# Sequence Models and Attention

CS280

Spring 2025

Angjoo Kanazawa

# Next few lectures

- Today: Transformers / Attention

- Next: Vision Transformers, DINO

- Next Week: Diffusion Models
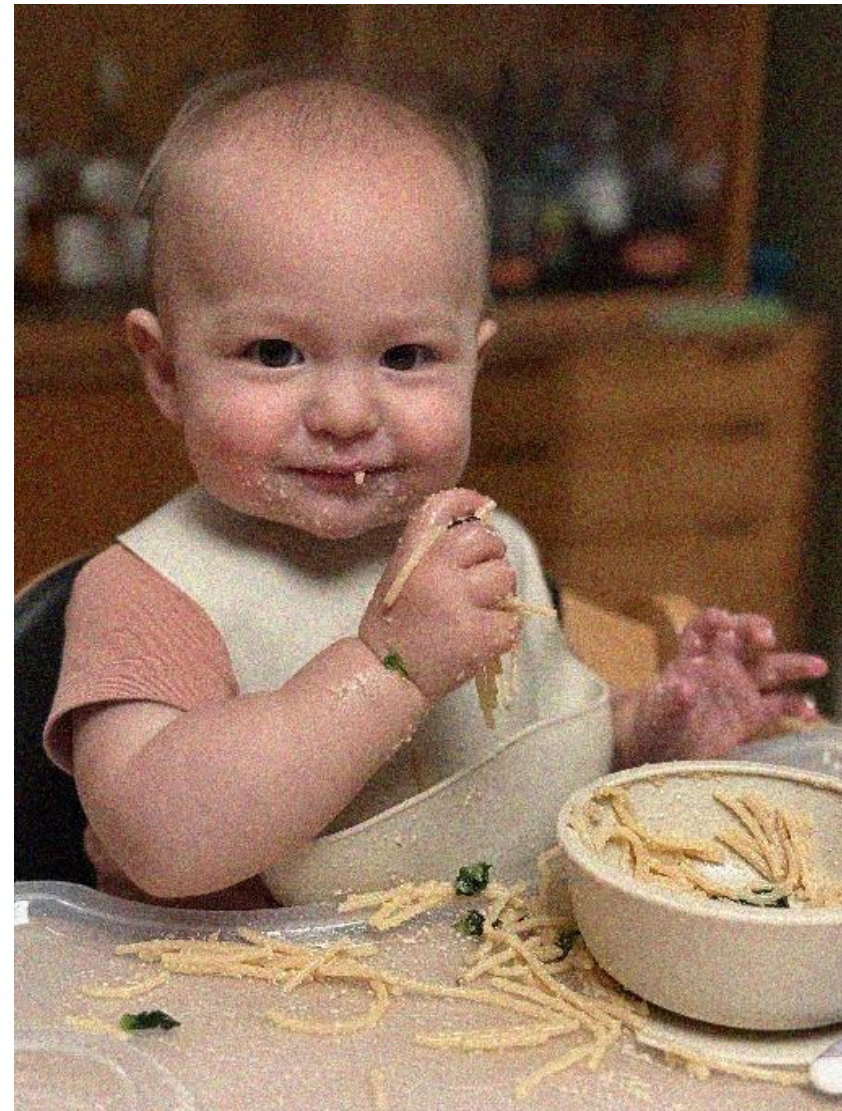
# What is Attention

A Computer Vision perspective

# Computer Vision History on filtering

- Gaussian Filter

- Anisotropic Diffusion

- Bilateral Filter

- Non-Local Means

# How to filter an image?

## Gaussian filter

## How to do edge preserving filter?

# Physical analog: Heat diffusion



Heat Equation

$$\frac{\partial u}{\partial t} = c\Delta u = \mathrm{div}(c\nabla u)$$

Moving from higher to lower concentration of signal u

# Physical analog: Heat diffusion



Heat Equation

$$\frac{\partial u}{\partial t} = c\Delta u = \text{div}(c\nabla u)$$

Solution is exactly convolution with a gaussian kernel!

Moving from higher to lower concentration of signal u

# Useful intuition for later…



Heat Equation $\dfrac{\partial u}{\partial t} = c\Delta u = \mathrm{div}(c\nabla u)$

SDE formulation (how each particle moves): $dx(t) = cdW(t)$

# How to do edge preserving filter?

How do we keep the two soups separate?

Why does gaussian filter destroy the edges?

# Anisotropic Diffusion Perona & Malik 1990

- Idea: Look to see if there is a wall, modulate diffusion across the edge!

- Anisotropic Diffusion:
  - Treats the edges in the image like this wall

  $$\frac{\partial u}{\partial t} = \text{div}(c(|\nabla u|)\nabla u)$$

  - Makes the diffusion process, edge dependent
  - ➔ **Data dependent** filtering
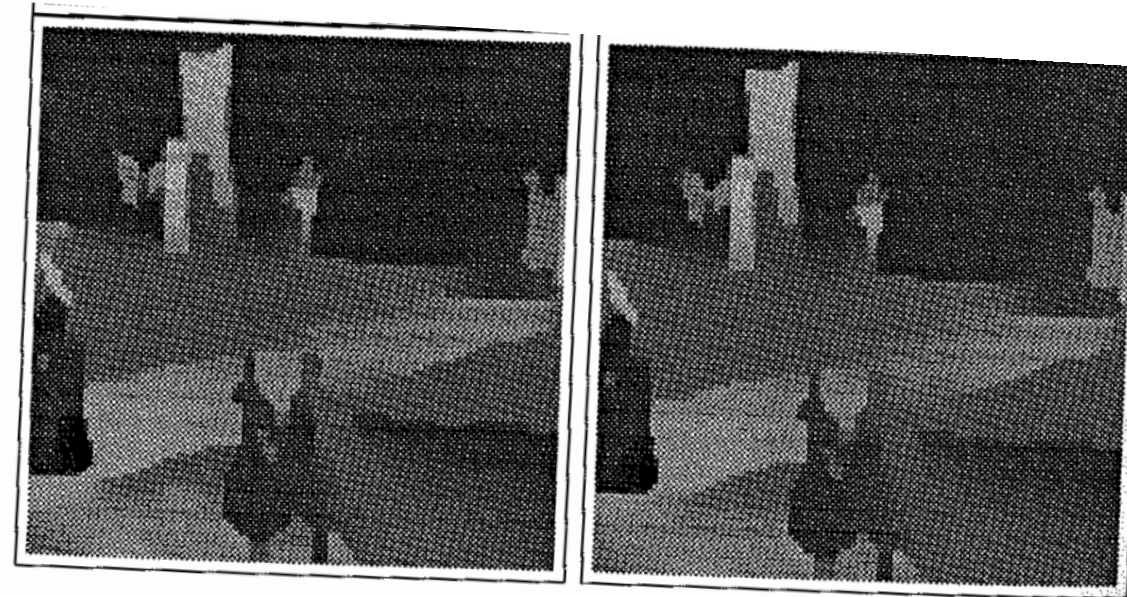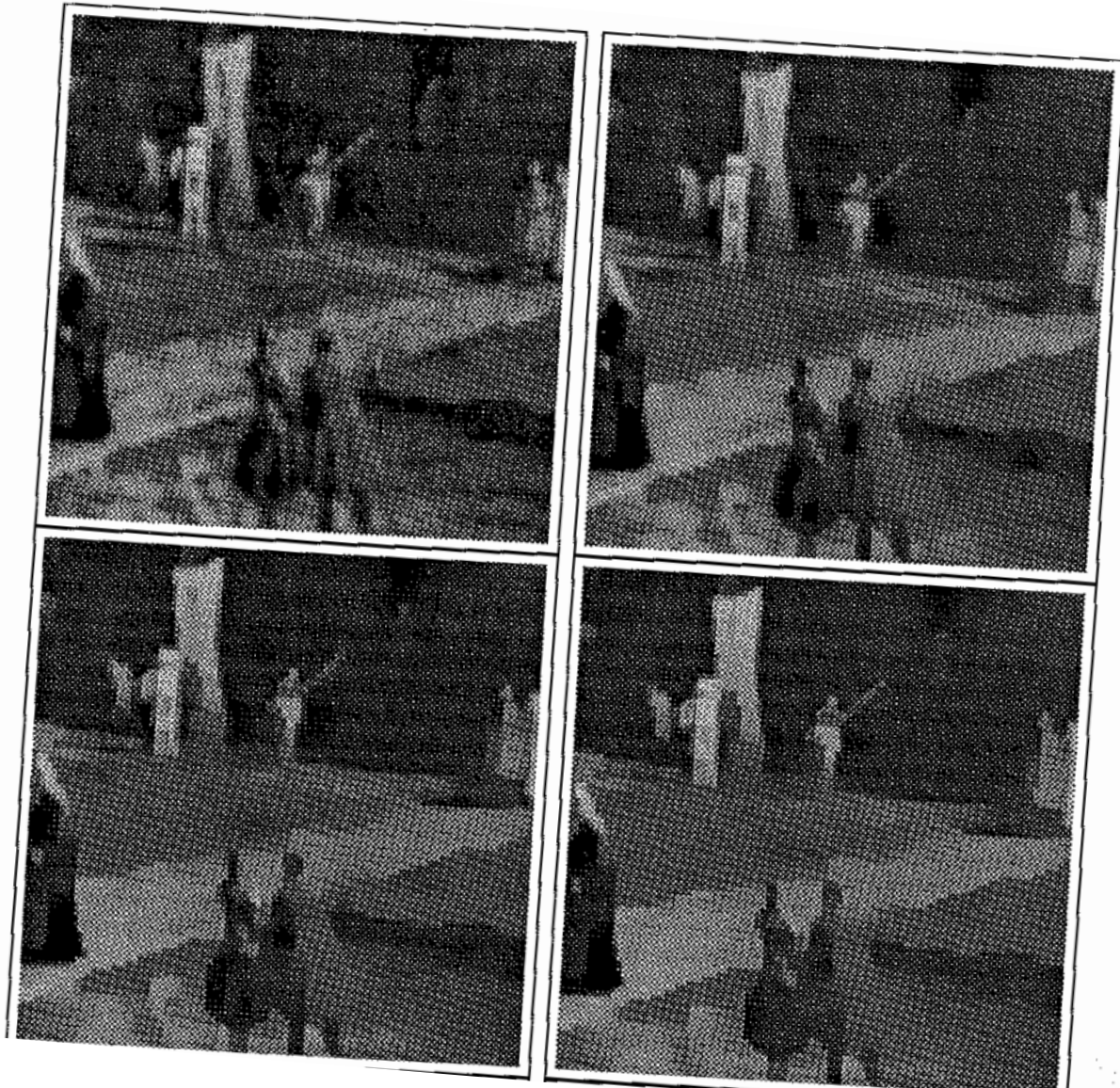
# Anisotropic Diffusion Results



Figure 3.7. Sequence of images produced by anisotropic diffusion. The code presented in figure 3.5 was run on the image at the top-left corner for 10, 20, 30, 60, 100 iterations. The original image has pixel values between 0 (black) and 255 (white) and had a size of $100 \times 100$ pixels. The coefficient K was set equal to $K = 10$.

# Bilateral Filter

Photo from wikipedia

# Bilateral Filter

Aurich and Weule 1995
Tomasi and Manduchi 1998...
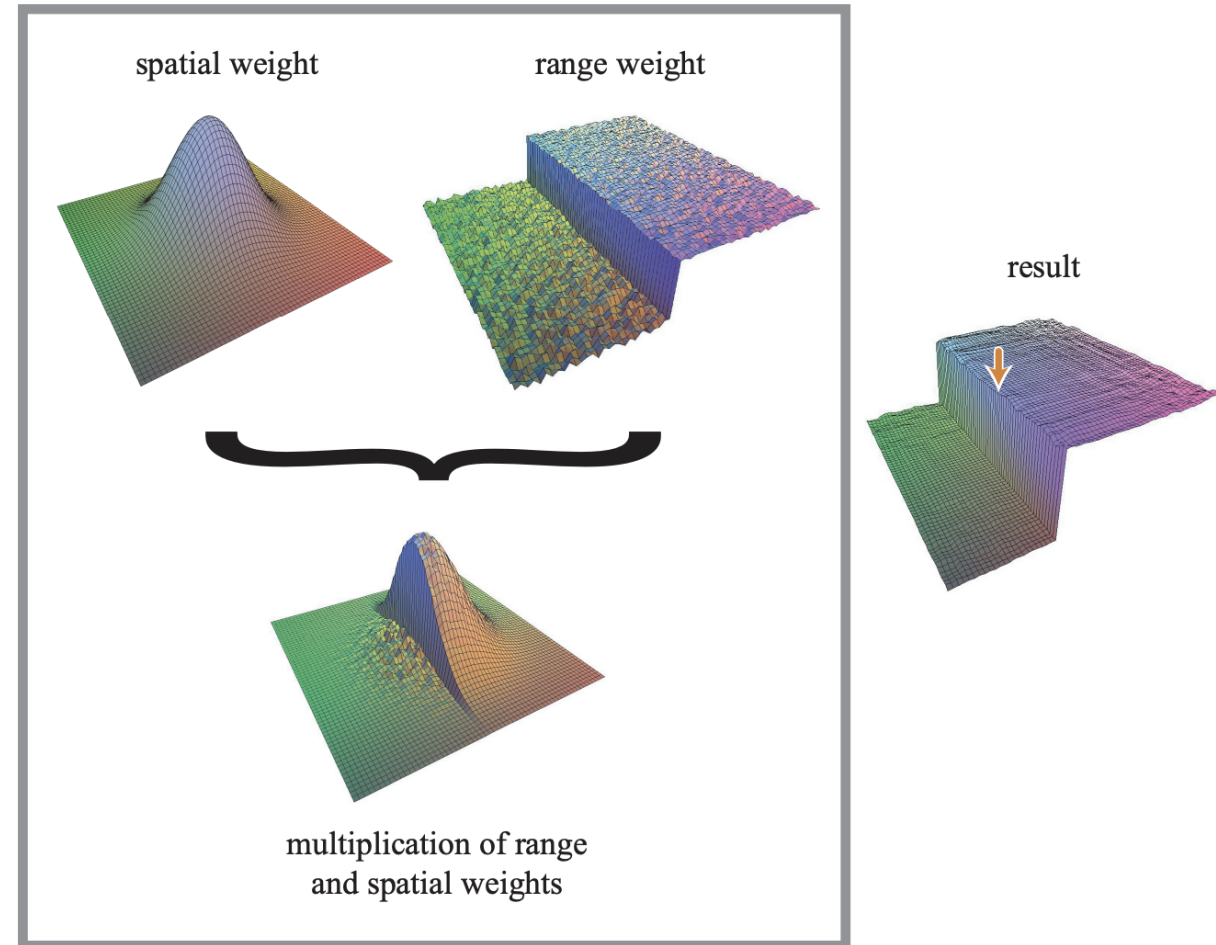
- Inspired by Anisotropic Diffusion

- Weight gaussian kernel with pixel similarity

- Again data-dependent

- This is core idea of ROPE

input

bilateral filter weights of the central pixel

spatial weight

range weight

result

multiplication of range and spatial weights

Distance Similarity

$$BF[I]_{\mathbf{p}} = \frac{1}{W_{\mathbf{p}}} \sum_{\mathbf{q} \in \mathcal{S}} G_{\sigma_s}(\|\mathbf{p} - \mathbf{q}\|) \, G_{\sigma_r}(I_{\mathbf{p}} - I_{\mathbf{q}}) \, I_{\mathbf{q}}$$

Pixel Similarity

Figure from paris et al.

