



NUS | **Computing**
National University
of Singapore

IT5005 Artificial Intelligence

Sirigina Rajendra Prasad
AY2025/2026: Semester 1

Transformers

Self-supervised Learning

- Model is trained on unlabeled data by having the model generate its labels from data itself
- **NLP Example:** Next-token prediction
 - Unlabeled data: *“Singapore is known for its efficient public transport system”*

“Singapore” ---> Target token: *“is”*

“Singapore is” ---> Target token: *“known”*

“Singapore is known” ---> Target token: *“for”*

“Singapore is known for” ---> Target token: *“its”*

“Singapore is known for its” ---> Target token: *“efficient”*

“Singapore is known for its efficient” ---> Target token: *“public”*

“Singapore is known for its efficient public” ---> Target token: *“transport”*

“Singapore is known for its efficient public transport” ---> Target token: *“system”*

Self-supervised Learning

- Model is trained on unlabeled data by having the model generate its labels from data itself
- **NLP Example:** Masked Language Model
 - Unlabeled data: *“Singapore is known for its efficient public transport system”*

“Singapore is _____ for its efficient public transport system”

Target token: *“**known**”*

“Singapore is known for its _____ public transport system”

Target token: *“**efficient**”*

Transformers: Motivation

- Language modeling and masked language modelling
 - Can be modelled using Markov chains
 - Transition probabilities capture the relation between words
- **Issue: Scaling**
 - Context window size is equivalent to order of Markov chains
 - Memory and computational complexity is exponential in context window size

Other alternatives for sequential data

- Convolutional Neural Networks (1-D kernels)
- Recurrent Neural Networks
- Long Short Term Memory (LSTM) and its variants (SOTA till 2017)

CNNs for Sequential Data

- 1-D kernels to capture temporal relations



CNNs can model the sequential data through 1-D kernels

Only local context and position-agnostic

Recurrent Neural Networks

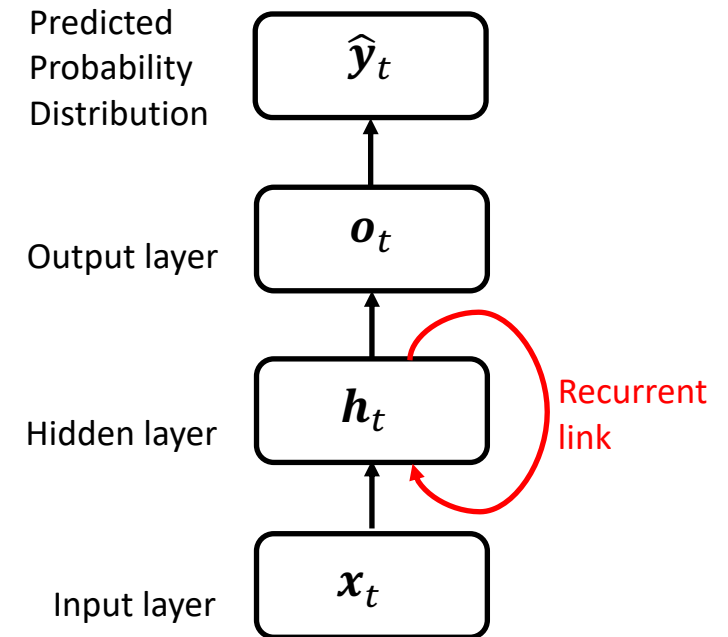
- Accept dynamic or sequence data of unequal length

x_t : Input at time t

h_t : Hidden layer at time t

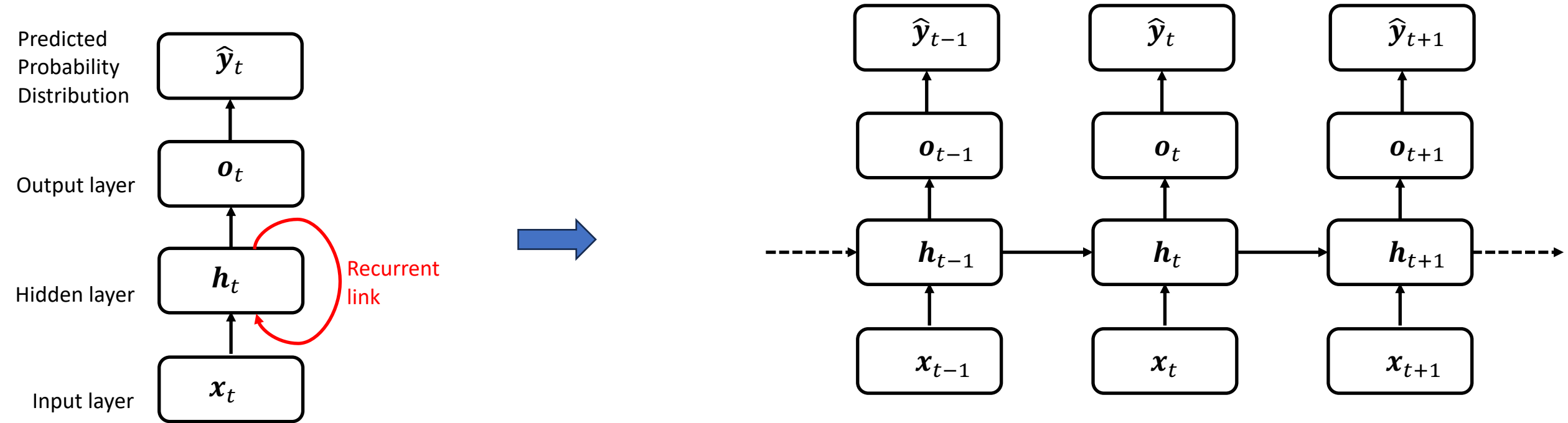
o_t : Output vector at time t

\hat{y}_t : Output of SoftMax layer at time t



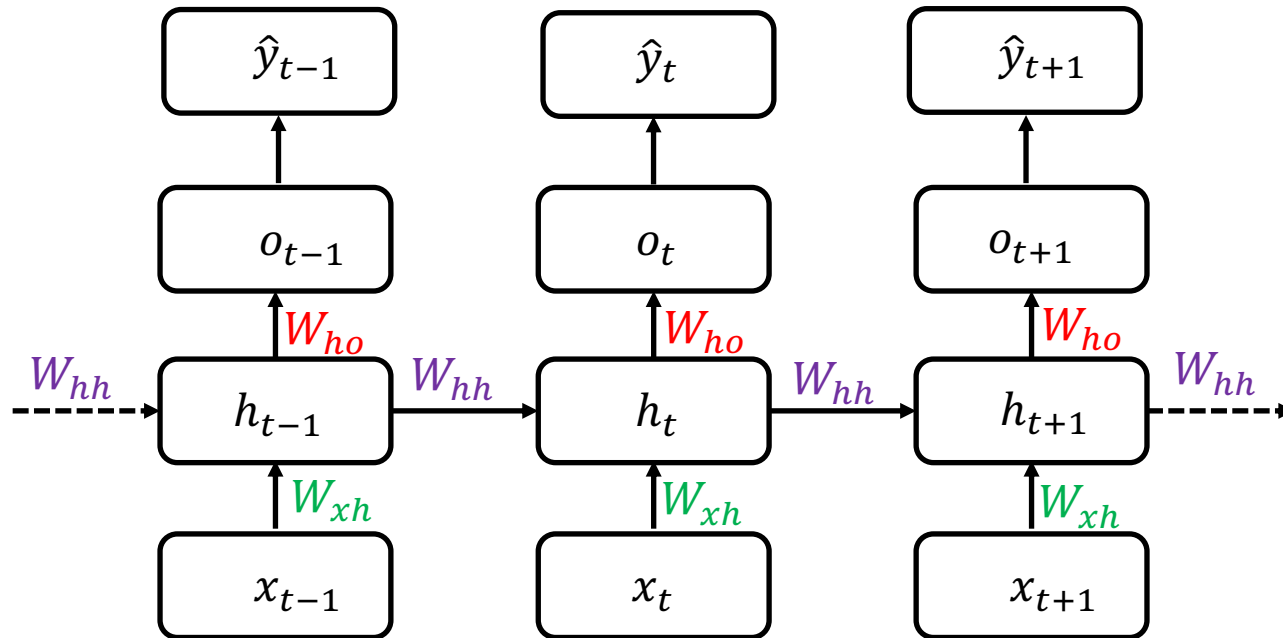
- Recurrent link
 - Output is feedback as input
- Recurrent link can be unrolled to handle variable length data

RNN: Unrolled



- Input sequences can be of different lengths
- Hidden layers from previous time instant provides context or memory

RNN: Unrolled



Reuse the same
weight matrix W_{xx}
at each time step

$$\mathbf{h}_t = \mathbf{f}(\mathbf{a}_t)$$

$$\text{where } \mathbf{a}_t = \mathbf{W}_{xh}\mathbf{x}_t + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{b}$$

$$\mathbf{W}_{xh} \in \mathbb{R}^{k \times d}$$

$$\mathbf{W}_{hh} \in \mathbb{R}^{k \times k}$$

$$\mathbf{b} \in \mathbb{R}^k$$

\mathbf{f} is Tanh or ReLU

$$\hat{\mathbf{y}}_t = g(\mathbf{o}_t)$$

$$\text{where } \mathbf{o}_t = \mathbf{W}_{ho}\mathbf{h}_t + \mathbf{c}$$

$$\mathbf{W}_{ho} \in \mathbb{R}^{|V| \times k}$$

$$\mathbf{c} \in \mathbb{R}^{|V|}$$

g is softmax

RNNs for Language Modelling

- RNNs model the joint probability $P(x_1, x_2, \dots, x_n)$, where x_1, x_2, \dots, x_n represent sequence of words.

$$P(x_1, x_2, \dots, x_n) = \prod_{t=1}^n P(x_t | x_{t-1}, x_{t-2}, \dots, x_1)$$

- Useful for many applications involving text/sentence generation
 - Prediction of next token, machine translation, speech recognition, question answering , etc.

Prediction of Next Token

During Training

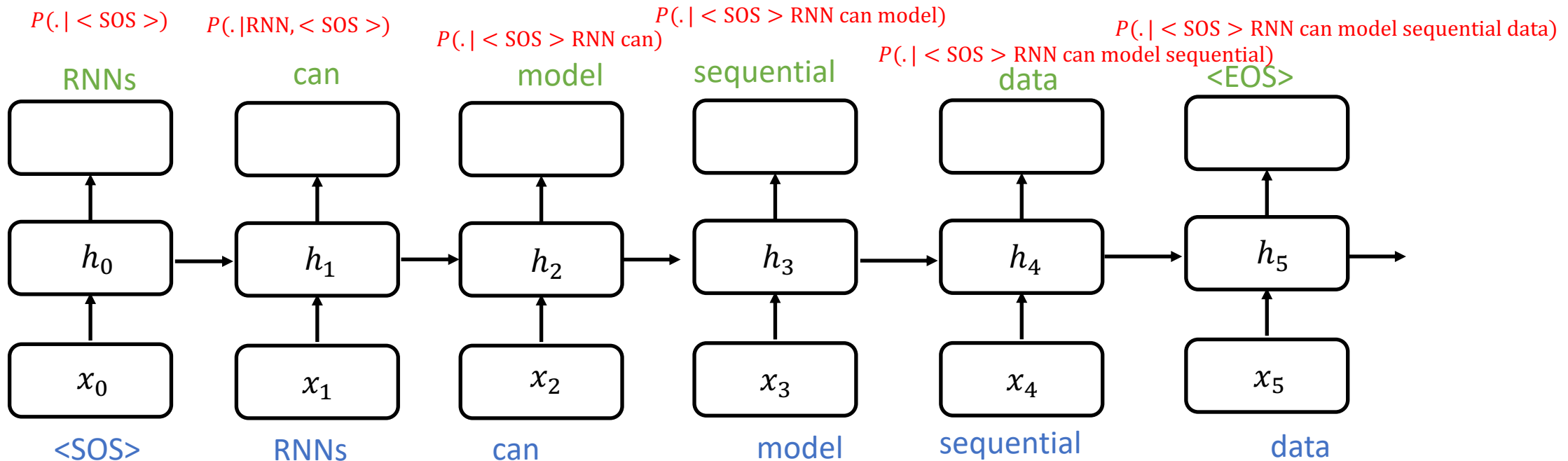
“RNNs can model the sequential data”

Input	Ground Truth	Output Distribution
$x_0 = \text{“<BOS>”}$	y_0 : <i>“RNNs”</i>	$P(. x_0)$
$x_{0:1} = \text{“<BOS> RNNs”}$	y_1 : <i>“can”</i>	$P(. x_{0:1})$
$x_{0:2} = \text{“<BOS> RNN can”}$	y_2 : <i>“model”</i>	$P(. x_{0:2})$
$x_{0:3} = \text{“<BOS> RNN can model”}$	y_3 : <i>“sequential”</i>	$P(. x_{0:3})$
$x_{0:4} = \text{“<BOS> RNN can model sequential”}$	y_4 : <i>“data”</i>	$P(. x_{0:4})$
$x_{0:5} = \text{“<BOS> RNN can model sequential data”}$	y_5 : <i>“<EOS>”</i>	$P(. x_{0:5})$

Prediction of Next Token

During Training

“RNNs can model the sequential data”



Input
Ground truth
Output distribution

Prediction of Next Token

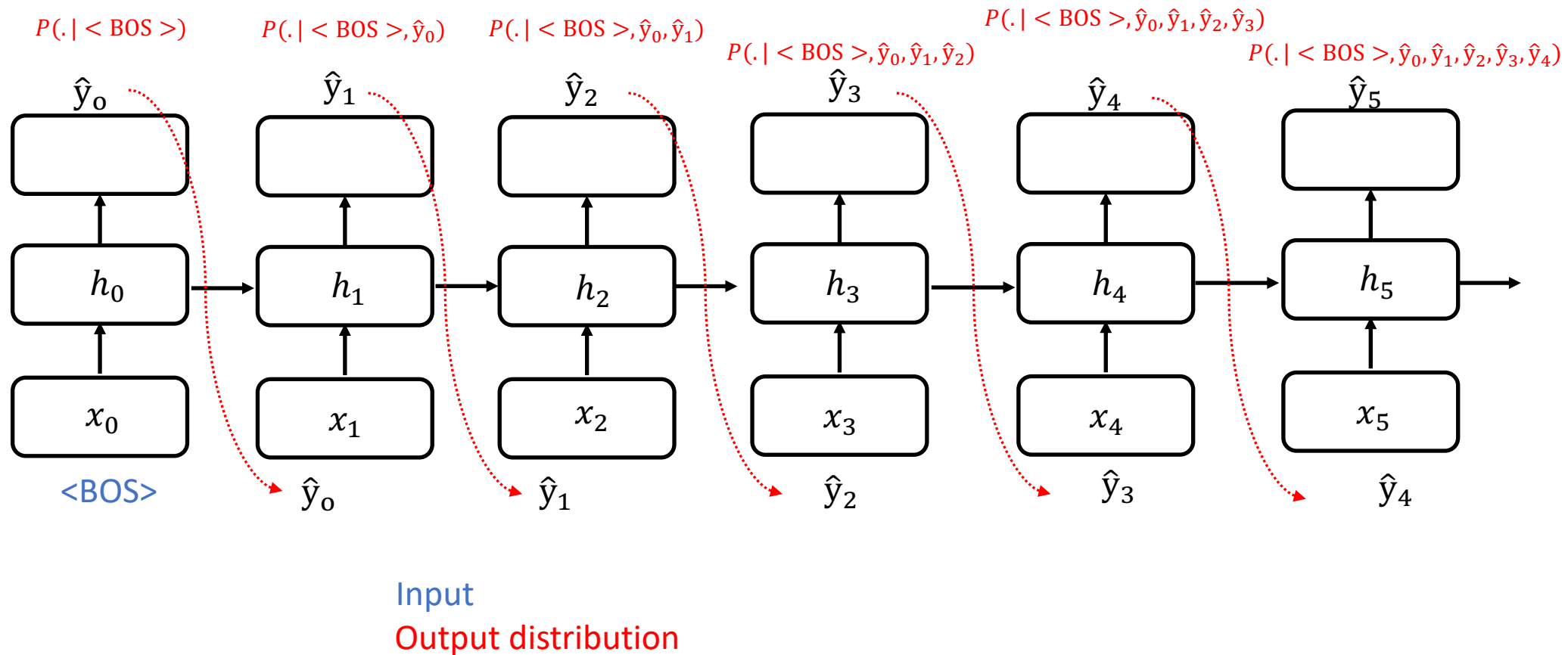
During Inference

“RNNs can model the sequential data”

Input	Predicted	Output Distribution
$x_0 = \text{<BOS>}$	\hat{y}_0	$P(. x_0)$
$x_{0:1} = \text{<BOS> } \hat{y}_0$	\hat{y}_1	$P(. x_{0:1})$
$x_{0:2} = \text{<BOS> } \hat{y}_0, \hat{y}_1$	\hat{y}_2	$P(. x_{0:2})$
$x_{0:3} = \text{<BOS> } \hat{y}_0, \hat{y}_1, \hat{y}_2$	\hat{y}_3	$P(. x_{0:3})$
$x_{0:4} = \text{<BOS> } \hat{y}_0, \hat{y}_1, \hat{y}_2, \hat{y}_3$	\hat{y}_4	$P(. x_{0:4})$
$x_{0:5} = \text{<BOS> } \hat{y}_0, \hat{y}_1, \hat{y}_2, \hat{y}_3, \hat{y}_4$	\hat{y}_5	$P(. x_{0:5})$

