

# Sequence Model

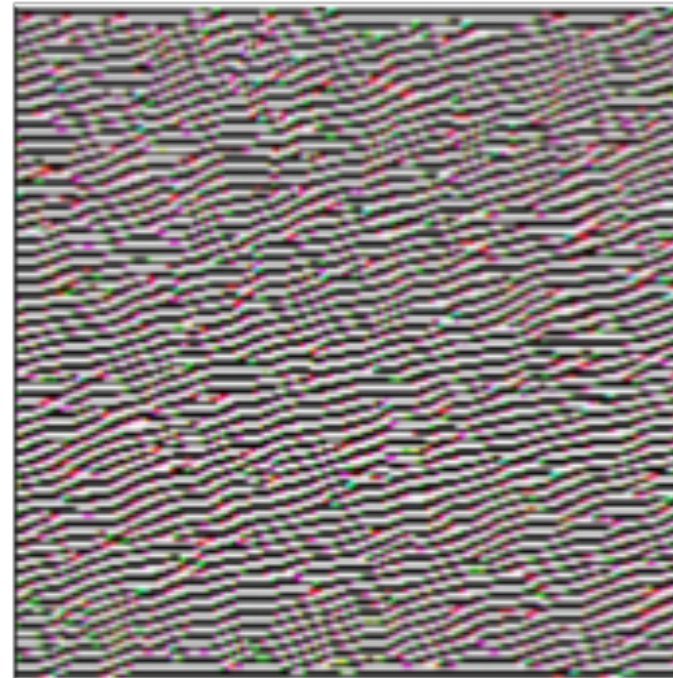
# **Announcements**

1. Makeup exam Dec 11
2. We will release the last reading quiz today

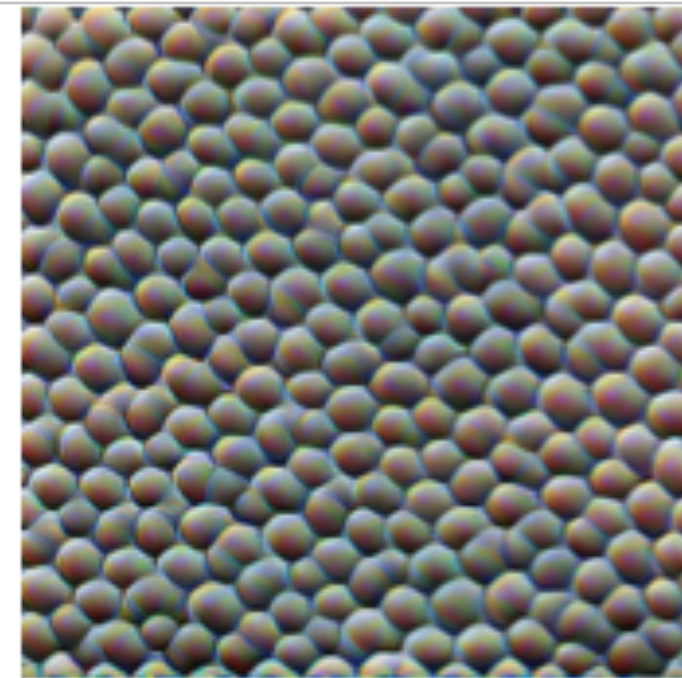
# Recap on Convolutional neural network

Learned feature representations in CNN

Edges



Textures



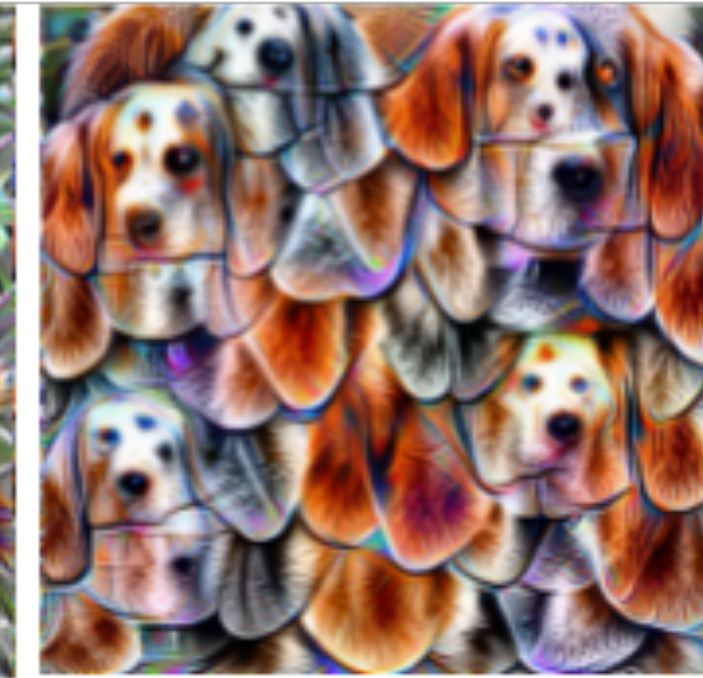
Patterns



Parts



Objects



# Objective today

Understanding neural network structures that are suitable for natural language (i.e., sequences of words)

