

# Computational Linguistics

## LT3233



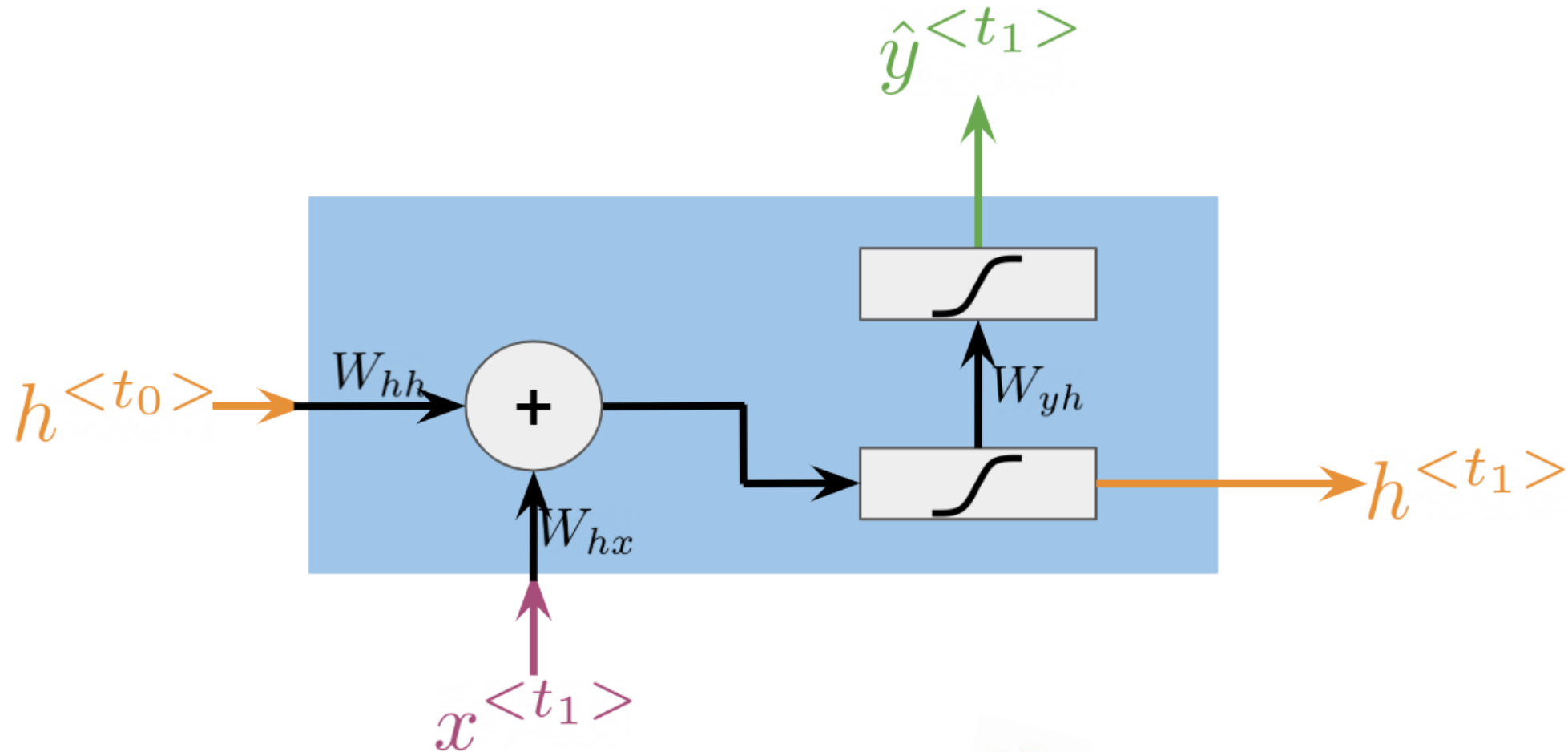
Jixing Li

### Lecture 13: Attention and Transformers

# Lecture plan

- Problems with RNN
- Attention
- Transformers
- Short break (15 mins)
- Hands-on exercises

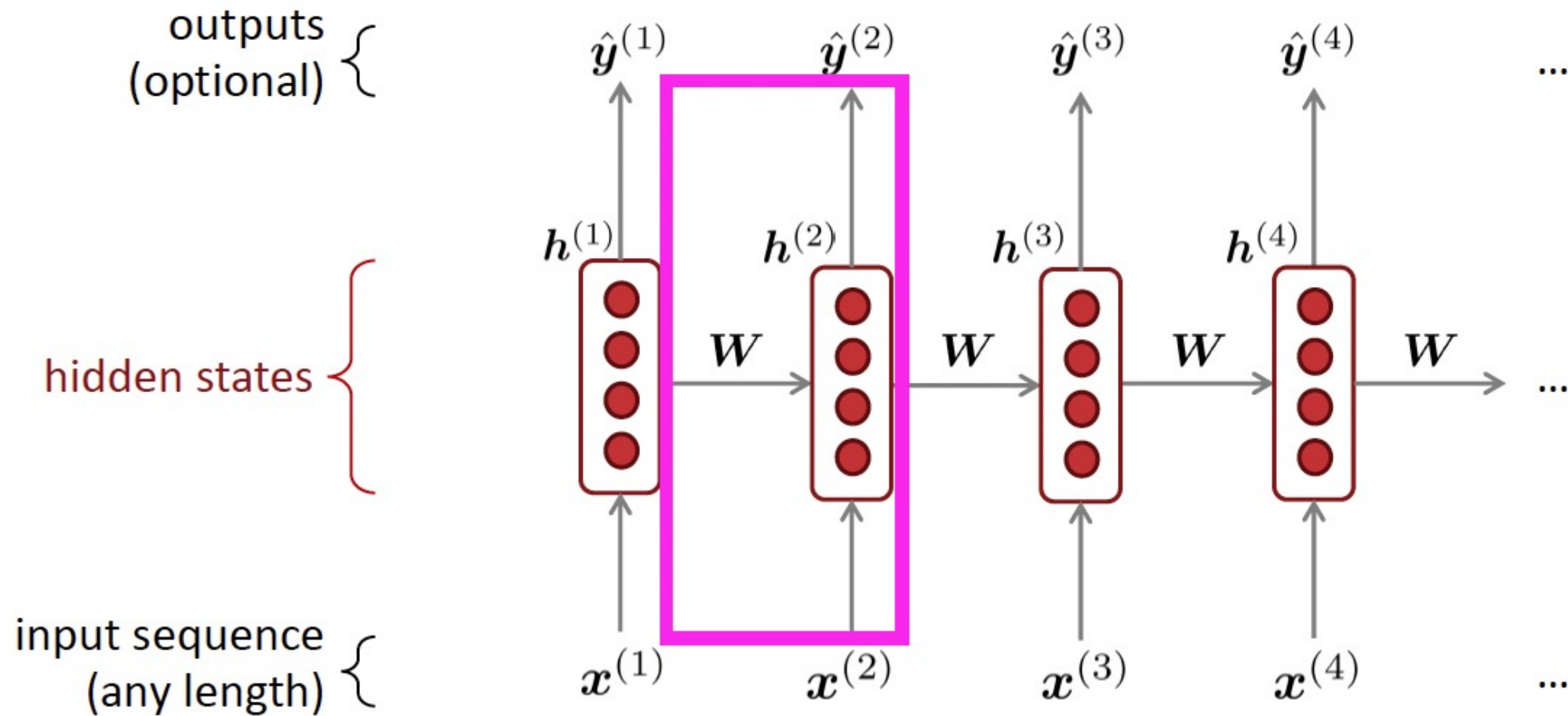
# Recurrent Neural Networks (RNN)