# Bilateral Filter Results



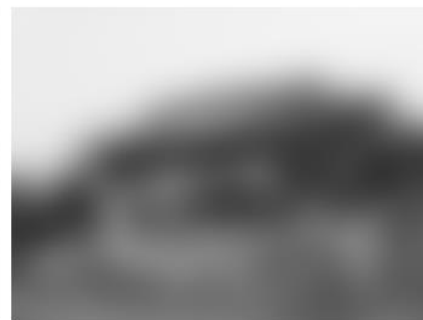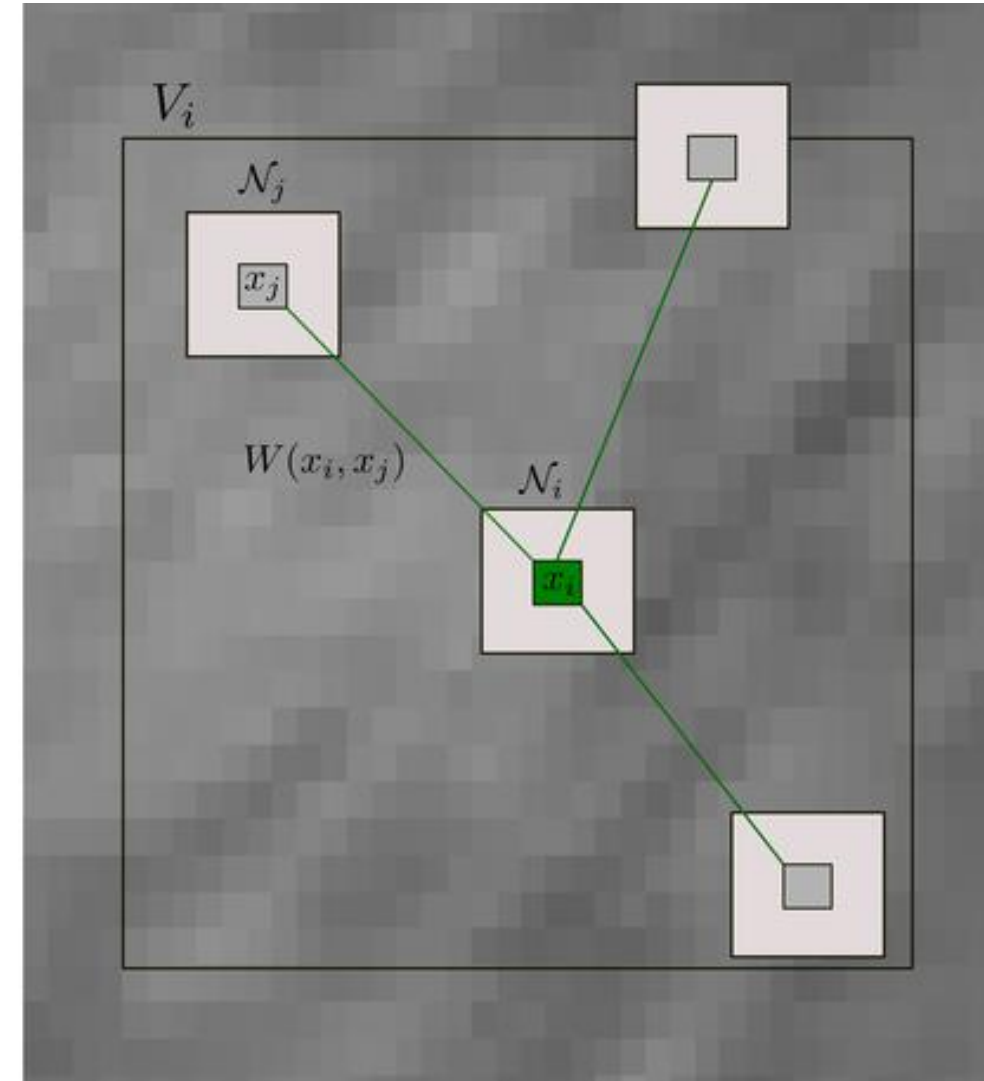|  $\sigma_s \backslash \sigma_r$ | 0.05 | 0.2 | 0.8 | $GB$ |
|---|---|---|---|---|
| 4 | | | | |
| 8 | | | | |
| 16 | | | | |

# Non-local means Buades, Coll, and Morel 2005

- Idea of a filter is to denoise by averaging similar pixels

- Why look at near by pixels only? The similar pixels **can be anywhere**!

- Idea: filter by average of similar patches, **from everywhere in the image**!

- Data-driven weights (similarity) with large field of view.

# Non-local means

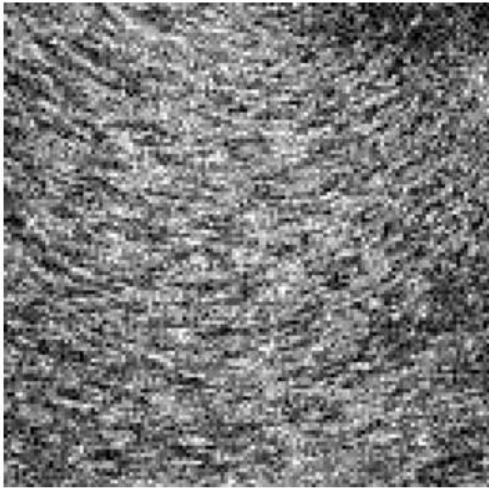$$NL[v](i) = \sum_{j \in I} w(i,j)v(j),$$

Weight is normalized affinity to all other pixels:

$$w(i,j) = \frac{1}{z_i} e^{-\|v_i - v_j\|/\sigma^2} \qquad z_i = \sum_j e^{-\|v_i - v_j\|/\sigma^2}$$
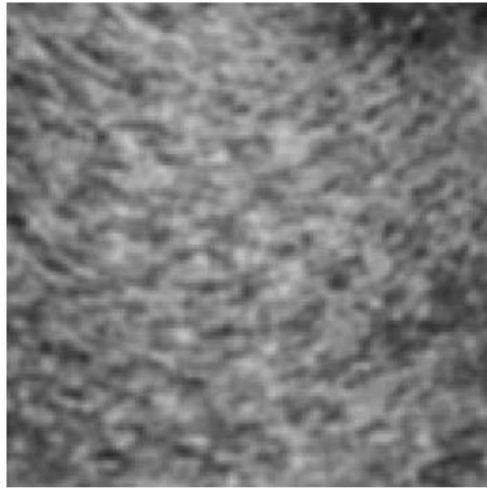
Softmax! In vectorized form, this is:

$$NL = softmax\big(dist(v_i, v_j)\big)V$$

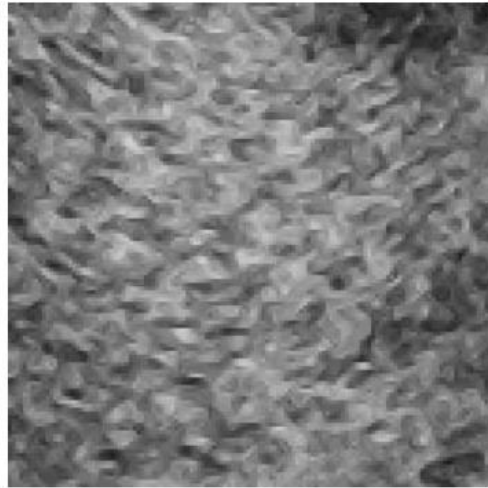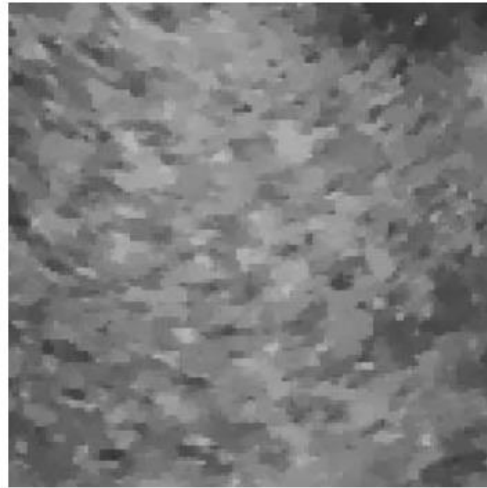# Results



| Input | Gaussian Filter | Anisotropic Filter | Total Variation | **Non-local means** |

# Aside Efros and Leung 1999

- Inspired non-local means



Synthesizing a pixel

non-parametric sampling

Input image

# Two key ideas

In contrast to ordinary convlution by a fixed kernel

1. **Data-driven kernel** → Weights conditioned on the data point
2. **Full field of view** → non-local connections

These are the key difference between convnets and transformers

**Convolution:**

$$x' = Wx + b$$

$$x_i' = \sum_{j}^{N} w_j x_j$$

$w$

Window K

**Attention:**

$$\mathrm{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

$$x_i' = \sum_{j}^{N} w_j(x_i, x_j)x_j$$

# Transformer Architecture

Excellent slides from Jiabin Huang (UMD) for the rest

https://www.youtube.com/watch?v=rcWMRA9E5RI

**Please explain transformers like I am five**



transformer

Sure! Let's break down what a transformer is in a way that's easy to understand.

Imagine you have a big book, and you want to summarize each chapter. A traditional way to do this might be to read each chapter from start to finish and write down the important points. But what if you could read the whole chapter at once and instantly understand which parts are important? That's kind of what a transformer does!

**Key Parts of a Transformer**

1. **Input Representation**: The text is turned into a form that the transformer can understand, usually a list of numbers.
2. **Attention Mechanism**: This is like a spotlight that shines on different parts of the text, helping the transformer focus on the important bits. The cool thing is, it can shine on multiple parts at once, not just one word at a time.
3. **Layers**: The transformer has several layers, and each layer helps it understand the text better. Think of these layers as levels in a video game, where each level makes you a bit stronger and better at the game.

**How It Works**

1. **Reading the Text**: When you give a transformer a sentence, it first converts each word into a number.
2. **Understanding Relationships**: The transformer then uses the attention mechanism to figure out how words in the sentence relate to each other. For example, in the sentence "The cat sat on the mat," it understands that "cat" and "sat" are related.
3. **Processing Layers**: The sentence goes through multiple layers, with each layer making the understanding of the sentence a bit clearer.
4. **Output**: Finally, the transformer gives its output. This could be a translation of the sentence into another language, a summary, or even the next word in a text you're writing.

**Why It's Cool**

Transformers are powerful because they can handle long sentences and understand context really

ChatGPT (OpenAI)

Beautiful, snowy Tokyo city is bustling. The camera moves through the bustling city street, following several people enjoying the beautiful snowy weather and shopping at nearby stalls. Gorgeous sakura petals are flying through the wind along with snowflakes.

transformer

Sora (OpenAI)

Which street was this photo taken on?

transformer

The photo was taken on Bourbon Street, which is located in the French Quarter of New Orleans, Louisiana.

Unified-IO (AI2)

transformer

🤔 **What?** 🤔 **How?** 🤔 **Why?**

**Sequence-to-Sequence model**

En

FR

How are you? → transformer → Comment ça va?
おげんきで
すか？

**Sequence-to-Sequence model**

**Sequence-to-Sequence model**

En

How are you? →

transformer

DE

你好嗎？r? es dir?

**Sequence-to-Sequence model**

En　How are you? → transformer → ZH　你好嗎？

# Sequence-to-Sequence model

ZH

你 好 吃 飽 嗎 … ? <end>

Encoders → Decoders

transformer

En

How are you?                    <start> 你

ZH

你 好 吃 飽 嗎 … ？ <end>

Encoders → Decoders

transformer

En

How are you?

<start> 你 好 嗎 ？

# Tokenization

Many words map to one token, but some don't: indivisible.

8607   4339   2472   311   832   4037   11   719   1063   1541   956   25   3687   23936   13

**One-hot encoding**

cat    dog    bear    cow    indiv

# tokens

$$
\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}
\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}
\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}
\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}
\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}
$$

Value 1 at 3687th entry

# TOKEN EMBEDDING

## One-hot encoding

# TOKEN EMBEDDING



Embedding Space

cat $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$
dog $\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$
bear $\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$
cow $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$
indiv $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$

Value 1 at $3687^{th}$ entry

# TOKEN EMBEDDING



Embedding Space

$$d \begin{bmatrix} 0.5 \\ 2.7 \\ 1.2 \\ \vdots \\ 0.2 \end{bmatrix} = d \begin{bmatrix} W_E \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

Embedded token          Embedding Matrix

# TOKEN EMBEDDING



Embedding Space

$$d \begin{bmatrix} 0.5 \\ 2.7 \\ 1.2 \\ \vdots \\ 0.2 \end{bmatrix} = d \begin{bmatrix} \phantom{x} & \phantom{x} & \phantom{x} & \phantom{x} & \cdots & \phantom{x} \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

dog

Embedded token        Embedding Matrix

# TOKEN EMBEDDING

Apple

dog

Embedding Space

$$
d \begin{bmatrix} 0.5 \\ 2.7 \\ 1.2 \\ \vdots \\ 0.2 \end{bmatrix} = d \begin{bmatrix} & & & & & \cdots & \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}
$$

Embedded token          Embedding Matrix

\# tokens

# TOKEN EMBEDDING

Apple

I bought an **apple** and an orange.

I bought an **apple** watch.

dog

$$d \begin{bmatrix} 0.5 \\ 2.7 \\ 1.2 \\ \vdots \\ 0.2 \end{bmatrix} = d \begin{bmatrix} & & & & & \cdots & & \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

\# tokens

Embedding Space

Embedded token

Embedding Matrix

Encoder

Embedded
Tokens

Token
Embedding

$W_E$ $W_E$ $W_E$ $W_E$ $W_E$ $W_E$ $W_E$

Tokens

**I** **bought** **an** **apple** **and** **an** **orange**

# Self-Attention

**apple**

Embedding Space

Embedded Tokens

Tokens

**I**  **bought**  **an**  **apple**  **and**  **an**  **orange**

# Self-Attention



Embedding Space

Embedded Tokens

Tokens

**I**  **bought**  **an**  **apple**  **and**  **an**  **orange**

# Self-Attention

**apple**

Embedding Space

Embedded Tokens

Tokens

**I**   **bought**   **an**   **apple**   **watch**

# Self-Attention



**apple**

Embedding Space

Embedded Tokens

Tokens

**I**  **bought**  **an**  **apple**  **watch**

# Self-Attention



Embedding Space

Embedded Tokens

Tokens

**I**  **bought**  **an**  **apple**  **watch**

# Self-Attention



Embedding Space

| Embedded Tokens | | | | |
| --- | --- | --- | --- | --- |
| **I** | **bought** | **an** | **apple** | **watch** |

Tokens

# Self-Attention

$$\alpha_{4,1} = \boldsymbol{x}_4^\top \boldsymbol{x}_1 \quad \alpha_{4,2} = \boldsymbol{x}_4^\top \boldsymbol{x}_2 \quad \alpha_{4,3} = \boldsymbol{x}_4^\top \boldsymbol{x}_3 \quad \alpha_{4,4} = \boldsymbol{x}_4^\top \boldsymbol{x}_4 \quad \alpha_{4,5} = \boldsymbol{x}_4^\top \boldsymbol{x}_5$$

Embedded
Tokens

$$\boldsymbol{x}_1 \in R^d \quad \boldsymbol{x}_2 \in R^d \quad \boldsymbol{x}_3 \in R^d \quad \boldsymbol{x}_4 \in R^d \quad \boldsymbol{x}_5 \in R^d$$

Tokens

**I**      **bought**      **an**      **apple**      **watch**

$\alpha'_{4,1}$    0.082      $\alpha'_{4,2}$    0.0495      $\alpha'_{4,3}$    0.0199      $\alpha'_{4,4}$    0.6034      $\alpha'_{4,5}$    0.2452

Softmax

$$\alpha'_{4,i} = \frac{\exp(\alpha_{4,i})}{\sum_j \exp(\alpha_{4,j})}$$

$\alpha_{4,1} = \boldsymbol{x}_4^\top \boldsymbol{x}_1$    $\alpha_{4,2} = \boldsymbol{x}_4^\top \boldsymbol{x}_2$    $\alpha_{4,3} = \boldsymbol{x}_4^\top \boldsymbol{x}_3$    $\alpha_{4,4} = \boldsymbol{x}_4^\top \boldsymbol{x}_4$    $\alpha_{4,5} = \boldsymbol{x}_4^\top \boldsymbol{x}_5$

Embedded Tokens    $\boldsymbol{x}_1$      $\boldsymbol{x}_2$      $\boldsymbol{x}_3$      $\boldsymbol{x}_4$      $\boldsymbol{x}_5$

Tokens    **I**      **bought**      **an**      **apple**      **watch**

$\alpha'_{4,1}$     $\alpha'_{4,2}$     $\alpha'_{4,3}$     $\alpha'_{4,4}$     $\alpha'_{4,5}$

Softmax

$\alpha_{4,1} = \boldsymbol{x}_4^\top \boldsymbol{x}_1$   $\alpha_{4,2} = \boldsymbol{x}_4^\top \boldsymbol{x}_2$   $\alpha_{4,3} = \boldsymbol{x}_4^\top \boldsymbol{x}_3$   $\alpha_{4,4} = \boldsymbol{x}_4^\top \boldsymbol{x}_4$   $\alpha_{4,5} = \boldsymbol{x}_4^\top \boldsymbol{x}_5$

Embedded Tokens

$\boldsymbol{x}_1$     $\boldsymbol{x}_2$     $\boldsymbol{x}_3$     $\boldsymbol{x}_4$     $\boldsymbol{x}_5$

Tokens

**I**     **bought**     **an**     **apple**     **watch**

Updated feature $\quad \boldsymbol{x}_4' \quad = \quad \alpha_{4,1}' \, \boldsymbol{x}_1 \quad + \quad \alpha_{4,2}' \, \boldsymbol{x}_2 \quad + \quad \alpha_{4,3}' \, \boldsymbol{x}_3 \quad + \quad \alpha_{4,4}' \, \boldsymbol{x}_4 \quad + \quad \alpha_{4,5}' \, \boldsymbol{x}_5$

$\alpha_{4,1}' \qquad\qquad \alpha_{4,2}' \qquad\qquad \alpha_{4,3}' \qquad\qquad \alpha_{4,4}' \qquad\qquad \alpha_{4,5}'$

watch

Softmax

I

$\alpha_{4,1} = \boldsymbol{x}_4^\top \boldsymbol{x}_1 \quad\Big|\quad \alpha_{4,2} = \boldsymbol{x}_4^\top \boldsymbol{x}_2 \quad\Big|\quad \alpha_{4,3} = \boldsymbol{x}_4^\top \boldsymbol{x}_3 \quad\Big|\quad \alpha_{4,4} = \boldsymbol{x}_4^\top \boldsymbol{x}_4 \quad\Big|\quad \alpha_{4,5} = \boldsymbol{x}_4^\top \boldsymbol{x}_5$

apple

an          bought

Embedded Tokens          $\boldsymbol{x}_1 \qquad\qquad \boldsymbol{x}_2 \qquad\qquad \boldsymbol{x}_3 \qquad\qquad \boldsymbol{x}_4 \qquad\qquad \boldsymbol{x}_5$

