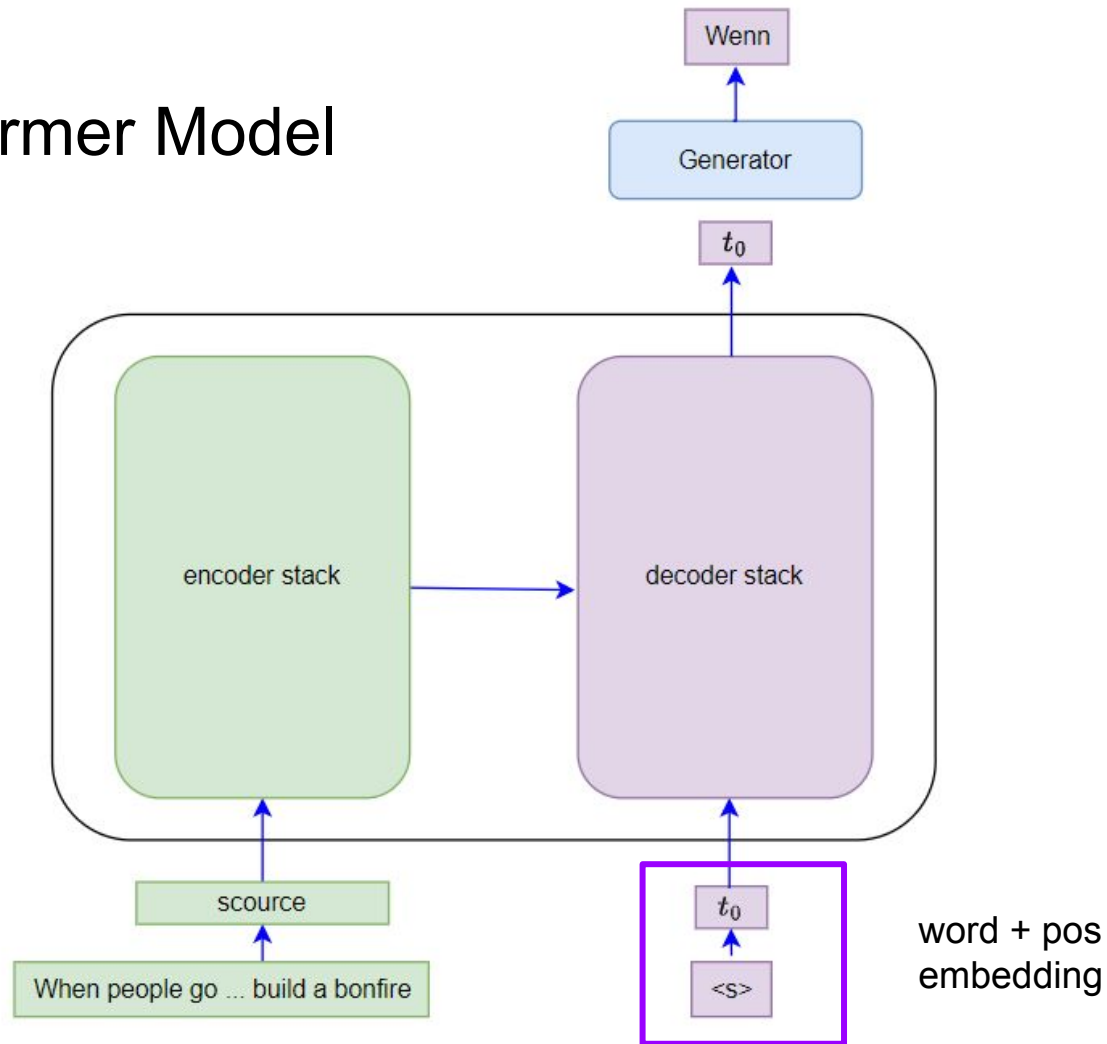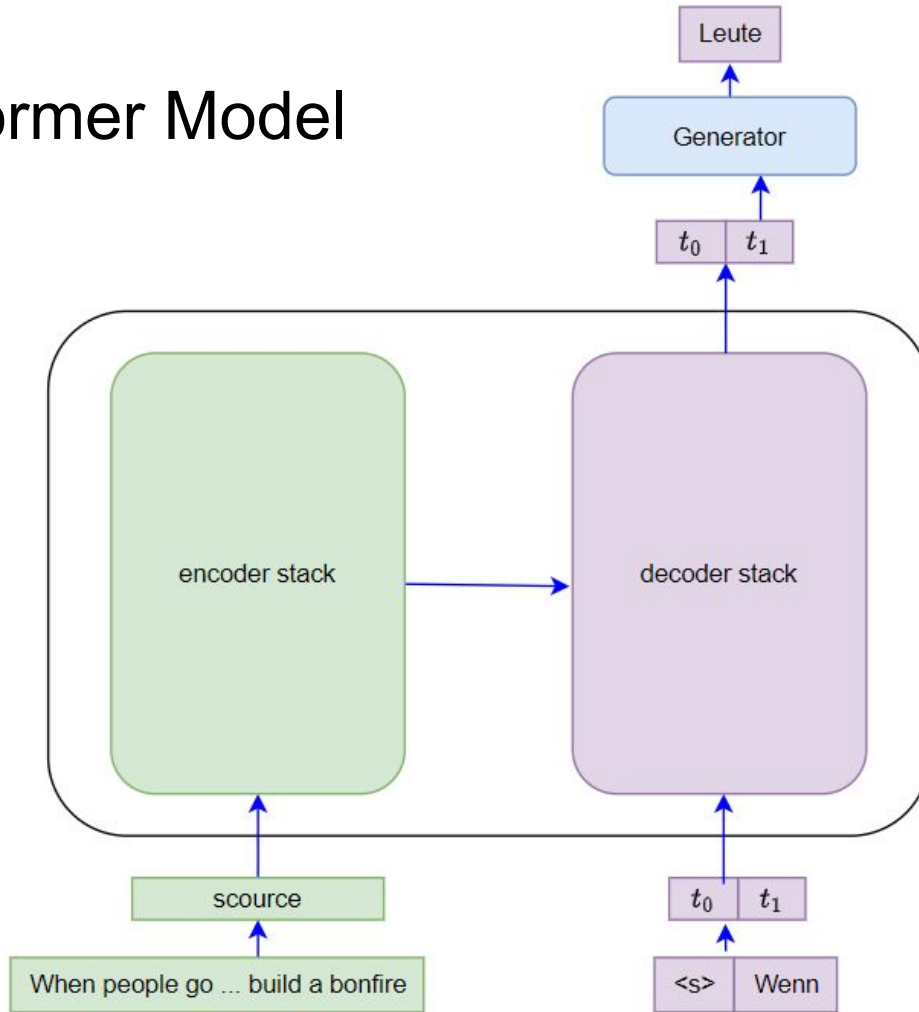# The Transformer Model

At Testing Time:



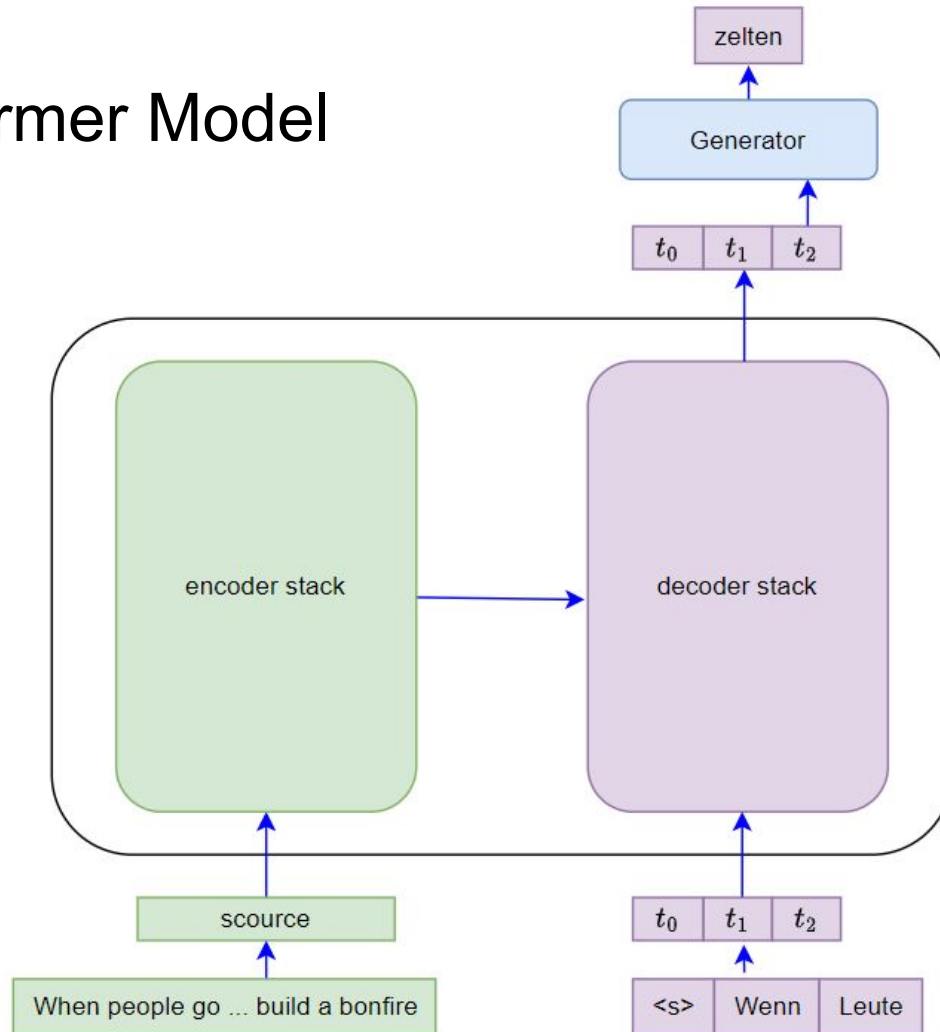word + pos embedding

# The Transformer Model
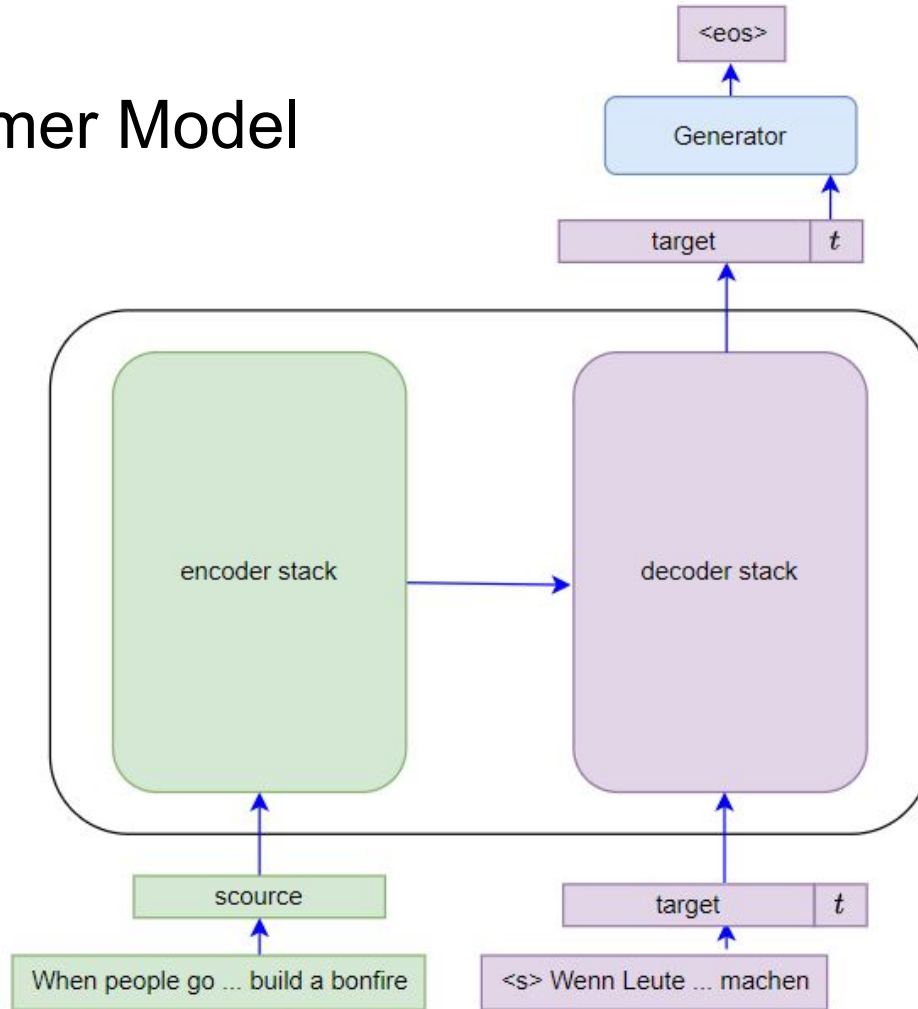
At Testing Time:

# The Transformer Model

At Testing Time:

