

Advanced Deep Learning

DATA.ML.230

Transformers

From: Simon J. D. Prince, Chapter 12 – Understanding Deep Learning, MIT Press (19 May 2025)

Motivation

Transformers were initially proposed for text processing tasks.

Example: Restaurant review

- “The restaurant refused to serve me a ham sandwich because it only cooks vegetarian food. In the end, they just gave me two slices of bread.”

Goal: design a network to process this text into a representation for downstream tasks:

- Classification to positive or negative
- Question answering: “Does the restaurant serve steak?”

Motivation

Transformers were initially proposed for text processing tasks.

Example: Restaurant review

- “The restaurant refused to serve me a ham sandwich because it only cooks vegetarian food. In the end, they just gave me two slices of bread.”

Parameter sharing → Varying size of input (or even sentences)

Goal: design a network to process this text into a representation for downstream tasks:

- Classification to positive or negative
- Question answering: “Does the restaurant serve steak?|

Motivation

Transformers were initially proposed for text processing tasks.

Example: Restaurant review

- “The **restaurant** refused to serve me a **ham sandwich** because **it** only cooks vegetarian food. In the end, they just gave me two slices of bread.”

Language can be ambiguous

Goal: design a network to process this text into a representation for downstream tasks:

- Classification to positive or negative
- Question answering: “Does the restaurant serve steak?”

Motivation

Transformers were initially proposed for text processing tasks.

Example: Restaurant review

- “The **restaurant** refused to serve me a **ham sandwich** because **it** only cooks vegetarian food. In the end, they just gave me two slices of bread.”

Attention



Language can be ambiguous

Goal: design a network to process this text into a representation for downstream tasks:

- Classification to positive or negative
- Question answering: “Does the restaurant serve steak?”

Dot-product self-attention

A self-attention block **sa[•]** takes N inputs x_1, \dots, x_N and returns N outputs:

Dot-product self-attention

A self-attention block **sa[•]** takes N inputs $\mathbf{x}_1, \dots, \mathbf{x}_N$ and returns N outputs:

- A set of *values* are computed for each input:

$$\mathbf{v}_m = \boldsymbol{\beta}_v + \boldsymbol{\Omega}_v \mathbf{x}_m$$

Dot-product self-attention

A self-attention block **sa[•]** takes N inputs $\mathbf{x}_1, \dots, \mathbf{x}_N$ and returns N outputs:

- A set of *values* are computed for each input:

$$\mathbf{v}_m = \beta_v + \Omega_v \mathbf{x}_m$$

Parameter sharing:
The same parameters
are applied to all inputs

This computation
scales linearly with the
sequence length N

Dot-product self-attention

A self-attention block $\text{sa}[\bullet]$ takes N inputs $\mathbf{x}_1, \dots, \mathbf{x}_N$ and returns N outputs:

- A set of *values* are computed for each input:

$$\mathbf{v}_m = \beta_v + \Omega_v \mathbf{x}_m$$

- Then, the n^{th} output $\text{sa}_n[\mathbf{x}_1, \dots, \mathbf{x}_N]$ is a weighted sum of all the values

$\mathbf{v}_1, \dots, \mathbf{v}_N$:

$$\text{sa}_n[\mathbf{x}_1, \dots, \mathbf{x}_N] = \sum_{m=1}^N a[\mathbf{x}_m, \mathbf{x}_n] \mathbf{v}_m$$

Dot-product self-attention

A self-attention block $\text{sa}[\bullet]$ takes N inputs $\mathbf{x}_1, \dots, \mathbf{x}_N$ and returns N outputs:

- A set of *values* are computed for each input:

$$\mathbf{v}_m = \beta_v + \Omega_v \mathbf{x}_m$$

- Then, the n^{th} output $\text{sa}_n[\mathbf{x}_1, \dots, \mathbf{x}_N]$ is a weighted sum of all the values

$\mathbf{v}_1, \dots, \mathbf{v}_N$:

$$\text{sa}_n[\mathbf{x}_1, \dots, \mathbf{x}_N] = \sum_{m=1}^N a[\mathbf{x}_m, \mathbf{x}_n] \mathbf{v}_m$$

Attention weights:

- Combine the values from different inputs
- N^2 in total (quadratic dependence on the sequence length N)
- Independent of the dimension of \mathbf{x}_n

Dot-product self-attention

A self-attention block $\text{sa}[\bullet]$ takes N inputs $\mathbf{x}_1, \dots, \mathbf{x}_N$ and returns N outputs:

- A set of *values* are computed for each input:

$$\mathbf{v}_m = \beta_v + \Omega_v \mathbf{x}_m$$

- Then, the n^{th} output $\text{sa}_n[\mathbf{x}_1, \dots, \mathbf{x}_N]$ is a weighted sum of all the values $\mathbf{v}_1, \dots, \mathbf{v}_N$:

$$\text{sa}_n[\mathbf{x}_1, \dots, \mathbf{x}_N] = \sum_{m=1}^N a[\mathbf{x}_m, \mathbf{x}_n] \mathbf{v}_m$$

- The scalar weight $a[\mathbf{x}_m, \mathbf{x}_n]$ is the attention that the n^{th} output pays to input \mathbf{x}_m
- The N weights $a[\bullet, \mathbf{x}_n]$ are non-negative and sum to one

Dot-product self-attention

A self-attention block $\text{sa}[\bullet]$ takes N inputs $\mathbf{x}_1, \dots, \mathbf{x}_N$ and returns N outputs:

- A set of *values* are computed for each input:

$$\mathbf{v}_m = \beta_v + \Omega_v \mathbf{x}_m$$

- Then, the n^{th} output $\text{sa}_n[\mathbf{x}_1, \dots, \mathbf{x}_N]$ is a weighted sum of all the values $\mathbf{v}_1, \dots, \mathbf{v}_N$:

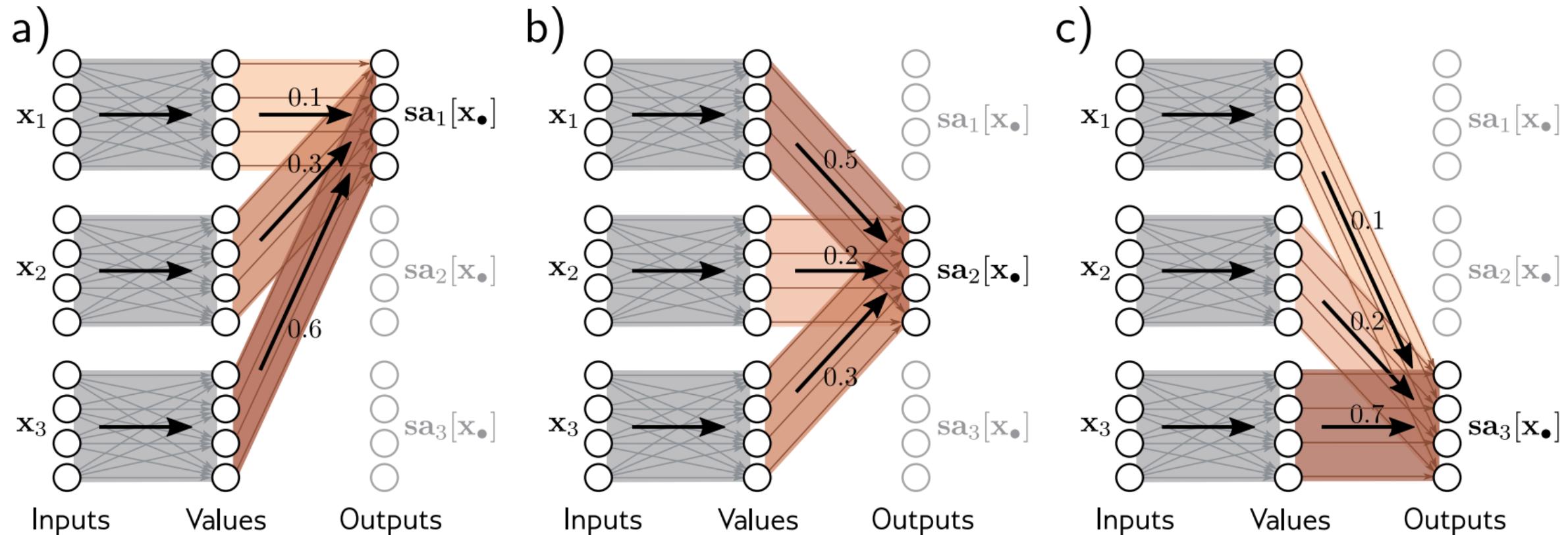
$$\text{sa}_n[\mathbf{x}_1, \dots, \mathbf{x}_N] = \sum_{m=1}^N a[\mathbf{x}_m, \mathbf{x}_n] \mathbf{v}_m$$

- The scalar weight $a[\mathbf{x}_m, \mathbf{x}_n]$ is the attention that the n^{th} output pays to input \mathbf{x}_m
- The N weights $a[\bullet, \mathbf{x}_n]$ are non-negative and sum to one

← Leads to non-linear mapping

Dot-product self-attention

A self-attention as routing



Attention weights

The overall self-attention computation is nonlinear:

- Apply two more linear transformations:

$$\text{Queries} \quad \mathbf{q}_n = \boldsymbol{\beta}_q + \boldsymbol{\Omega}_q \mathbf{x}_n$$

$$\text{Keys} \quad \mathbf{k}_m = \boldsymbol{\beta}_k + \boldsymbol{\Omega}_k \mathbf{x}_m$$

Attention weights

The overall self-attention computation is nonlinear:

- Apply two more linear transformations:

$$\text{Queries} \quad \mathbf{q}_n = \boldsymbol{\beta}_q + \boldsymbol{\Omega}_q \mathbf{x}_n$$

$$\text{Keys} \quad \mathbf{k}_m = \boldsymbol{\beta}_k + \boldsymbol{\Omega}_k \mathbf{x}_m$$

- Compute the dot products between the queries and keys and pass the results through a softmax function:

$$a[\mathbf{x}_m, \mathbf{x}_n] = \text{softmax}_m [\mathbf{k}_m^T \mathbf{q}_n] = \frac{\exp [\mathbf{k}_m^T \mathbf{q}_n]}{\sum_{m'=1}^N \exp [\mathbf{k}_{m'}^T \mathbf{q}_n]}$$

Attention weights

The overall self-attention computation is nonlinear:

- Apply two more linear transformations:

$$\text{Queries} \quad \mathbf{q}_n = \boldsymbol{\beta}_q + \boldsymbol{\Omega}_q \mathbf{x}_n$$

$$\text{Keys} \quad \mathbf{k}_m = \boldsymbol{\beta}_k + \boldsymbol{\Omega}_k \mathbf{x}_m$$

- Compute the dot products between the queries and keys and pass the results through a softmax function:

$$a[\mathbf{x}_m, \mathbf{x}_n] = \text{softmax}_m [\mathbf{k}_m^T \mathbf{q}_n] = \frac{\exp [\mathbf{k}_m^T \mathbf{q}_n]}{\sum_{m'=1}^N \exp [\mathbf{k}_{m'}^T \mathbf{q}_n]}$$

The weights $a[\mathbf{x}_\bullet, \mathbf{x}_n]$ sum to 1

Attention weights

The overall self-attention computation is nonlinear:

- Apply two more linear transformations:

$$\text{Queries} \quad \mathbf{q}_n = \boldsymbol{\beta}_q + \boldsymbol{\Omega}_q \mathbf{x}_n$$

$$\text{Keys} \quad \mathbf{k}_m = \boldsymbol{\beta}_k + \boldsymbol{\Omega}_k \mathbf{x}_m$$

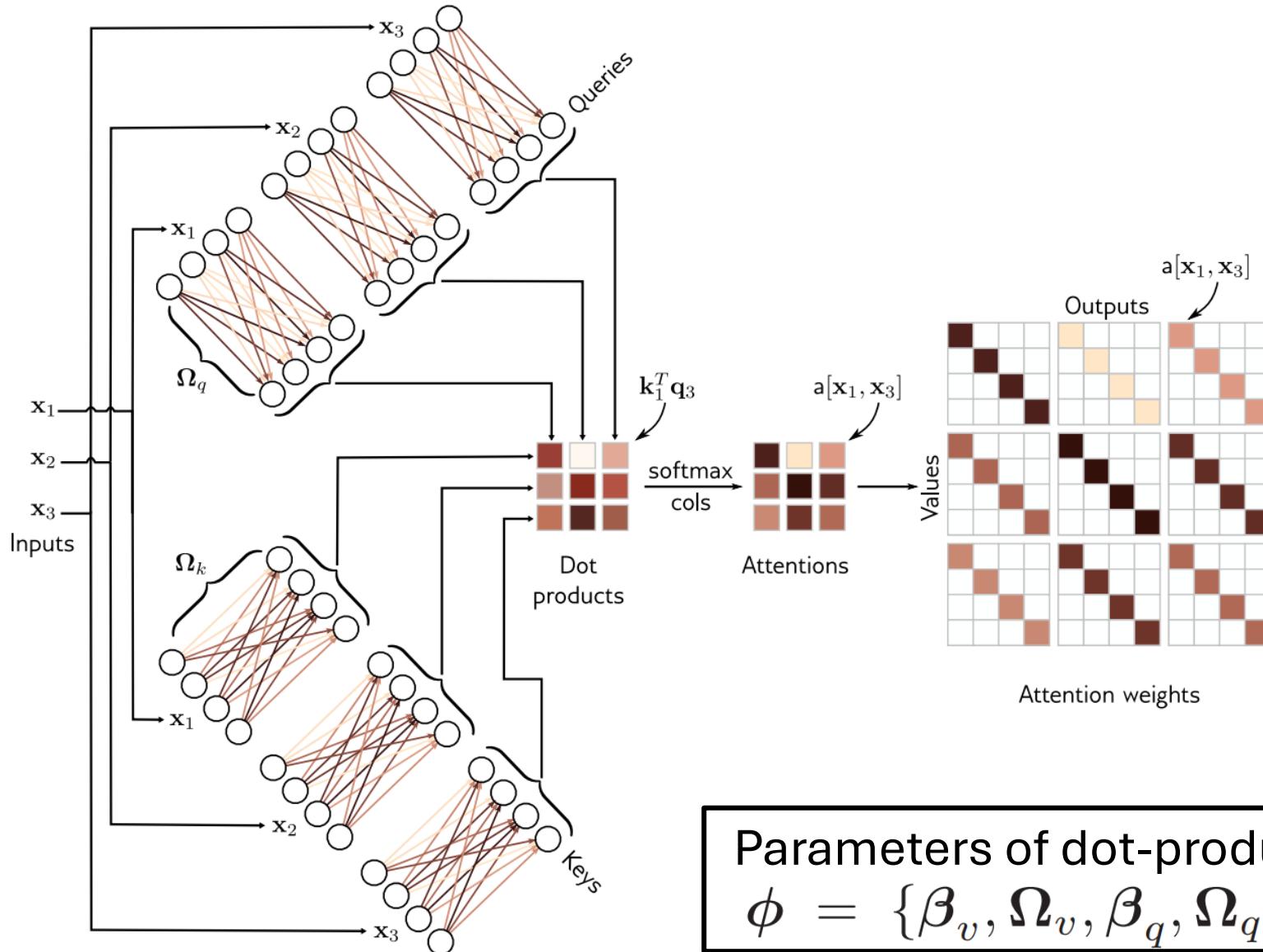
- Compute the dot products between the queries and keys and pass the results through a softmax function:

$$a[\mathbf{x}_m, \mathbf{x}_n] = \text{softmax}_m [\mathbf{k}_m^T \mathbf{q}_n] = \frac{\exp [\mathbf{k}_m^T \mathbf{q}_n]}{\sum_{m'=1}^N \exp [\mathbf{k}_{m'}^T \mathbf{q}_n]}$$