Tokens          I          bought          an          apple          watch

Updated feature $\quad \boldsymbol{x}'_4 \;=\; \alpha'_{4,1}\,\boldsymbol{x}_1 \;+\; \alpha'_{4,2}\,\boldsymbol{x}_2 \;+\; \alpha'_{4,3}\,\boldsymbol{x}_3 \;+\; \alpha'_{4,4}\,\boldsymbol{x}_4 \;+\; \alpha'_{4,5}\,\boldsymbol{x}_5$

**delicious apple**

$\alpha'_{4,1} \qquad\qquad \alpha'_{4,2} \qquad\qquad \alpha'_{4,3} \qquad\qquad \alpha'_{4,4} \qquad\qquad \alpha'_{4,5}$

Softmax

$\alpha_{4,1} = \boldsymbol{x}_4^\top \boldsymbol{x}_1 \quad \alpha_{4,2} = \boldsymbol{x}_4^\top \boldsymbol{x}_2 \quad \alpha_{4,3} = \boldsymbol{x}_4^\top \boldsymbol{x}_3 \quad \alpha_{4,4} = \boldsymbol{x}_4^\top \boldsymbol{x}_4 \quad \alpha_{4,5} = \boldsymbol{x}_4^\top \boldsymbol{x}_5$

Embedded
Tokens $\qquad \boldsymbol{x}_1 \qquad\qquad \boldsymbol{x}_2 \qquad\qquad \boldsymbol{x}_3 \qquad\qquad \boldsymbol{x}_4 \qquad\qquad \boldsymbol{x}_5$

Tokens $\qquad$ **I** $\qquad$ **bought** $\qquad$ **an** $\qquad$ **apple** $\qquad$ **watch**

Updated feature

$$\boldsymbol{x}'_4 = \alpha'_{4,1}\, v_1 + \alpha'_{4,2}\, v_2 + \alpha'_{4,3}\, v_3 + \alpha'_{4,4}\, v_4 + \alpha'_{4,5}\, v_5$$

Updated feature $\boldsymbol{x}'_4 = W^O\left( \alpha'_{4,1} v_1 + \alpha'_{4,2} v_2 + \alpha'_{4,3} v_3 + \alpha'_{4,4} v_4 + \alpha'_{4,5} v_5 \right)$

$$= \sum_i \alpha'_{4,1} \; W^O \; W^V \; \boldsymbol{x}_i$$

$W^O \in R^{d \times d_v}$

$W^Q \in R^{d_k \times d}$

$W^K \in R^{d_k \times d}$

$W^V \in R^{d_v \times d}$

$\alpha'_{4,1}$    $\alpha'_{4,2}$    $\alpha'_{4,3}$    $\alpha'_{4,4}$    $\alpha'_{4,5}$

Softmax

$\alpha_{4,1} = k_1^\top q_4$   $\alpha_{4,2} = k_2^\top q_4$   $\alpha_{4,3} = k_3^\top q_4$   $\alpha_{4,4} = k_4^\top q_4$   $\alpha_{4,5} = k_5^\top q_4$

$k_1$   $v_1$    $k_2$   $v_2$    $k_3$   $v_3$   $q_4$   $k_4$   $v_4$    $k_5$   $v_5$

$W^K$ $W^V$   $W^K$ $W^V$   $W^K$ $W^V$   $W^Q$ $W^K$ $W^V$   $W^K$ $W^V$

Embedded
Tokens

$\boldsymbol{x}_1$    $\boldsymbol{x}_2$    $\boldsymbol{x}_3$    $\boldsymbol{x}_4$    $\boldsymbol{x}_5$

Tokens

**I**    **bought**    **an**    **apple**    **watch**

Updated feature $x'_4 = W^O ( \alpha'_{4,1} v_1 + \alpha'_{4,2} v_2 + \alpha'_{4,3} v_3 + \alpha'_{4,4} v_4 + \alpha'_{4,5} v_5 )$

$$= \sum_i \alpha'_{4,1} ( W^O W^V ) x_i$$

$W^O \in R^{d \times d_v}$

$W^Q \in R^{d_k \times d}$

$W^K \in R^{d_k \times d}$

$W^V \in R^{d_v \times d}$

$\alpha'_{4,1} \qquad \alpha'_{4,2} \qquad \alpha'_{4,3} \qquad \alpha'_{4,4} \qquad \alpha'_{4,5}$

Softmax

$\alpha_{4,1} = k_1^\top q_4 \quad \alpha_{4,2} = k_2^\top q_4 \quad \alpha_{4,3} = k_3^\top q_4 \quad \alpha_{4,4} = k_4^\top q_4 \quad \alpha_{4,5} = k_5^\top q_4$

$k_1 \quad v_1 \qquad k_2 \quad v_2 \qquad k_3 \quad v_3 \quad q_4 \quad k_4 \quad v_4 \qquad k_5 \quad v_5$

$W^K \quad W^V \qquad W^K \quad W^V \qquad W^K \quad W^V \quad W^Q \quad W^K \quad W^V \qquad W^K \quad W^V$

Embedded Tokens

$x_1 \qquad x_2 \qquad x_3 \qquad x_4 \qquad x_5$

Tokens

I        bought        an        apple        watch

Updated feature $\boldsymbol{x}'_2 = W^O\left(\alpha'_{2,1}\ v_1 + \alpha'_{2,2}\ v_2 + \alpha'_{2,3}\ v_3 + \alpha'_{2,4}\ v_4 + \alpha'_{2,5}\ v_5\right)$

$W^O \in R^{d \times d_v}$

$W^Q \in R^{d_k \times d}$

$W^K \in R^{d_k \times d}$

$W^V \in R^{d_v \times d}$

$\alpha'_{2,1}$    $\alpha'_{2,2}$    $\alpha'_{2,3}$    $\alpha'_{2,4}$    $\alpha'_{2,5}$

Softmax

$\alpha_{2,1} = k_1^\top q_2$   $\alpha_{2,2} = k_2^\top q_2$   $\alpha_{2,3} = k_3^\top q_2$   $\alpha_{2,4} = k_4^\top q_2$   $\alpha_{2,5} = k_5^\top q_2$

$k_1$   $v_1$   $q_2$   $k_2$   $v_2$    $k_3$   $v_3$    $k_4$   $v_4$    $k_5$   $v_5$

$W^K$   $W^V$   $W^Q$   $W^K$   $W^V$    $W^K$   $W^V$    $W^K$   $W^V$    $W^K$   $W^V$

Embedded Tokens

$\boldsymbol{x}_1$    $\boldsymbol{x}_2$    $\boldsymbol{x}_3$    $\boldsymbol{x}_4$    $\boldsymbol{x}_5$

Tokens

**I**    **bought**    **an**    **apple**    **watch**

$$\begin{matrix} \alpha_{1,1} & \alpha_{2,1} & \alpha_{3,1} & \alpha_{4,1} & \alpha_{5,1} \\ \alpha_{1,2} & \alpha_{2,2} & \alpha_{3,2} & \alpha_{4,2} & \alpha_{5,2} \\ \alpha_{1,3} & \alpha_{2,3} & \alpha_{3,3} & \alpha_{4,3} & \alpha_{5,3} \\ \alpha_{1,4} & \alpha_{2,4} & \alpha_{3,4} & \alpha_{4,4} & \alpha_{5,4} \\ \alpha_{1,5} & \alpha_{2,5} & \alpha_{3,5} & \alpha_{4,5} & \alpha_{5,5} \end{matrix} = \begin{bmatrix} k_1^\top \\ k_2^\top \\ k_3^\top \\ k_4^\top \\ k_5^\top \end{bmatrix} \begin{matrix} q_1 & q_2 & q_3 & q_4 & q_5 \end{matrix}$$



Embedded Tokens

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

Tokens

**I**    **bought**    **an**    **apple**    **watch**

$$\frac{1}{\sqrt{d_k}}$$

$\otimes$

Softmax

$$\begin{bmatrix} \alpha_{1,1} & \alpha_{2,1} & \alpha_{3,1} & \alpha_{4,1} & \alpha_{5,1} \\ \alpha_{1,2} & \alpha_{2,2} & \alpha_{3,2} & \alpha_{4,2} & \alpha_{5,2} \\ \alpha_{1,3} & \alpha_{2,3} & \alpha_{3,3} & \alpha_{4,3} & \alpha_{5,3} \\ \alpha_{1,4} & \alpha_{2,4} & \alpha_{3,4} & \alpha_{4,4} & \alpha_{5,4} \\ \alpha_{1,5} & \alpha_{2,5} & \alpha_{3,5} & \alpha_{4,5} & \alpha_{5,5} \end{bmatrix} = \begin{bmatrix} k_1^\top \\ k_2^\top \\ k_3^\top \\ k_4^\top \\ k_5^\top \end{bmatrix} \begin{bmatrix} q_1 & q_2 & q_3 & q_4 & q_5 \end{bmatrix}$$

$$\begin{bmatrix} \alpha'_{1,1} & \alpha'_{2,1} & \alpha'_{3,1} & \alpha'_{4,1} & \alpha'_{5,1} \\ \alpha'_{1,2} & \alpha'_{2,2} & \alpha'_{3,2} & \alpha'_{4,2} & \alpha'_{5,2} \\ \alpha'_{1,3} & \alpha'_{2,3} & \alpha'_{3,3} & \alpha'_{4,3} & \alpha'_{5,3} \\ \alpha'_{1,4} & \alpha'_{2,4} & \alpha'_{3,4} & \alpha'_{4,4} & \alpha'_{5,4} \\ \alpha'_{1,5} & \alpha'_{2,5} & \alpha'_{3,5} & \alpha'_{4,5} & \alpha'_{5,5} \end{bmatrix}$$

$$k = [k^1, k^2, \cdots, k^{d_k}]^\top \qquad E[k^i] = E[q^i] = 0$$
$$q = [q^1, q^2, \cdots, q^{d_k}]^\top \qquad \mathrm{Var}[k^i] = \mathrm{Var}[q^i] = 1$$

$$k^\top q = \sum_{i=1}^{d_k} k^i q^i \qquad \mathrm{Var}[k^\top q] = d_k$$

$q_1$ $k_1$ $v_1$  $q_2$ $k_2$ $v_2$  $q_3$ $k_3$ $v_3$  $q_4$ $k_4$ $v_4$  $q_5$ $k_5$ $v_5$

$W^Q$ $W^K$ $W^V$   $W^Q$ $W^K$ $W^V$   $W^Q$ $W^K$ $W^V$   $W^Q$ $W^K$ $W^V$   $W^Q$ $W^K$ $W^V$

Embedded Tokens

$\boldsymbol{x}_1$ $\boldsymbol{x}_2$ $\boldsymbol{x}_3$ $\boldsymbol{x}_4$ $\boldsymbol{x}_5$

Tokens

**I** **bought** **an** **apple** **watch**

$$\frac{1}{\sqrt{d_k}}$$

$$\otimes$$

Softmax

$$\begin{bmatrix} \alpha_{1,1} & \alpha_{2,1} & \alpha_{3,1} & \alpha_{4,1} & \alpha_{5,1} \\ \alpha_{1,2} & \alpha_{2,2} & \alpha_{3,2} & \alpha_{4,2} & \alpha_{5,2} \\ \alpha_{1,3} & \alpha_{2,3} & A_{3,3} & \alpha_{4,3} & \alpha_{5,3} \\ \alpha_{1,4} & \alpha_{2,4} & \alpha_{3,4} & \alpha_{4,4} & \alpha_{5,4} \\ \alpha_{1,5} & \alpha_{2,5} & \alpha_{3,5} & \alpha_{4,5} & \alpha_{5,5} \end{bmatrix} = \begin{bmatrix} k_1^\top \\ k_2^\top \\ K_3^\top \\ k_4^\top \\ k_5^\top \end{bmatrix} \begin{bmatrix} q_1 & q_2 & Q & q_4 & q_5 \end{bmatrix}$$

$$\begin{bmatrix} \alpha'_{1,1} & \alpha'_{2,1} & \alpha'_{3,1} & \alpha'_{4,1} & \alpha'_{5,1} \\ \alpha'_{1,2} & \alpha'_{2,2} & \alpha'_{3,2} & \alpha'_{4,2} & \alpha'_{5,2} \\ \alpha'_{1,3} & \alpha'_{2,3} & A'_{3,3} & \alpha'_{4,3} & \alpha'_{5,3} \\ \alpha'_{1,4} & \alpha'_{2,4} & \alpha'_{3,4} & \alpha'_{4,4} & \alpha'_{5,4} \\ \alpha'_{1,5} & \alpha'_{2,5} & \alpha'_{3,5} & \alpha'_{4,5} & \alpha'_{5,5} \end{bmatrix}$$

| $q_1$ | $k_1$ | $v_1$ | $q_2$ | $k_2$ | $v_2$ | $q_3$ | $k_3$ | $v_3$ | $q_4$ | $k_4$ | $v_4$ | $q_5$ | $k_5$ | $v_5$ |

| $W^Q$ | $W^K$ | $W^V$ | $W^Q$ | $W^K$ | $W^V$ | $W^Q$ | $W^K$ | $W^V$ | $W^Q$ | $W^K$ | $W^V$ | $W^Q$ | $W^K$ | $W^V$ |

Embedded Tokens

$$\boldsymbol{x}_1 \qquad \boldsymbol{x}_2 \qquad \boldsymbol{x}_3 \qquad \boldsymbol{x}_4 \qquad \boldsymbol{x}_5$$

Tokens

**I**     **bought**     **an**     **apple**     **watch**

$$\frac{1}{\sqrt{d_k}}$$

$$\otimes$$

Softmax

$$\begin{bmatrix} \alpha_{1,1} & \alpha_{2,1} & \alpha_{3,1} & \alpha_{4,1} & \alpha_{5,1} \\ \alpha_{1,2} & \alpha_{2,2} & \alpha_{3,2} & \alpha_{4,2} & \alpha_{5,2} \\ \alpha_{1,3} & \alpha_{2,3} & A_{3,3} & \alpha_{4,3} & \alpha_{5,3} \\ \alpha_{1,4} & \alpha_{2,4} & \alpha_{3,4} & \alpha_{4,4} & \alpha_{5,4} \\ \alpha_{1,5} & \alpha_{2,5} & \alpha_{3,5} & \alpha_{4,5} & \alpha_{5,5} \end{bmatrix} = \begin{bmatrix} k_1^\top \\ k_2^\top \\ K^\top \\ k_4^\top \\ k_5^\top \end{bmatrix} \begin{bmatrix} q_1 & q_2 & Q & q_4 & q_5 \end{bmatrix}$$