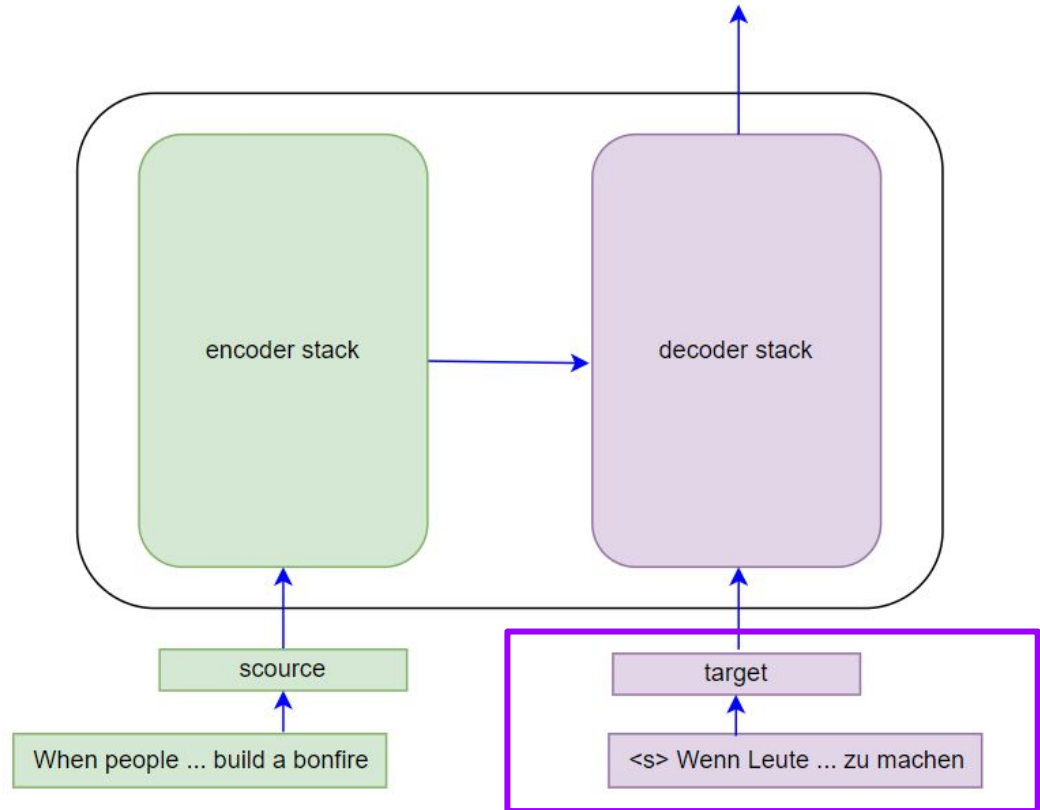# The Transformer Model

At Testing Time:
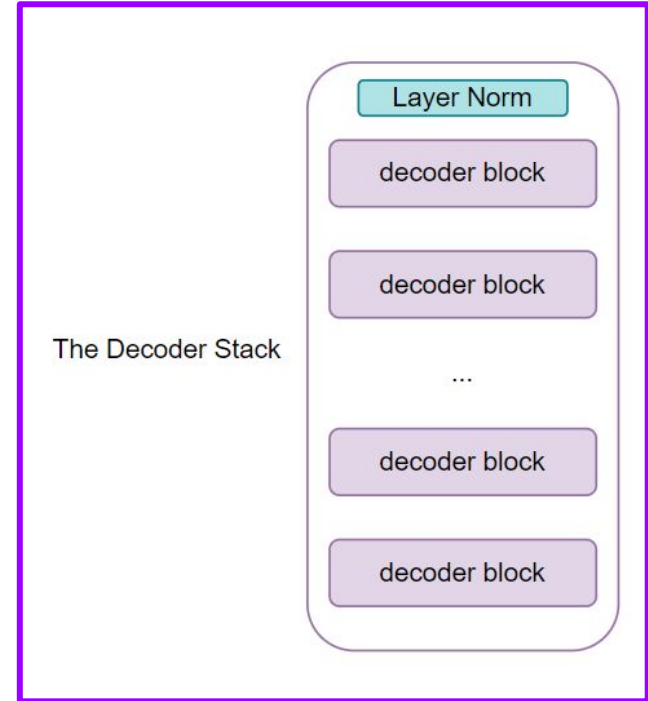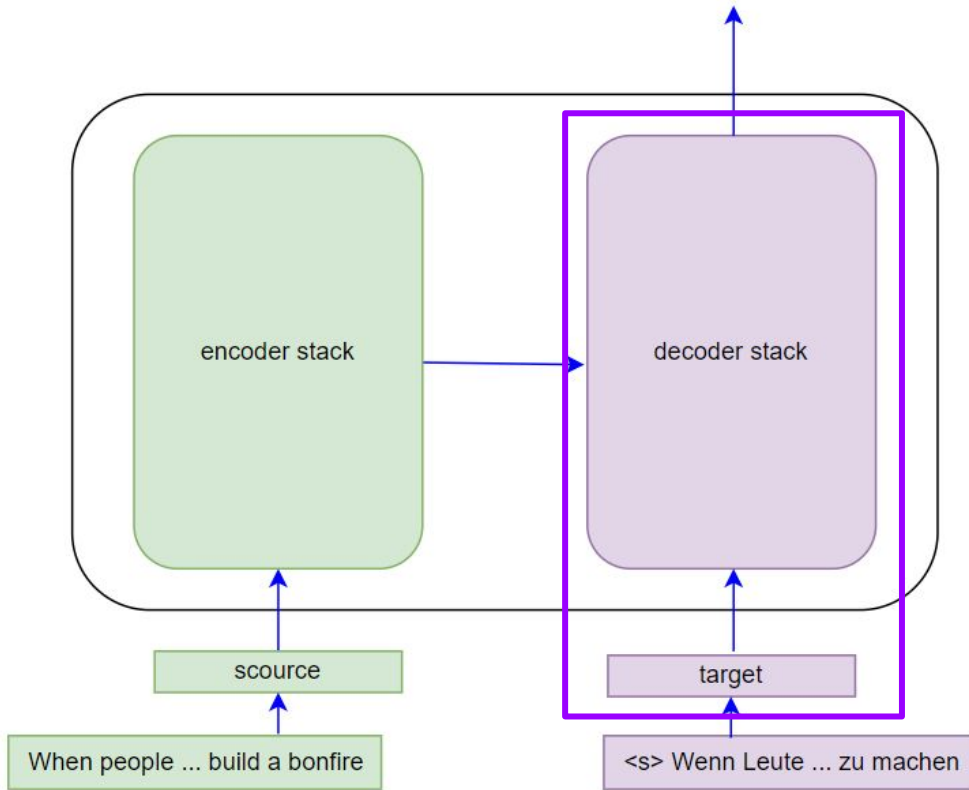
# Target Input Embedding

Add one special character at beginning
To shift every word one position behind
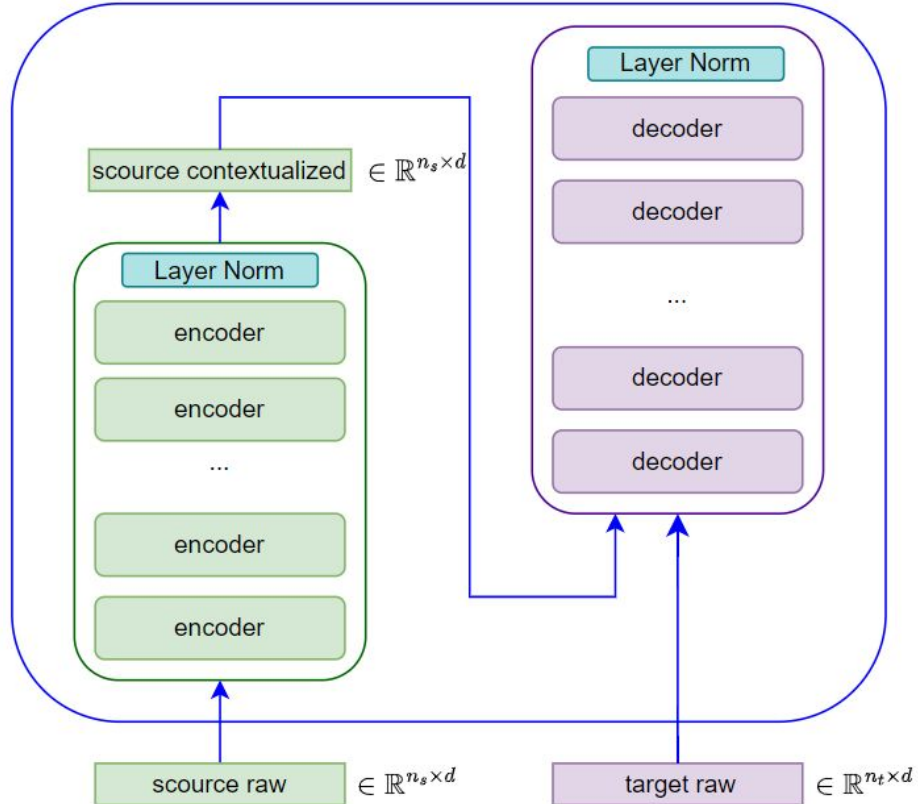
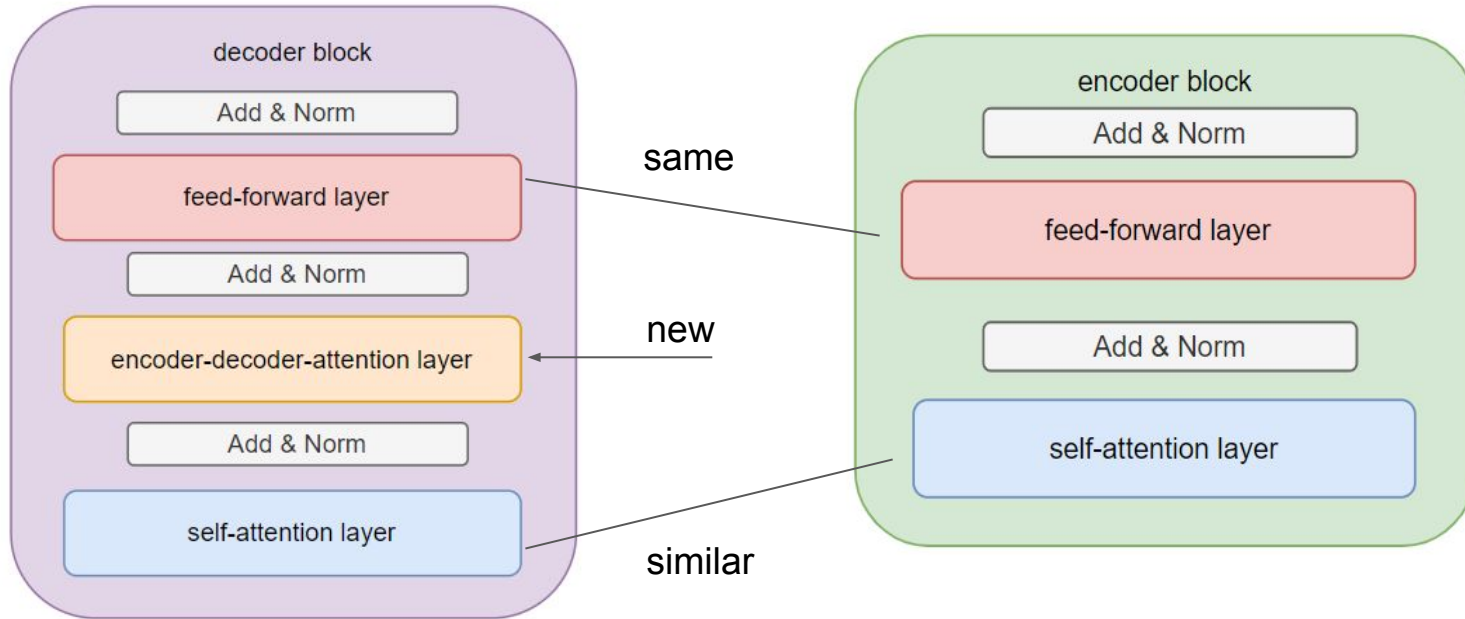1. word embedding