**simple/vanilla/Elman RNN**



# RNN basic structure



# RNN language model

output distribution

$$\hat{y}^{(t)} = \text{softmax}(U\mathbf{h}^{(t)} + \mathbf{b}_2) \in \mathbb{R}^{|V|}$$

hidden states

$$\mathbf{h}^{(t)} = \sigma(\mathbf{W}_h\mathbf{h}^{(t-1)} + \mathbf{W}_e\mathbf{e}^{(t)} + \mathbf{b}_1)$$

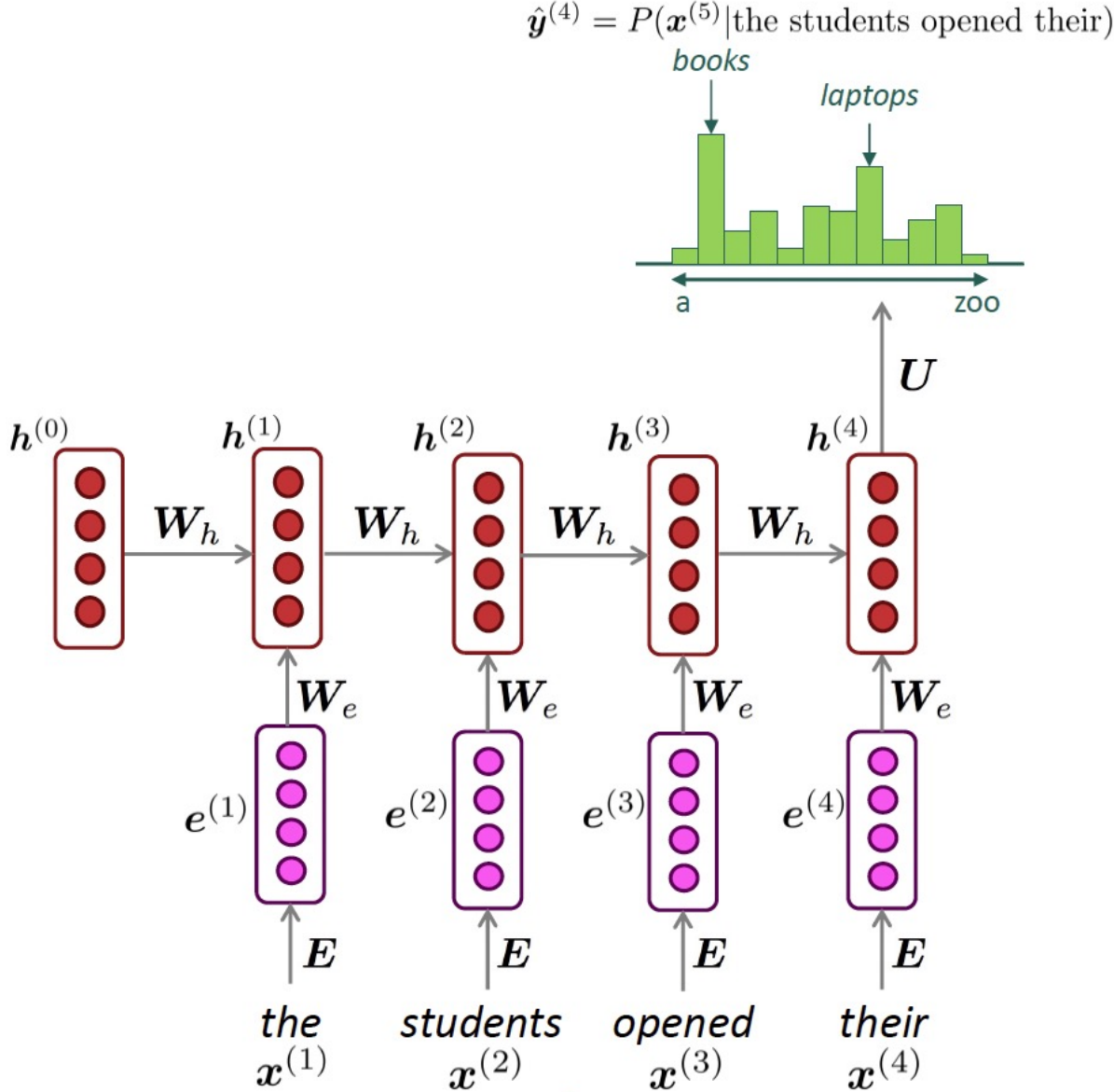
$\mathbf{h}^{(0)}$  is the initial hidden state

word embeddings

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

words / one-hot vectors

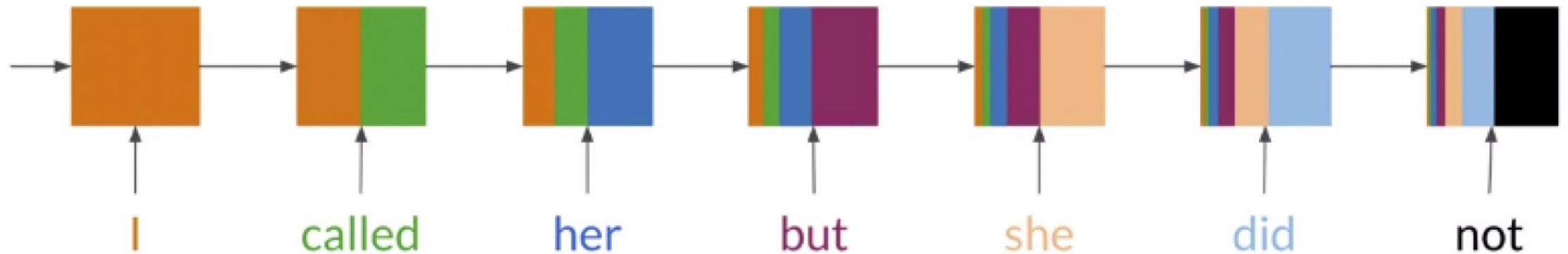
$$\mathbf{x}^{(t)} \in \mathbb{R}^{|V|}$$