# Outline today

1. Word-2-Vec embedding and positional embedding

2. Attention model

3. Putting things together: the Transformer model

# Example: autocompletion

e.g., I went to the climbing gym and I \_\_\_\_

**A Language model is a conditional probability model:**

$$y_1 \sim P(Y = \cdot \mid x_1, \dots, x_n) \in \mathbb{R}^{100k}$$

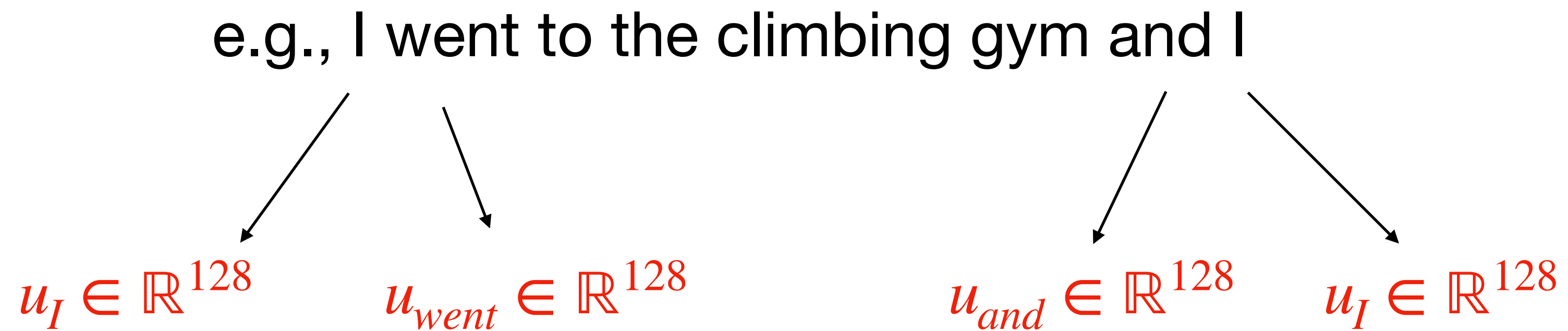
$$y_2 \sim P(Y = \cdot \mid x_1, \dots, x_n, y_1)$$

$$y_m \sim P(Y = \cdot \mid x_1, \dots, x_n, y_1, \dots, y_{m-1})$$

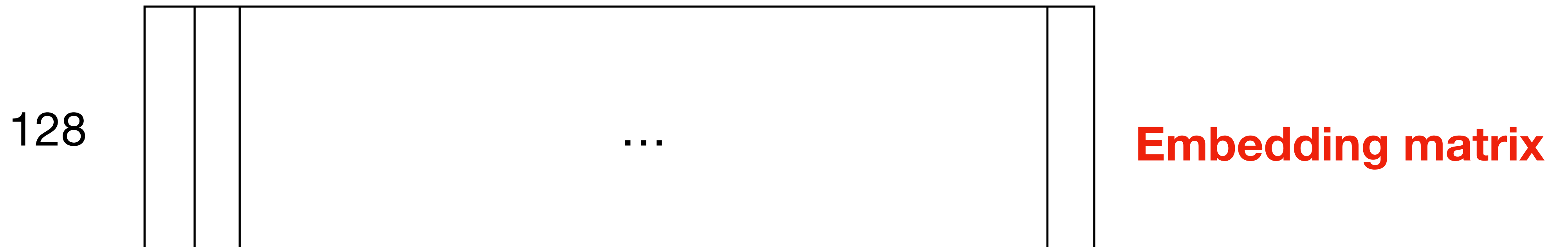


# Word to Vector Embedding

ML models only take vectors of real numbers as inputs...



Size of the English vocabulary (e.g., 100k)



# Positional embedding

Order of the words and their positions matter...

e.g., When I say Transformer in ML, I do not mean the transformer in the movies

$$\begin{aligned} & \swarrow \\ u_{\text{transformer}} & \in \mathbb{R}^{128} \\ & + p_4 \in \mathbb{R}^{128} \end{aligned}$$

$$\begin{aligned} & \swarrow \\ u_{\text{transformer}} + p_{13} & \in \mathbb{R}^{128} \end{aligned}$$

---

Create positional embedding using sin functions

High frequency

$$p_t = \begin{bmatrix} \sin(t/c_1) \\ \sin(t/c_2) \\ \vdots \\ \sin(t/c_{128}) \end{bmatrix}$$