Keys compete to contribute to the result

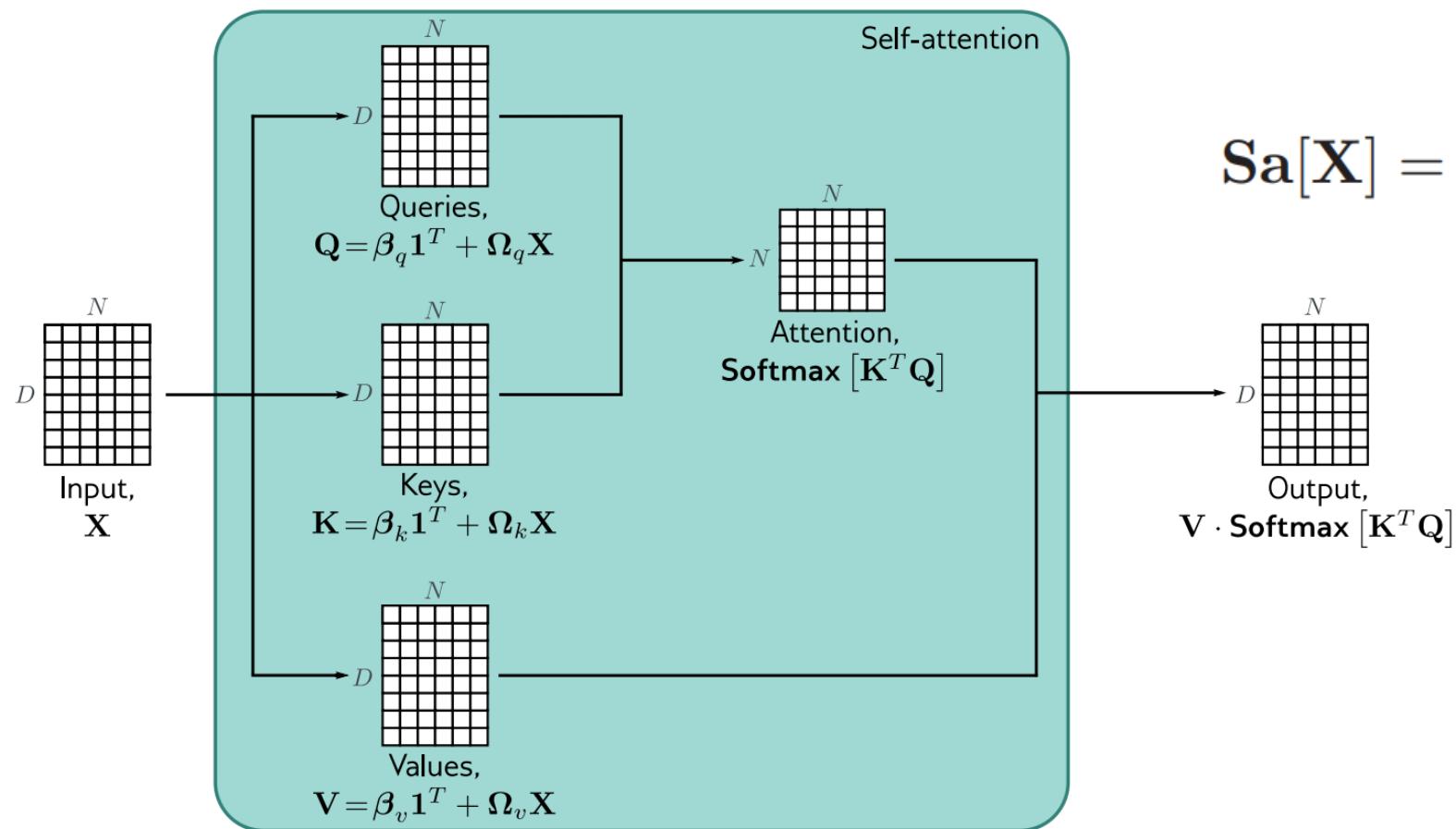
The weights $a[\mathbf{x}_\bullet, \mathbf{x}_n]$ sum to 1

Dot-product self-attention



Attention weights

Matrix form:

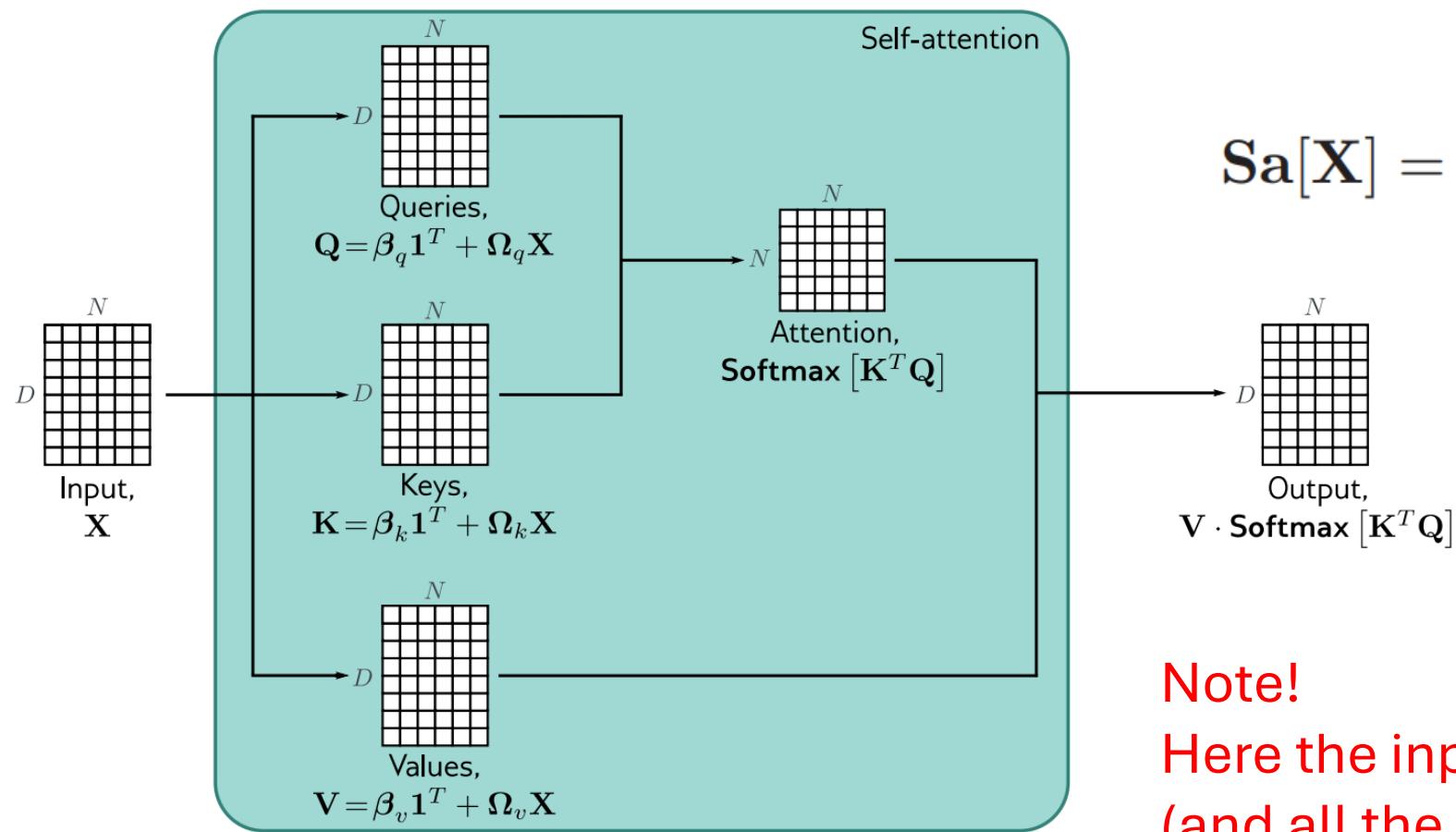


$$\begin{aligned}\mathbf{V}[X] &= \beta_v \mathbf{1}^T + \Omega_v X \\ \mathbf{Q}[X] &= \beta_q \mathbf{1}^T + \Omega_q X \\ \mathbf{K}[X] &= \beta_k \mathbf{1}^T + \Omega_k X\end{aligned}$$

$$\mathbf{Sa}[X] = \mathbf{V}[X] \cdot \text{Softmax}[\mathbf{K}[X]^T \mathbf{Q}[X]]$$

Attention weights

Matrix form:



$$\begin{aligned}\mathbf{V}[X] &= \beta_v \mathbf{1}^T + \Omega_v X \\ \mathbf{Q}[X] &= \beta_q \mathbf{1}^T + \Omega_q X \\ \mathbf{K}[X] &= \beta_k \mathbf{1}^T + \Omega_k X\end{aligned}$$

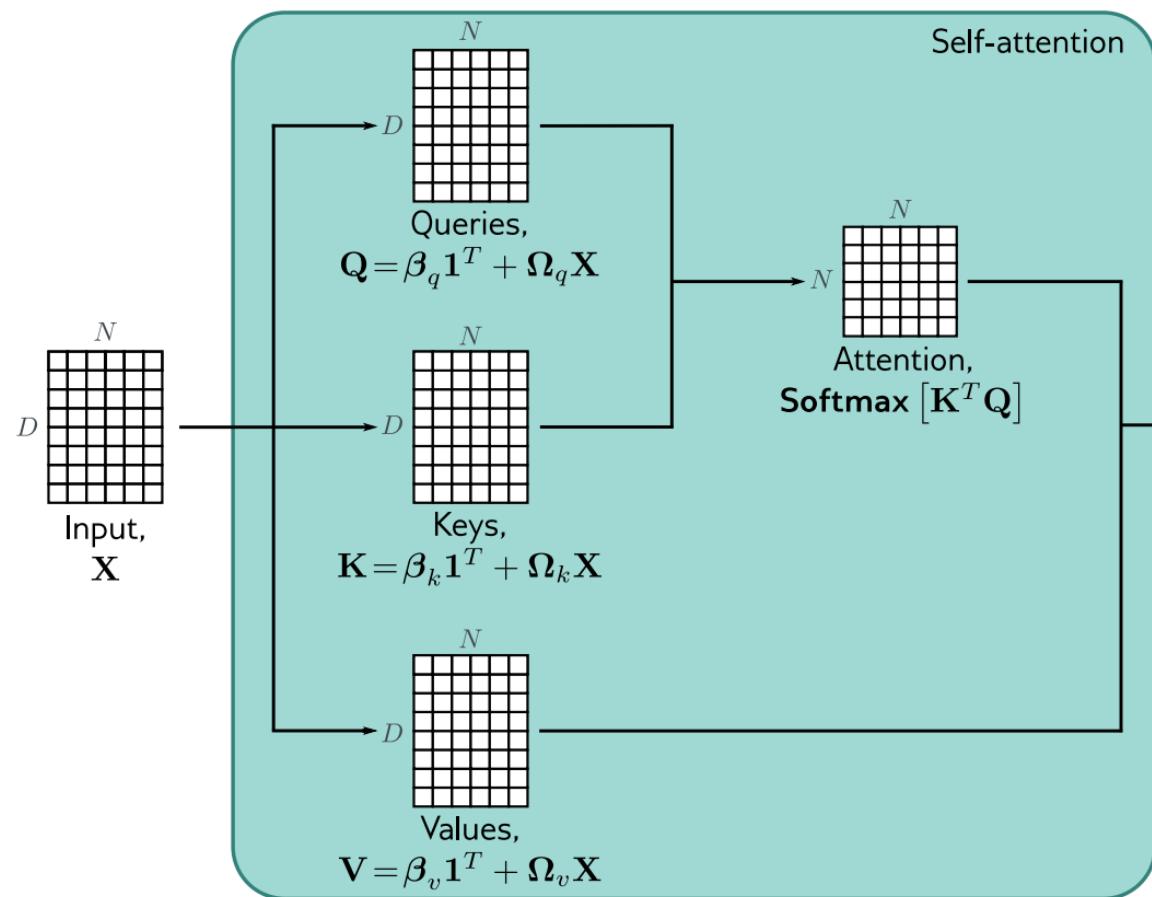
$$\mathbf{Sa}[X] = \mathbf{V}[X] \cdot \text{Softmax} [\mathbf{K}[X]^T \mathbf{Q}[X]]$$

Note!

Here the input X is a column-matrix. Often, X (and all the rest matrices) are row-matrices

Attention weights

Matrix form:



$$\begin{aligned}\mathbf{V}[X] &= \beta_v \mathbf{1}^T + \Omega_v X \\ \mathbf{Q}[X] &= \beta_q \mathbf{1}^T + \Omega_q X \\ \mathbf{K}[X] &= \beta_k \mathbf{1}^T + \Omega_k X\end{aligned}$$

$$Sa[X] = \mathbf{V}[X] \cdot \text{Softmax} [\mathbf{K}[X]^T \mathbf{Q}[X]]$$



Scaled dot-product attention

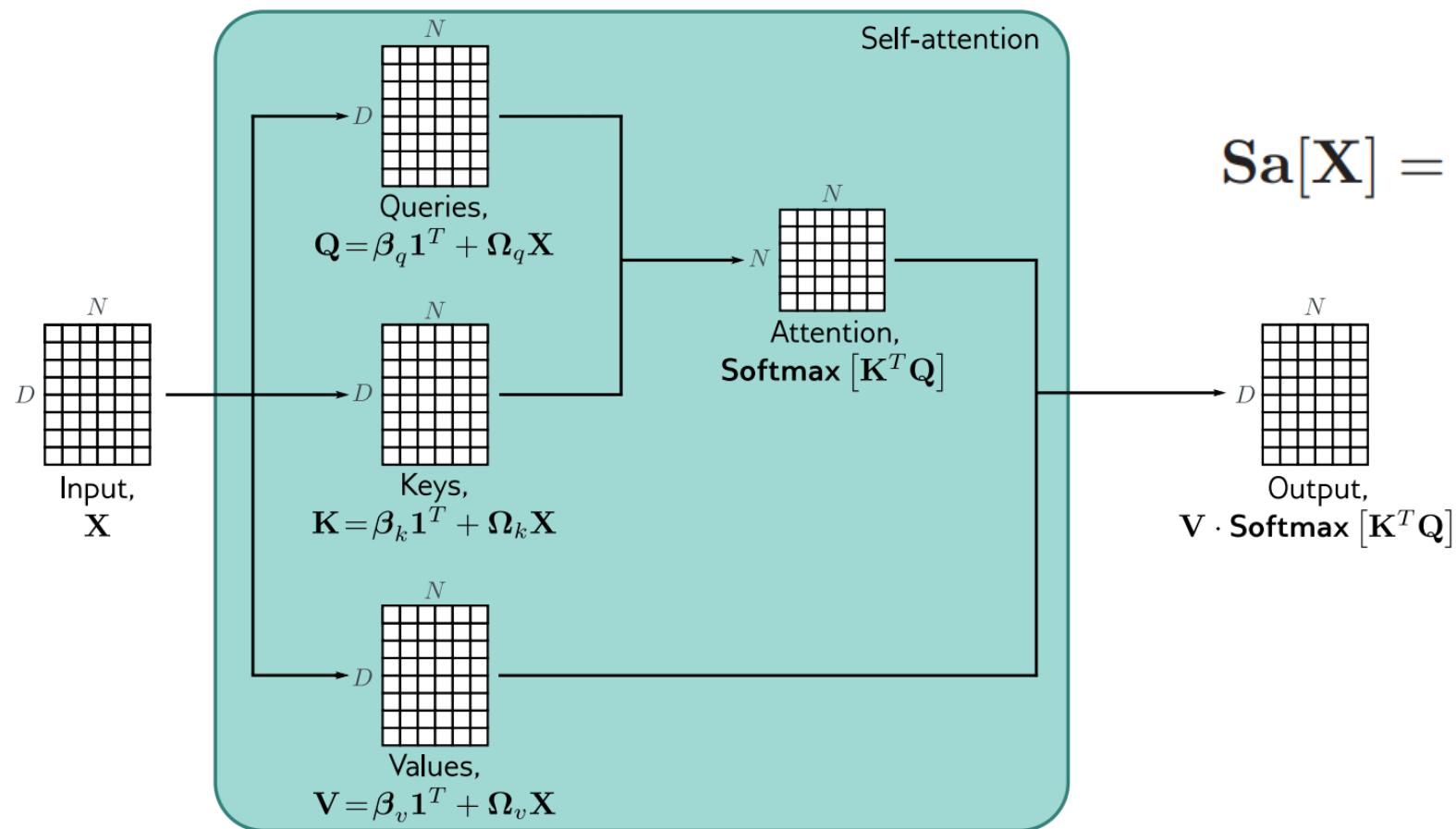
$$Sa[X] = \mathbf{V} \cdot \text{Softmax} \left[\frac{\mathbf{K}^T \mathbf{Q}}{\sqrt{D_q}} \right]$$

Note!