# RNN language model

## Advantages:

- Can process any length input
- Computation for step  $t$  can (in theory) use information from many steps back
- Model size doesn't increase for longer input context: Same weights applied on every timestep

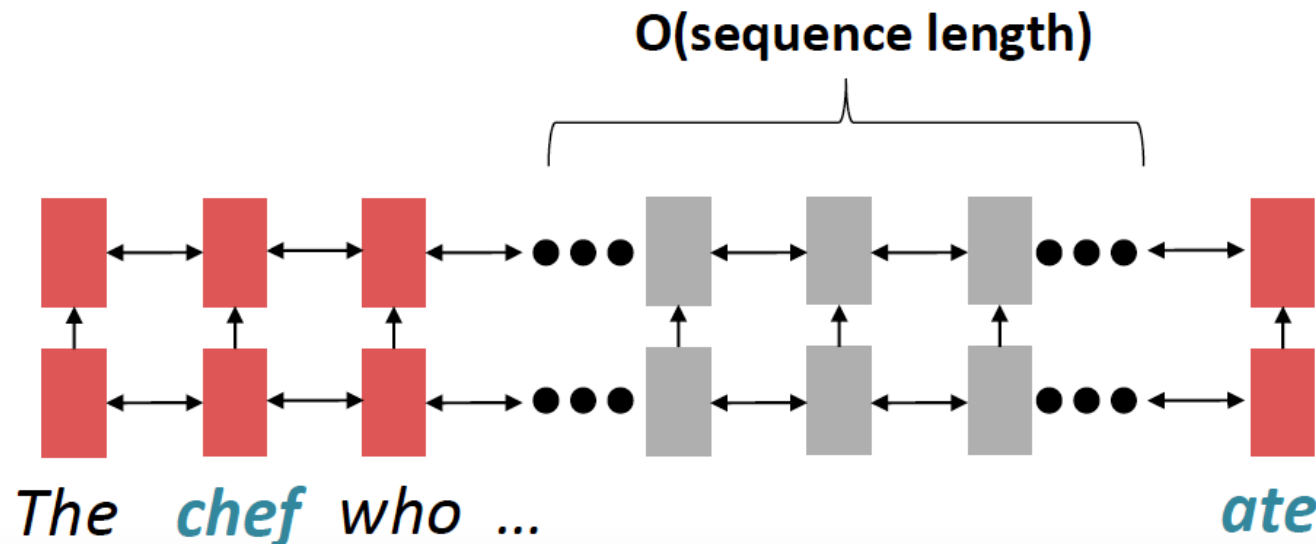


# RNN language model

## Disadvantages:

RNNs take  **$O(\text{sequence length})$**  steps for distant word pairs to interact.

- Hard to learn long-distance dependencies (vanishing gradients)
- Forward and backward passes have  **$O$**  unparallelizable operations:
  - future RNN hidden states can't be computed before past RNN hidden states have been computed



Info of *chef* has gone through  $O(\text{sequence length})$  many layers!

# Attention

---

## 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

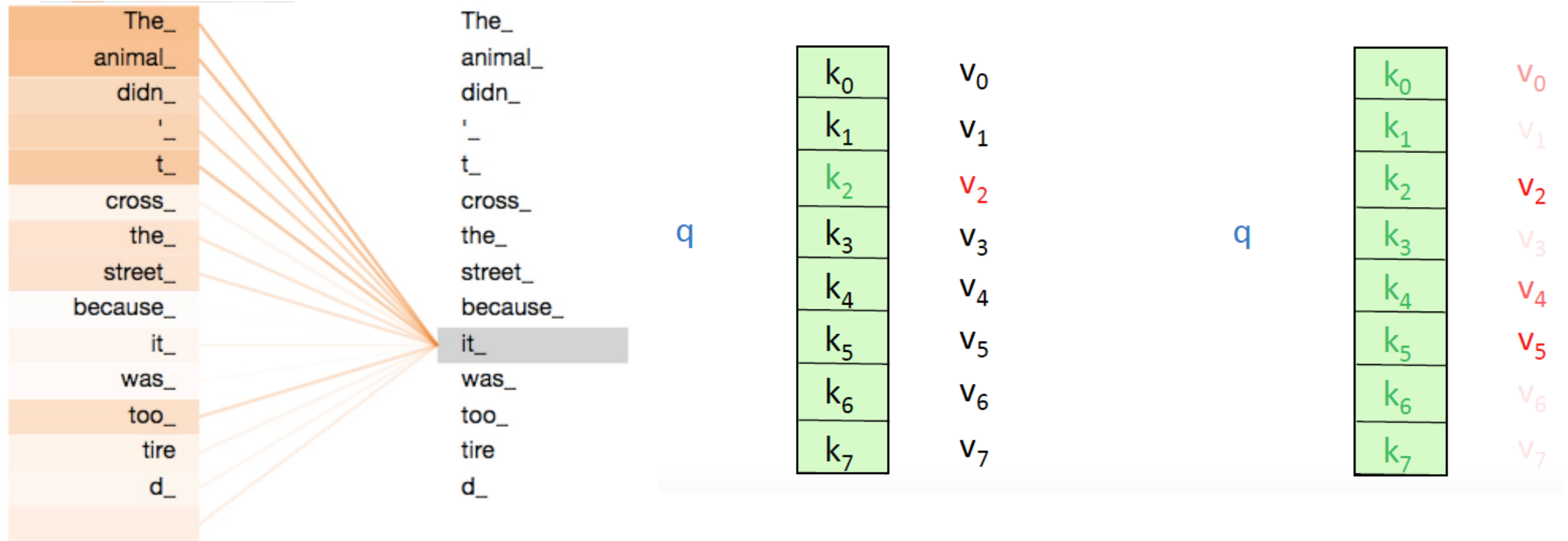
Idea: Compare an item of interest to a collection of other items in a way that reveals their relevance in the current context.