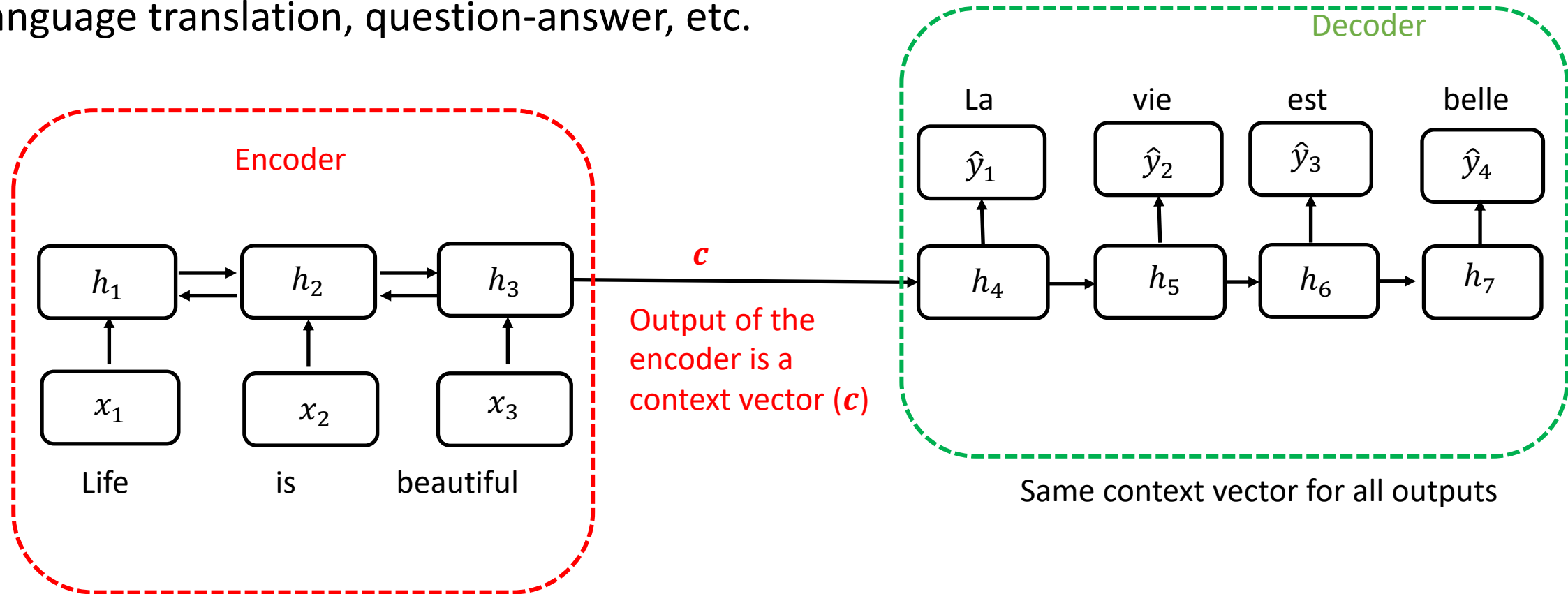
Prediction of Next Token

During Inference



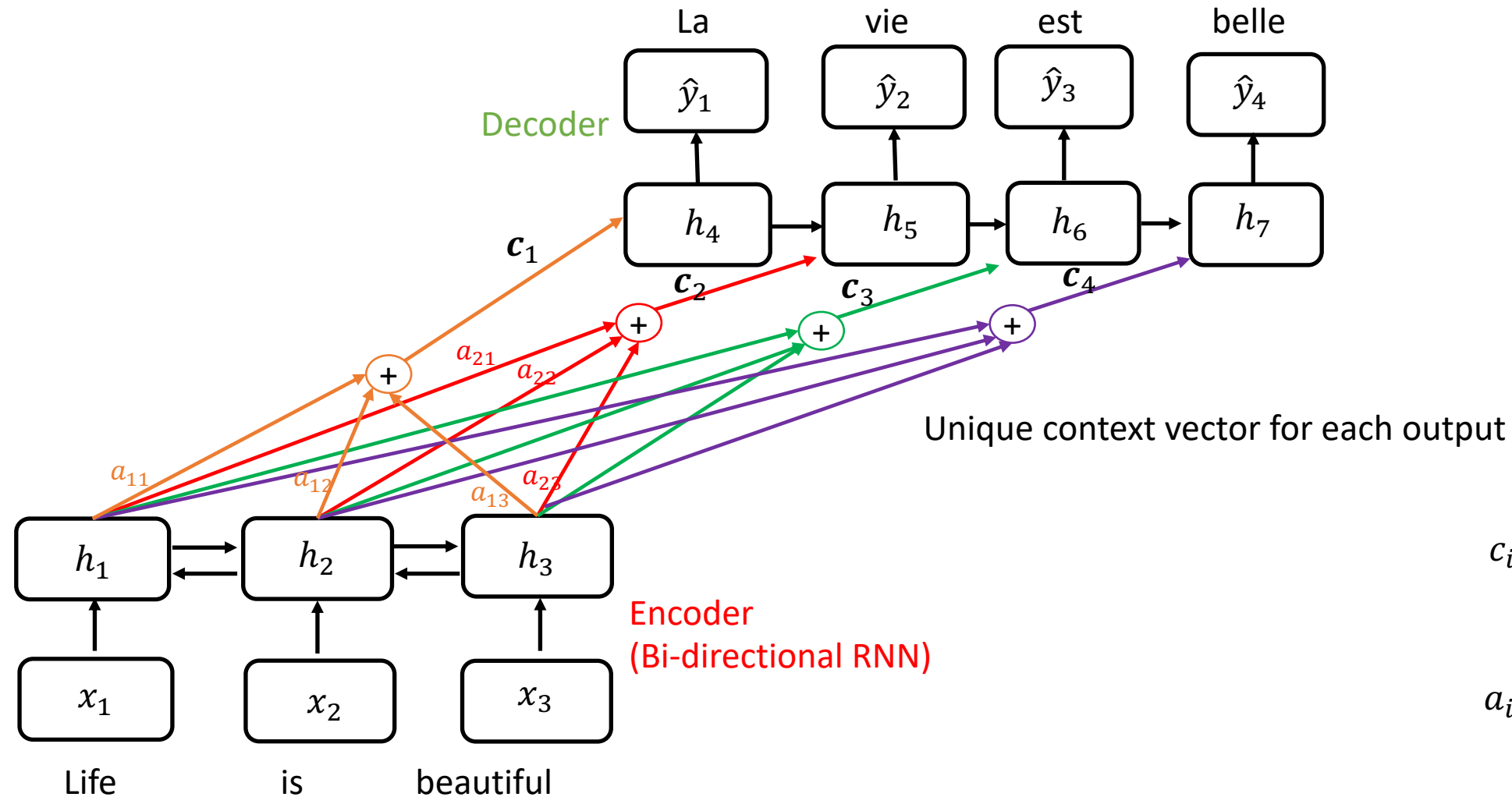
Sequence-to-Sequence Model or Encoder-Decoder Model

- RNNs process the entire input *sequentially* and then produce the output
- **Applications:**
 - language translation, question-answer, etc.



Context vector (c) is called "bottleneck"

Variants of RNNs: RNNs with Attention



$$c_i = \sum_{j=1}^T a_{ij} h_j$$

$$a_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^T \exp(e_{ik})}$$

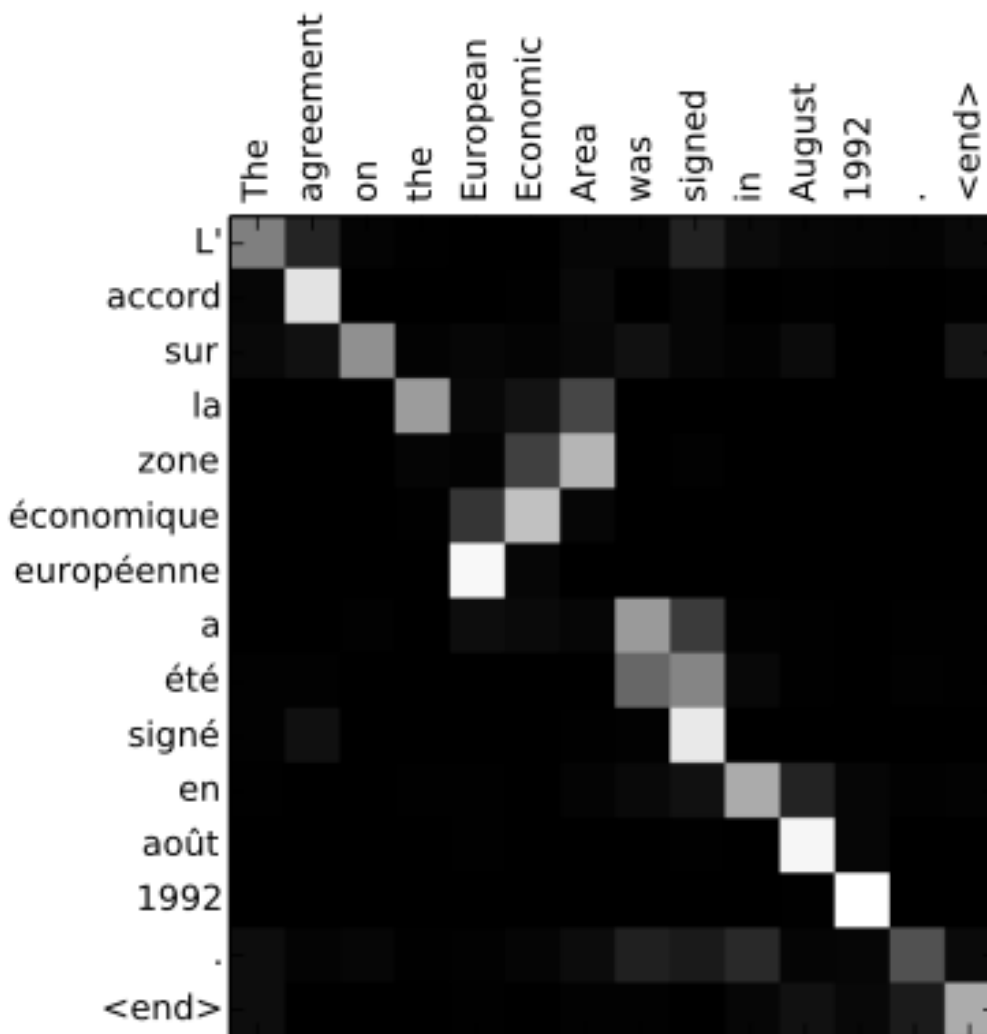
$$e_{ij} = a(y_i, h_j)$$

Neural Machine Translation by Jointly Learning to Align and Translate

Code: https://docs.pytorch.org/tutorials/intermediate/seq2seq_translation_tutorial.html

Variants of RNN: RNN with Attention

Only for your reading



$$\text{Attention Matrix} = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \dots & a_{mn} \end{bmatrix}$$

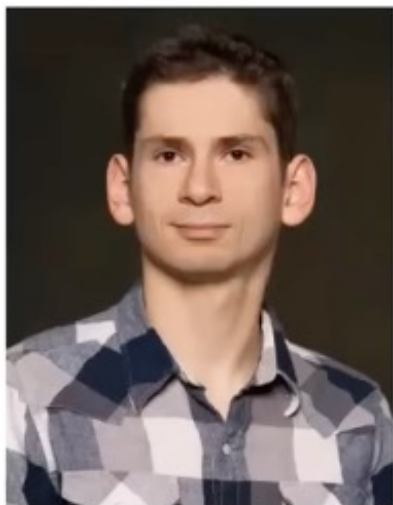
$$c_i = \sum_{j=1}^T a_{ij} h_j$$

$$a_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^T \exp(e_{ik})}$$

$$e_{ij} = a(y_i, h_j)$$

Neural Machine Translation by Jointly Learning to Align and Translate

Code: https://docs.pytorch.org/tutorials/intermediate/seq2seq_translation_tutorial.html



Dzmitry Bahdanau

RE: history of “attention” (from email correspondence)



So I started thinking about how to avoid the bottleneck between encoder and decoder RNN. My first idea was to have a model with two "cursors", one moving through the source sequence (encoded by a BiRNN) and another one moving through the target sequence. The cursor trajectories would be marginalized out using dynamic programming. KyungHyun Cho recognized this as an equivalent to Alex Graves' RNN Transducer model. Following that, I may have also read Graves' hand-writing recognition paper. The approach looked inappropriate for machine translation though.

The above approach with cursors would be too hard to implement in the remaining 5 weeks of my internship. So I tried instead something simpler - two cursors moving at the same time synchronously (effectively hard-coded diagonal attention). That sort of worked, but the approach lacked elegance.

So one day I had this thought that it would be nice to enable the decoder RNN to learn to search where to put the cursor in the source sequence. This was sort of inspired by translation exercises that learning English in my middle school involved. Your gaze shifts back and forth between source and target sequence as you translate. I expressed the soft search as softmax and then weighted averaging of BiRNN states. It worked great from the very first try to my great excitement. I called the architecture RNNSearch, and we rushed to publish an ArXiv paper as we knew that Ilya and co at Google are somewhat ahead of us with their giant 8 GPU LSTM model (RNN Search still ran on 1 GPU).

As it later turned out, the name was not great. The better name (attention) was only added by Yoshua to the conclusion in one of the final passes.

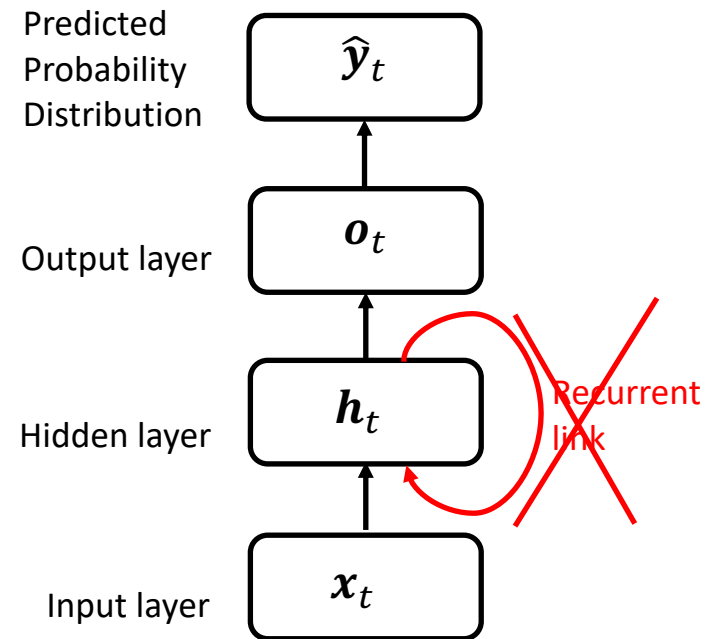


Issues with Sequential Models

- CNNs
 - Smaller context windows
 - Lack of positional information
- RNNs and LSTMs
 - Difficult to achieve parallelism due to sequential processing
 - Vanishing/exploding gradients for larger sequences

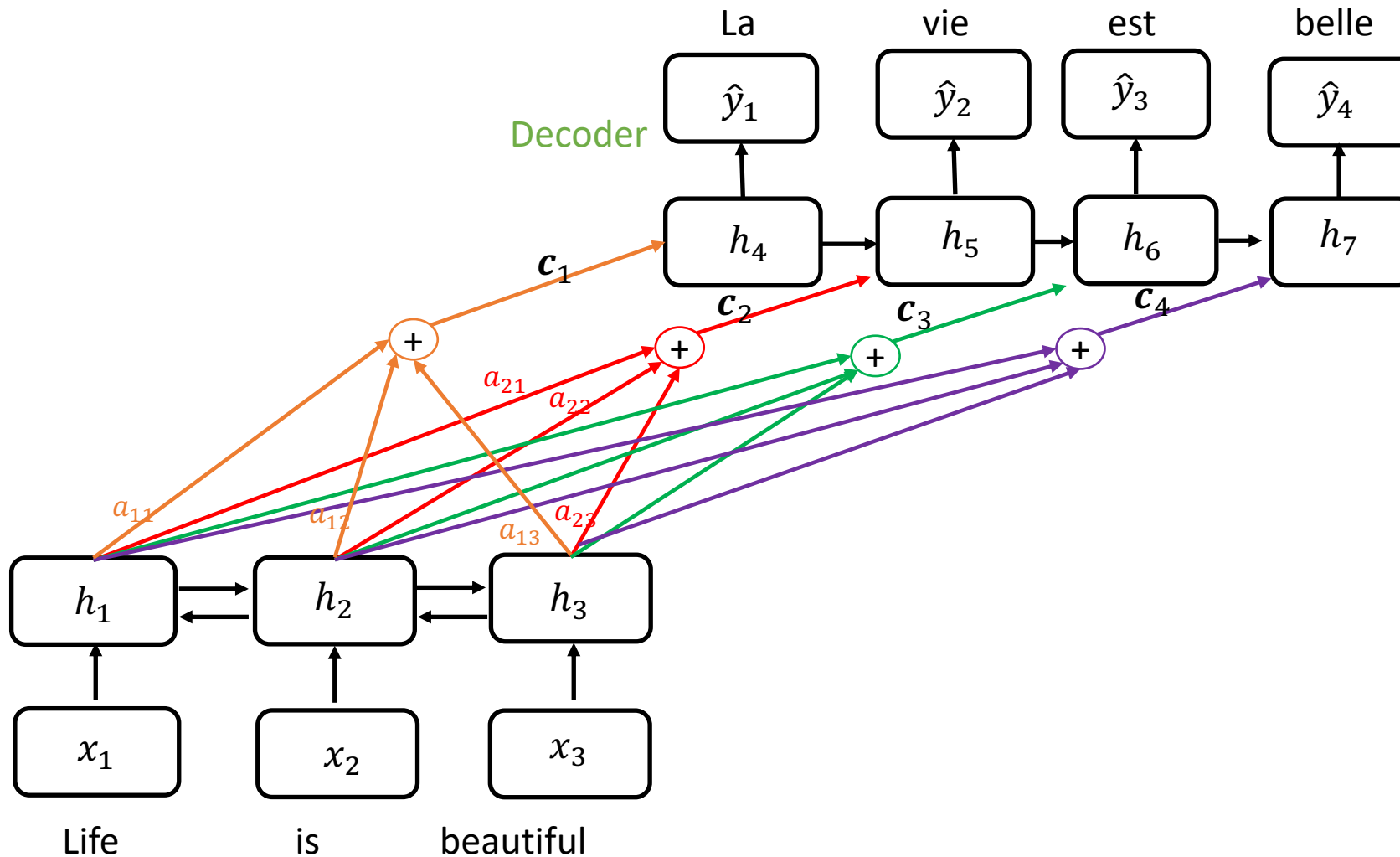
How transformers overcome this issue?

- Recursion Elimination



How transformers overcome this issue?

- Parallel processing all inputs simultaneously



Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Łukasz Kaiser*
Google Brain
lukaszkaizer@google.com

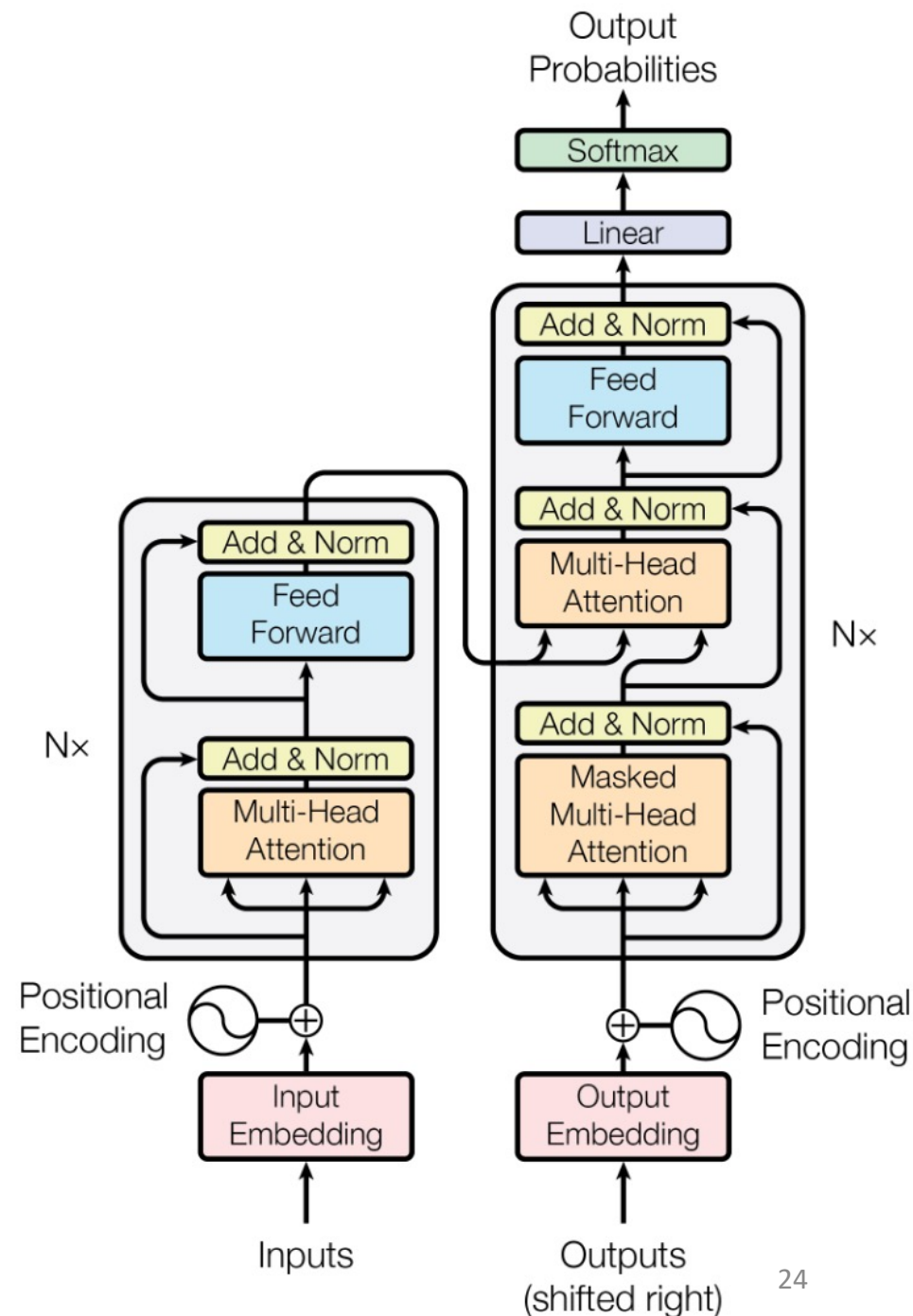
Illia Polosukhin* ‡
illia.polosukhin@gmail.com

Abstract

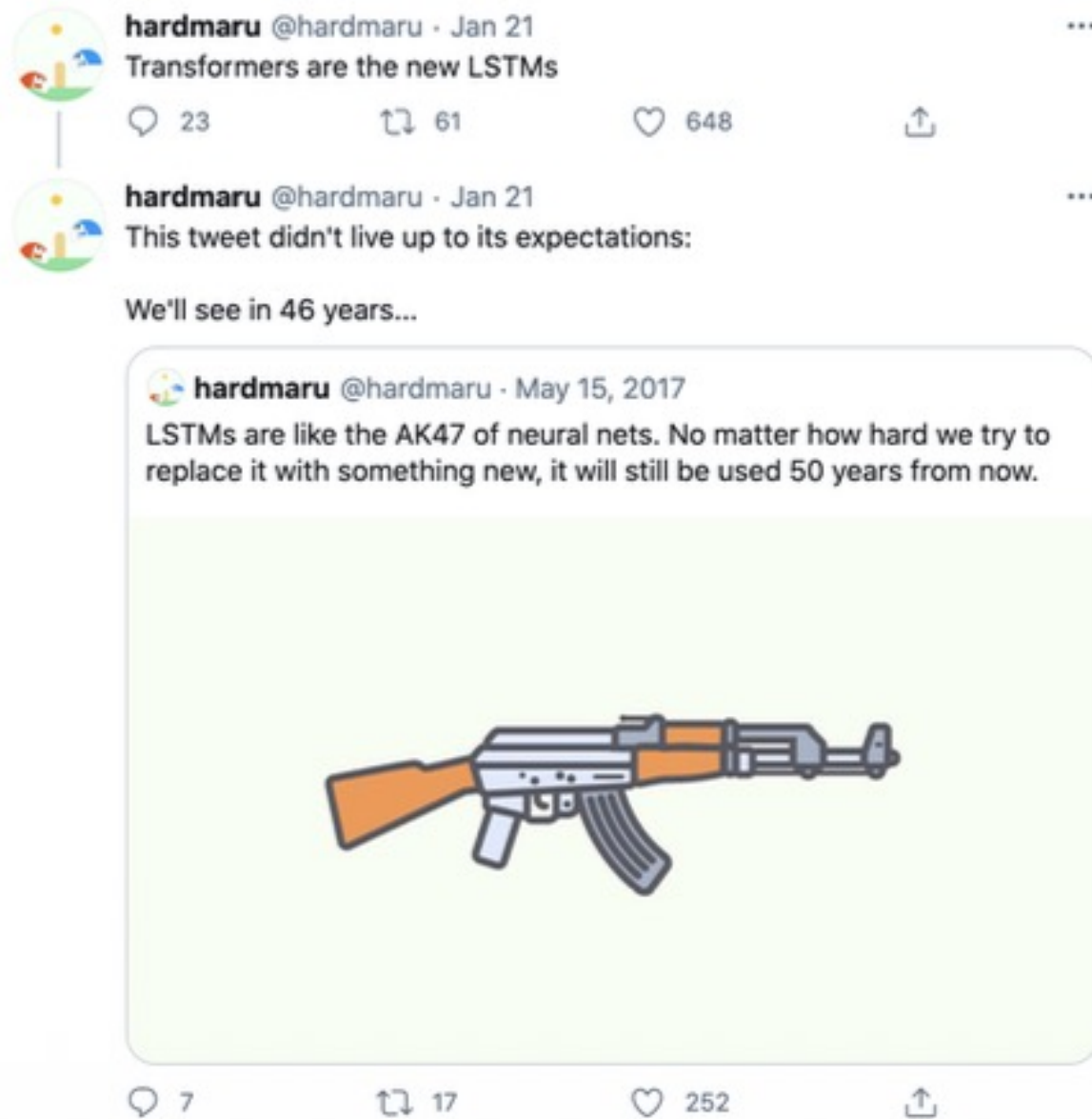
The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

Of course, we also need

1. Skip connections for deeper networks
2. Layer normalization
3. Positional encoding for explicit capture of sequence information



Transformers



Slide credit: Prof. You Yang (CS5242)

Transformers

- Transformers are friendly to modern hardware like TPUs and GPUs
 - We need a high parallelism to make full use of the computing units in TPUs and GPUs (e.g. thousands of threads)



Slide credit: Prof. You Yang (CS5242)

Transformer is transforming everything in deep learning



- All the attention of deep learning is now on Attention!