2. positional embedding

# The Decoder Stack
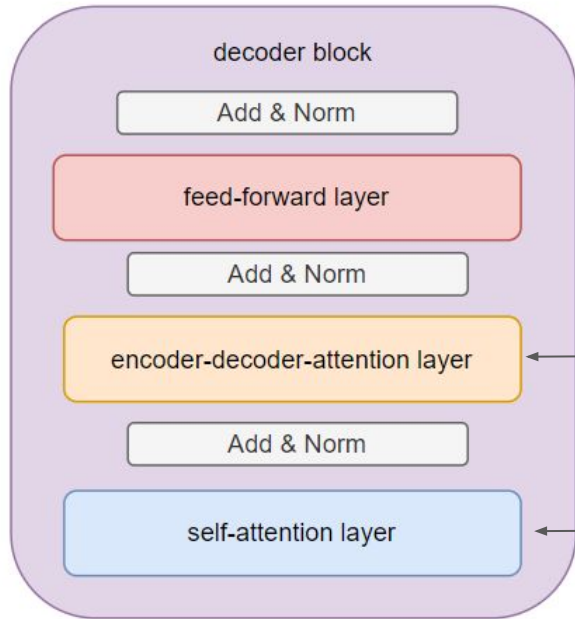
# The Decoder Stack

Training Time:

# The Decoder Block

# The Decoder Block



decoder block

Add & Norm

feed-forward layer

Add & Norm

encoder-decoder-attention layer — self-attention over target embedding and contextualized source embedding

Add & Norm

self-attention layer — self-attention over target embedding

target raw $\in \mathbb{R}^{n_t \times d}$

scource contextualized $\in \mathbb{R}^{n_s \times d}$

target raw $\in \mathbb{R}^{n_t \times d}$

# The Decoder Block



decoder block

| Add & Norm |
| feed-forward layer |
| Add & Norm |
| encoder-decoder-attention layer |
| Add & Norm |
| self-attention layer |

But uses a special attention mask

Input Embedding $\in \mathbb{R}^d$

$x_1$ $x_2$

sequence length = m

target embedded

Multi-headed Attention

$W_Q$ $W_Q$ ... $W_Q$
$W_K$ $W_K$ $W_K$
$W_V$ $W_V$ $W_V$

$W_O$

Output Embedding $\in \mathbb{R}^d$

$x_1$ $x_2$

sequence length = m

# The Decoder Block Self-attention Mask

| Good | morning | [PAD] |
|------|---------|-------|

Recall: use attention mask to avoid attention to [PAD] token

| 0 | 0 | -inf |
|---|---|------|
| 0 | 0 | -inf |
| 0 | 0 | -inf |

Not only [PAD] token, also can be used to mask any position in attention matrix we don't want

# The Decoder Block Self-attention Mask

Every word in decoder self-attention will only
attend to words before it and itself

| Ich | weiß | nicht |
|-----|------|-------|

The reason is, **during test time**,
when we have no ground-truth target embedding given,
we will predict each word **one-by-one**, not together

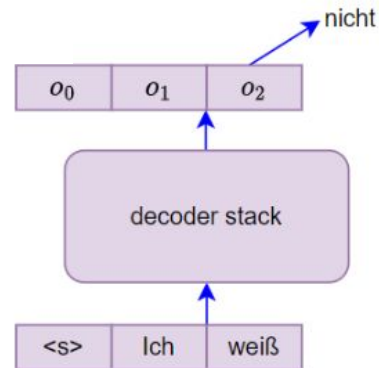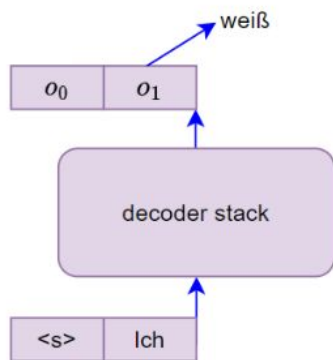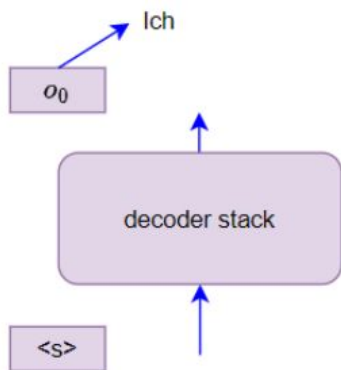| $x_1 \rightarrow x_1$ | $x_1 \rightarrow x_2$ | $x_1 \rightarrow x_3$ |
|-----|------|-------|
| $x_2 \rightarrow x_1$ | $x_2 \rightarrow x_2$ | $x_2 \rightarrow x_3$ |
| $x_3 \rightarrow x_1$ | $x_3 \rightarrow x_2$ | $x_3 \rightarrow x_3$ |

<s>              → Ich
<s> Ich          → weiß
<s> Ich weiß   → nicht

green means attends,
gray means masked

Each word can only see previous words in test time
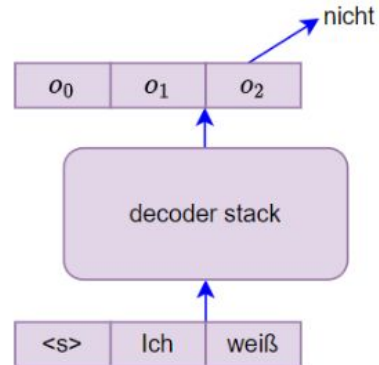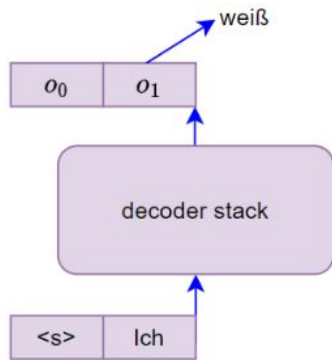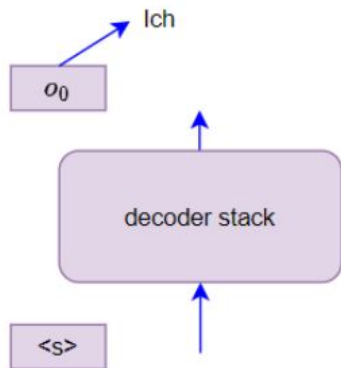So we mimic the same thing at training time

# The Decoder Block Self-attention Mask



Testing

Training
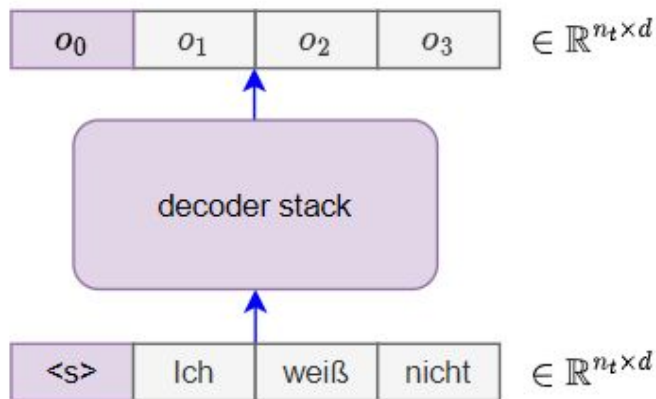
possible,
but want less
training time

# The Decoder Block Self-attention Mask

| $x_1 \rightarrow x_1$ | $x_1 \rightarrow x_2$ | $x_1 \rightarrow x_3$ |
|---|---|---|
| $x_2 \rightarrow x_1$ | $x_2 \rightarrow x_2$ | $x_2 \rightarrow x_3$ |
| $x_3 \rightarrow x_1$ | $x_3 \rightarrow x_2$ | $x_3 \rightarrow x_3$ |

Training Time:

| $o_0$ | $o_1$ | $o_2$ | $o_3$ | $\in \mathbb{R}^{n_t \times d}$ |

decoder stack

| <s> | Ich | weiß | nicht | $\in \mathbb{R}^{n_t \times d}$ |

Testing Time:

Ich

| $o_0$ |

decoder stack

| <s> |

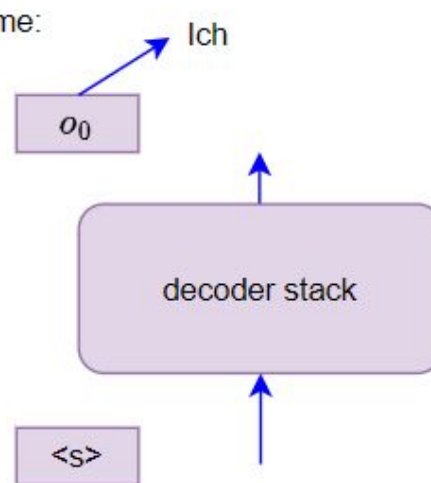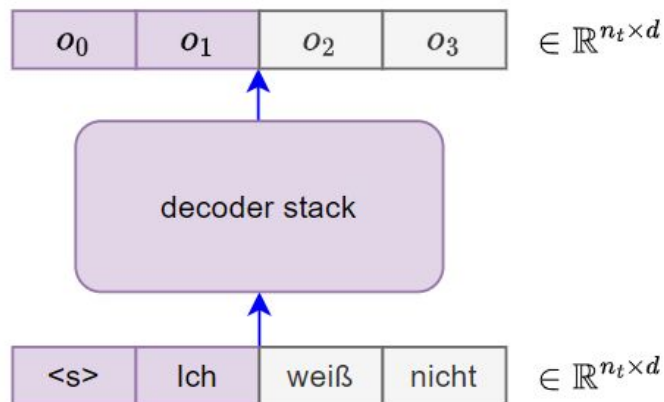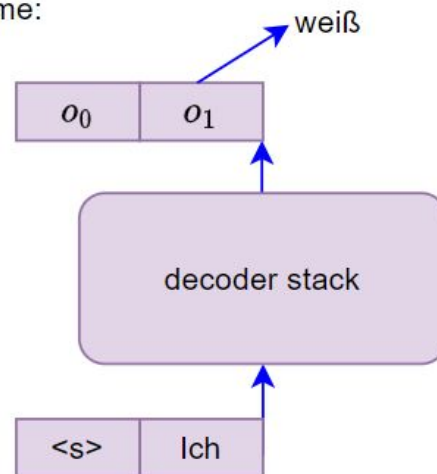$o_0$ only attends to start of sentence character