Attention is the core innovation in **Transformer** architecture which revolutionizes deep learning



# Attention intuition

"*The animal didn't cross the street because it was too tired.*"



To look up a **value**, we compare a **query** against **keys** in a table. Each **query** matches each **key** to varying degrees. We return a sum of **values** weighted by the **query-key** match.

# Self-Attention

1. For each word  $x_i$ , calculate its **query**, **key**, and **value**.

$$q_i = W^Q x_i \quad k_i = W^K x_i \quad v_i = W^V x_i$$

2. Calculate attention score between **query** and **keys**

$$e_{ij} = q_i \cdot k_j$$

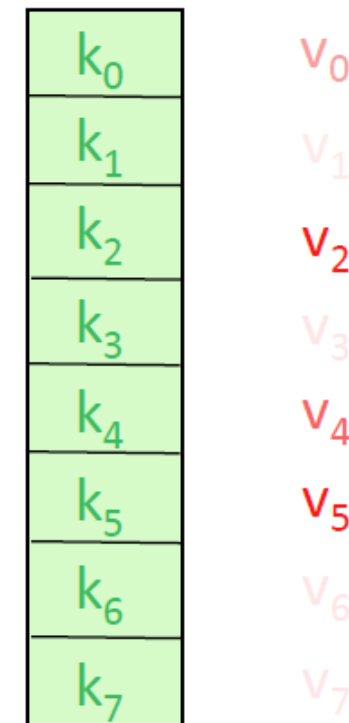
q

3. Take the softmax to normalize attention scores.

$$\alpha_{ij} = \text{softmax}(e_{ij}) = \frac{\exp(e_{ij})}{\sum_k \exp(e_{ik})}$$

4. Take a weighted sum of **values**.

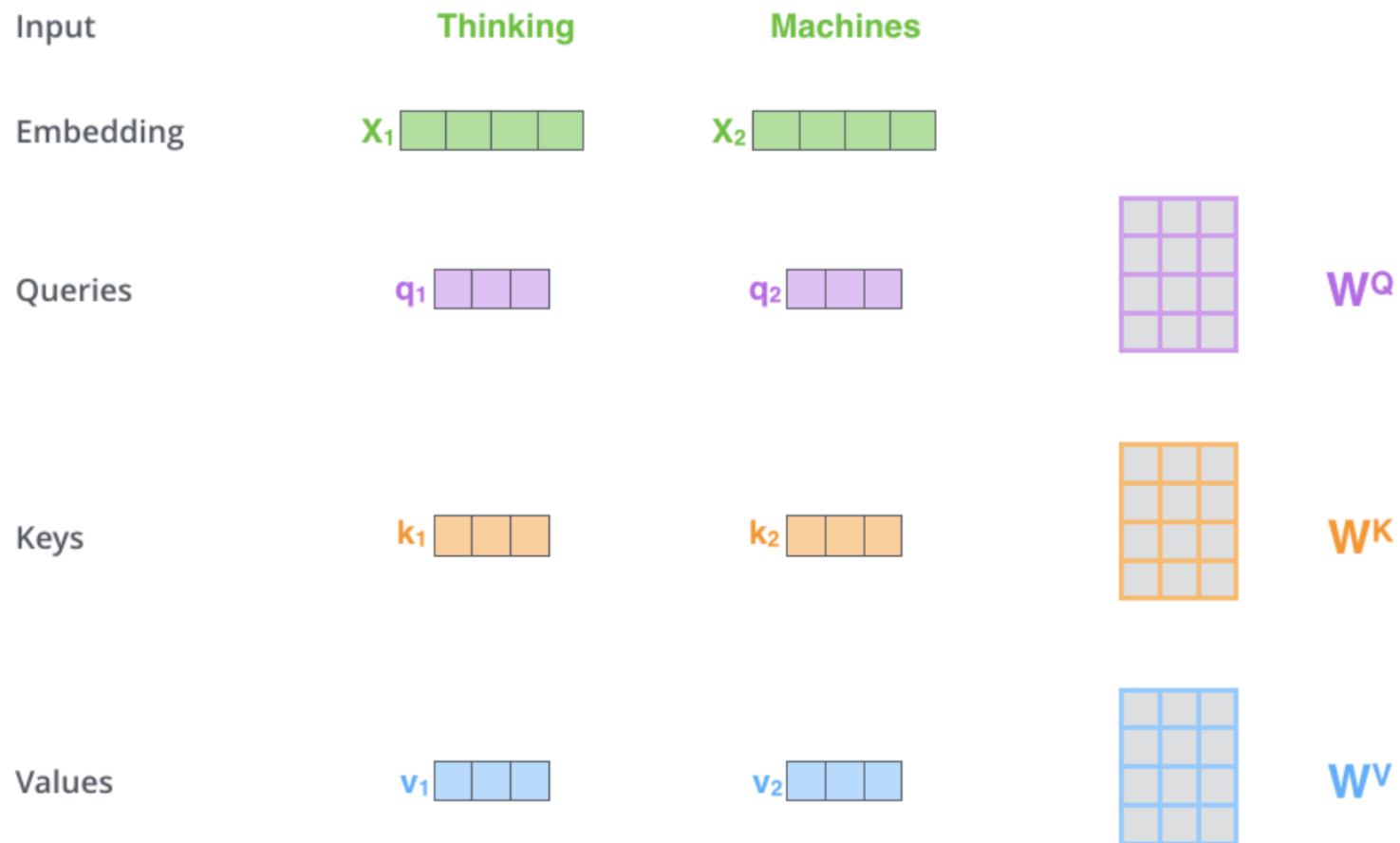
$$\text{Output}_i = \sum_j \alpha_{ij} v_j$$