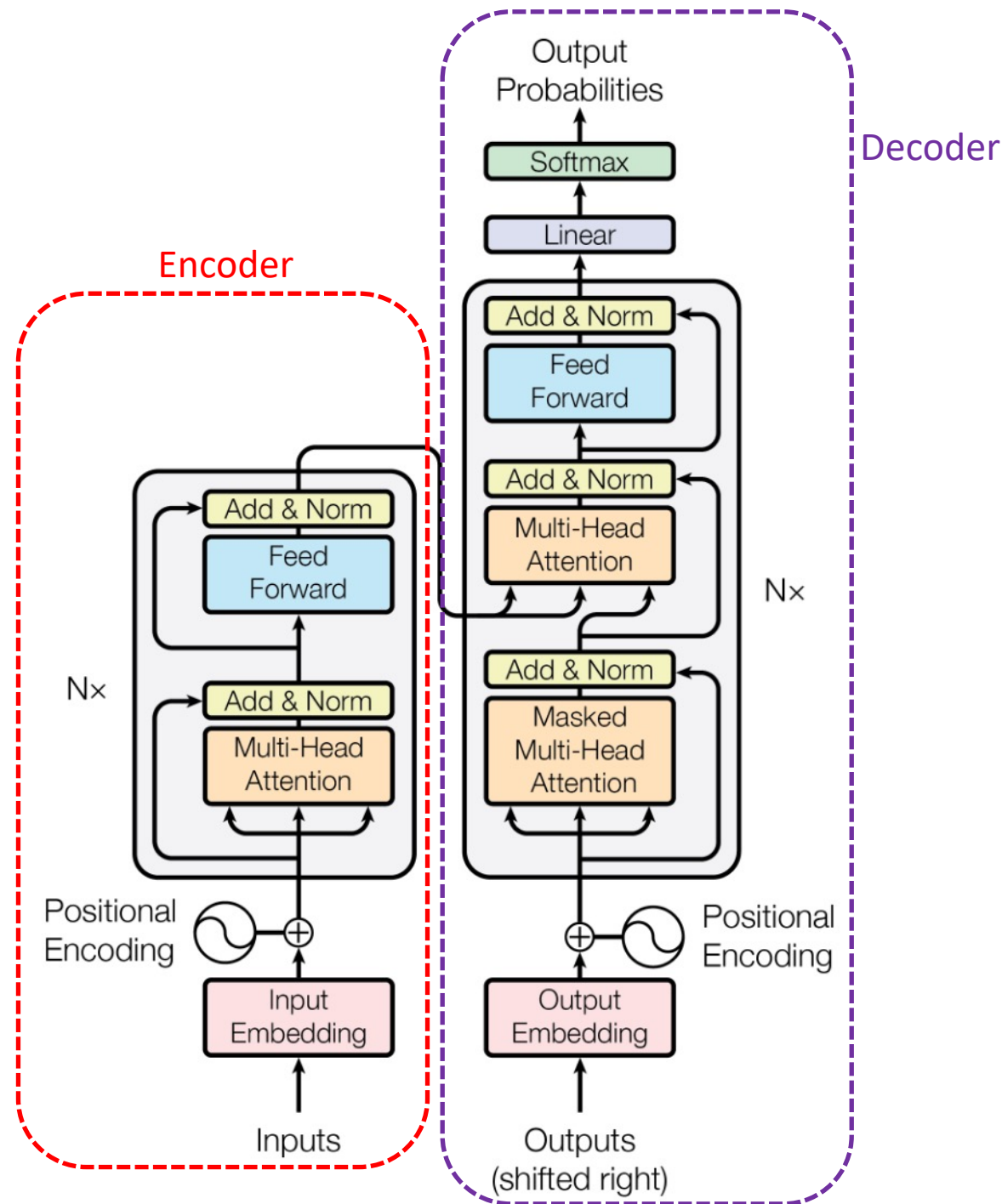
Transformers

- Prior to the transformers, research is fragmented across domains
- Transformer is the unifying force
 - Single architecture for all tasks

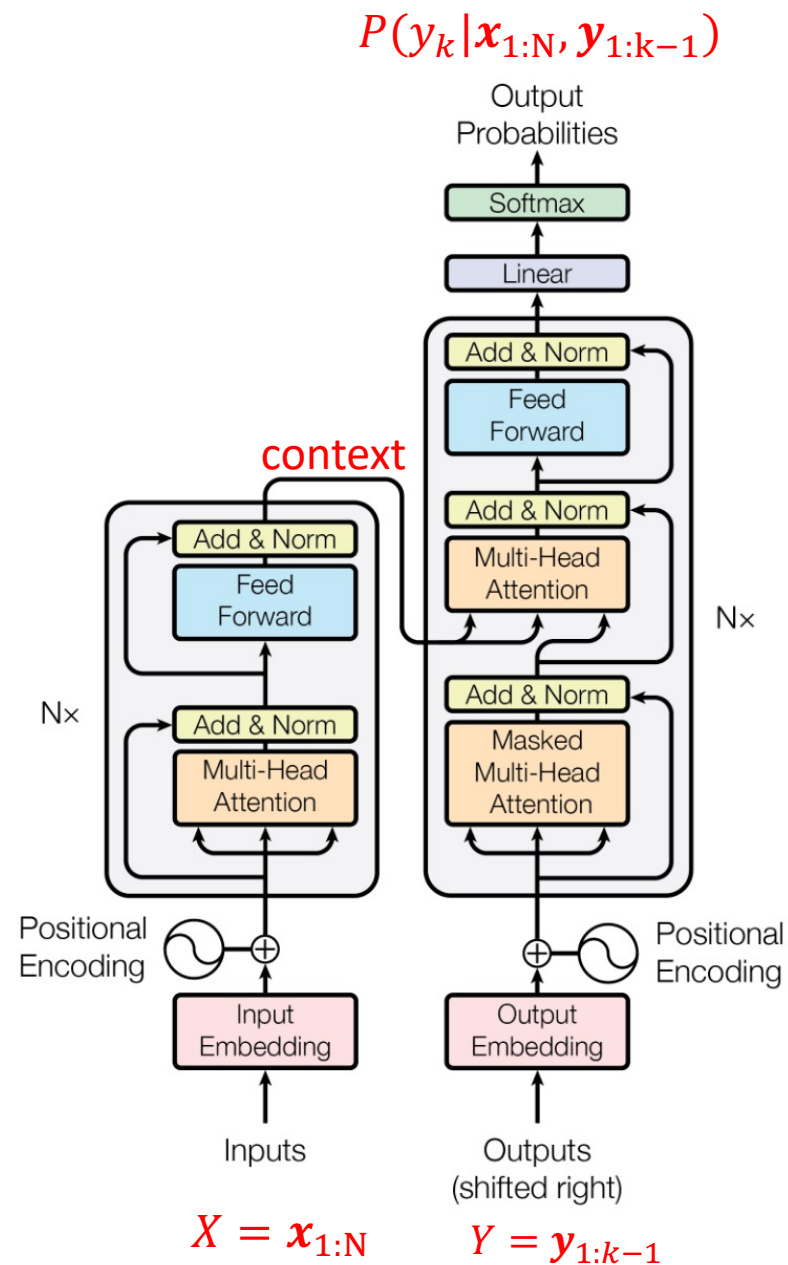


generated using OpenAI's GPT Image 1 model.

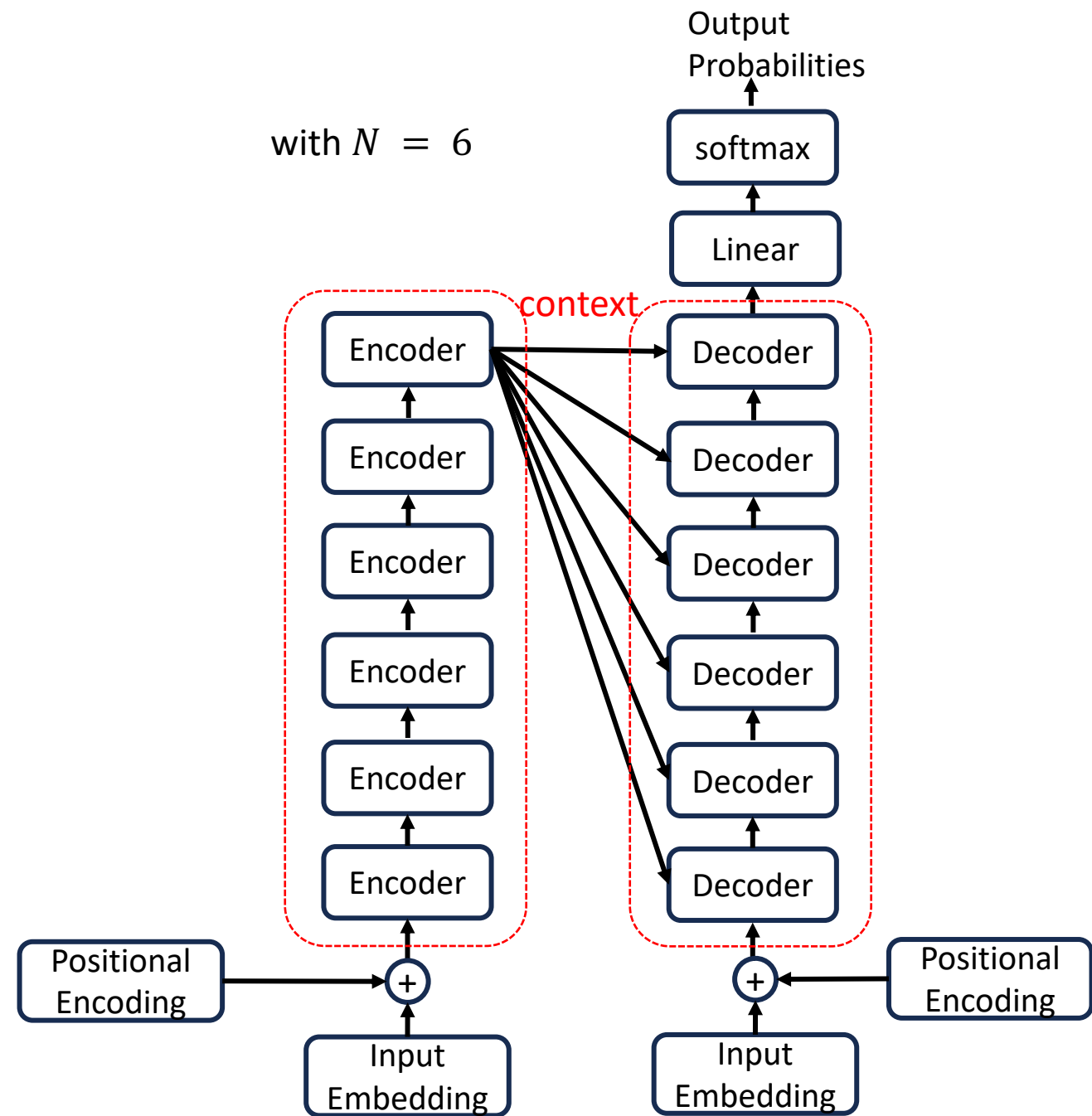
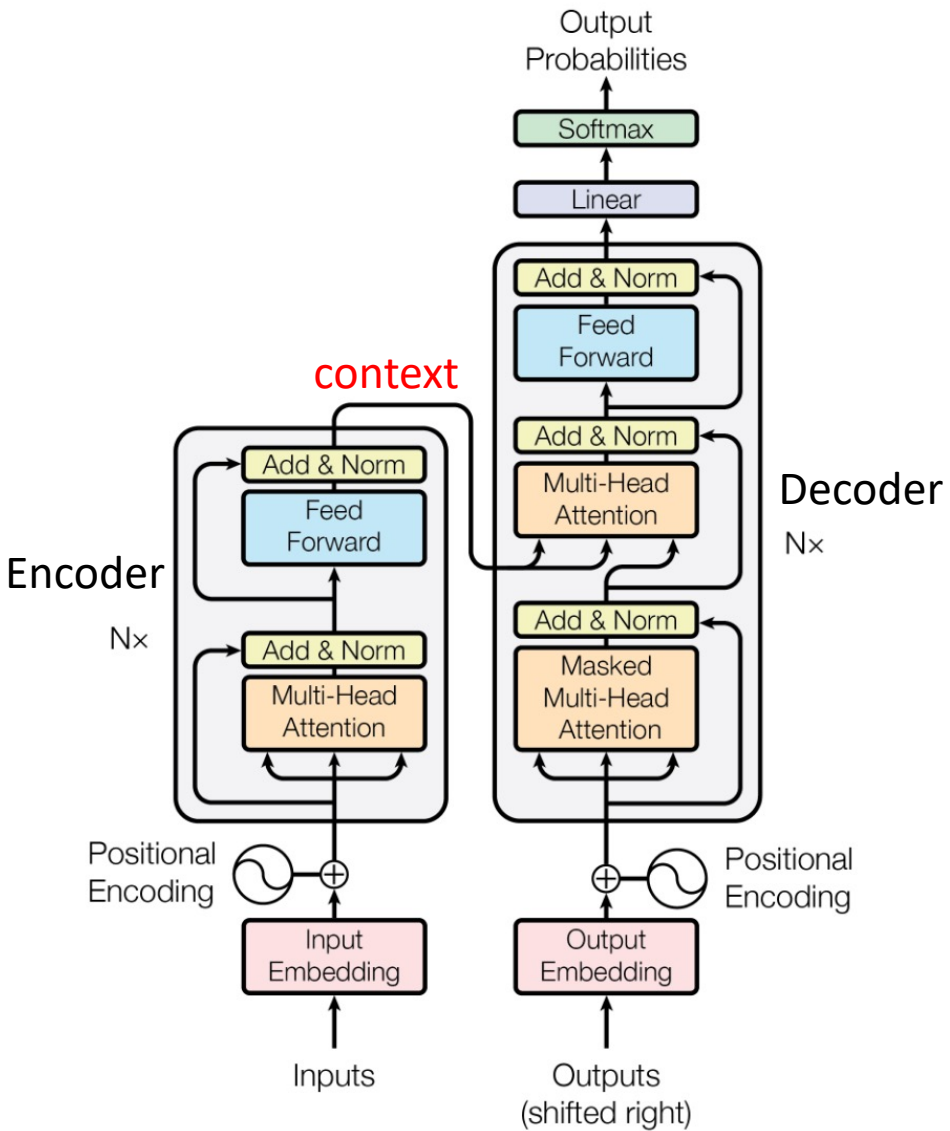
Transformer Architecture



Transformer Architecture

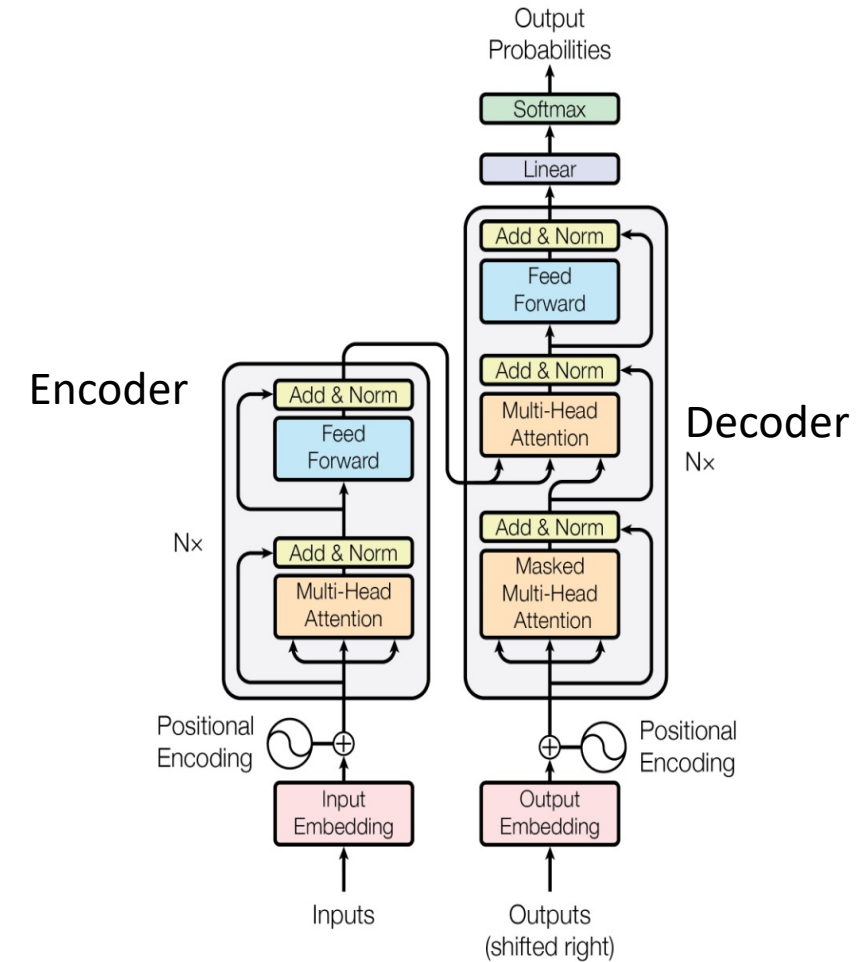


Transformer Architecture



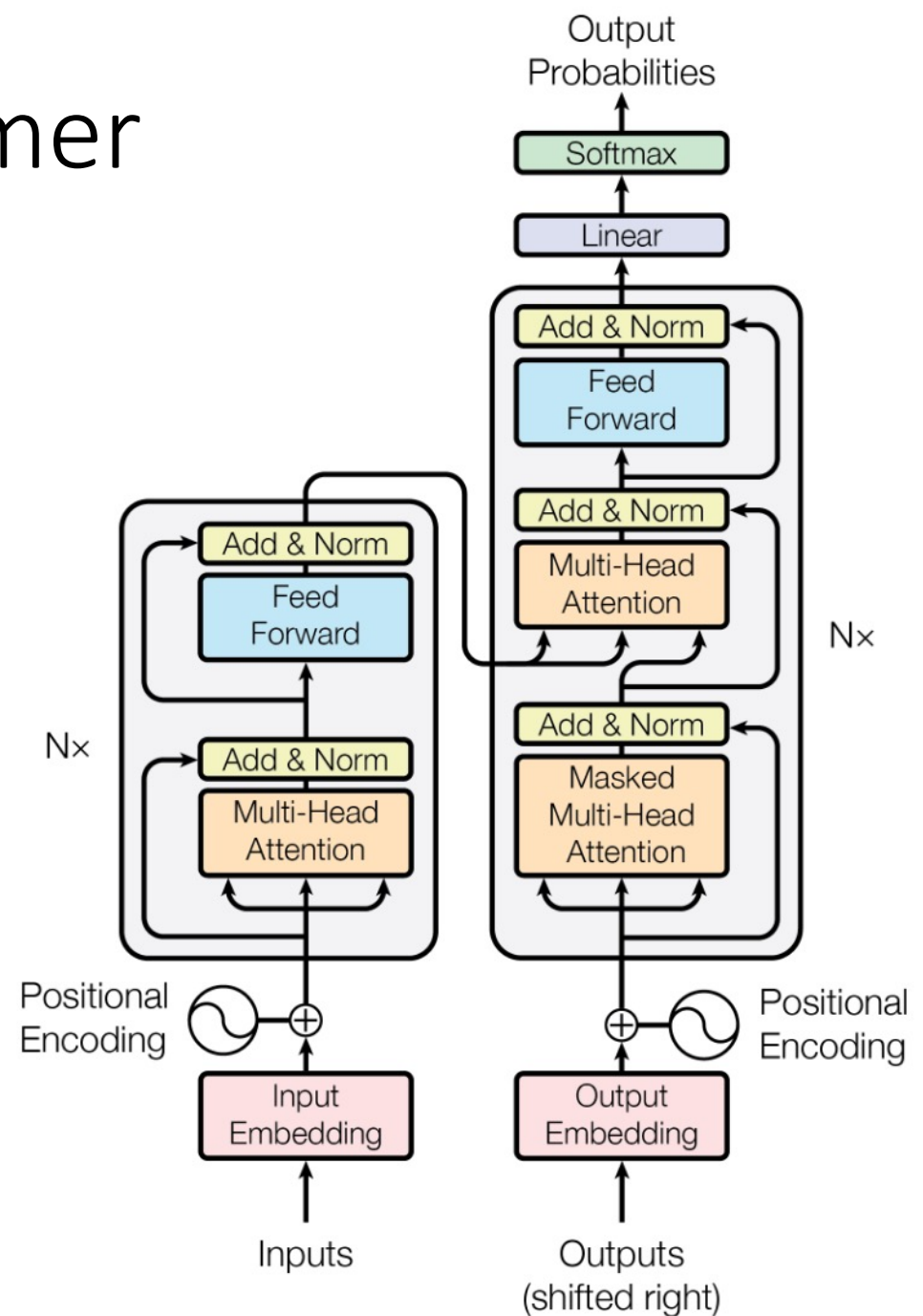
Variants of Transformer Architectures

- Both Encoder and Decoder
 - Language Translation
- Encoder-only
 - Eg: BERT
(Bidirectional encoder representations from transformers)
- Decoder-only
 - Eg: GPT (Generative Pre-trained Transformer)