Here the input X is a column-matrix. Often, X (and all the rest matrices) are row-matrices

Attention weights

Matrix form:



$$\begin{aligned}\mathbf{V}[\mathbf{X}] &= \beta_v \mathbf{1}^T + \Omega_v \mathbf{X} \\ \mathbf{Q}[\mathbf{X}] &= \beta_q \mathbf{1}^T + \Omega_q \mathbf{X} \\ \mathbf{K}[\mathbf{X}] &= \beta_k \mathbf{1}^T + \Omega_k \mathbf{X}\end{aligned}$$

$$\mathbf{Sa}[\mathbf{X}] = \mathbf{V}[\mathbf{X}] \cdot \text{Softmax} \left[\mathbf{K}[\mathbf{X}]^T \mathbf{Q}[\mathbf{X}] \right]$$



Scaled dot-product attention

$$\begin{array}{c}N \\ \text{Output,} \\ \mathbf{V} \cdot \text{Softmax} [\mathbf{K}^T \mathbf{Q}]\end{array}$$

$$\mathbf{Sa}[\mathbf{X}] = \mathbf{V} \cdot \text{Softmax} \left[\frac{\mathbf{K}^T \mathbf{Q}}{\sqrt{D_q}} \right]$$

Self-attention:

Keys and Queries (defining the attention) depend on the same input as the Values

Positional encodings

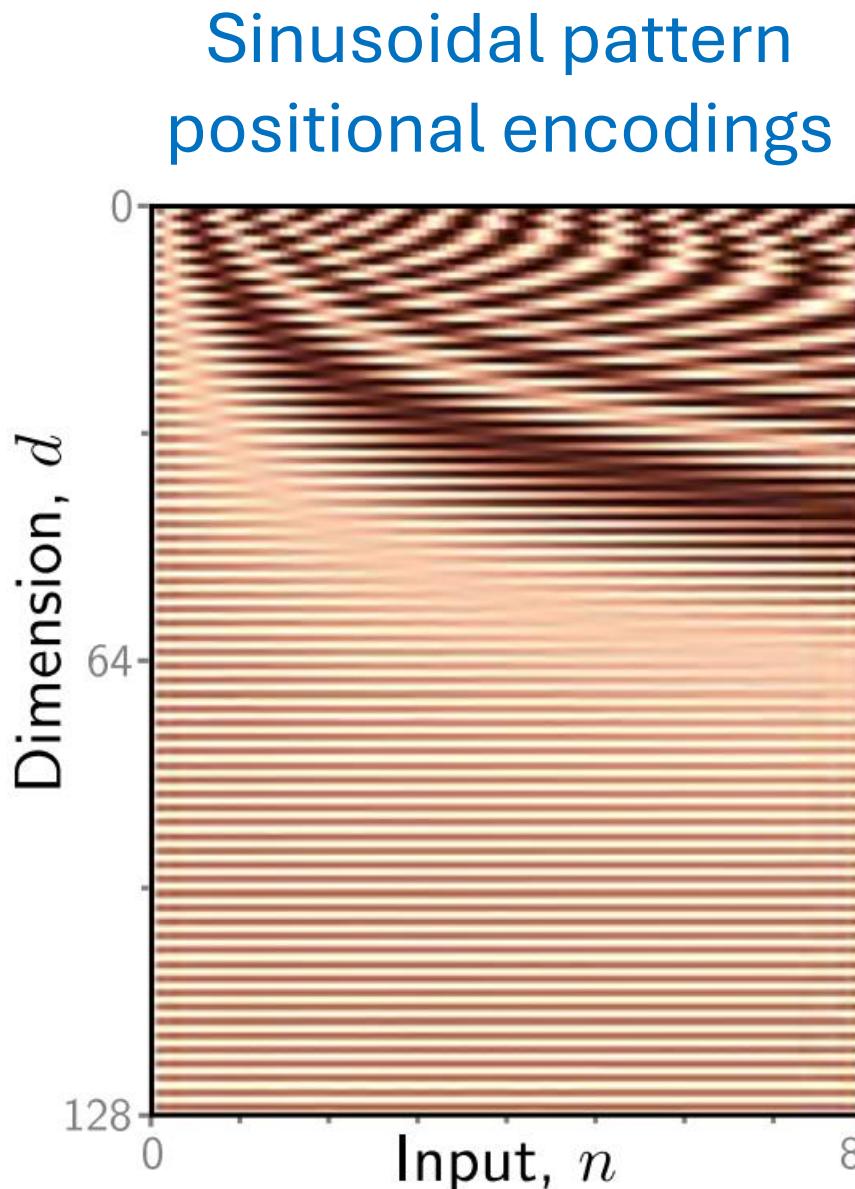
The self-attention mechanism does not consider the order of the inputs x_n :

- Multiple equivalent attentions for different permutations of the data
- Order of the inputs is important, e.g., order of words in a sentence:
 - The woman ate the racoon
 - The racoon ate the woman

Positional encodings

Absolute positional encoding:

- Add a matrix Π encoding the positional information to the input X
- Π can be predefined or learned
- Π can be added in the input or at every network layer
- Sometimes Π is added to X in the computation of queries and keys, but not to the values



Positional encodings

Relative positional encoding:

- Sometimes the absolute position of an input is less important compared to the relative position with respect to other inputs
- Each element of the attention matrix:
 - A particular offset between key position a and query position b

Process:

- Learn a parameter $\pi_{a,b}$ for each offset
- Use it to modify the attention matrix by altering it (e.g., addition, multiplication)

Scaled dot-product self-attention

The dot products in the attention computation can have large magnitudes:

- the largest value completely dominates the softmax function calculation
- Normalize the input to the softmax function

$$\text{Sa}[\mathbf{X}] = \mathbf{V} \cdot \text{Softmax} \left[\frac{\mathbf{K}^T \mathbf{Q}}{\sqrt{D_q}} \right]$$

Multiple heads

- Calculate multiple self-attentions in parallel:

$$\mathbf{V}_h = \beta_{vh} \mathbf{1}^T + \Omega_{vh} \mathbf{X}$$

$$\mathbf{Q}_h = \beta_{qh} \mathbf{1}^T + \Omega_{qh} \mathbf{X}$$

$$\mathbf{K}_h = \beta_{kh} \mathbf{1}^T + \Omega_{kh} \mathbf{X}$$

$$\mathbf{Sa}_h[\mathbf{X}] = \mathbf{V}_h \cdot \text{Softmax} \left[\frac{\mathbf{K}_h^T \mathbf{Q}_h}{\sqrt{D_q}} \right]$$

Each self-attention is called *head*, and learns different parameters $\{\Omega_{\cdot h}, \beta_{\cdot h}\}$

Multiple heads

- Calculate multiple self-attentions in parallel:

$$\mathbf{V}_h = \beta_{vh} \mathbf{1}^T + \Omega_{vh} \mathbf{X}$$

$$\mathbf{Q}_h = \beta_{qh} \mathbf{1}^T + \Omega_{qh} \mathbf{X}$$

$$\mathbf{K}_h = \beta_{kh} \mathbf{1}^T + \Omega_{kh} \mathbf{X}$$

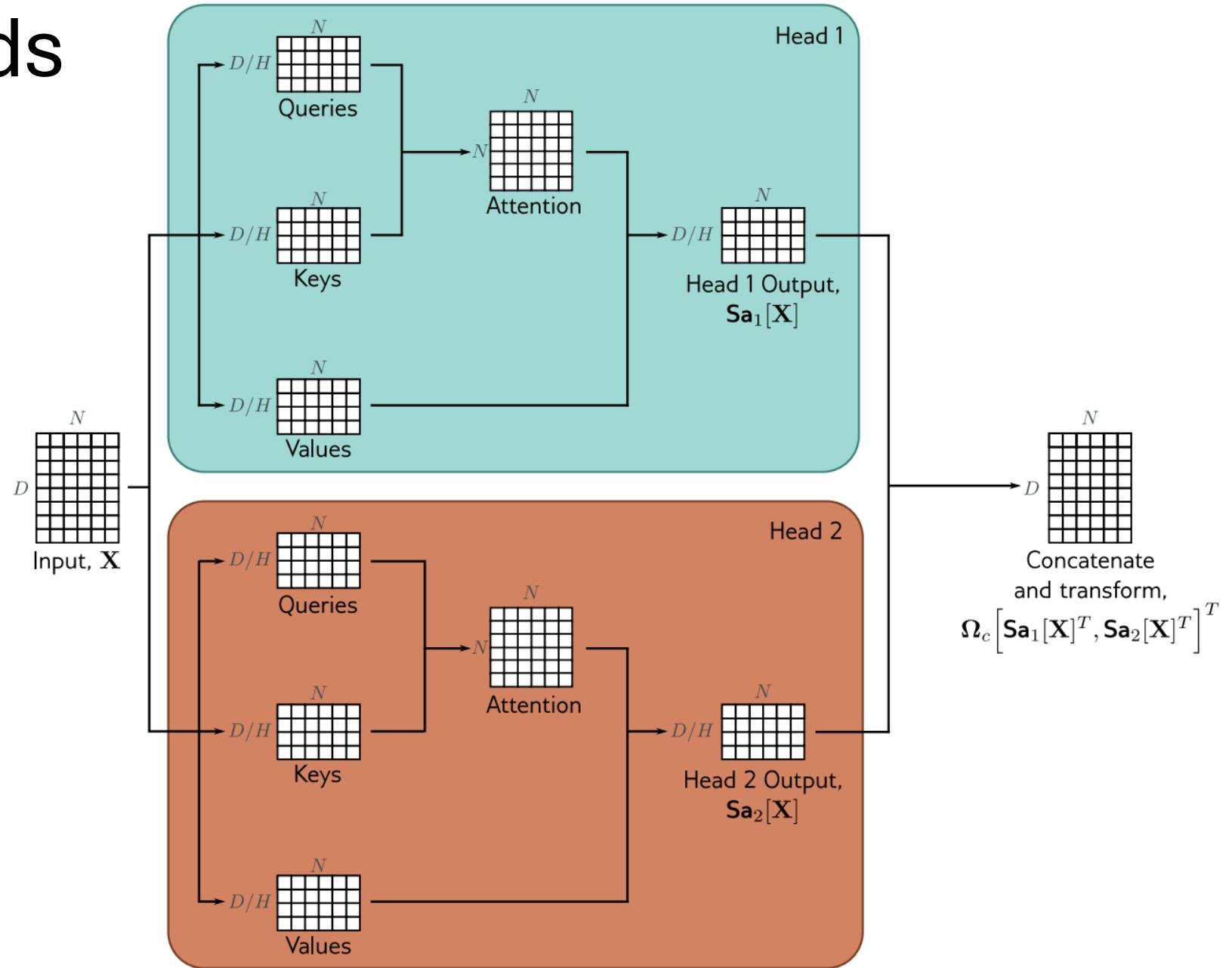
$$\mathbf{Sa}_h[\mathbf{X}] = \mathbf{V}_h \cdot \text{Softmax} \left[\frac{\mathbf{K}_h^T \mathbf{Q}_h}{\sqrt{D_q}} \right]$$

- Concatenate the outputs and perform a linear transformation for fusing the information:

$$\mathbf{MhSa}[\mathbf{X}] = \Omega_c \left[\mathbf{Sa}_1[\mathbf{X}]^T, \mathbf{Sa}_2[\mathbf{X}]^T, \dots, \mathbf{Sa}_H[\mathbf{X}]^T \right]^T$$

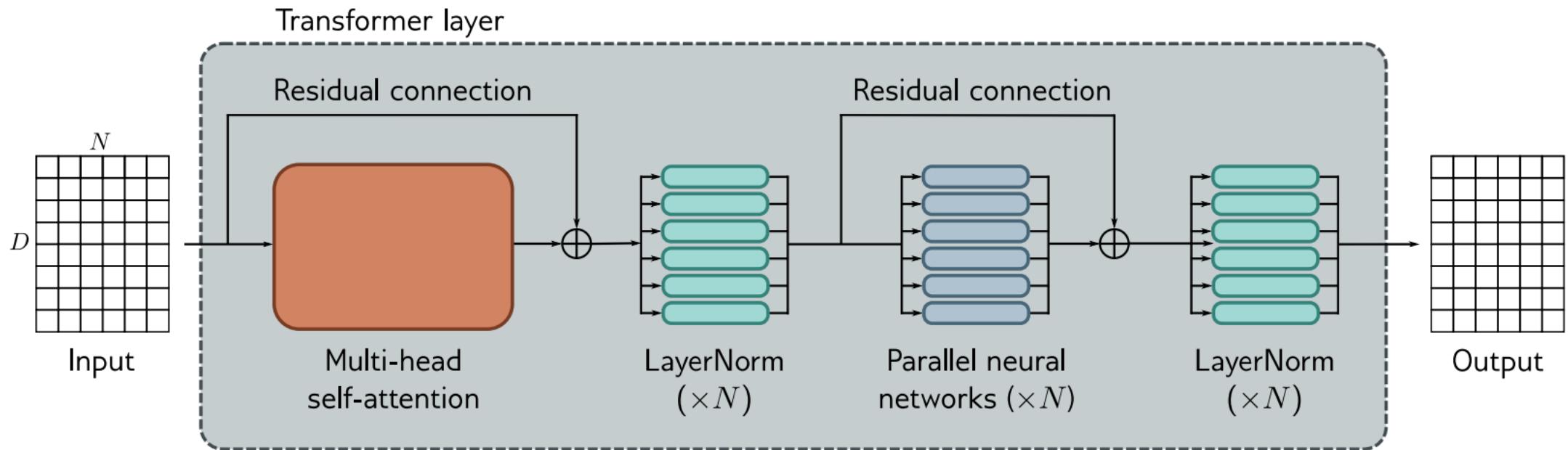
Each self-attention is called *head*, and learns different parameters $\{\Omega_{\cdot h}, \beta_{\cdot h}\}$

Multiple heads



Transformer layer

Calculations:

$$\mathbf{X} \leftarrow \mathbf{X} + \text{MhSa}[\mathbf{X}]$$
$$\mathbf{X} \leftarrow \text{LayerNorm}[\mathbf{X}]$$
$$\mathbf{x}_n \leftarrow \mathbf{x}_n + \text{mlp}[\mathbf{x}_n]$$
$$\mathbf{X} \leftarrow \text{LayerNorm}[\mathbf{X}]$$


Transformer layer

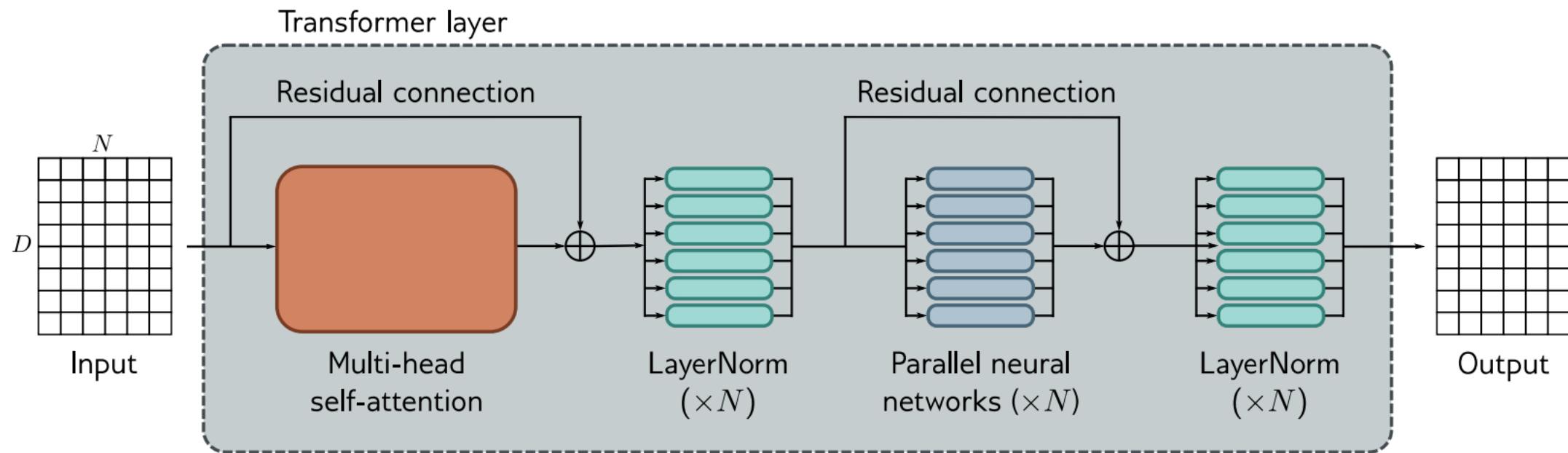
Calculations: $\mathbf{X} \leftarrow \mathbf{X} + \text{MhSa}[\mathbf{X}]$

$\mathbf{X} \leftarrow \text{LayerNorm}[\mathbf{X}]$

$\mathbf{x}_n \leftarrow \mathbf{x}_n + \text{mlp}[\mathbf{x}_n]$

$\mathbf{X} \leftarrow \text{LayerNorm}[\mathbf{X}]$

LayerNorm is similar to BatchNorm but normalizes each embedding in each batch element separately using statistics calculated across its D embedding dimensions



Transformers for Natural Language Processing

Tokenization:

- Splits the text into smaller constituent units (called *tokens*) from a vocabulary of possible tokens

Transformers for Natural Language Processing

Tokenization:

- Splits the text into smaller constituent units (called *tokens*) from a vocabulary of possible tokens
- Difficulties:
 - To include all possible words, an enormous (infinite...) vocabulary is needed
 - Unclear how to handle punctuation (e.g., question mark)
 - Unclear how to clarify/encode different versions of a word (e.g., a different token?)