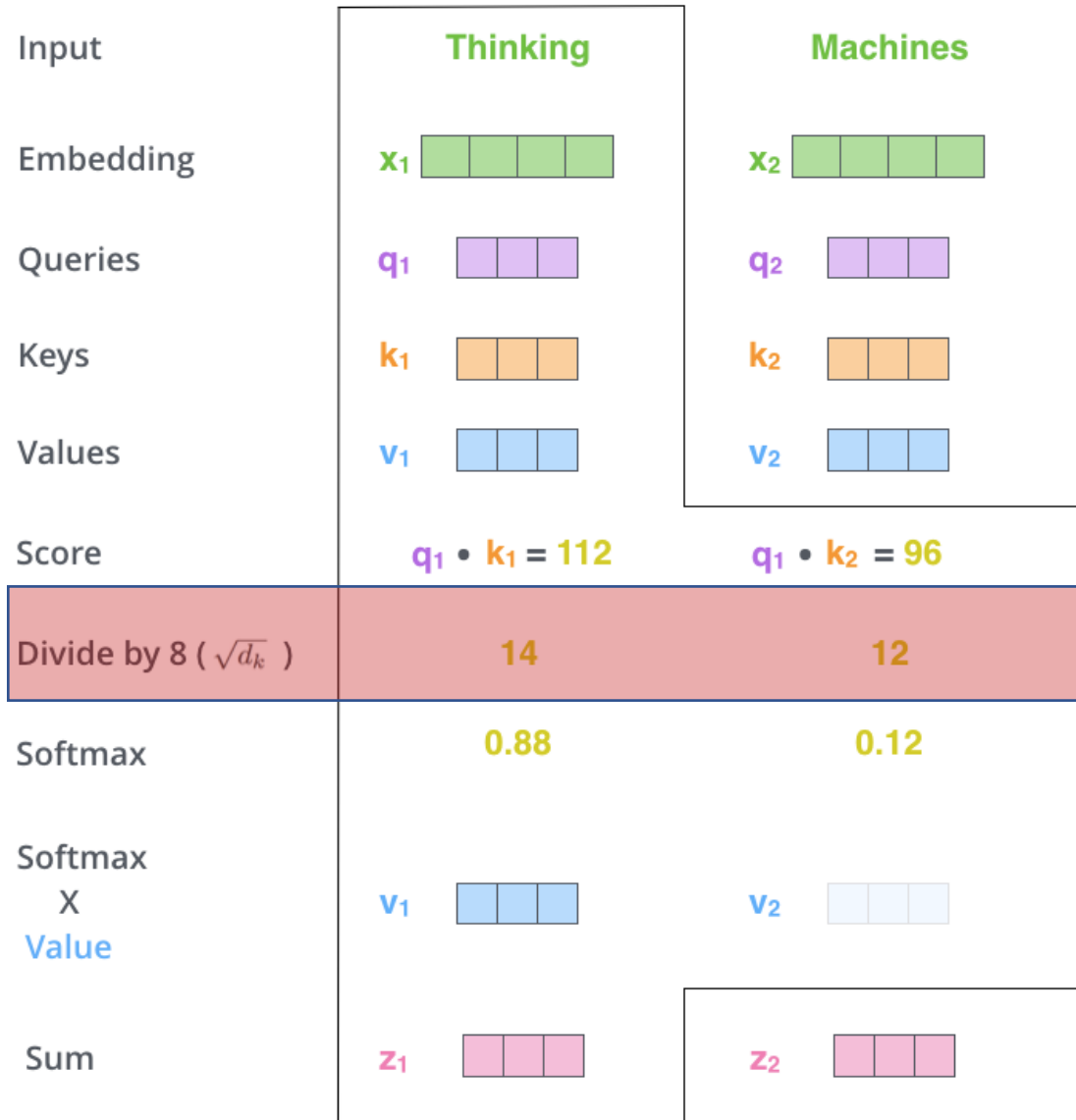
# Illustrated self-attention

For each word  $x_i$ , calculate its **query**, **key**, and **value**.

$$q_i = W^Q x_i \quad k_i = W^K x_i \quad v_i = W^V x_i$$



# Illustrated self-attention



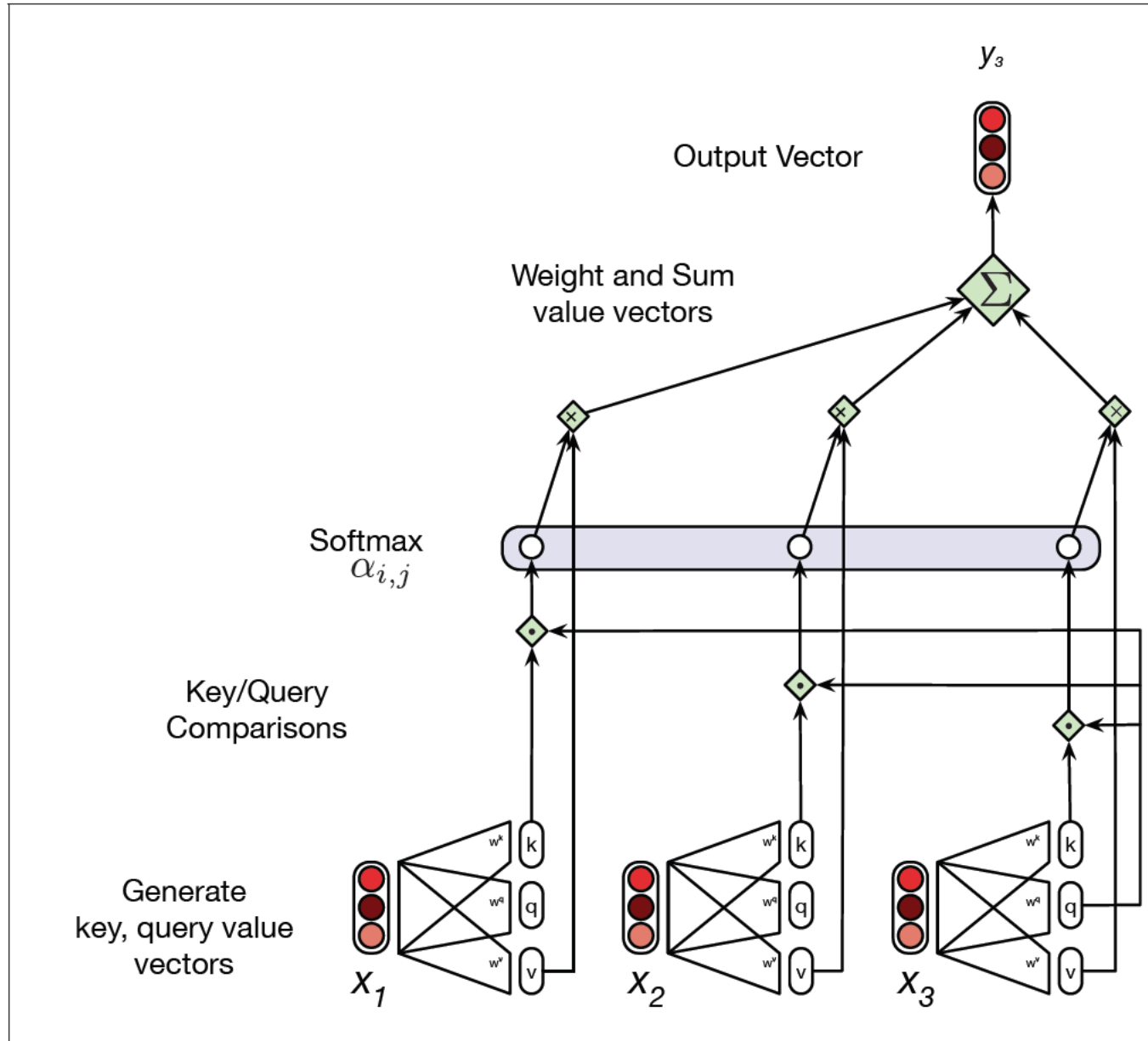
Calculate attention score between query and keys.  $e_{ij} = q_i \cdot k_j$