$$Q = W^Q \ x_1 \ x_2 \ x_3 \ x_4 \ x_5$$

$$K = W^K \ x_1 \ x_2 \ x_3 \ x_4 \ x_5$$
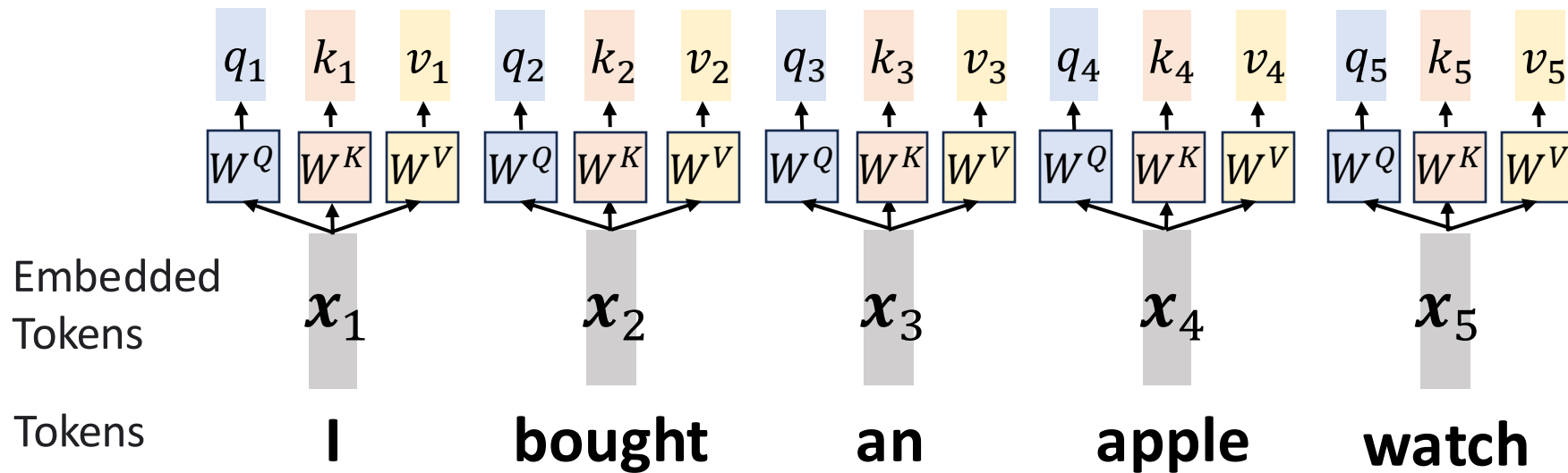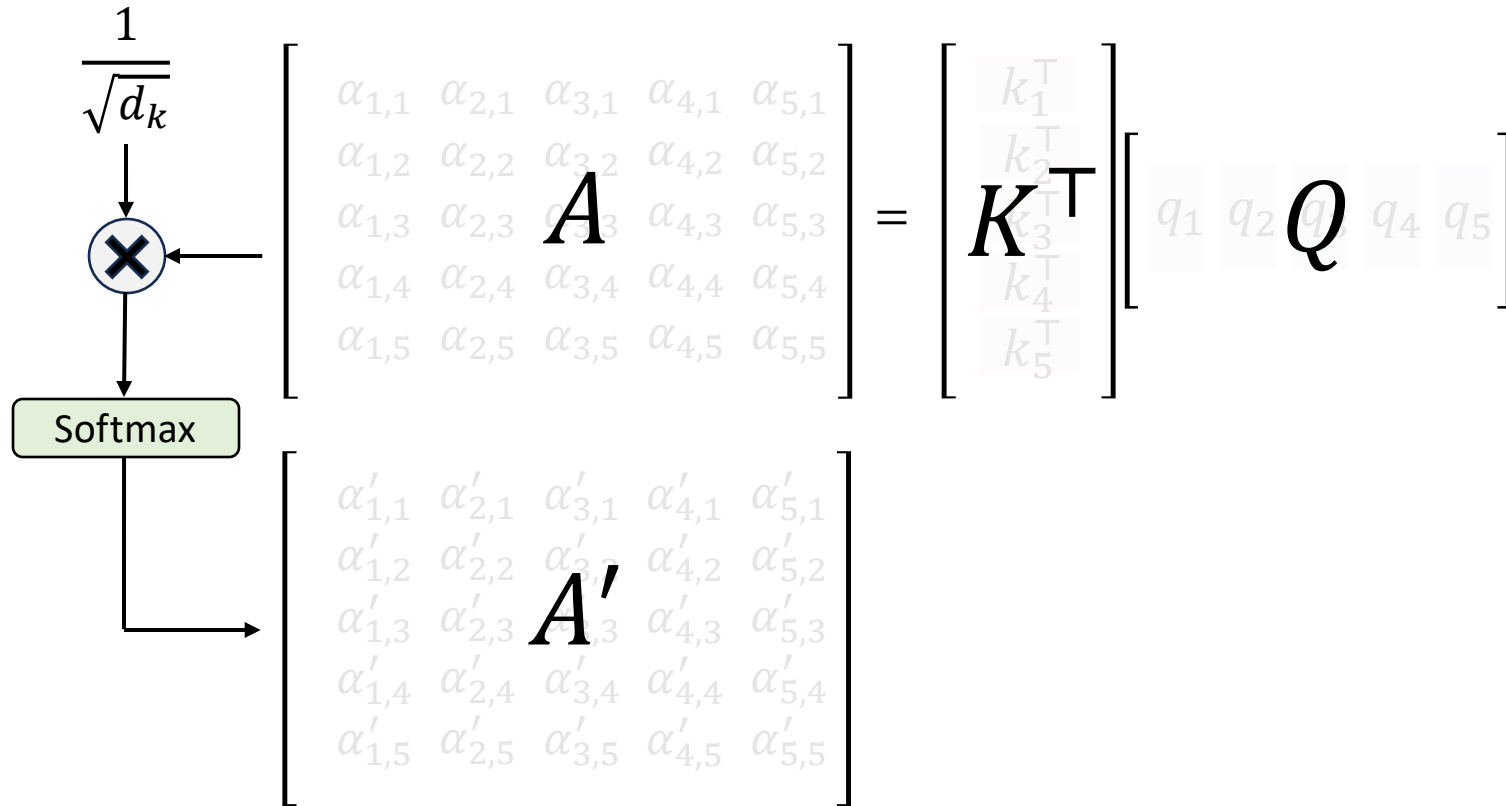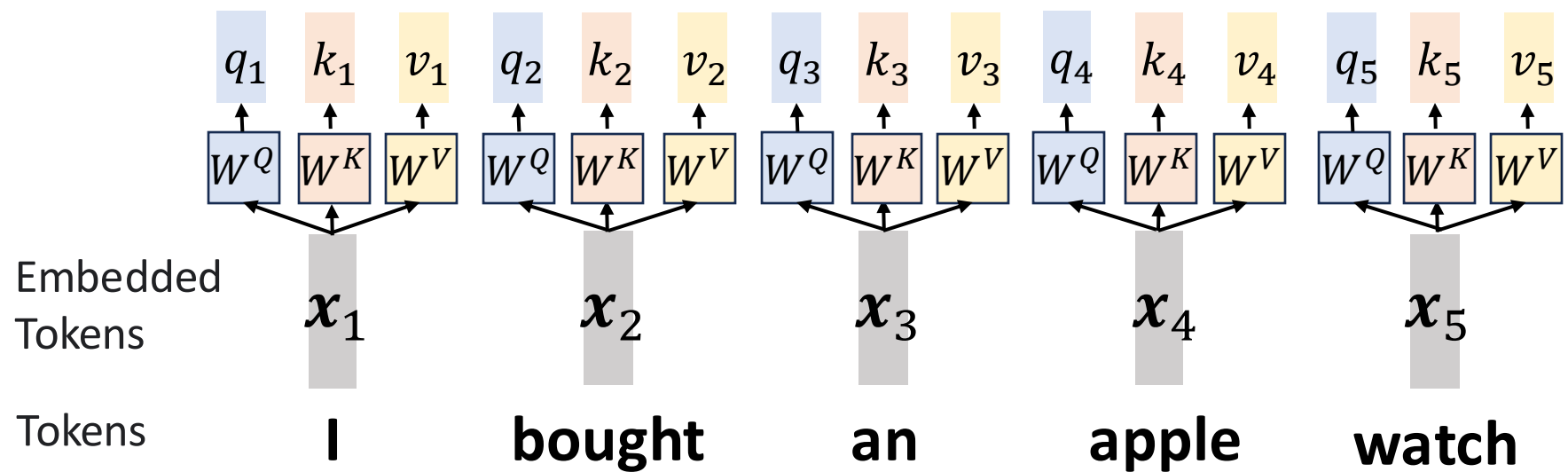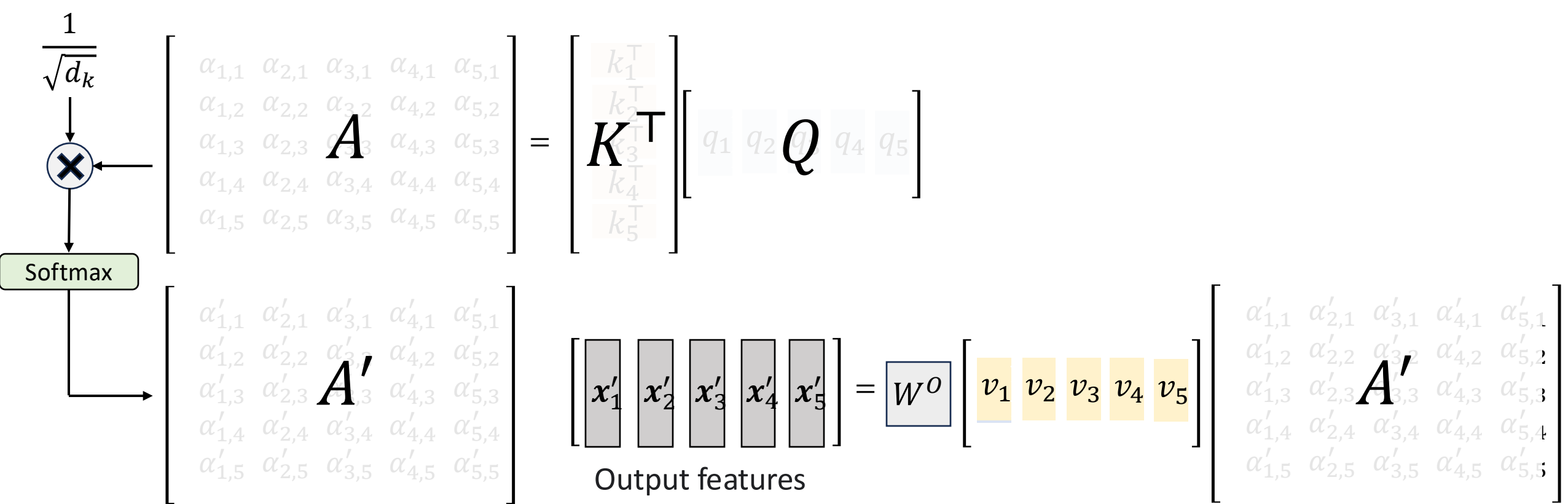
$$V = W^V \ x_1 \ x_2 \ x_3 \ x_4 \ x_5$$

$$\begin{bmatrix} \alpha'_{1,1} & \alpha'_{2,1} & \alpha'_{3,1} & \alpha'_{4,1} & \alpha'_{5,1} \\ \alpha'_{1,2} & \alpha'_{2,2} & \alpha'_{3,2} & \alpha'_{4,2} & \alpha'_{5,2} \\ \alpha'_{1,3} & \alpha'_{2,3} & A' & \alpha'_{4,3} & \alpha'_{5,3} \\ \alpha'_{1,4} & \alpha'_{2,4} & \alpha'_{3,4} & \alpha'_{4,4} & \alpha'_{5,4} \\ \alpha'_{1,5} & \alpha'_{2,5} & \alpha'_{3,5} & \alpha'_{4,5} & \alpha'_{5,5} \end{bmatrix}$$

$$\begin{bmatrix} x'_1 & x'_2 & x'_3 & x'_4 & x'_5 \end{bmatrix} = W^O \begin{bmatrix} v_1 & v_2 & V & v_4 & v_5 \end{bmatrix} \begin{bmatrix} \alpha'_{1,1} & \alpha'_{2,1} & \alpha'_{3,1} & \alpha'_{4,1} & \alpha'_{5,1} \\ \alpha'_{1,2} & \alpha'_{2,2} & \alpha'_{3,2} & \alpha'_{4,2} & \alpha'_{5,2} \\ \alpha'_{1,3} & \alpha'_{2,3} & A' & \alpha'_{4,3} & \alpha'_{5,3} \\ \alpha'_{1,4} & \alpha'_{2,4} & \alpha'_{3,4} & \alpha'_{4,4} & \alpha'_{5,4} \\ \alpha'_{1,5} & \alpha'_{2,5} & \alpha'_{3,5} & \alpha'_{4,5} & \alpha'_{5,5} \end{bmatrix}$$

Output features

$q_1 \ k_1 \ v_1$  $q_2 \ k_2 \ v_2$  $q_3 \ k_3 \ v_3$  $q_4 \ k_4 \ v_4$  $q_5 \ k_5 \ v_5$

$W^Q \ W^K \ W^V$  $W^Q \ W^K \ W^V$  $W^Q \ W^K \ W^V$  $W^Q \ W^K \ W^V$  $W^Q \ W^K \ W^V$

Embedded Tokens

$x_1$  $x_2$  $x_3$  $x_4$  $x_5$
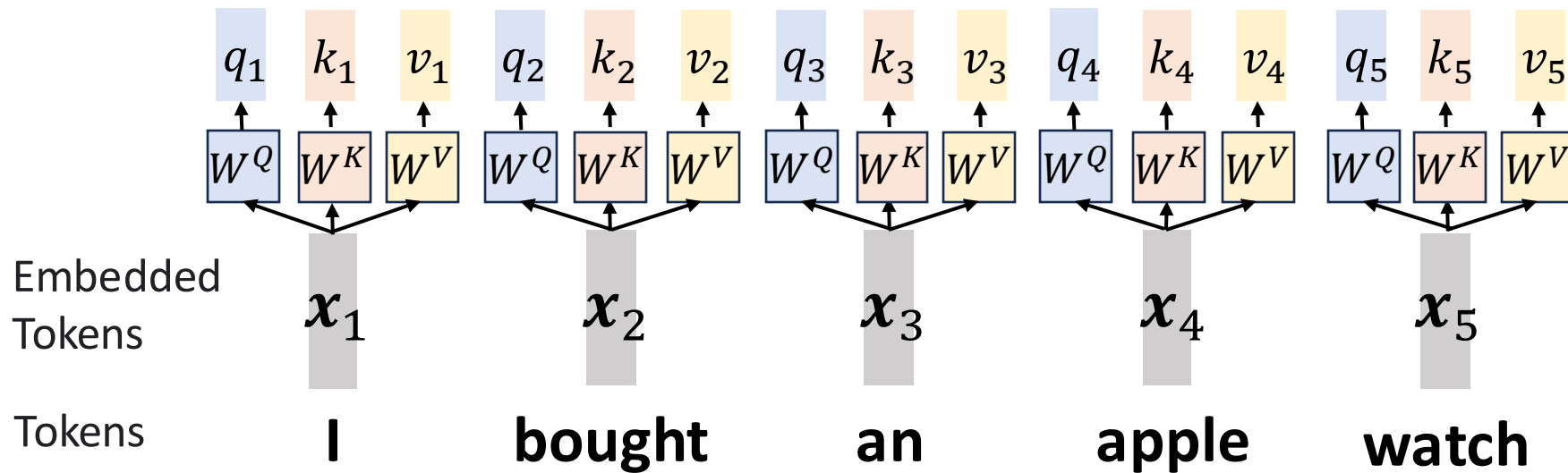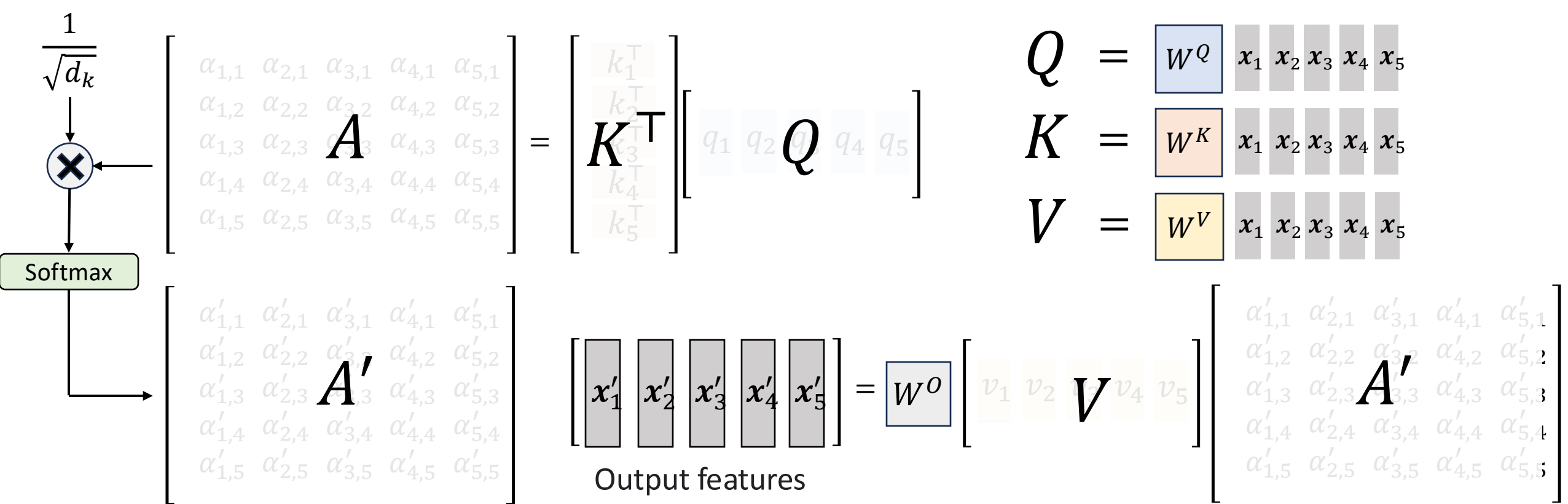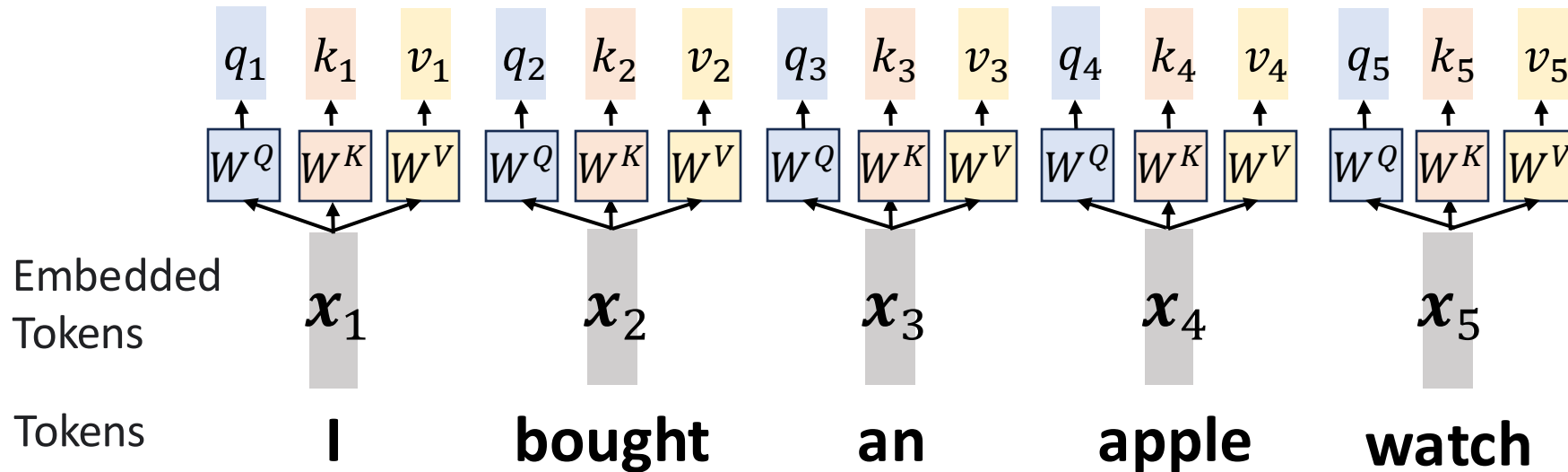
Tokens