# The Decoder Block Self-attention Mask

| $x_1 \rightarrow x_1$ | $x_1 \rightarrow x_2$ | $x_1 \rightarrow x_3$ |
|---|---|---|
| $x_2 \rightarrow x_1$ | $x_2 \rightarrow x_2$ | $x_2 \rightarrow x_3$ |
| $x_3 \rightarrow x_1$ | $x_3 \rightarrow x_2$ | $x_3 \rightarrow x_3$ |

Training Time:

| $o_0$ | $o_1$ | $o_2$ | $o_3$ | $\in \mathbb{R}^{n_t \times d}$ |

decoder stack

| <s> | Ich | weiß | nicht | $\in \mathbb{R}^{n_t \times d}$ |

Testing Time:

weiß

| $o_0$ | $o_1$ |

decoder stack

| <s> | Ich |

$o_1$  attends to precious words

# The Decoder Block Self-attention Mask

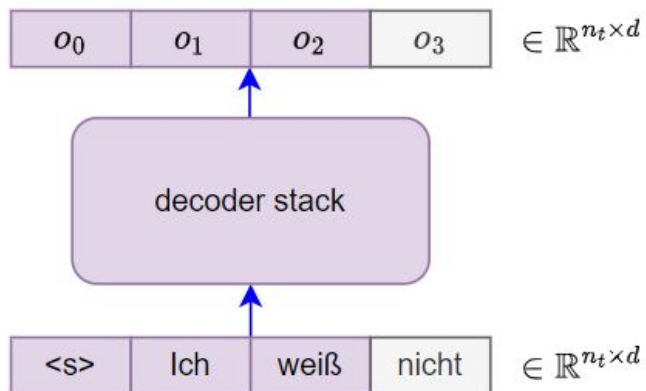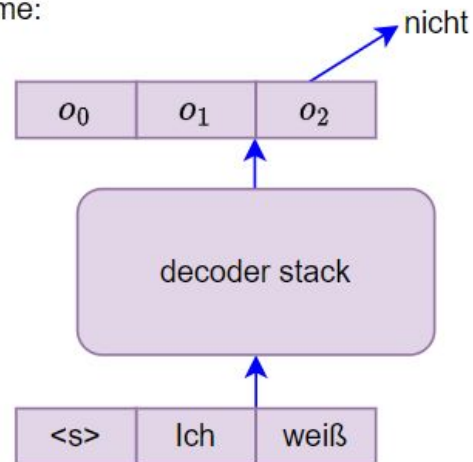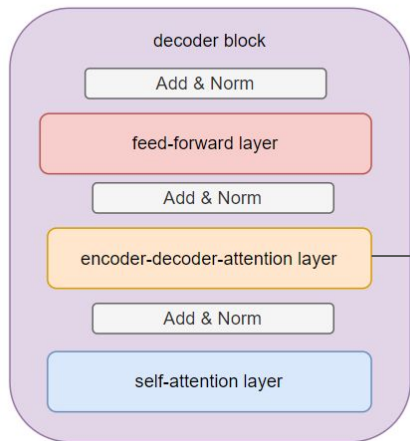| | | |
|---|---|---|
| $x_1 \rightarrow x_1$ | $x_1 \rightarrow x_2$ | $x_1 \rightarrow x_3$ |
| $x_2 \rightarrow x_1$ | $x_2 \rightarrow x_2$ | $x_2 \rightarrow x_3$ |
| $x_3 \rightarrow x_1$ | $x_3 \rightarrow x_2$ | $x_3 \rightarrow x_3$ |

Training Time:

| $o_0$ | $o_1$ | $o_2$ | $o_3$ | $\in \mathbb{R}^{n_t \times d}$ |

decoder stack

| <s> | Ich | weiß | nicht | $\in \mathbb{R}^{n_t \times d}$ |

Testing Time:

nicht

| $o_0$ | $o_1$ | $o_2$ |

decoder stack

| <s> | Ich | weiß |

$o_2$   attends to precious words

# The encoder-decoder attention



But compures QKV differently

# The encoder-decoder attention

Attention Head

Input Embedding $\in \mathbb{R}^d$

$x_1$   $x_2$

sequence length = m

$W_Q$

$W_K$

$W_V$

$q_1$
$q_2$
$Q \in \mathbb{R}^{m \times n}$

$k_1$
$k_2$
$K \in \mathbb{R}^{m \times n}$

$v_1$
$v_2$
$V \in \mathbb{R}^{m \times n}$

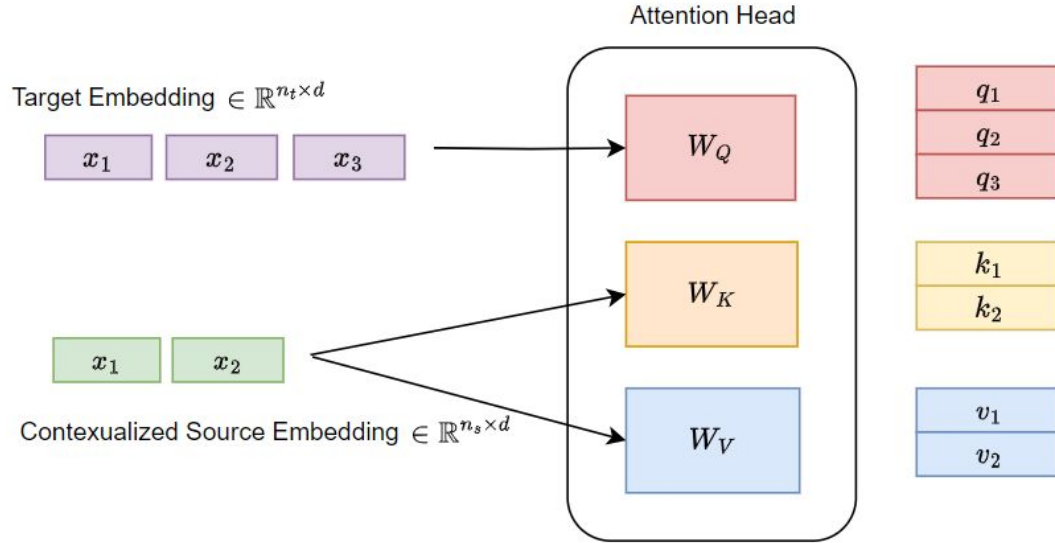In encoder-self-attention, QKV are computed from source embedding