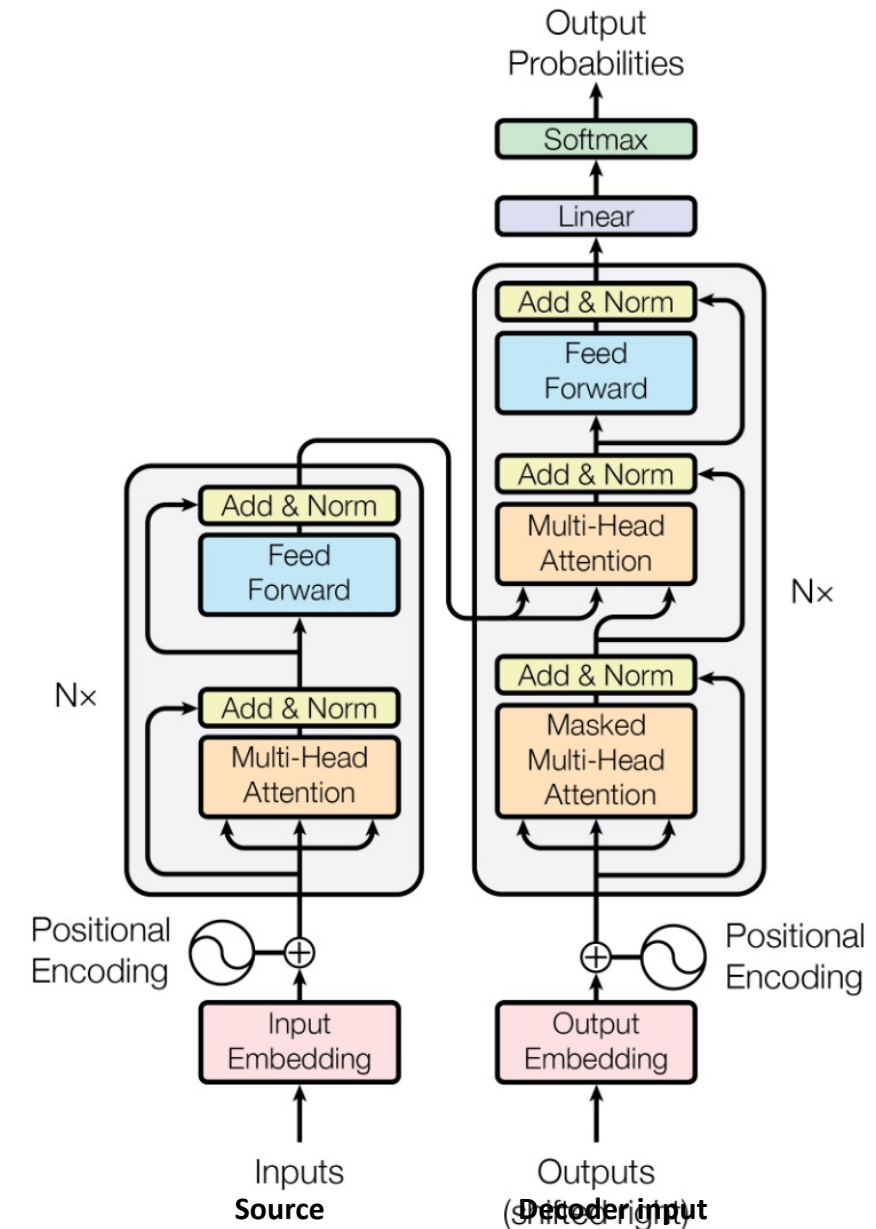
Building Blocks of Transformer

1. Tokenizer
2. Token Embedding
3. Positional Encoder
4. Multi-head Attention
 - a. Self Attention
 - b. Cross Attention
 - c. Masked Attention
5. Skip Connection
6. Normalization
7. Feed Forward Network
8. Linear Layer
9. Softmax

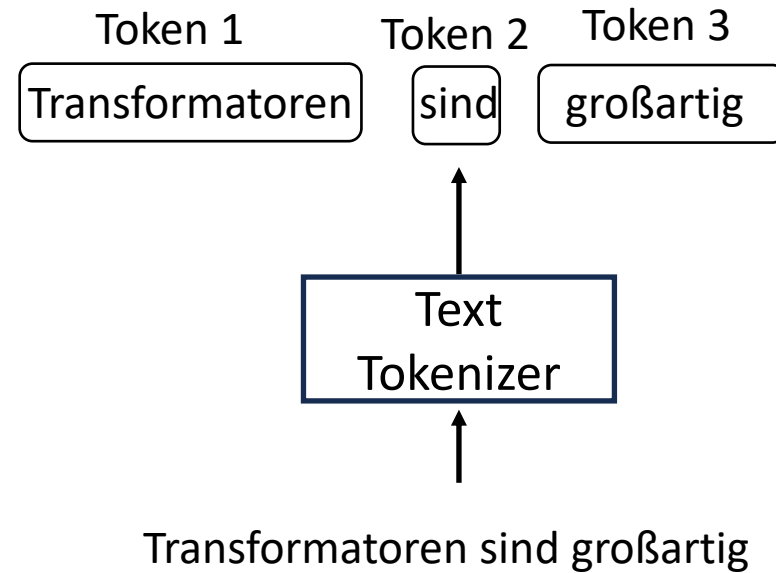
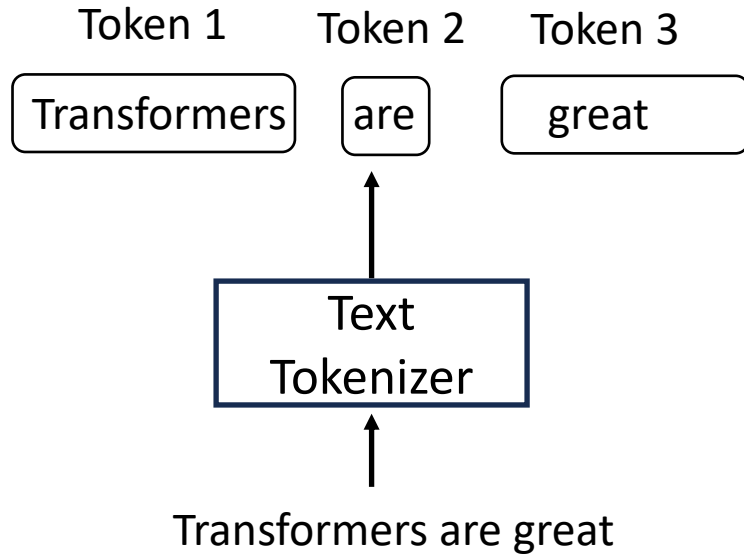


Language Translation: Example

- **Source:** “Transformers are great”
- **Target:** “Transformatoren sind großartig”
- During Training:
 - **Source (English):**
 - ["Transformers", "are", "great"]
 - **Target (German):**
 - ["Transformatoren", "sind", "großartig", "</eos>"]
 - **Decoder input (shifted-right targets):**
 - ["<bos>", "Transformatoren", "sind", "großartig"]



Tokenizer

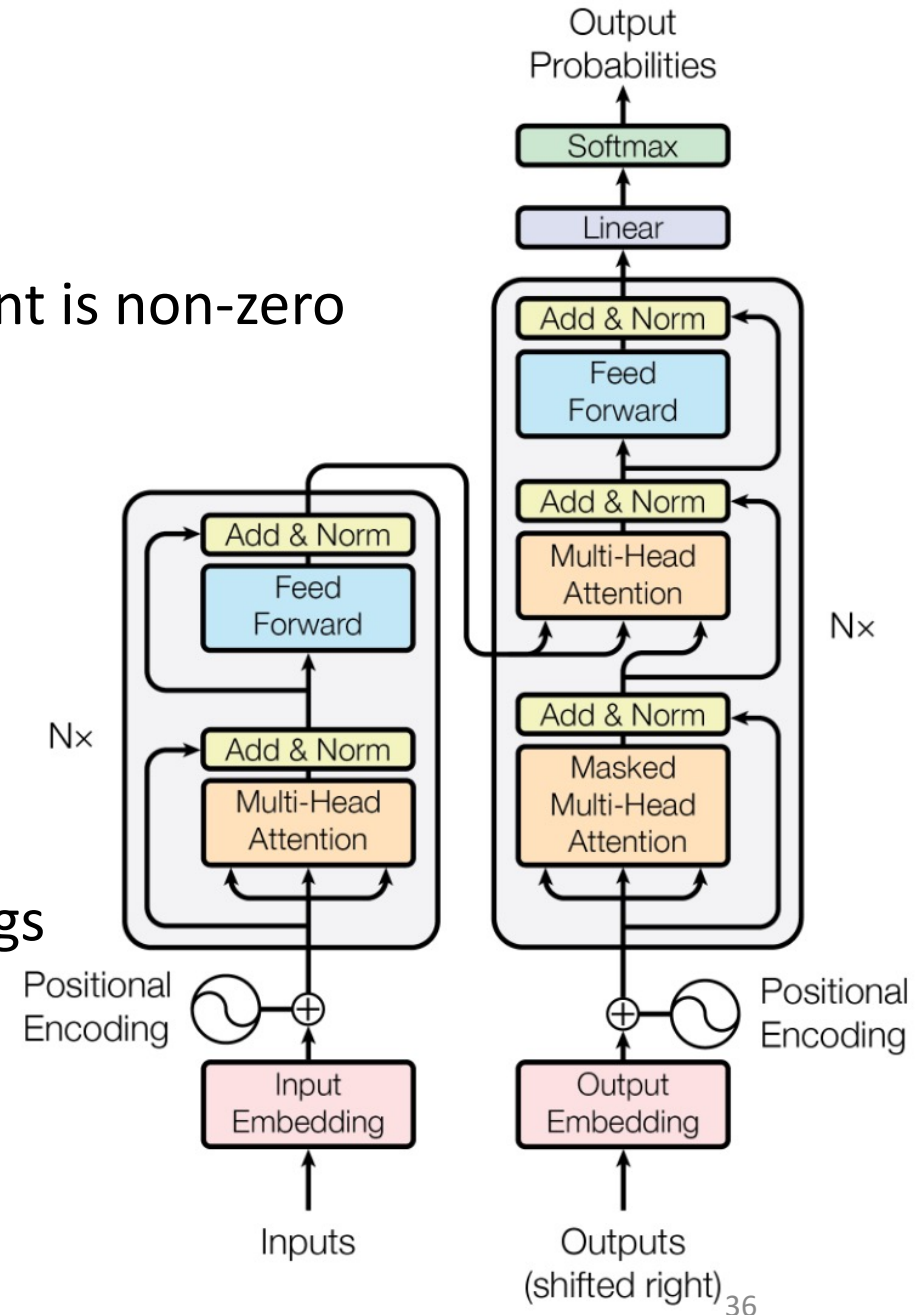


- Natural language
 - *tiktoken* (for use in OpenAI's models), etc.

<https://github.com/openai/tiktoken>

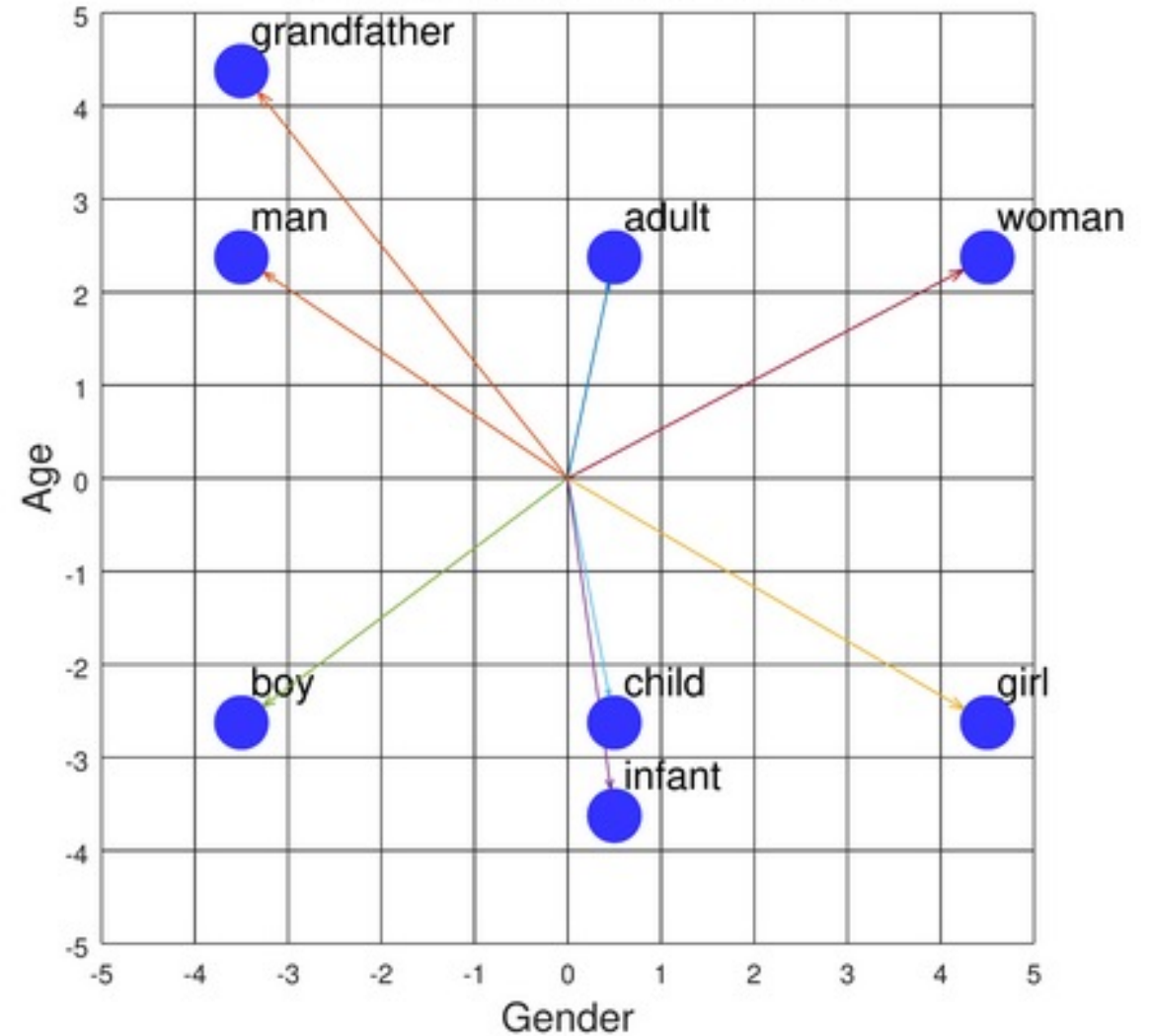
Token Representation

- One-hot encoding
 - Represent each token as a vector with only one element is non-zero
 - Index of non-zero element (one) represent the token
 - Size of the vector is same as size of vocabulary
 - Leads to sparse representation
- Token Embedding
 - Maps each token to dense real vector
 - Can be learned or can use pretrained token embeddings
 - Example: GloVe and Word2Vec



Token Embedding

- Word embeddings capture semantic relations
 - Clustering based on semantics
- Word embeddings are **static**
 - Lack contextual information
- Self-Attention captures contextual information



<https://www.cs.cmu.edu/~dst/WordEmbeddingDemo/tutorial.html>

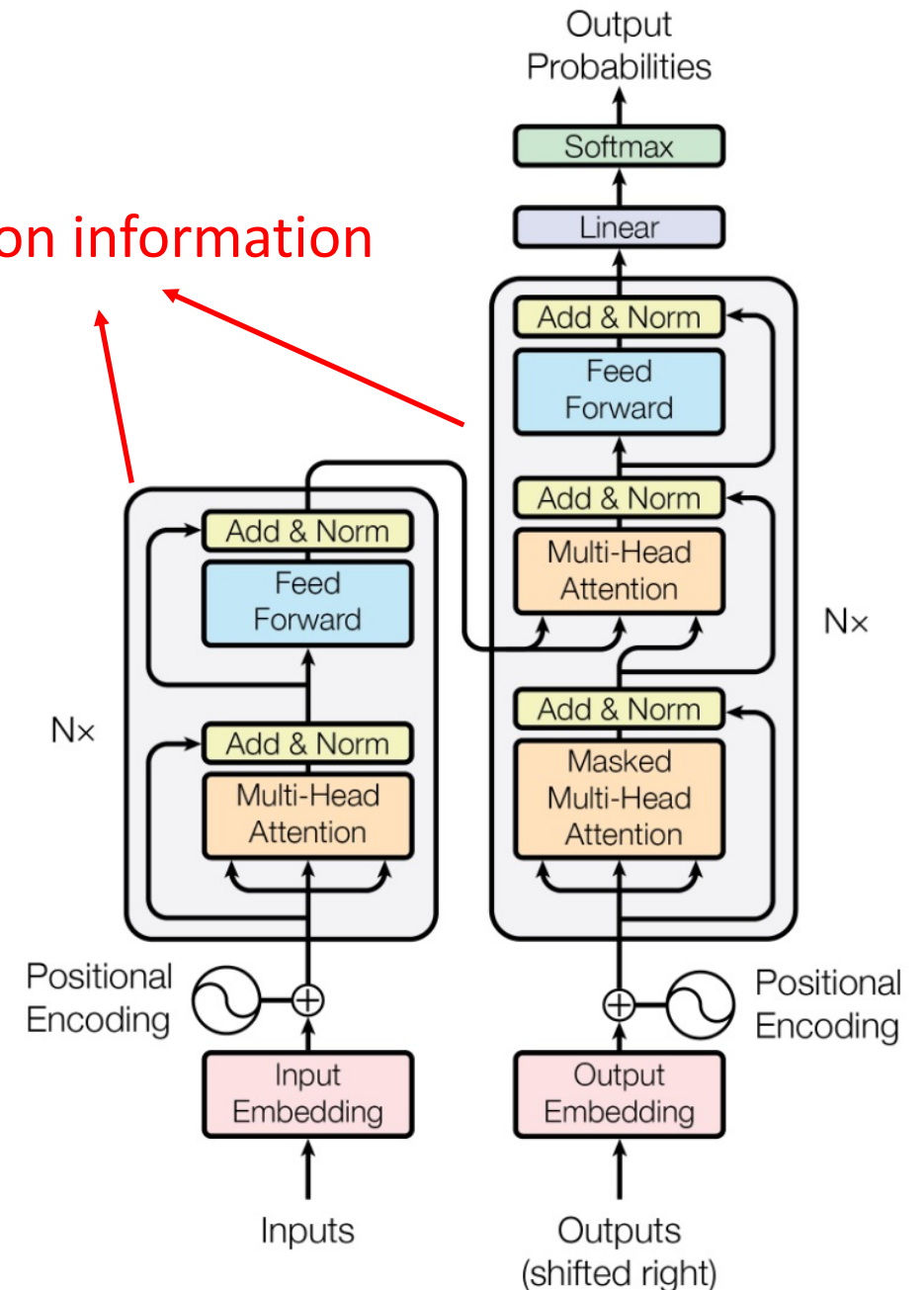
Token Embedding: Example

- “*The phone didn’t connect to **the network** because **it** was down.*”
- “***The phone** didn’t connect to the network because **it** was broken.*”
- Semantically both “**it**”s are same
 - Token embeddings would be the same for both
- But given the context they refer to different tokens
 - Self-attention layer fixes this issue

Why Positional Encoding?

