# CS-E4890: Deep Learning
# Q&A Session Assignment 5 - Transformers

Lukas Prediger
31.03.2023

## Your tasks

You will have to implement the same translation task as in the previous assignment, but this time using a transformer architecture. This is broken down into the implementation steps

1. the `collate` function for converting input data batches to PyTorch tensors,
2. `EncoderBlock` and `Encoder` modules,
3. `DecoderBlock` and `Decoder` modules,
4. the training loop,
5. the `translate` function to perform translation of any source sentence without knowing the target.

1. Recap of the transformer architecture
2. Batching, padding and `collate`
3. Common implementation pitfalls & tips
4. Your questions

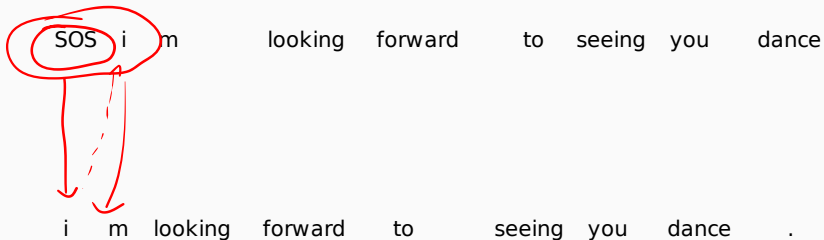# Transformer Recap

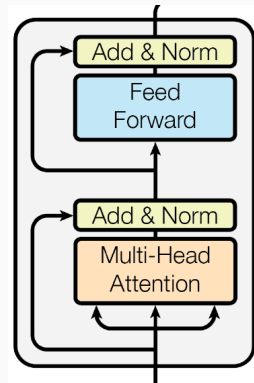je suis impatiente de te voir danser .

*ENCODER*

*EMBEDDED MEANING*

*DECODER*

i m looking forward to seeing you dance .

i  m  looking  forward  to  seeing  you ...

ENCODER

DECODER

i  m  looking  forward  to  seeing  you        dance   .

SOS   i   m        looking   forward        to   seeing   you      dance

i    m   looking    forward      to         seeing   you    dance      .

je  suis  impatiente  de  te  voir  danser  .

}WORD EMBEDDING

}POSITIONAL ENCODING

i think that i could …

je    suis    (impatiente)    de    te    voir    danser    .



$\vec{s}_1$   $\vec{s}_2$        $\vec{s}_3$

$$\vec{s}_1^{\,'} = \sum_{j=1}^{8} w_{1j} \, \vec{s}_j \quad \longleftarrow \text{VALUES}$$

$$\sum_{j=1}^{8} w_{1j} \overset{!}{=} 1$$

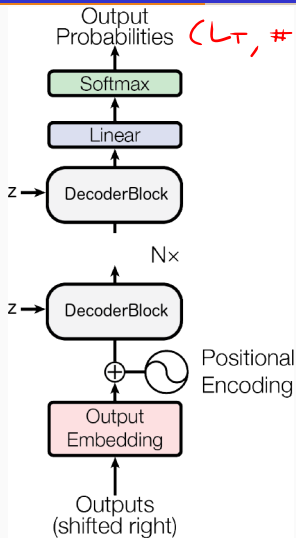$$w_{1j} \approx \text{SIMILARITY}(\vec{s}_1, \vec{s}_j) = \vec{s}_1^{\,T} \vec{s}_j =: a_{1j}$$

$$\underset{\text{QUERY}}{\uparrow} \quad \underset{\text{KEY}}{\uparrow}$$

$$w_{1j} = \frac{e^{a_{1j}}}{\sum e^{a_{1i}}}$$

KEYS

QUERIES

| | | $w_{1,1}$ | $w_{1,2}$ | $w_{1,3}$ | $w_{1,4}$ | $w_{1,5}$ | $w_{1,6}$ | $w_{1,7}$ | $w_{1,8}$ | $w_{1,9}$ | $w_{1,10}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|

$\vec{s}_1$    $w_{1,1}$   $w_{1,2}$   $w_{1,3}$    $w_{1,4}$   $w_{1,5}$   $w_{1,6}$    $w_{1,7}$   $w_{1,8}$   $w_{1,9}$   $w_{1,10}$

$\vec{s}_2$    $w_{2,1}$   $w_{2,2}$   $w_{2,3}$    —    —    —

$\vec{s}_3$    $w_{3,1}$   $w_{3,2}$   —   —   —

$V \rightarrow \vec{s}_1' \quad . \quad \vec{s}_8'$

SOS i m looking forward to seeing you dance

i m looking forward to seeing you dance .