In decoder-self-attention, QKV are computed from target embedding
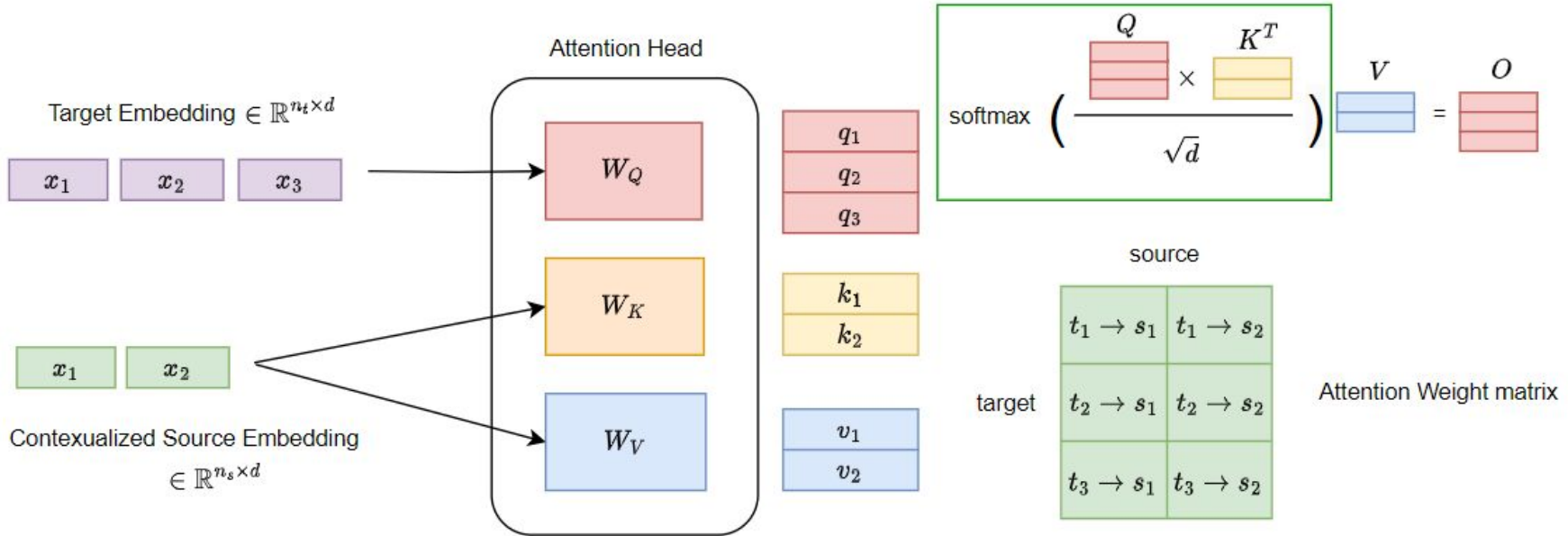
# The encoder-decoder attention



In encoder-decoder-attention,
Q is computed from target embedding
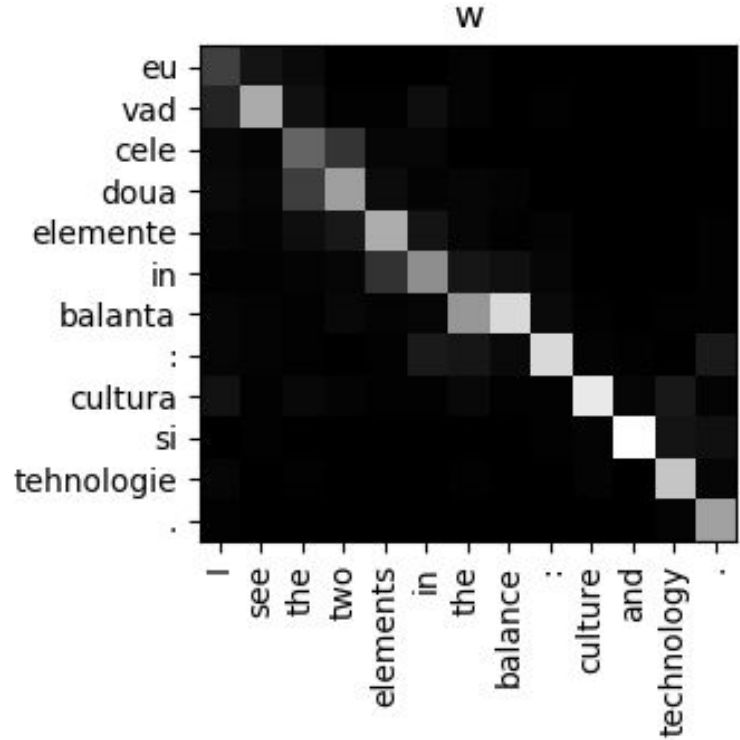KV are computed from contextualized source embedding
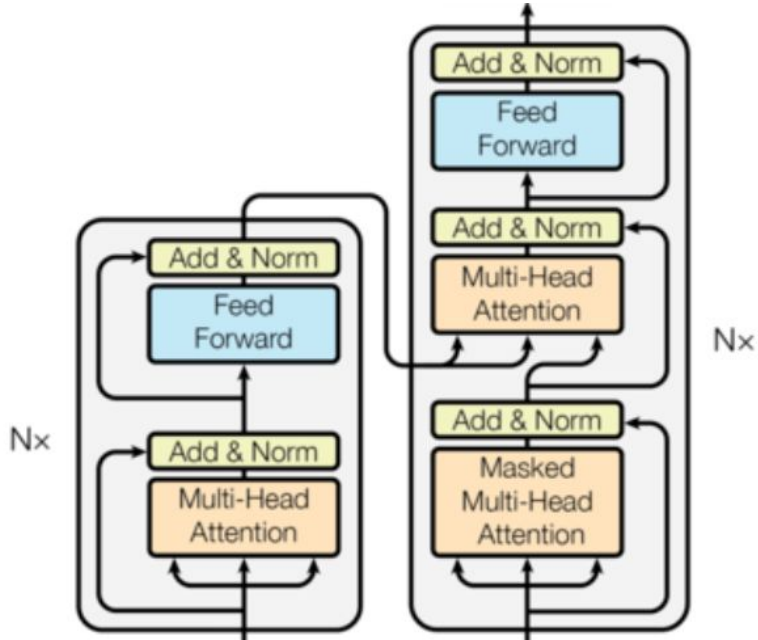
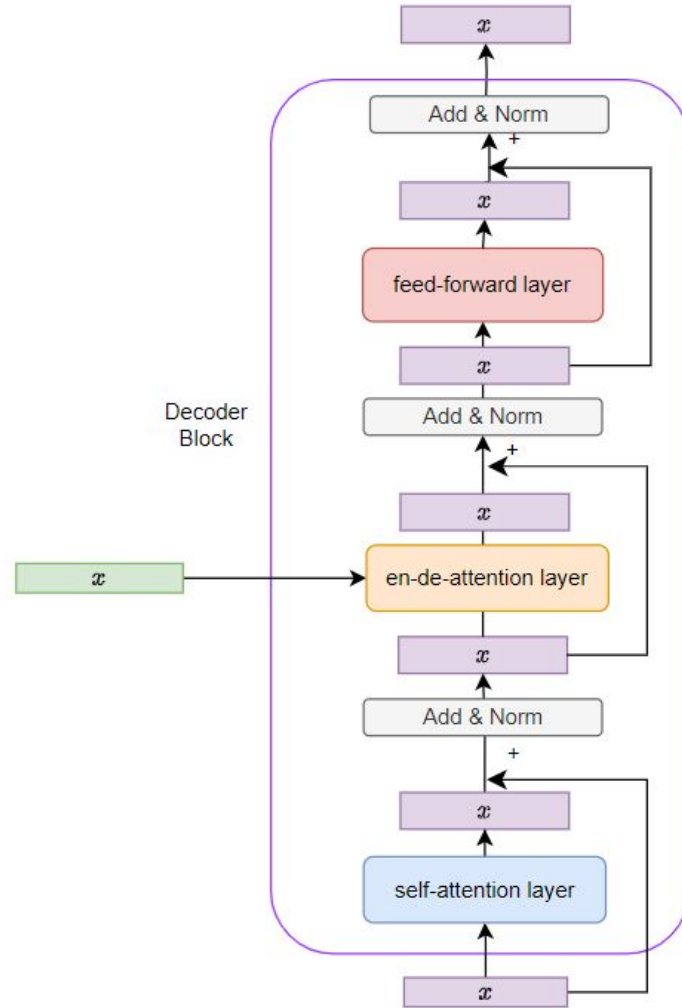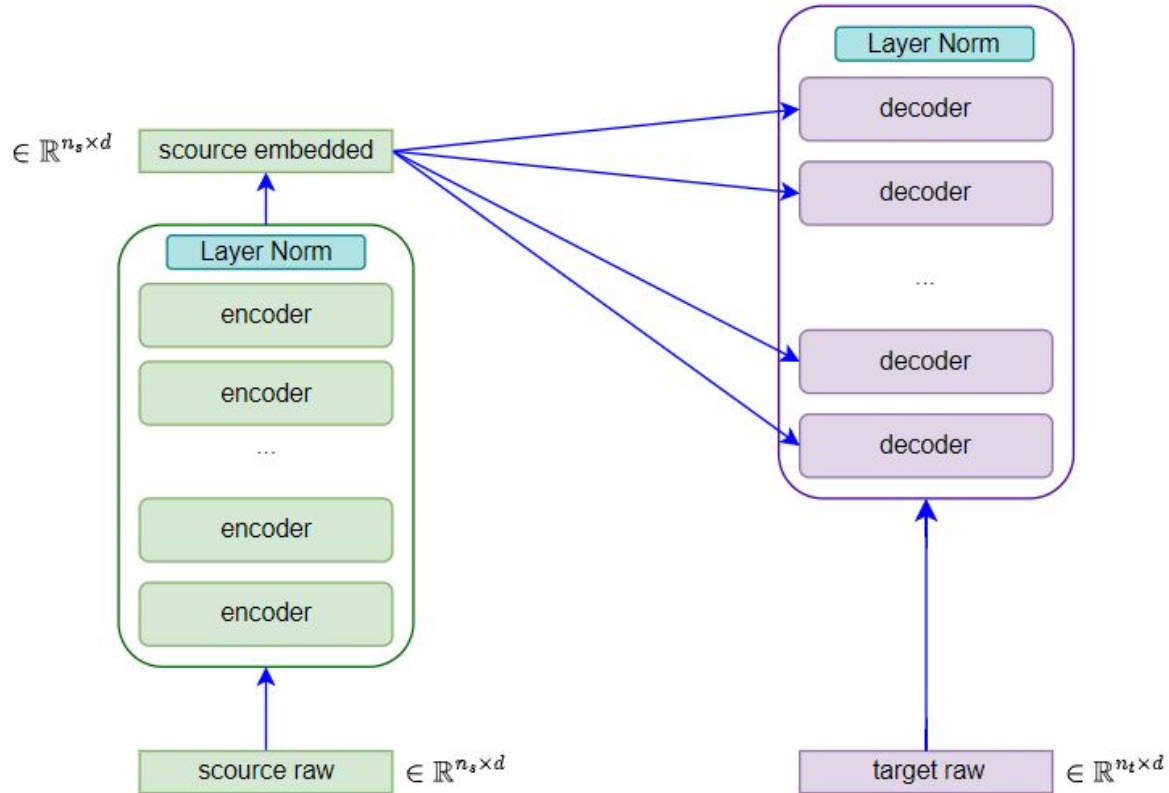# The encoder-decoder attention