**I**  **bought**  **an**  **apple**  **watch**

**Single-head** attention

$$\text{Attention}(Q, K, V) = V \, \text{softmax}\left(\frac{K^\top Q}{\sqrt{d_k}}\right)$$

$Q = W^Q \; x_1 \; x_2 \; x_3 \; x_4 \; x_5$

$K = W^K \; x_1 \; x_2 \; x_3 \; x_4 \; x_5$

$V = W^V \; x_1 \; x_2 \; x_3 \; x_4 \; x_5$

$q_1 \quad k_1 \quad v_1 \qquad q_2 \quad k_2 \quad v_2 \qquad q_3 \quad k_3 \quad v_3 \qquad q_4 \quad k_4 \quad v_4 \qquad q_5 \quad k_5 \quad v_5$

$W^Q \; W^K \; W^V \qquad W^Q \; W^K \; W^V \qquad W^Q \; W^K \; W^V \qquad W^Q \; W^K \; W^V \qquad W^Q \; W^K \; W^V$

Embedded Tokens

$x_1 \qquad x_2 \qquad x_3 \qquad x_4 \qquad x_5$
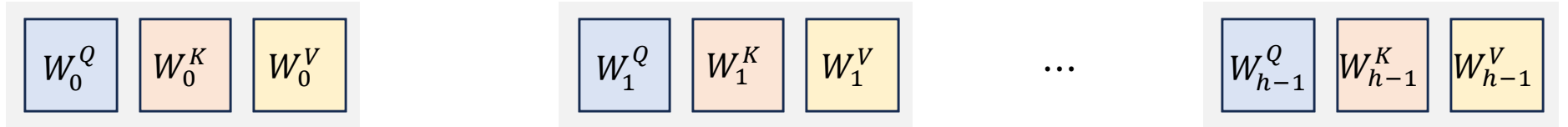
Tokens

**I**     **bought**     **an**     **apple**     **watch**

**Single-head** attention

$$\text{Attention}(Q, K, V) = V \text{ softmax}\left(\frac{K^\top Q}{\sqrt{d_k}}\right)$$

$$Q = \boxed{W^Q} \; \boxed{x_1} \; \boxed{x_2} \; \boxed{x_3} \; \boxed{x_4} \; \boxed{x_5}$$

$$K = \boxed{W^K} \; \boxed{x_1} \; \boxed{x_2} \; \boxed{x_3} \; \boxed{x_4} \; \boxed{x_5}$$

$$V = \boxed{W^V} \; \boxed{x_1} \; \boxed{x_2} \; \boxed{x_3} \; \boxed{x_4} \; \boxed{x_5}$$

$\boxed{W_0^Q} \; \boxed{W_0^K} \; \boxed{W_0^V}$   $\boxed{W_1^Q} \; \boxed{W_1^K} \; \boxed{W_1^V}$   $\ldots$   $\boxed{W_{h-1}^Q} \; \boxed{W_{h-1}^K} \; \boxed{W_{h-1}^V}$

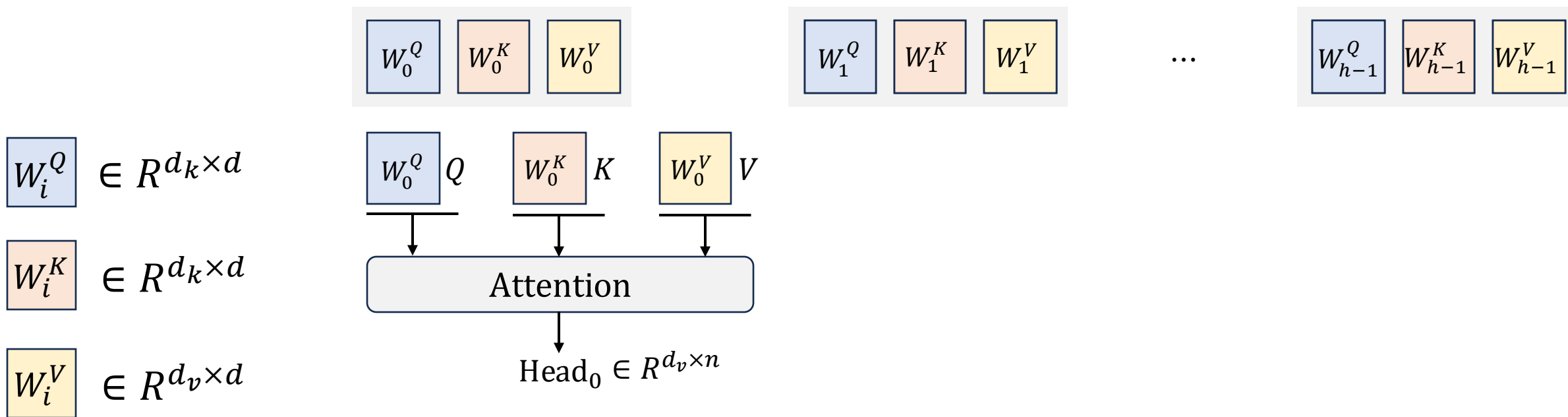$\boxed{W_i^Q} \in R^{d_k \times d}$

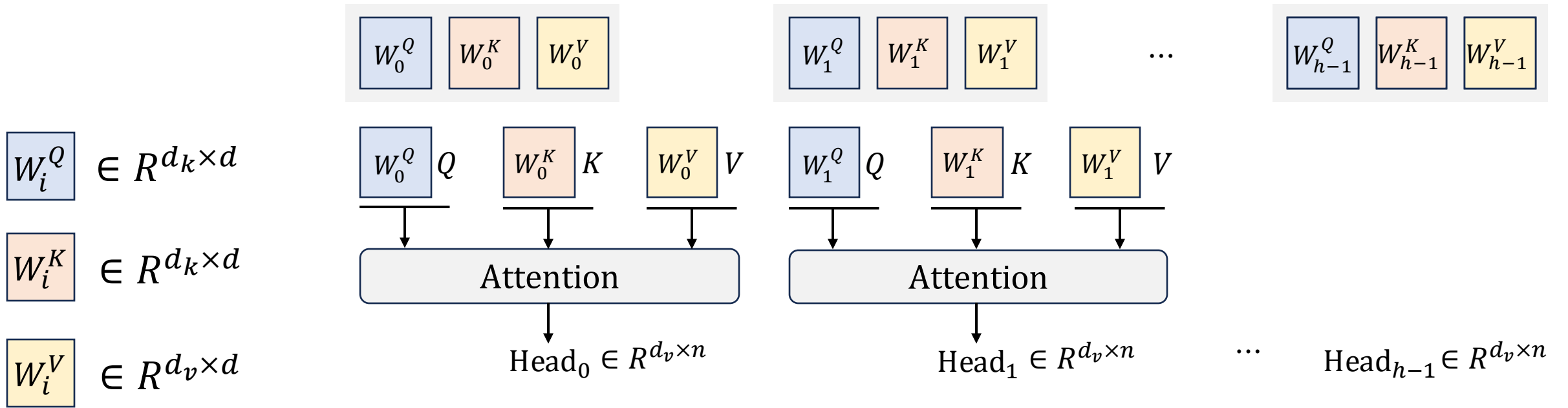$\boxed{W_i^K} \in R^{d_k \times d}$

$\boxed{W_i^V} \in R^{d_v \times d}$

**Single-head** attention

$$\text{Attention}(Q, K, V) = V \, \text{softmax}\left(\frac{K^\top Q}{\sqrt{d_k}}\right)$$

$Q \;=\; \boxed{W^Q}\;\; \boxed{x_1}\;\boxed{x_2}\;\boxed{x_3}\;\boxed{x_4}\;\boxed{x_5}$

$K \;=\; \boxed{W^K}\;\; \boxed{x_1}\;\boxed{x_2}\;\boxed{x_3}\;\boxed{x_4}\;\boxed{x_5}$

$V \;=\; \boxed{W^V}\;\; \boxed{x_1}\;\boxed{x_2}\;\boxed{x_3}\;\boxed{x_4}\;\boxed{x_5}$

$\boxed{W_0^Q}\;\boxed{W_0^K}\;\boxed{W_0^V}$  $\boxed{W_1^Q}\;\boxed{W_1^K}\;\boxed{W_1^V}$  $\cdots$  $\boxed{W_{h-1}^Q}\;\boxed{W_{h-1}^K}\;\boxed{W_{h-1}^V}$

$\boxed{W_i^Q} \in R^{d_k \times d}$

$\boxed{W_i^K} \in R^{d_k \times d}$

$\boxed{W_i^V} \in R^{d_v \times d}$

$\boxed{W_0^Q}\,Q \qquad \boxed{W_0^K}\,K \qquad \boxed{W_0^V}\,V$
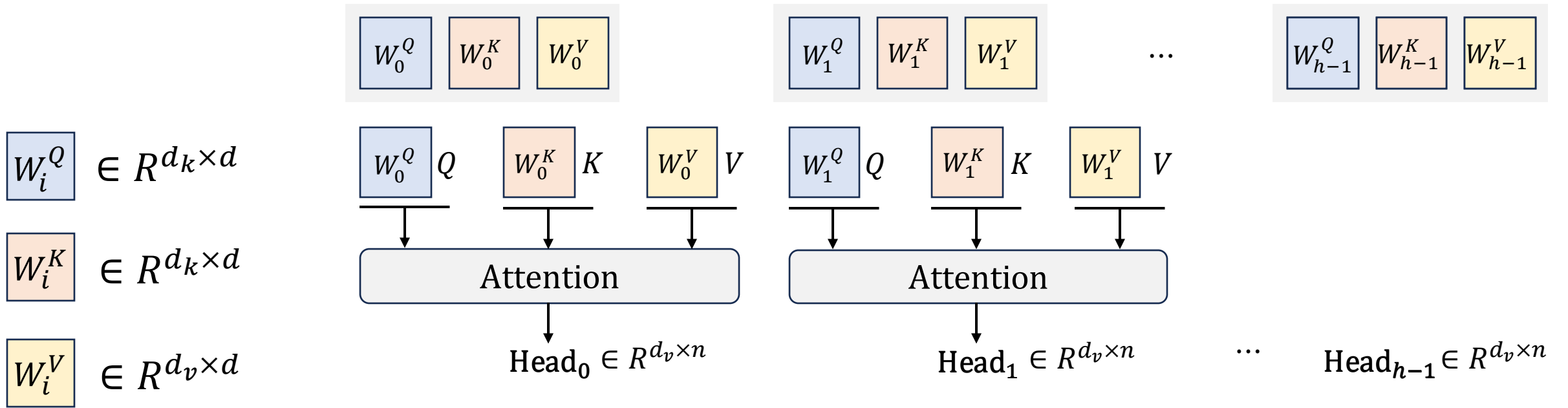
Attention

$\text{Head}_0 \in R^{d_v \times n}$

**Single-head** attention

$$\text{Attention}(Q, K, V) = V \, \text{softmax}\left(\frac{K^\top Q}{\sqrt{d_k}}\right)$$

$$Q = W^Q \; x_1 \; x_2 \; x_3 \; x_4 \; x_5$$

$$K = W^K \; x_1 \; x_2 \; x_3 \; x_4 \; x_5$$

$$V = W^V \; x_1 \; x_2 \; x_3 \; x_4 \; x_5$$

$W_0^Q \quad W_0^K \quad W_0^V$

$W_1^Q \quad W_1^K \quad W_1^V$ ... $W_{h-1}^Q \quad W_{h-1}^K \quad W_{h-1}^V$

$W_i^Q \in R^{d_k \times d}$

$W_i^K \in R^{d_k \times d}$

$W_i^V \in R^{d_v \times d}$

$W_0^Q \, Q \qquad W_0^K \, K \qquad W_0^V \, V$

$W_1^Q \, Q \qquad W_1^K \, K \qquad W_1^V \, V$

Attention

Attention

$\text{Head}_0 \in R^{d_v \times n}$

$\text{Head}_1 \in R^{d_v \times n}$ ... $\text{Head}_{h-1} \in R^{d_v \times n}$

**Single-head** attention

$$\text{Attention}(Q, K, V) = V \, \text{softmax}\left(\frac{K^\top Q}{\sqrt{d_k}}\right)$$

$Q \;=\; W^Q \; x_1 \; x_2 \; x_3 \; x_4 \; x_5$

$K \;=\; W^K \; x_1 \; x_2 \; x_3 \; x_4 \; x_5$

$V \;=\; W^V \; x_1 \; x_2 \; x_3 \; x_4 \; x_5$

$W_0^Q \quad W_0^K \quad W_0^V$ $\qquad$ $W_1^Q \quad W_1^K \quad W_1^V$ $\quad \cdots \quad$ $W_{h-1}^Q \quad W_{h-1}^K \quad W_{h-1}^V$

$W_i^Q \in R^{d_k \times d}$

$W_i^K \in R^{d_k \times d}$

$W_i^V \in R^{d_v \times d}$

$W_0^Q Q \qquad W_0^K K \qquad W_0^V V$ $\qquad$ $W_1^Q Q \qquad W_1^K K \qquad W_1^V V$

Attention $\qquad$ Attention

$\text{Head}_0 \in R^{d_v \times n}$ $\qquad$ $\text{Head}_1 \in R^{d_v \times n}$ $\quad \cdots \quad$ $\text{Head}_{h-1} \in R^{d_v \times n}$

**Single-head** attention

$$\text{Attention}(Q, K, V) = V \text{ softmax}\left(\frac{K^\top Q}{\sqrt{d_k}}\right)$$

**Multi-head** attention

$$Q = \boxed{W^Q} \; \boxed{x_1} \; \boxed{x_2} \; \boxed{x_3} \; \boxed{x_4} \; \boxed{x_5}$$

$$K = \boxed{W^K} \; \boxed{x_1} \; \boxed{x_2} \; \boxed{x_3} \; \boxed{x_4} \; \boxed{x_5}$$

$$V = \boxed{W^V} \; \boxed{x_1} \; \boxed{x_2} \; \boxed{x_3} \; \boxed{x_4} \; \boxed{x_5}$$

$\boxed{W^O} \in R^{d \times hd_v}$

$\boxed{W_i^Q} \in R^{d_k \times d}$

$\boxed{W_i^K} \in R^{d_k \times d}$

$\boxed{W_i^V} \in R^{d_v \times d}$



$$\boxed{W_0^Q} \; \boxed{W_0^K} \; \boxed{W_0^V} \qquad \boxed{W_1^Q} \; \boxed{W_1^K} \; \boxed{W_1^V} \qquad \cdots \qquad \boxed{W_{h-1}^Q} \; \boxed{W_{h-1}^K} \; \boxed{W_{h-1}^V}$$

$\boxed{W_0^Q} Q \quad \boxed{W_0^K} K \quad \boxed{W_0^V} V \qquad \boxed{W_1^Q} Q \quad \boxed{W_1^K} K \quad \boxed{W_1^V} V$

Attention      Attention

$\text{Head}_0 \in R^{d_v \times n}$      $\text{Head}_1 \in R^{d_v \times n}$    $\cdots$    $\text{Head}_{h-1} \in R^{d_v \times n}$

$$\text{MultiHeadedAttention}(Q, K, V) = \boxed{W^O} \begin{bmatrix} \text{Head}_0 \\ \text{Head}_1 \\ \vdots \\ \text{Head}_{h-1} \end{bmatrix}$$

# Feed Forward Network (FFN)

$$FFN(\boldsymbol{x}) = \boldsymbol{W}_2 \text{ReLU}(\boldsymbol{W}_1 \boldsymbol{x} + \boldsymbol{b}_1) + \boldsymbol{b}_2$$

# Feed Forward Network (FFN)

$$FFN(\boldsymbol{x}) = \boldsymbol{W}_2 \text{ReLU}(\boldsymbol{W}_1 \boldsymbol{x} + \boldsymbol{b}_1) + \boldsymbol{b}_2$$

Encoder #1

Feed Forward

Feed Forward

Feed Forward

Feed Forward

Feed Forward

Multi-head Self-Attention

Embedded Tokens

$\boldsymbol{x}_1$    $\boldsymbol{x}_2$    $\boldsymbol{x}_3$    $\boldsymbol{x}_4$    $\boldsymbol{x}_5$

Tokens

I    bought    an    apple    watch

# 💡 Positional encoding

# 💡 Positional encoding

| Position $k$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $2^3$ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | ← Slow oscillating |
| $2^2$ | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | |
| $2^1$ | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 | |
| $2^0$ | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | ← Fast oscillating |

Dimension

Positional embedding $d$

Embedded Tokens $d$    $x_1$    $x_2$    $x_3$    $x_4$    $x_5$

Tokens    **I**    **bought**    **an**    **apple**    **watch**

# 💡 Positional encoding

Position $k$

Angular frequency

$w_i = N^{-2i/d}$

$N = 100,000$

$$d \begin{bmatrix} \sin(w_0 k) \\ \cos(w_0 k) \\ \sin(w_1 k) \\ \cos(w_1 k) \\ \vdots \\ \vdots \\ \sin\left(w_{\frac{d}{2}-1} k\right) \\ \cos\left(w_{\frac{d}{2}-1} k\right) \end{bmatrix}$$

← Fast oscillating

← Slow oscillating

Positional embedding · $d$

Embedded Tokens · $d$    $\boldsymbol{x}_1$      $\boldsymbol{x}_2$      $\boldsymbol{x}_3$      $\boldsymbol{x}_4$      $\boldsymbol{x}_5$

