

# Language Modeling: RNNs, Attention, Transformers

CSE 5525: Foundations of Speech and Language Processing  
<https://shocheen.github.io/cse-5525-spring-2026/>



**THE OHIO STATE UNIVERSITY**

---

Sachin Kumar (kumar.1145@osu.edu)

# Logistics

- Homework 2 is due date next week.
  - Any thoughts, questions, concerns?
  
- Final project: have you formed teams already?
  - Everyone has access to Tinker? Did you get the credits?
  - A project proposal will be due Feb 20.
  - Will post sample project proposals and final reports on teams later today.

# Recap

- What is a language model?
  - $P(\text{text})$
  - Probability distribution over a sequence (or words or “tokens”).
  - Apply chain rule to model  $p(\text{current word} \mid \text{prefix})$
- How to evaluate a language model?
  - Perplexity
  - Other extrinsic metrics.
- N-gram language models
  - Markov assumption – current word’s probability only depends on previous  $n-1$  words
  - Compute probability by counting n-gram frequencies
  - Smoothing to avoid 0 probabilities.
- Feedforward Neural Language Model
  - Still n-gram but make it neural – word embeddings + MLP + classification (over the vocabulary).
  - Need to make unreasonable assumptions and lose information from the long context.

We want to model

$$P(X_1, \dots, X_t)$$

Apply chain rule

$$P(X_1)P(X_2|X_1) \dots P(X_t|X_1, \dots, X_{t-1})$$

## A variant conditional Language Model

$$P(X_1, \dots, X_N \mid Y_1, \dots, Y_M)$$

additional input

Conditional Language Model

# LMs w/ Recurrent Neural Nets

- Core idea: apply **a model repeatedly**

outputs { **output distribution**  
 $\hat{y}^{(t)} = \text{softmax}(U h^{(t)} + b_2) \in \mathbb{R}^{|V|}$

hidden states

{  $h^{(t)} = \sigma(W_h h^{(t-1)} + W_e e^{(t)} + b_1)$   
 $h^{(0)}$  is the initial hidden state

Input

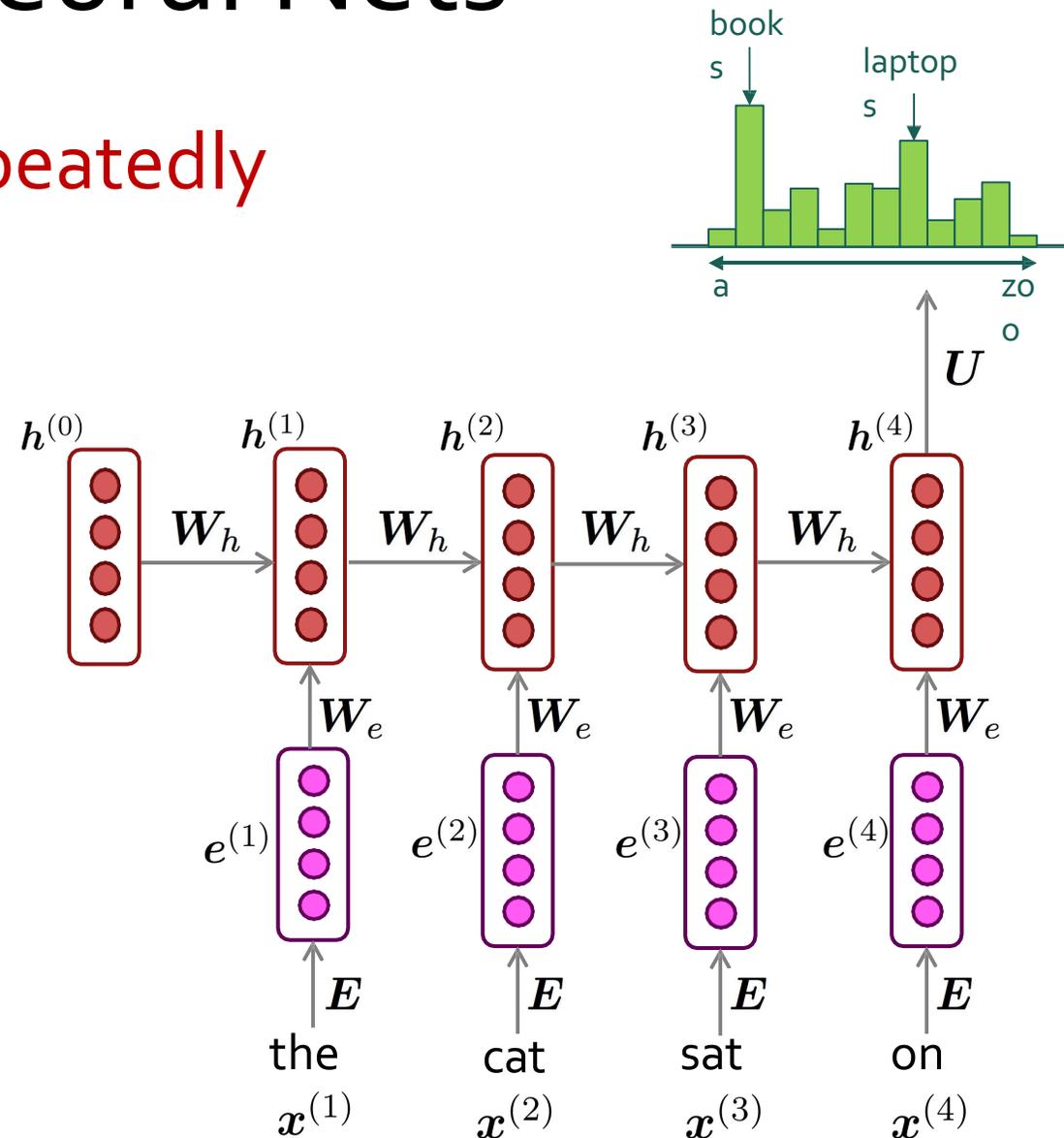
embedding

{ **word embeddings**

$$e^{(t)} = E x^{(t)}$$

words / one-hot vectors

$$x^{(t)} \in \mathbb{R}^{|V|}$$



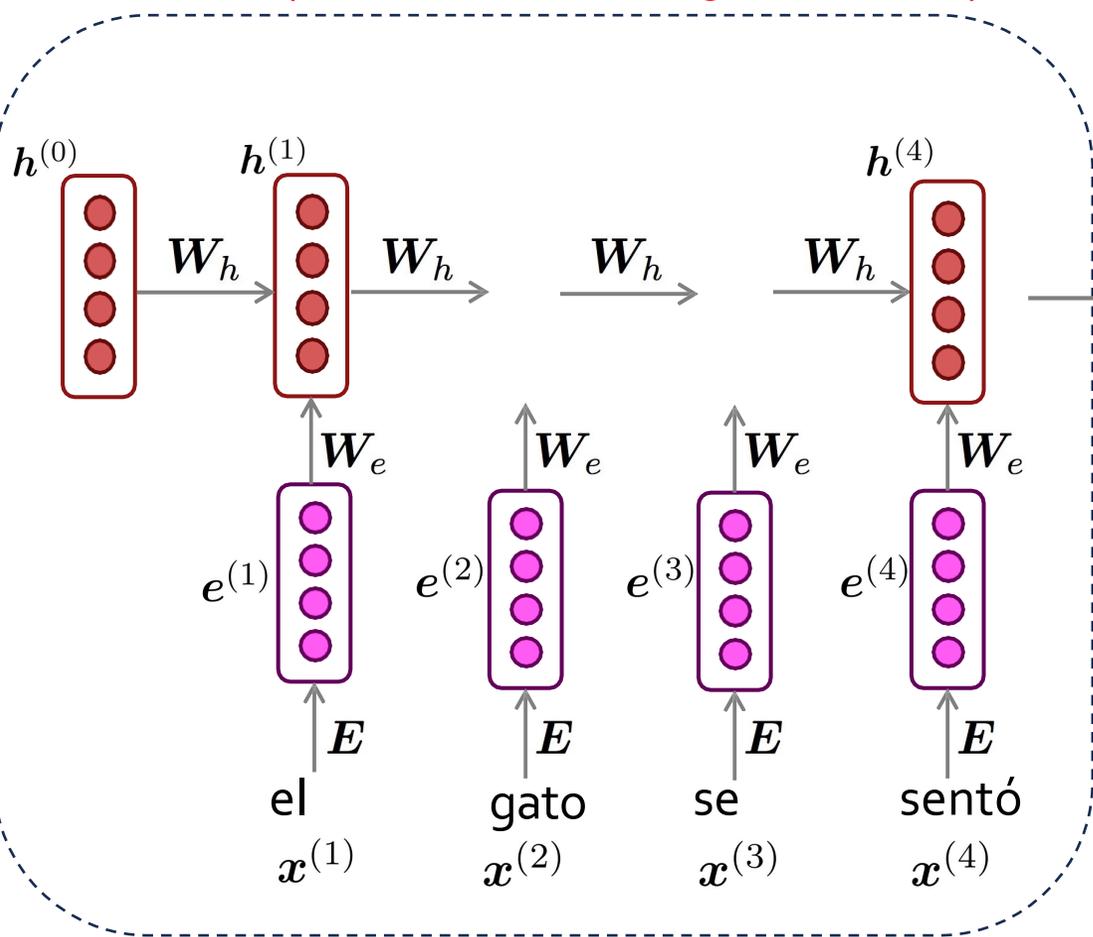
# Recurrent Neural Networks

- Applied to sequential data iteratively.
  - $h_t = f(h_{t-1}, x_t; \theta)$
  - there are many ways to define  $f$  (we will only talk about simple RNNs)
  - Note this  $\theta$  is shared across all the items in the sequence
- Why RNNs
  - They allow modeling infinite context (in theory)
  - They can retain sequential information as opposed to bag of words models
- Intuitively, at every hidden state, the model encodes all the necessary information required to predict the next token at that position
  - At least that's the hope

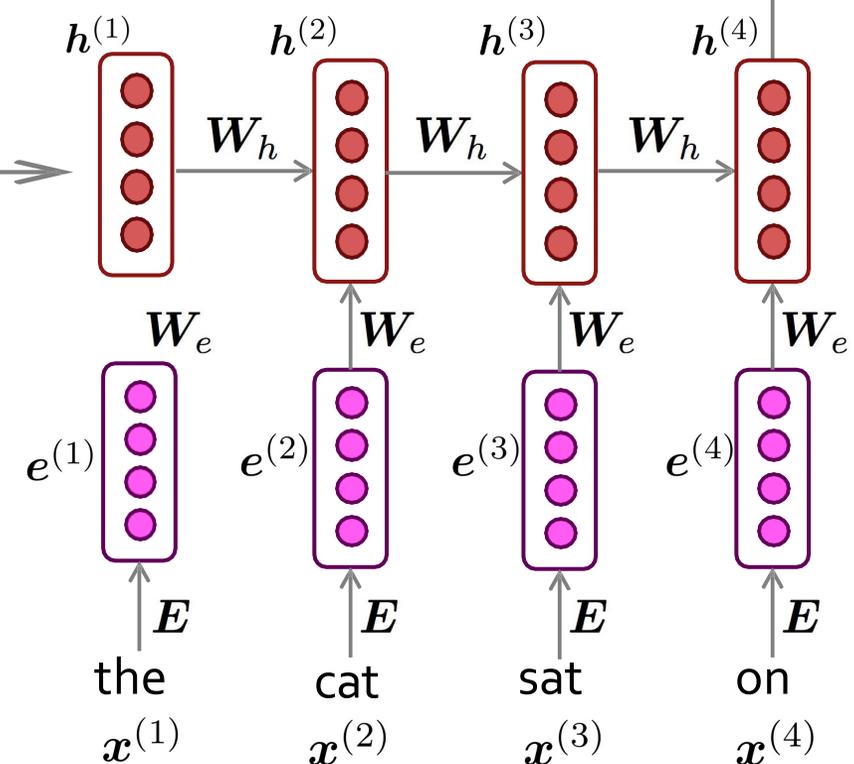
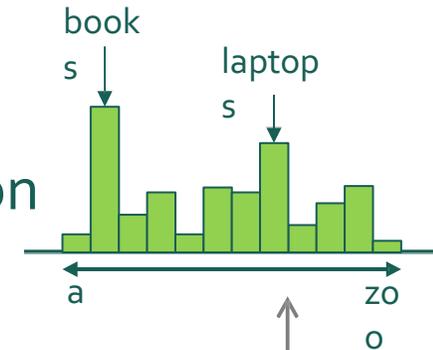
# Conditional LMs with RNNs

Two RNNs – encoder and decoder

Encoder (different set of weight matrices)



output distribution



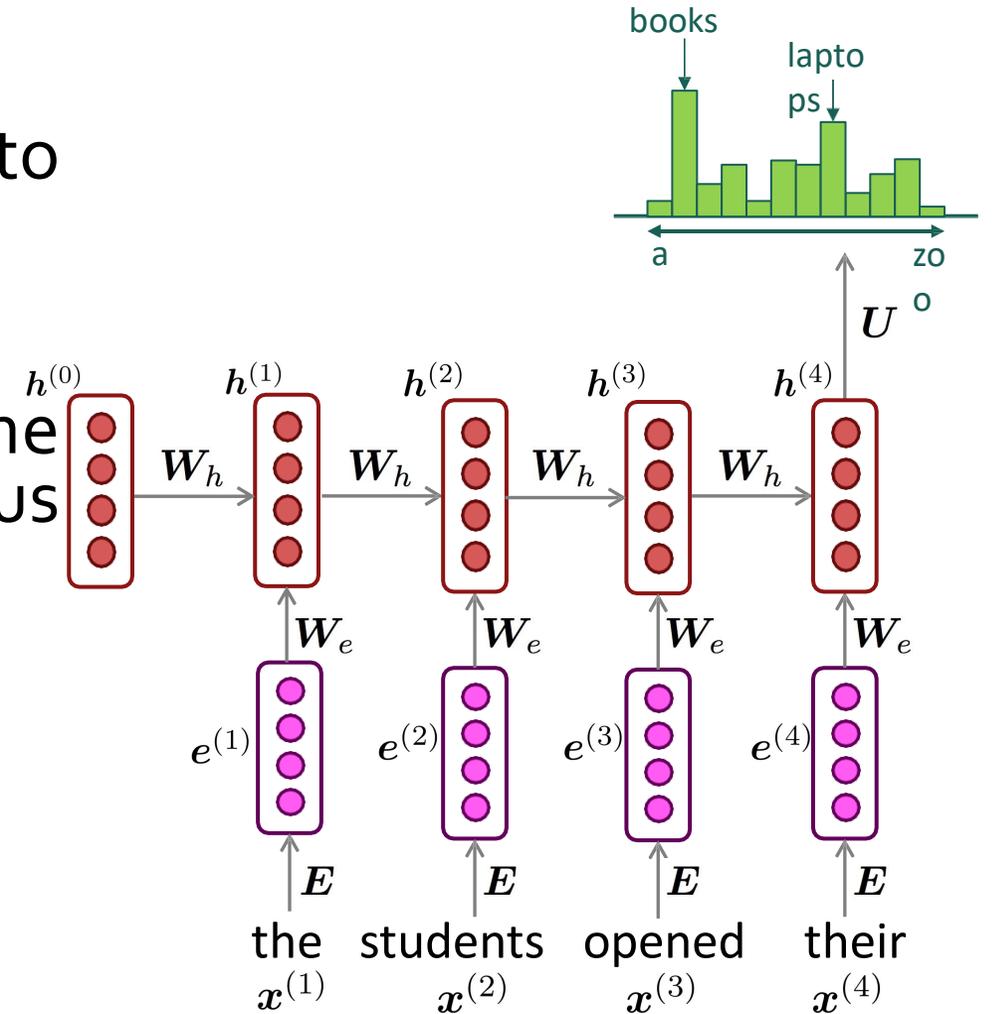
Decoder

# How to train RNNs?

- Using our favorite algorithm: gradient descent.
- Loss: classification loss at every step  $i$  (using cross-entropy loss)
- But backpropagation is applied over and over to the same parameters  $\theta$ 
  - Also known as backpropagation through time (BPTT)
- Issues with RNNs
  - Gradients can explode or vanish.
  - Solution: modify optimization algorithms / architectures (e.g. LSTMs) [won't discuss in this course, look at readings]

# Other issues with RNNs

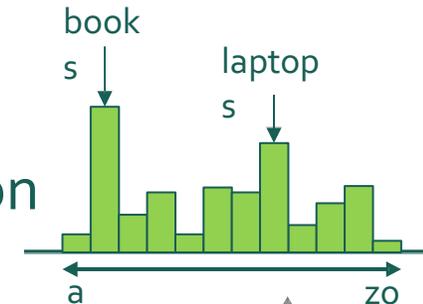
- Recurrent computation is **slow**, difficult to parallelize.
- Each hidden state is expected to store the entire information from the previous context
  - Is it even possible?



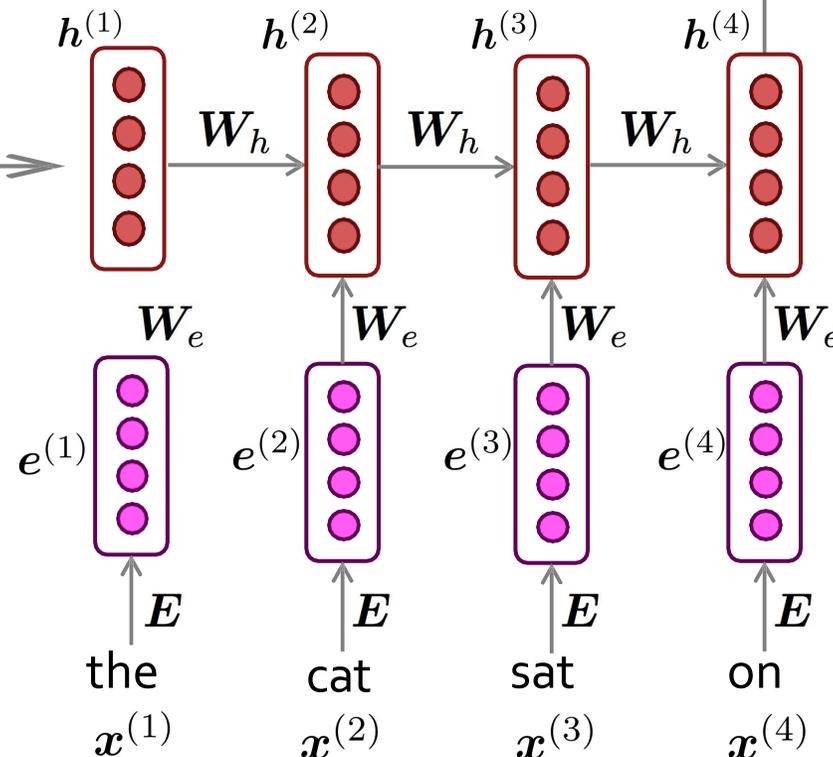
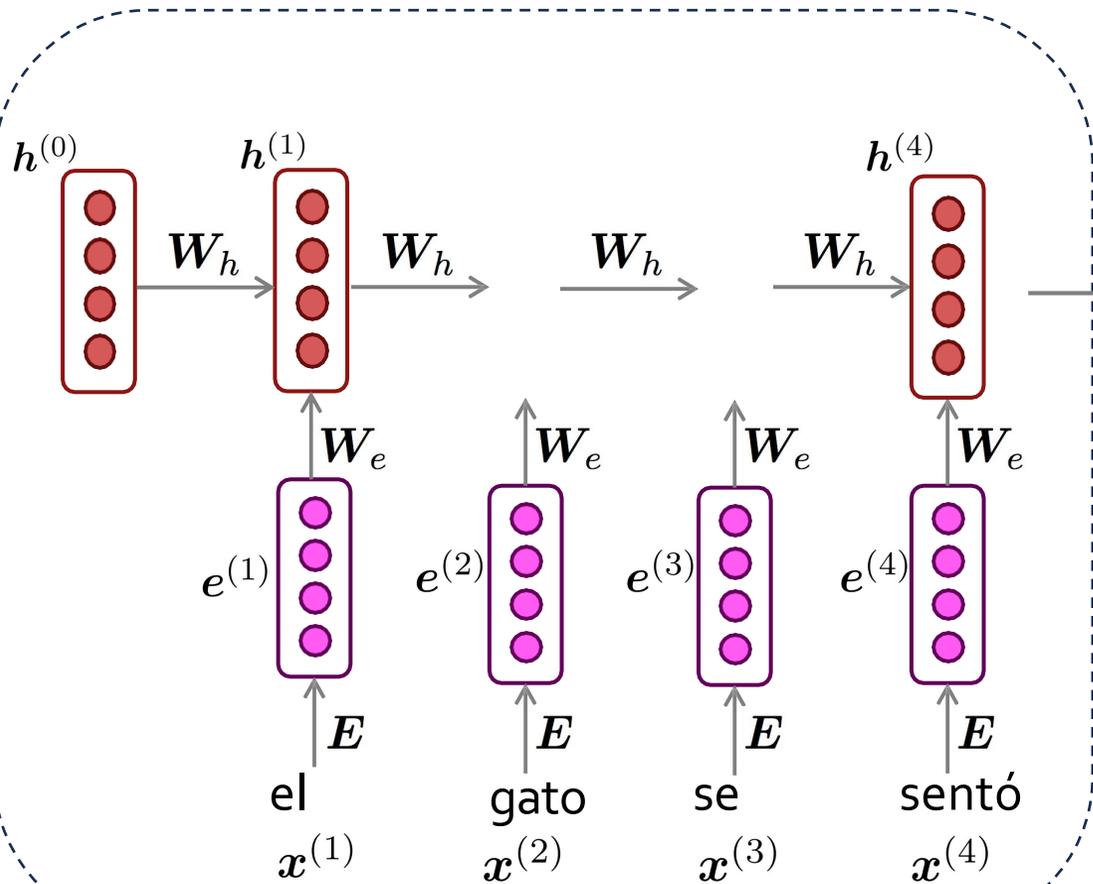
# Machine Translation with RNNs

Read the source only once, generate translation from memory

output distribution



Encoder

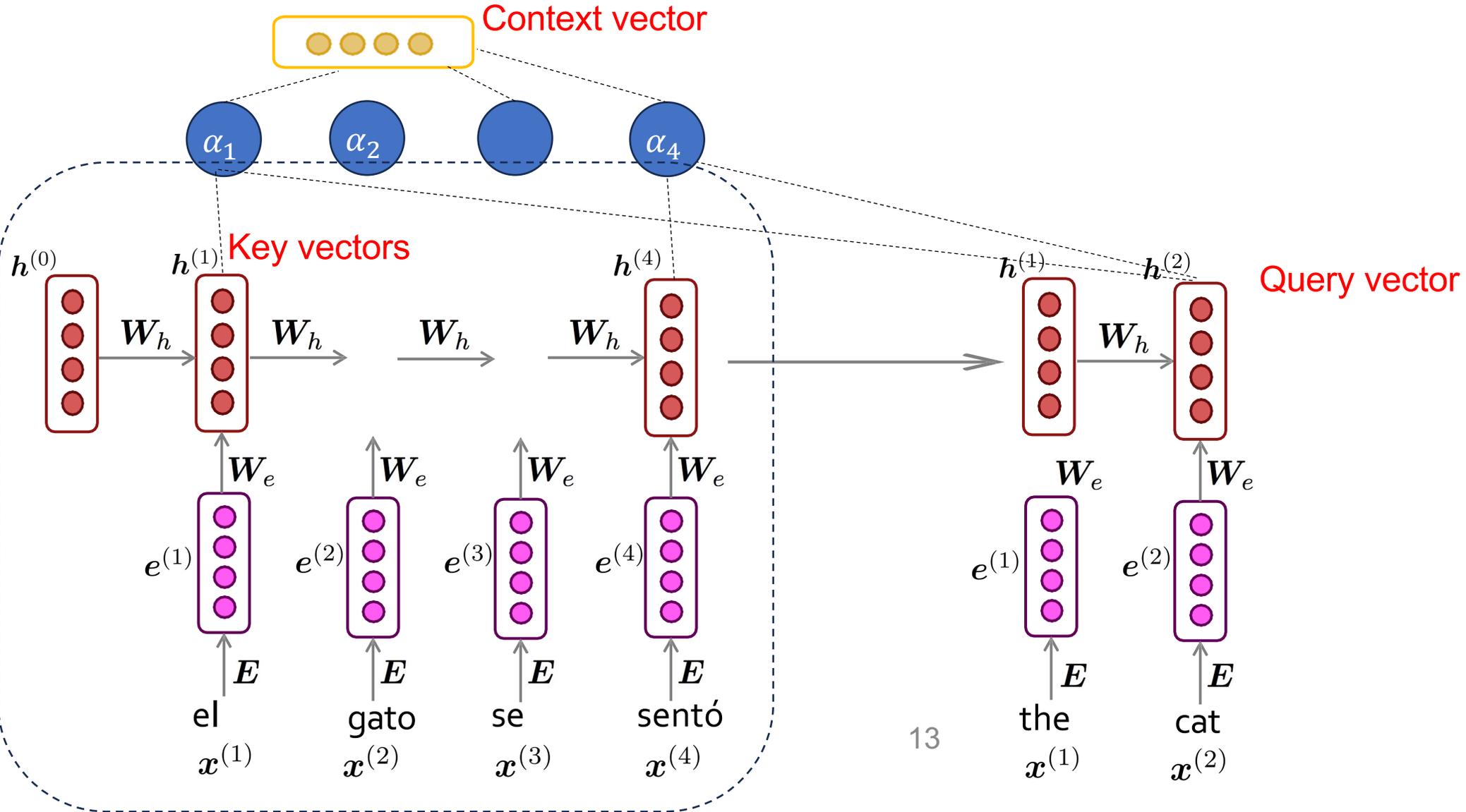


Decoder

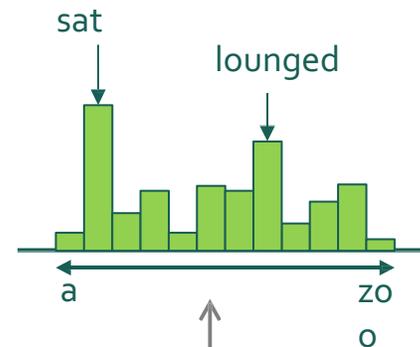
# Solution: Attention

- What if the decoder at each step pays “attention” to a distribution of all of encoder’s hidden states?
- Intuition: when we (humans) translate a sentence, we don’t just consume the original sentence then regurgitate in a new language; we continuously look back at the original while focusing on different parts