# The encoder-decoder attention

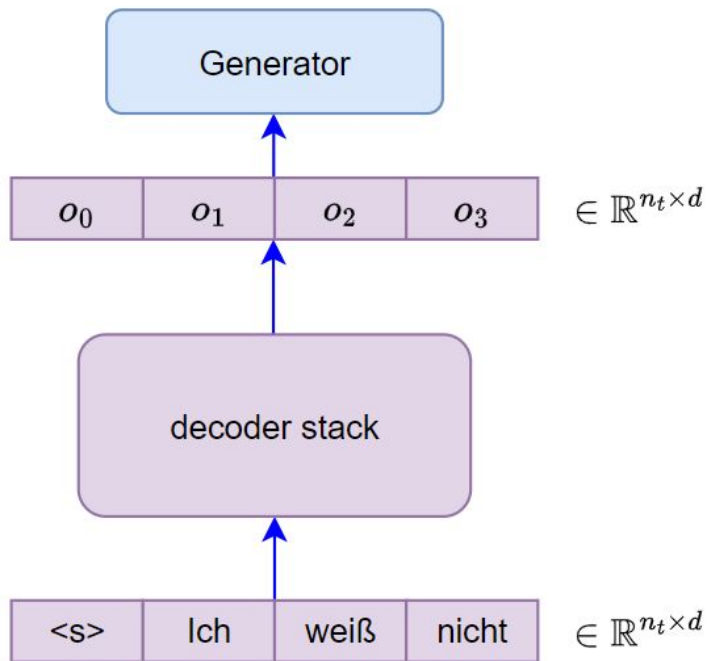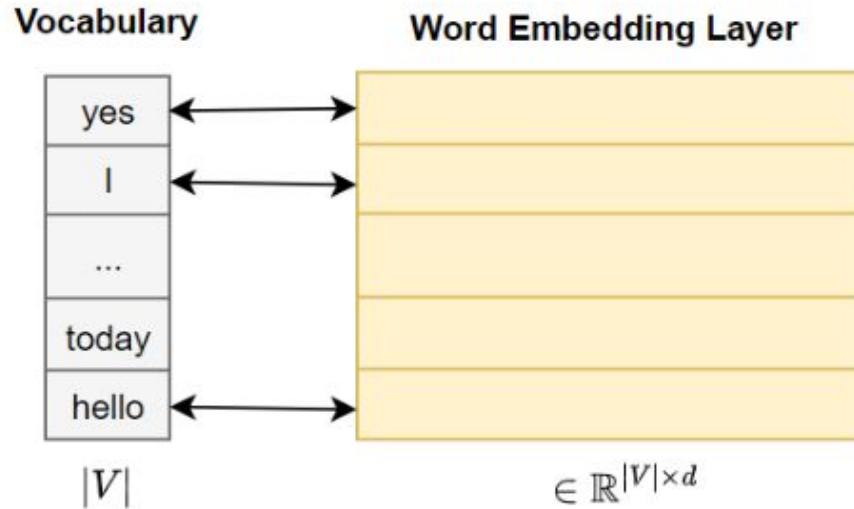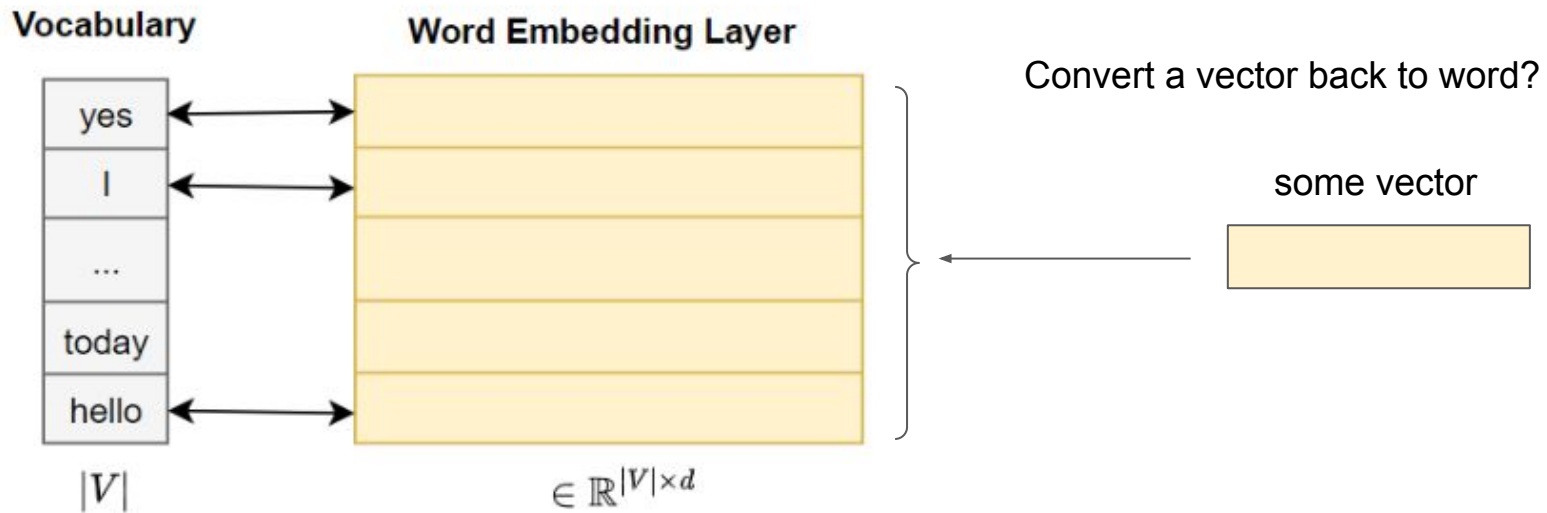# The Decoder Block



In original paper

# The Decoder Block

# The Generator

# The Generator



Word embedding: for a fixed word convert to a fixed embedding

# The Generator



Vocabulary

Word Embedding Layer

| yes |
| I |
| ... |
| today |
| hello |

$|V|$

$\in \mathbb{R}^{|V| \times d}$

Convert a vector back to word?