$A \in \mathbb{R}^{L_T \times L_T}$

$A \oplus M$

KEYS

QUERIES

$\vec{t}_1 \quad \vec{t}_2 \quad \vec{t}_3 \quad \vec{t}_4 \quad \vec{t}_5 \quad \vec{t}_6 \quad \vec{t}_7 \quad \vec{t}_8 \quad \vec{t}_9 \quad \vec{t}_{10} \quad \vec{t}_{11}$

$\vec{t}_1 \quad A_{1,1} \quad A_{1,2} \quad A_{1,3} \quad A_{1,4} \quad A_{1,5} \quad A_{1,6} \quad A_{1,7} \quad A_{1,8} \quad A_{1,9} \quad A_{1,10} \quad A_{1,11}$

$\vec{t}_2 \quad A_{2,1} \quad A_{2,2} \quad A_{2,3} \quad A_{2,4} \quad A_{2,5} \quad A_{2,6} \quad A_{2,7} \quad A_{2,8} \quad A_{2,9} \quad A_{2,10} \quad A_{2,11}$

$\vec{t}_3 \quad A_{3,1} \quad A_{3,2} \quad A_{3,3} \quad A_{3,4} \quad A_{3,5} \quad A_{3,6} \quad A_{3,7} \quad A_{3,8} \quad A_{3,9} \quad A_{3,10} \quad A_{3,11}$

$A \in \mathbb{P}^{L_T \times L_T}$

$\vdots$

12

$\vec{S}_1$ --- $\vec{S}_8$ ← KEYS, VALUES

$\vec{t}_1$ --- $\vec{t}_T$ ← QUERYS

**Batching, padding,** `collate`

Transformers are (also) about efficiency: We want to parallelise as much as possible.

Since sequences have different lengths, we need to pad them in order to create a tensor!

je    suis    impatiente    de    te    voir    danser    .    *PAD*

il    est    toujours    en train    de    se    plaindre    .

# Padding and Attention

je    suis    impatiente    de    te    voir    danser    .    PAD

$(B, L)$



1   2   3   4   5   6   7   8   9

$P_1$    ———    ———    ———    ———    ———    ———    $P_9$

$w_{19} \overset{!}{=} 0$    $w_{ij} \overset{!}{=} 0$

$$P_j = \begin{cases} 0 & \text{if padding} \\ 1 & \text{otherwise} \end{cases}$$

$$w_{ij} = \frac{e^{\alpha_{ij}} \cdot P_j}{\sum_{k=1} \left( e^{a_{ik}} \cdot P_k \right)}$$

15

# Common implementation pitfalls

We provide tr.PositionalEncoding to you for the positional encoding step.

tr.PositionalEncoding adds the positional encoding to the input it is provided *internally*

you do not need to perform the addition step yourself!

## Including padded positions in the loss...

... incentivises the model to become really good at predicting padding.

At the expense of becoming less good at predicting interesting things.

*Solution:* PyTorch loss functions (`nll_loss`, `cross_entropy`, ...) have a `ignore_index` parameter.

`nn.Embedding` similarly has a `padding_idx` parameter.

... causes the decoder to learn to predict the *current* instead of the *next* word.

## Some other tips

- `nn.Sequential` simplifies MLP implementation: https://pytorch.org/docs/1.10/generated/torch.nn.Sequential.html#torch.nn.Sequential.
- `nn.ModuleList` helps with blocks in Encoder and Decoder: https://pytorch.org/docs/1.10/generated/torch.nn.Sequential.html#torch.nn.ModuleList.
- Consider creating additional `nn.Module` classes for reusable sub-blocks.

# Questions?

# Room for questions