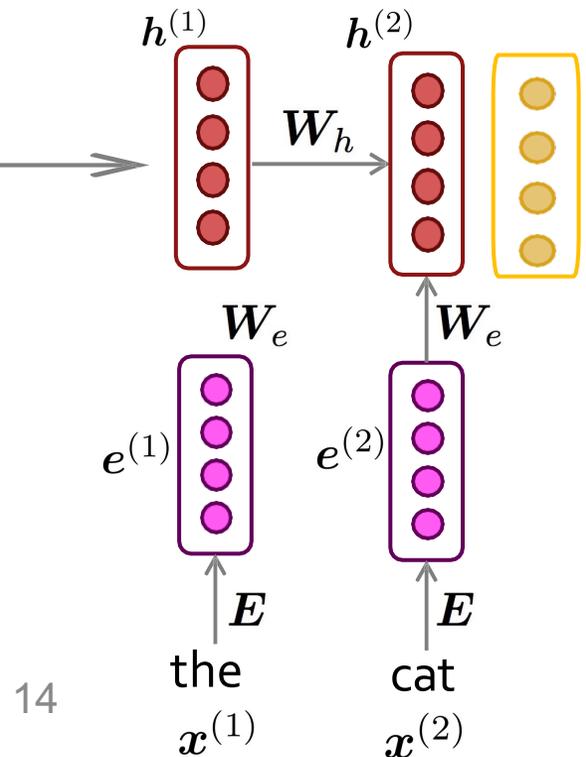
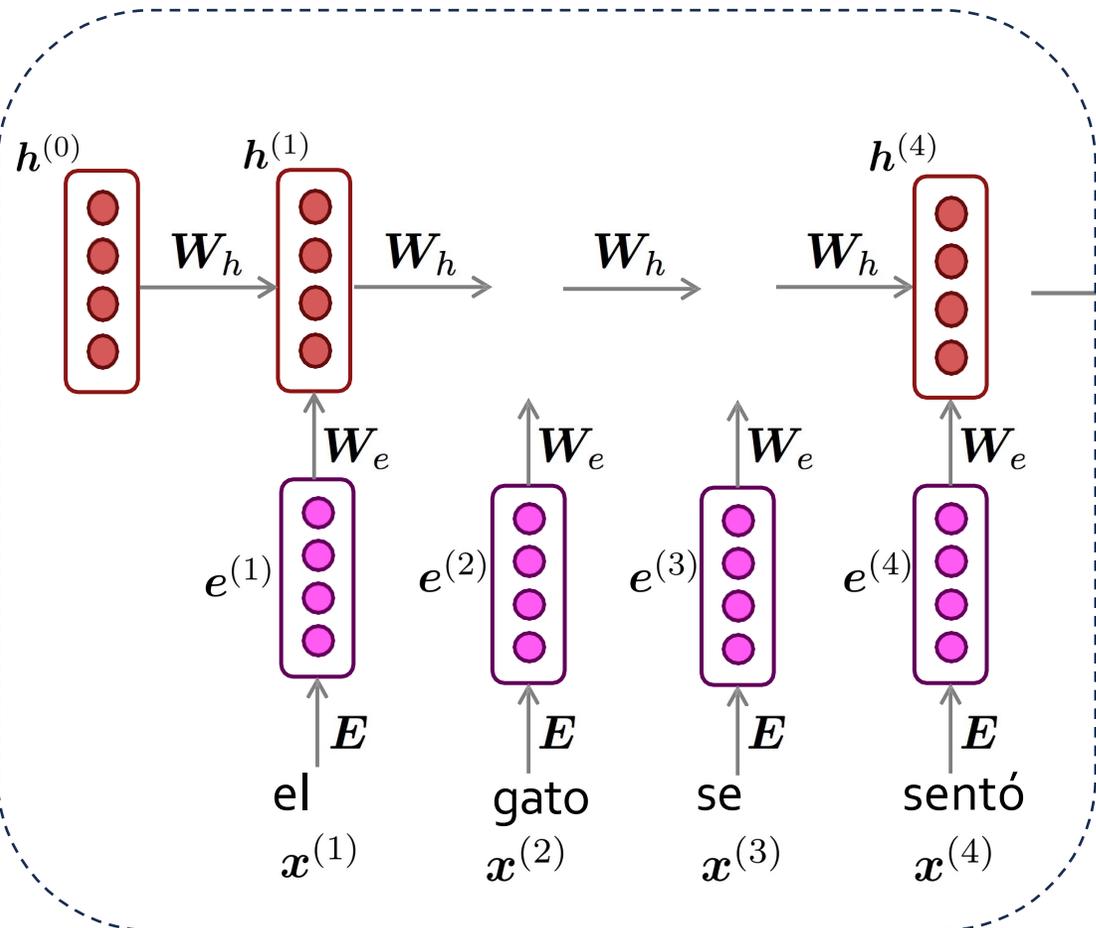
# RNNs with Attention



# RNNs with Attention



$[h^{(2)}, c^{(2)}]$



14

14  
14

# RNNs with Attention

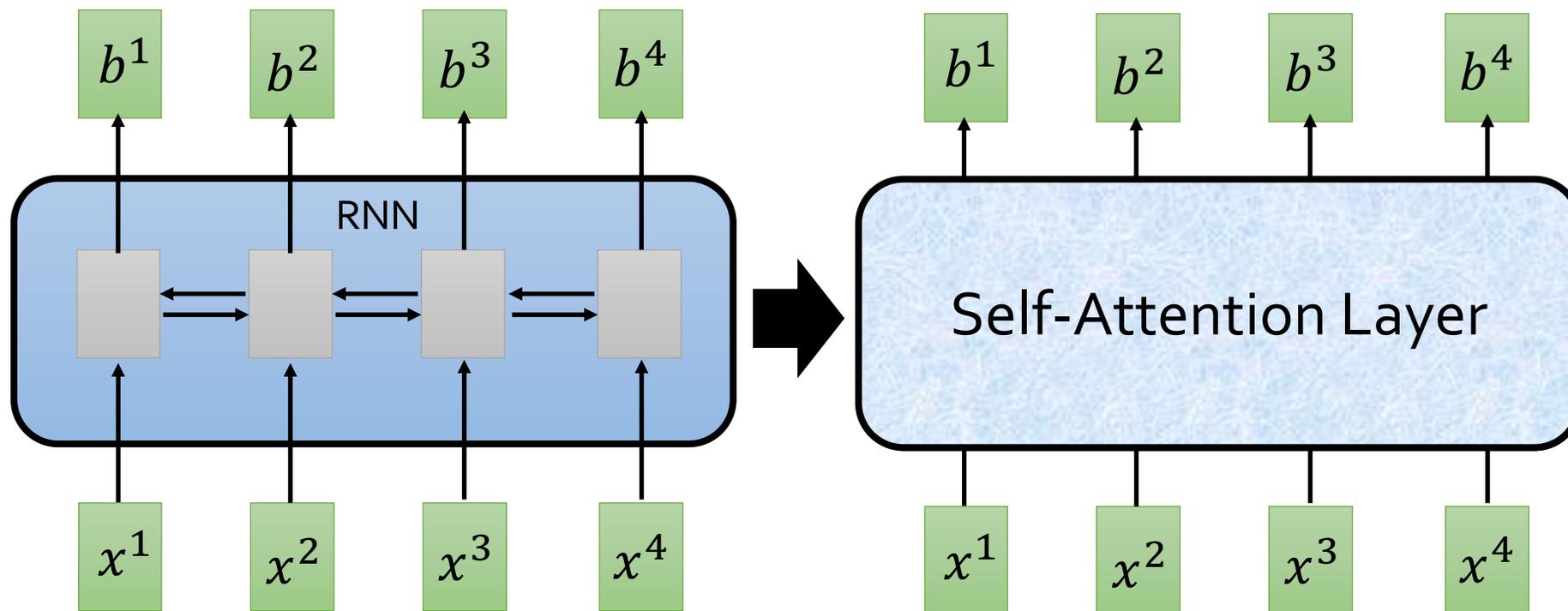
- Attention allowed modelling longer context and obtain higher performance
- But
  - It is still slow because of linear computation in RNN
  - It still has gradient vanishing/exploding issues (RNN variants can resolve this)
- Solution: what if we removed the RNN component and only use attention
  - Attention is all you need (Vaswani et al 2017)

# Transformers

- Replace the linear part of RNNs with **self-attention**
- Introduce **residual connections + layernorm** to improve gradient flow (avoid gradient vanishing issues)
- Introduce **positional embeddings** to encode sequential order

# Self-Attention

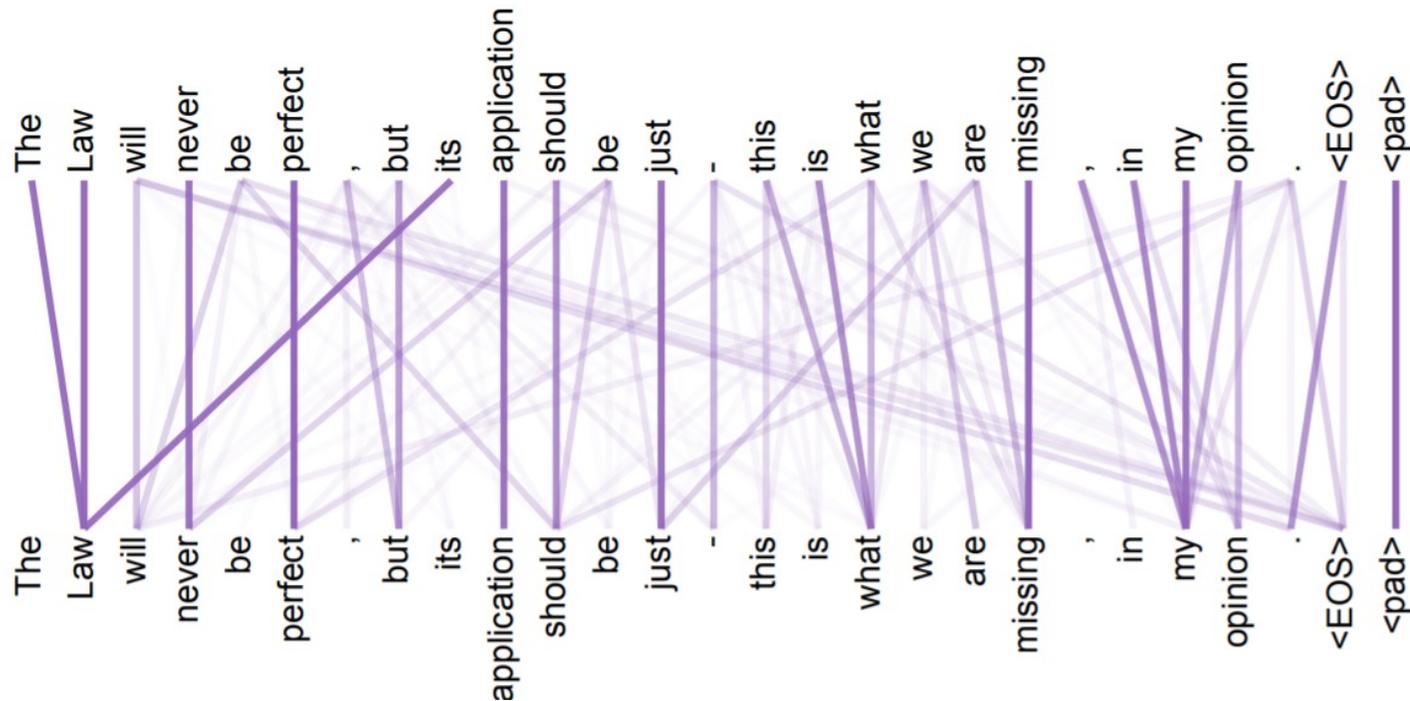
- $b^t$  is obtained based on the whole input sequence.
- can be parallelly computed.



Idea: replace any thing done by RNN with **self-attention**.

# Attention

- Core idea: on each step, *use direct connection to focus (“attend”) on a particular part* of the context
  - Kind of similar to deep averaging networks but a “weighted average”



# Outline

-  Self-Attention
-  Transformer Encoder
-  Transformer Decoder
-  Language Modeling With  
Transformers

# Outline

-  Self-Attention
-  Transformer Encoder
-  Transformer Decoder
-  Language Modeling With  
Transformers

# Defining Self-Attention

- **Terminology:**
  - **Query:** to match others
  - **Key:** to be matched
  - **Value:** information to be extracted
- **Definition:** Given a set of vector **values**, and a vector **query**, *attention* is a technique to compute a weighted sum of the **value**, dependent on the **keys**. [In RNNs, the key and value vectors were the same]

$q$ : query (to match others)

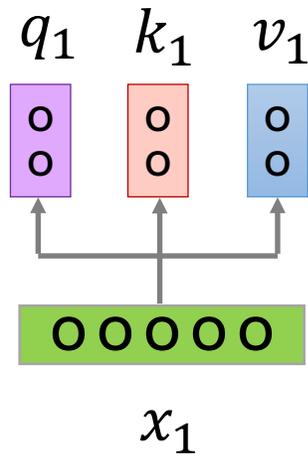
$$q_t = W^q x_t$$

$k$ : key (to be matched)

$$k_t = W^k x_t$$

$v$ : value (information to be extracted)

$$v_t = W^v x_t$$



The

$q$ : query (to match others)

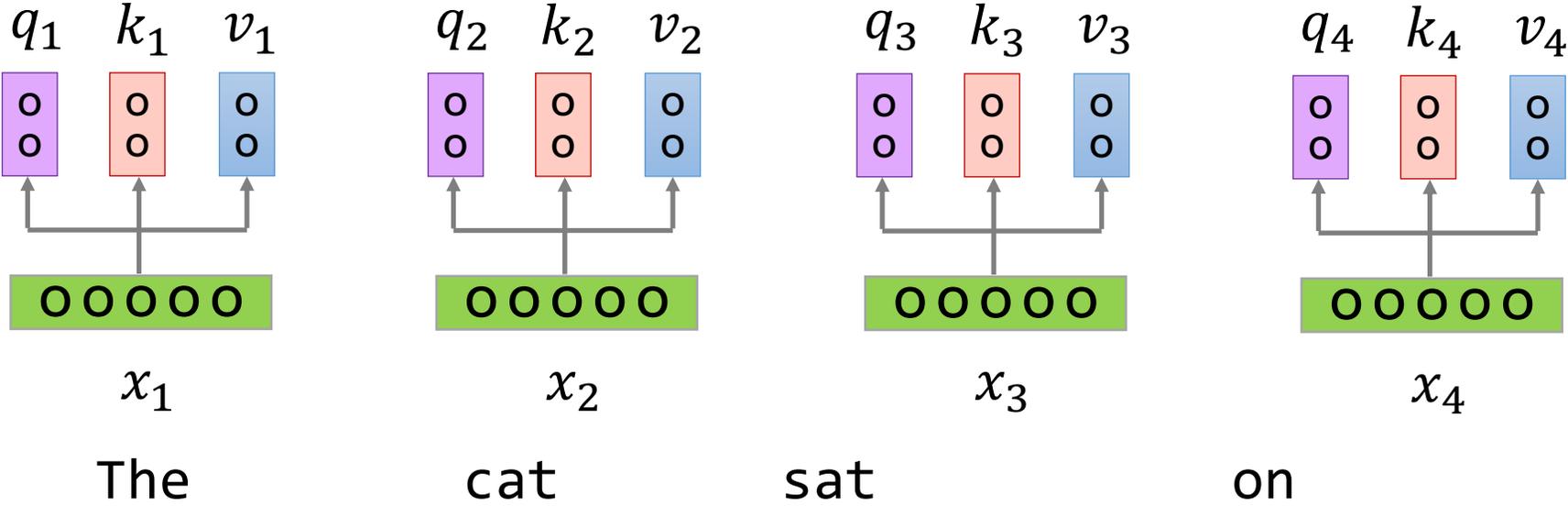
$$q_t = W^q x_t$$

$k$ : key (to be matched)

$$k_t = W^k x_t$$

$v$ : value (information to be extracted)

$$v_t = W^v x_t$$



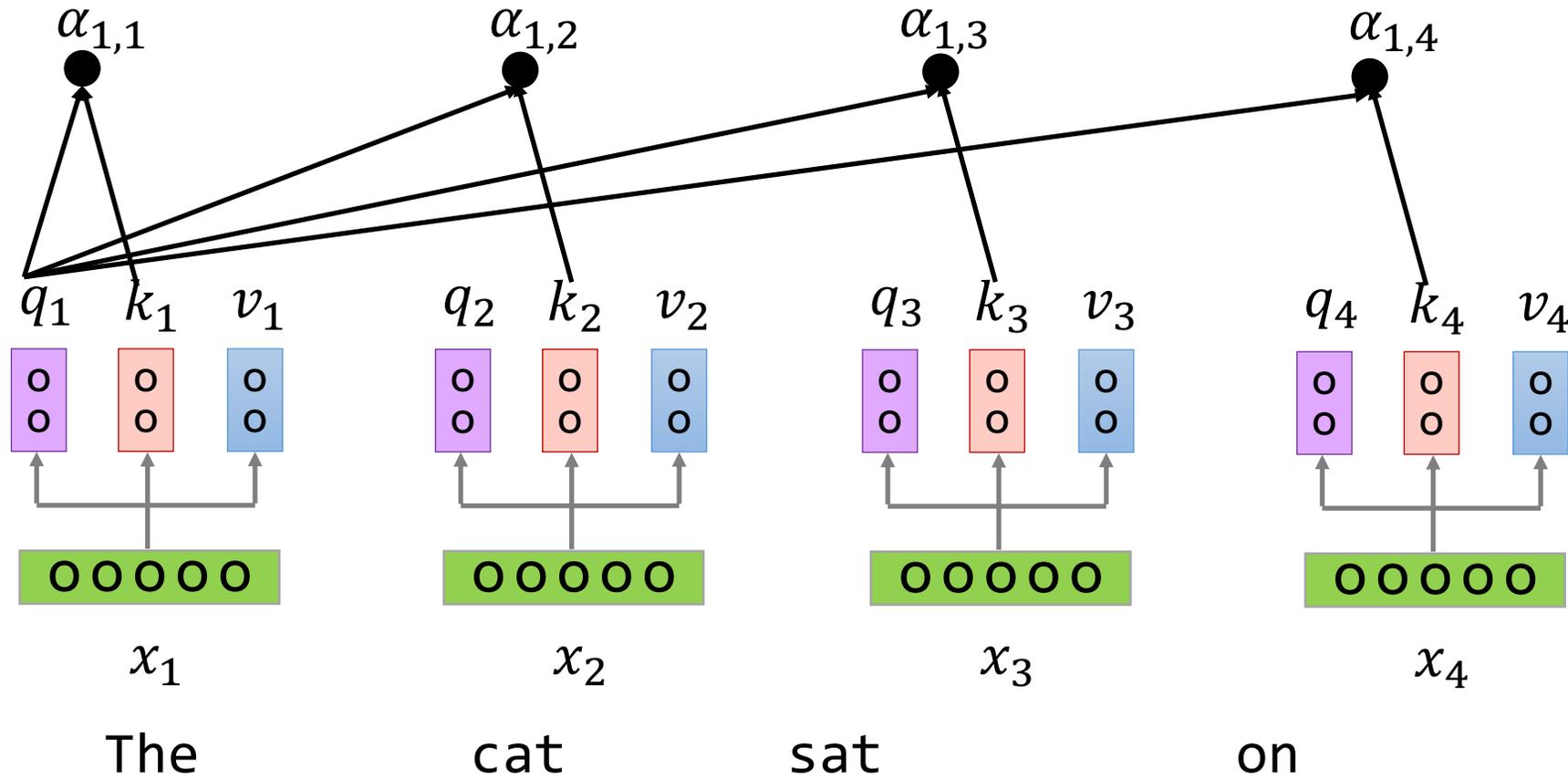
$$\alpha_{1,t} = \underbrace{q^1 \cdot k^t / \alpha}_{\text{Scaled dot product}}$$

$q$ : query (to match others)

$k$ : key (to be matched)

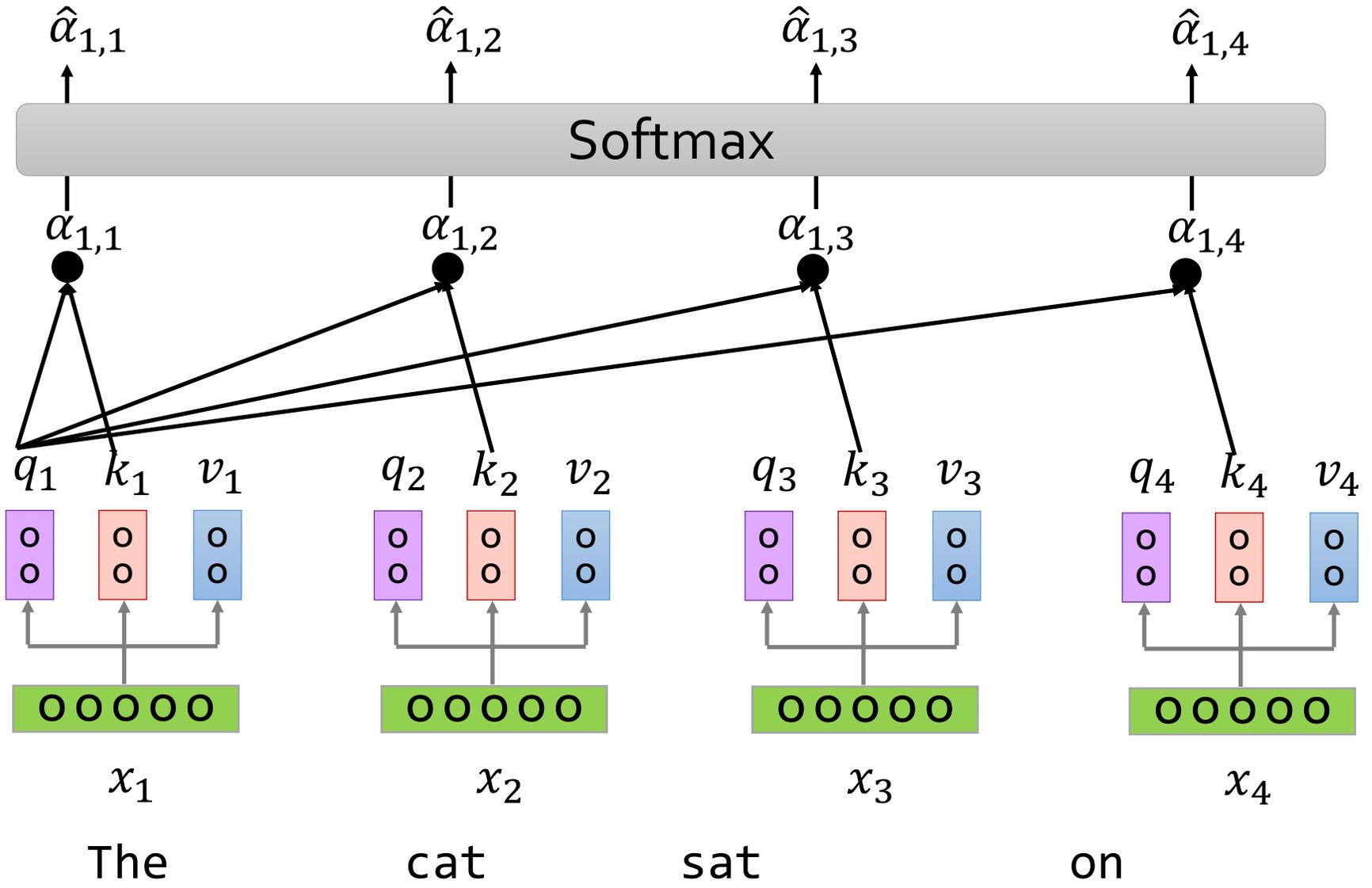
$v$ : value (information to be extracted)

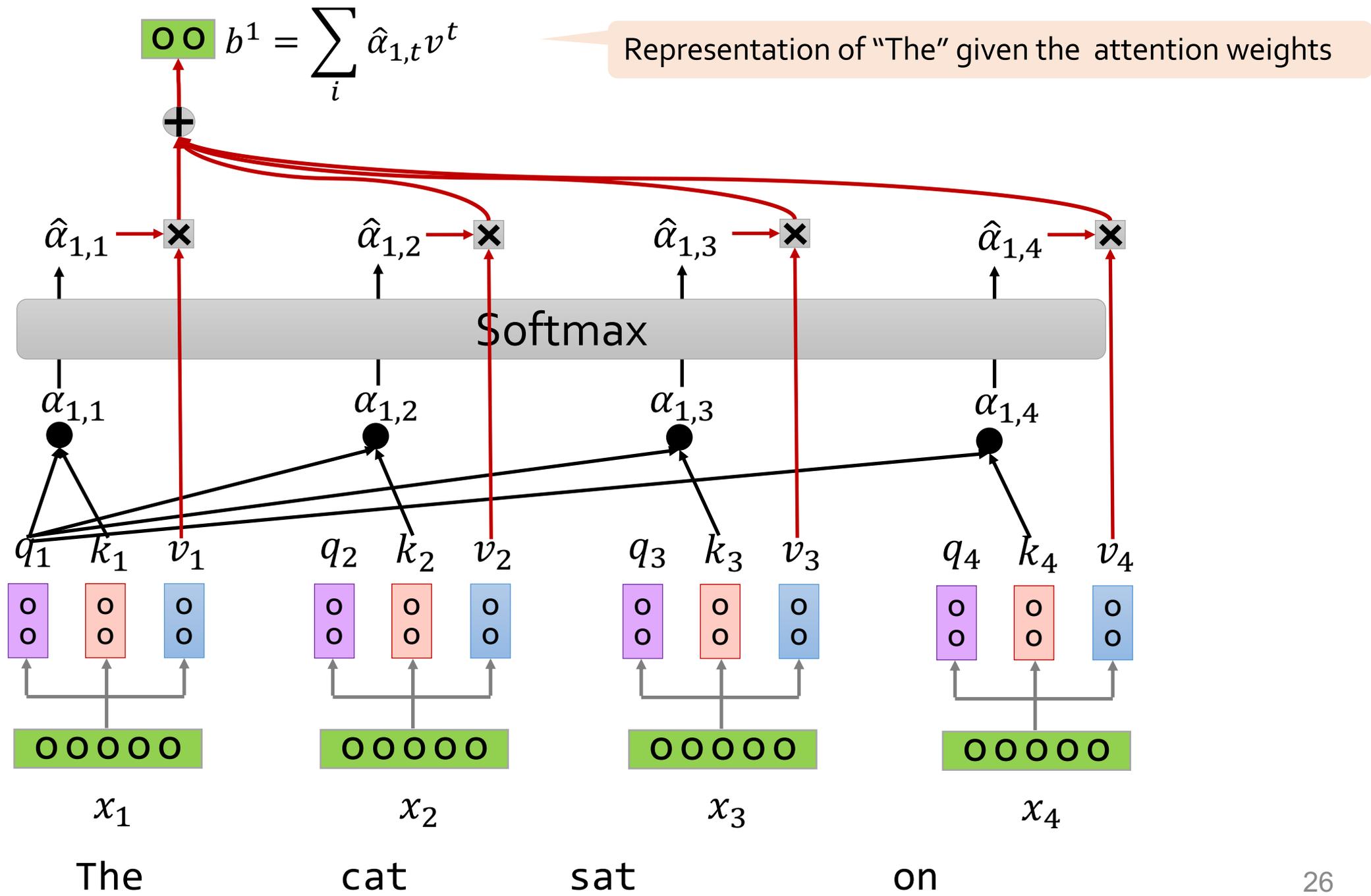
How much should "The" attend to other positions?