Tokens    **I**     **bought**     **an**     **apple**     **watch**

# 💡 Positional encoding

👍 **Normalized Range**

👍 **Unique identifier, unlimited length**

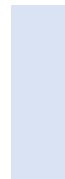👍 **Relative positions as linear transform**

Position $k$

Angular frequency

$$w_i = N^{-2i/d}$$

$N = 100{,}000$

$$d \begin{bmatrix} \sin(w_0 k) \\ \cos(w_0 k) \\ \sin(w_1 k) \\ \cos(w_1 k) \\ \vdots \\ \sin\left(w_{\frac{d}{2}-1} k\right) \\ \cos\left(w_{\frac{d}{2}-1} k\right) \end{bmatrix}$$

$$\begin{bmatrix} \sin(w_i(k + \Delta k)) \\ \cos(w_i(k + \Delta k)) \end{bmatrix} = \begin{bmatrix} \sin(w_i k)\cos(w_i \Delta k) + \cos(w_i k)\sin(w_i \Delta k) \\ \cos(w_i k)\cos(w_i \Delta k) - \sin(w_i k)\sin(w_i \Delta k) \end{bmatrix}$$

Positional embedding

$d$



Embedded Tokens

$d$ $\boldsymbol{x}_1$ $\boldsymbol{x}_2$ $\boldsymbol{x}_3$ $\boldsymbol{x}_4$ $\boldsymbol{x}_5$

Tokens  **I**  **bought**  **an**  **apple**  **watch**

# 💡 Positional encoding

👍 **Normalized Range**

👍 **Unique identifier, unlimited length**

👍 **Relative positions as linear transform**

Position $k$

Angular frequency

$w_i = N^{-2i/d}$

$N = 100,000$

$$d \begin{bmatrix} \sin(w_0 k) \\ \cos(w_0 k) \\ \sin(w_1 k) \\ \cos(w_1 k) \\ \vdots \\ \sin\left(w_{\frac{d}{2}-1} k\right) \\ \cos\left(w_{\frac{d}{2}-1} k\right) \end{bmatrix}$$

$$\begin{bmatrix} \sin(w_i(k + \Delta k)) \\ \cos(w_i(k + \Delta k)) \end{bmatrix} = \begin{bmatrix} \sin(w_i k)\cos(w_i \Delta k) + \cos(w_i k)\sin(w_i \Delta k) \\ \cos(w_i k)\cos(w_i \Delta k) - \sin(w_i k)\sin(w_i \Delta k) \end{bmatrix}$$

$$= \begin{bmatrix} \cos(w_i \Delta k) & \sin(w_i \Delta k) \\ -\sin(w_i \Delta k) & \cos(w_i \Delta k) \end{bmatrix} \begin{bmatrix} \sin(w_i k) \\ \cos(w_i k) \end{bmatrix}$$

Positional embedding
$d$   $\boldsymbol{P}_1$     $\boldsymbol{P}_2$     $\boldsymbol{P}_3$     $\boldsymbol{P}_4$     $\boldsymbol{P}_5$     $P_{k+\Delta k} = MP_k$

Embedded Tokens
$d$   $\boldsymbol{x}_1$     $\boldsymbol{x}_2$     $\boldsymbol{x}_3$     $\boldsymbol{x}_4$     $\boldsymbol{x}_5$

Tokens    **I**     **bought**     **an**     **apple**     **watch**

# 💡 Positional encoding

Position $k$

Angular frequency

$w_i = N^{-2i/d}$

$N = 100{,}000$

$$d \begin{bmatrix} \sin(w_0 k) \\ \cos(w_0 k) \\ \sin(w_1 k) \\ \cos(w_1 k) \\ \vdots \\ \vdots \\ \sin\left(w_{\frac{d}{2}-1} k\right) \\ \cos\left(w_{\frac{d}{2}-1} k\right) \end{bmatrix}$$

Sinusoidal positional encoding

Relative positional encoding

KERPLE          RoPE          CoPE

NoPE          YaRN          FIRE

Positional embedding     $d$     $\boldsymbol{P}_1$          $\boldsymbol{P}_2$          $\boldsymbol{P}_3$          $\boldsymbol{P}_4$          $\boldsymbol{P}_5$

Embedded Tokens     $d$     $\boldsymbol{x}_1$          $\boldsymbol{x}_2$          $\boldsymbol{x}_3$          $\boldsymbol{x}_4$          $\boldsymbol{x}_5$

Tokens          I          bought          an          apple          watch

Encoder #1

Multi-head Self-Attention

Feed Forward

Embedded Tokens

$d$

$\boldsymbol{x}_1$ $\boldsymbol{P}_1$   $\boldsymbol{x}_2$ $\boldsymbol{P}_2$   $\boldsymbol{x}_3$ $\boldsymbol{P}_3$   $\boldsymbol{x}_4$ $\boldsymbol{P}_4$   $\boldsymbol{x}_5$ $\boldsymbol{P}_5$

Tokens

**I**   **bought**   **an**   **apple**   **watch**

| | | | | |
|---|---|---|---|---|
| Feed Forward | Feed Forward | Feed Forward | Feed Forward | Feed Forward |

Encoder #2

Multi-head Self-Attention

| | | | | |
|---|---|---|---|---|
| Feed Forward | Feed Forward | Feed Forward | Feed Forward | Feed Forward |

Encoder #1

Multi-head Self-Attention

STACK MORE LAYERS

NEURAL NETWORKS

LAYERS

LAYERS

Embedded Tokens

$d$

$x_1$ $P_1$   $x_2$ $P_2$   $x_3$ $P_3$   $x_4$ $P_4$   $x_5$ $P_5$

Tokens

**I**   **bought**   **an**   **apple**   **watch**

# 💡 **Residual connection**



Feed Forward

Multi-head Self-Attention

Embedded Tokens

$d$

$\boldsymbol{x}_1$ $\boldsymbol{P}_1$  $\boldsymbol{x}_2$ $\boldsymbol{P}_2$  $\boldsymbol{x}_3$ $\boldsymbol{P}_3$  $\boldsymbol{x}_4$ $\boldsymbol{P}_4$  $\boldsymbol{x}_5$ $\boldsymbol{P}_5$

Tokens

**I**   **bought**   **an**   **apple**   **watch**

💡 **Residual connection**

Feed Forward

Feed Forward

Feed Forward

Feed Forward

Feed Forward

Multi-head Self-Attention

Embedded Tokens

$d$

$\boldsymbol{x}_1$ $\boldsymbol{P}_1$    $\boldsymbol{x}_2$ $\boldsymbol{P}_2$    $\boldsymbol{x}_3$ $\boldsymbol{P}_3$    $\boldsymbol{x}_4$ $\boldsymbol{P}_4$    $\boldsymbol{x}_5$ $\boldsymbol{P}_5$

Tokens

**I**        **bought**        **an**        **apple**        **watch**

Residual connection

Layer normalization

💡 **Residual connection**

💡 **Layer normalization**

$$\text{LayerNorm}(\boldsymbol{x}) =$$

$$\gamma\left(\frac{\boldsymbol{x} - \text{mean}(\boldsymbol{x})}{\sqrt{\text{Variance}(\mathbf{x}) + \epsilon}}\right) + \beta$$

$\gamma, \beta \in R$

Learnable parameters

| LayerNorm | LayerNorm | LayerNorm | LayerNorm | LayerNorm |

⊕

| Feed Forward | Feed Forward | Feed Forward | Feed Forward | Feed Forward |

| LayerNorm | LayerNorm | LayerNorm | LayerNorm | LayerNorm |

⊕

Multi-head Self-Attention

⊕

Embedded Tokens

$d$  $\boldsymbol{x}_1$ $\boldsymbol{P}_1$   $\boldsymbol{x}_2$ $\boldsymbol{P}_2$   $\boldsymbol{x}_3$ $\boldsymbol{P}_3$   $\boldsymbol{x}_4$ $\boldsymbol{P}_4$   $\boldsymbol{x}_5$ $\boldsymbol{P}_5$

Tokens

**I**     **bought**     **an**     **apple**     **watch**

💡 **Residual connection**

💡 **Layer normalization**

$$\mathrm{LayerNorm}(\boldsymbol{x}) =$$

$$\gamma \left( \frac{\boldsymbol{x} - \mathrm{mean}(\boldsymbol{x})}{\sqrt{\mathrm{Variance}(\mathbf{x}) + \epsilon}} \right) + \beta$$

$$\gamma, \beta \in R$$

Learnable parameters



[Xiong et al. 2020]
On Layer Normalization in the Transformer Architecture

Encoder #6

Encoder #5

Encoder #4

Encoder #3

Encoder #2

Encoder #1

Embedded Tokens

$d$

$x_1$ $P_1$   $x_2$ $P_2$   $x_3$ $P_3$   $x_4$ $P_4$   $x_5$ $P_5$

Tokens

**I**   **bought**   **an**   **apple**   **watch**

| Encoder #6 |
| Encoder #5 |
| Encoder #4 |
| Encoder #3 |
| Encoder #2 |
| Encoder #1 |

Embedded Tokens

$d$

$x_1$ $P_1$    $x_2$ $P_2$    $x_3$ $P_3$    $x_4$ $P_4$

Tokens

**How**    **are**    **you**    **?**

Encoder #6

⋮

Encoder #1

Embedded
Tokens

$d$

$x_1$ $P_1$    $x_2$ $P_2$    $x_3$ $P_3$    $x_4$ $P_4$

Tokens

**How**    **are**    **you**    **?**

你 好 吃 飽 嗎 … ？ <end>

Softmax

Linear

Decoder #$L$

Decoder #1

$e_1$ $e_2$ $e_3$ $e_4$

Encoder #$L$

Encoder #1

$d$ $x_1$ $P_1$ $x_2$ $P_2$ $x_3$ $P_3$ $x_4$ $P_4$

How are you ?

$d$ $z_1$ $P_1$ $z_2$ $P_2$

<start> 你

| $e_1$ | $e_2$ | $e_3$ | $e_4$ |

Encoder #$L$

⋮

Encoder #1

$x_1$ $P_1$　$x_2$ $P_2$　$x_3$ $P_3$　$x_4$ $P_4$

**How**　**are**　**you**　**?**

Multi-head Self-Attention

Decoder #1

$z_1$ $P_1$　$z_2$ $P_2$　$z_3$ $P_3$　$z_4$ $P_4$

**<start>**　你　好　嗎

Decoder #1

Multi-head Self-Attention

$d$

$\boldsymbol{z}_1$ $\boldsymbol{P}_1$ $\boldsymbol{z}_2$ $\boldsymbol{P}_2$ $\boldsymbol{z}_3$ $\boldsymbol{P}_3$ $\boldsymbol{z}_4$ $\boldsymbol{P}_4$

<start>　　你　　　好　　　嗎

Decoder #1

Multi-head Self-Attention

$d$

$\boldsymbol{z}_1$ $\boldsymbol{P}_1$ $\boldsymbol{z}_2$ $\boldsymbol{P}_2$ $\boldsymbol{z}_3$ $\boldsymbol{P}_3$ $\boldsymbol{z}_4$ $\boldsymbol{P}_4$
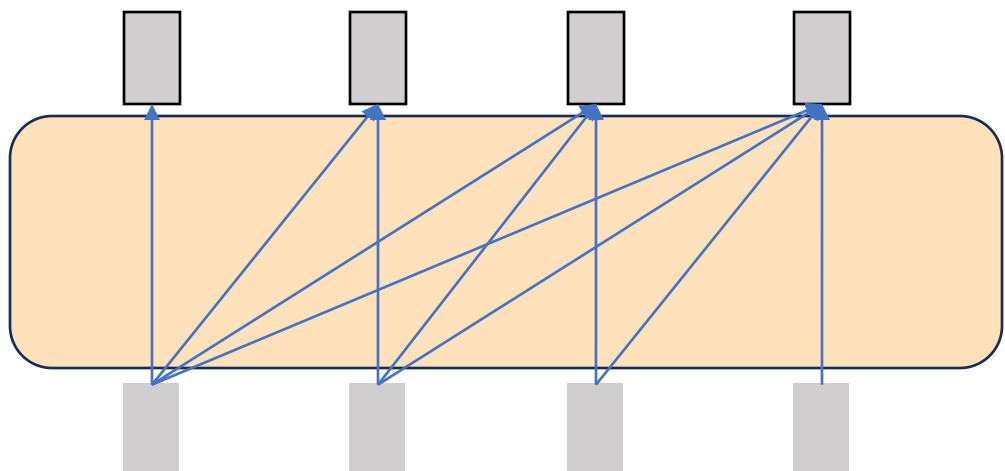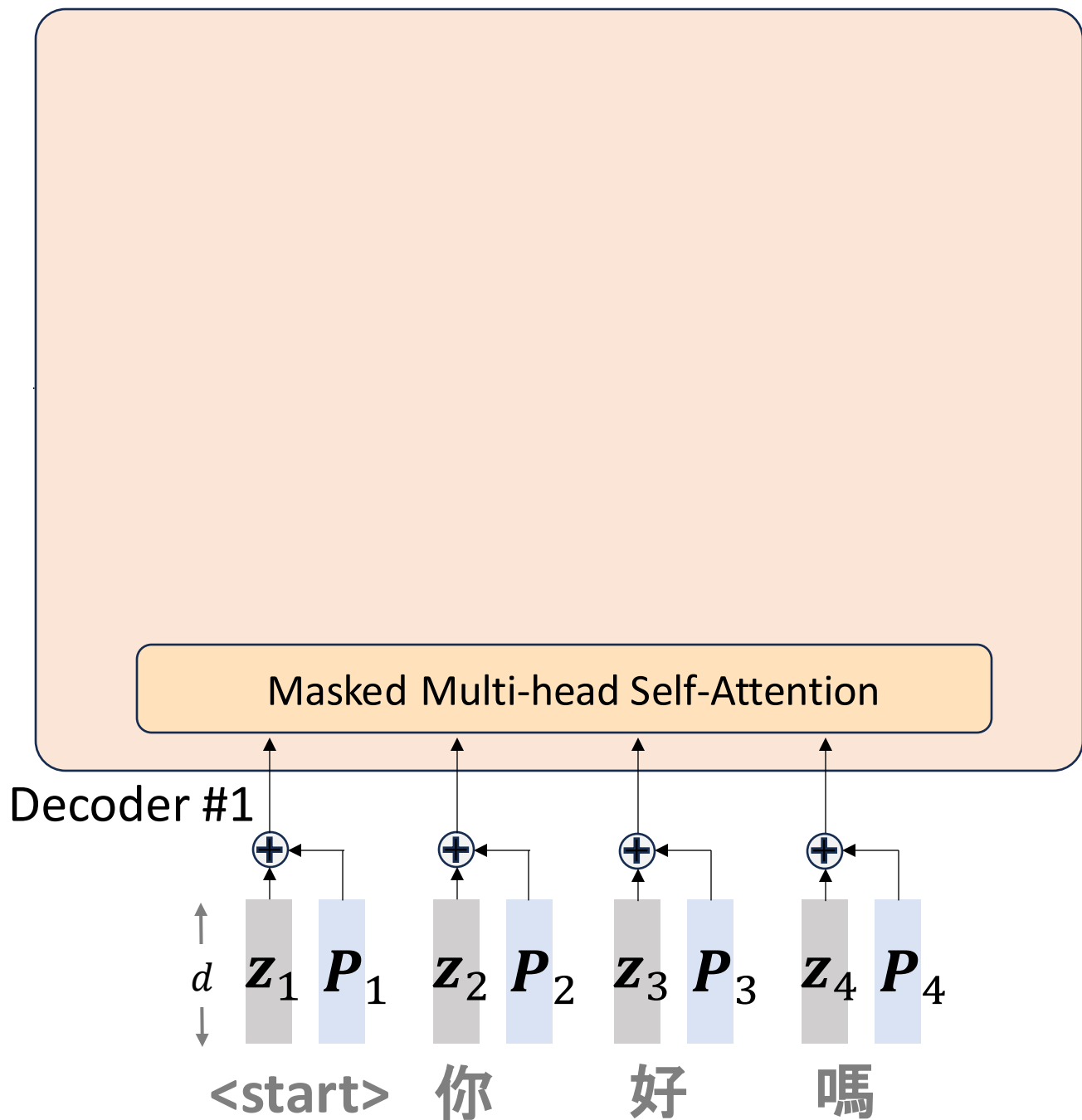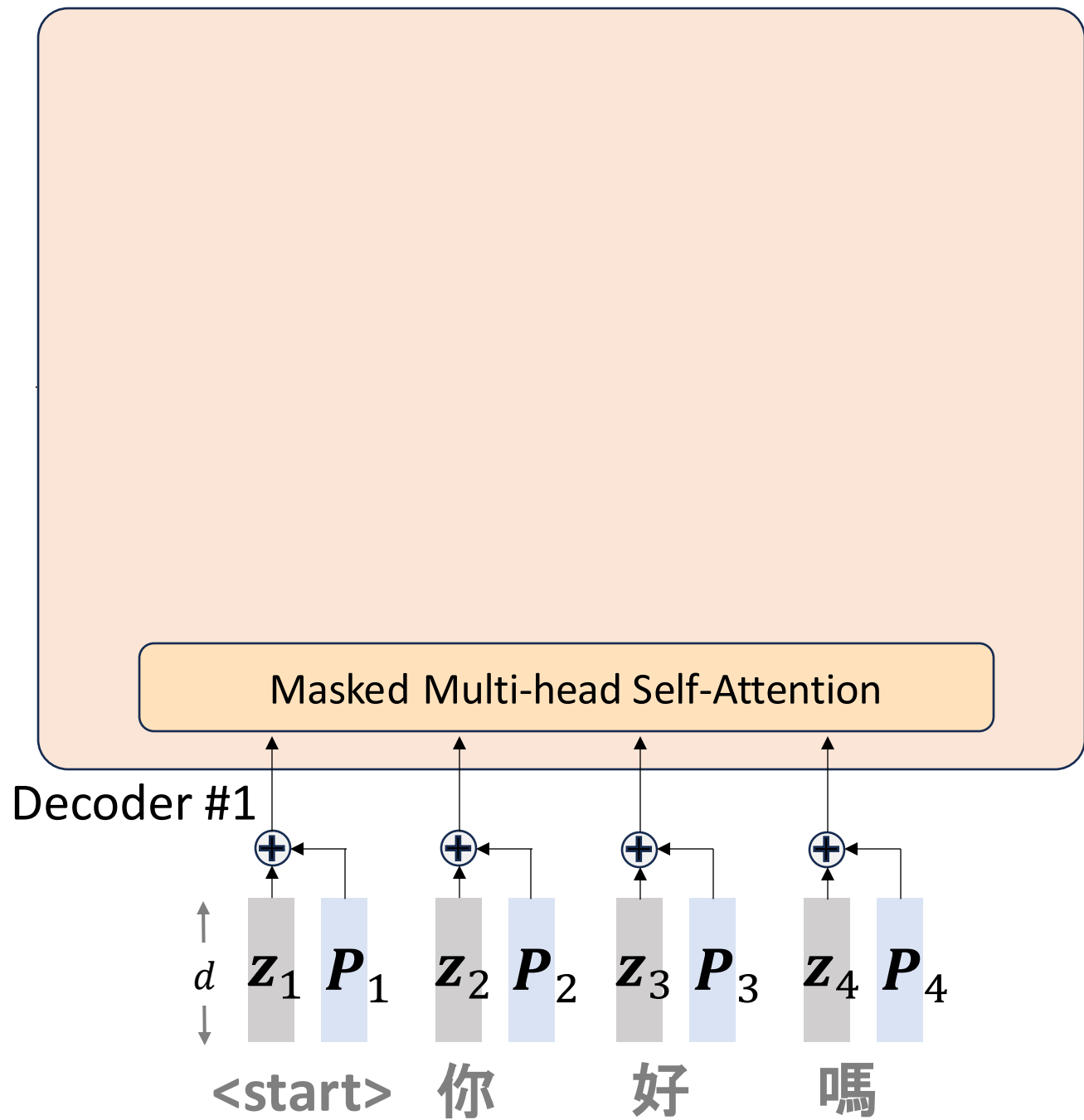
<start> 你 好 嗎