some vector

# The Generator

## Vocabulary

"yes"

| yes |
|-----|
| I |
| ... |
| today |
| hello |

$|V|$

## Word Embedding Layer

$\in \mathbb{R}^{|V| \times d}$
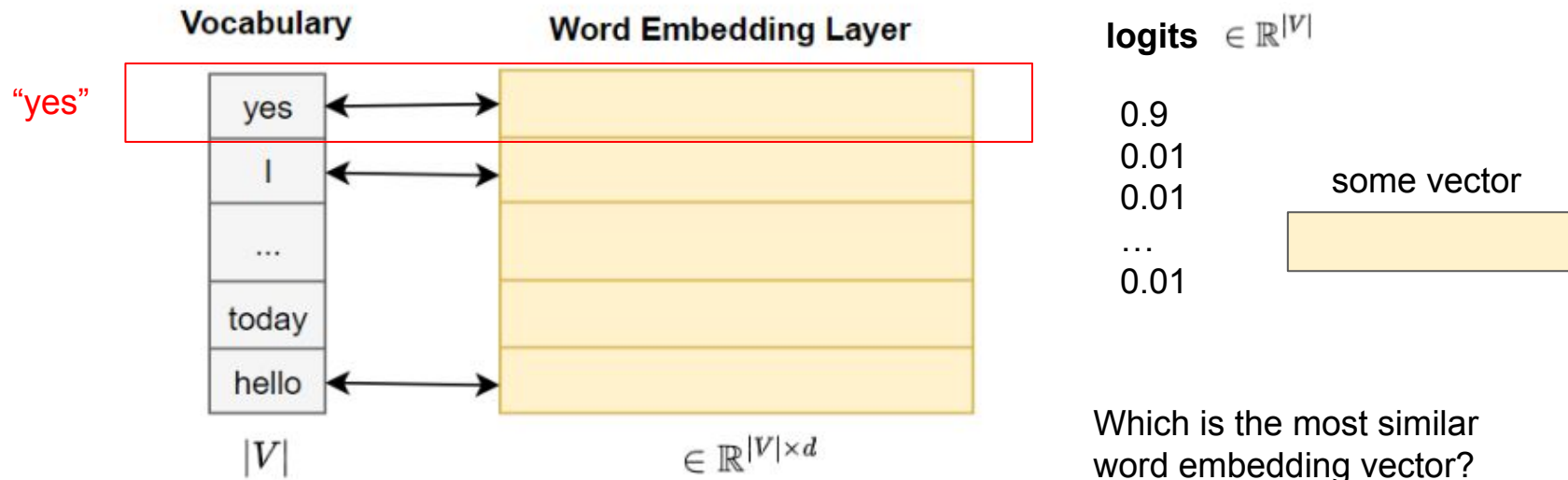
**logits** $\in \mathbb{R}^{|V|}$

0.9
0.01
0.01
…
0.01

some vector

Which is the most similar word embedding vector?

dot product with every word embedding vector
do softmax over result to get logits

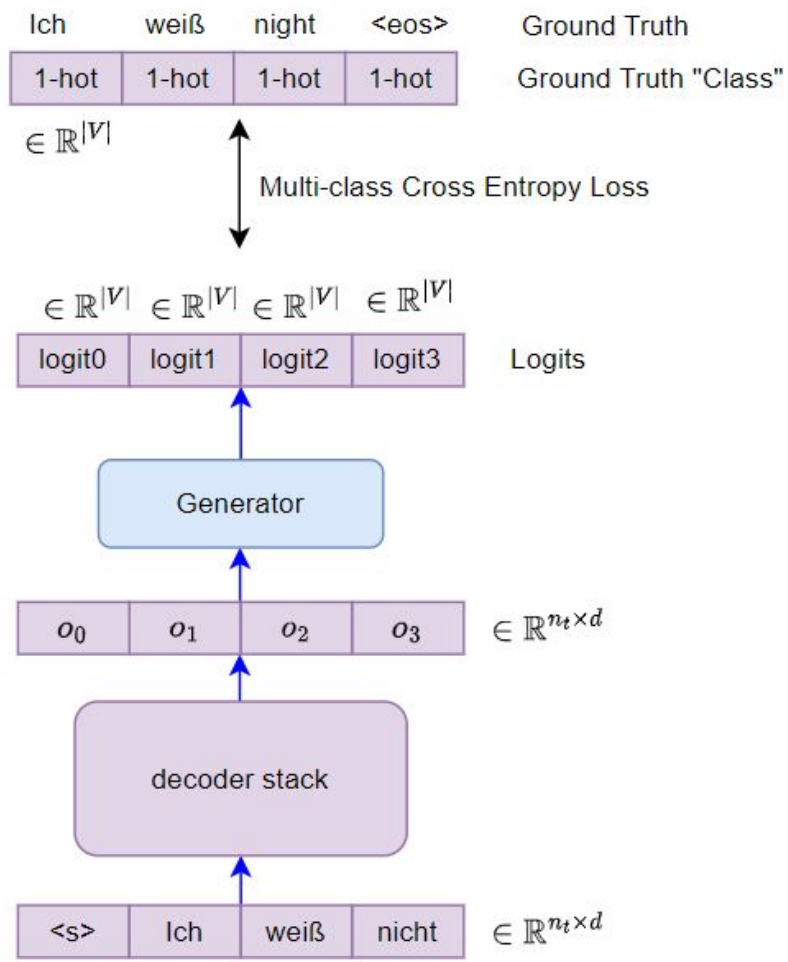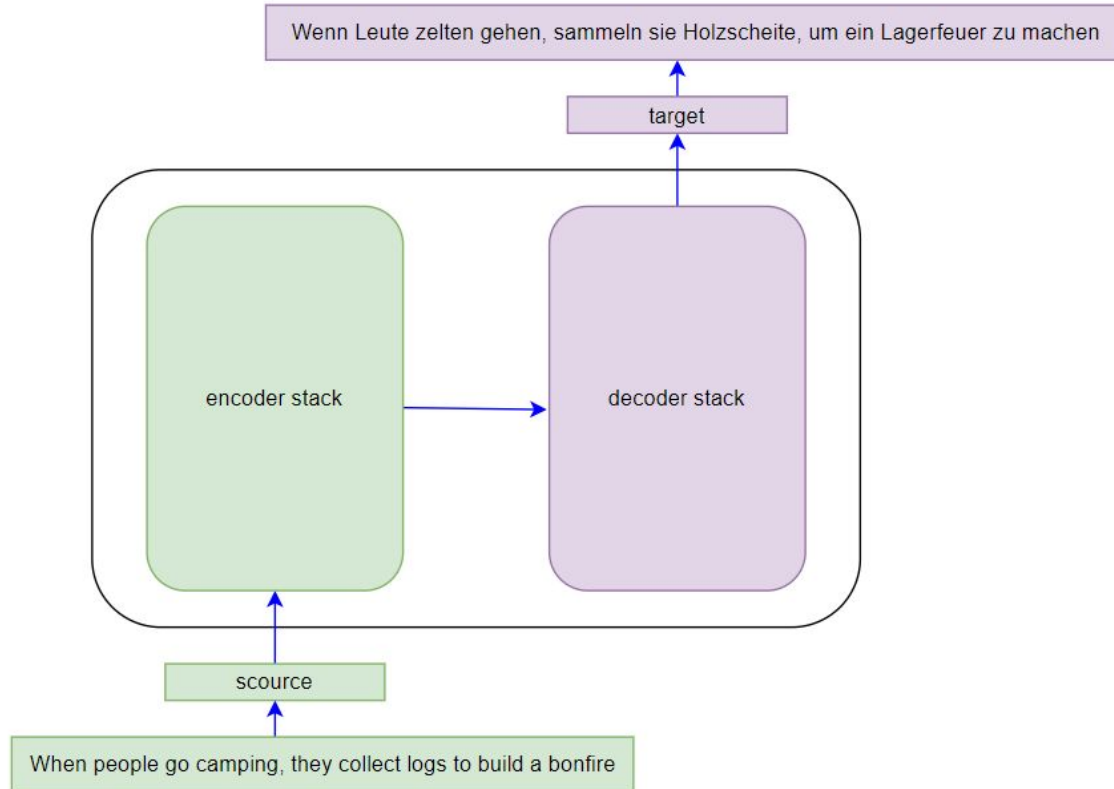# The Training Loss

| ground truth | logit |
|---|---|
| 1 | 0.9 |
| 0 | 0.01 |
| 0 | 0.01 |
| … | … |
| 0 | 0.01 |

# The Transformer Model

# References

- The annotated transformer http://nlp.seas.harvard.edu/annotated-transformer
- A Gentle Introduction to Positional Encoding in Transformer Models, Part 1 https://machinelearningmastery.com/a-gentle-introduction-to-positional-encoding-in-transformer-models-part-1/#:~:text=What%20Is%20Positional%20Encoding%3F,item's%20position%20in%20transformer%20models.
- Database analogy modified from Dive into Deep Learning CH11.1

# References

- Dzmitry Bahdanau, Kyunghyun Cho and Yoshua Bengio. 2014. Neural Machine Translation by Jointly Learning to Align and Translate. In International Conference on Learning Representations.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser and Illia Polosukhin. 2017. Attention is All you Need. In Advances in Neural Information Processing Systems. Curran Associates, Inc..
- Kaiming He, Xiangyu Zhang, Shaoqing Ren and Jian Sun. 2015. Deep Residual Learning for Image Recognition
- Jimmy Lei Ba, Jamie Ryan Kiros and Geoffrey E. Hinton. 2016. Layer Normalization