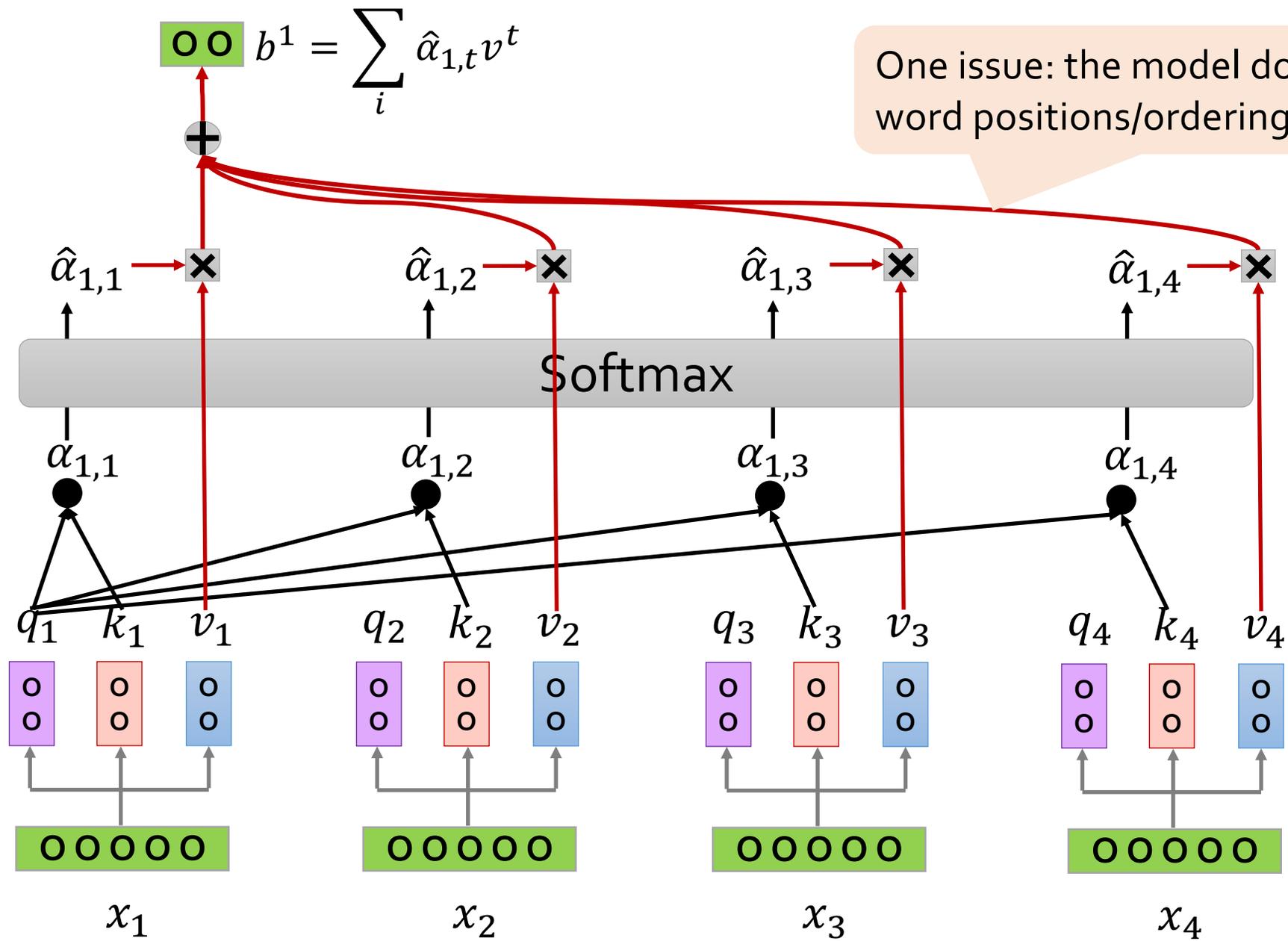


$$\sigma(z)_t = \frac{\exp(z_t)}{\sum_j \exp(z_j)}$$

How much should "The" attend to other positions?







One issue: the model doesn't know word positions/ordering.

The

cat

sat

on

# Self-Attention

**Step 1:** Our Self-Attention Head I has just 3 weight matrices  $W_q$ ,  $W_k$ ,  $W_v$  in total. **These same 3 weight matrices** are multiplied by each  $x_i$  to create all vectors:

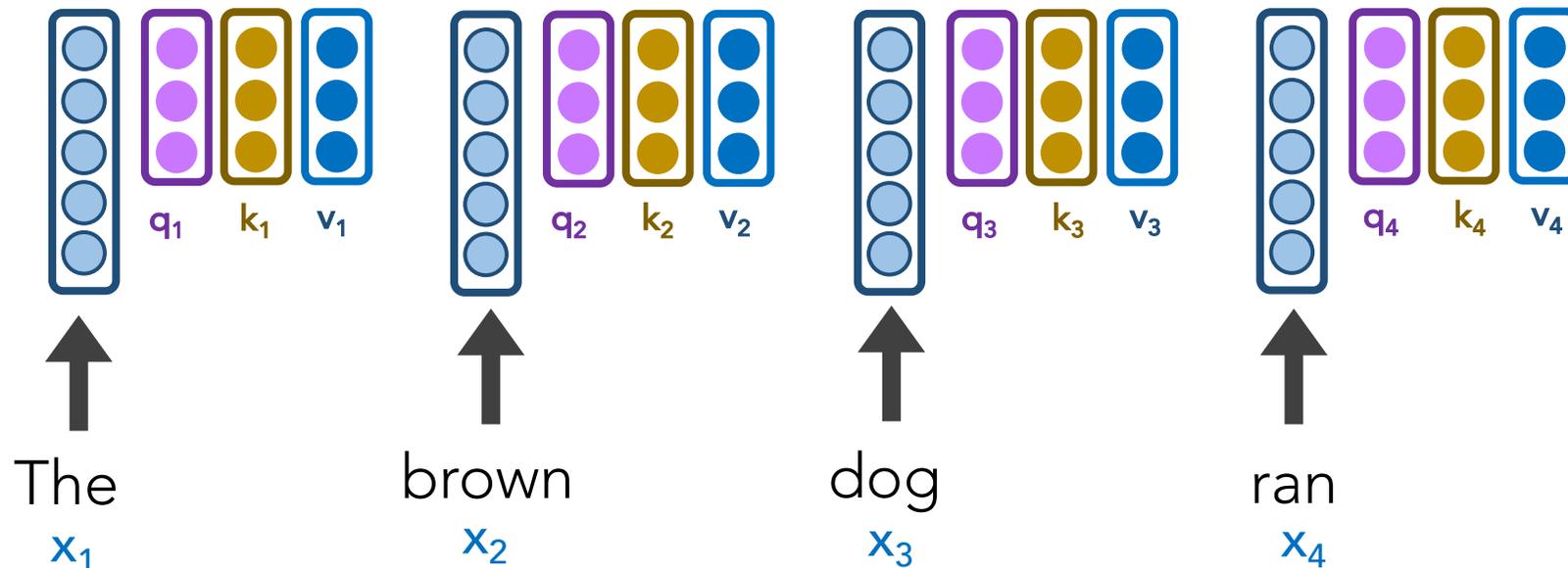
$$q_i = W_q x_i$$

$$k_i = W_k x_i$$

$$v_i = W_v x_i$$

Under the hood, each  $x_i$  has 3 small, associated vectors. For example,  $x_1$  has:

- Query  $q_1$
- Key  $k_1$
- Value  $v_1$

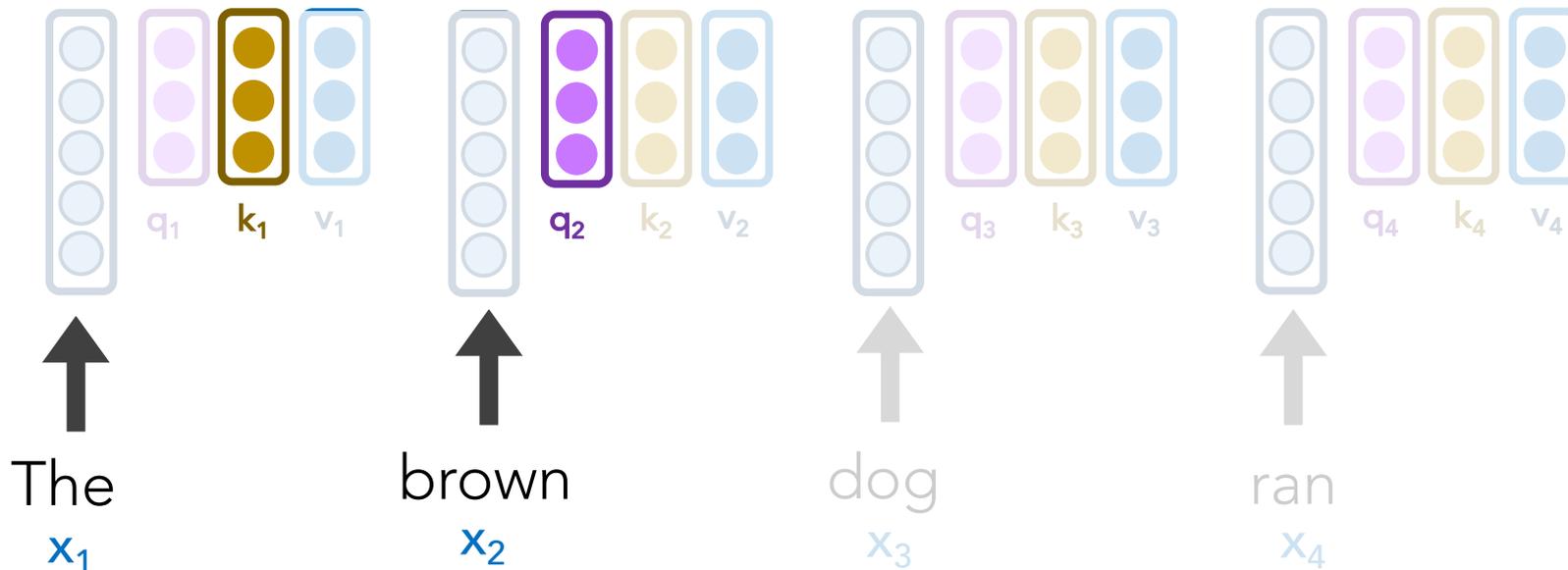


# Self-Attention

---

**Step 2:** For word  $x_2$ , let's calculate the scores  $s_1, s_2, s_3, s_4$ , which represent how much attention to pay to each respective "word"  $v_i$

$$s_1 = q_2 \cdot k_1 = 92$$

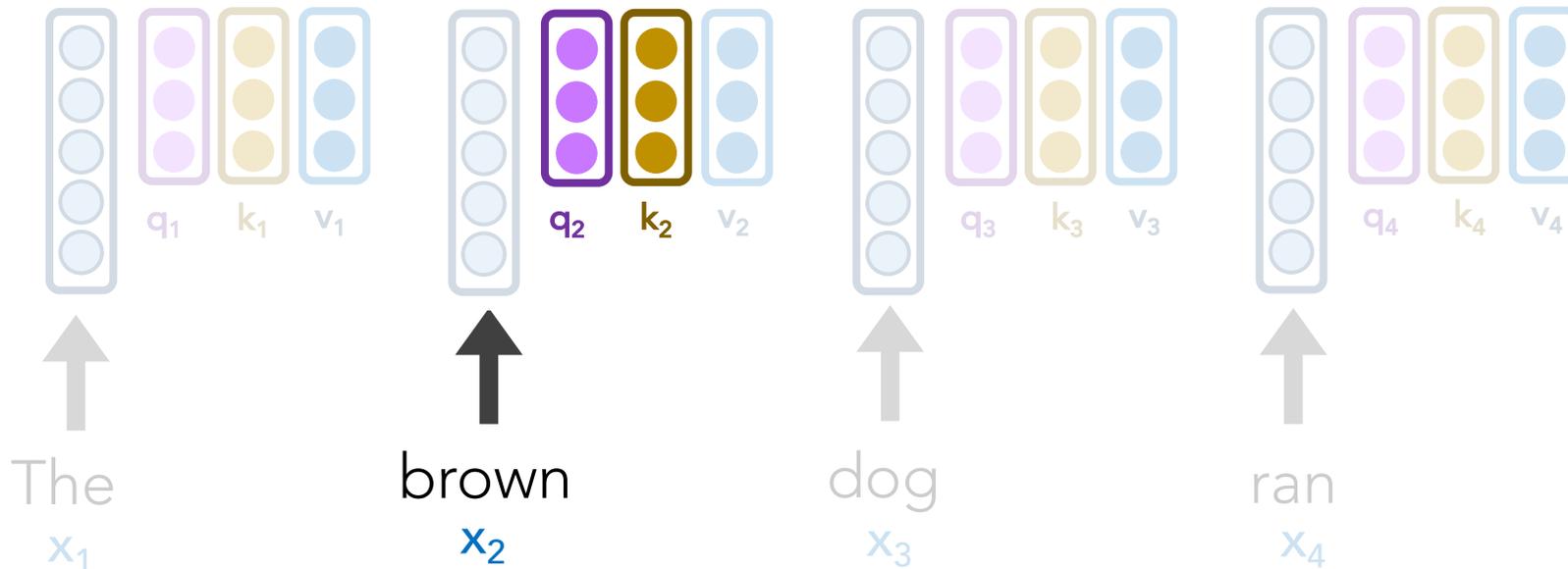


# Self-Attention

**Step 2:** For word  $x_2$ , let's calculate the scores  $s_1, s_2, s_3, s_4$ , which represent how much attention to pay to each respective "word"  $v_i$

$$s_2 = q_2 \cdot k_2 = 124$$

$$s_1 = q_2 \cdot k_1 = 92$$



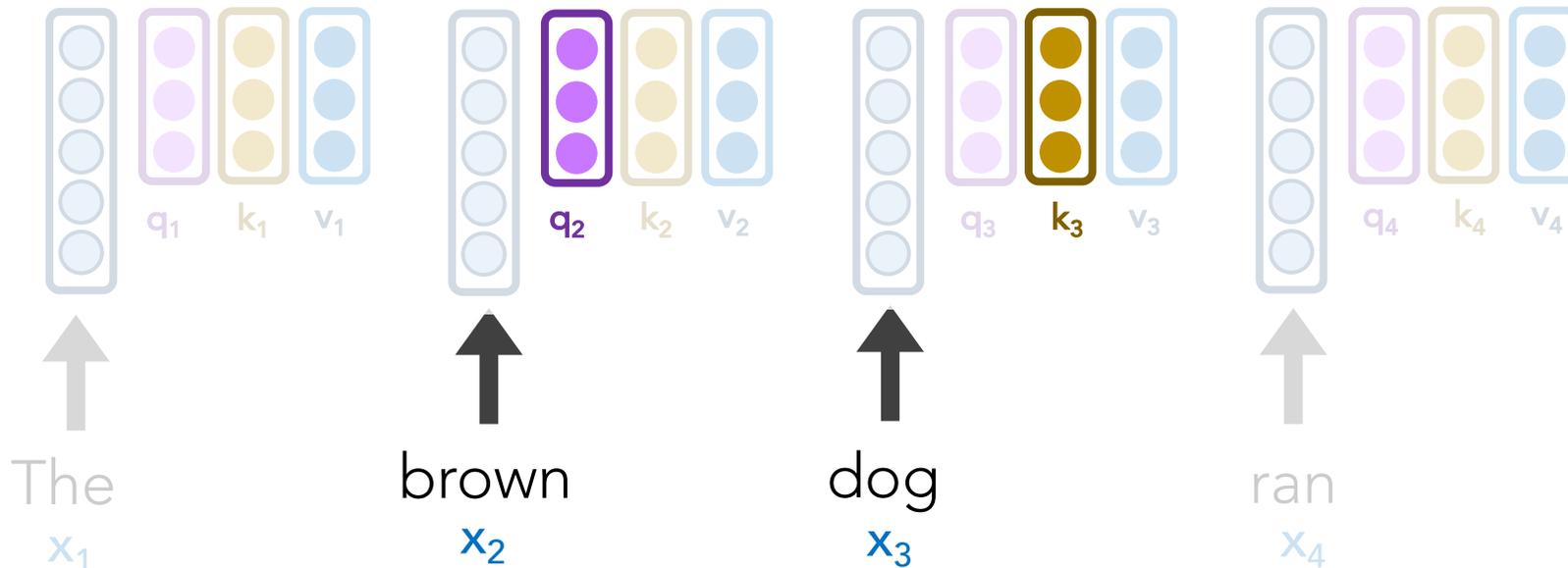
# Self-Attention

**Step 2:** For word  $x_2$ , let's calculate the scores  $s_1, s_2, s_3, s_4$ , which represent how much attention to pay to each respective "word"  $v_i$

$$s_3 = q_2 \cdot k_3 = 22$$

$$s_2 = q_2 \cdot k_2 = 124$$

$$s_1 = q_2 \cdot k_1 = 92$$



# Self-Attention

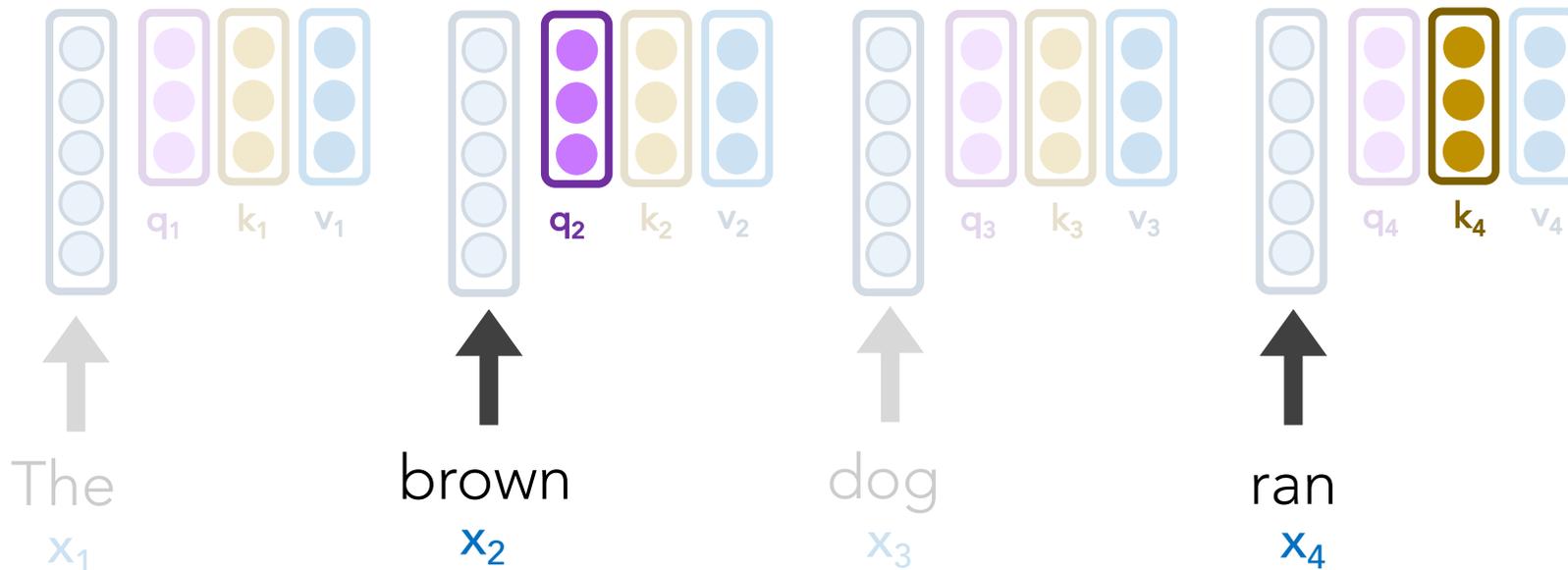
**Step 2:** For word  $x_2$ , let's calculate the scores  $s_1, s_2, s_3, s_4$ , which represent how much attention to pay to each respective "word"  $v_i$

$$s_4 = q_2 \cdot k_4 = 8$$

$$s_3 = q_2 \cdot k_3 = 22$$

$$s_2 = q_2 \cdot k_2 = 124$$

$$s_1 = q_2 \cdot k_1 = 92$$



# Self-Attention

**Step 3:** Our scores  $s_1, s_2, s_3, s_4$  don't sum to 1. Let's divide by  $\sqrt{\text{len}(k_i)}$  and **softmax** it

$$s_4 = q_2 \cdot k_4 = 8$$

$$a_4 = \sigma(s_4/8) = 0$$

$$s_3 = q_2 \cdot k_3 = 22$$

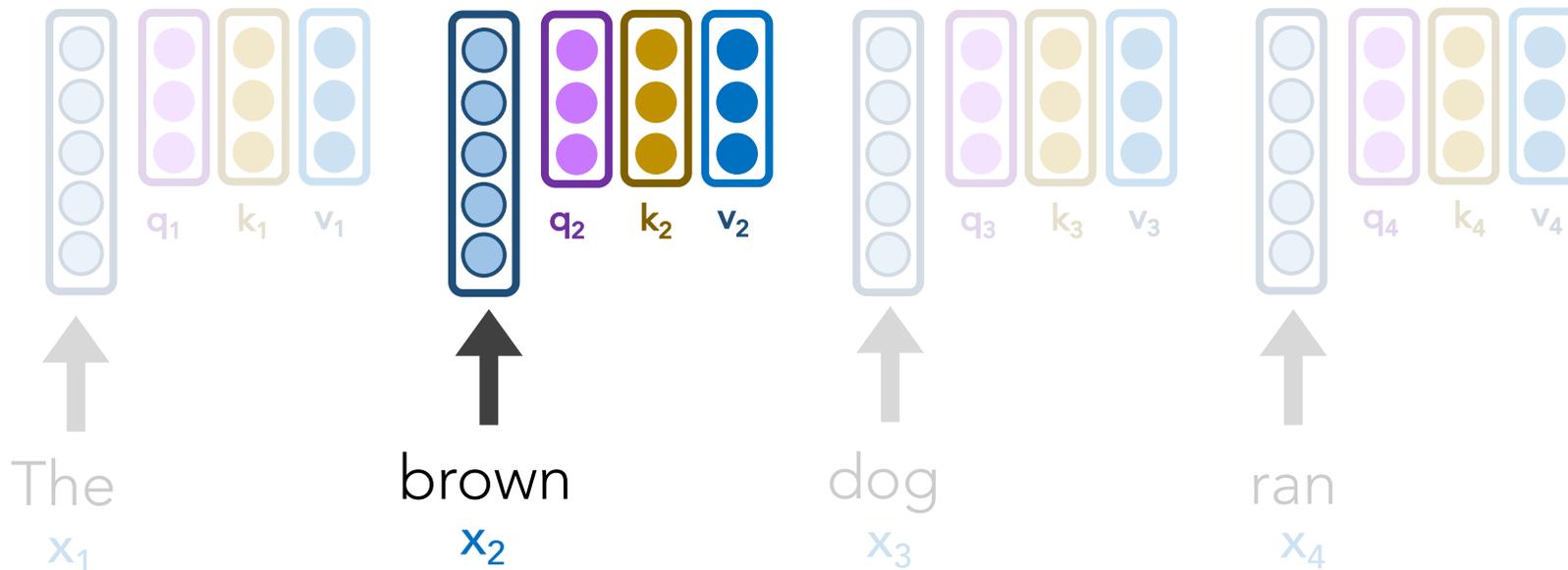
$$a_3 = \sigma(s_3/8) = .01$$

$$s_2 = q_2 \cdot k_2 = 124$$

$$a_2 = \sigma(s_2/8) = .91$$

$$s_1 = q_2 \cdot k_1 = 92$$

$$a_1 = \sigma(s_1/8) = .08$$



# Self-Attention

**Step 3:** Our scores  $s_1, s_2, s_3, s_4$  don't sum to 1. Let's divide by  $\sqrt{\text{len}(k_i)}$  and **softmax** it

$$s_4 = \mathbf{q}_2 \cdot \mathbf{k}_4 = 8$$

$$\mathbf{a}_4 = \sigma(s_4/8) = 0$$

$$s_3 = \mathbf{q}_2 \cdot \mathbf{k}_3 = 22$$

$$\mathbf{a}_3 = \sigma(s_3/8) = .01$$

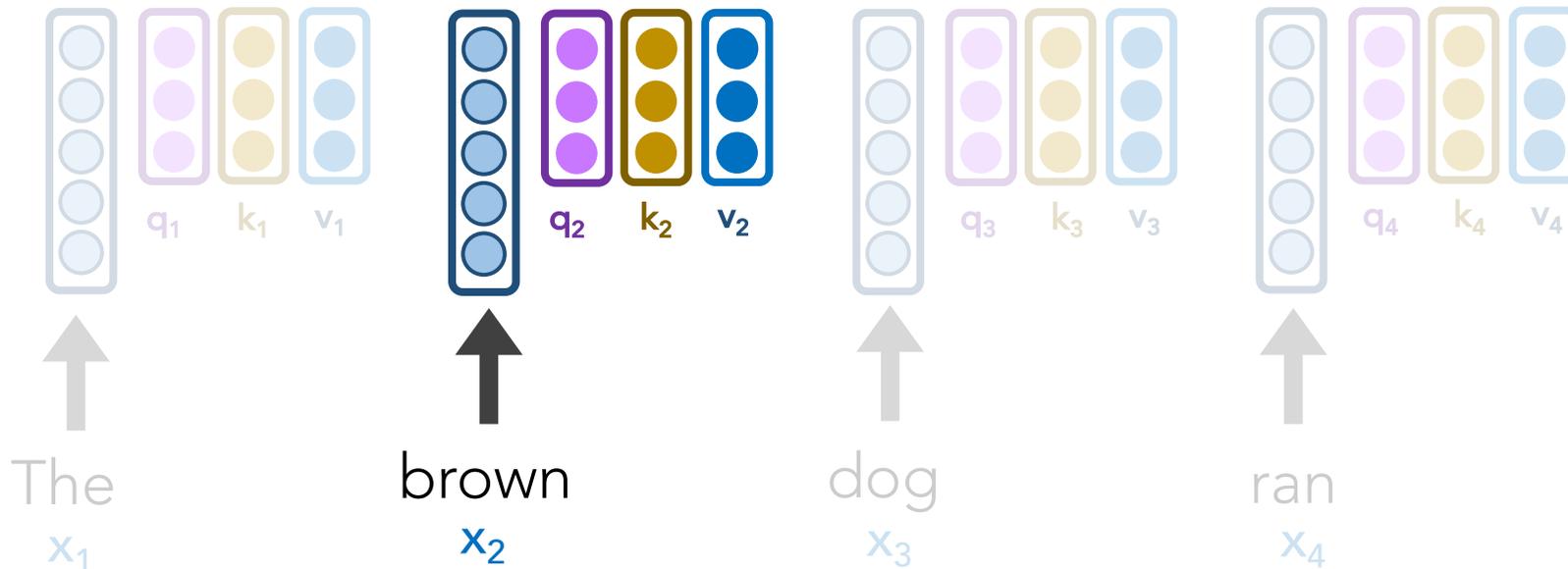
$$s_2 = \mathbf{q}_2 \cdot \mathbf{k}_2 = 124$$

$$\mathbf{a}_2 = \sigma(s_2/8) = .91$$

$$s_1 = \mathbf{q}_2 \cdot \mathbf{k}_1 = 92$$

$$\mathbf{a}_1 = \sigma(s_1/8) = .08$$

Instead of these  $\mathbf{a}_i$  values directly weighting our original  $\mathbf{x}_i$  word vectors, they directly weight our  $\mathbf{v}_i$  vectors.