Take the softmax to normalize attention scores.

Take a weighted sum of values.

$$Output_i = \sum_j \alpha_{ij} v_j$$

# Illustrated self-attention

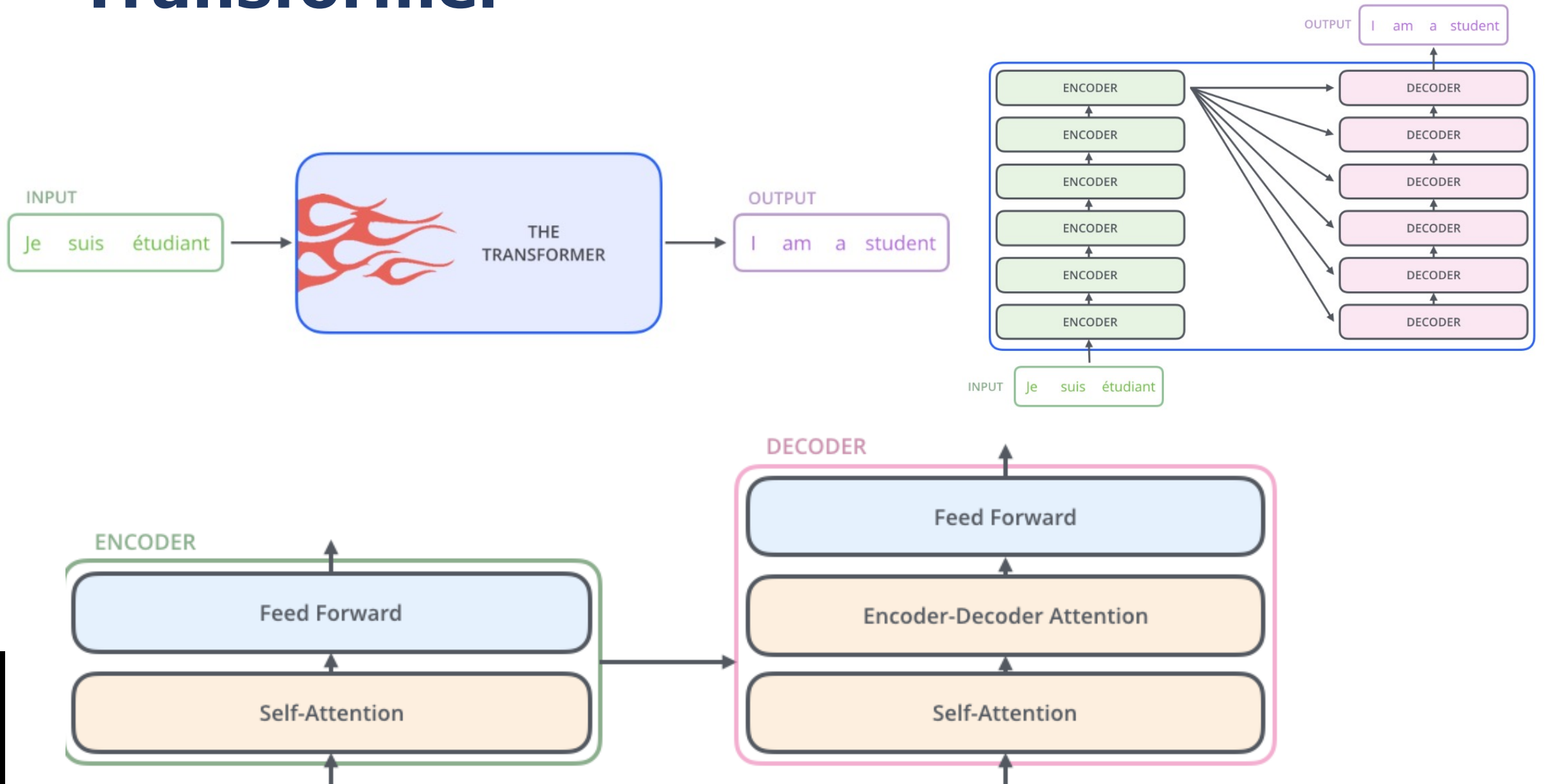


Calculating the value of  $y_3$ , the third element of a sequence using self-attention.