Low frequency



# Summary so far

**We turn words into vectors of real numbers**

e.g., When I say Transformer in ML, I do not mean the transformer in the movies


$$u_{\text{transformer}} + p_4$$

$$u_{\text{transformer}} + p_{13} \in \mathbb{R}^{128}$$

Feature of the word + feature of the position

# Outline today

1. Word-2-Vec embedding and positional embedding

2. Attention model

3. Putting things together: the Transformer model

# Motivation

e.g., When I say Transformer in ML, I do not mean the transformer in the movies

e.g., When I say Transformer, I literally mean the transformer in the movies

**Contextual feature: feature of a word should depend on the context around it**

# Self-attention

I went to the climbing gym

Word-2-vec + positional

$x_1$   $x_2$   $x_3$   $x_4$   $x_5$   $x_6$

$(q_1, k_1, v_1)$   $(q_5, k_5, v_5)$   $(q_6, k_6, v_6)$

$k_1^\top q_5$  ...  $k_i^\top q_5$  ...  $k_6^\top q_5$

Softmax:  $p_{i,5} = \exp(k_i^\top q_5) / \sum_{j=1}^6 \exp(k_j^\top q_5)$

Self-attention layer

$p_{1,5}, p_{2,5}, \dots, p_{6,5}$

$x'_5 = p_{1,5}v_1 + p_{2,5}v_2 + \dots + p_{6,5}v_6$

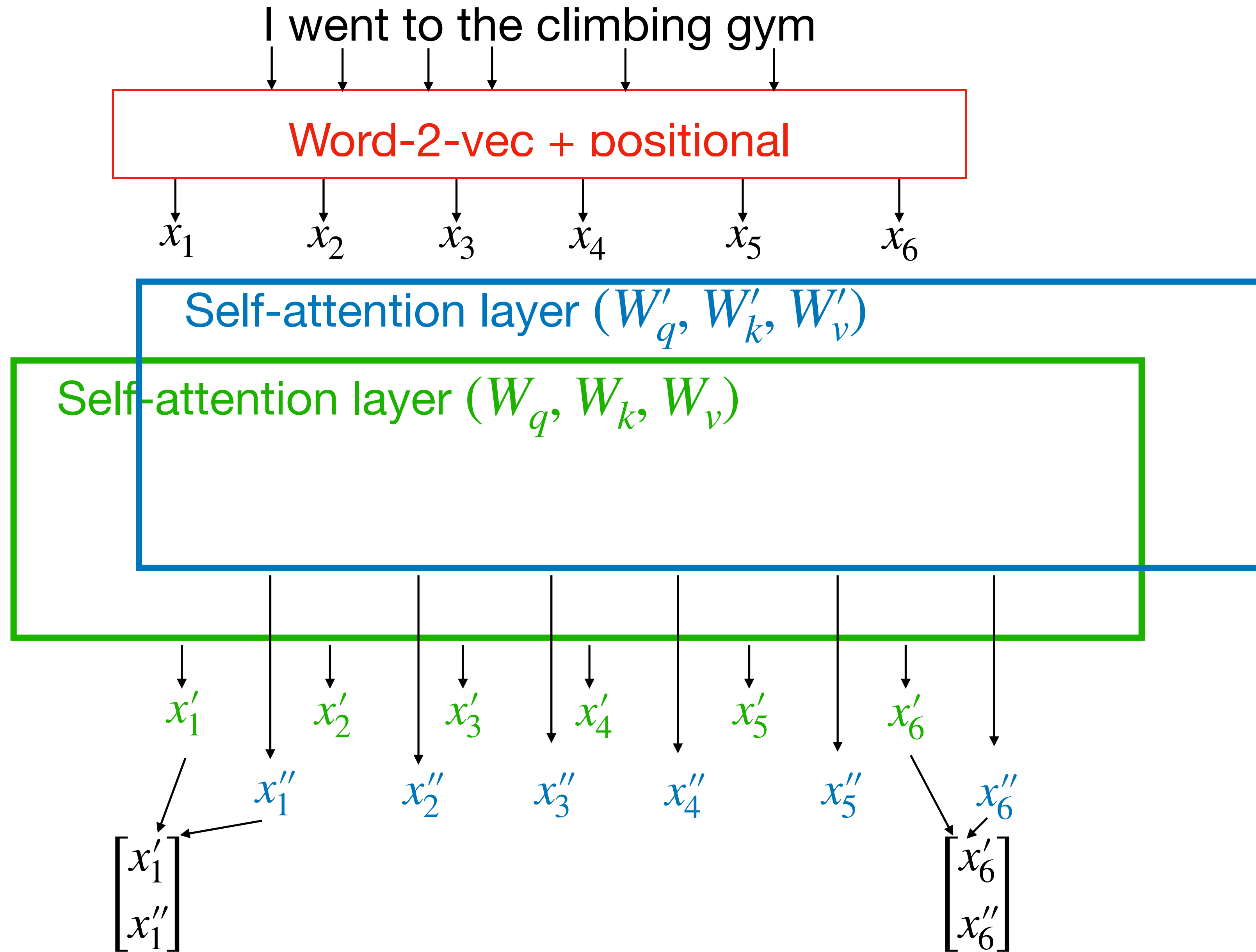
Attention head:  
three matrices:

$W_q, W_k, W_v$

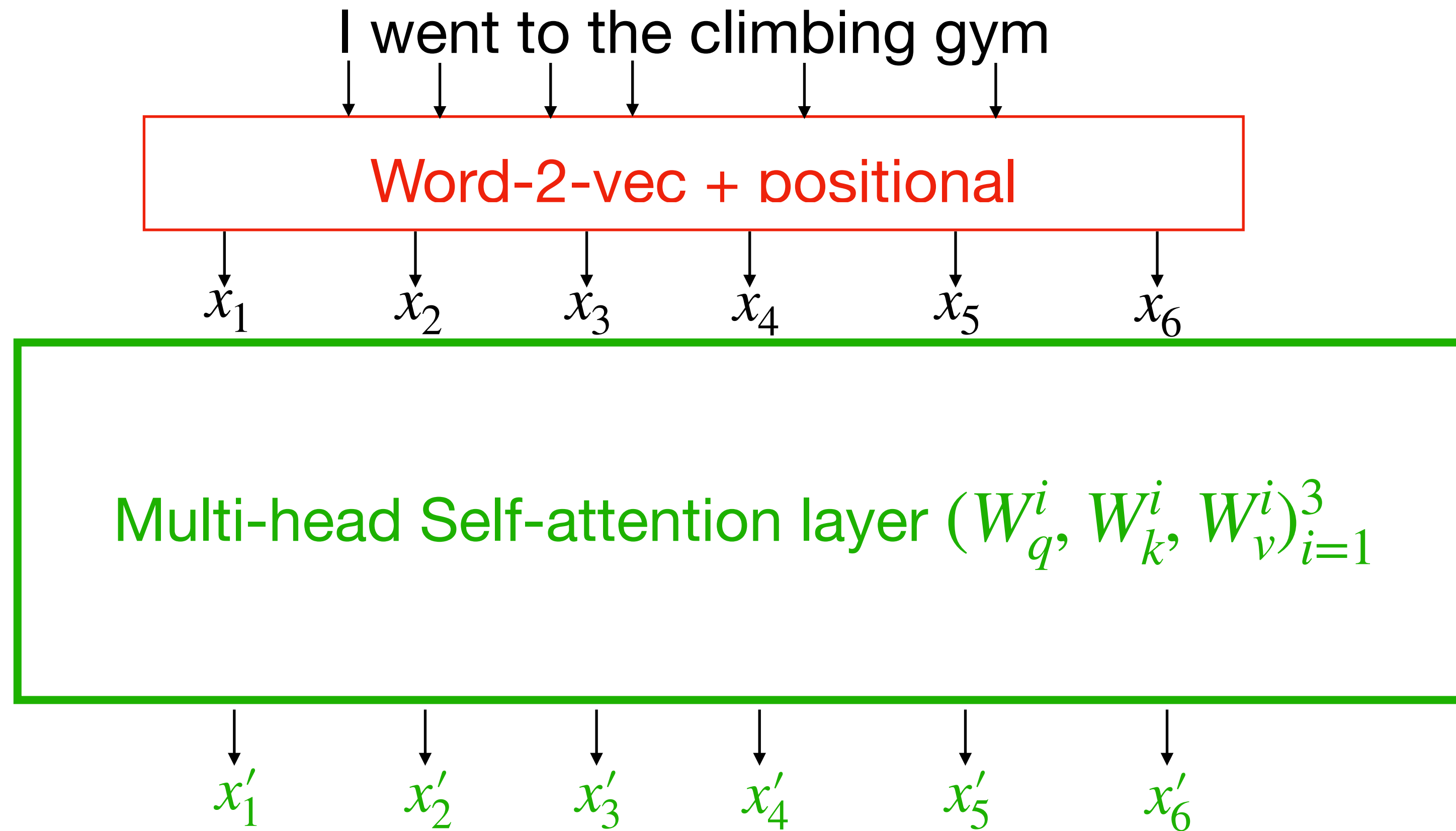
$q = W_q x$   $k = W_k x$   $v = W_v x$

Query key value

# Multi-head self-attention