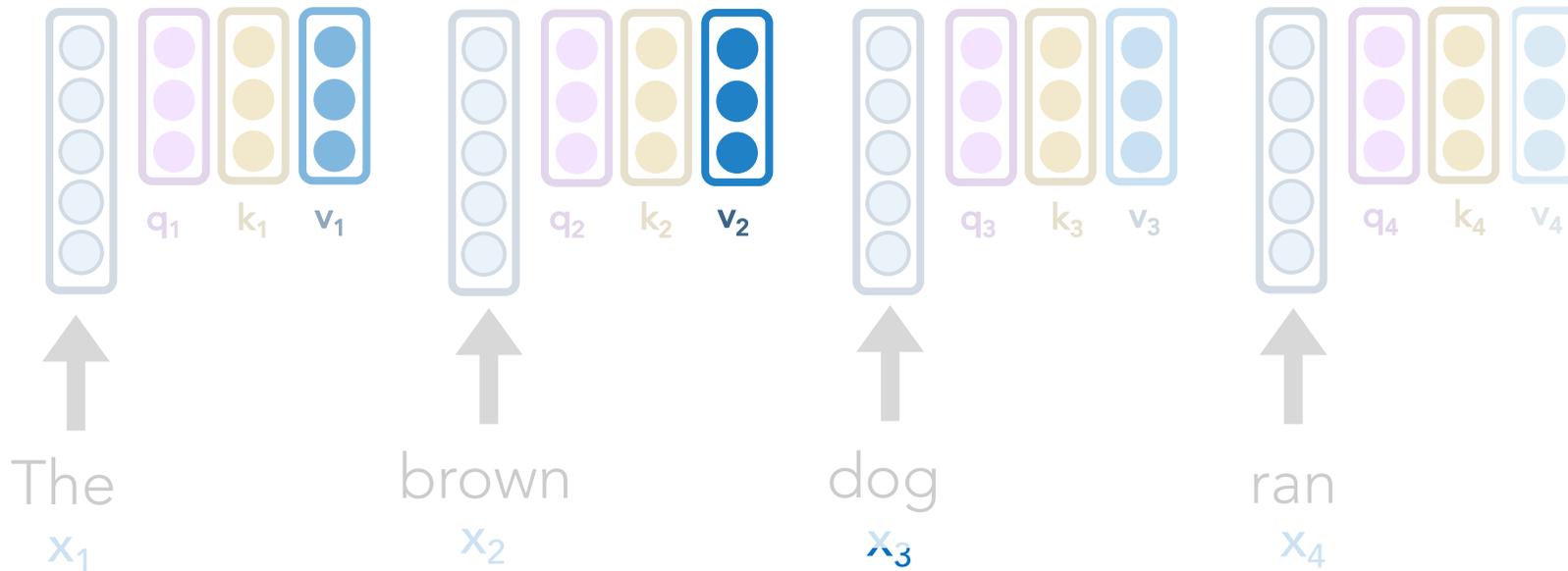


# Self-Attention

Step 4: Let's weight our  $v_i$  vectors and simply sum them up!

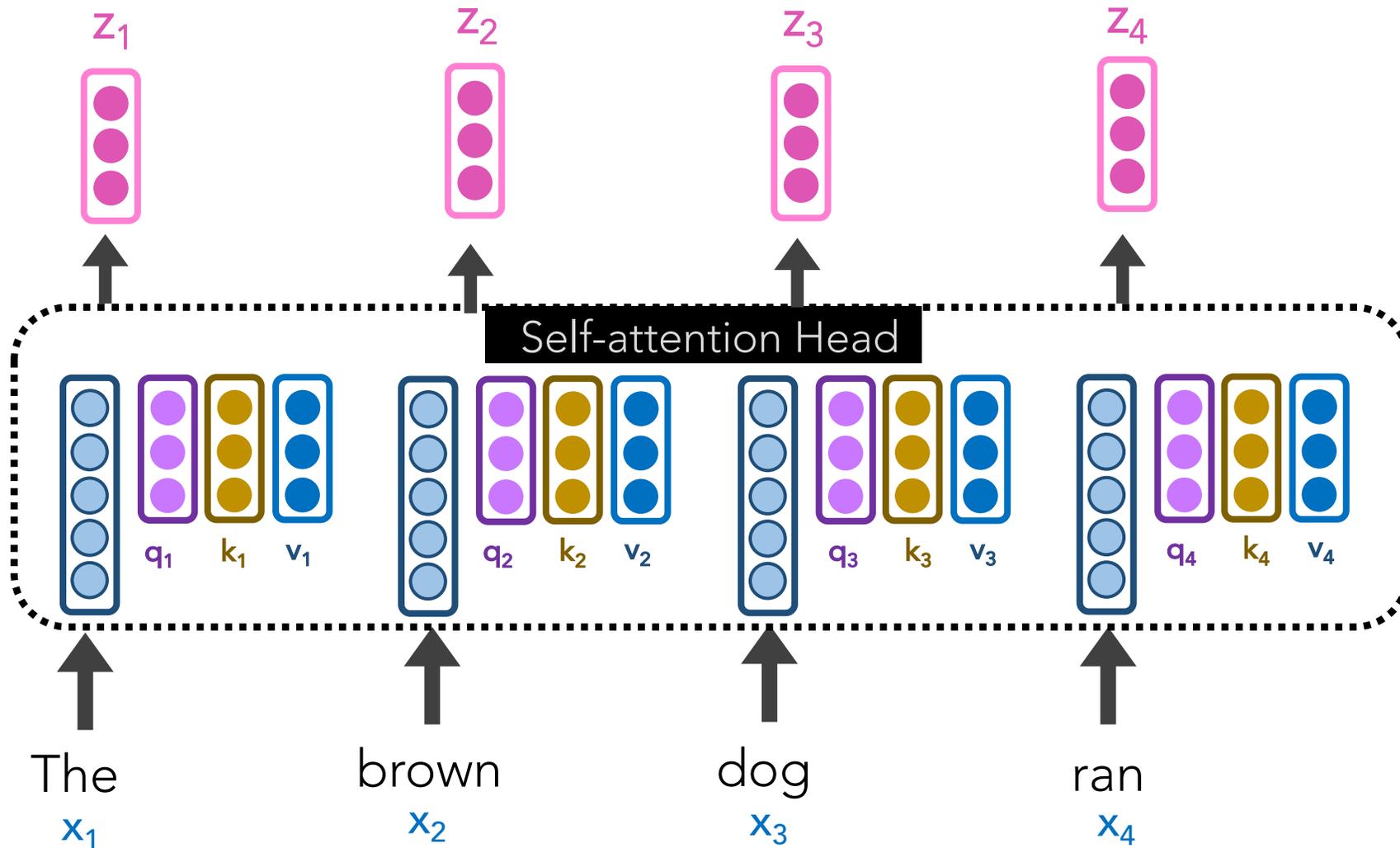


$$\begin{aligned}z_2 &= \mathbf{a}_1 \cdot \mathbf{v}_1 + \mathbf{a}_2 \cdot \mathbf{v}_2 + \mathbf{a}_3 \cdot \mathbf{v}_3 + \mathbf{a}_4 \cdot \mathbf{v}_4 \\ &= 0.08 \cdot \mathbf{v}_1 + 0.91 \cdot \mathbf{v}_2 + 0.01 \cdot \mathbf{v}_3 + 0 \cdot \mathbf{v}_4\end{aligned}$$



# Self-Attention

Tada! Now we have great, new representations  $z_i$  via a self-attention head

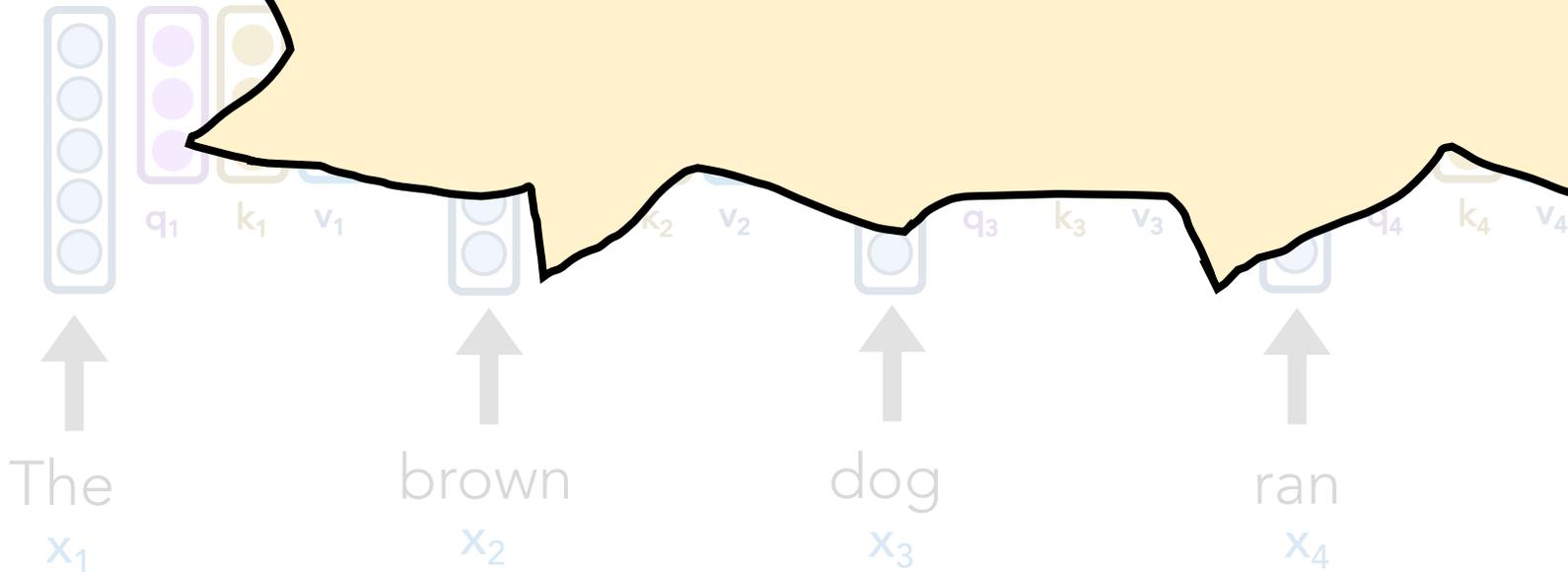


# Self-Attention

Tada! No

Takeaway:

**Self-Attention** is powerful; allows us to create great, context-aware representations



# Self-Attention

$$b = \text{softmax} \left( \frac{QK^T}{\alpha} \right) V$$



hardmaru  
@hardmaru



The most important formula in deep learning after 2018

## Self-Attention

**What is self-attention?** Self-attention calculates a weighted average of feature representations with the weight proportional to a similarity score between pairs of representations. Formally, an input sequence of  $n$  tokens of dimensions  $d$ ,  $X \in \mathbf{R}^{n \times d}$ , is projected using three matrices  $W_Q \in \mathbf{R}^{d \times d_q}$ ,  $W_K \in \mathbf{R}^{d \times d_k}$ , and  $W_V \in \mathbf{R}^{d \times d_v}$  to extract feature representations  $Q$ ,  $K$ , and  $V$ , referred to as query, key, and value respectively with  $d_k = d_q$ . The outputs  $Q$ ,  $K$ ,  $V$  are computed as

$$Q = XW_Q, \quad K = XW_K, \quad V = XW_V. \quad (1)$$

So, self-attention can be written as,

$$S = D(Q, K, V) = \text{softmax} \left( \frac{QK^T}{\sqrt{d_q}} \right) V, \quad (2)$$

where softmax denotes a *row-wise* softmax normalization function. Thus, each element in  $S$  depends on all other elements in the same row.

9:08 PM · Feb 9, 2021 · Twitter Web App

553 Retweets 42 Quote Tweets 3,338 Likes

# Outline

-  Self-Attention
-  Transformer Encoder
-  Transformer Decoder
-  Language Modeling With  
Transformers

# Outline

 Self-Attention

 Transformer Encoder

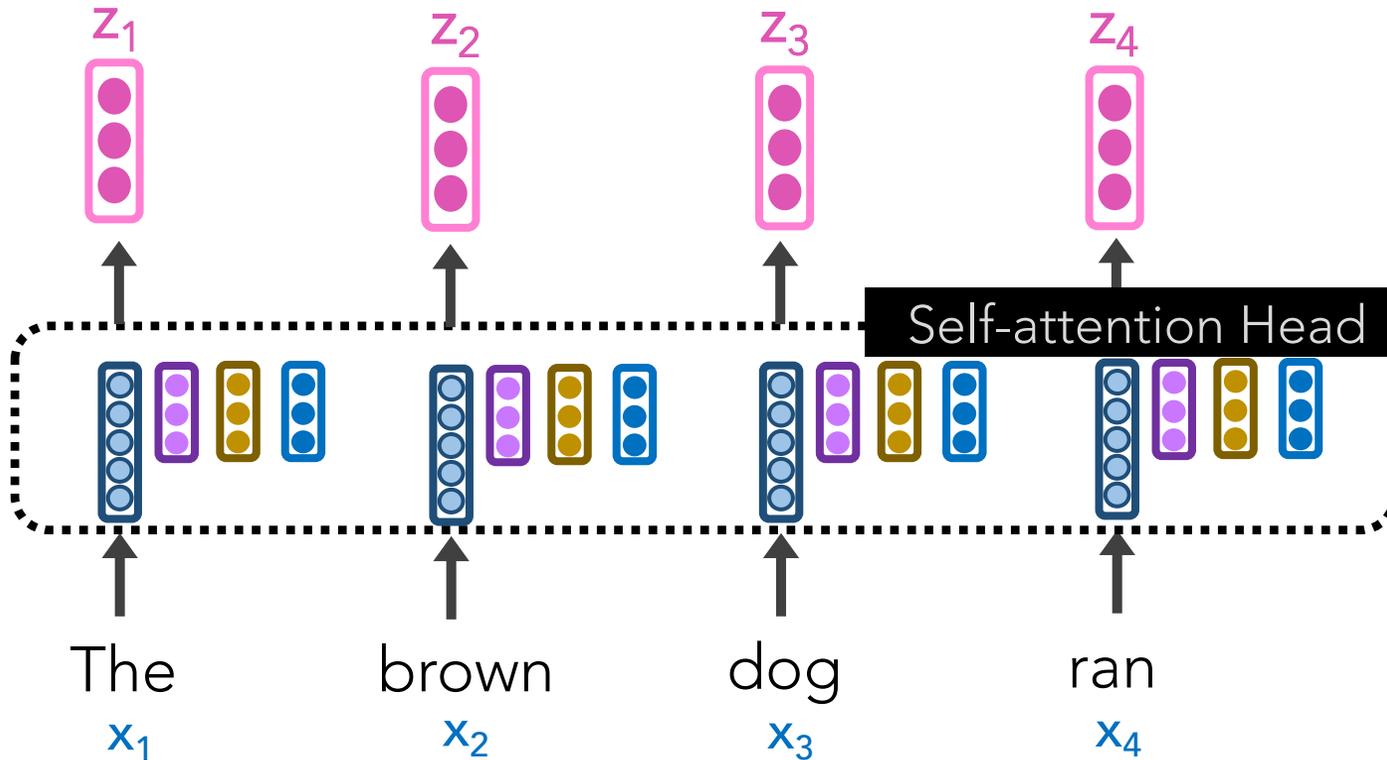
 Transformer Decoder

 Language Modeling With  
Transformers

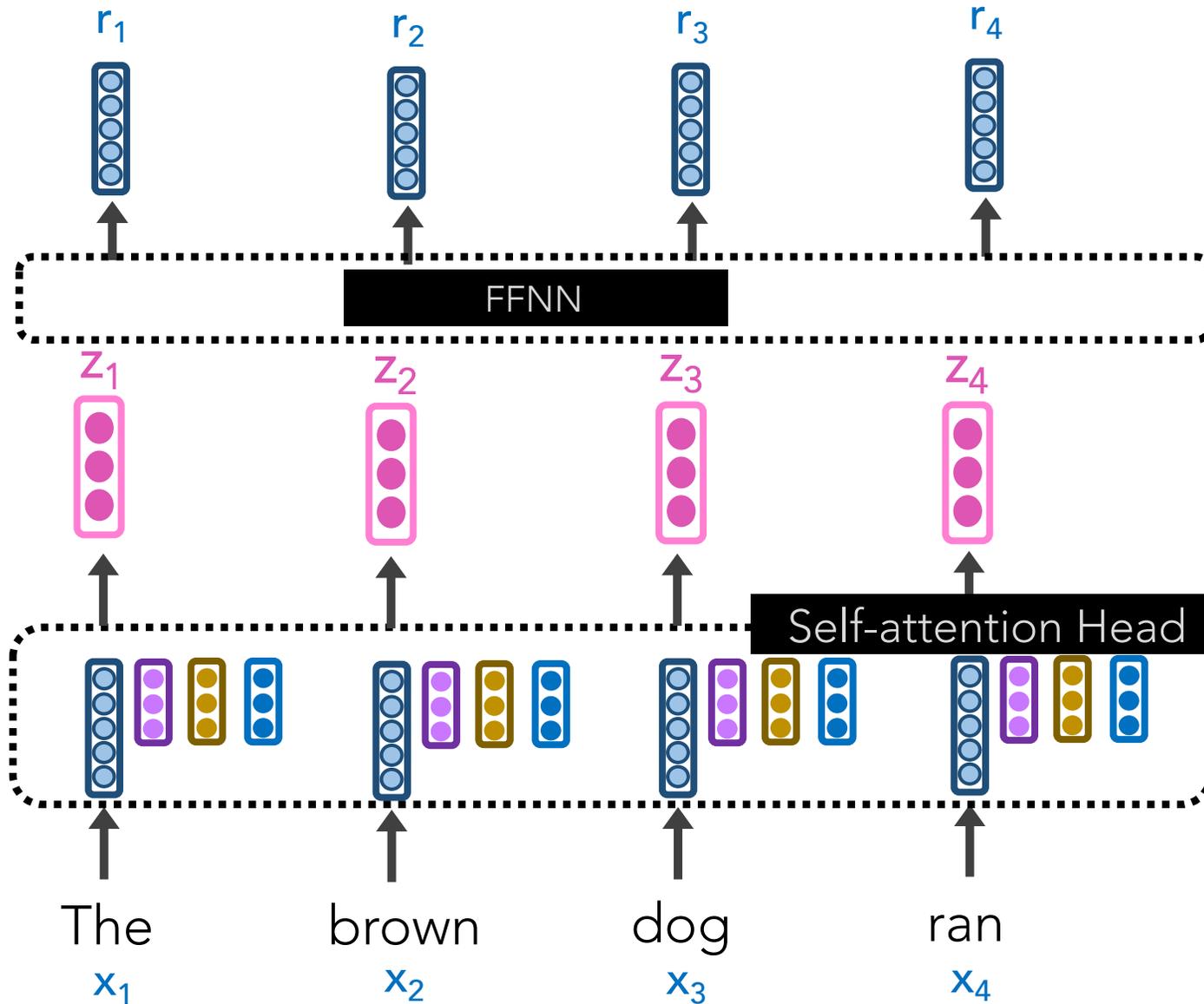
# Self-Attention

---

Let's further pass each  $z_i$  through a FeedForward NN

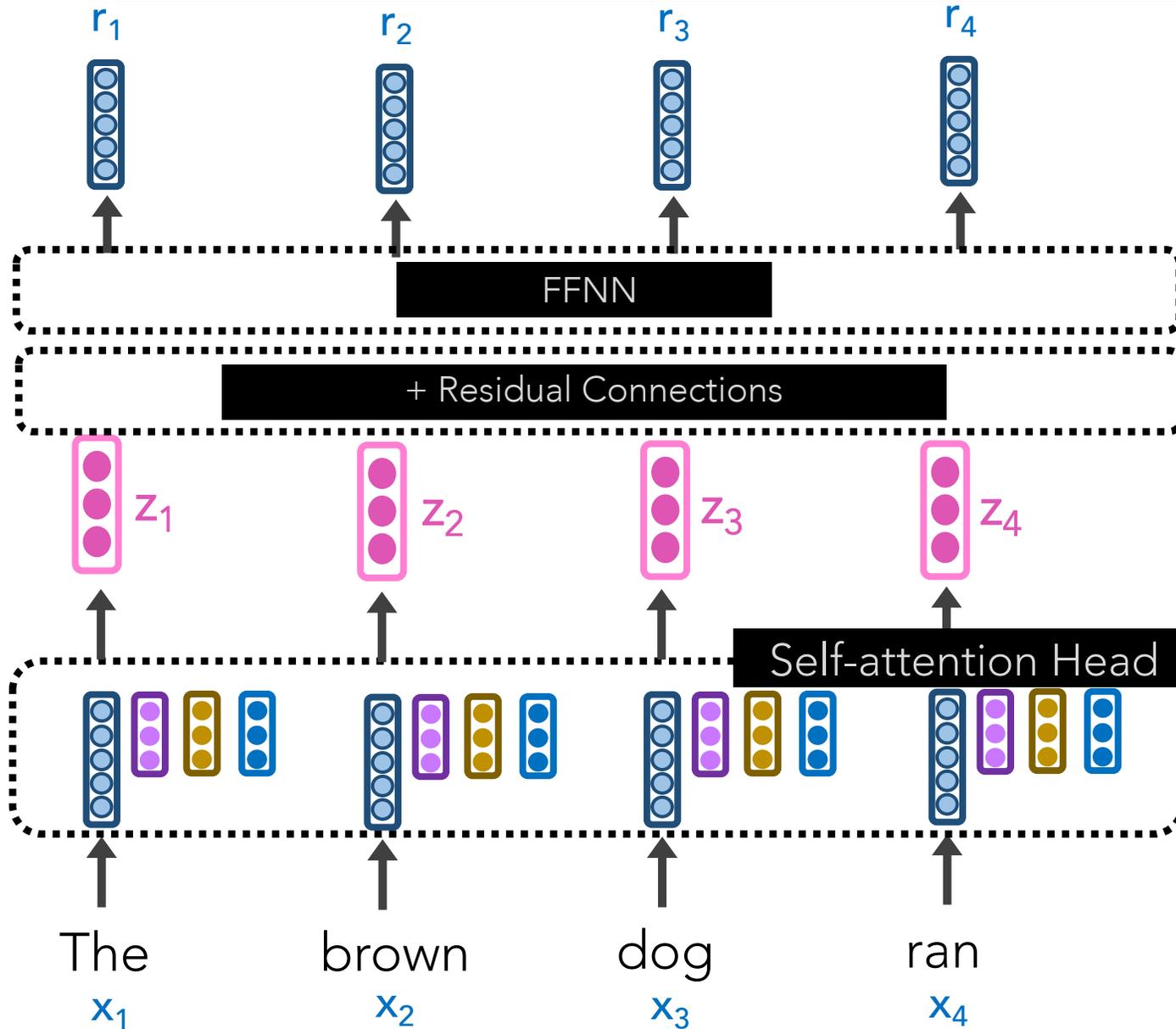


# Self-Attention + FFNN



Let's further pass each  $z_i$  through a Feed Forward NN

# Self-Attention + FFNN + Residual Connections

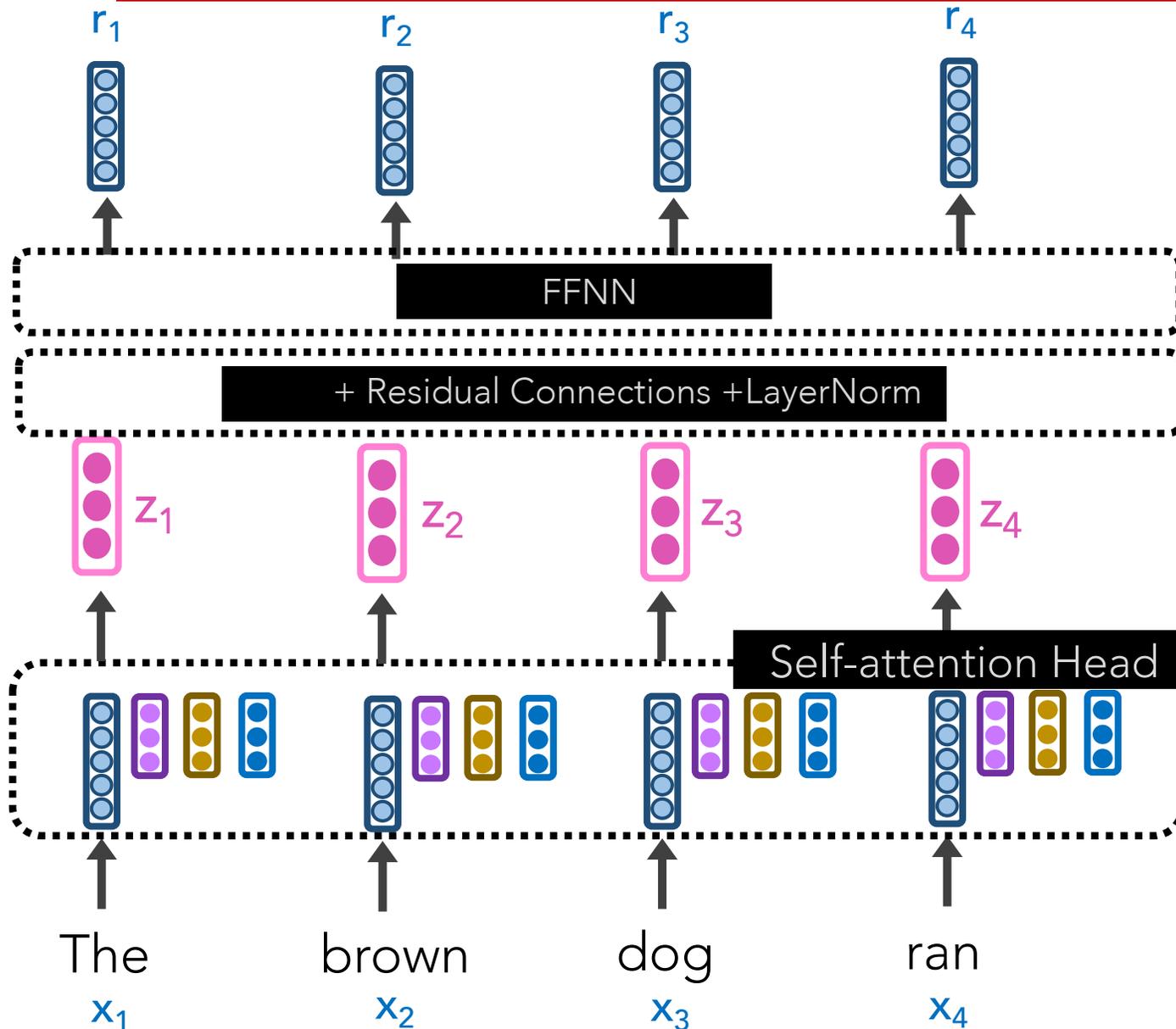


Let's further pass each  $z_i$  through a FFNN

We add a **residual connection** to help ensure relevant info is getting forward passed.

$$v = z + x$$

# Self-Attention + FFNN + Residual Connections



Let's further pass each  $z_i$  through a FFNN

We add **residual connection** to help ensure relevant info is getting forward passed.

$$v = z + x$$

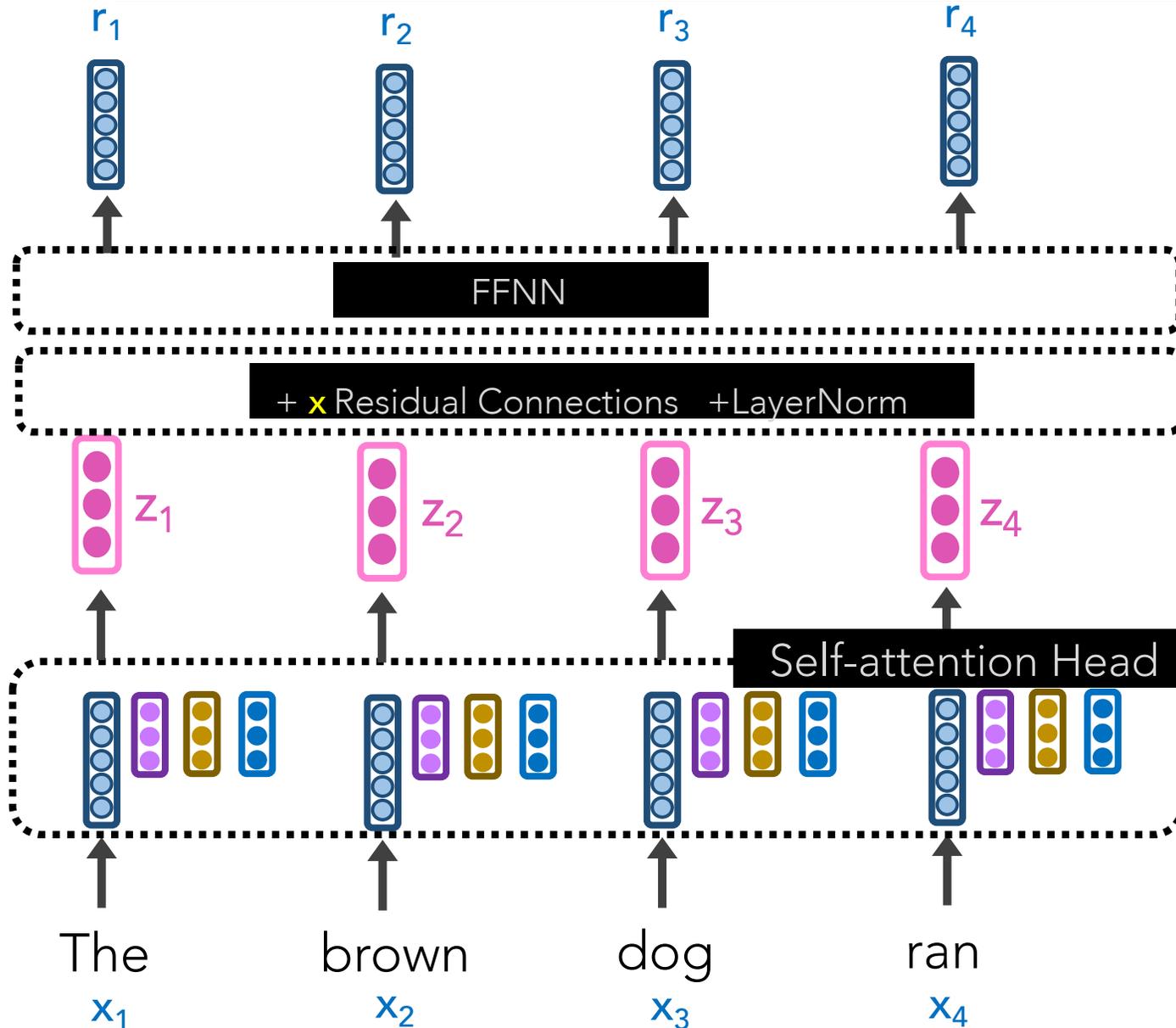
We perform **LayerNorm** to stabilize the network and allow for proper gradient flow. You should do this after the FFNN, too.