Transformers for Natural Language Processing

Tokenization:

- Splits the text into smaller constituent units (called *tokens*) from a vocabulary of possible tokens
- Difficulties:
 - To include all possible words, an enormous (infinite...) vocabulary is needed
 - Unclear how to handle punctuation (e.g., question mark)
 - Unclear how to clarify/encode different versions of a word (e.g., a different token?)
- Compromise between letters and full words
 - Include both common words and word fragments
 - Those can be used to create larger and less frequent words

Transformers for Natural Language Processing

Example sub-word tokenizer:

- a) The tokens are initially just the characters and whitespace (represented by an underscore), and their frequencies are displayed in the table.
- b) At each iteration, the sub-word tokenizer looks for the most commonly occurring adjacent pair of tokens (in this case, se) and merges them
- c) At the second iteration, the algorithm merges e and the whitespace character_.
- d) After 22 iterations, the tokens consist of a mix of letters, word fragments, and commonly occurring words.
- e) If we continue this process indefinitely, the tokens eventually represent the full words

In a real situation, the algorithm would terminate when the vocabulary size (number of tokens) reached a predetermined value

a)

a_sailor_went_to_sea_sea_sea_
to_see_what_he_could_see_see_see_
but_all_that_he_could_see_see_see_see_
was_the_bottom_of_the_deep_blue_sea_sea_sea_sea_

_	e	s	a	t	o	h	l	u	b	d	w	c	f	i	m	n	p	r
33	28	15	12	11	8	6	6	4	3	3	3	2	1	1	1	1	1	1

b)

a_sailor_went_to_sea_sea_sea_
to_see_what_he_could_see_see_see_
but_all_that_he_could_see_see_see_see_
was_the_bottom_of_the_deep_blue_sea_sea_sea_

_	e	se	a	t	o	h	l	u	b	d	w	c	s	f	i	m	n	p	r
33	15	13	12	11	8	6	6	4	3	3	3	2	2	1	1	1	1	1	1

c)

a_sailor_went_to_sea_sea_sea_
to_see_what_he_could_see_see_see_
but_all_that_he_could_see_see_see_see_
was_the_bottom_of_the_deep_blue_sea_sea_sea_

_	se	a	e	t	o	h	l	u	b	d	e	w	c	s	f	i	m	n	p	r
21	13	12	12	11	8	6	6	4	3	3	3	3	2	2	1	1	1	1	1	1

d)

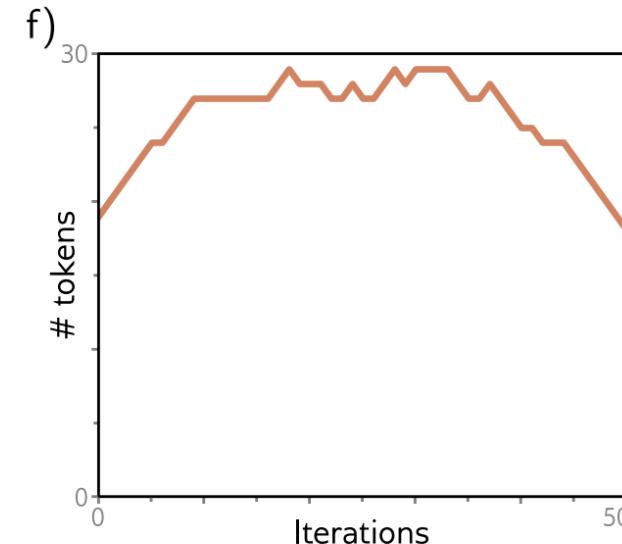
see_sea_e_b_l_w_a_could_hat_he_o_t_t_the_to_u_a_d_f_m_n_p_s_sailor_to

see_	sea_	e	b	l	w	a	could_	hat_	he_	o	t	t	the_	to_	u	a	d	f	m	n	p	s	sailor_	to
7	6	4	3	3	3	3	2	2	2	2	2	2	2	2	2	1	1	1	1	1	1	1	1	1

e)

see_sea_could_he_the_a_all_blue_bottom_but_deep_of_sailor_that_to_was_went_what_

see_	sea_	could_	he_	the_	a	all_	blue_	bottom_	but_	deep_	of_	sailor_	that_	to_	was_	went_	what_
7	6	2	2	2	1	1	1	1	1	1	1	1	1	1	1	1	1



Transformers for Natural Language Processing

Embeddings:

- Each token in the vocabulary V is mapped to a unique *word embedding* stored in a matrix $\Omega_e \in \mathbb{R}^{D \times |\mathcal{V}|}$

Transformers for Natural Language Processing

Embeddings:

- Each token in the vocabulary V is mapped to a unique *word embedding* stored in a matrix $\Omega_e \in \mathbb{R}^{D \times |V|}$

$$\Omega_e \in \mathbb{R}^{D \times |V|}$$

1024 30,000

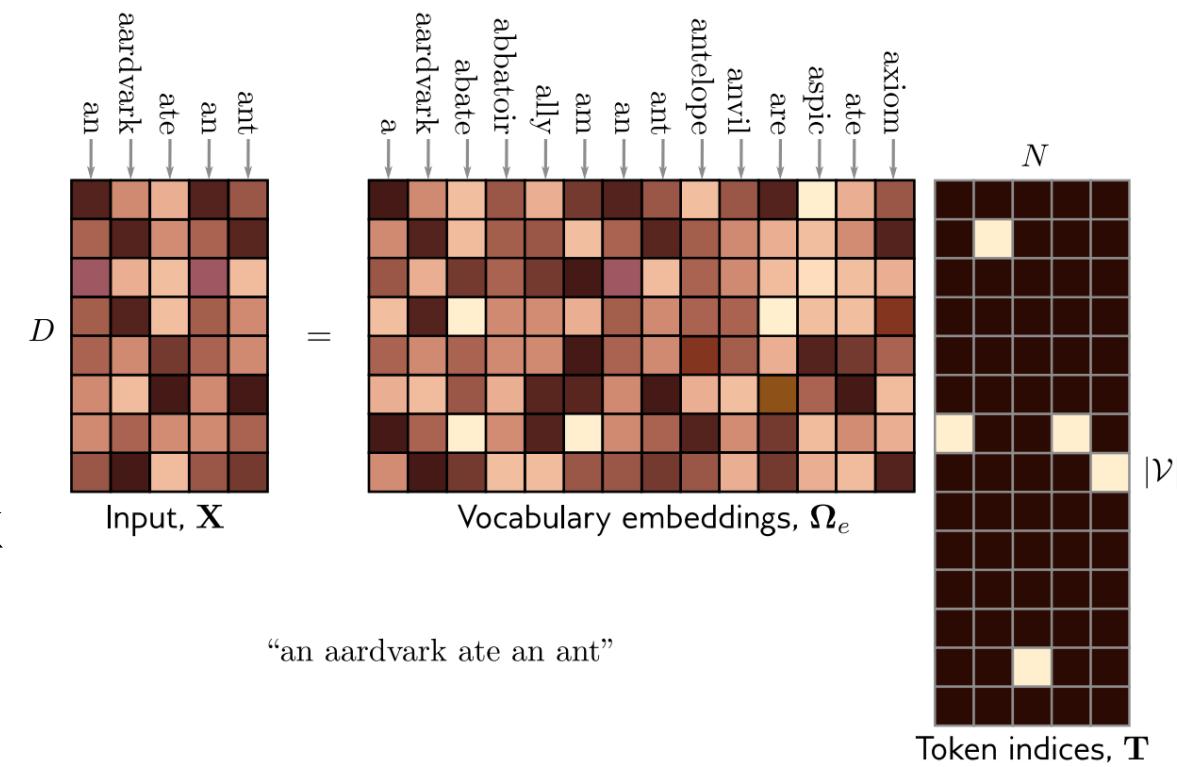
Transformers for Natural Language Processing

Embeddings:

- Each token in the vocabulary V is mapped to a unique *word embedding* stored in a matrix $\Omega_e \in \mathbb{R}^{D \times |V|}$

Learning embeddings:

- Create the one-hot vector-based matrix $T \in \mathbb{R}^{|V| \times N}$
- Learn the word embedding $X = \Omega_e T$ like any other network parameter



Transformers for Natural Language Processing

Transformer models:

- Encoder: transforms the text embeddings to a representation that can support different tasks
- Decoder: predicts the next token to continue the input text
- Encoder-decoder: used in sequence-to-sequence tasks (e.g., machine translation)

Encoder model example (BERT)

The Bidirectional Encoder Representations from Transformers (BERT):

- Vocabulary: 30,000 tokens
- Word embeddings: 1024 dimensions
- Structure: 24 Transformer layers, each with 16 heads
 - 64-dimensional queries, keys, values
 - MLP layer of 4096 units

Encoder model example (BERT)

The Bidirectional Encoder Representations from Transformers (BERT):

- Vocabulary: 30,000 tokens
- Word embeddings: 1024 dimensions
- Structure: 24 Transformer layers, each with 16 heads

64-dimensional queries, keys, values

MLP layer of 4096 units



≈340 million
parameters

Encoder model example (BERT)

The Bidirectional Encoder Representations from Transformers (BERT):

- Vocabulary: 30,000 tokens
- Word embeddings: 1024 dimensions
- Structure: 24 Transformer layers, each with 16 heads
 - 64-dimensional queries, keys, values
 - MLP layer of 4096 units
- Training using *Transfer Learning*:
 - *Pre-training* using *self-supervision* on a large dataset
 - *Fine-tuning* on a smaller labelled dataset to adapt to a downstream task

Encoder model example (BERT)

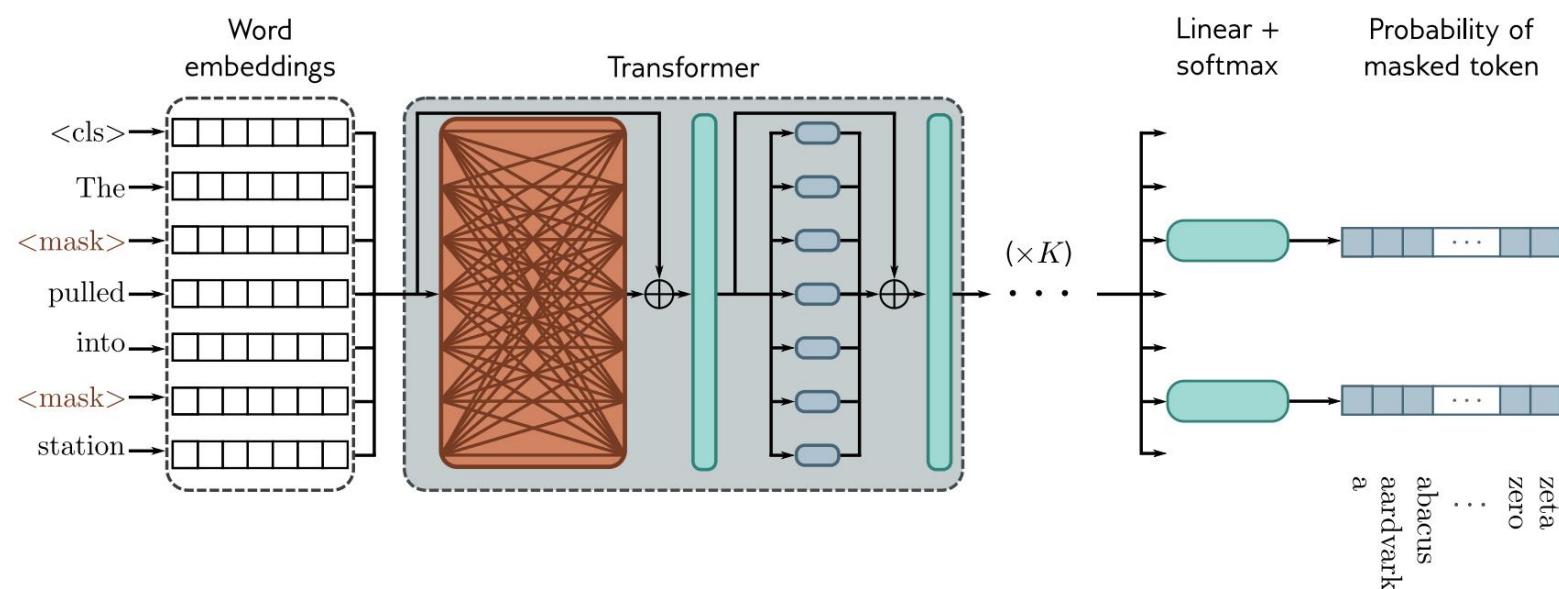
Pre-training using self-supervision:

- Allows for using enormous amounts of data without needing labels

Tasks: 1) Predict missing words from sentences

2) Predict whether two sentences were originally adjacent

- Maximum input length of 512 tokens
- Batch size of 256
- Dataset size: 3.3-billion-word corpus
- 1 million training updates (~50 epochs)



Encoder model example (BERT)

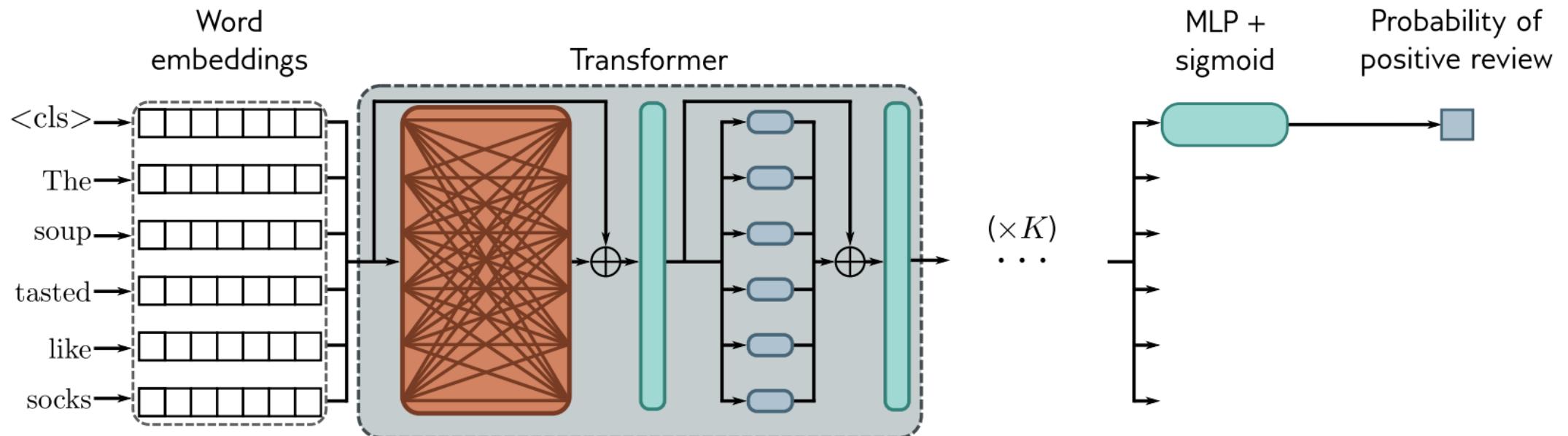
Fine-tuning:

- Initialize the model with the parameters obtained by the pre-training
- Add an extra layer at the end to fit the type and output of the downstream task
- Adjust/update the network parameters with a labeled dataset (some parameters may be frozen / not updated)

Encoder model example (BERT)

Downstream tasks:

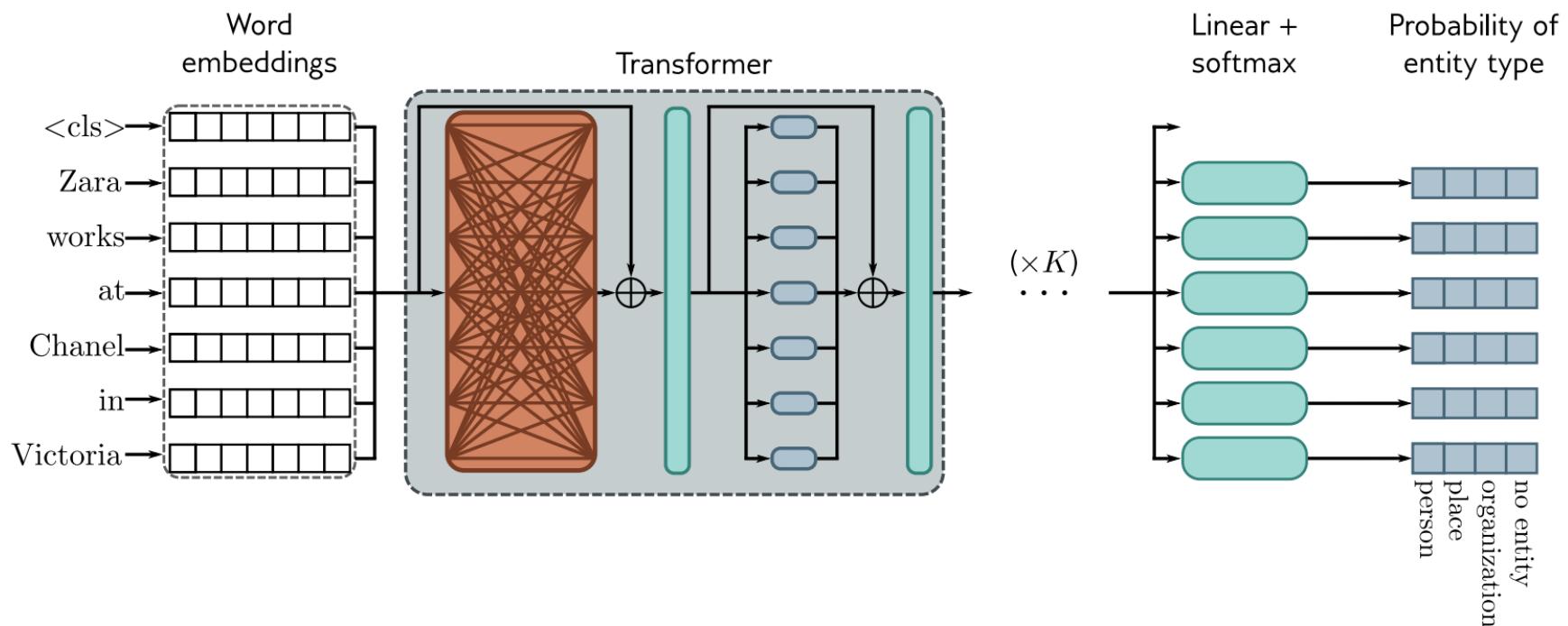
- Text classification: Use a special output `<cls>` token of the final Transformer layer as input to a classification layer



Encoder model example (BERT)

Downstream tasks:

- Word classification (named entity recognition): Classify each output word embedding from the last Transformer layer using a softmax layer

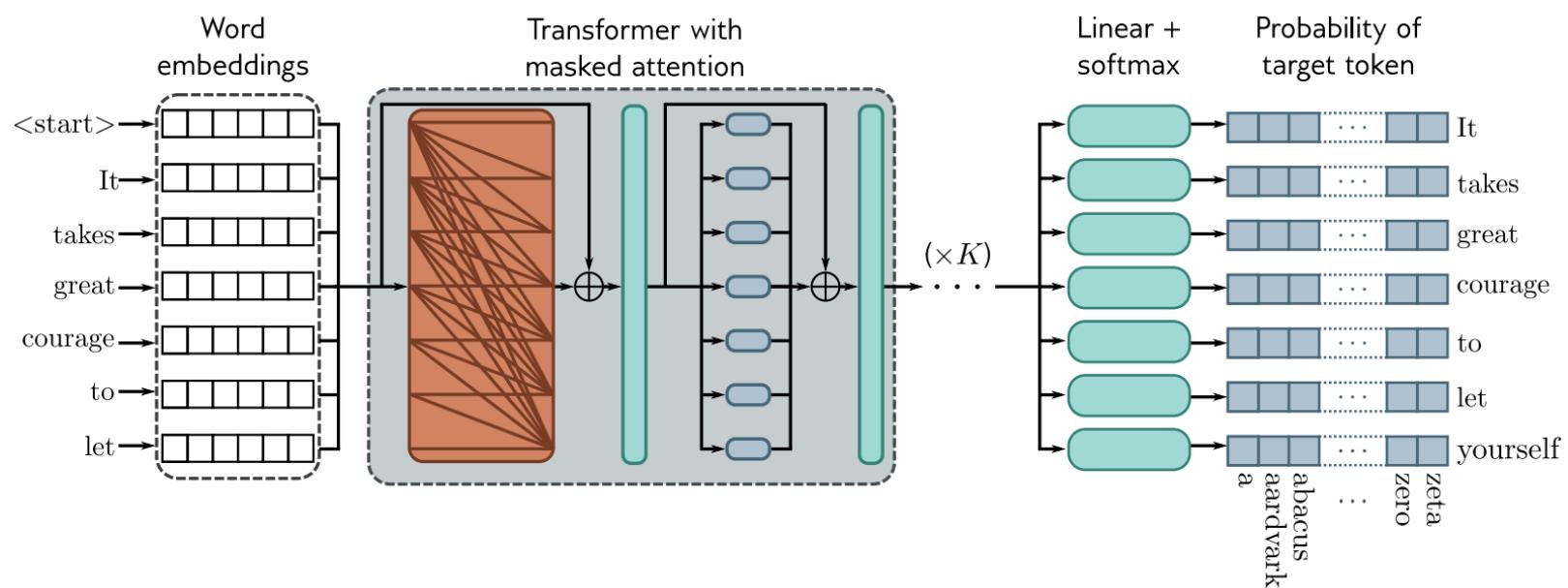


Decoder model example (GPT3)

The Generative Pre-trained Transformer version 3 (GPT3):

- Generates the next token in a sequence

- Sequence length: 2048 tokens
- Total batch size: 3.2 million tokens
- Architecture:
 - 96 Transformer layers
 - word embedding size 12,288
 - 96 heads in the self-attention
 - Value/Key dimension: 128
- Parameter count: 176 billion
- Dataset size: 300-billion-word corpus



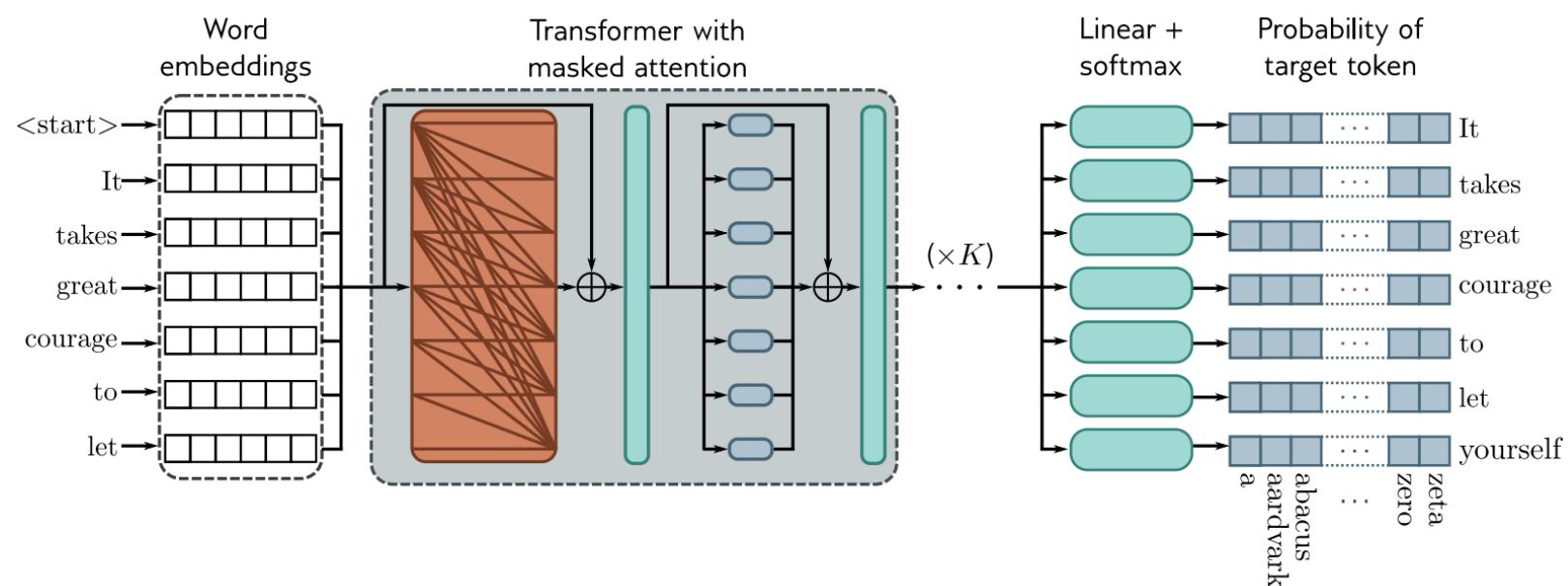
Decoder model example (GPT3)

The Generative Pre-trained Transformer version 3 (GPT3):

- Generates the next token in a sequence

It can generate a coherent text passage by feeding the extended sequence back into the model (*autoregressive model*)

- Sequence length: 2048 tokens
- Total batch size: 3.2 million tokens
- Architecture:
 - 96 Transformer layers
 - word embedding size 12,288
 - 96 heads in the self-attention
 - Value/Key dimension: 128
- Parameter count: 176 billion
- Dataset size: 300-billion-word corpus



Decoder model example (GPT3)

Example: Sentence **It takes great courage to let yourself appear weak**

$$\begin{aligned} Pr(\text{It takes great courage to let yourself appear weak}) &= \\ Pr(\text{It}) \times Pr(\text{takes}|\text{It}) \times Pr(\text{great}|\text{It takes}) \times Pr(\text{courage}|\text{It takes great}) \times \\ Pr(\text{to}|\text{It takes great courage}) \times Pr(\text{let}|\text{It takes great courage to}) \times \\ Pr(\text{yourself}|\text{It takes great courage to let}) \times \\ Pr(\text{appear}|\text{It takes great courage to let yourself}) \times \\ Pr(\text{weak}|\text{It takes great courage to let yourself appear}). \end{aligned}$$

Decoder model example (GPT3)

Example: Sentence **It takes great courage to let yourself appear weak**

$$\begin{aligned} Pr(\text{It takes great courage to let yourself appear weak}) &= \\ Pr(\text{It}) \times Pr(\text{takes}|\text{It}) \times Pr(\text{great}|\text{It takes}) \times Pr(\text{courage}|\text{It takes great}) \times \\ Pr(\text{to}|\text{It takes great courage}) \times Pr(\text{let}|\text{It takes great courage to}) \times \\ Pr(\text{yourself}|\text{It takes great courage to let}) \times \\ Pr(\text{appear}|\text{It takes great courage to let yourself}) \times \\ Pr(\text{weak}|\text{It takes great courage to let yourself appear}). \end{aligned}$$

- Predict the conditional probabilities $Pr(t_n|t_1, \dots, t_{n-1})$ of each token given all the prior tokens.
- Thus, it can calculate the joint probability:

$$Pr(t_1, t_2, \dots, t_N) = Pr(t_1) \prod_{n=2}^N Pr(t_n|t_1, \dots, t_{n-1})$$

Decoder model example (GPT3)

Masked self-attention:

- Using the full sentence for learning the probability of each token is not appropriate

E.g., computing $\log [Pr(\text{great}|\text{It takes})]$ would have access to both:

- The answer **great**
- The right context (rest of the sentence) **courage to let yourself appear weak**

Decoder model example (GPT3)

Masked self-attention:

- Using the full sentence for learning the probability of each token is not appropriate

E.g., computing $\log [Pr(\text{great}|\text{It takes})]$ would have access to both:

- The answer **great**
- The right context (rest of the sentence) **courage to let yourself appear weak**
- Observation → Tokens interact only in the self-attention layers



Make the attention to the answer and the context equal to zero

Decoder model example (GPT3)

Masked self-attention:

- Using the full sentence for learning the probability of each token is not appropriate

E.g., computing $\log [Pr(\text{great}|\text{It takes})]$ would have access to both:

- The answer **great**
- The right context (rest of the sentence) **courage to let yourself appear weak**
- Observation → Tokens interact only in the self-attention layers



Make the attention to the answer and the context equal to zero



Set the corresponding elements in the self-attention to $-\infty$

Decoder model example (GPT3)

Masked self-attention:

- Using the full sentence for learning the probability of each token is not appropriate

E.g., computing $\log [Pr(\text{great}|\text{It takes})]$ would have access to both:

- The answer **great**
- The right context (rest of the sentence) **courage to let yourself appear weak**
- Observation → Tokens interact only in the self-attention layers

Masked
self-attention

Each token attends only
to the previous tokens

Make the attention to the answer and the context equal to zero



Set the corresponding elements in the self-attention to $-\infty$

Decoder model example (GPT3)

Generating text:

- Feed with an input sequence text (e.g., a special <start> token)
- The network outputs the probabilities over possible subsequent tokens

Decoder model example (GPT3)

Generating text:

- Feed with an input sequence text (e.g., a special `<start>` token)
- The network outputs the probabilities over possible subsequent tokens
- Pick the next token of the sequence:
 - Pick the most likely token, or
 - Sample from the output probability distribution
- Introduce the extended sequence to the network

Decoder model example (GPT3)

Generating text:

- Feed with an input sequence text (e.g., a special <start> token)
- The network outputs the probabilities over possible subsequent tokens
- Pick the next token of the sequence:
 - Pick the most likely token, or
 - Sample from the output probability distribution
- Introduce the extended sequence to the network

Prior embeddings do not depend
on subsequent ones due to the
masked self-attention



Earlier computations can be
recycled as we generate
subsequent tokens (KV-cache)

Encoder-decoder model example

Machine translation from English to French:

- The Encoder receives the sentence in English and outputs a representation for each input token

Encoder-decoder model example

Machine translation from English to French:

- The Encoder receives the sentence in English and outputs a representation for each input token
- The Decoder:
 - receives the ground truth translation in French and passes it through a series of Transformer layers (that use masked self-attention) to predict the following word at each position
 - each of its layers attends to the output of the Encoder

Encoder-decoder model example

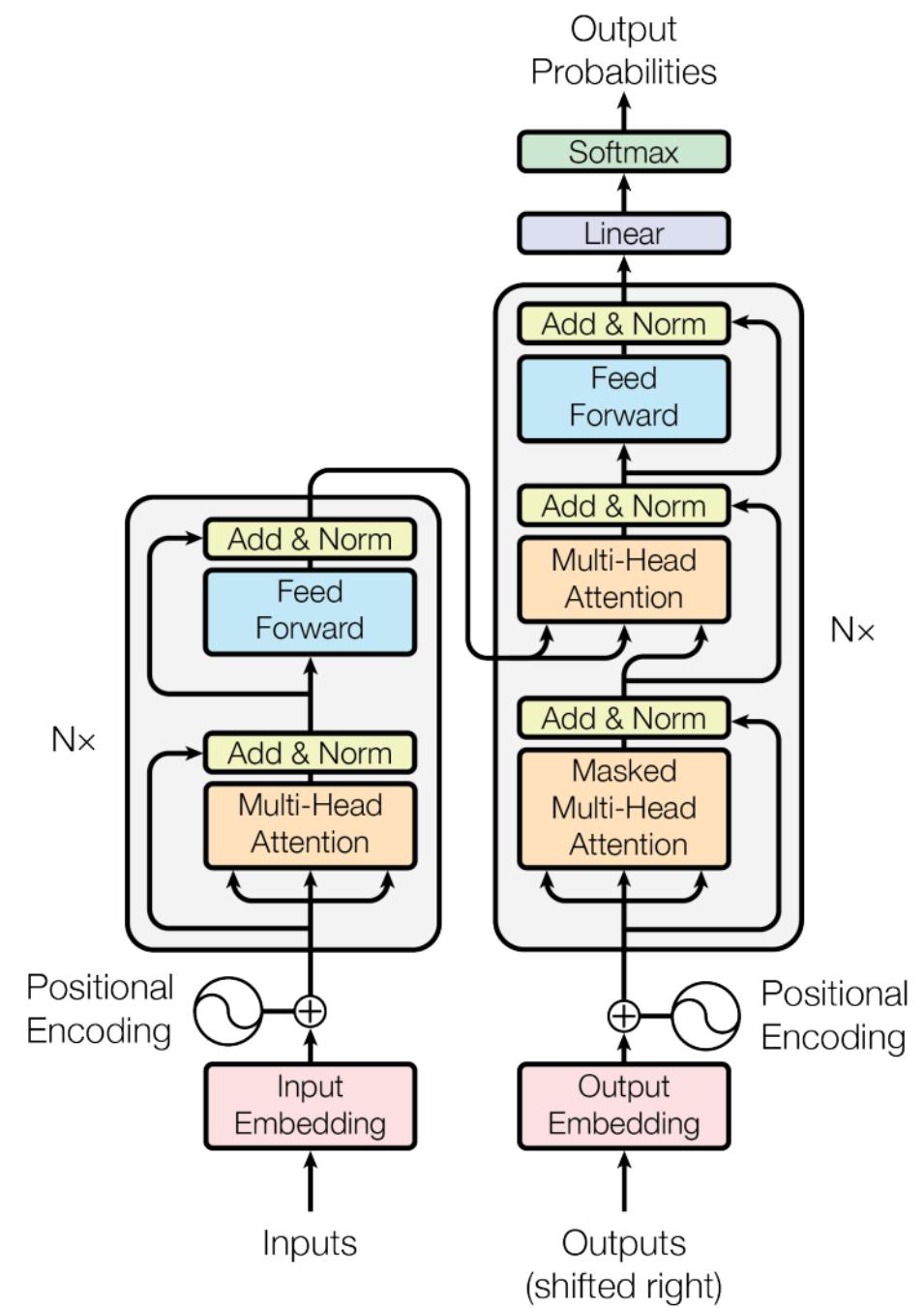
Machine translation from English to French:

- The Encoder receives the sentence in English and outputs a representation for each input token
- The Decoder:
 - receives the ground truth translation in French and passes it through a series of Transformer layers (that use masked self-attention) to predict the following word at each position
 - each of its layers attends to the output of the Encoder

Each French word is conditioned on:

- the previous output words
- the source English sentence

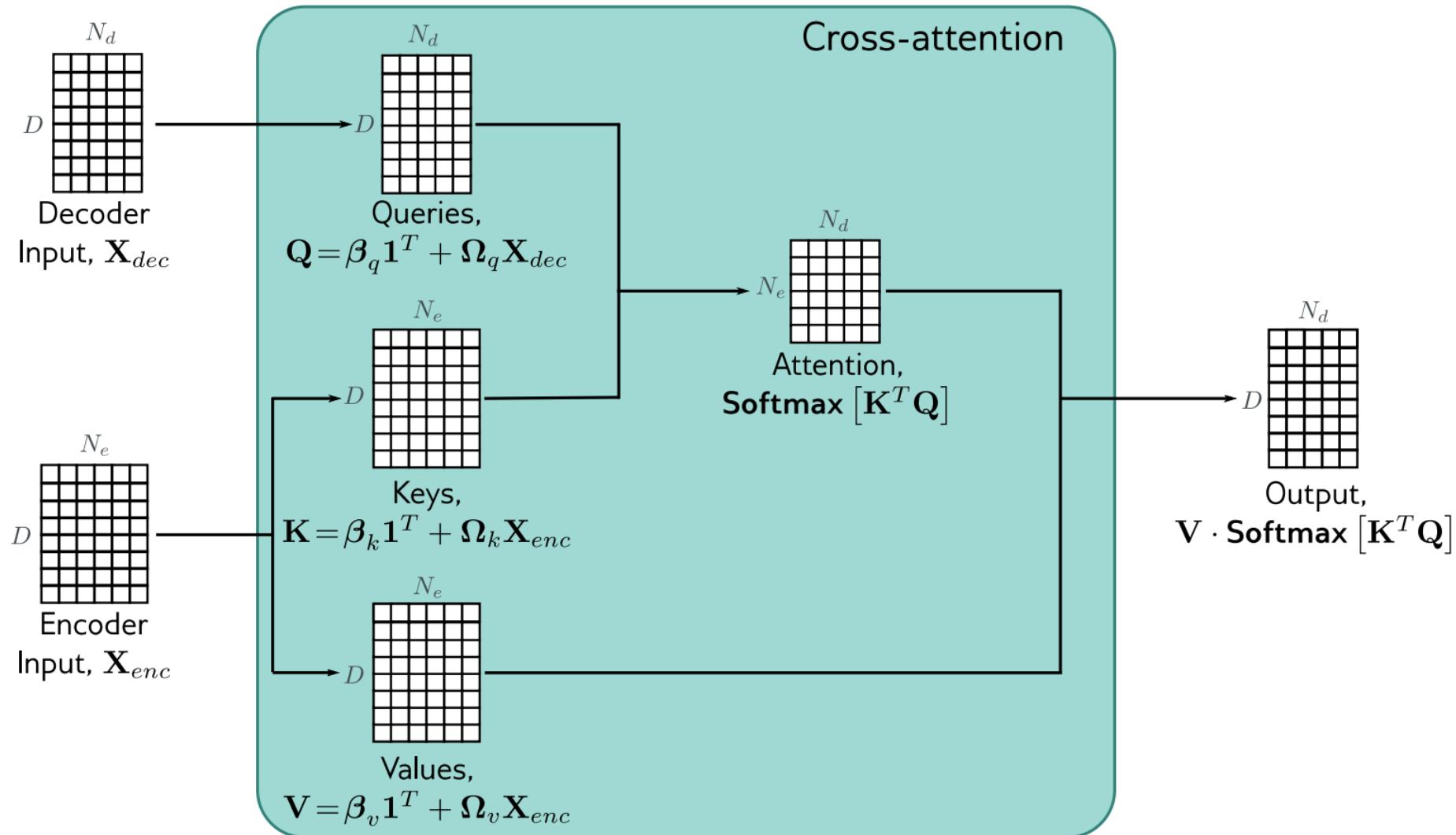
Cross-Attention



Vaswani et al., “Attention is all you need”, NIPS 2017

Figure 1: The Transformer - model architecture.

Cross-Attention



Processing long sequences

Computational complexity
of self-attention scales
quadratically with the
length of the sequence
→ impractical for long
sequences

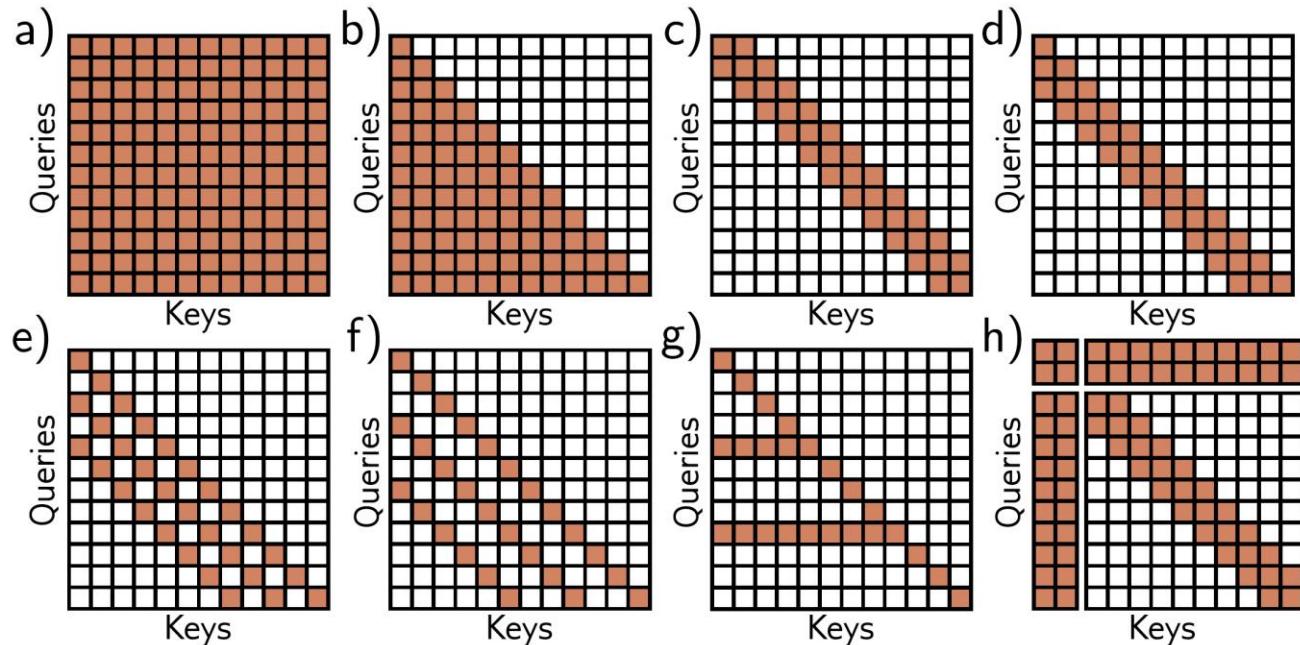


Figure 12.15 Interaction matrices for self-attention. a) In an encoder, every token interacts with every other token, and computation expands quadratically with the number of tokens. b) In a decoder, each token only interacts with the previous tokens, but complexity is still quadratic. c) Complexity can be reduced by using a convolutional structure (encoder case). d) Convolutional structure for decoder case. e–f) Convolutional structure with dilation rate of two and three (decoder case). g) Another strategy is to allow selected tokens to interact with all the other tokens (encoder case) or all the previous tokens (decoder case pictured). h) Alternatively, global tokens can be introduced (left two columns and top two rows). These interact with all of the tokens as well as with each other.

Processing long sequences

Note!
Here the input (and all the rest matrices) is row-matrix

Example model: Nyströmformer approximates the Scaled Dot-Product Attention with low-rank approximation:

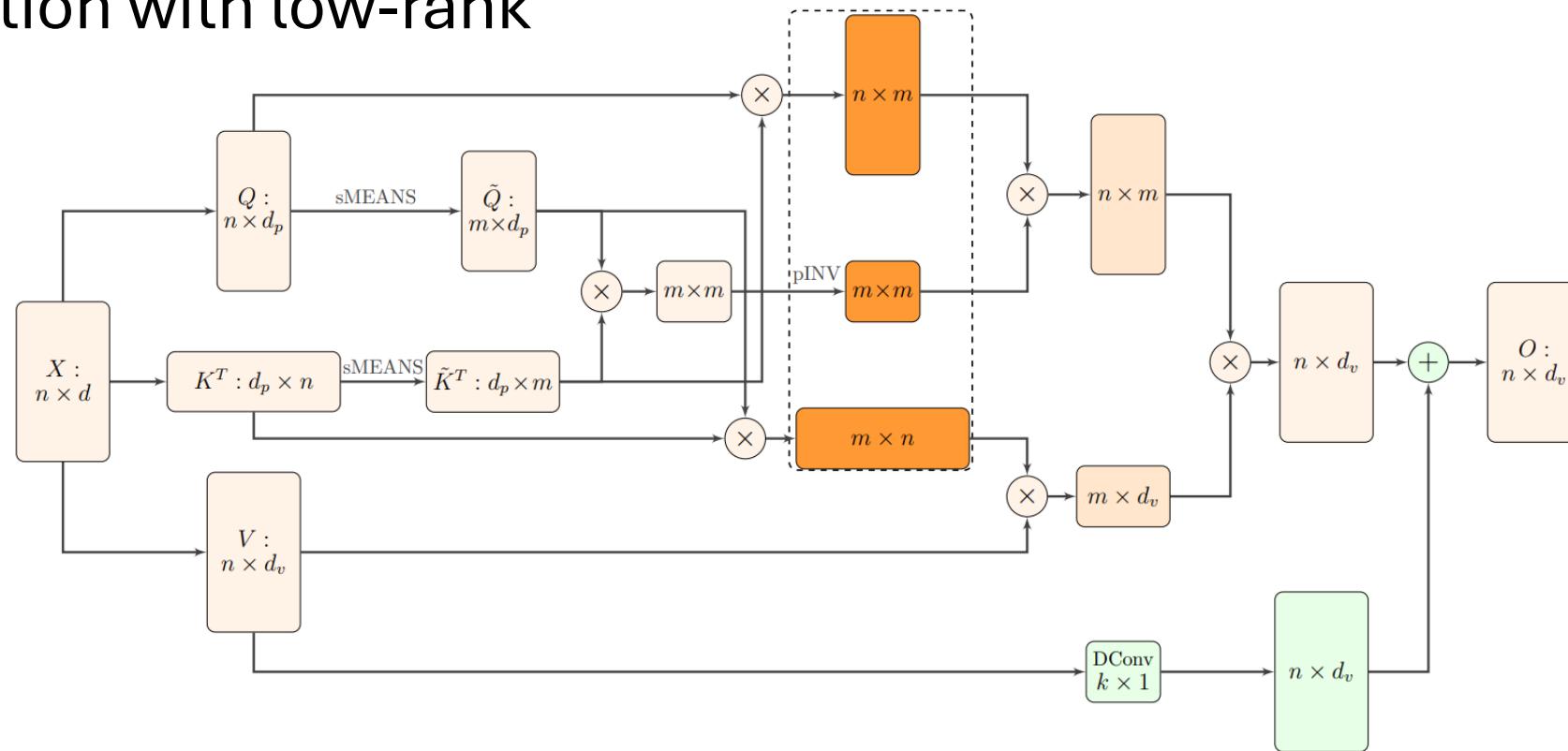
$$S \left(\frac{QK^T}{\sqrt{d}} \right) \approx S \left(\frac{Q\tilde{K}^T}{\sqrt{d}} \right) S \left(\frac{\tilde{Q}\tilde{K}^T}{\sqrt{d}} \right)^\dagger S \left(\frac{\tilde{Q}K^T}{\sqrt{d}} \right)$$

where Ω^\dagger is the Moore-Penrose pseudo-inverse of matrix Ω , and \tilde{Q} and \tilde{K} are matrices formed by sets of m landmarks, computed as the segment means of the matrices Q and K , respectivel

Processing long sequences

Note!
Here the input (and all the rest matrices) is row-matrix

Example model: Nyströmformer approximates the Scaled Dot-Product Attention with low-rank



Transformers for images

ImageGPT is a Transformer decoder:

- Builds an autoregressive model:
 - Input: partial image
 - Output: next pixel value
- Largest model (6.8 billion parameters) can operate on 64x64 pixel images and nine-bit color space quantization, so the system ingests (and predicts) one of 512 possible tokens at each position

Transformers for images

ImageGPT

Image generation

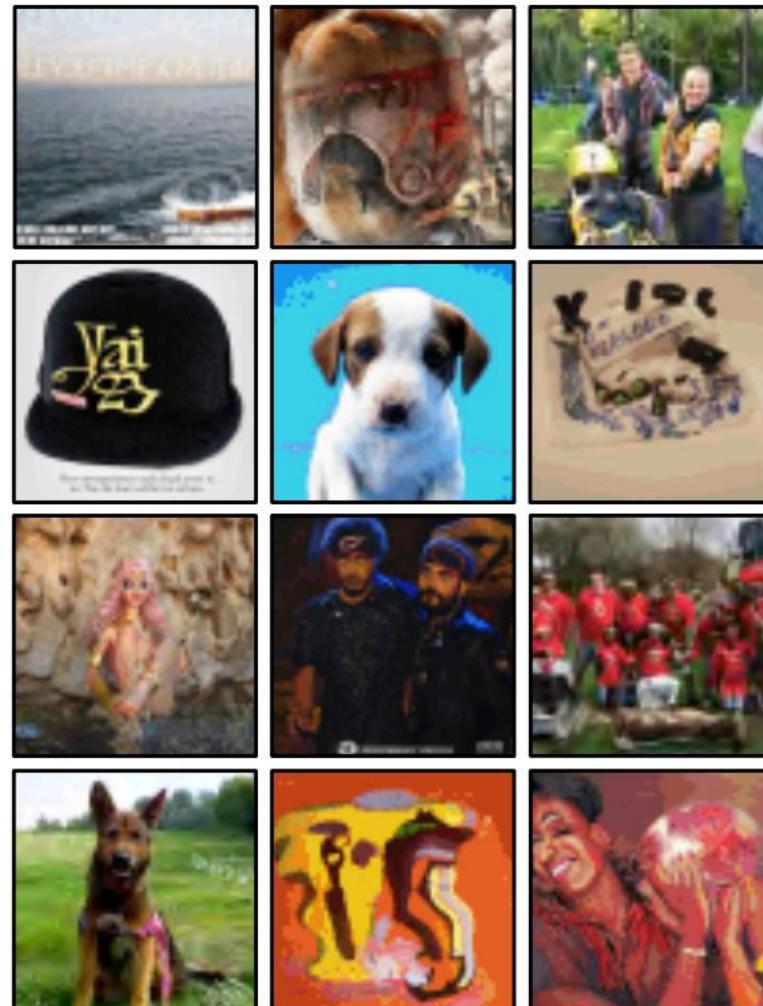
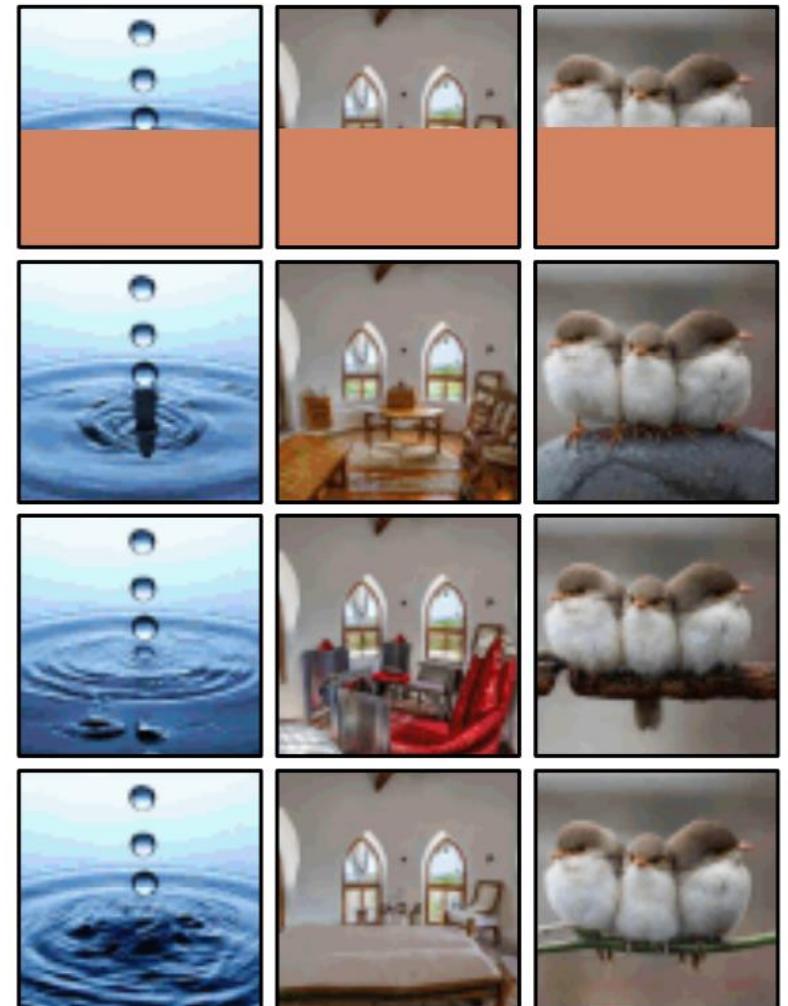