# Summary so far



Contextual features: e.g.,  $x'_4$  encodes information from all words



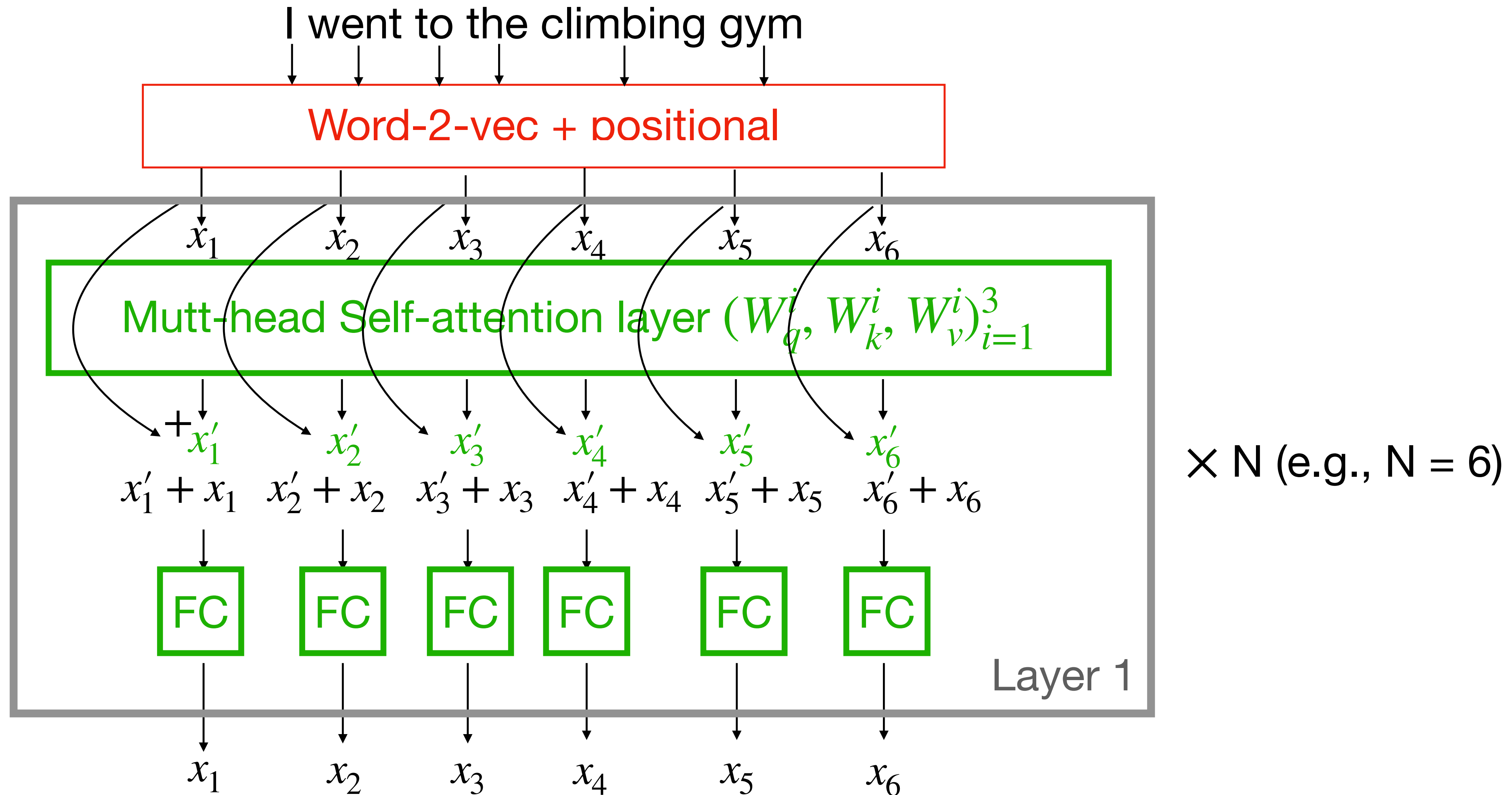