# Stabilizing Gradient Flow: Residual Connection and LayerNorm

- Residual connection:  $y = f(x) + x$ 
  - $f$  might be a complex function and gives small gradients wrt  $x$ , adding  $x$  back to  $f(x)$  gives higher values of the gradient
- Layer Normalization (LayerNorm):
  - Another way to prevent vanishing gradients

$$y = \frac{x - \mathbf{E}[x]}{\sqrt{\mathbf{Var}[x] + \epsilon}} * \gamma + \beta$$

# Self-Attention + FFNN



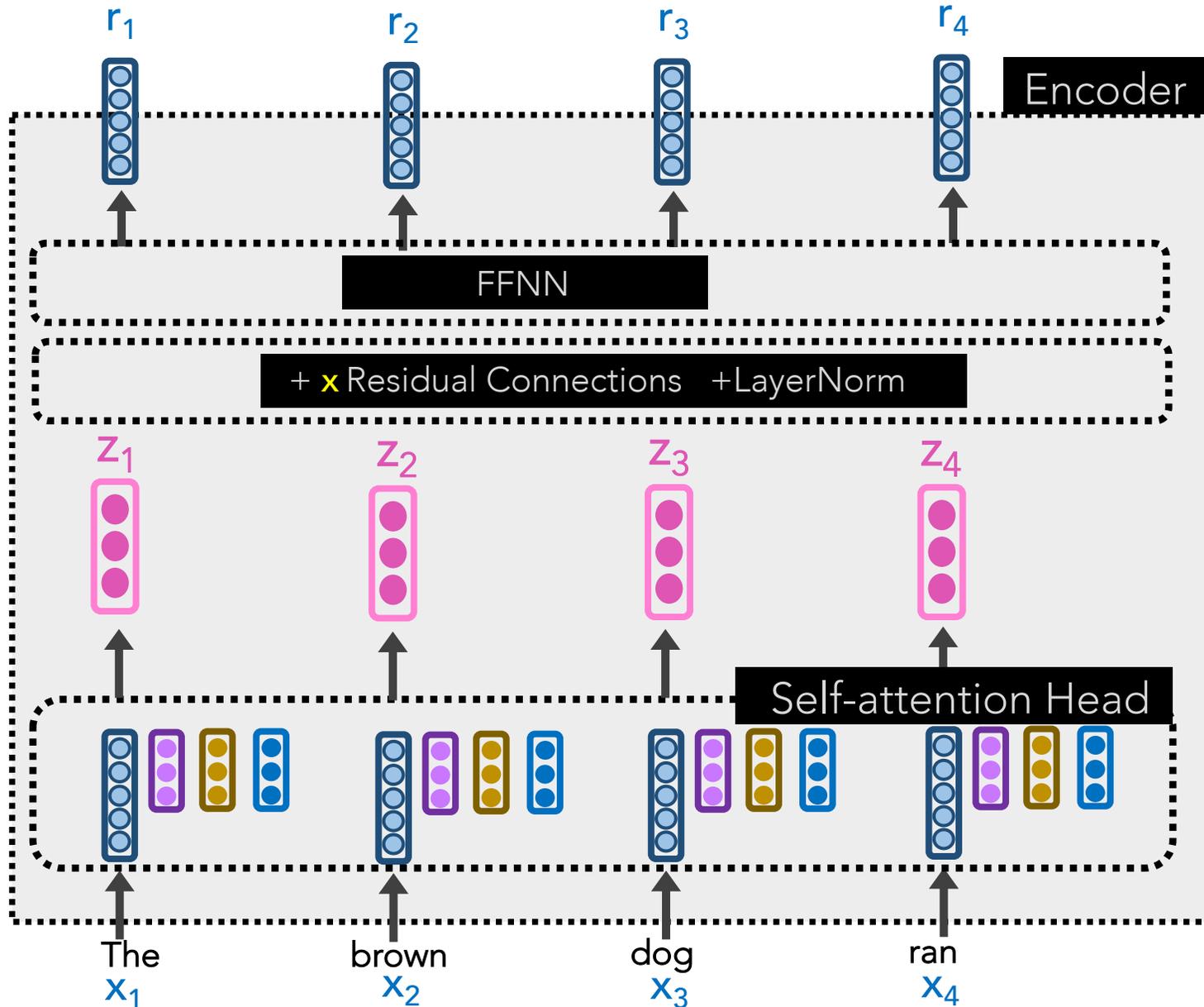
Let's further pass each  $z_i$  through a FFNN

We concat w/ a **residual connection** to help ensure relevant info is getting forward passed.

We perform **LayerNorm** to stabilize the network and allow for proper gradient flow. You should do this after the FFNN, too.

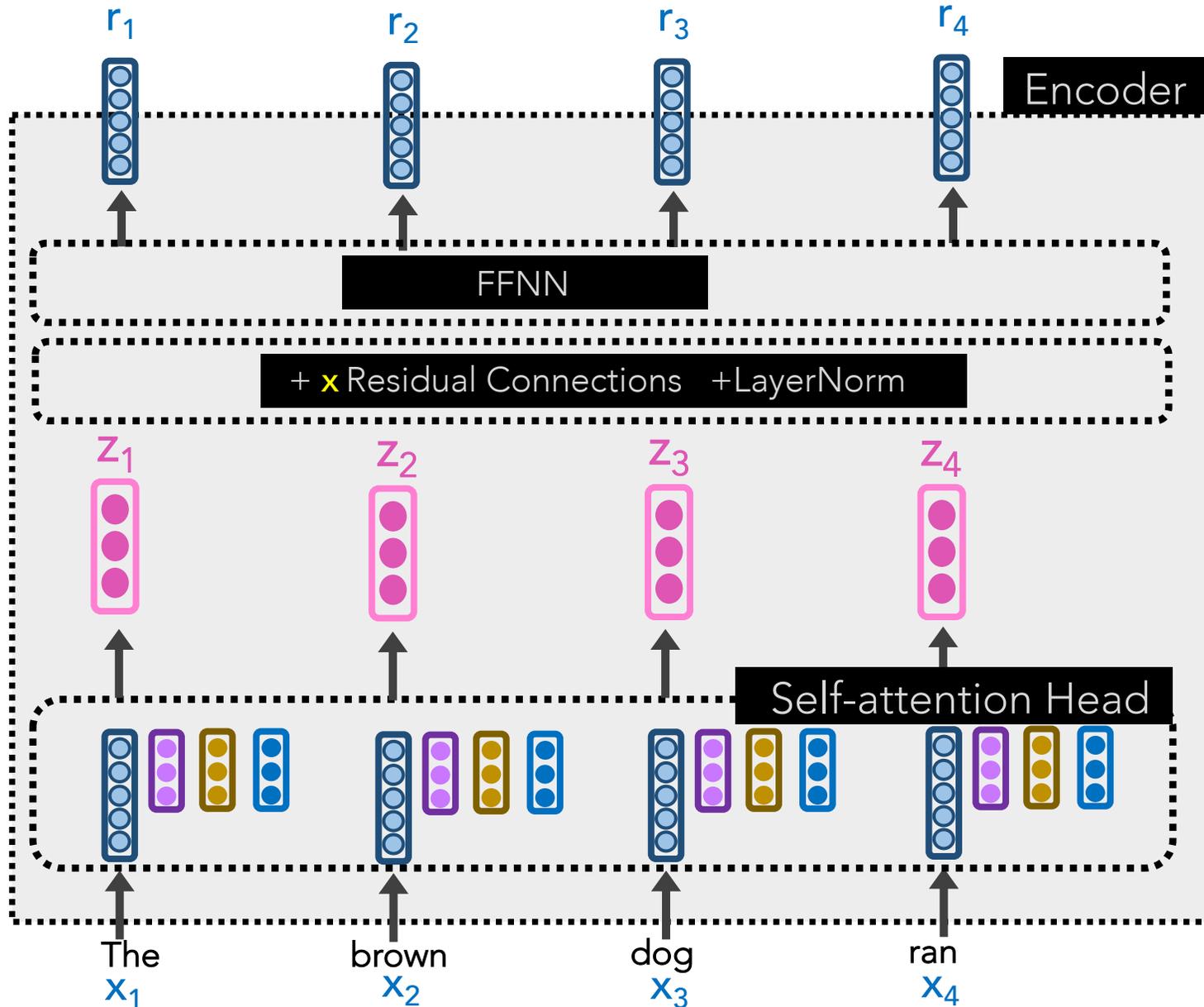
Each  $z_i$  can be computed in **parallel**, unlike RNNs!

# Transformer Encoder



Yay! Our  $r_i$  vectors are our new representations, and this entire process is called a **Transformer Encoder**

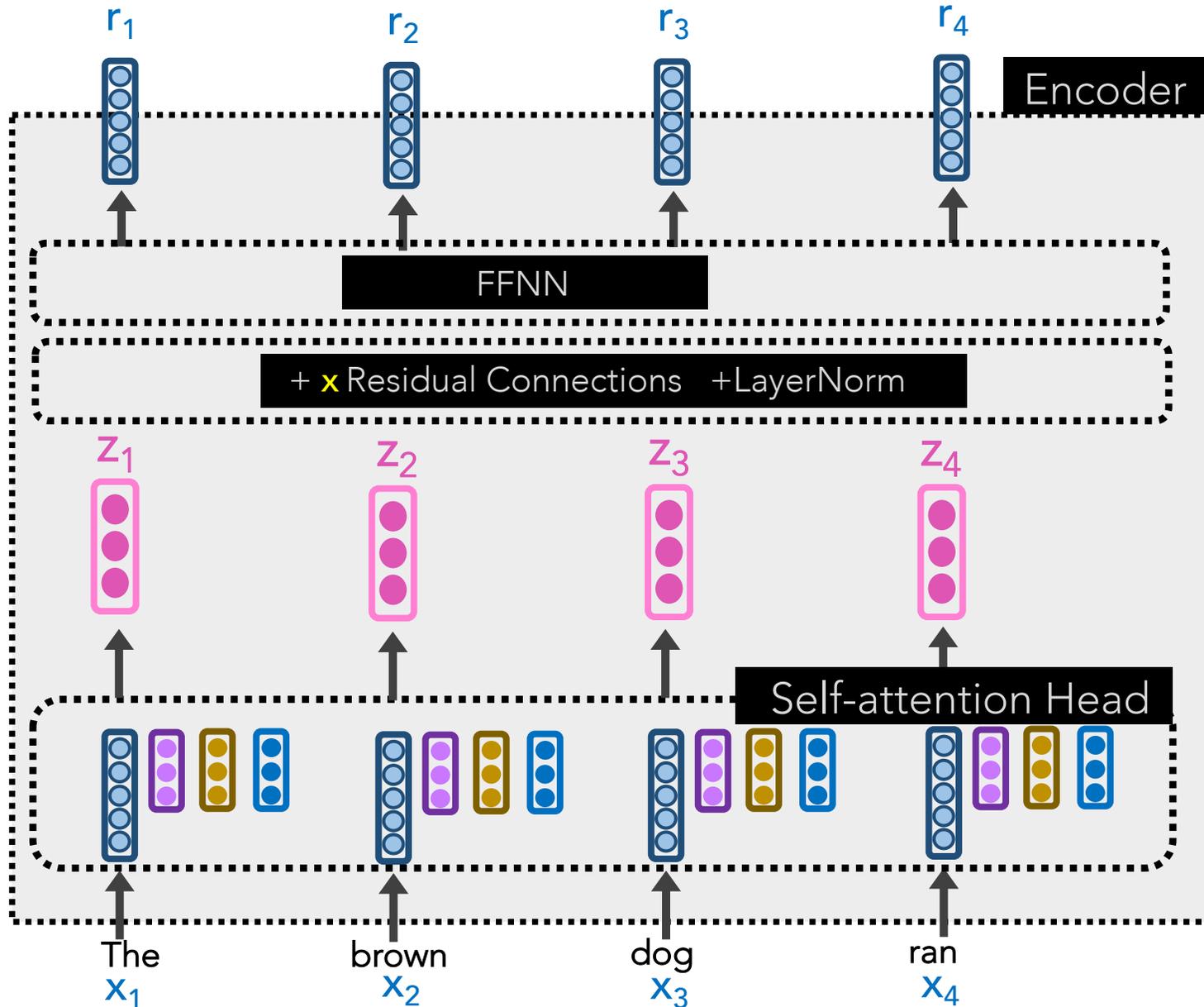
# Transformer Encoder



Yay! Our  $r_i$  vectors are our new representations, and this entire process is called a **Transformer Encoder**

**Problem:** there is no concept of positionality. Words are weighted as if a "bag of words"

# Transformer Encoder



Yay! Our  $r_i$  vectors are our new representations, and this entire process is called a **Transformer Encoder**

**Problem:** there is no concept of positionality. Words are weighted as if a "bag of words"

**Solution:** add to each input word  $x_i$  a **positional encoding**

Input to the model is now  $x_i + pos_i$

# How to encode position information?

- Self attention doesn't have a way to know whether an input token comes before or after another
  - Position is important in sequence modeling in NLP
- A way to introduce position information is add individual position encodings to the input for each position in the sequence

$$x_i = x_i + pos_i$$

Where  $pos_t$  is a position vector

# Properties of a good positional embedding

- It should output a unique encoding for each time-step (word's position in a sentence)
- Distance between any two time-steps should be consistent across sentences with different lengths.
  - The cat sat on the mat
  - The happy cat sat on the mat
- Our model should generalize to longer sentences without any efforts. Its values should be bounded.

# Absolute position embeddings

- Define a maximum context length you model can encode: say 1000 tokens.
  - Create a separate embedding table for each position.
  - Each index 1, 2, 3, ... gets an embedding.
  - Learn the embeddings with the model.
- Issues with Learned positions embeddings:
  - Maximum length that can be presented is limited (what if I get a 2000 token input)
  - Difficult to encode relative positions
    - The cat sat on the mat
    - The happy cat sat on the mat

# Functional (and fixed) position embeddings

## Sinusoidal embeddings

$$\vec{p}_t^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k \cdot t), & \text{if } i = 2k \\ \cos(\omega_k \cdot t), & \text{if } i = 2k + 1 \end{cases}$$

where

$$\omega_k = \frac{1}{10000^{2k/d}}$$

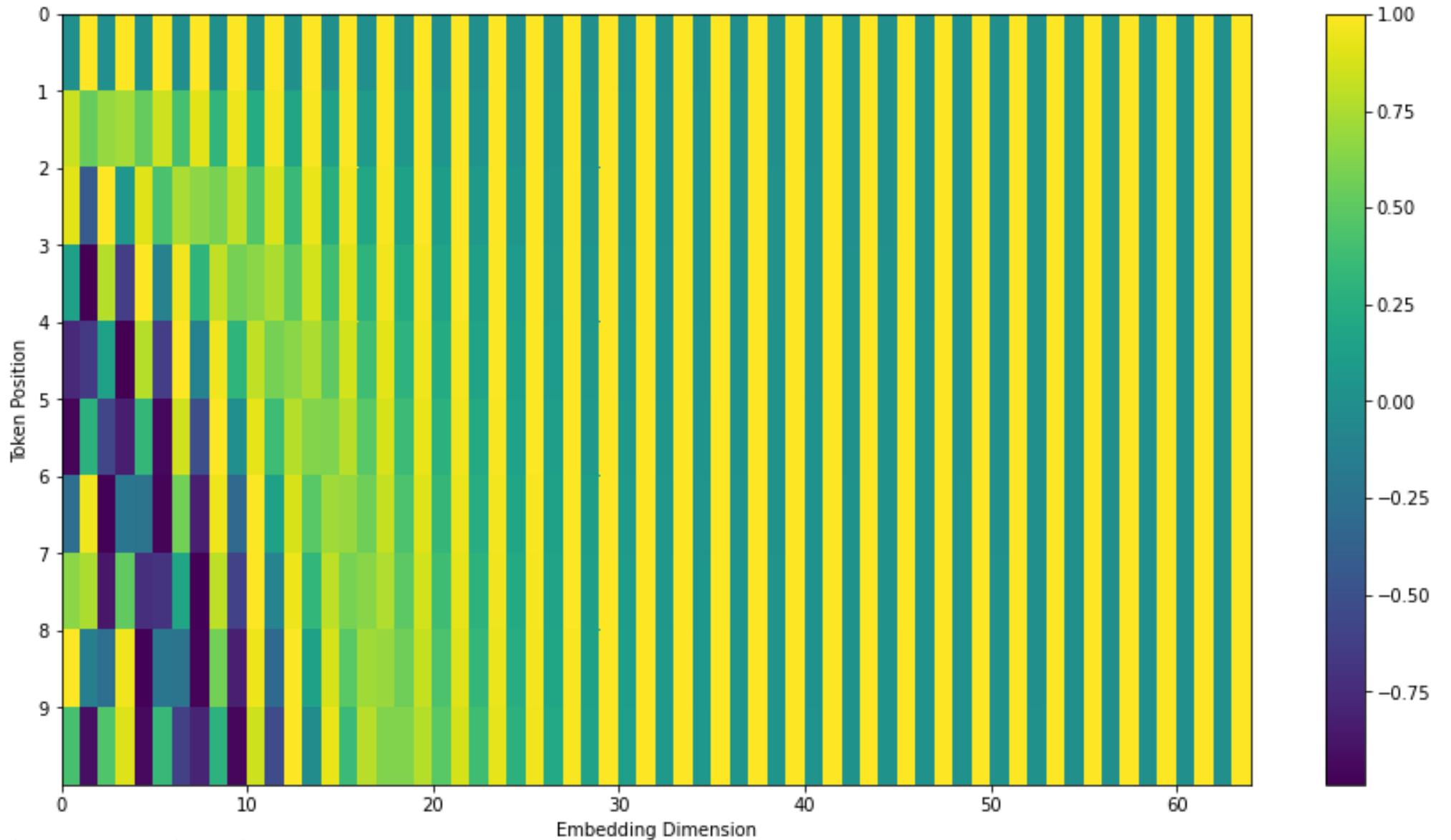
$$\vec{p}_t = \begin{bmatrix} \sin(\omega_1 \cdot t) \\ \cos(\omega_1 \cdot t) \\ \\ \sin(\omega_2 \cdot t) \\ \cos(\omega_2 \cdot t) \\ \\ \vdots \\ \\ \sin(\omega_{d/2} \cdot t) \\ \cos(\omega_{d/2} \cdot t) \end{bmatrix}$$

The frequencies are decreasing along the vector dimension. It forms a geometric progression on the wavelengths.

# Sinusoidal Embeddings: Intuition

0 :	0	0	0	0	8 :	1	0	0	0
1 :	0	0	0	1	9 :	1	0	0	1
2 :	0	0	1	0	10 :	1	0	1	0
3 :	0	0	1	1	11 :	1	0	1	1
4 :	0	1	0	0	12 :	1	1	0	0
5 :	0	1	0	1	13 :	1	1	0	1
6 :	0	1	1	0	14 :	1	1	1	0
7 :	0	1	1	1	15 :	1	1	1	1

# Position Encodings



# Variants of Positional Embeddings

- Rotary Positional Embeddings (RoPE): [\[2104.09864\] RoFormer: Enhanced Transformer with Rotary Position Embedding \(arxiv.org\)](#)
- AliBi: [\[2108.12409\] Train Short, Test Long: Attention with Linear Biases Enables Input Length Extrapolation \(arxiv.org\)](#)
- No embeddings(!?): [\[2203.16634\] Transformer Language Models without Positional Encodings Still Learn Positional Information \(arxiv.org\)](#)

# Summary So Far

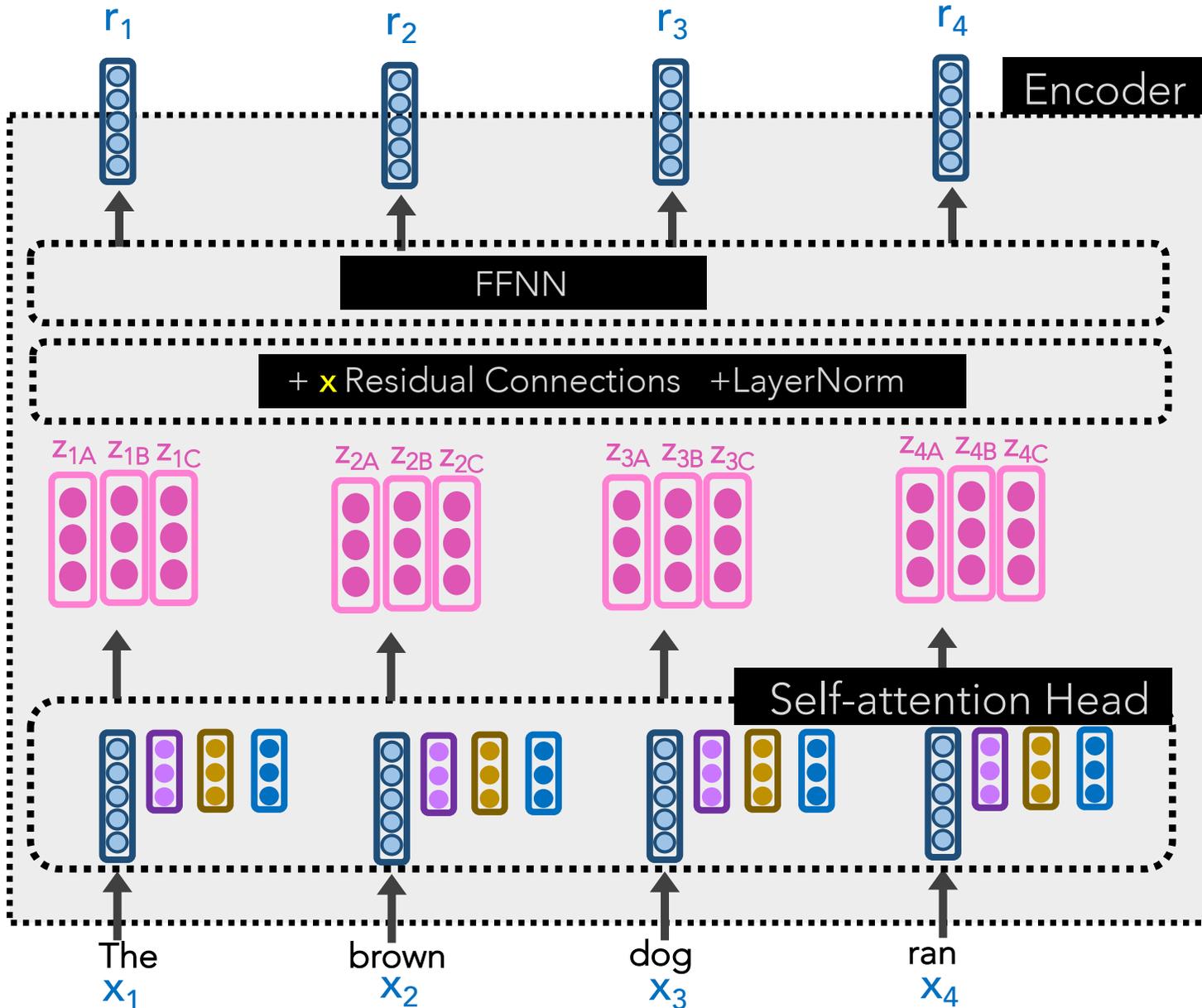
- Self Attention (Query, Key, Value)
- Residual Connections, Layernorm, FFN layers in between
- Positional Embeddings
- Combine all that and we get a transformer encoder.

A **Self-Attention Head** has just one set of **query/key/value** weight matrices  $w_q, w_k, w_v$

Words can relate in many ways, so it's restrictive to rely on just one Self-Attention Head in the system.