Image completion



Transformers for images

Vision Transformer (ViT) is a Transformer encoder:

- Divides the image into 16x16 patches
- Each patch is mapped to an input embedding (learnable)
- Resulting embeddings are introduced to the Transformer network
- Standard 1D positional encodings are learned

Transformers for images

Vision Transformer (ViT) is a Transformer encoder:

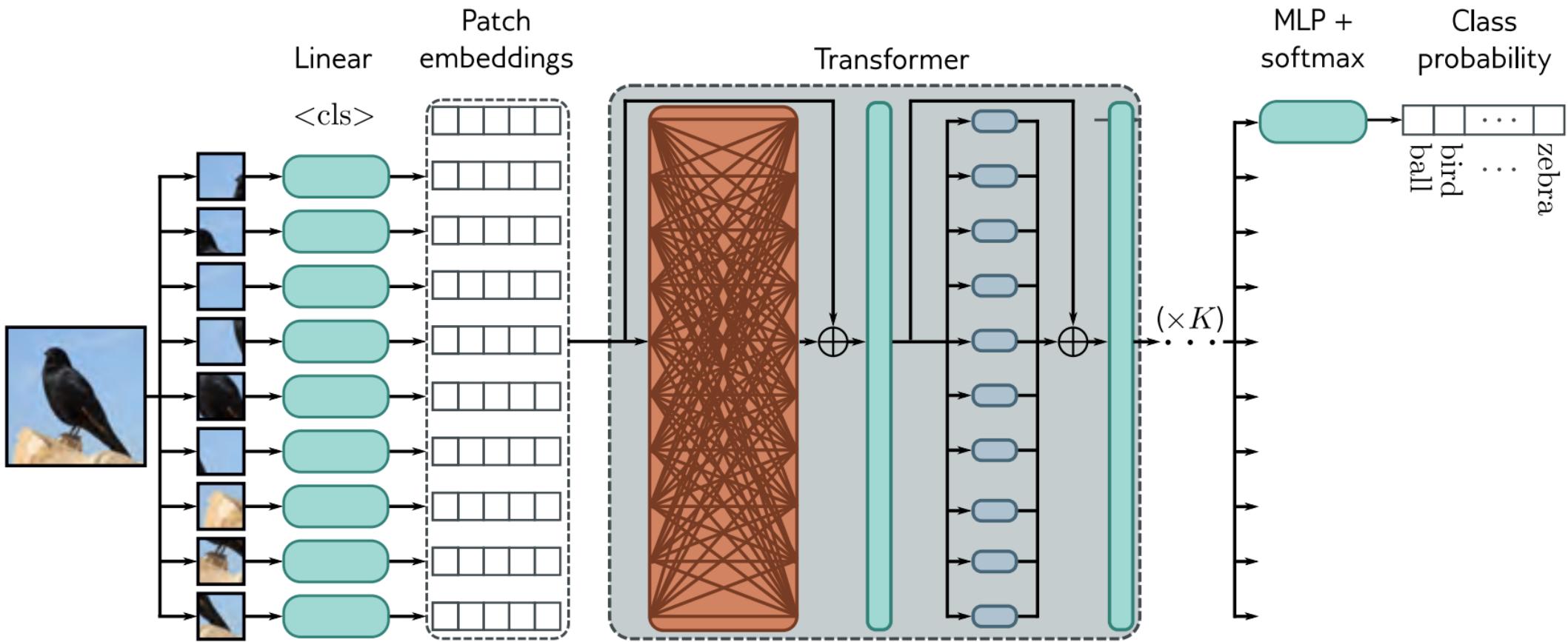
- Divides the image into 16x16 patches
- Each patch is mapped to an input embedding (learnable)
- Resulting embeddings are introduced to the Transformer network
- Standard 1D positional encodings are learned

Used supervised pre-training on a dataset with

- 303 million labeled images
- 18,000 classes

Transformers for images

Vision Transformer (ViT)

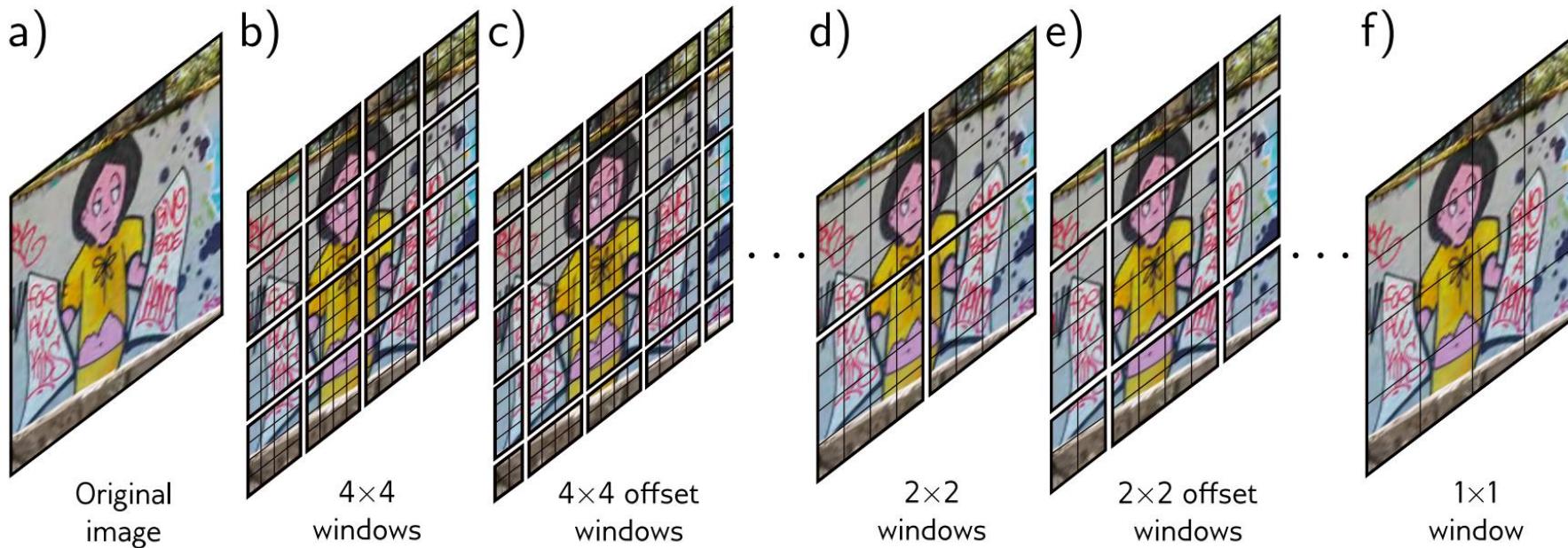


Transformers for images

Multi-scale Vision Transformer:

- Start with small high-resolution patches
- Gradually enlarge the receptive field, decrease spatial resolution

Example: Shifted-window Transformer (Swin Transformer)



Transformers for videos

Video Vision Transformer (ViViT)

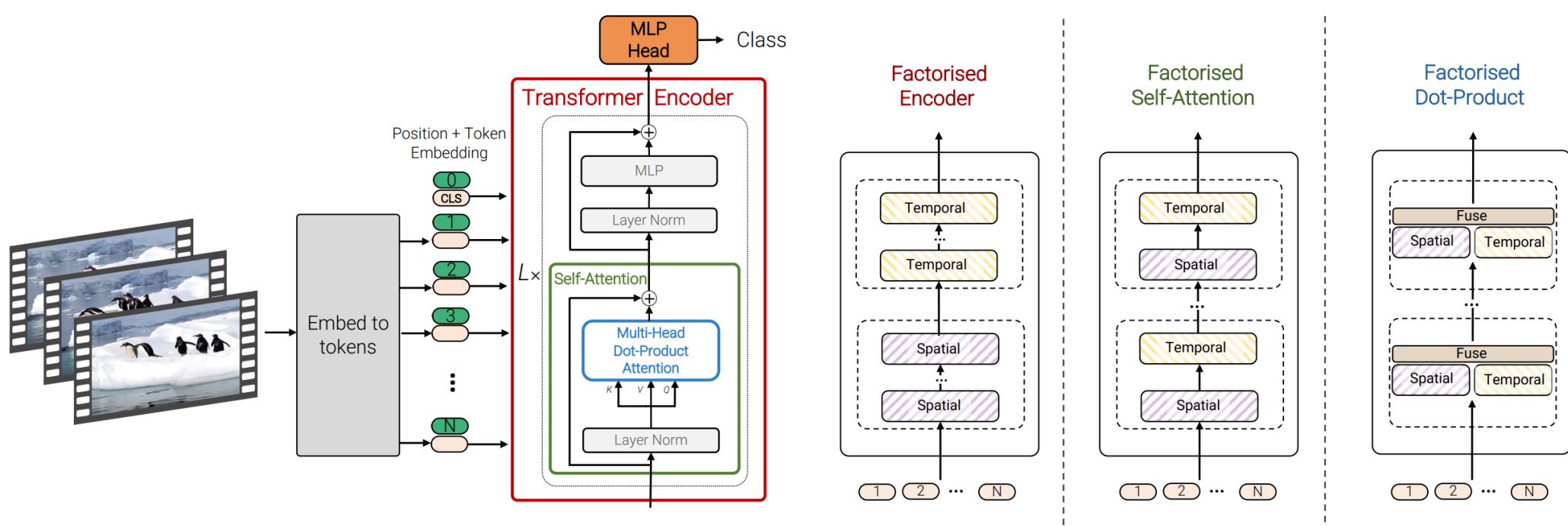


Figure 1: We propose a pure-transformer architecture for video classification, inspired by the recent success of such models for images [18]. To effectively process a large number of spatio-temporal tokens, we develop several model variants which factorise different components of the transformer encoder over the spatial- and temporal-dimensions. As shown on the right, these factorisations correspond to different attention patterns over space and time.

Transformers for videos

Video Vision Transformer (ViViT)

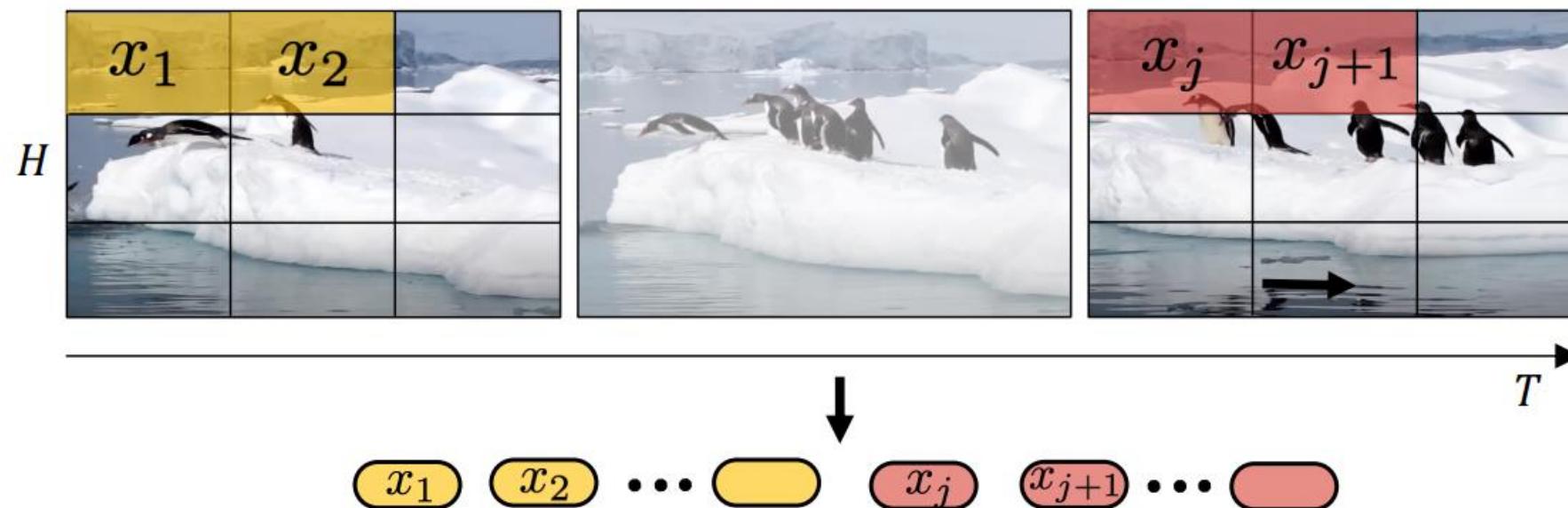


Figure 2: Uniform frame sampling: We simply sample n_t frames, and embed each 2D frame independently following ViT [17].

Transformers for videos

Video Vision Transformer (ViViT)

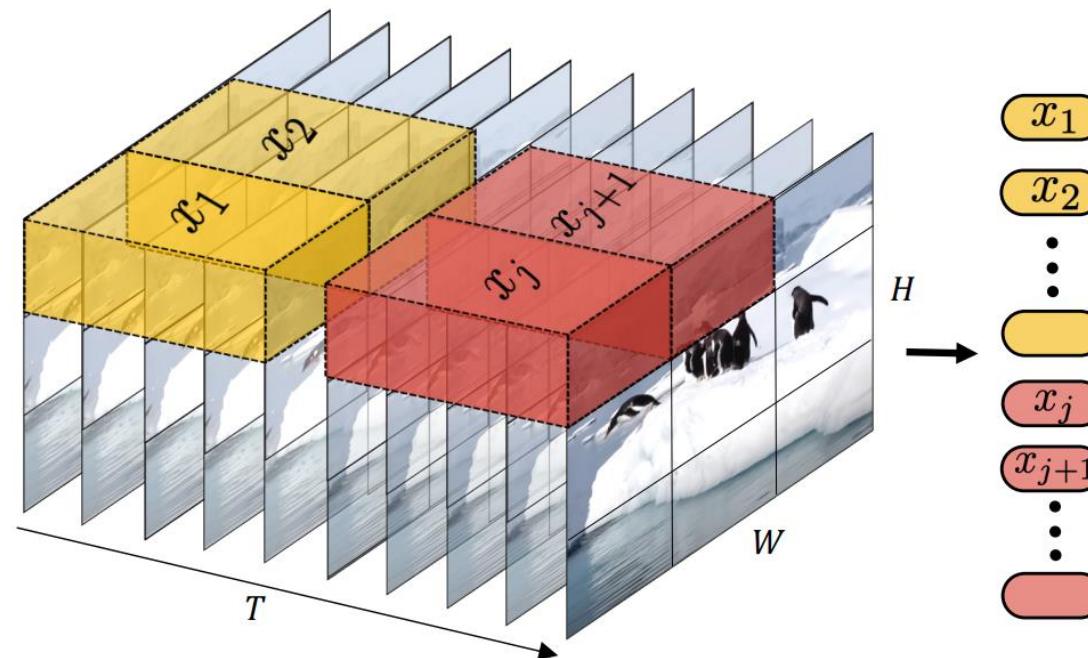


Figure 3: Tubelet embedding. We extract and linearly embed non-overlapping tubelets that span the spatio-temporal input volume.