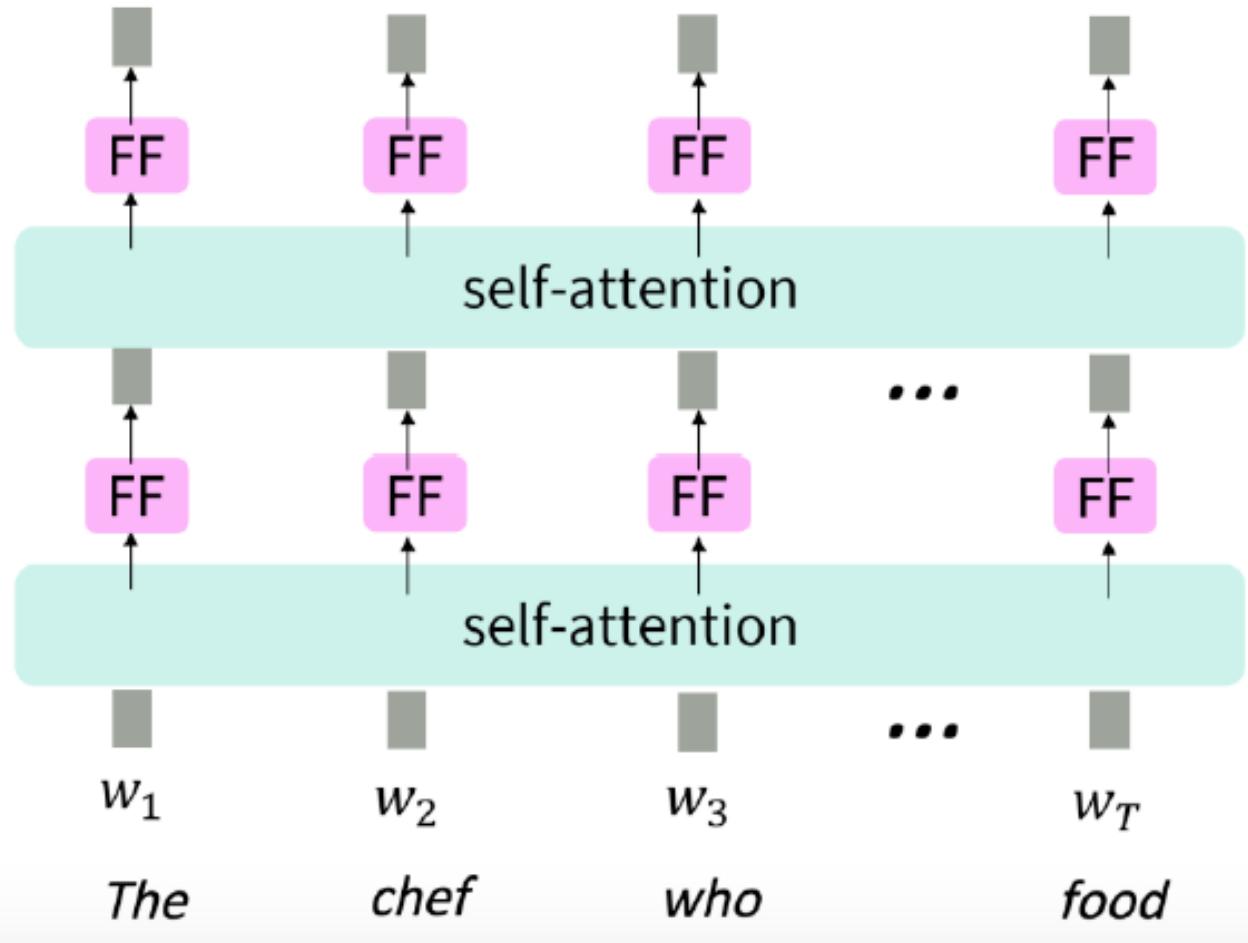
# Summary

- Attention is an added layer that lets a model **focus on what's important**
- **Queries, Values, and Keys** are used for information retrieval inside the Attention layer
- Attention has become the powerful, flexible, general way pointer and memory manipulation in all deep learning models.

# Transformer



# Feed-forward layer



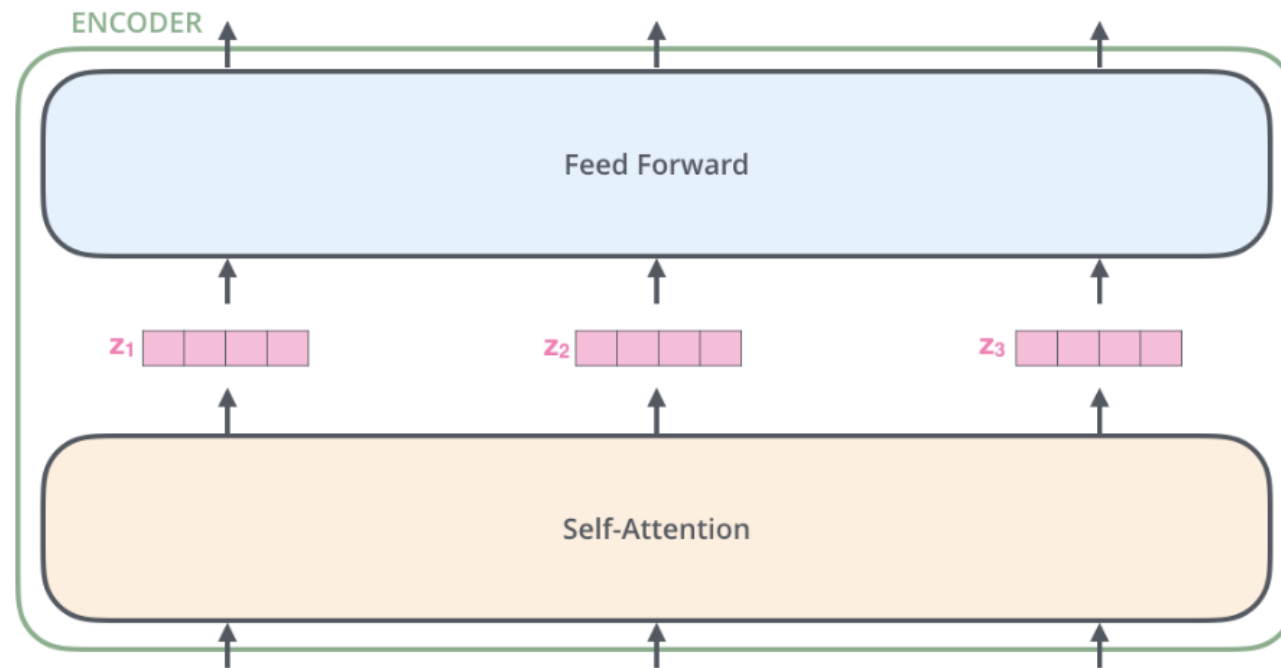
Since there are **no element-wise non-linearities**, self-attention is simply performing a re-averaging of the value vectors.