Let's create Multi-headed Self-Attention

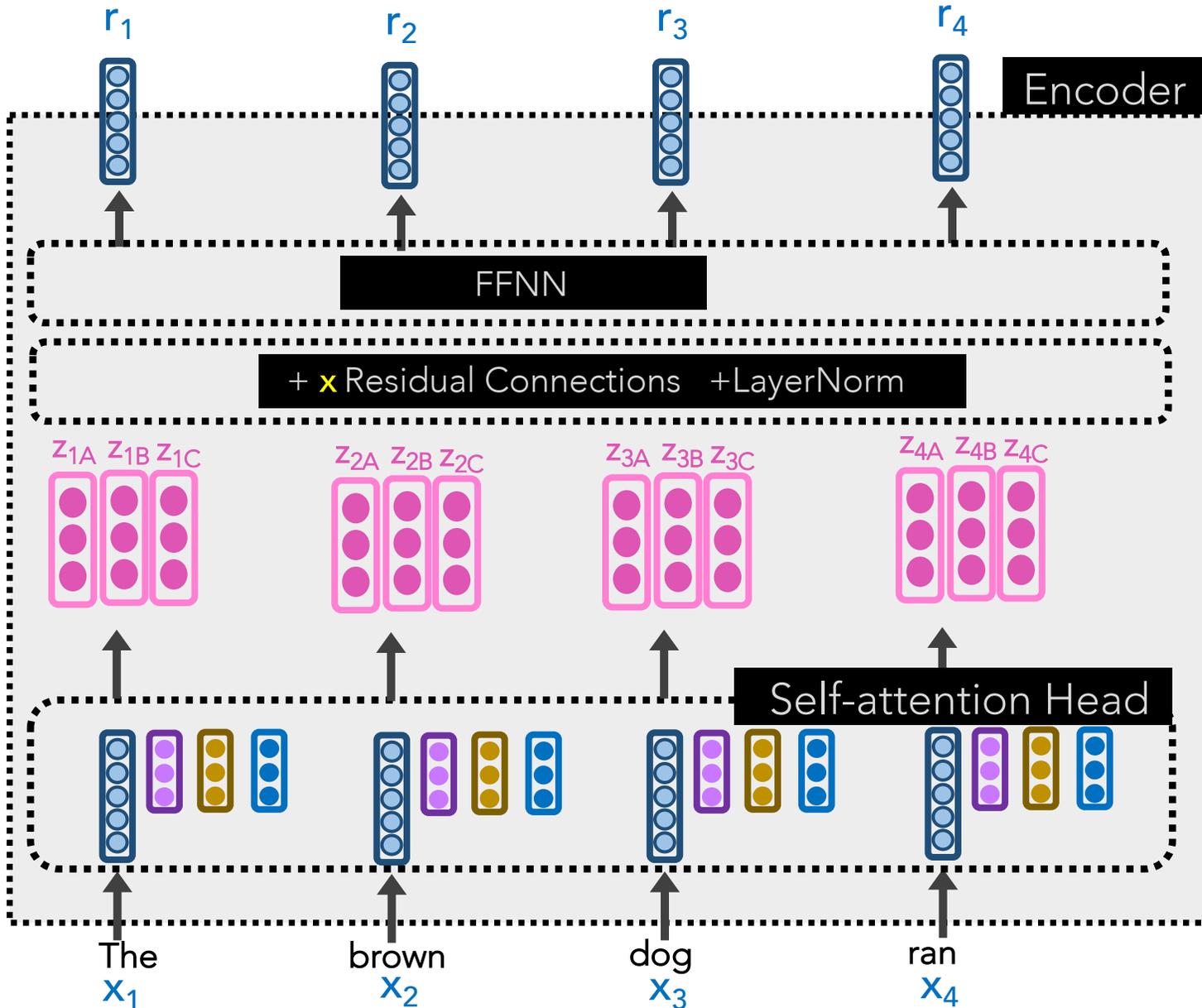
# Multi-head Attention



Each **Self-Attention Head** produces a  $z_i$  vector using query, key, and value vectors

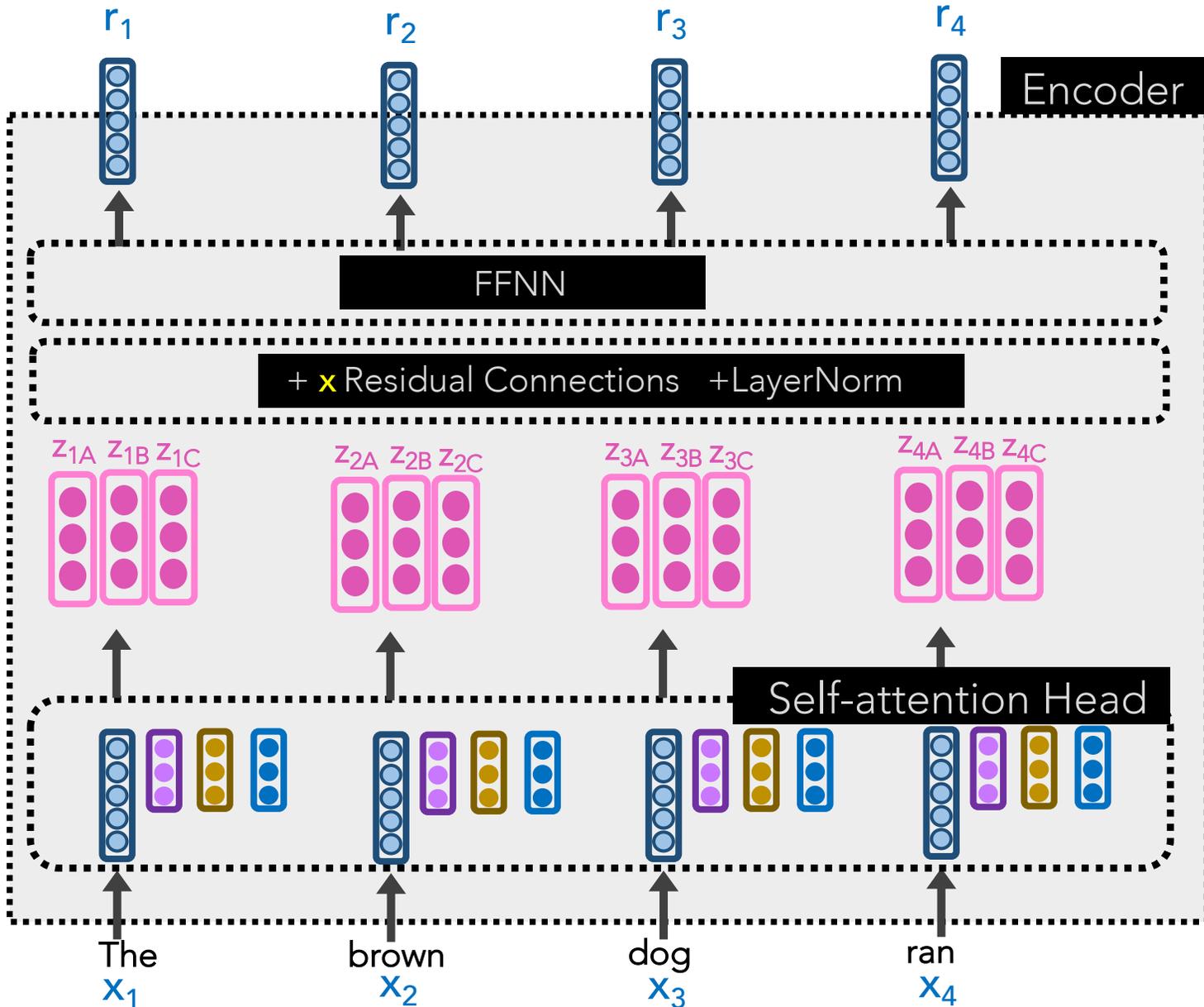
We can, in parallel, use **multiple heads** and concat the  $z_i$ 's. For each input create multiple query, key, and value vectors

# Transformer Encoder

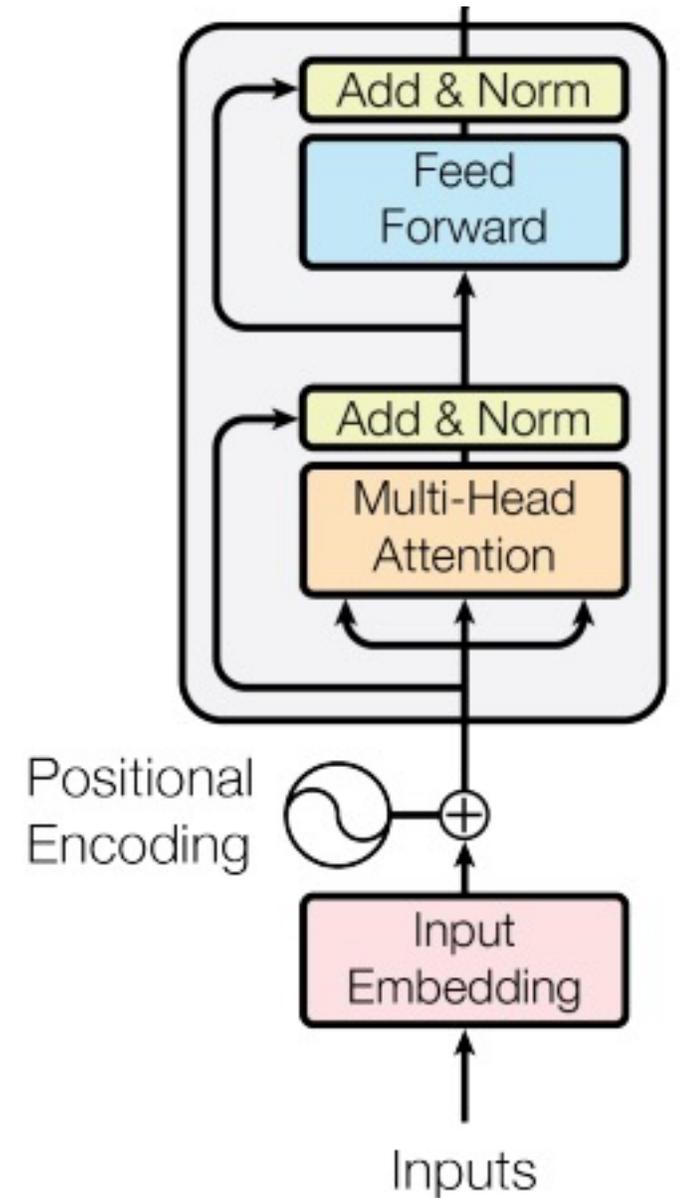


To recap: all of this looks fancy, but ultimately it's just producing a very good **contextualized embedding**  $r_i$  of each word  $x_i$

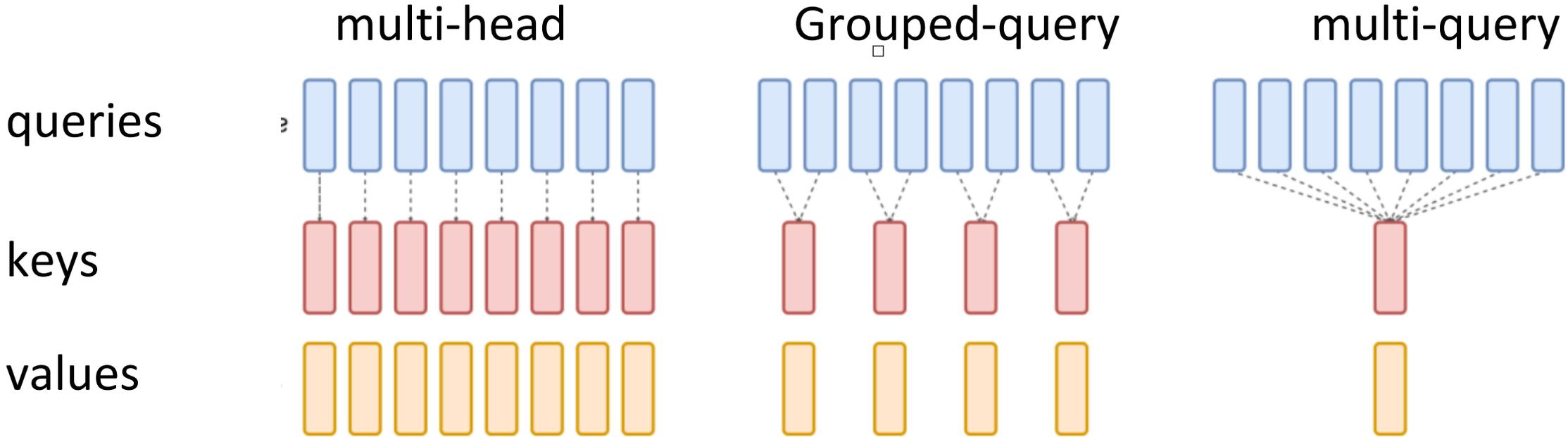
# Transformer Encoder



=

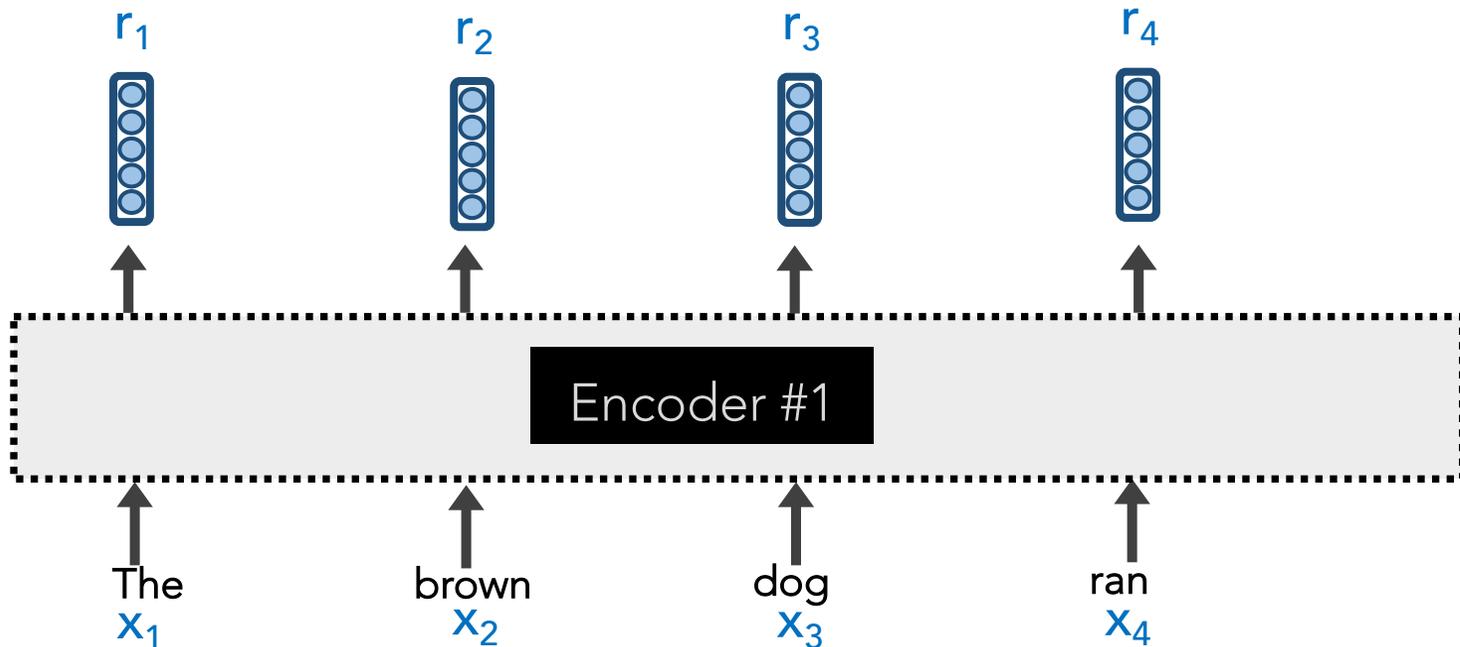


# Variants of multi-head attention attention



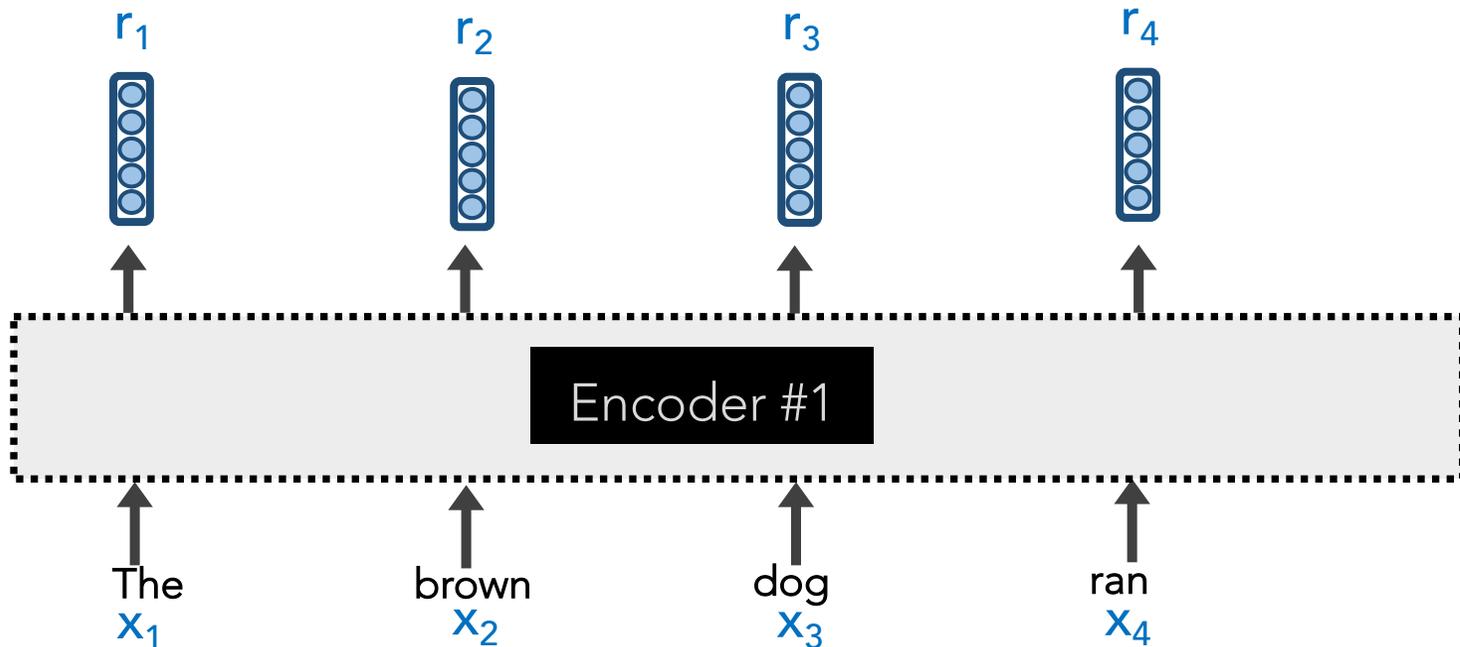
# Transformer Encoder

To recap: all of this looks fancy, but ultimately it's just producing a very good **contextualized embedding**  $r_i$  of each word  $x_i$



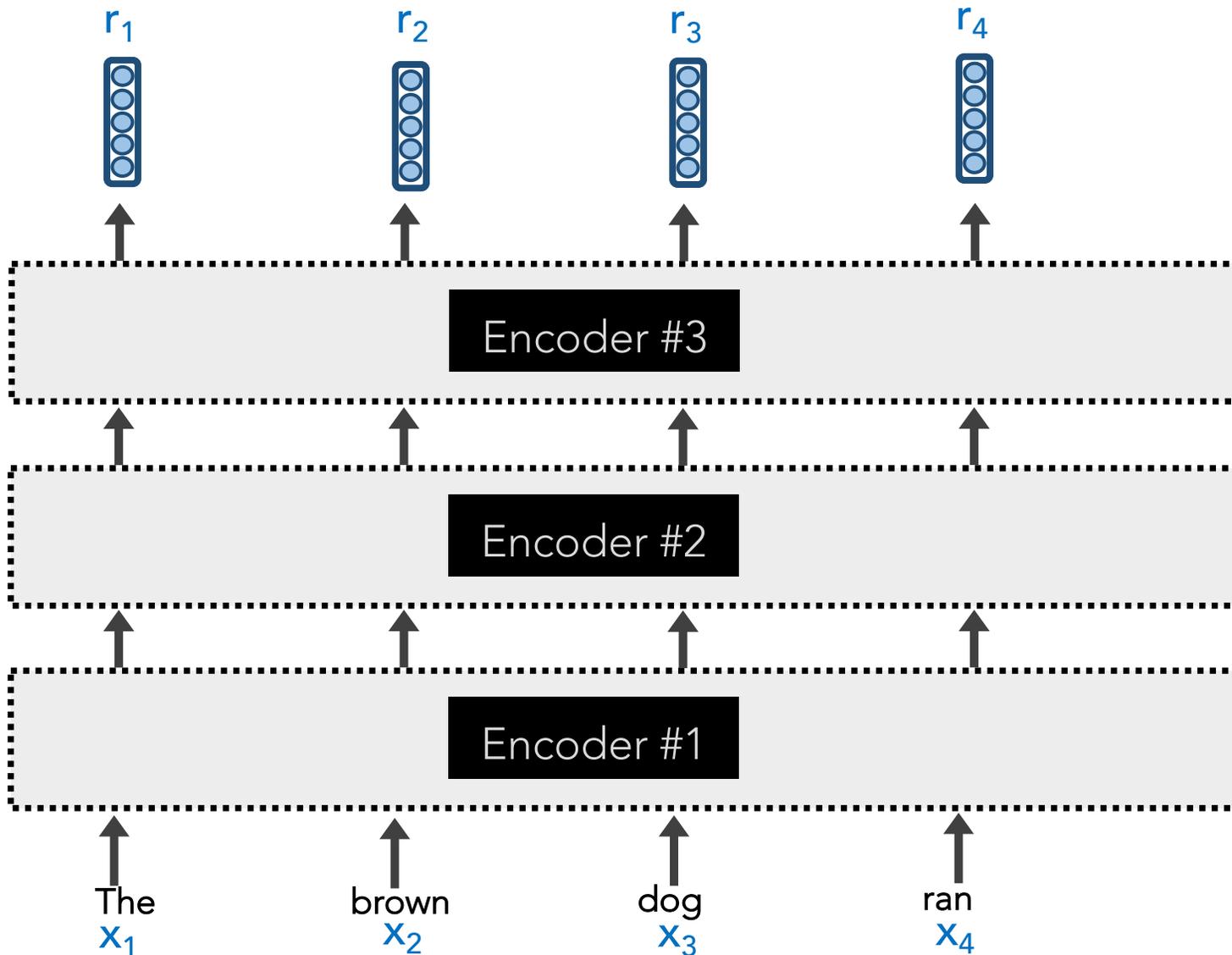
# Transformer Encoder

To recap: all of this looks fancy, but ultimately it's just producing a very good **contextualized embedding**  $r_i$  of each word  $x_i$



Why stop with just 1 **Transformer Encoder**?  
We could stack several!

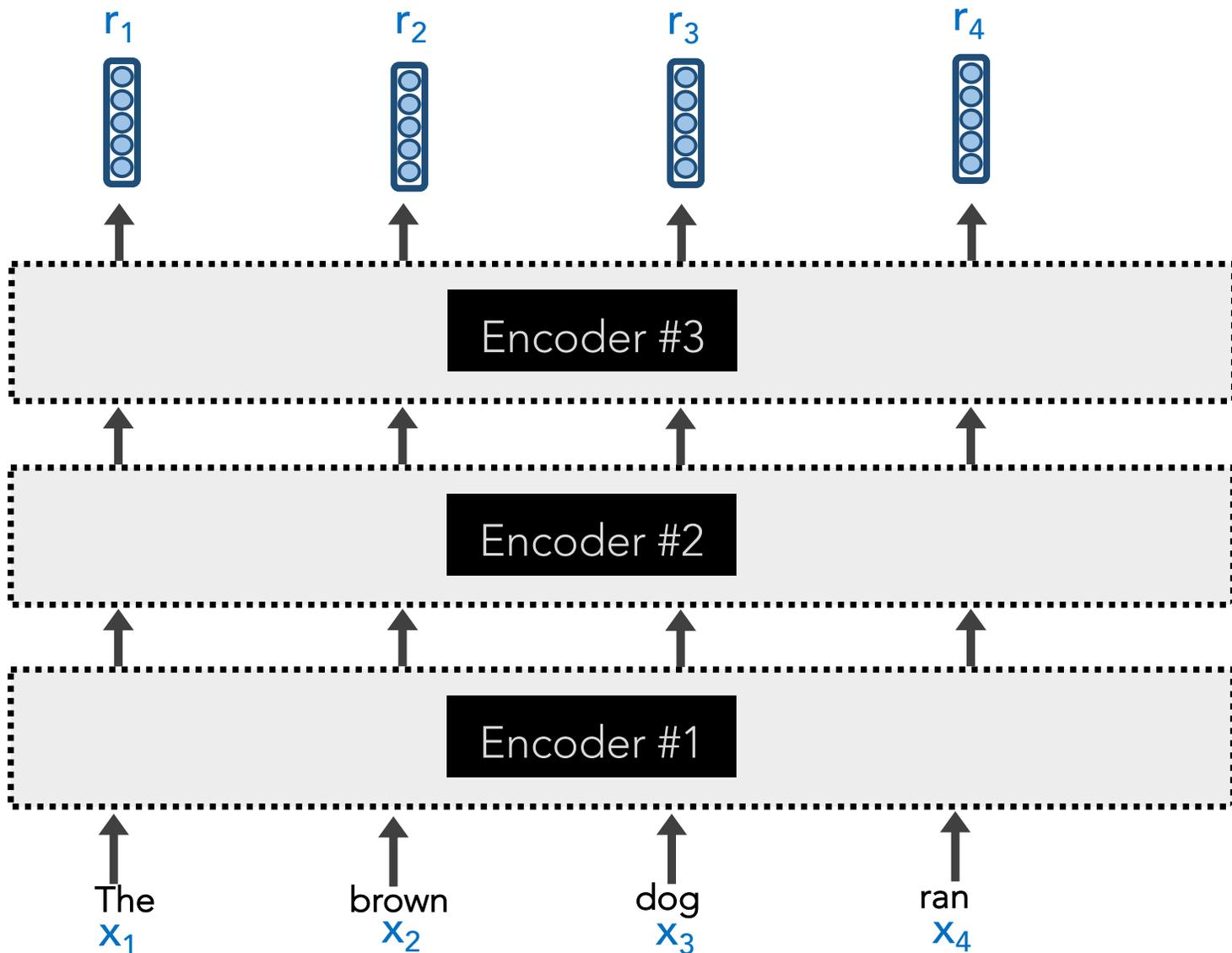
# Transformer Encoder



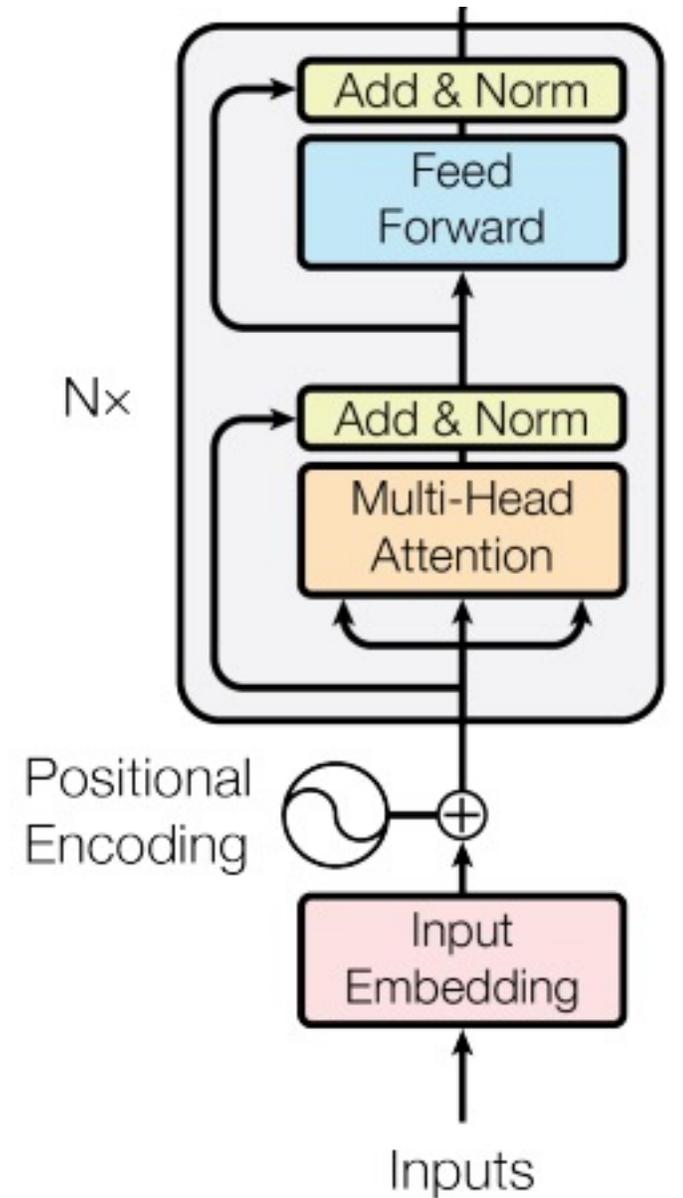
To recap: all of this looks fancy, but ultimately it's just producing a very good **contextualized embedding**  $r_i$  of each word  $x_i$

Why stop with just 1 **Transformer Encoder**?  
We could stack several!