# Outline today

1. Word-2-Vec embedding and positional embedding

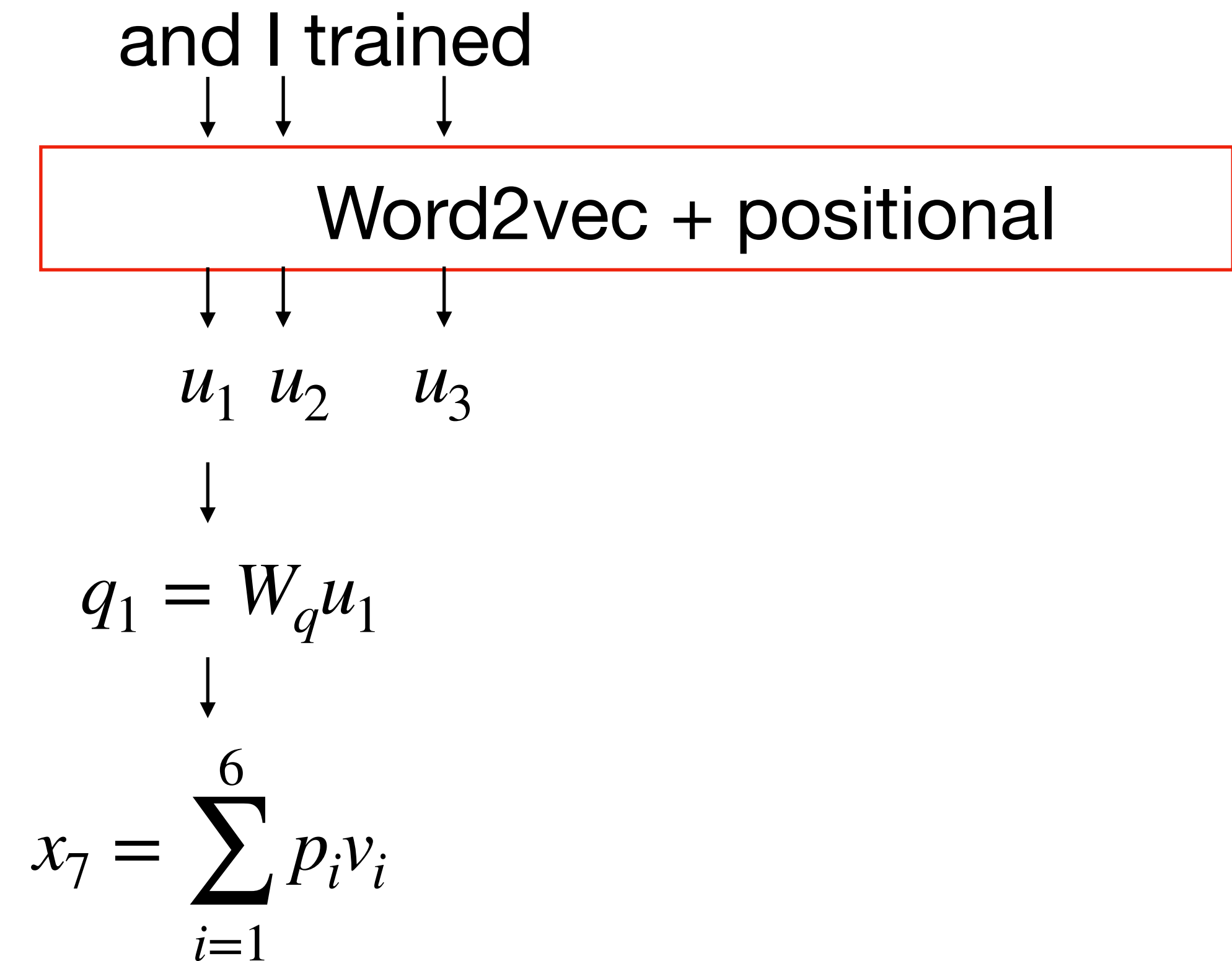
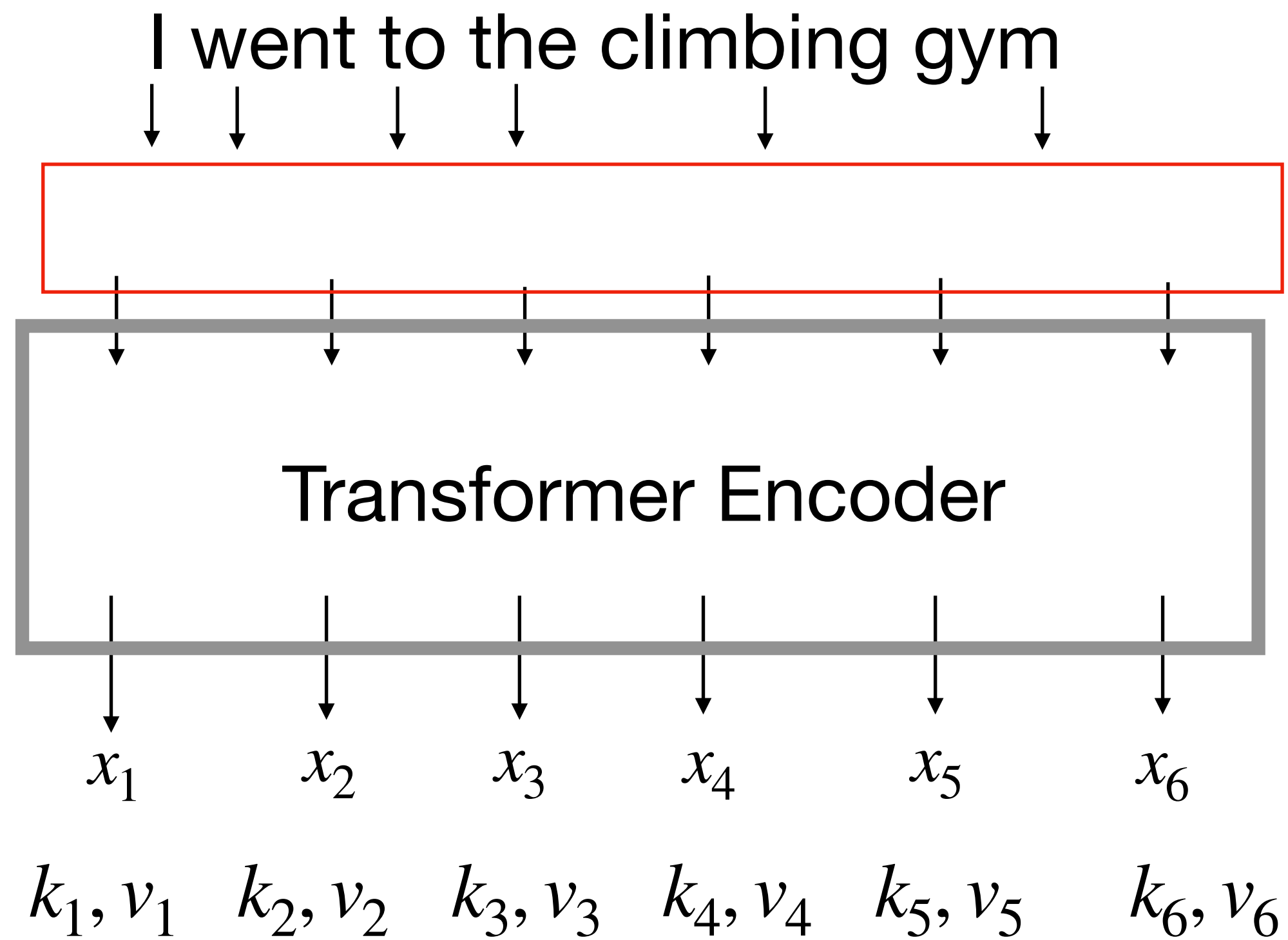
2. Attention model