→ Apply a feedforward layer to the output of attention, providing **non-linear activation (and additional expressive power)**.

$$W_2 * \text{ReLU}(W_1 \times \text{output}_i + b_1) + b_2$$



# Positional embeddings

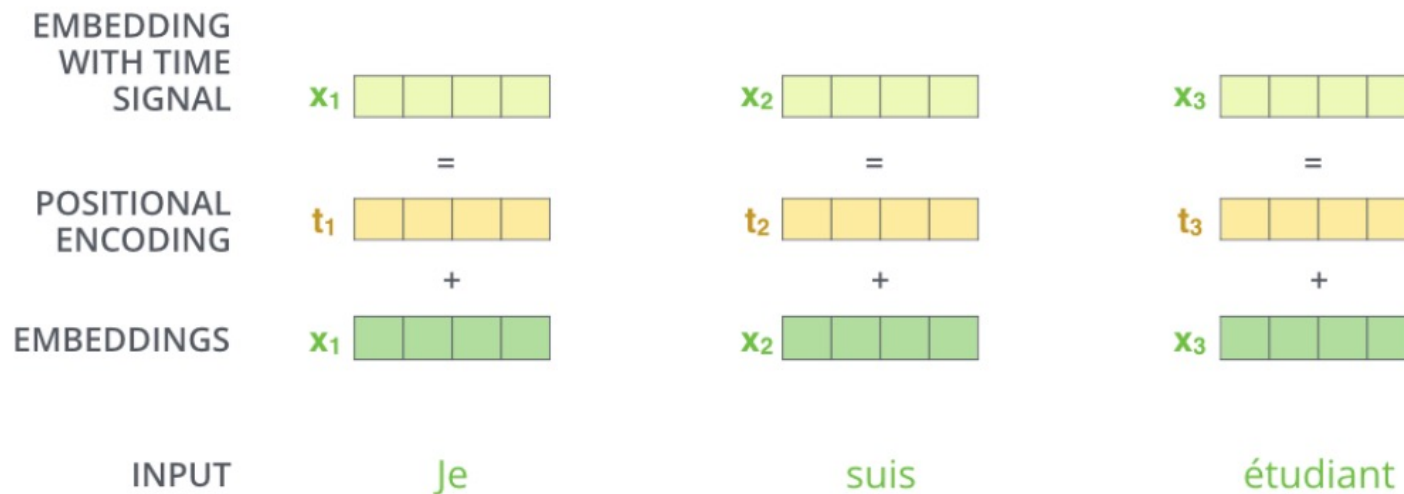


**Problem:**

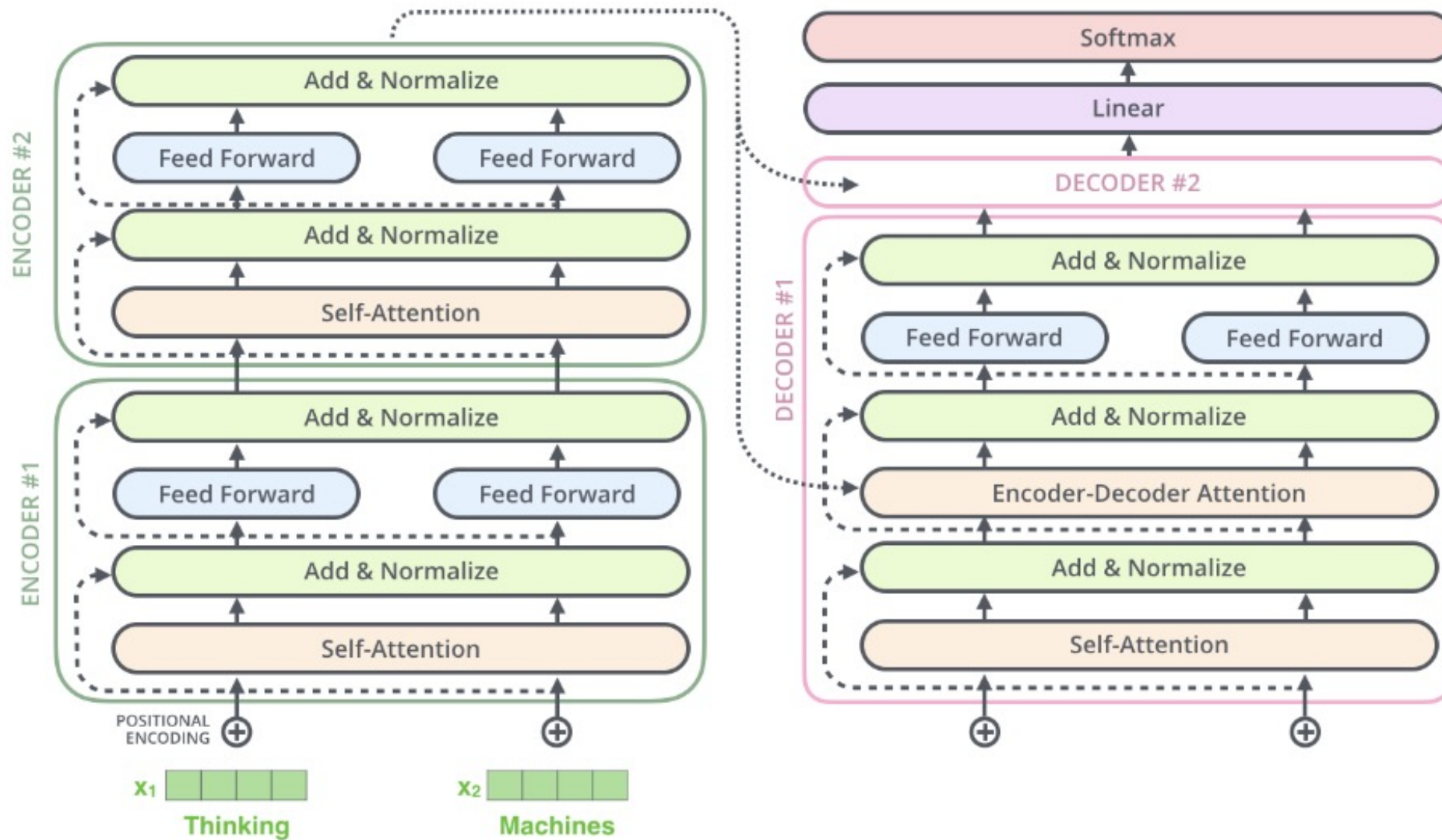
No word order information!

**Solution:**

Add positional embeddings



# Transformer block



# Transformer advantages

- Number of unparallelizable operations does not increase with sequence length.
- Each "word" interacts with each other

## State-of-the-art transformers:

Radford, A., et al. (2018)  
Open AI

Devlin, J., et al. (2018)  
Google AI Language

GPT-2: Generative Pre-training for  
Transformer

BERT: Bidirectional Encoder  
Representations from Transformers

# To do

- Optional reading: **SLP** Ch9
- Review course materials