# Transformer Encoder



=



The original Transformer model was intended for Machine Translation, so it had **Decoders**, too

# Outline

 Self-Attention

 Transformer Encoder

 Transformer Decoder

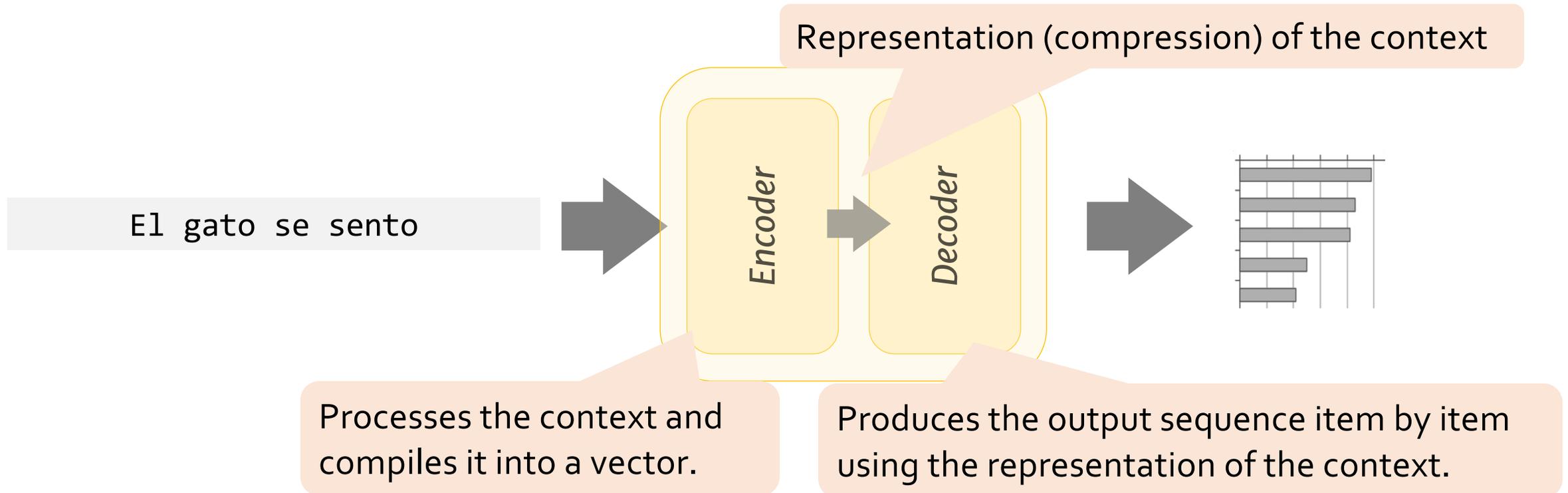
 Language Modeling With  
Transformers

# Outline

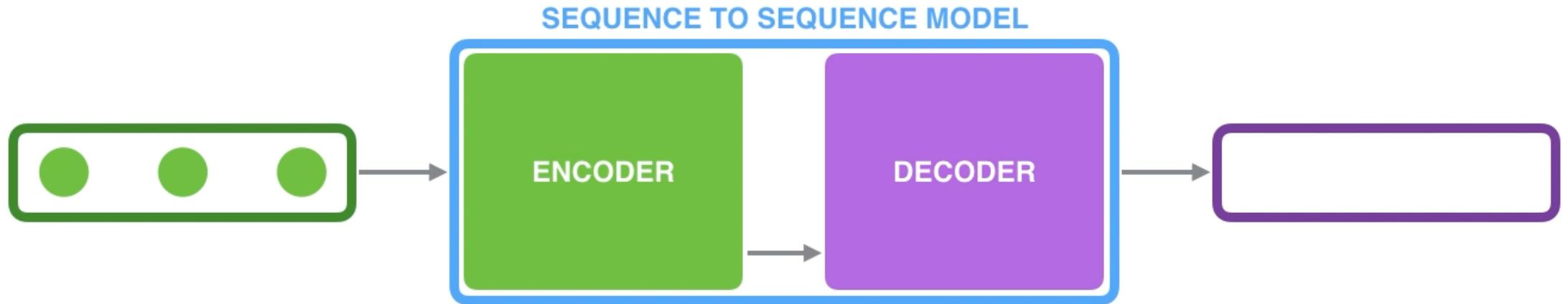
-  Self-Attention
-  Transformer Encoder
-  Transformer Decoder
-  Language Modeling With  
Transformers

# Encoder-Decoder Architectures

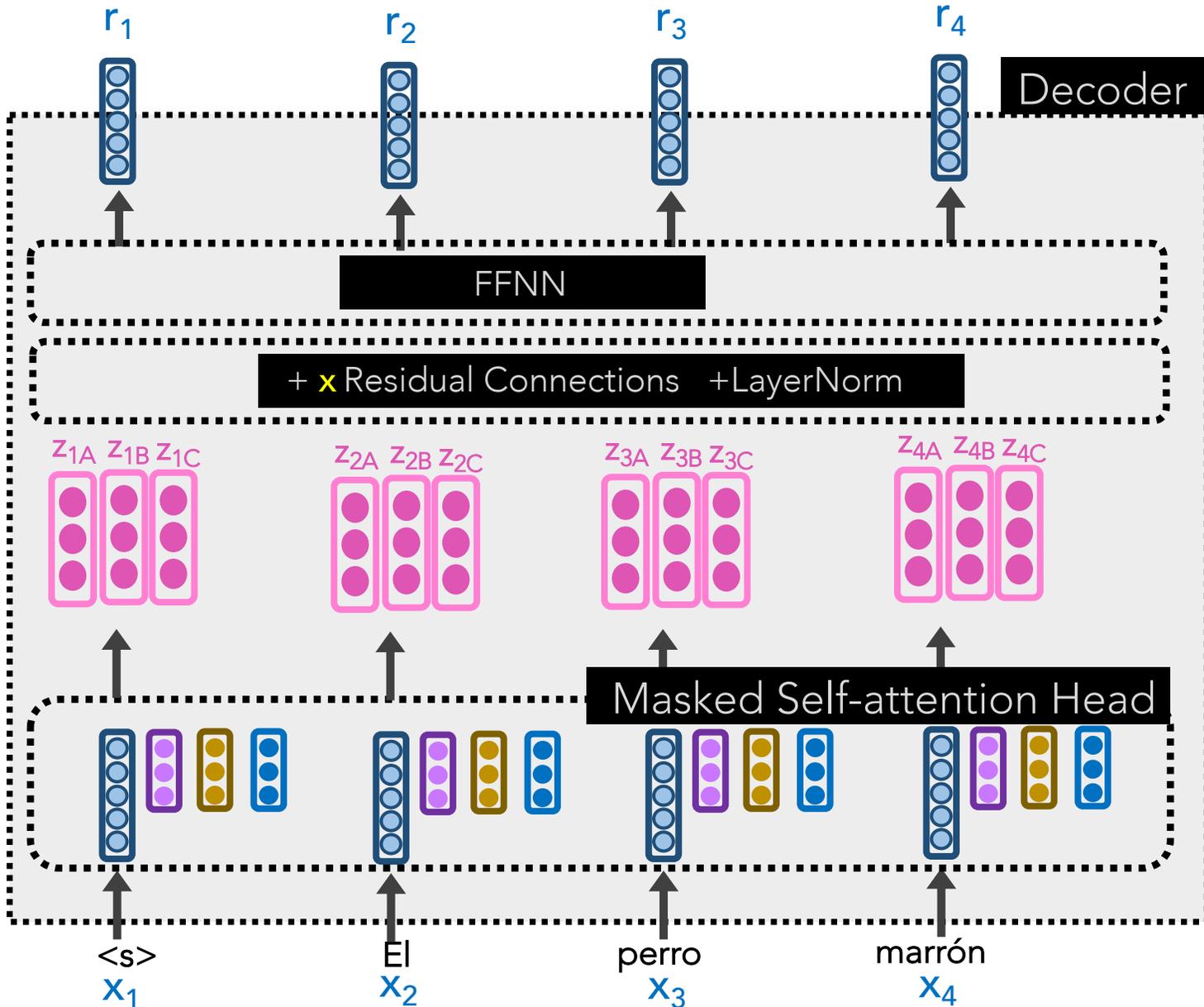
- Original transformer had two sub-models.



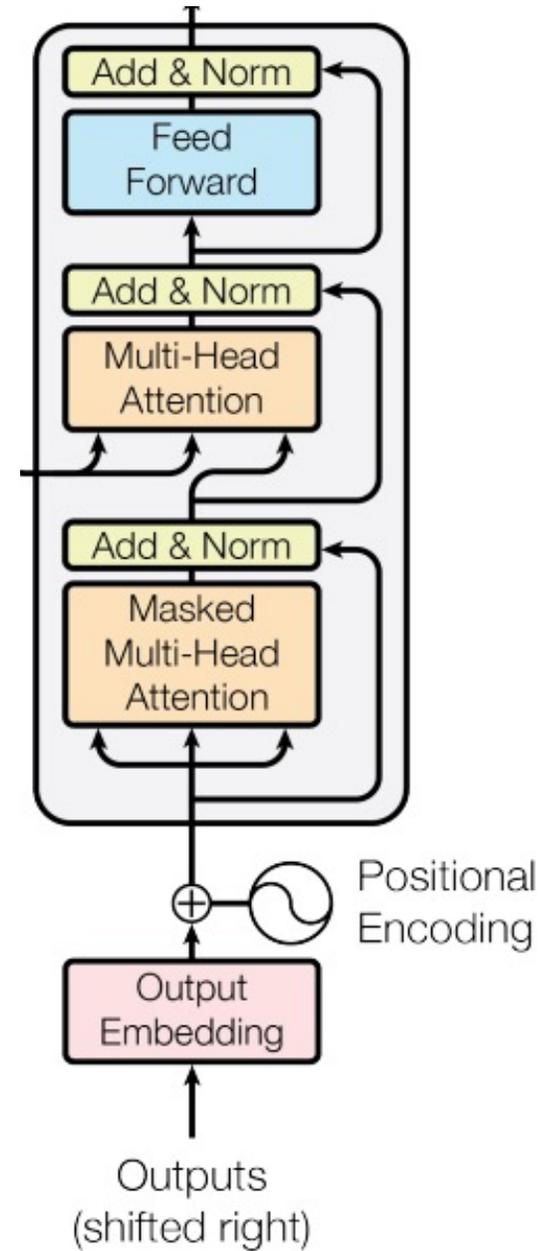
# Encoder-Decoder Architectures



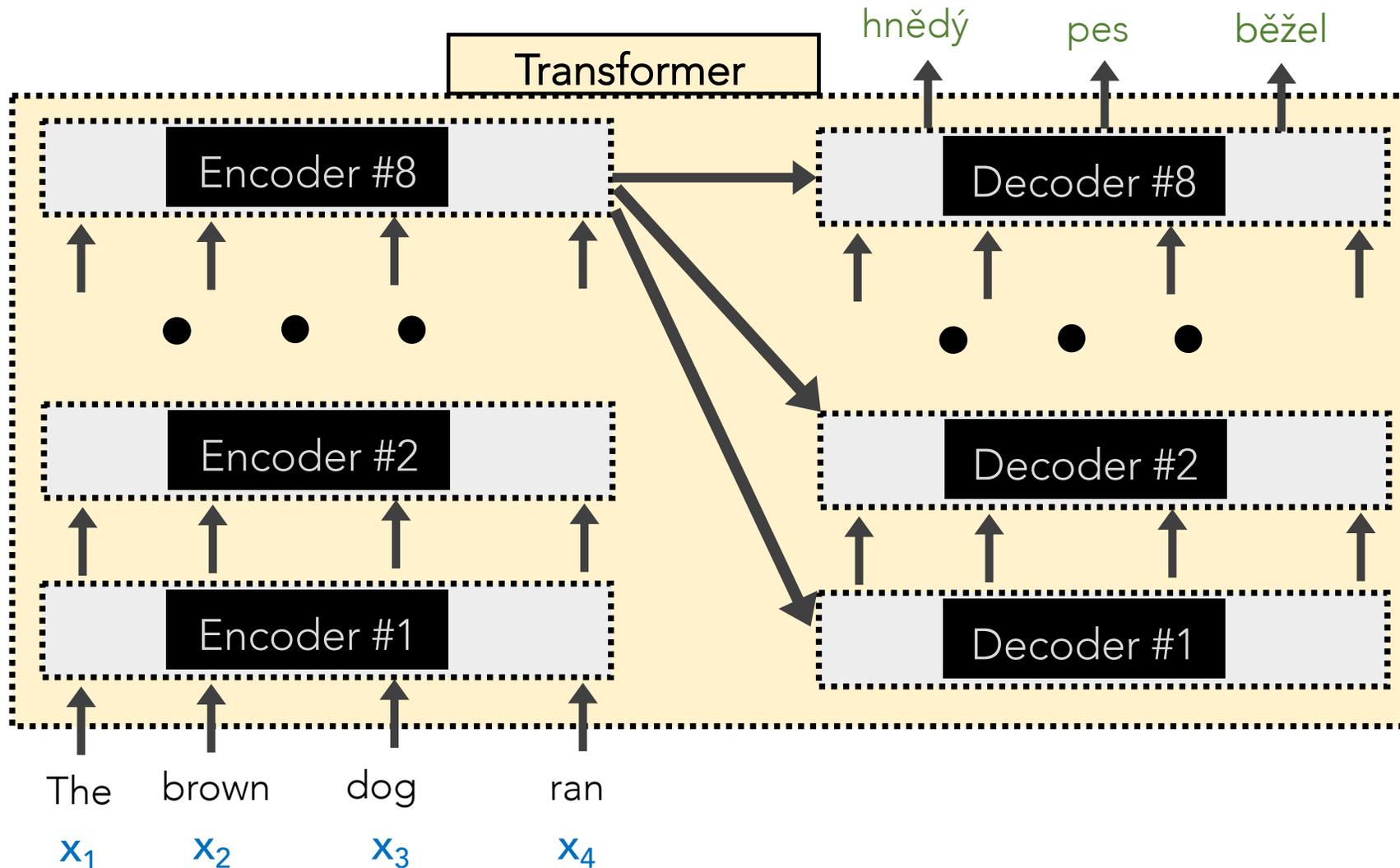
# Transformer Decoder



=



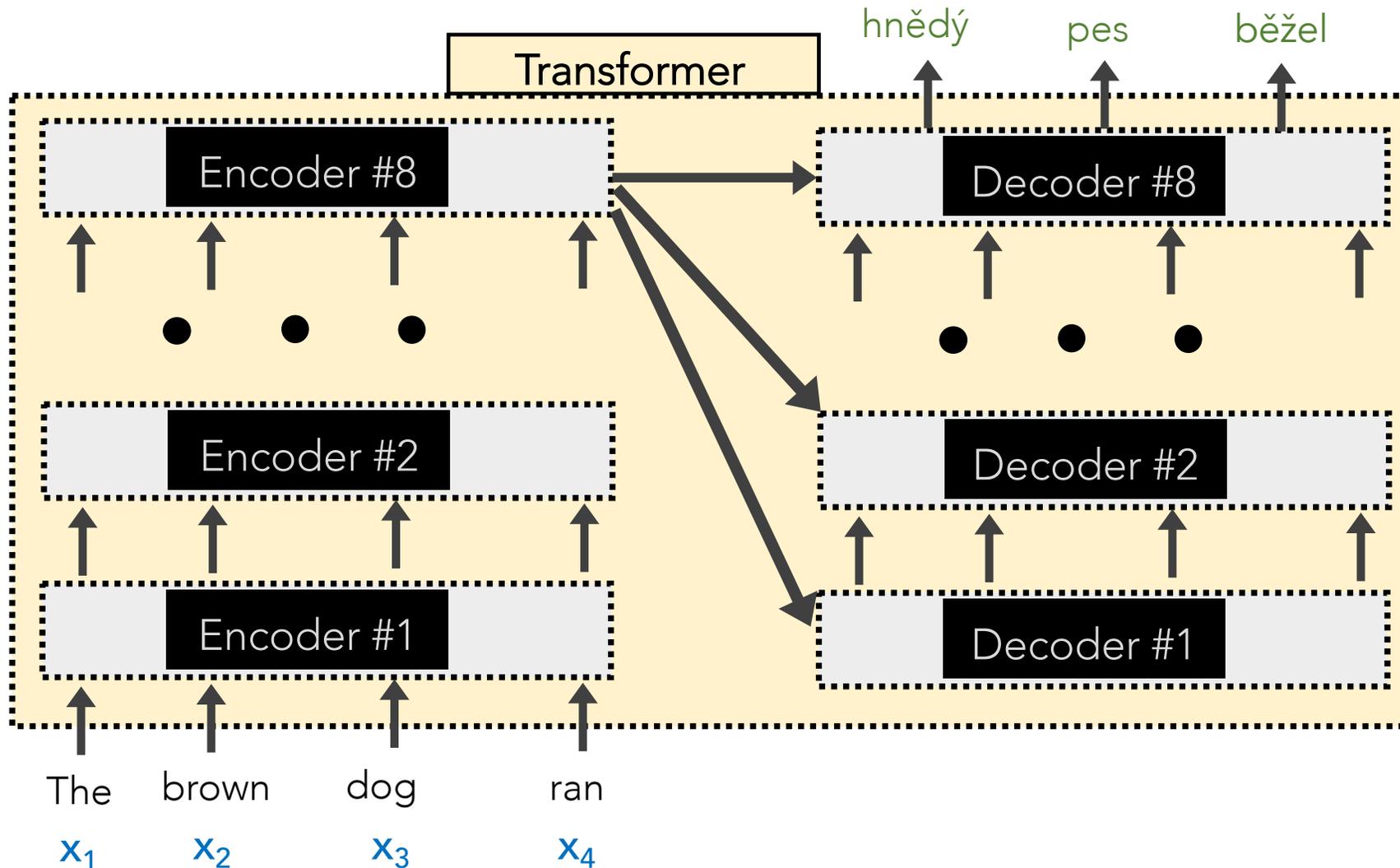
# Transformer Encoders and Decoders



Transformer Encoders produce **contextualized embeddings** of each word

Transformer Decoders generate new sequences of text

# Transformer Encoders and Decoders

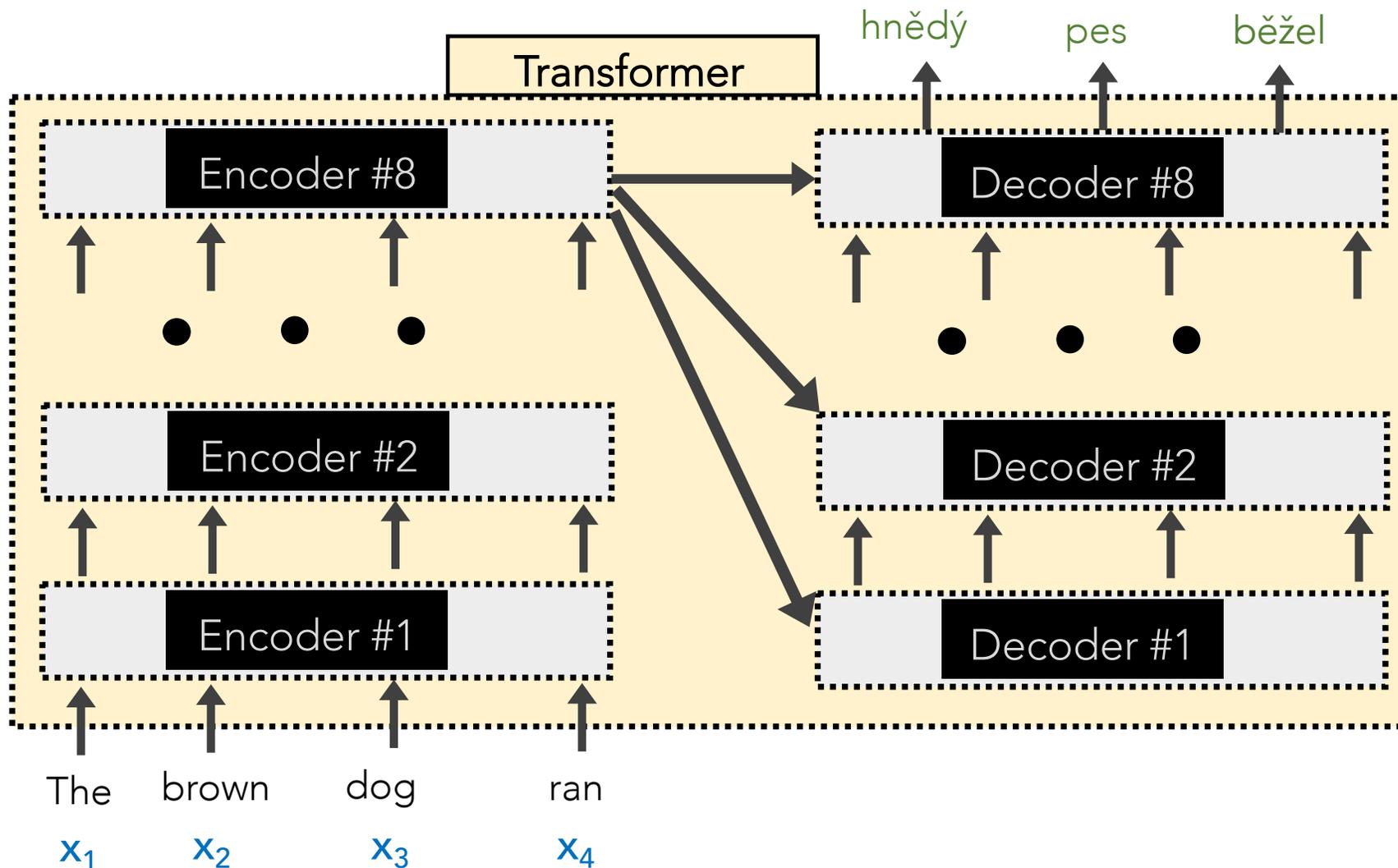


## NOTE

Transformer Decoders are identical to the Encoders, except they have an additional **Attention Head** in between the Self-Attention and FFNN layers.

This additional **Attention Head** focuses on parts of the encoder's representations.

# Transformer Encoders and Decoders

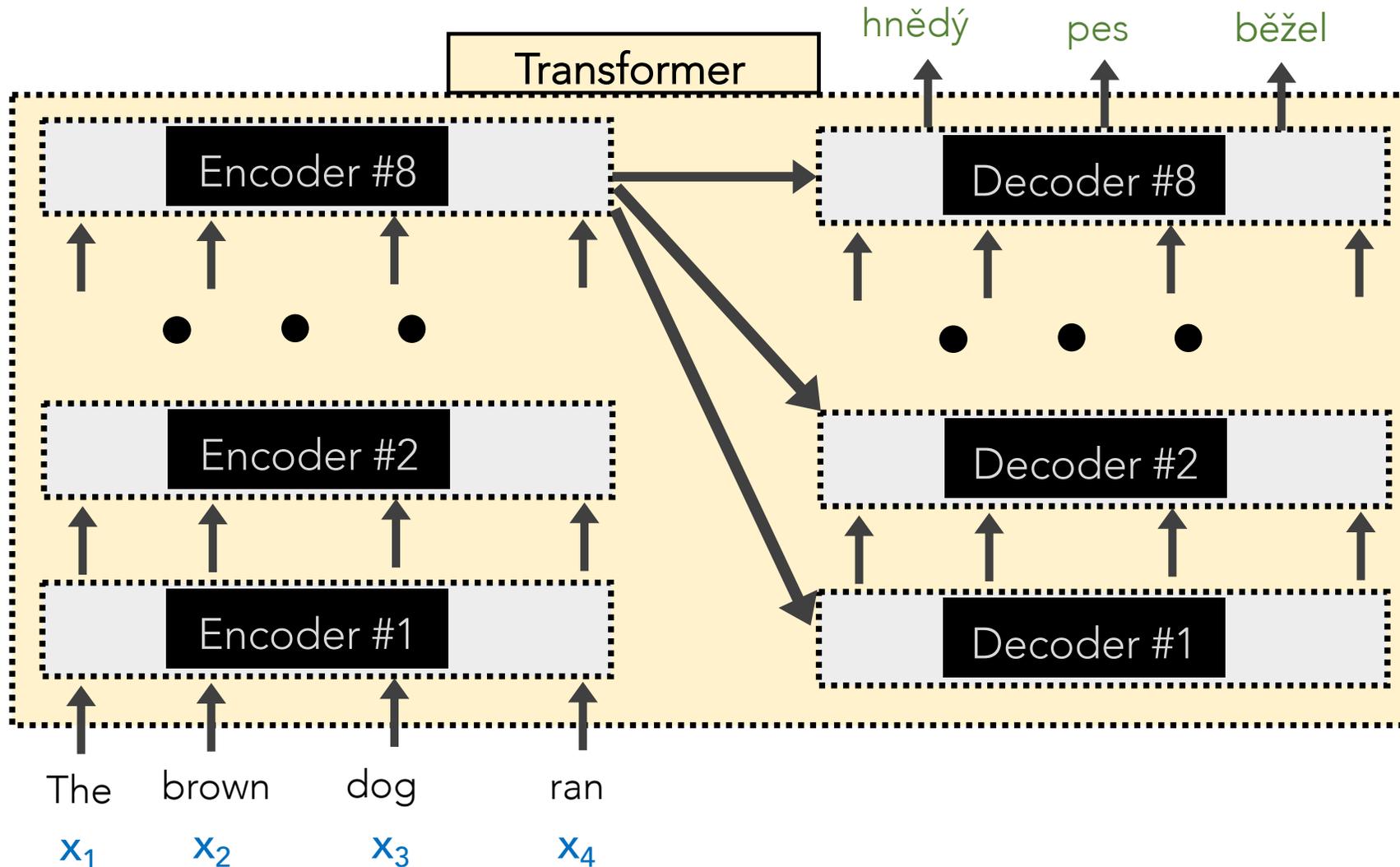


## NOTE

The **query** vector for a Transformer **Decoder's Attention Head** (not Self-Attention Head) is from the output of the previous decoder layer.