3. Putting things together: the Transformer model

# The Transformer model: encoder

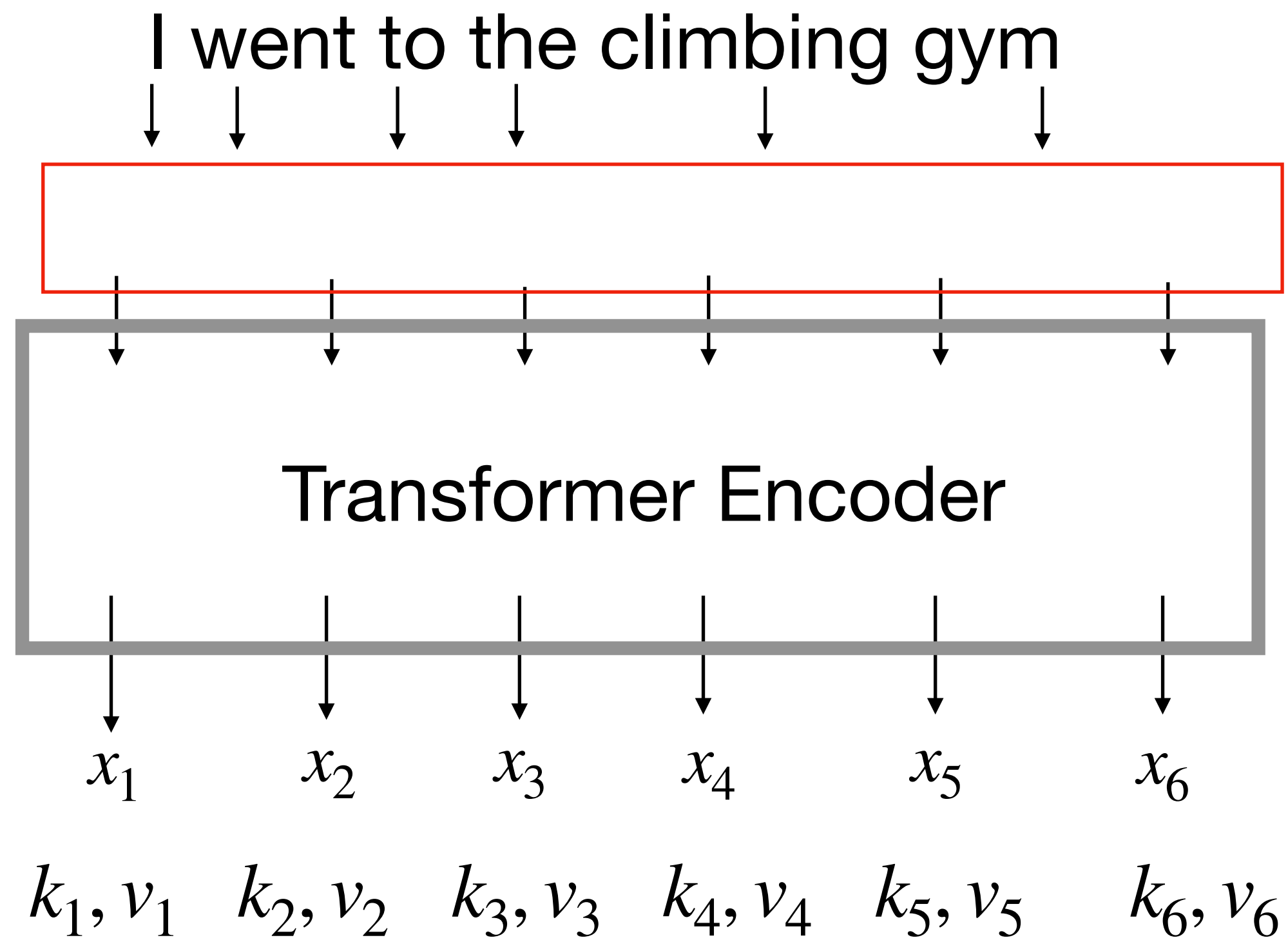


# The Transformer model: decoder

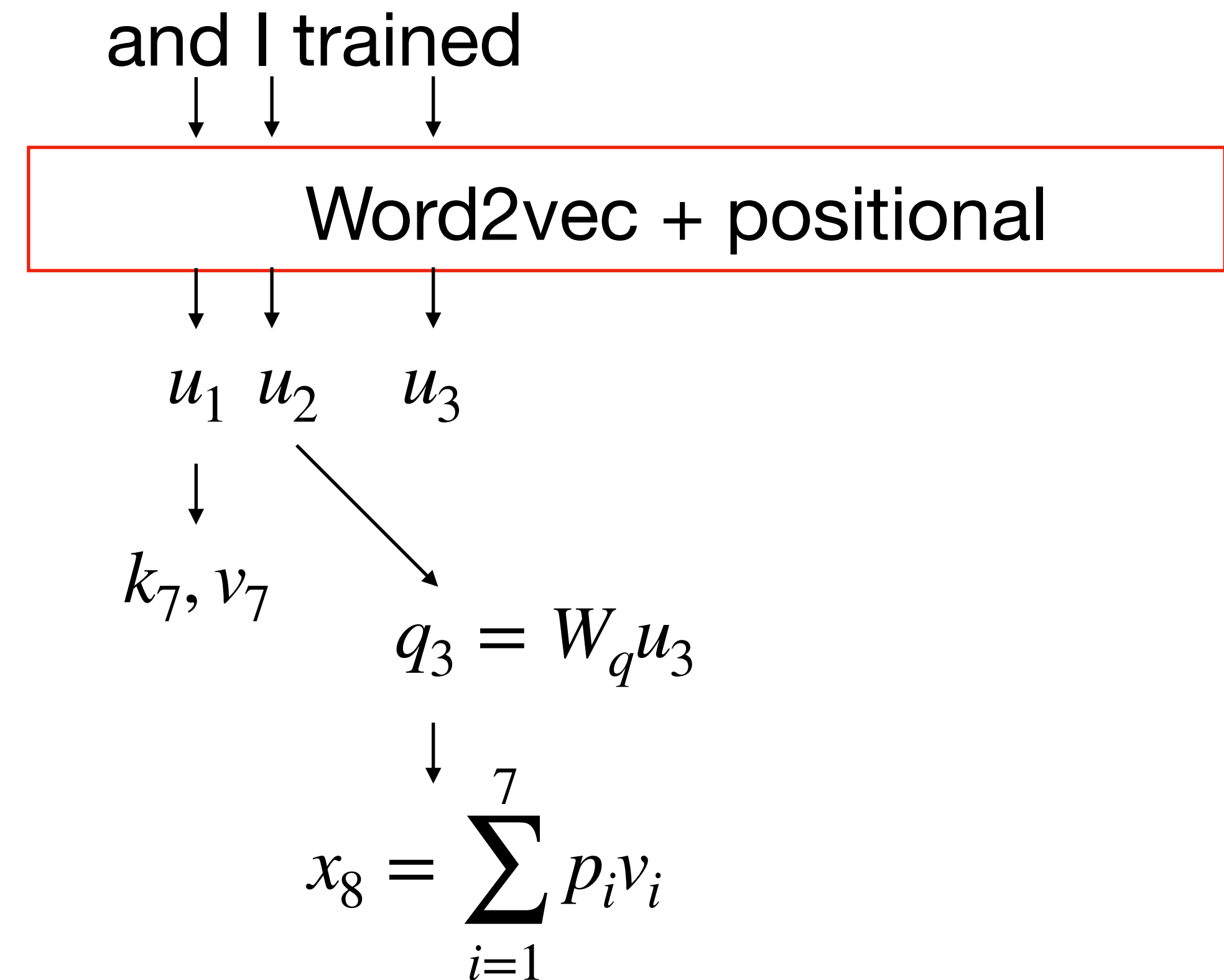


cross-attention ( $W_q, W_k, W_v$ )

# The Transformer model: decoder

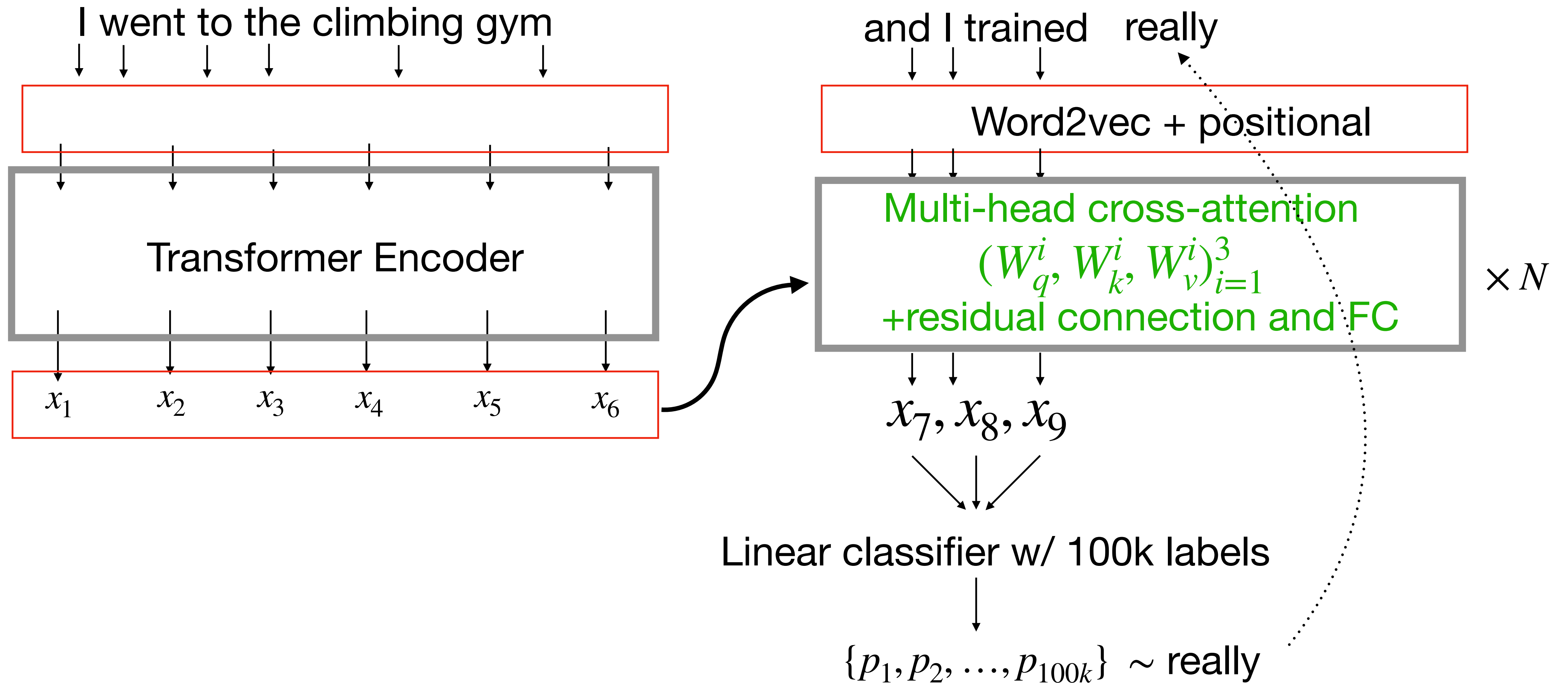


cross-attention ( $W_q, W_k, W_v$ )



Note: we do not pay attention to future words

# The Transformer model: decoder



# Take home task:

Check out the the original paper (not too hard to read!)

---

## 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