- *Feed Forward* and *Attention* layers process **set** of vectors
 - Set has no inherent order
 - **Ignores** position information or order of data
- Need explicit injection of positional information
- Adding positional encoding to data helps capture relative positional information of words or tokens

ignore position information



Positional Encoding

- Positional Encoding (PE) vector is unique for a given position in a sequence

$$PE(p, 2i) = \sin(w_i p)$$

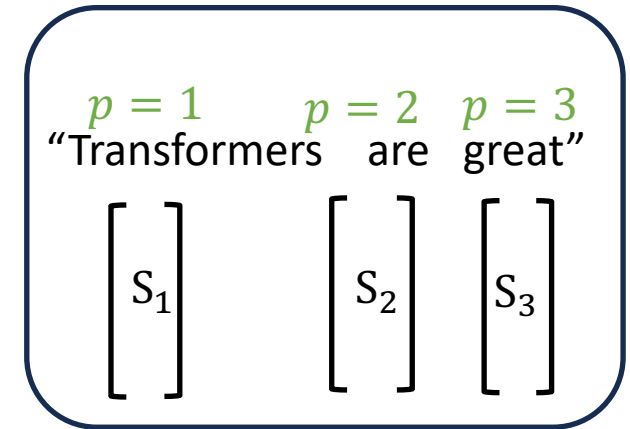
$$PE(p, 2i + 1) = \cos(w_i p)$$

where

$$w_i = \frac{1}{10000^{\frac{2i}{d}}}$$

p is the position of the embedding vector in the sequence of embedding vectors

i is the position within the embedding vector



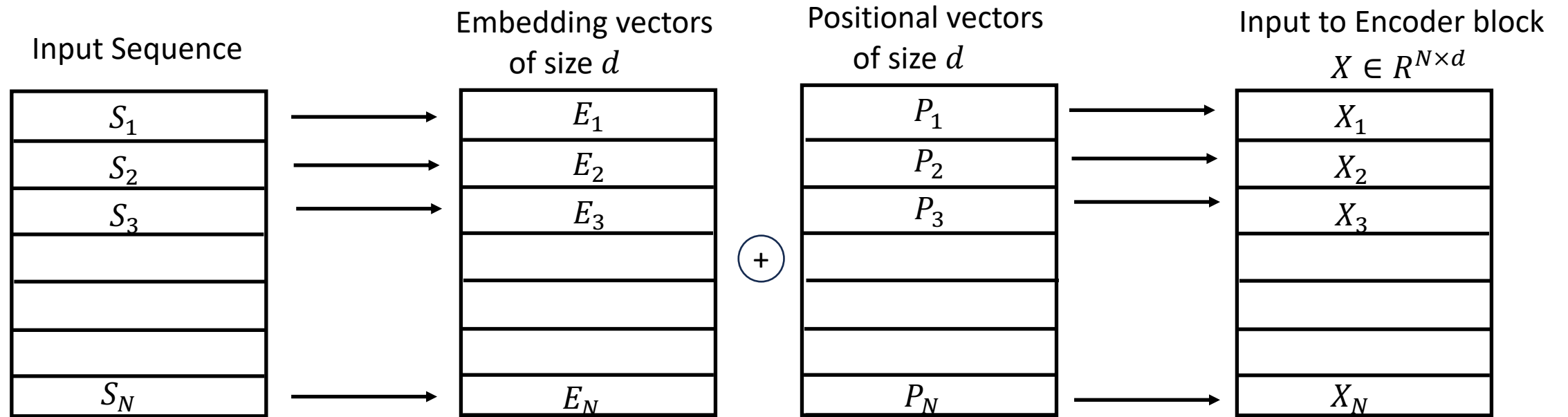
Positional Encoding Vector

$$\text{PE}(p) = \begin{bmatrix} \sin(\omega_1 p) \\ \cos(\omega_1 p) \\ \sin(\omega_2 p) \\ \cos(\omega_2 p) \\ \vdots \\ \sin(\omega_{d/2} p) \\ \cos(\omega_{d/2} p) \end{bmatrix}$$

where $p = pos$, $d = d_{model}$, and $\omega_k = \frac{1}{10000^{2(k-1)/d}}$, for $k = 1, 2, \dots, d/2$.

Number of learnable parameters in positional embedding?

Positional Encoding



How PE helps?

- Dot product of positional encoding vectors reveals the relative distance between two positions
- Given two positions p_1 and p_2 , PE vectors for these two positions:

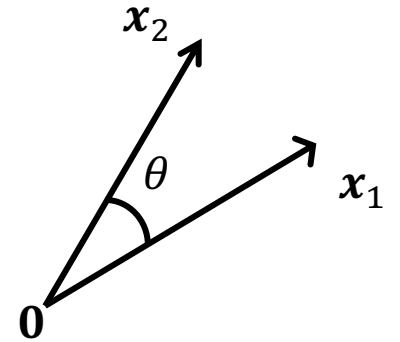
$$\text{PE}(p_1) = \begin{bmatrix} \sin(\omega_1 p_1) \\ \cos(\omega_1 p_1) \\ \sin(\omega_2 p_1) \\ \cos(\omega_2 p_1) \\ \vdots \\ \sin(\omega_{d/2} p_1) \\ \cos(\omega_{d/2} p_1) \end{bmatrix} \quad \text{PE}(p_2) = \begin{bmatrix} \sin(\omega_1 p_2) \\ \cos(\omega_1 p_2) \\ \sin(\omega_2 p_2) \\ \cos(\omega_2 p_2) \\ \vdots \\ \sin(\omega_{d/2} p_2) \\ \cos(\omega_{d/2} p_2) \end{bmatrix}$$

- Inner product of these two PE vectors:

$$\begin{aligned} \text{PE}(p_1) \cdot \text{PE}(p_2) &= \sum_{k=1}^{d/2} \left(\sin(\omega_k p_1) \sin(\omega_k p_2) + \cos(\omega_k p_1) \cos(\omega_k p_2) \right) \\ &= \sum_{k=1}^{d/2} \cos(\omega_k (p_1 - p_2)) \end{aligned} \quad \longrightarrow \quad \begin{array}{l} \text{function of distance} \\ \text{between two vectors} \end{array}$$

Attention: Intuition

- Measures similarity or relation between two vectors
- Inner or Dot product provides a measure of similarity
 - $\mathbf{x}_1 \cdot \mathbf{x}_2 = \mathbf{x}_1^T \mathbf{x}_2 = \mathbf{x}_2^T \mathbf{x}_1 = \langle \mathbf{x}_1, \mathbf{x}_2 \rangle$
 - $\langle \mathbf{x}_1, \mathbf{x}_2 \rangle = \|\mathbf{x}_1\|_2 \|\mathbf{x}_2\|_2 \cos \theta$
- Inner or dot product between the vectors represent relation between them



What does *Attention* focus on?

“The book that John gave to Mary was interesting.”

“Le livre que John a donné à Mary était intéressant.”

What does *Attention* focus on?

“The book that John gave to Mary was interesting.”

“Le livre que John a donné à Mary était intéressant.”

Self-Attention (Encoder)

- Determiner-noun binding
“The” → “book”
- Subject-verb agreement
“book” ↔ “was”
- Long-range dependency
“book” ↔ “that John gave to Mary”
- Entity tracking
Tracks “John”, “Mary” throughout
- Word order mapping
Helps rearrange for French syntax



Self-attention block at *Decoder* captures similar relations within the sentence in French

Cross-Attention (Decoder)

- Article choice:
“The” ↔ “Le”
- Noun alignment
“book” ↔ “livre”
- Subject-verb agreement:
“book” ↔ “était”
- Verb decomposition
“gave to” ↔ “a donné à”

Why multiple heads?

- To simultaneously capture multiple types of relationships

“The book that John gave to Mary was interesting.”

“Le livre que John a donné à Mary était intéressant.”

- **Head 1:** Subject–verb agreement
Focuses on “book” \leftrightarrow “was/était”
- **Head 2:** Long-range dependency
Links “book” \leftrightarrow “that John gave to Mary”
- **Head 3:** Entity tracking
Tracks “John”, “Mary” throughout
- **Head 4:** Local context
Focuses on “gave to” \leftrightarrow “a donné à”
- **Head 5:** Word order mapping
Helps rearrange for French syntax

Query (W^Q), Key (W^K) and Value (W^V) Parameters

- Model uses W^Q , W^K and W^V parameters to capture these relationships
 - $W^Q \in R^{d \times d}$
 - $W^K \in R^{d \times d}$
 - $W^V \in R^{d \times d}$
- Projection of the input embedding onto W^Q converts the input to a query
- Projection of the input embedding onto W^K converts the input to a key
- Projection of the input embedding onto W^V converts the input to a value

Differentiable parameters and learnable
Hence the name *differentiable* dictionary

How to think of Query (Q), Key (K) and Value (V)?

- **Query:**

“What do I need right now?”

- **Key:**

“How should others find me?”

- **Value:**

“What information do I provide if they pick me?”

Example for Self-Attention

*“The **book** that John gave to Mary **was** interesting.”*

“Le livre que John a donné à Mary était intéressant.”

- Example: Resolving subject-verb
 - **Query** (W^Q) while processing the token “**was**” learns to ask “**Find the subject?**”
 - **Key** (W^K) of “**book**” learns to say “***I am a singular noun and may be a subject***”
 - **Value** (W^V) of the “**book**” provides the semantic content

Multi-Head Attention (*MHA*)

- $MHA(Q, K, V) = HW_0$

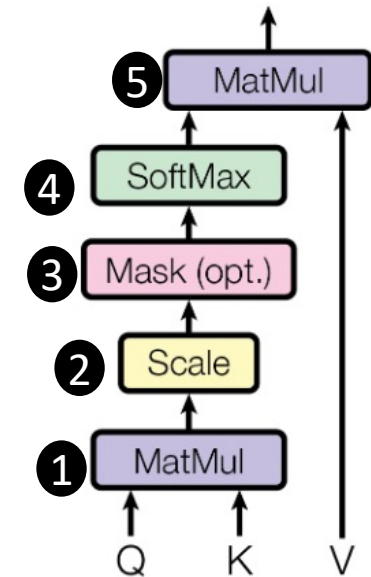
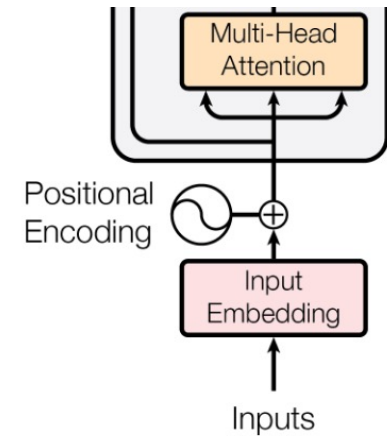
where $H = \text{concat}(\text{head}_1, \dots, \text{head}_h)$

$$\text{head}_i = \text{softmax} \left(\underbrace{\frac{(QW_i^Q)(KW_i^K)^T}{\sqrt{d_K}} \circ M}_{\text{5 MatMul}} \right) (VW_i^V)$$

Annotations in the diagram:

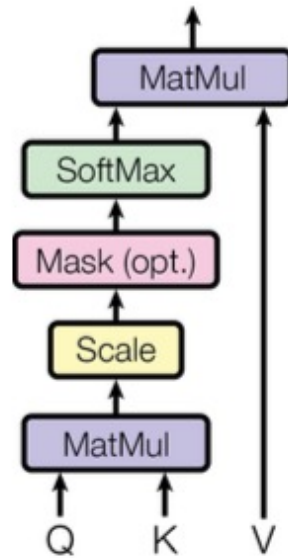
- ① MatMul: $(QW_i^Q)(KW_i^K)^T$
- ② Scale: $\sqrt{d_K}$
- ③ Mask (optional): M
- ④ Softmax: softmax
- ⑤ MatMul: (VW_i^V)

W s are learnable parameters

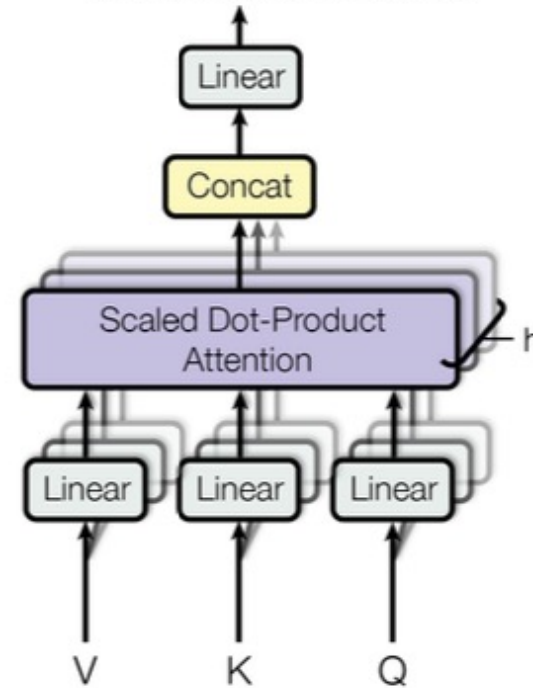


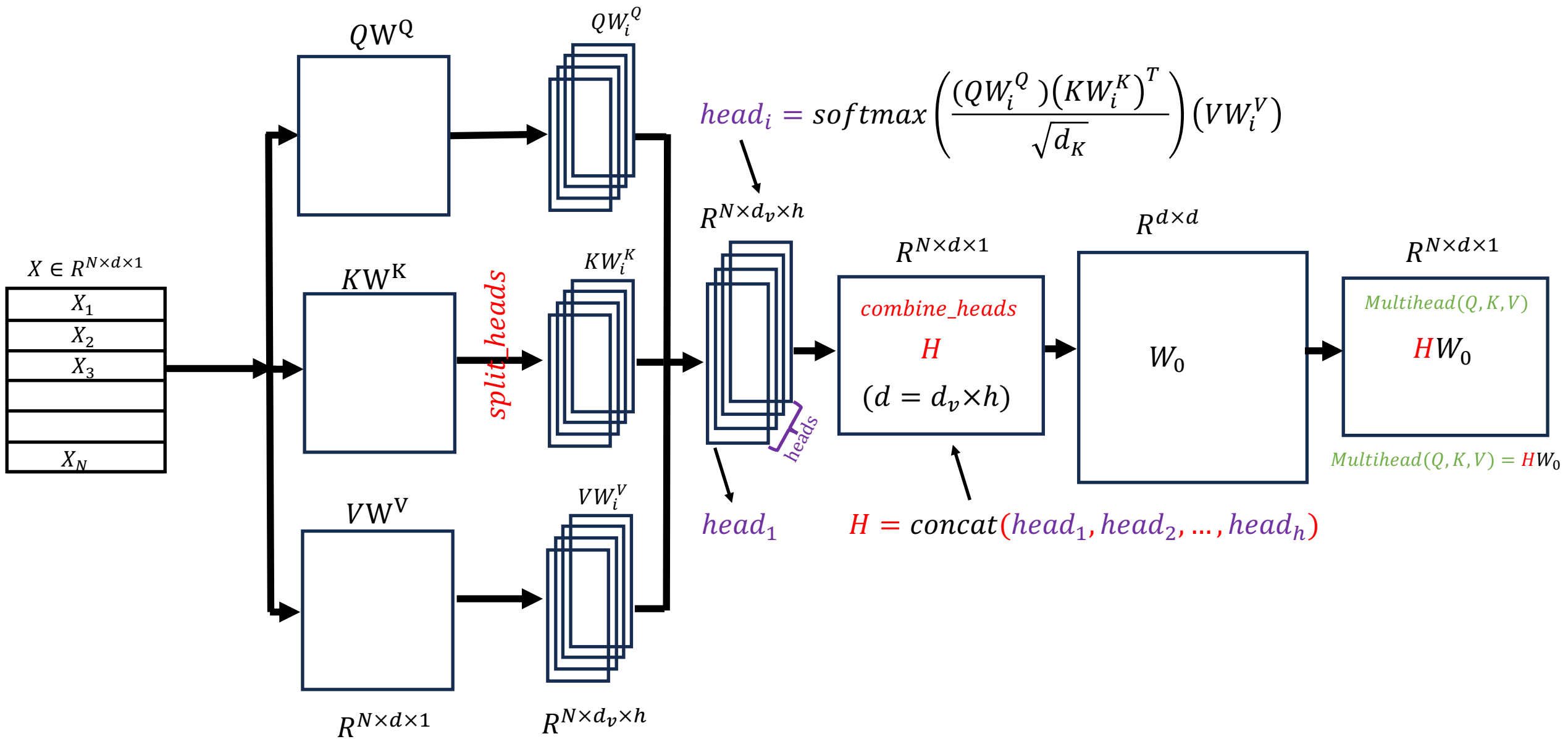
Multi-Head Attention (*MHA*)

Scaled Dot-Product Attention



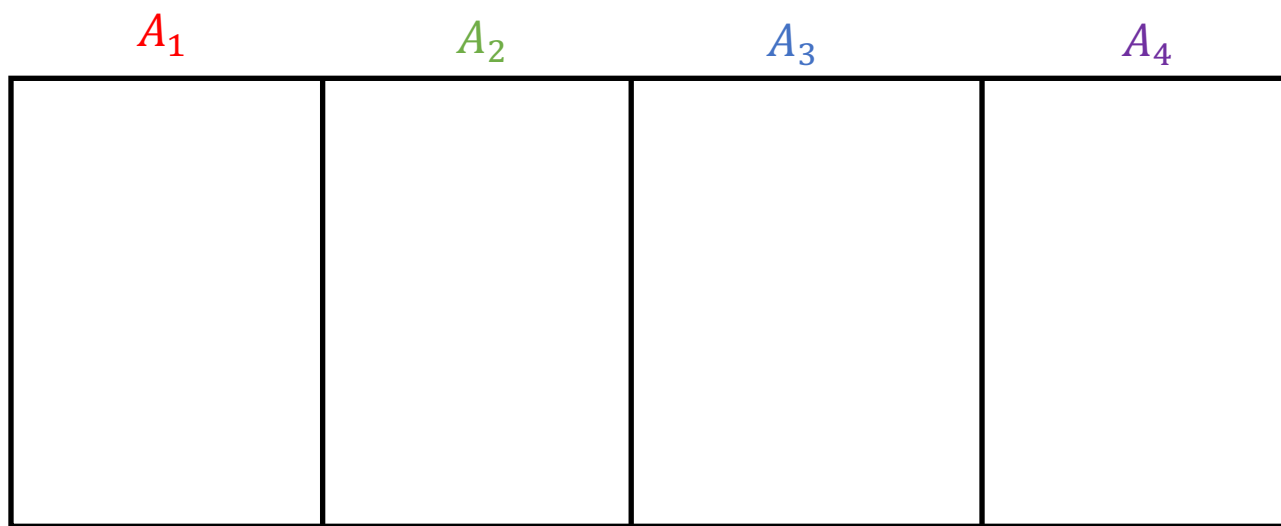
Multi-Head Attention





$$d_v = \frac{d}{h}$$

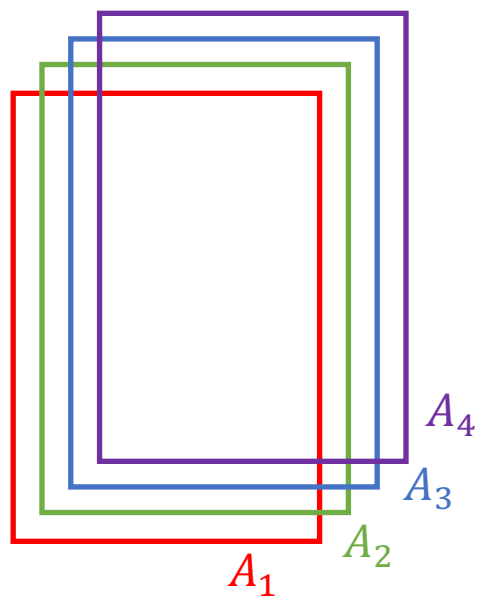
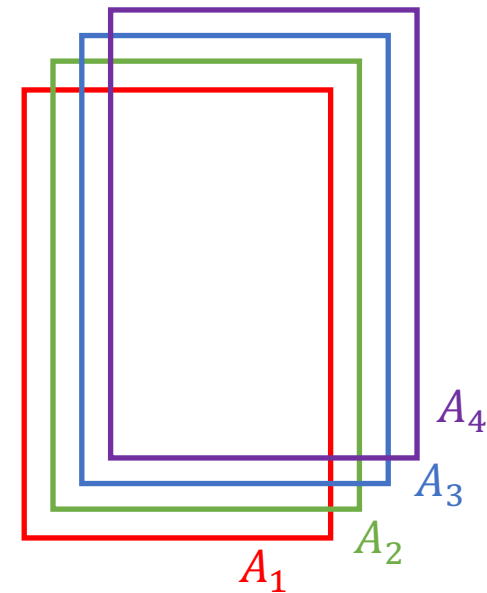
Refer next slide for *split_heads* and *combine_heads*