However, the **key** and **value** vectors are from the **Transformer Encoders'** outputs.

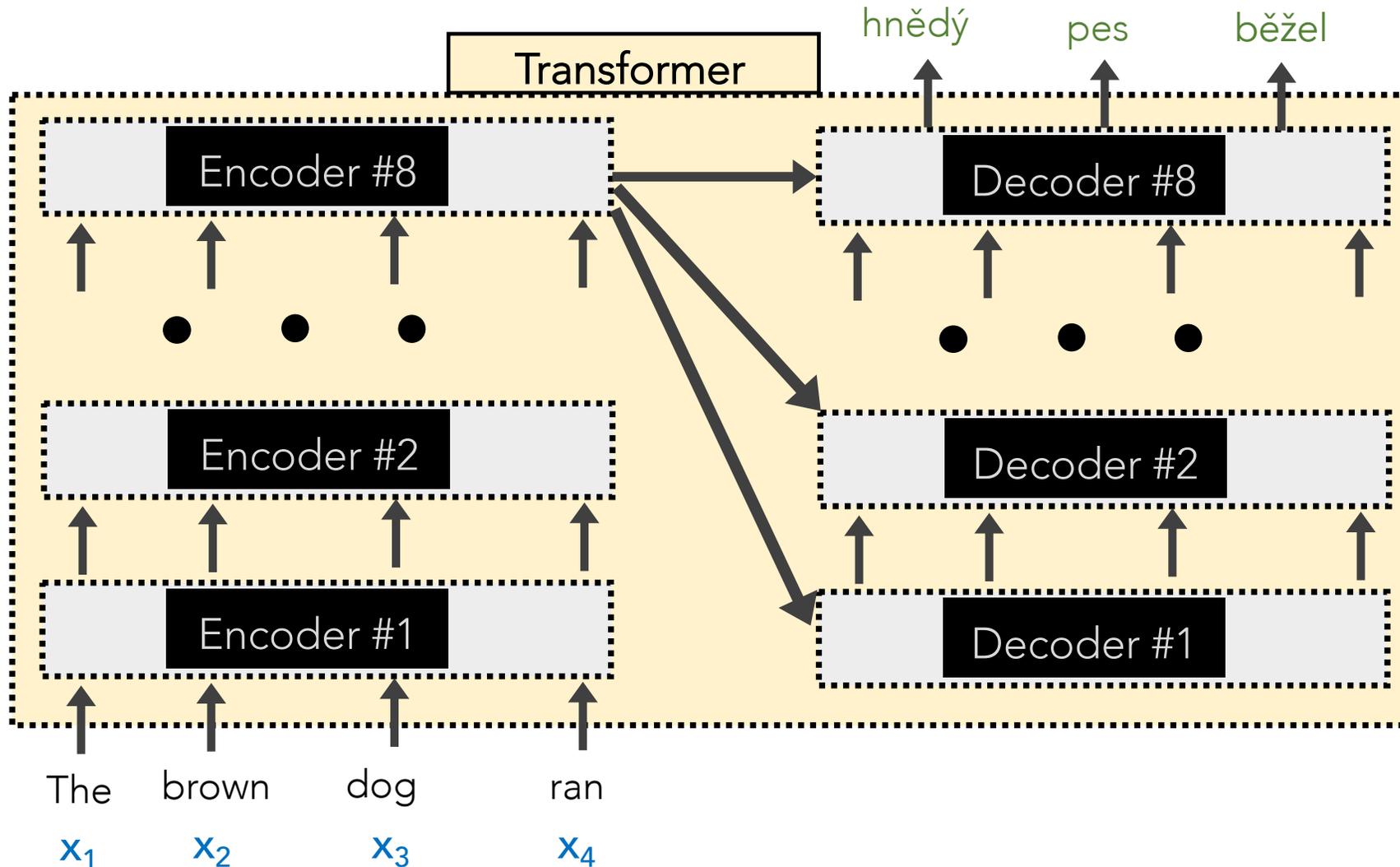
# Transformer Encoders and Decoders



## NOTE

The **query**, **key**, and **value** vectors for a Transformer **Decoder's Self-Attention Head** (not Attention Head) are all from the output of the previous decoder layer.

# Transformer Encoders and Decoders



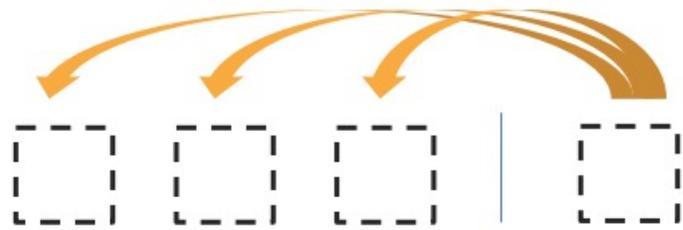
## IMPORTANT

The Transformer **Decoders** have **positional embeddings**, too, just like the **Encoders**.

Critically, each position is **only allowed to attend to the previous indices**. This *masked Attention* preserves it as being an auto-regressive LM.

# Transformer [Vaswani et al. 2017]

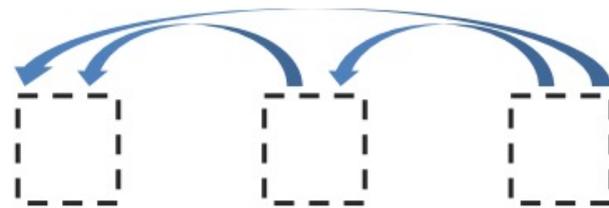
- An **encoder-decoder** architecture built with **attention** modules.
- 3 forms of attention



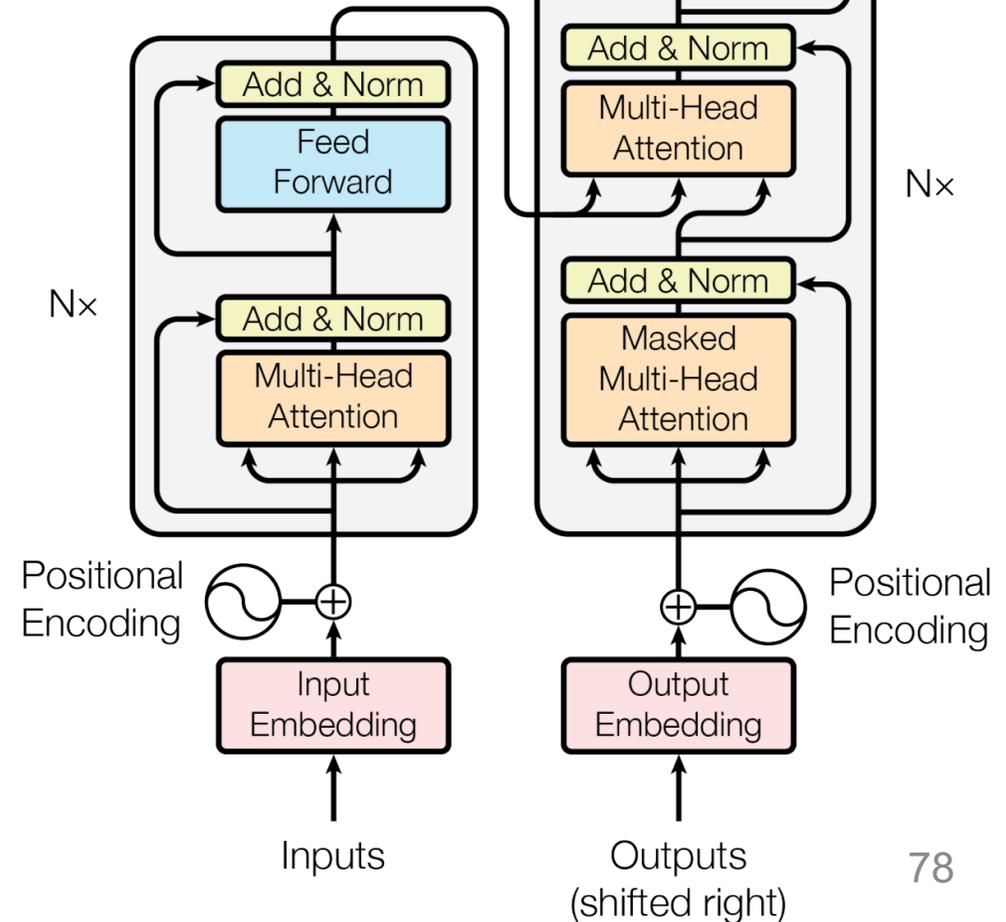
Encoder-Decoder Attention



Encoder Self-Attention

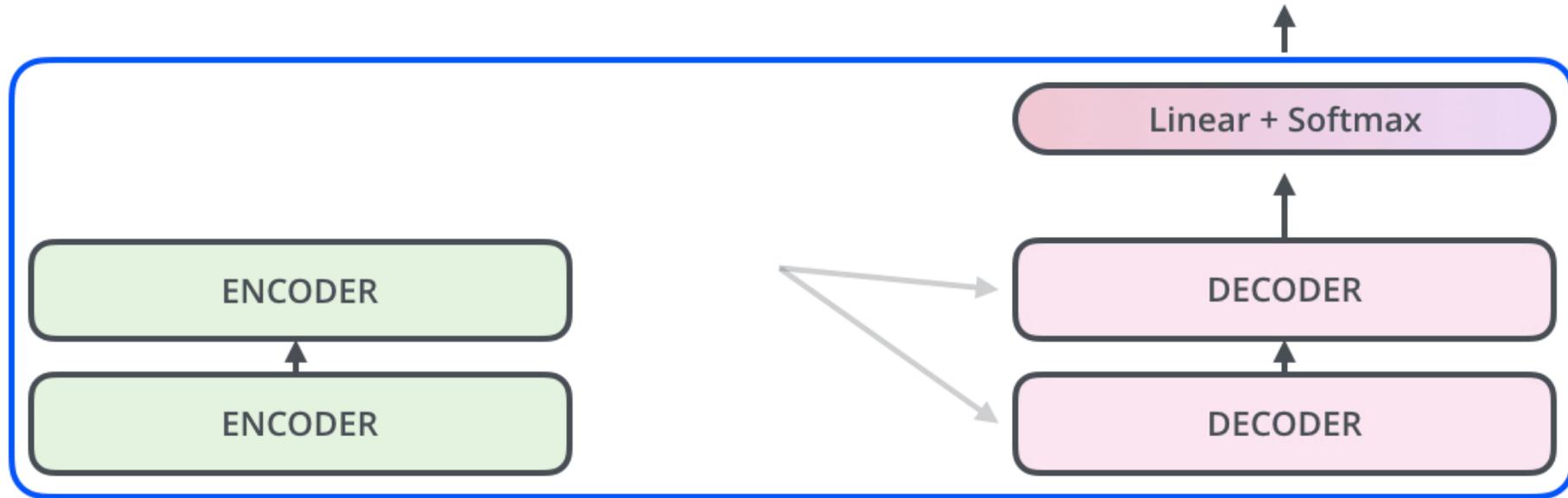


Masked Decoder Self-Attention

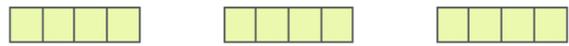


Decoding time step: 1 2 3 4 5 6

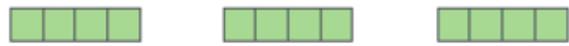
OUTPUT



EMBEDDING WITH TIME SIGNAL



EMBEDDINGS

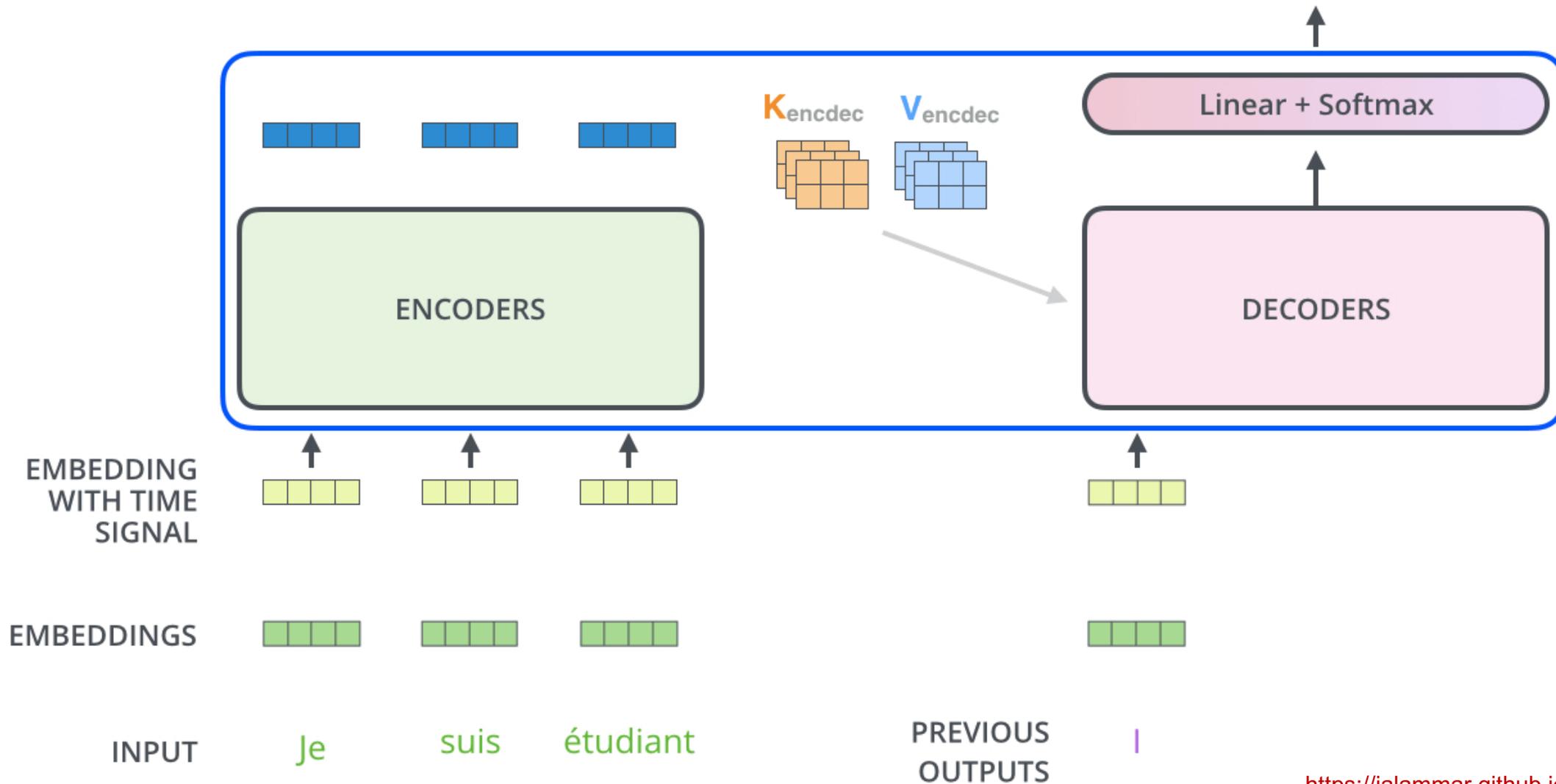


INPUT

Je suis étudiant

Decoding time step: 1 2 3 4 5 6

OUTPUT |



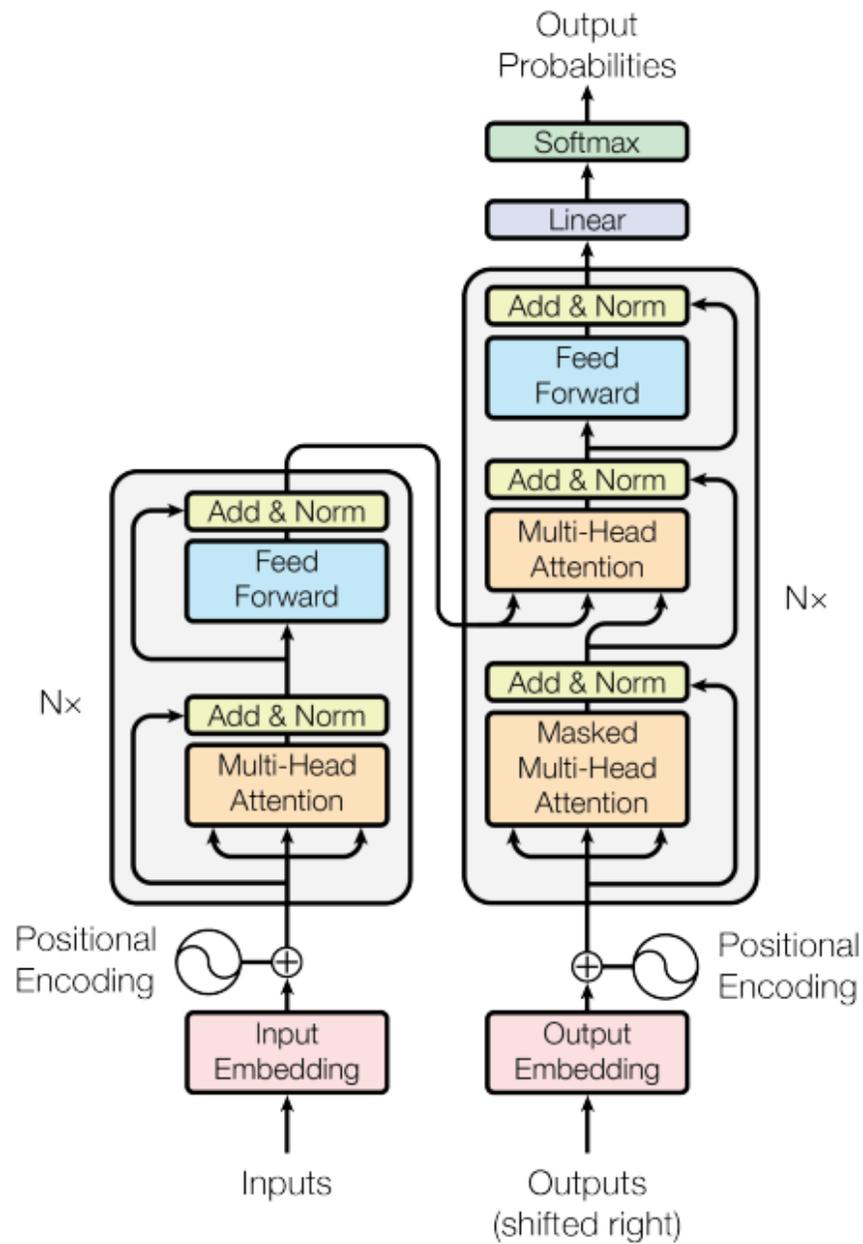


Figure 1: The Transformer - model architecture.

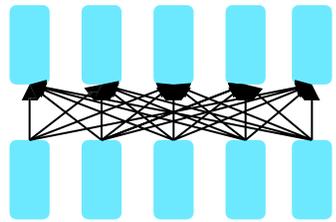
# Impact of Transformers

- Let to better predictive models of language ala GPTs!

Model	Layers	Heads	Perplexity
LSTMs (Grave et al., 2016)	-	-	40.8
QRNNs (Merity et al., 2018)	-	-	33.0
Transformer	16	16	19.8

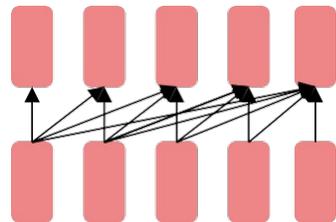
# Impact of Transformers

- A building block for a variety of LMs



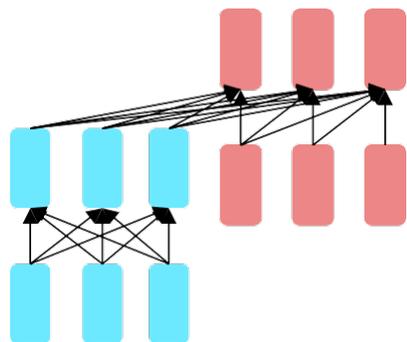
**Encoders**

- ❖ **Examples:** BERT, RoBERTa, SciBERT.
- ❖ Captures bidirectional context. How do we pretrain them?



**Decoders**

- ❖ **Examples:** GPT-2, GPT-3, Llama models, and many many more
- ❖ Other name: **causal or auto-regressive language model**
- ❖ Nice to generate from; can't condition on future words



**Encoder-  
Decoders**

- ❖ **Examples:** Transformer, T5, BART
- ❖ What's the best way to pretrain them?

# Transformer LMs + Scale = LLMs

- 2 main dimensions:
- Model size, pretraining data size

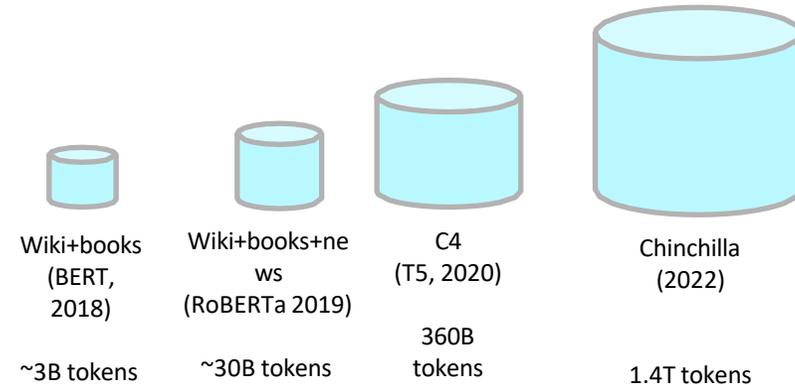
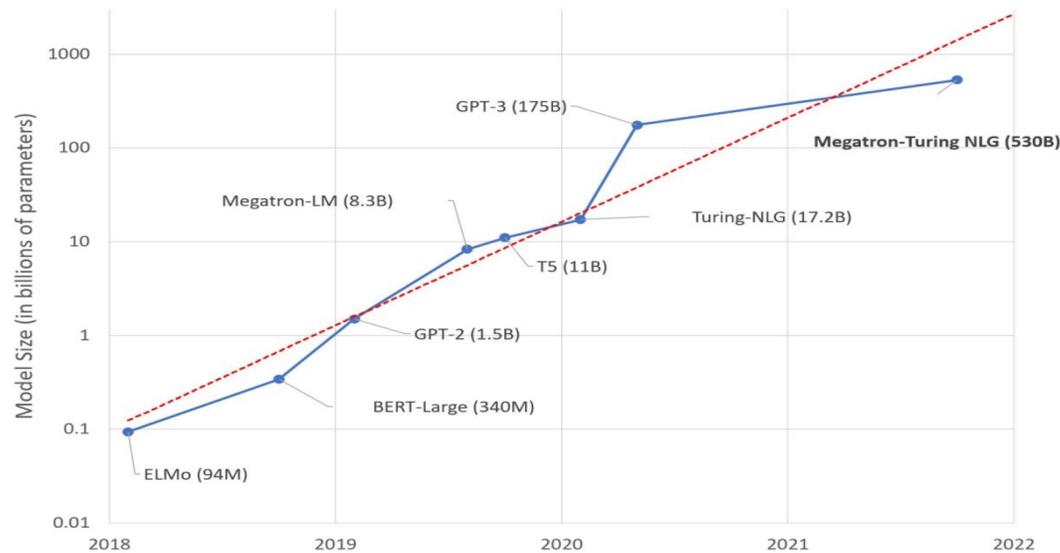
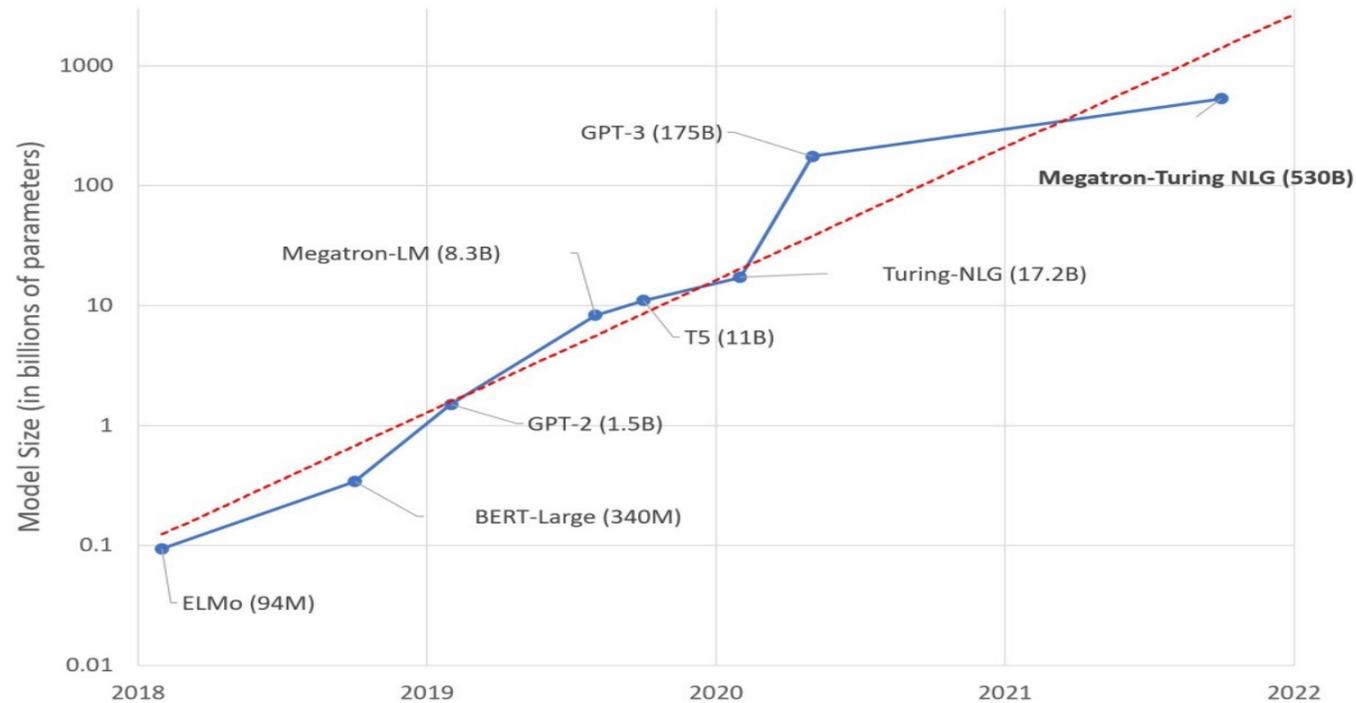


Photo credit: <https://www.microsoft.com/en-us/research/blog/using-deepspeed-and-megatron-to-train-megatron-turing-nlg-530b-the-worlds-largest-and-most-powerful-generative-language-model/>

# Large Language Models

- Not only they improved performance on many NLP tasks, but exhibited new capabilities



# Transformers - Summary

- Self-attention + positional embedding + others = NLP go brr
- Much faster to train than any previous architectures, much easier to scale
- Perform on par or better than previous RNN based models
  - Ease of scaling allows to extract much better performance