$$\text{Attention}(Q, K, V) = V \ \text{softmax}\left(\frac{K^{\top}Q}{\sqrt{d_k}}\right)$$

$$\text{MaskedAttention}(Q, K, V) = V \ \text{softmax}\left(\frac{K^{\top}Q}{\sqrt{d_k}} + M\right)$$

$$M = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ -\infty & 0 & 0 & 0 & 0 \\ -\infty & -\infty & 0 & 0 & 0 \\ -\infty & -\infty & -\infty & 0 & 0 \\ -\infty & -\infty & -\infty & -\infty & 0 \end{bmatrix}$$

Multi-head Self-Attention

Decoder #1

$d$

$\boldsymbol{z}_1$ $\boldsymbol{P}_1$ $\quad$ $\boldsymbol{z}_2$ $\boldsymbol{P}_2$ $\quad$ $\boldsymbol{z}_3$ $\boldsymbol{P}_3$ $\quad$ $\boldsymbol{z}_4$ $\boldsymbol{P}_4$

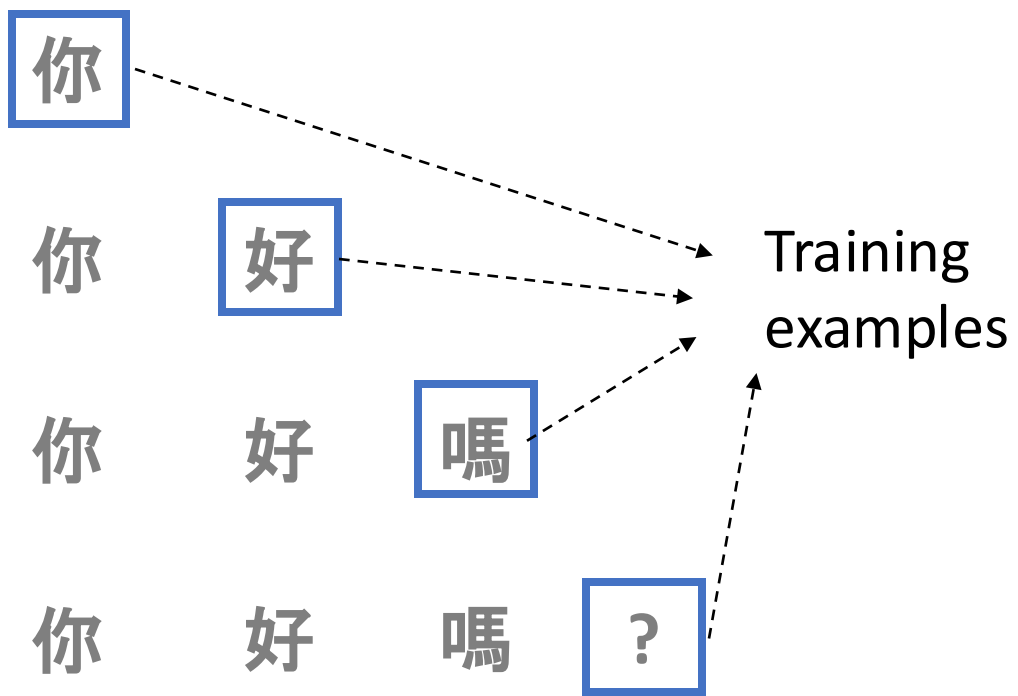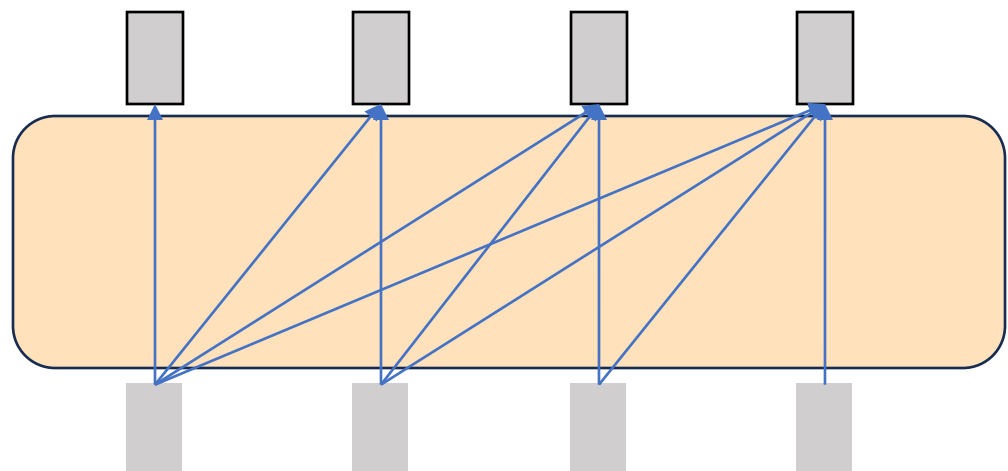&lt;start&gt; $\quad$ 你 $\quad$ 好 $\quad$ 嗎

$$\text{Attention}(Q, K, V) = V \, \text{softmax}\left(\frac{K^\top Q}{\sqrt{d_k}}\right)$$
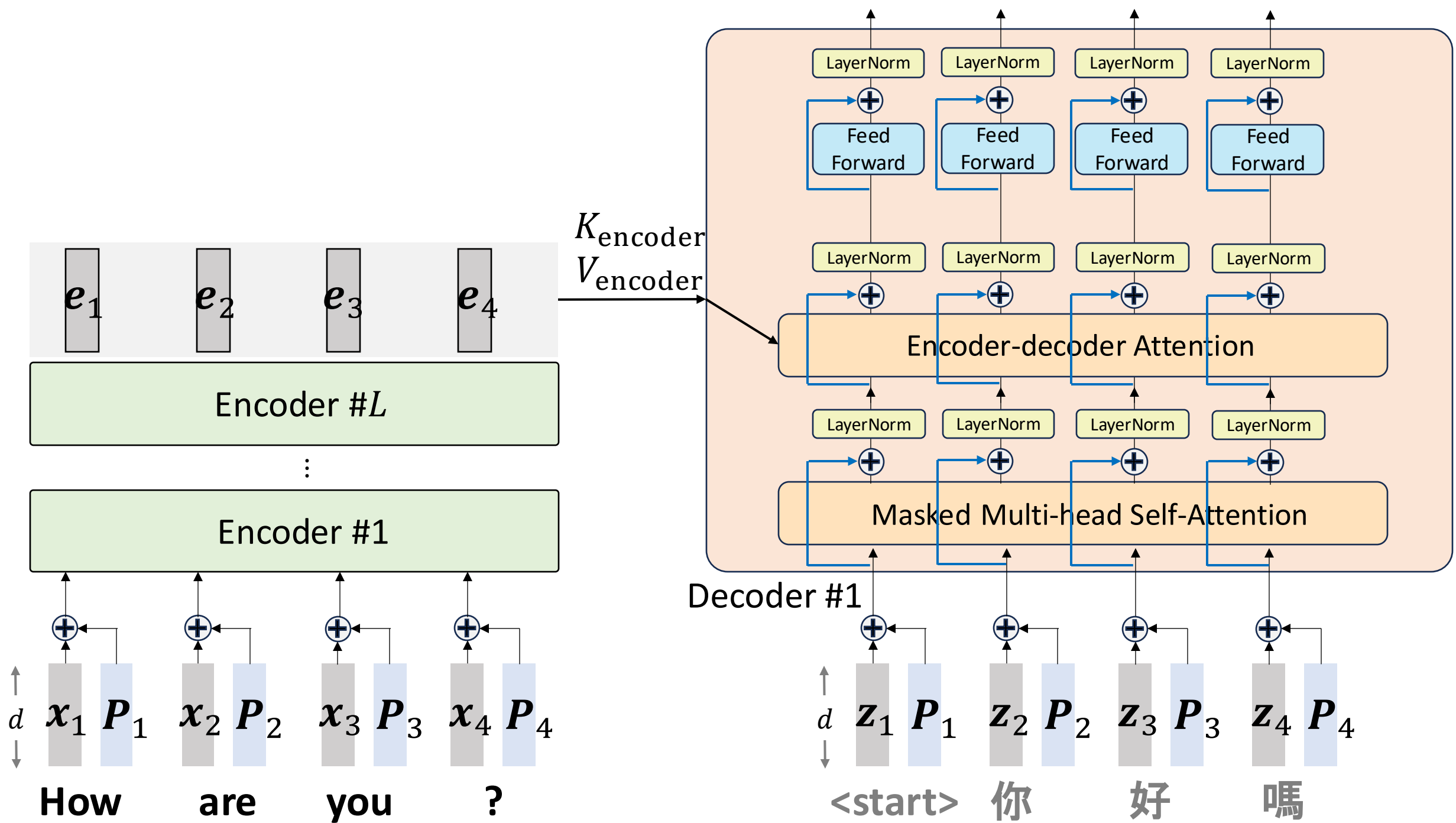
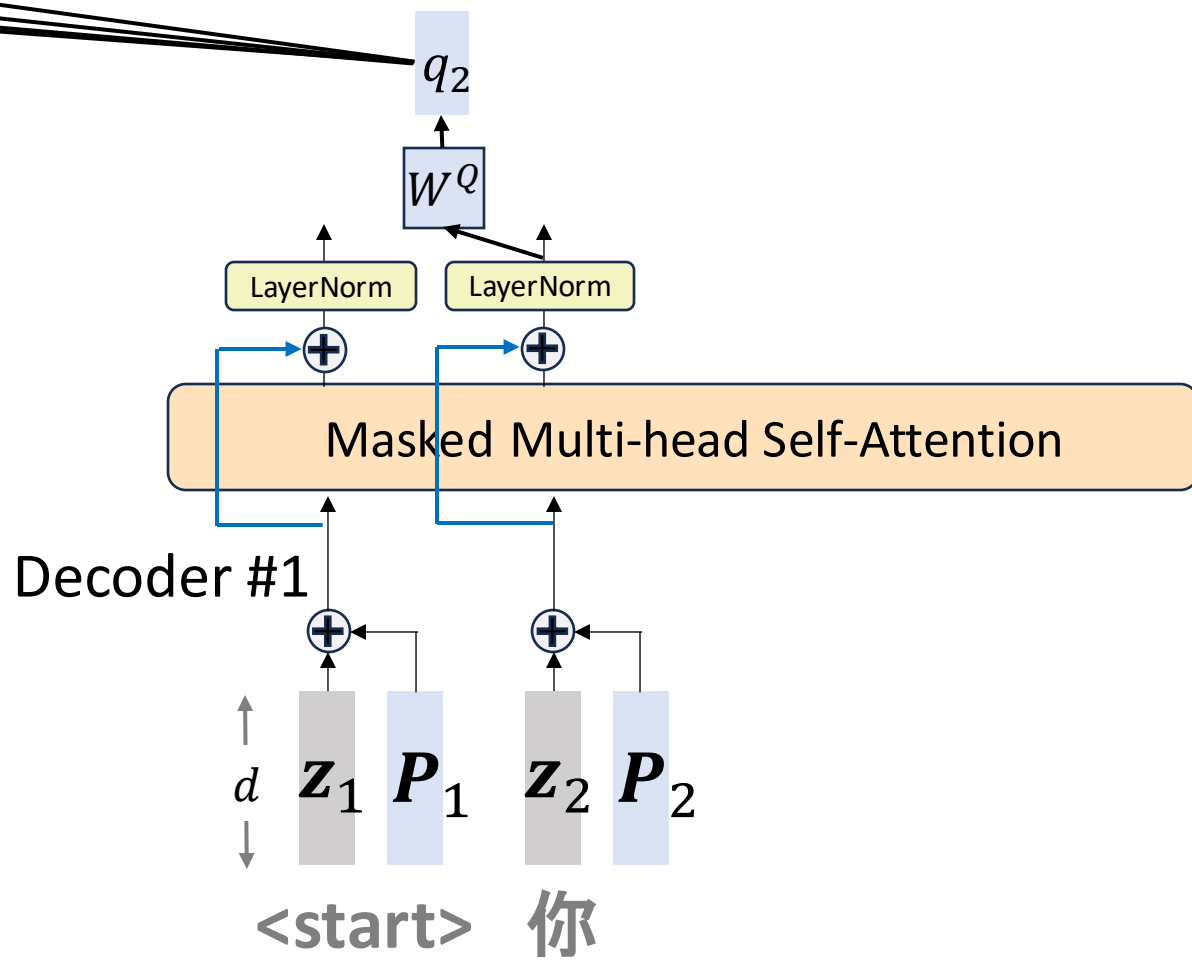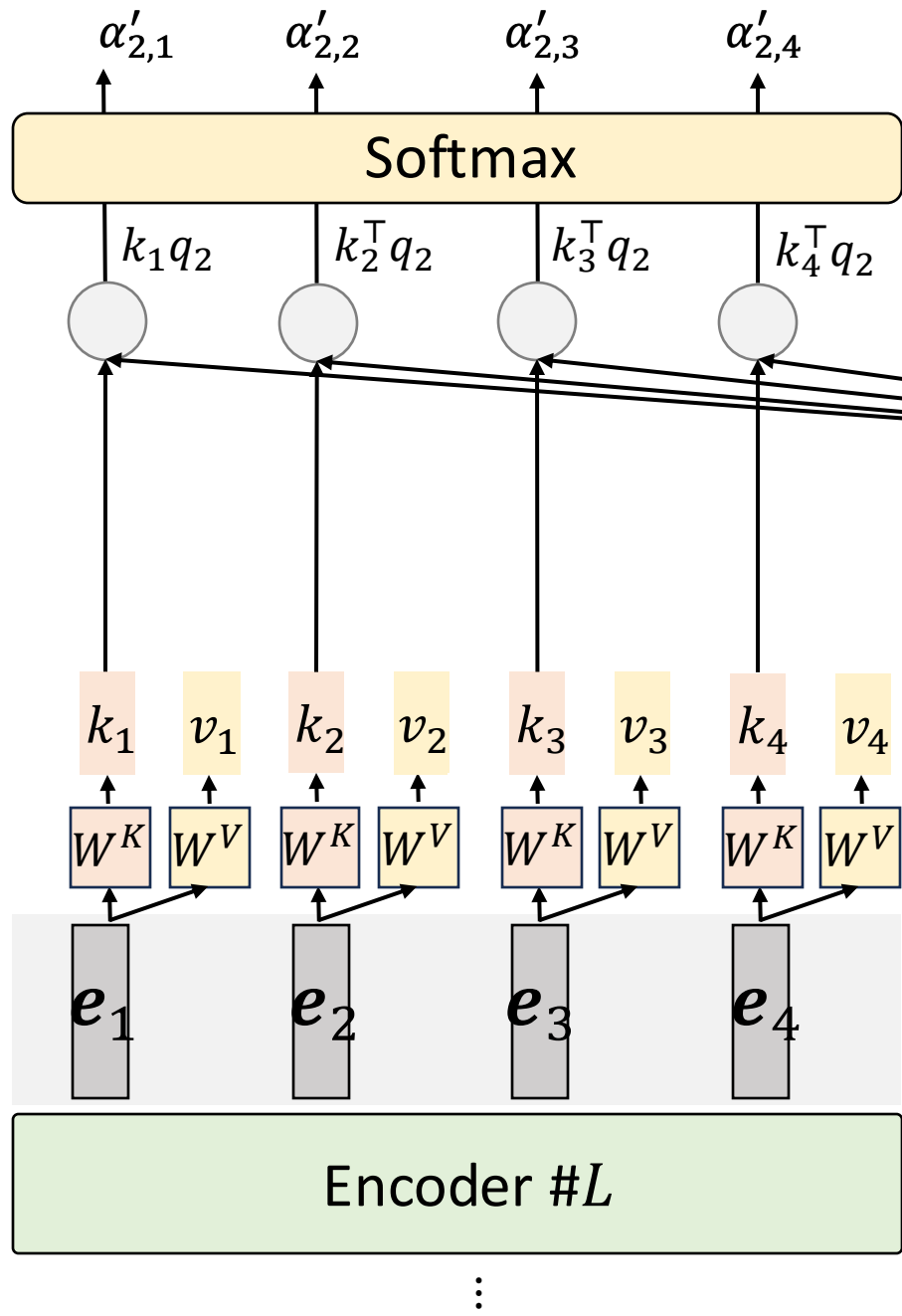$$\text{MaskedAttention}(Q, K, V) = V \, \text{softmax}\left(\frac{K^\top Q}{\sqrt{d_k}} + M\right)$$

$$M = \begin{array}{|c|c|c|c|c|}
\hline
0 & 0 & 0 & 0 & 0 \\
\hline
-\infty & 0 & 0 & 0 & 0 \\
\hline
-\infty & -\infty & 0 & 0 & 0 \\
\hline
-\infty & -\infty & -\infty & 0 & 0 \\
\hline
-\infty & -\infty & -\infty & -\infty & 0 \\
\hline
\end{array}$$