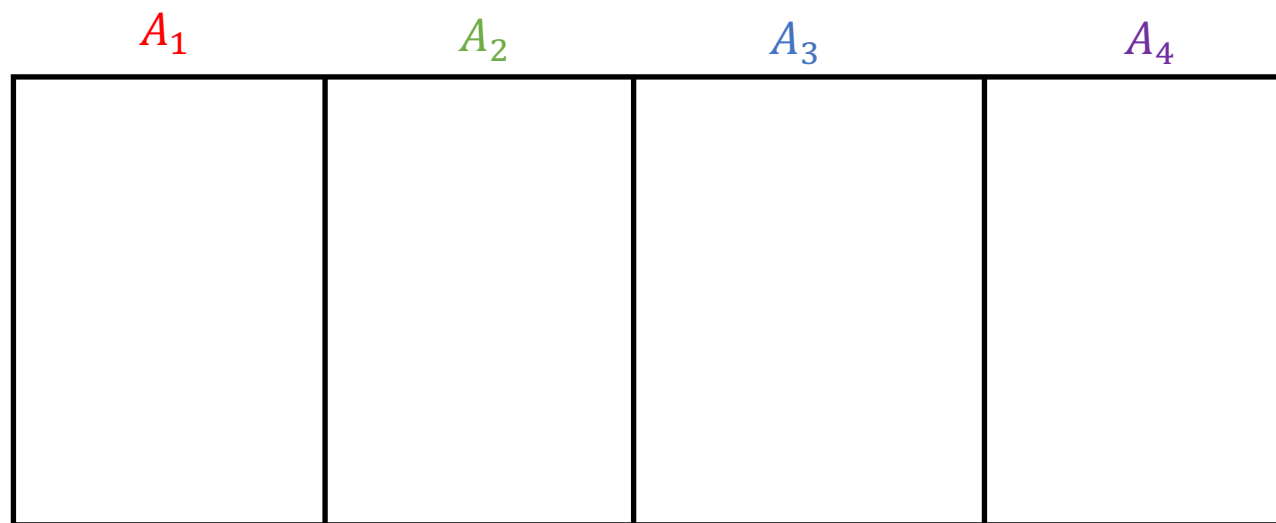
split_heads



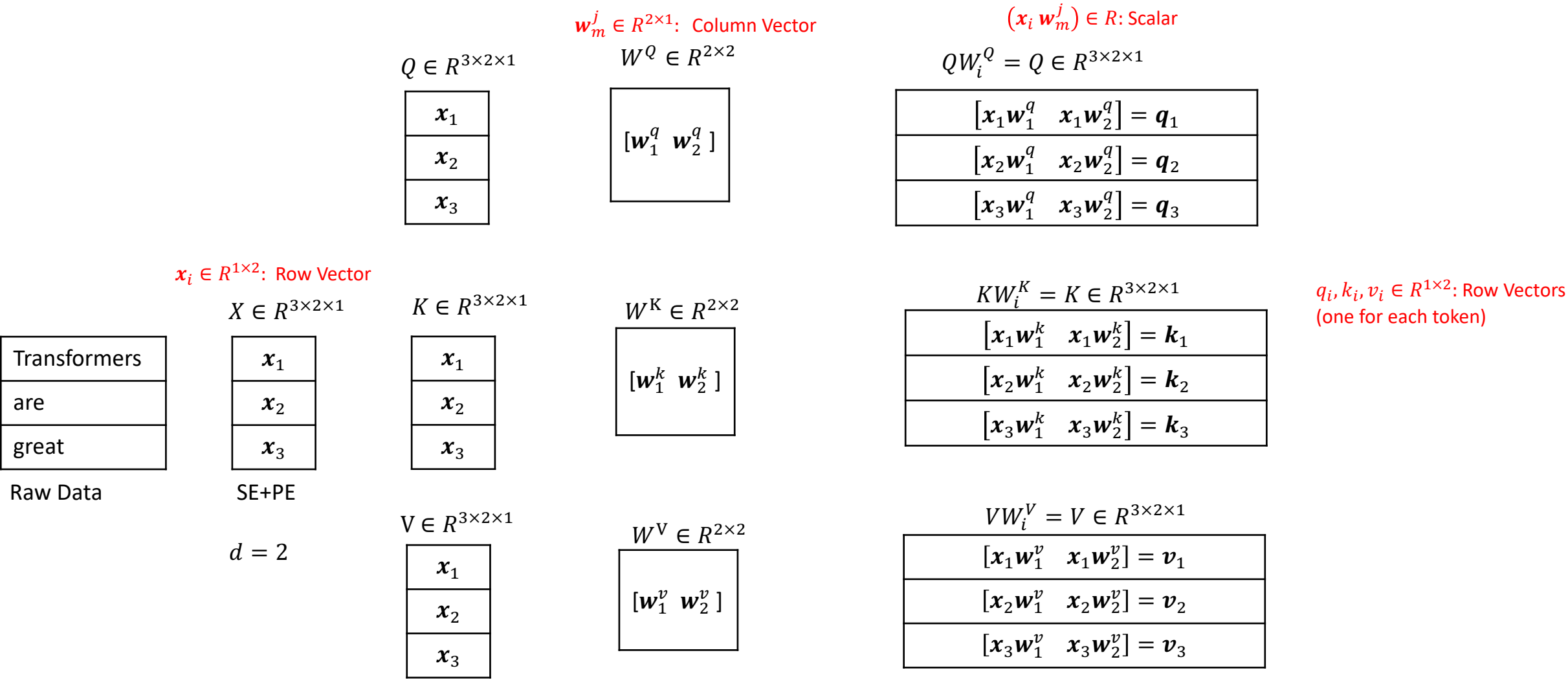
- Split along embedding dimension
- Each row corresponds to a token



combine_heads



Scaled Dot-Product Attention: Example with Single Head



$$\text{head}_i = \text{softmax} \left(\frac{(QW_i^Q)(KW_i^K)^T}{\sqrt{d_K}} \right) (VW_i^V)$$

$$\frac{(QW_i^Q)(KW_i^K)^T}{\sqrt{d_k}} = \begin{bmatrix} \frac{\mathbf{q}_1 \mathbf{k}_1^T}{\sqrt{d_1}} & \frac{\mathbf{q}_1 \mathbf{k}_2^T}{\sqrt{d_1}} & \frac{\mathbf{q}_1 \mathbf{k}_3^T}{\sqrt{d_1}} \\ \frac{\mathbf{q}_2 \mathbf{k}_1^T}{\sqrt{d_2}} & \frac{\mathbf{q}_2 \mathbf{k}_2^T}{\sqrt{d_2}} & \frac{\mathbf{q}_2 \mathbf{k}_3^T}{\sqrt{d_2}} \\ \frac{\mathbf{q}_3 \mathbf{k}_1^T}{\sqrt{d_3}} & \frac{\mathbf{q}_3 \mathbf{k}_2^T}{\sqrt{d_3}} & \frac{\mathbf{q}_3 \mathbf{k}_3^T}{\sqrt{d_3}} \end{bmatrix}$$

$$\text{softmax} \left(\frac{(QW_i^Q)(KW_i^K)^T}{\sqrt{d_k}} \right) = \begin{bmatrix} \text{softmax} \left(\begin{bmatrix} \frac{\mathbf{q}_1 \mathbf{k}_1^T}{\sqrt{d_1}} & \frac{\mathbf{q}_1 \mathbf{k}_2^T}{\sqrt{d_1}} & \frac{\mathbf{q}_1 \mathbf{k}_3^T}{\sqrt{d_1}} \end{bmatrix} \right) \\ \text{softmax} \left(\begin{bmatrix} \frac{\mathbf{q}_2 \mathbf{k}_1^T}{\sqrt{d_2}} & \frac{\mathbf{q}_2 \mathbf{k}_2^T}{\sqrt{d_2}} & \frac{\mathbf{q}_2 \mathbf{k}_3^T}{\sqrt{d_2}} \end{bmatrix} \right) \\ \text{softmax} \left(\begin{bmatrix} \frac{\mathbf{q}_3 \mathbf{k}_1^T}{\sqrt{d_3}} & \frac{\mathbf{q}_3 \mathbf{k}_2^T}{\sqrt{d_3}} & \frac{\mathbf{q}_3 \mathbf{k}_3^T}{\sqrt{d_3}} \end{bmatrix} \right) \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}$$

Attention Scores

	Transformers	are	great
Transformers	a_{11}	a_{12}	a_{13}
are	a_{21}	a_{22}	a_{23}
great	a_{31}	a_{32}	a_{33}

Visualization of Attention Scores From “Attention is all you need”

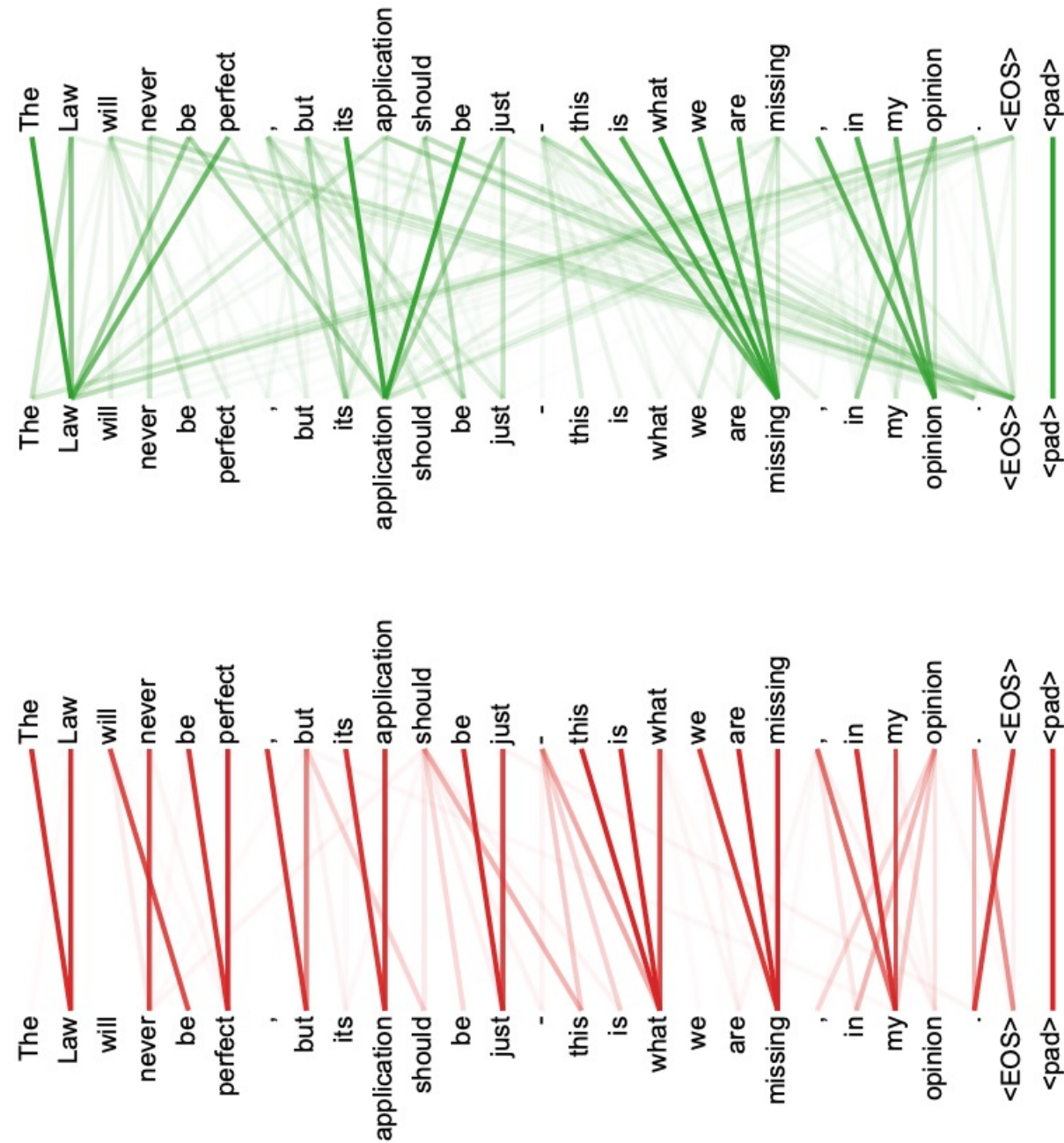


Figure 5: Many of the attention heads exhibit behaviour that seems related to the structure of the sentence. We give two such examples above, from two different heads from the encoder self-attention at layer 5 of 6. The heads clearly learned to perform different tasks.

Softmax - Recap

$$\mathbf{z} = [z_1, z_2, z_3, \dots, z_n]$$

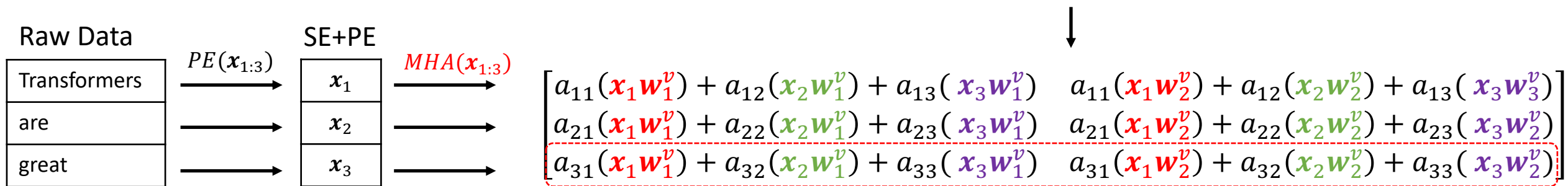
$$\textit{softmax}(\mathbf{z}) = \left[\frac{\exp(z_1)}{\sum_i \exp(z_i)}, \frac{\exp(z_2)}{\sum_i \exp(z_i)}, \frac{\exp(z_3)}{\sum_i \exp(z_i)}, \dots, \frac{\exp(z_n)}{\sum_i \exp(z_i)} \right]$$

$$head_i = softmax \left(\frac{(QW_i^Q)(KW_i^K)^T}{\sqrt{d_K}} \right) (VW_i^V)$$

$$softmax \left(\frac{(QW_i^Q)(KW_i^K)^T}{\sqrt{d_k}} \right) \begin{bmatrix} \mathbf{x}_1 \mathbf{w}_1^v & \mathbf{x}_1 \mathbf{w}_2^v \\ \mathbf{x}_2 \mathbf{w}_1^v & \mathbf{x}_2 \mathbf{w}_2^v \\ \mathbf{x}_3 \mathbf{w}_1^v & \mathbf{x}_3 \mathbf{w}_2^v \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \mathbf{w}_1^v & \mathbf{x}_1 \mathbf{w}_2^v \\ \mathbf{x}_2 \mathbf{w}_1^v & \mathbf{x}_2 \mathbf{w}_2^v \\ \mathbf{x}_3 \mathbf{w}_1^v & \mathbf{x}_3 \mathbf{w}_2^v \end{bmatrix}$$

Attention Scores

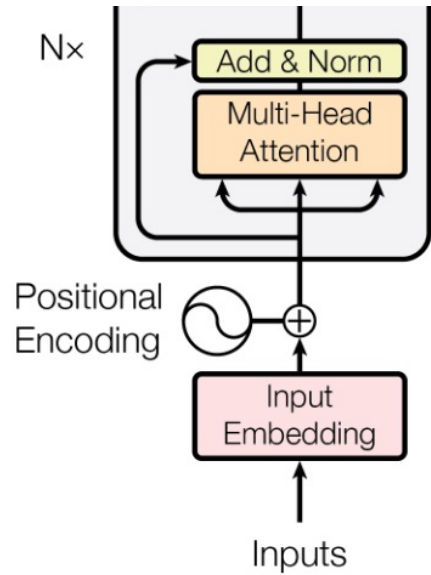
Each row vector includes **semantic**, **positional** and **contextual** information of a word



Embedding of “great” is a linear combination of embeddings of “Transformers”, “are”, and “great”

$$a_{ij} = \frac{\exp(q_i k_j^T)}{\sum_m \exp(q_i k_m^T)}$$

Add and Norm after MHA

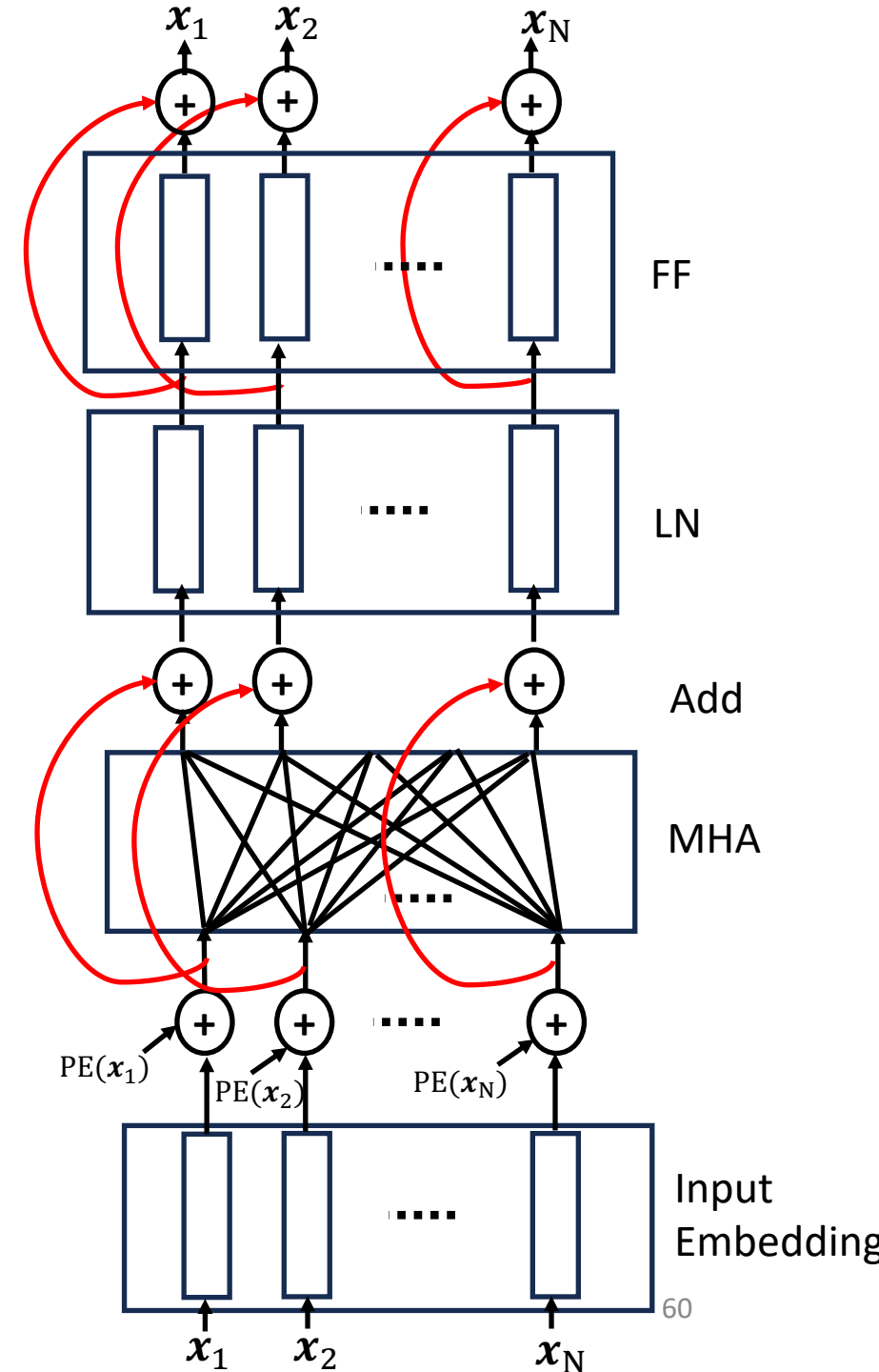


Layer Normalization

Residual connection

$$Output = LN(MHA(x_{1:N}) + x_{1:N})$$

LN is applied for each embedding vector independently (refer next slide)