Transformers for videos

Video Vision Transformer (ViViT)

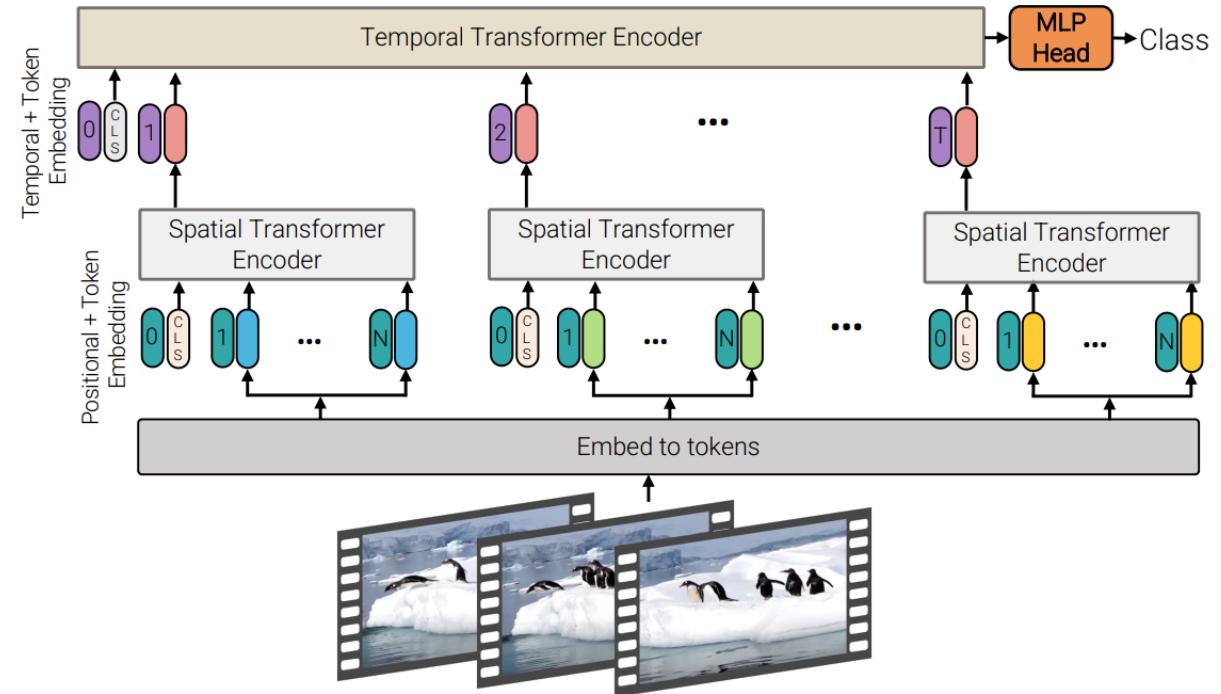


Figure 4: Factorised encoder (Model 2). This model consists of two transformer encoders in series: the first models interactions between tokens extracted from the same temporal index to produce a latent representation per time-index. The second transformer models interactions between time steps. It thus corresponds to a “late fusion” of spatial- and temporal information.

Transformers for videos

Video Vision Transformer (ViViT)

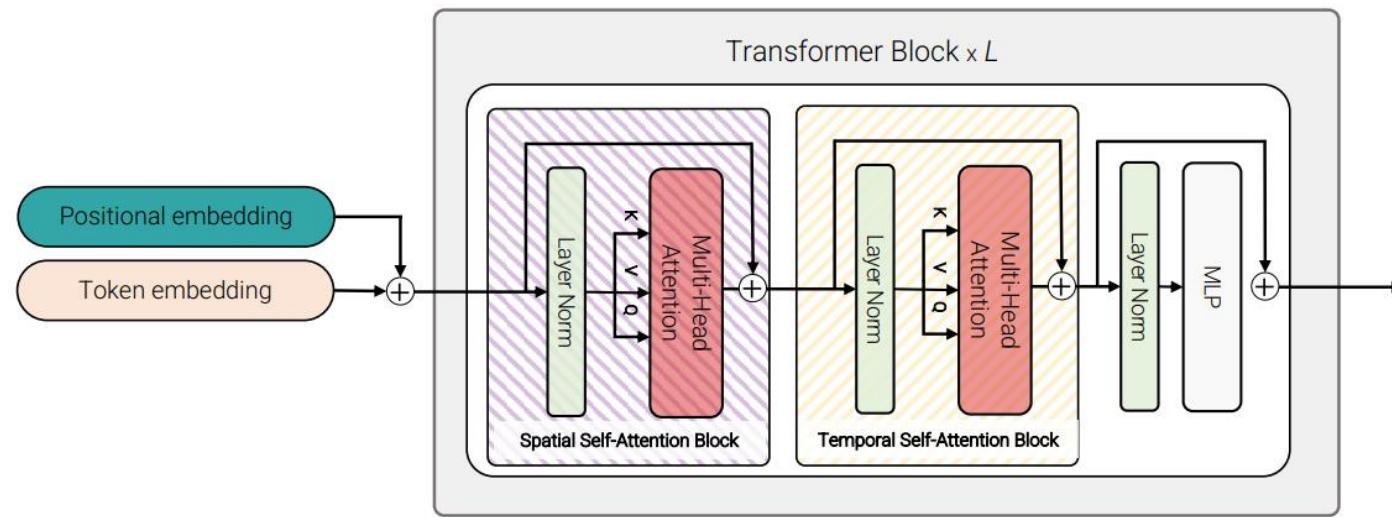


Figure 5: Factorised self-attention (Model 3). Within each transformer block, the multi-headed self-attention operation is factorised into two operations (indicated by striped boxes) that first only compute self-attention spatially, and then temporally.

Continual Inference Networks

A deep neural network, which:

- is capable of continual step inference without computational redundancy,
- is capable of batch inference corresponding to a non-continual neural network,
- produces identical outputs for batch- and step inference given identical receptive fields,
- uses one set of trainable parameters for both batch and step inference

Continual Transformer Encoder

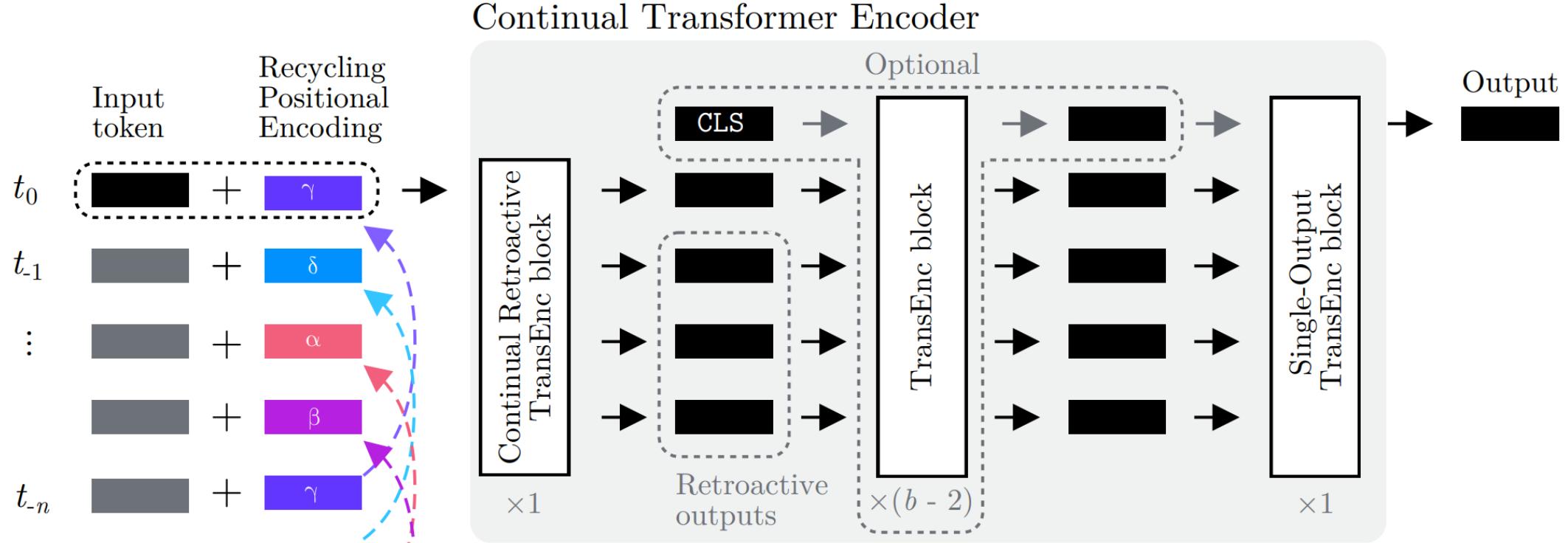


Figure 1: **Multi-block Continual Transformer Encoder with Recycling Positional Encoding.** For $b > 2$ blocks, regular Transformer Encoder blocks can be added between an initial Continual Retroactive block and a final Single-Output block. A class-token may be used after the initial block.

Continual Transformer Encoder

Regular scaled dot-product attention:

- Given query, key, and value matrices $Q, K, V \in \mathbb{R}^{n \times d}$, it can be defined as:

$$\text{Att}(Q, K, V) = D^{-1} A V \quad A = \exp(QK^\top / \sqrt{d}) \quad D = \text{diag}(A \mathbf{1}_n^\top)$$

where $A, D \in \mathbb{R}^{n \times n}$ and $\mathbf{1}_{\textcolor{brown}{n}}$ is a row vector of ones.

Continual Transformer Encoder

Continual retroactive scaled dot-product attention:

We can compute $\mathbf{D}^{-1} \mathbf{A} \mathbf{V}$ in a step-wise manner using the latest query, key, and value steps, $\mathbf{q}_{\text{new}}, \mathbf{k}_{\text{new}}, \mathbf{v}_{\text{new}} \in \mathbb{R}^{1 \times d}$, alongside appropriately cached partial results. The softmax normalisation with \mathbf{D}^{-1} can be efficiently implemented via column-aligned element-wise multiplications (denoted by \odot hereafter) of a column-vector $\mathbf{d} = \mathbf{A} \mathbb{1}_n^\top$. If we cache the $n - 1$ values for the prior step tokens, i.e. $\mathbf{d}_{\text{mem}} = \mathbf{A}_{\text{prev}}^{(-n+1:-1)} \mathbb{1}_{n-1}^\top$, alongside \mathbf{Q} and \mathbf{K} , we can define the step update as:

$$\begin{aligned}\mathbf{d}^{(-n+1:-1)} &= \mathbf{d}_{\text{mem}}^{(-n+2:0)} - \exp(\mathbf{Q}_{\text{mem}} \mathbf{k}_{\text{old}}^\top) + \exp(\mathbf{Q}_{\text{mem}} \mathbf{k}_{\text{new}}^\top) \\ \mathbf{d}^{(0)} &= \exp\left(\frac{\mathbf{q}_{\text{new}}}{\sqrt{d}} (\mathbf{K}_{\text{mem}} \parallel \mathbf{k}_{\text{new}})^\top\right) \mathbb{1}_n^\top,\end{aligned}$$

where \mathbf{Q}_{mem} (\mathbf{K}_{mem}) are the $n - 1$ prior query (key) tokens, \mathbf{k}_{old} is the key from n steps ago, and \parallel denotes concatenation of matrices along the first dimension. Negative indices indicate prior time-steps.

Continual Transformer Encoder

Continual retroactive scaled dot-product attention:

An update for \mathbf{AV} can likewise be defined as a function of the $n - 1$ prior values \mathbf{AV}_{mem} :

$$\begin{aligned}\mathbf{AV}^{(-n+1:-1)} &= \mathbf{AV}_{\text{mem}}^{(-n+2:0)} - \exp(\mathbf{Q}_{\text{mem}} \mathbf{k}_{\text{old}}^\top) \mathbf{v}_{\text{old}} + \exp(\mathbf{Q}_{\text{mem}} \mathbf{k}_{\text{new}}^\top) \mathbf{v}_{\text{new}} \\ \mathbf{AV}^{(0)} &= \exp\left(\frac{\mathbf{q}_{\text{new}}}{\sqrt{d}} (\mathbf{K}_{\text{mem}} \parallel \mathbf{k}_{\text{new}})^\top\right) (\mathbf{V}_{\text{mem}} \parallel \mathbf{v}_{\text{new}}).\end{aligned}$$

Finally, we compute the Continual Retroactive Attention output in the usual manner:

$$\text{CoReAtt}(\mathbf{q}_{\text{new}}, \mathbf{k}_{\text{new}}, \mathbf{v}_{\text{new}}) = \mathbf{d}^{-1} \odot \mathbf{AV}.$$

Continual Transformer Encoder

Continual single-output scaled dot-product attention:

Both the Regular and Continual Retroactive Dot-Product Attentions produce attention outputs for the current step, as well as $n - 1$ retroactively updated steps. In cases where retroactive updates are not needed, we can simplify the computation greatly via a Continual Single-Output Dot-Product Attention (*CoSiAtt*). In essence, the regular SDA is reused, but prior values of \mathbf{k} and \mathbf{v} are cached between steps (as in (Ott et al., 2019)), and only the attention corresponding to a single query token \mathbf{q} is computed:

$$CoSiAtt(\mathbf{q}, \mathbf{k}_{\text{new}}, \mathbf{v}_{\text{new}}) = \mathbf{a} (\mathbf{V}_{\text{mem}} \parallel \mathbf{v}_{\text{new}}) / \mathbf{a} \mathbb{1}_n^\top, \quad \mathbf{a} = \exp \left(\frac{\mathbf{q}}{\sqrt{d}} (\mathbf{K}_{\text{mem}} \parallel \mathbf{k}_{\text{new}})^\top \right).$$

Continual Transformer Encoder

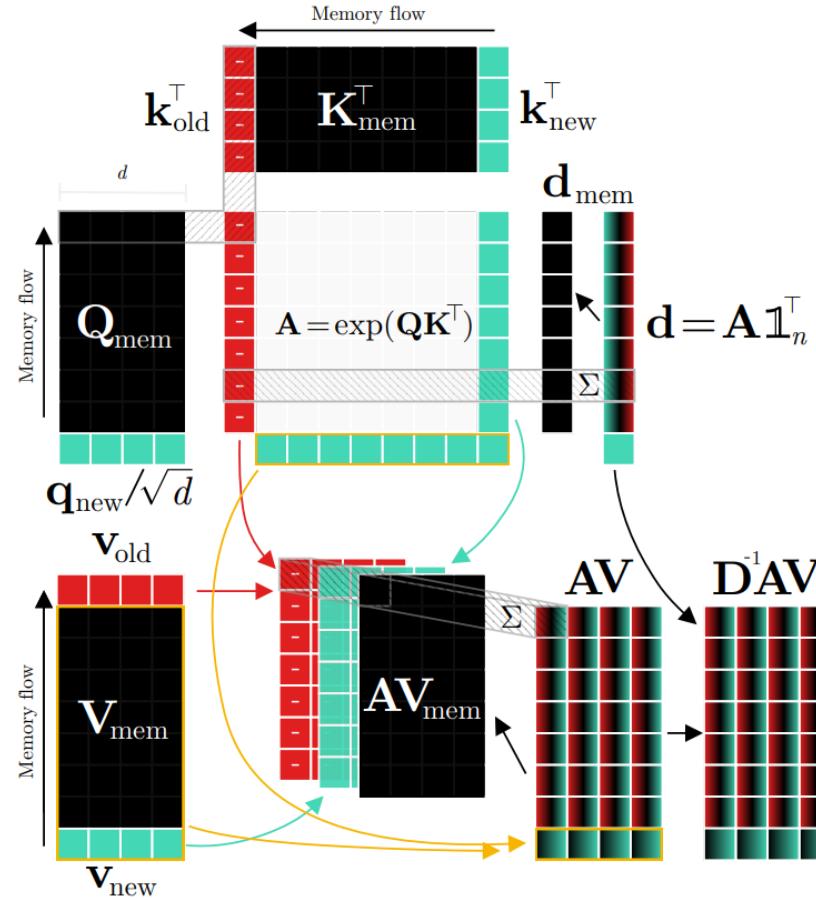


Figure 4: **Continual Retroactive Dot-Product Attention.** The query (Q), key (K), and value (V) matrices are aggregated over time by caching the step vectors \mathbf{q}_{new} , \mathbf{k}_{new} , and \mathbf{v}_{new} in each their FIFO queue (denoted by \square_{mem}). During each step, only the entries of A associated with \mathbf{q}_{new} , \mathbf{k}_{new} and the oldest K step, \mathbf{k}_{old} are computed. The diagonal entries of the row-normalisation matrix D as well as the AV can be updated retroactively by subtracting features corresponding to \mathbf{k}_{old} and adding features related to \mathbf{k}_{new} to the cached outputs of the previous step, D_{mem} and AV_{mem} , respectively.