Layer Normalization (LN)

- $LN(x_i) = \gamma \circ \left(\frac{x_i - \mu_i}{\sigma_i} \right) + \beta$

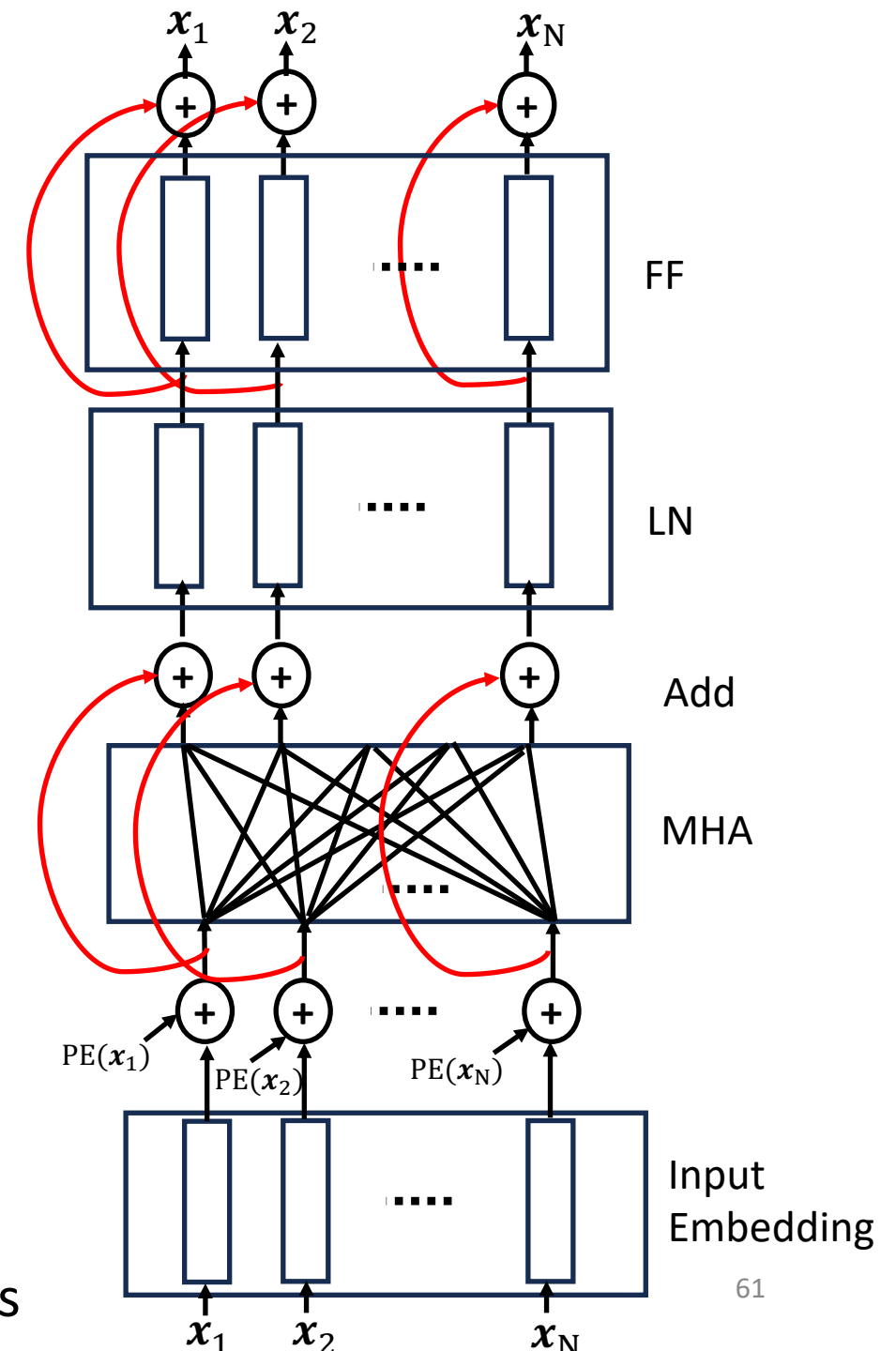
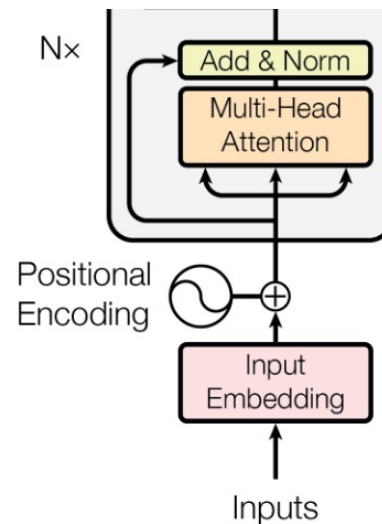
where

$$\mu_i = \text{mean}(x_i)$$

$$\sigma_i = \text{std}(x_i)$$

$$x_i, \gamma, \beta \in \mathbb{R}^d$$

γ and β are learnable parameters shared across the tokens



Add and Norm after MHA

Post-Layer Norm:

Layer norm is applied after MHA block

Layer Normalization

$$\text{Output} = \text{LN}(\text{MHA}(\mathbf{x}_{1:N}) + \mathbf{x}_{1:N})$$

Residual connection

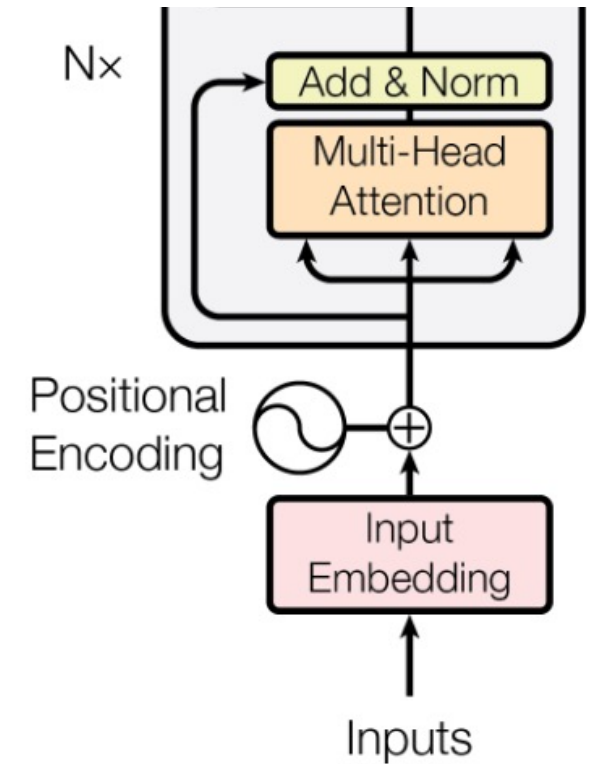
Pre-Layer Norm:

Layer norm is applied before MHA block

Layer Normalization

$$\text{Output} = \text{MHA}(\text{LN}(\mathbf{x}_{1:N})) + \mathbf{x}_{1:N}$$

Residual connection

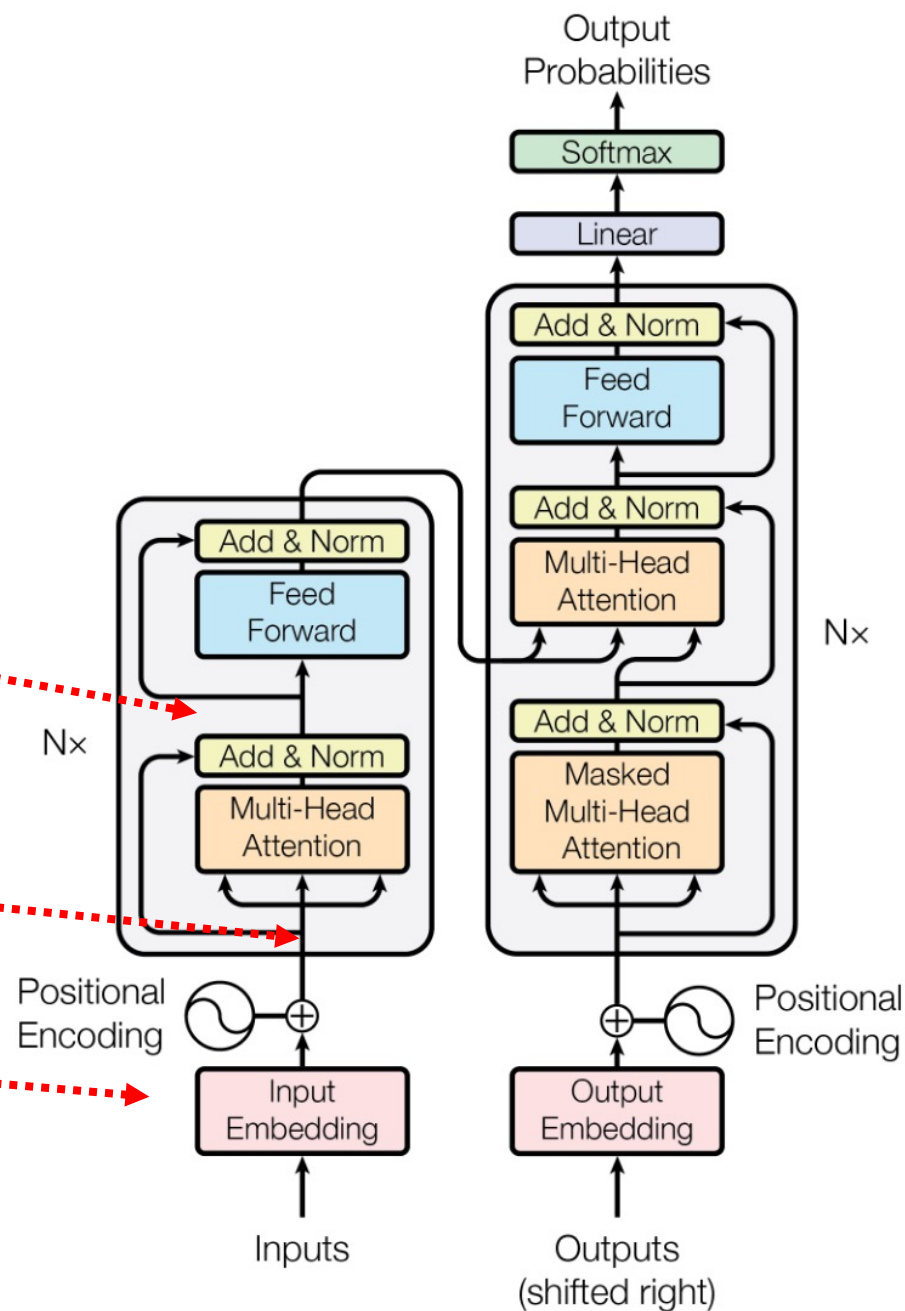


Semantic Embedding to Contextual Embedding

Embeddings with Semantic, Positional and Contextual Information

Embeddings with Semantic and Positional Information

Embeddings with Semantic Information



Positionwise Feedforward Layer

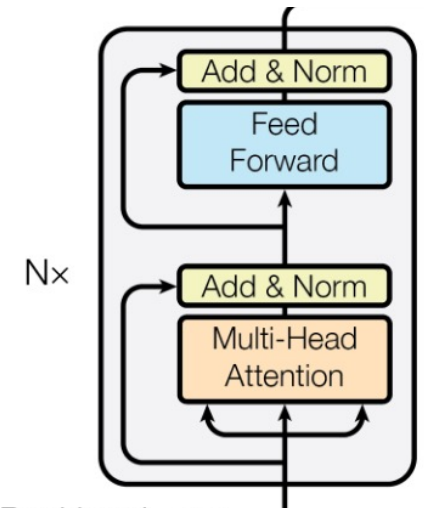
- A two-layer MLP with ReLU activation, skip connection and layer normalization

Layer Normalization Linear function ReLU activation Residual Connection

$$MLP(\mathbf{x}_i) = LN(\mathbf{f}^{[2]}(\mathbf{g}^{[1]}(\mathbf{f}^{[1]}(\mathbf{x}_i))) + \mathbf{x}_i)$$
$$\mathbf{f}^{[j]}(\mathbf{x}_i) = (\mathbf{W}^{[j]})^T \mathbf{x}_i + \mathbf{b}^{[j]}$$

- Each embedding vector is transformed independently
- Example with ReLU activation:

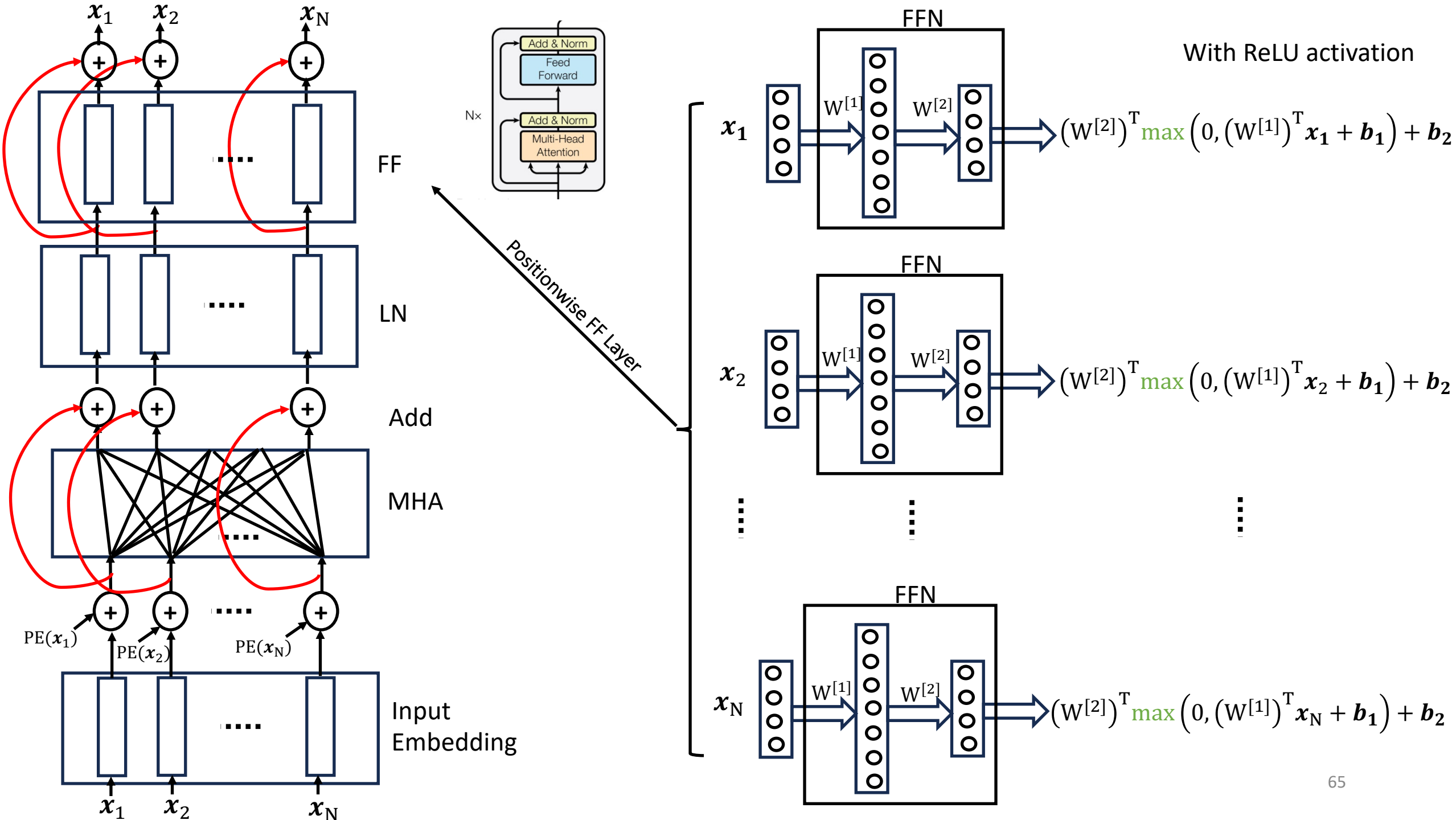
$$MLP(\mathbf{x}_i) = LN\left((\mathbf{W}^{[2]})^T \max\left(0, (\mathbf{W}^{[1]})^T \mathbf{x}_i + \mathbf{b}_1\right) + \mathbf{b}_2 + \mathbf{x}_i\right)$$



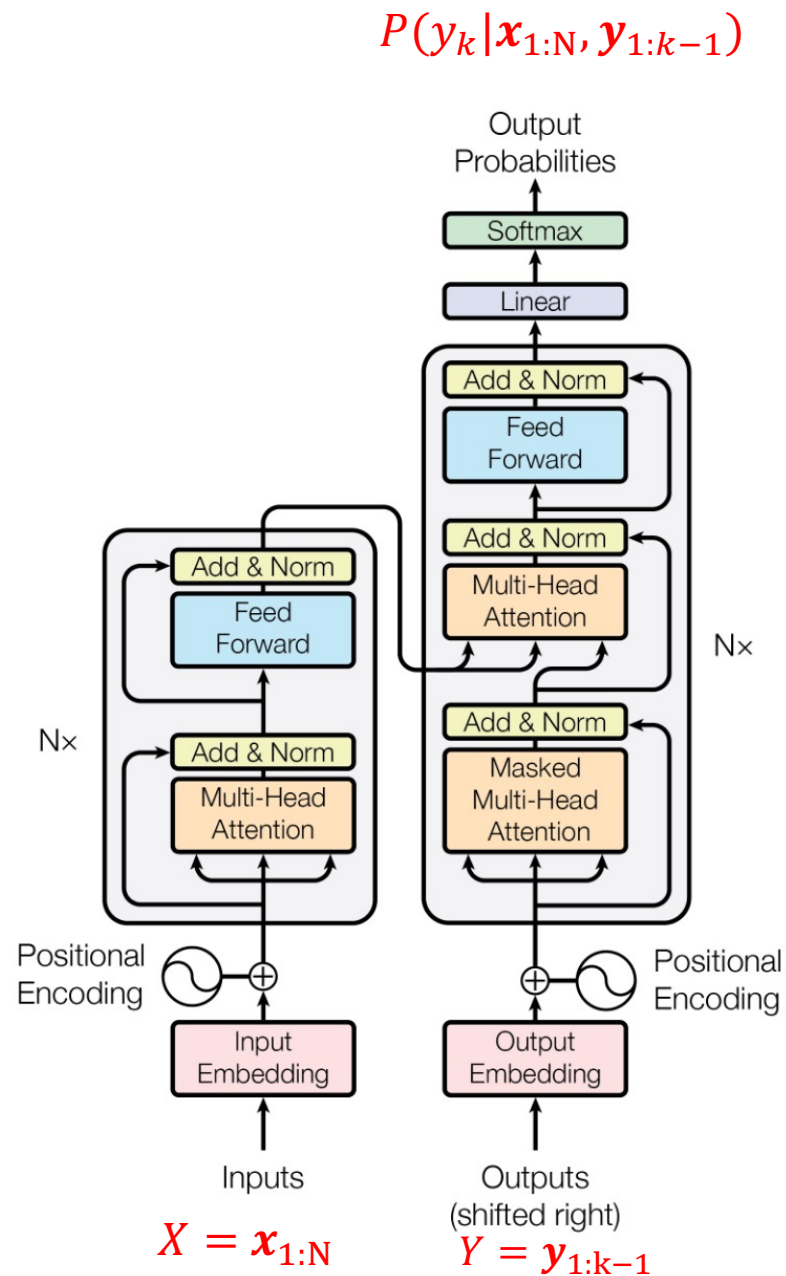
$$\mathbf{W}^{[1]} \in \mathbb{R}^{d \times d_f}$$

$$\mathbf{W}^{[2]} \in \mathbb{R}^{d_f \times d}$$

$\mathbf{W}^{[j]}$ are shared across the tokens

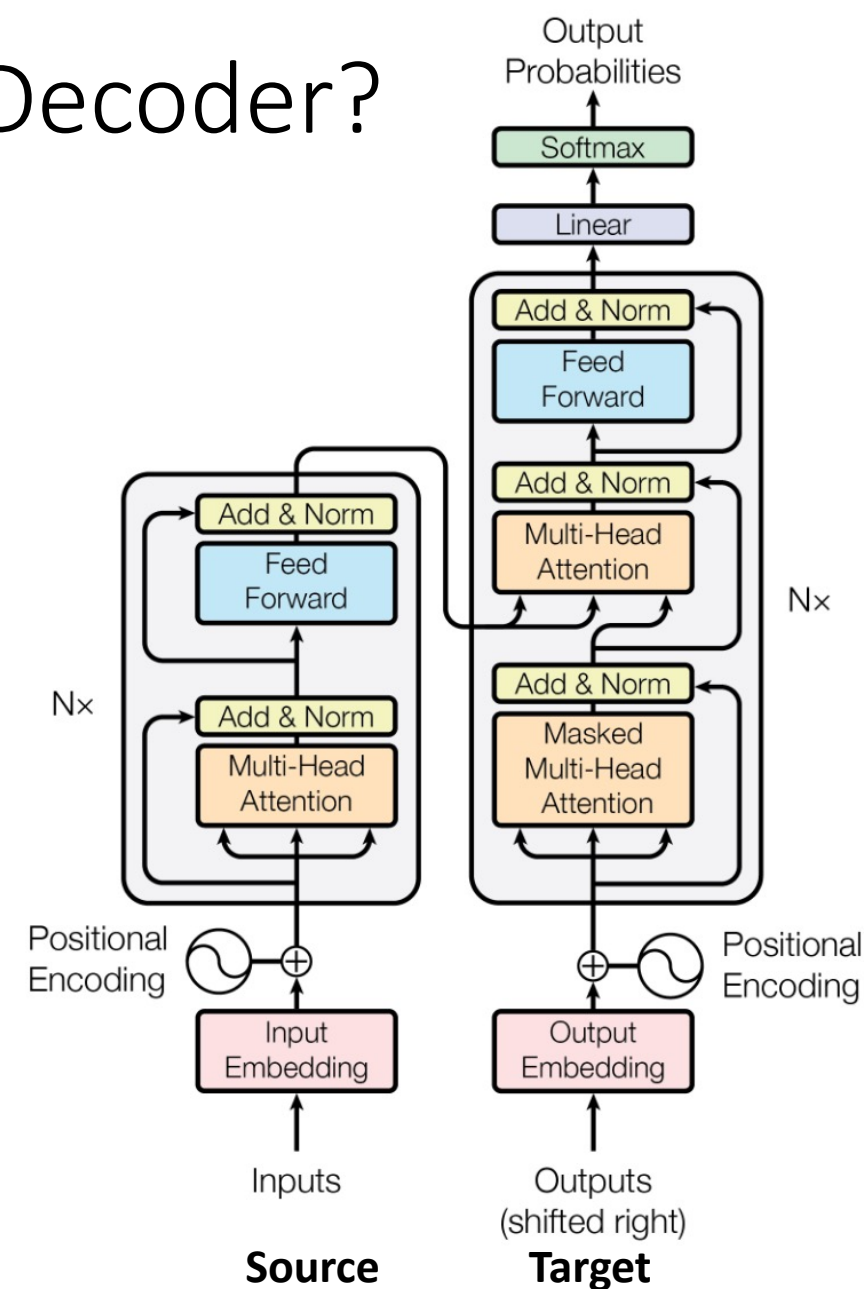


Why Masked Multi-head Attention at Decoder?



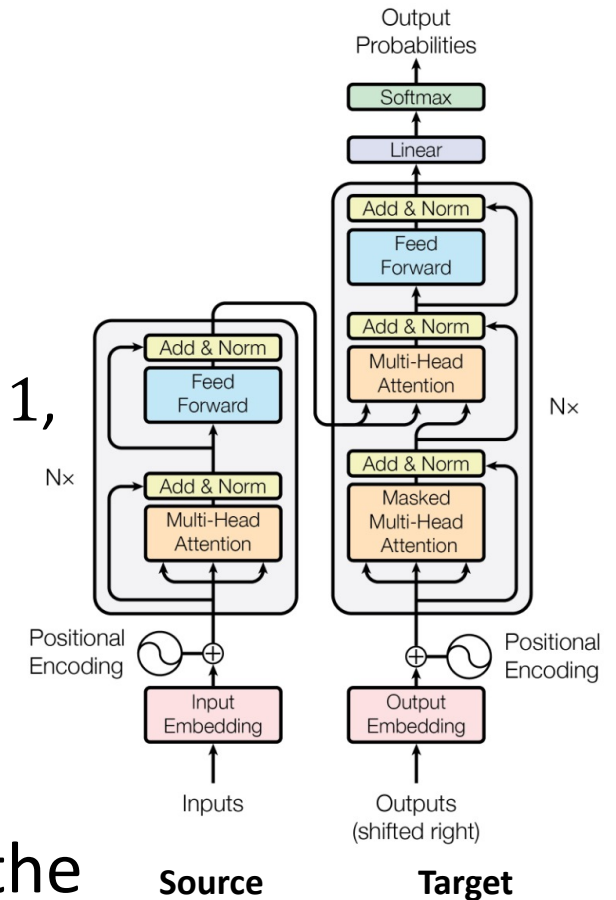
Why Masked Multi-head Attention at Decoder?

- During Training:
 - **Source (English):**
 - ["Transformers", "are", "great"]
 - **Target (German):**
 - ["Transformatoren", "sind", "großartig", "</eos>"]
 - **Decoder input (shifted-right targets):**
 - ["<bos>", "Transformatoren", "sind", "großartig"]
- Encoder is bi-directional (no masking)
 - Each word looks at the past and future words to build context
 - The encoder processes the entire source sequence at once



Why Masked Multi-head Attention at Decoder?

- Decoder generates tokens sequentially (one per step)
- What the decoder sees and predicts?
 - At each target position t , the decoder input is everything up to $t - 1$,
 - Decoder predicts the token at t .
 - $t = 1$: input "<bos>" → predict "Transformatoren"
 - $t = 2$: input "<bos> Transformatoren" → predict "sind"
 - $t = 3$: input "<bos> Transformatoren sind" → predict "großartig"
- We don't want decoder to look at the future words to build the context



Masked Multi-head Attention

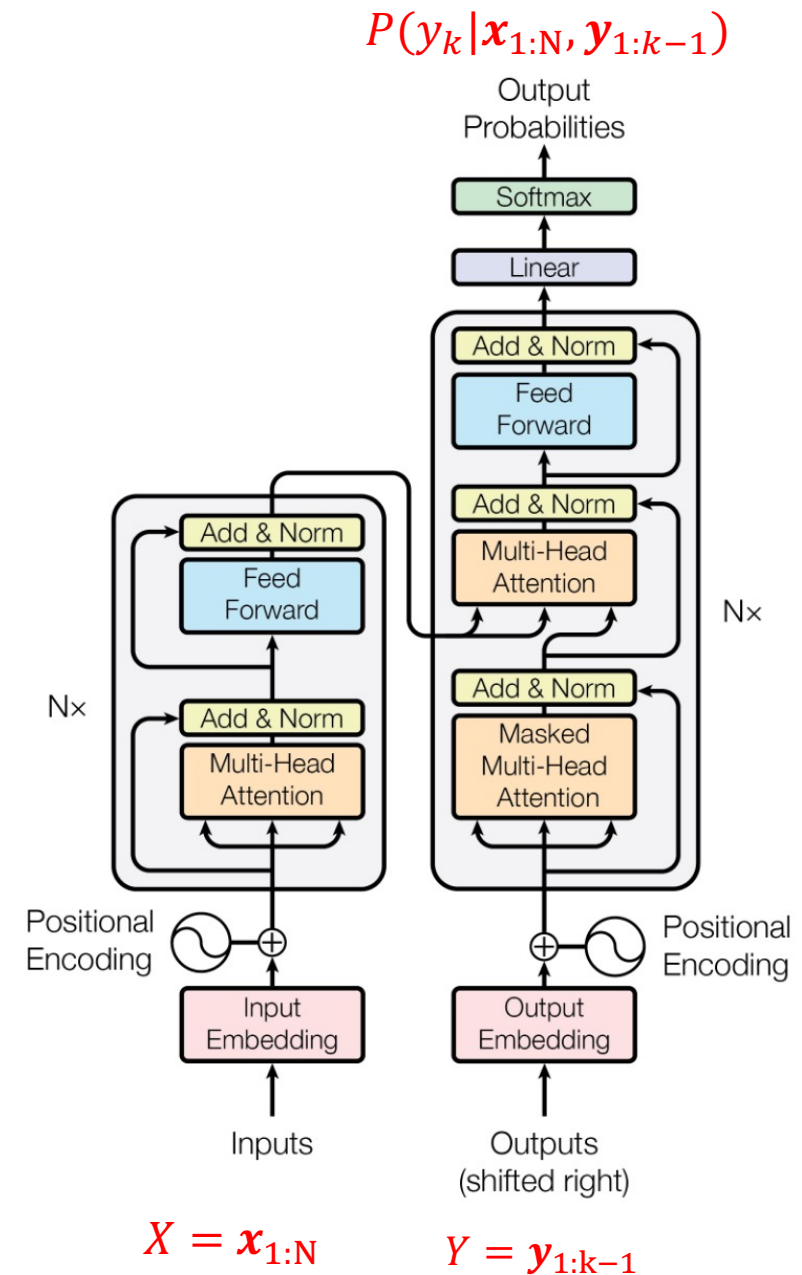
- $head_i = softmax \left(\frac{(QW_i^Q)(KW_i^K)^T}{\sqrt{d_K}} \circ \mathbf{M} \right) (VW_i^V)$

where $M_{ij} = \begin{cases} 1, & \text{if attention between } i \text{ and } j \\ -\infty, & \text{if no attention} \end{cases}$

- Example Mask (M):

$$M = \begin{bmatrix} 1 & -\infty & -\infty & -\infty & -\infty \\ 1 & 1 & -\infty & -\infty & -\infty \\ 1 & 1 & 1 & -\infty & -\infty \\ 1 & 1 & 1 & 1 & -\infty \\ 1 & 1 & 1 & 1 & 1 \end{bmatrix}$$

- Masking ensures causal property is not violated
 - Don't predict current output based on future outputs



$$\text{head}_i = \text{softmax} \left(\frac{(QW_i^Q)(KW_i^K)^T}{\sqrt{d_K}} \circ \mathbf{M} \right) (VW_i^V)$$

$$\frac{(QW_i^Q)(KW_i^K)^T}{\sqrt{d_k}} = \begin{bmatrix} \frac{\mathbf{q}_1 \mathbf{k}_1^T}{\sqrt{d_1}} & -\infty & -\infty \\ \frac{\mathbf{q}_2 \mathbf{k}_1^T}{\sqrt{d_2}} & \frac{\mathbf{q}_2 \mathbf{k}_2^T}{\sqrt{d_2}} & -\infty \\ \frac{\mathbf{q}_3 \mathbf{k}_1^T}{\sqrt{d_3}} & \frac{\mathbf{q}_3 \mathbf{k}_2^T}{\sqrt{d_3}} & \frac{\mathbf{q}_3 \mathbf{k}_3^T}{\sqrt{d_3}} \end{bmatrix}$$

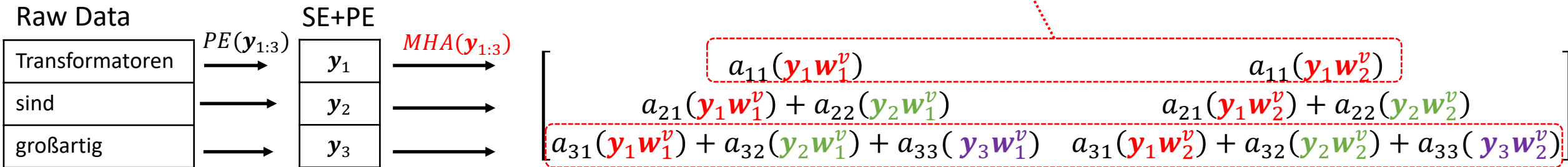
$$\text{softmax} \left(\frac{(QW_i^Q)(KW_i^K)^T}{\sqrt{d_k}} \right) = \begin{bmatrix} \text{softmax} \left(\begin{bmatrix} \frac{\mathbf{q}_1 \mathbf{k}_1^T}{\sqrt{d_1}} & -\infty & -\infty \end{bmatrix} \right) \\ \text{softmax} \left(\begin{bmatrix} \frac{\mathbf{q}_2 \mathbf{k}_1^T}{\sqrt{d_2}} & \frac{\mathbf{q}_2 \mathbf{k}_2^T}{\sqrt{d_2}} & -\infty \end{bmatrix} \right) \\ \text{softmax} \left(\begin{bmatrix} \frac{\mathbf{q}_3 \mathbf{k}_1^T}{\sqrt{d_3}} & \frac{\mathbf{q}_3 \mathbf{k}_2^T}{\sqrt{d_3}} & \frac{\mathbf{q}_3 \mathbf{k}_3^T}{\sqrt{d_3}} \end{bmatrix} \right) \end{bmatrix} = \begin{bmatrix} a_{11} & 0 & 0 \\ a_{21} & a_{22} & 0 \\ a_{31} & a_{32} & a_{33} \end{bmatrix}$$

	Transformatoren	sind	großartig
Transformatoren	a_{11}	0	0
sind	a_{21}	a_{22}	0
großartig	a_{31}	a_{32}	a_{33}

$$\text{softmax}\left(\frac{(QW_i^Q)(KW_i^K)^T}{\sqrt{d_k}}\right) \begin{bmatrix} \mathbf{y}_1 \mathbf{w}_1^v & \mathbf{y}_1 \mathbf{w}_2^v \\ \mathbf{y}_2 \mathbf{w}_1^v & \mathbf{y}_2 \mathbf{w}_2^v \\ \mathbf{y}_3 \mathbf{w}_1^v & \mathbf{y}_3 \mathbf{w}_2^v \end{bmatrix} = \begin{bmatrix} a_{11} & 0 & 0 \\ a_{21} & a_{22} & 0 \\ a_{31} & a_{32} & a_{33} \end{bmatrix} \begin{bmatrix} \mathbf{y}_1 \mathbf{w}_1^v & \mathbf{y}_1 \mathbf{w}_2^v \\ \mathbf{y}_2 \mathbf{w}_1^v & \mathbf{y}_2 \mathbf{w}_2^v \\ \mathbf{y}_3 \mathbf{w}_1^v & \mathbf{y}_3 \mathbf{w}_2^v \end{bmatrix}$$

Attention Scores

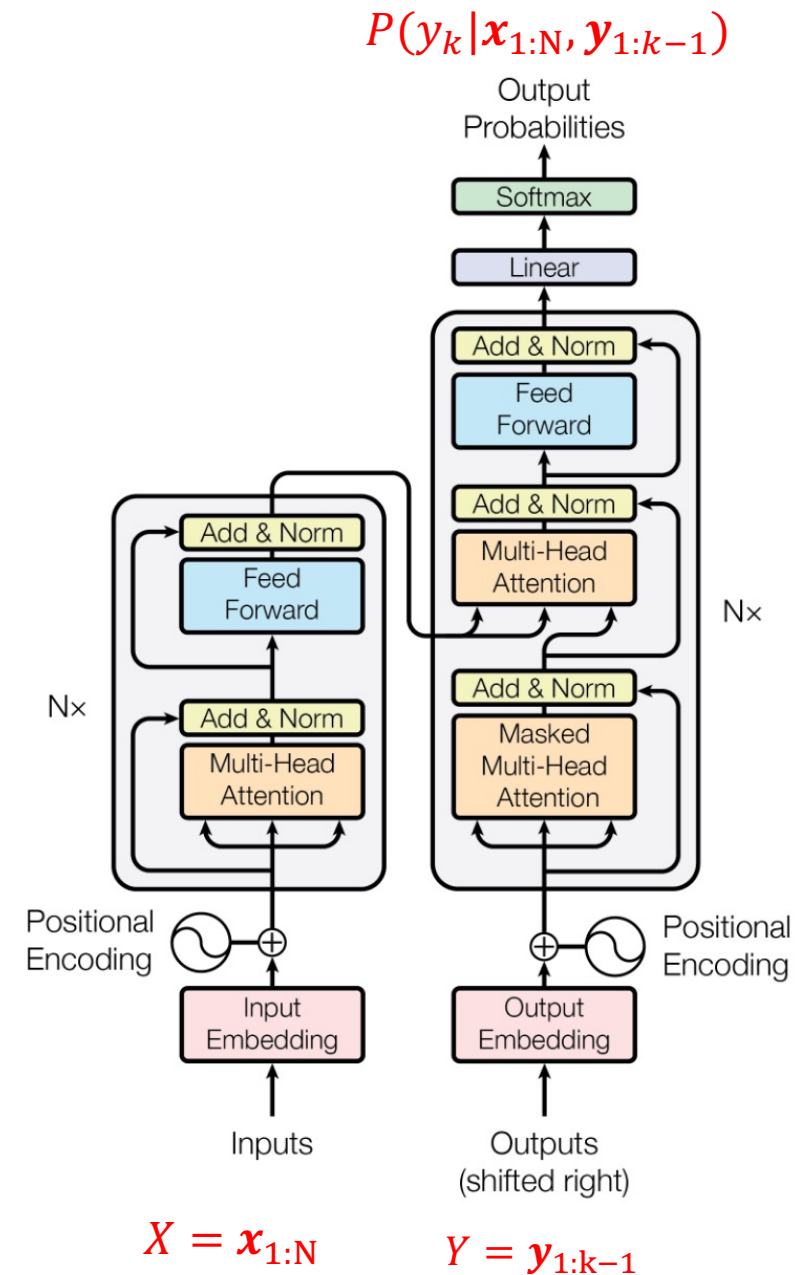
Embedding of “Transformatoren” does not depend on “sind”, and “großartig”



Embedding of “großartig” is a linear combination of embeddings of “Transformatoren”, “sind”, and “großartig”

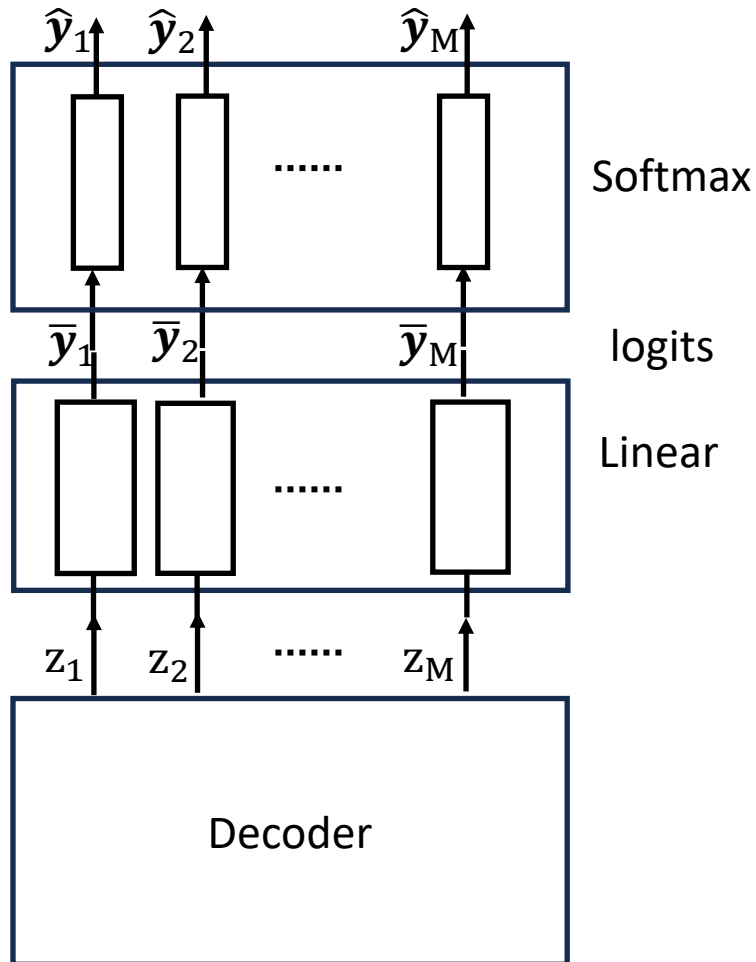
Cross-Attention Layer

- Same as self-attention layer, except that
 - K and V are from encoder
 - Q is from output sequence
- Why no masking?
 - Decoding starts only after entire encoder input is processed
 - Decoder's queries attend to the encoder's keys/values
 - Attending to any source position does not leak future target information
 - Therefore, no causal mask is needed in cross-attention.
- Streaming or online tasks
 - Eg: Live translation, etc.
 - We may need masking



Output Layer

- Let $Z \in R^{M \times d}$ be the output of decoder, where M is the output sequence length and d is the embedding dimension
 - $Z = [z_1 \ z_2 \ \dots \ z_M]$, where $z \in R^{d \times 1}$



$$\hat{y}_i = \text{softmax}(\bar{y}_i)$$

$$\hat{y}_i \in R^{|V| \times 1}$$

$$\bar{y}_i \in R^{|V| \times 1}$$

$$\bar{y}_i = (W^0)^T z_i$$

$$W^0 \in R^{d \times |V|}$$