Decoder #1

Masked Multi-head Self-Attention

$d$   $\boldsymbol{z}_1$   $\boldsymbol{P}_1$   $\boldsymbol{z}_2$   $\boldsymbol{P}_2$   $\boldsymbol{z}_3$   $\boldsymbol{P}_3$   $\boldsymbol{z}_4$   $\boldsymbol{P}_4$

**\<start\>**    你    好    嗎

Training examples

Decoder #1

Masked Multi-head Self-Attention

$z_1$ $P_1$ $z_2$ $P_2$ $z_3$ $P_3$ $z_4$ $P_4$

$d$

<start> 你 好 嗎

$$\alpha'_{2,1} \quad \alpha'_{2,2} \quad \alpha'_{2,3} \quad \alpha'_{2,4}$$

Softmax

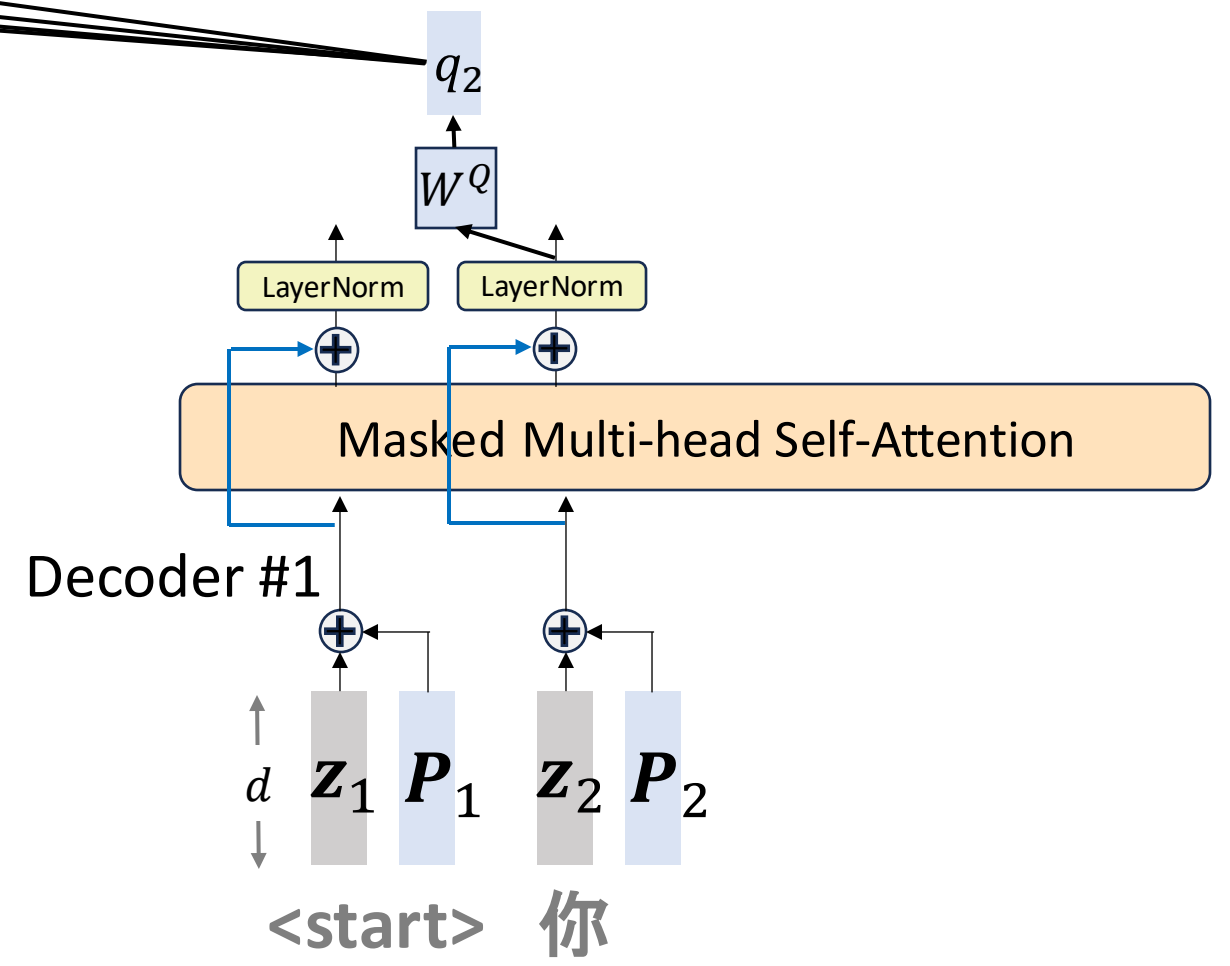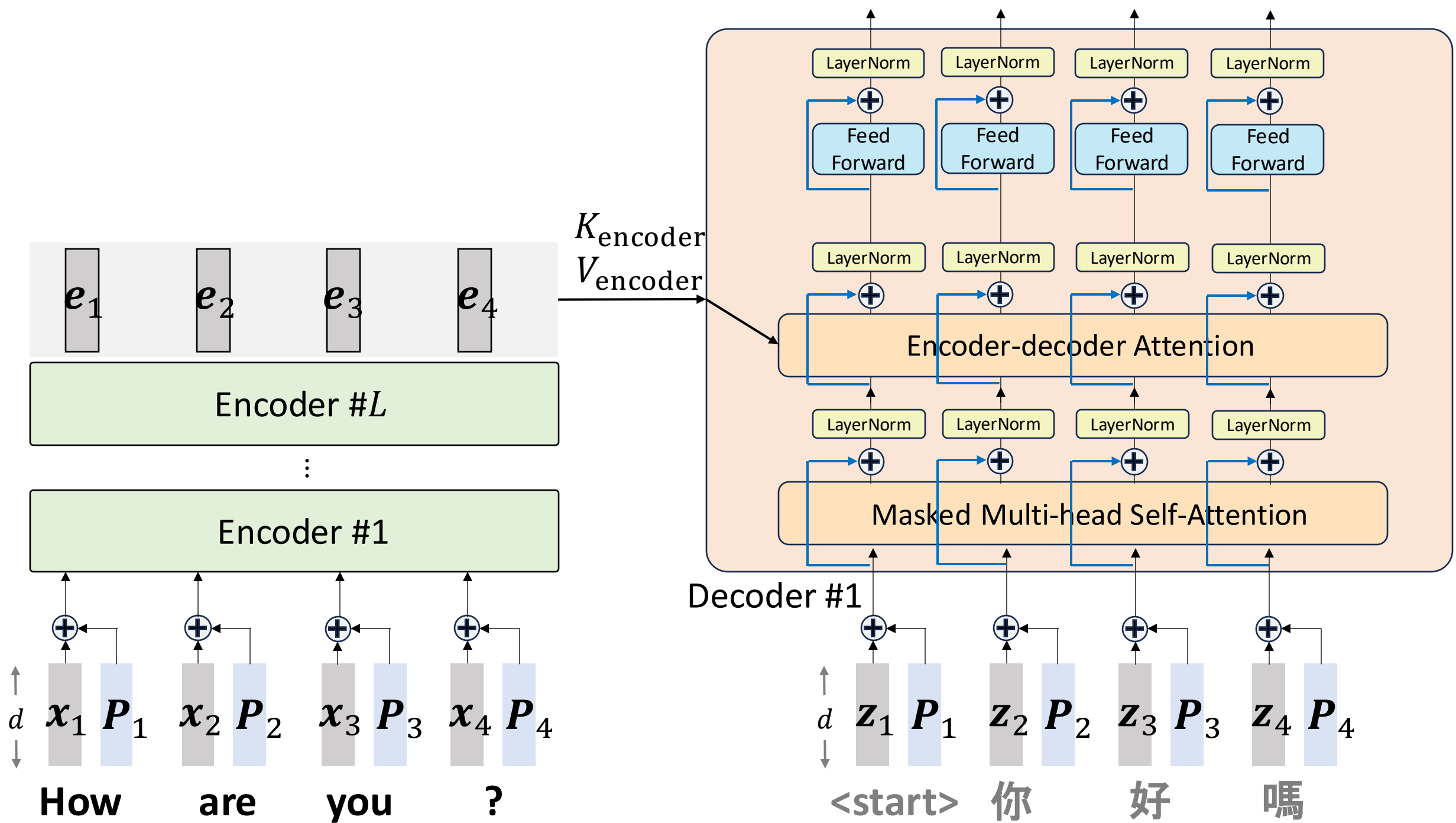$$k_1 q_2 \quad k_2^\top q_2 \quad k_3^\top q_2 \quad k_4^\top q_2$$

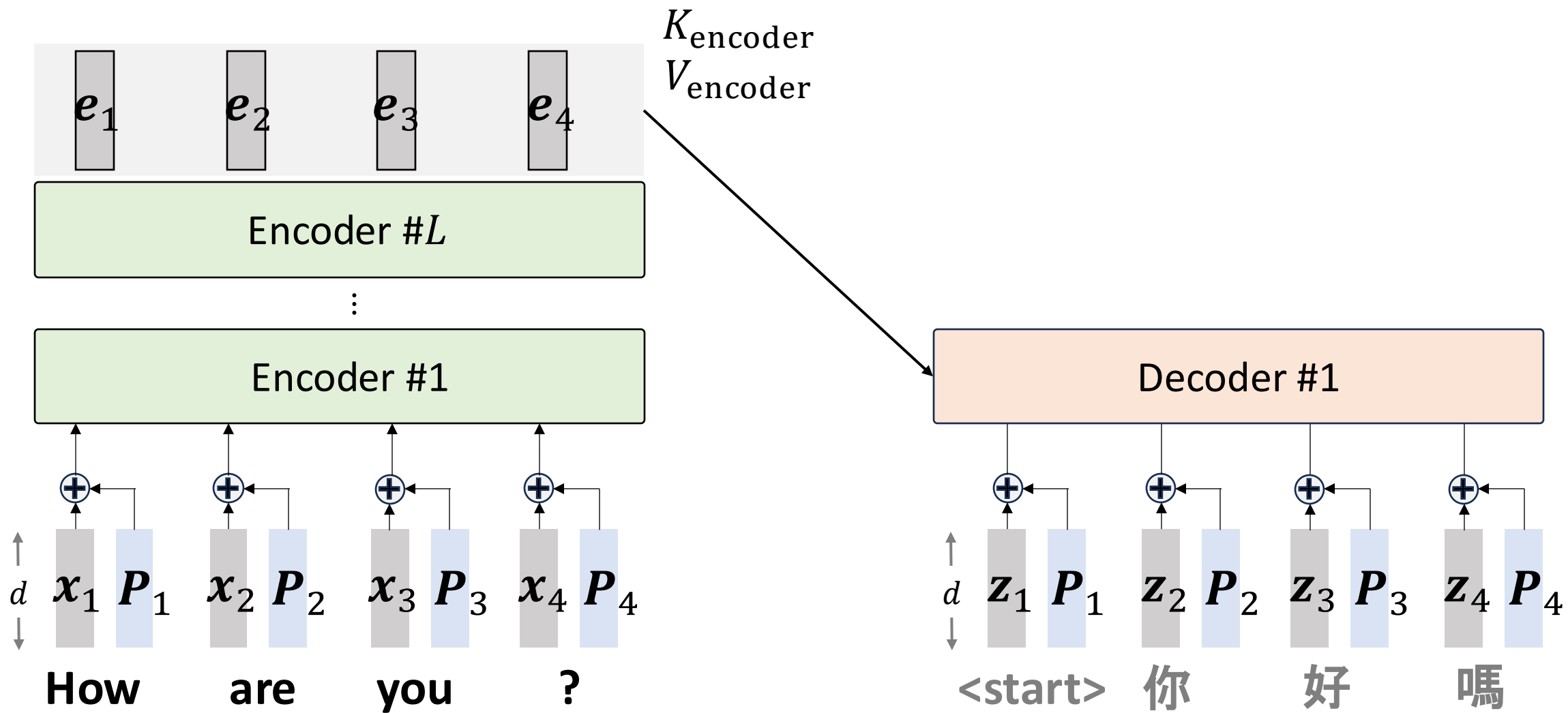$$\mathbf{z}'_2 = W^O( \alpha'_{2,1} v_1 + \alpha'_{2,2} v_2 + \alpha'_{2,3} v_3 + \alpha'_{2,4} v_4 )$$

Cross-attention

Encoder-decoder attention

$q_2$

$W^Q$

LayerNorm | LayerNorm

$+$ | $+$

Masked Multi-head Self-Attention

$$k_1 \quad v_1 \quad k_1 \quad v_2 \quad k_1 \quad v_3 \quad k_1 \quad v_4$$

$$W^K \quad W^V \quad W^K \quad W^V \quad W^K \quad W^V \quad W^K \quad W^V$$

$$\mathbf{e}_1 \quad \mathbf{e}_2 \quad \mathbf{e}_3 \quad \mathbf{e}_4$$

Decoder #1

Encoder #L

$$d \quad \mathbf{z}_1 \quad \mathbf{P}_1 \quad \mathbf{z}_2 \quad \mathbf{P}_2$$

**&lt;start&gt;** 你

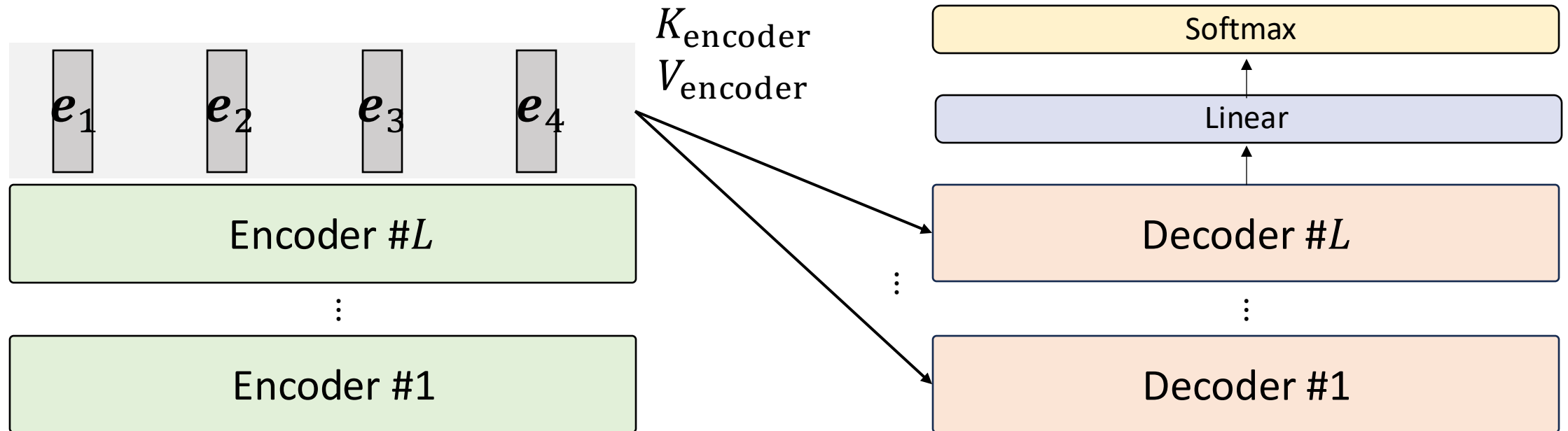(ignore the scaling $1/\sqrt{d_k}$ here for simplicity)
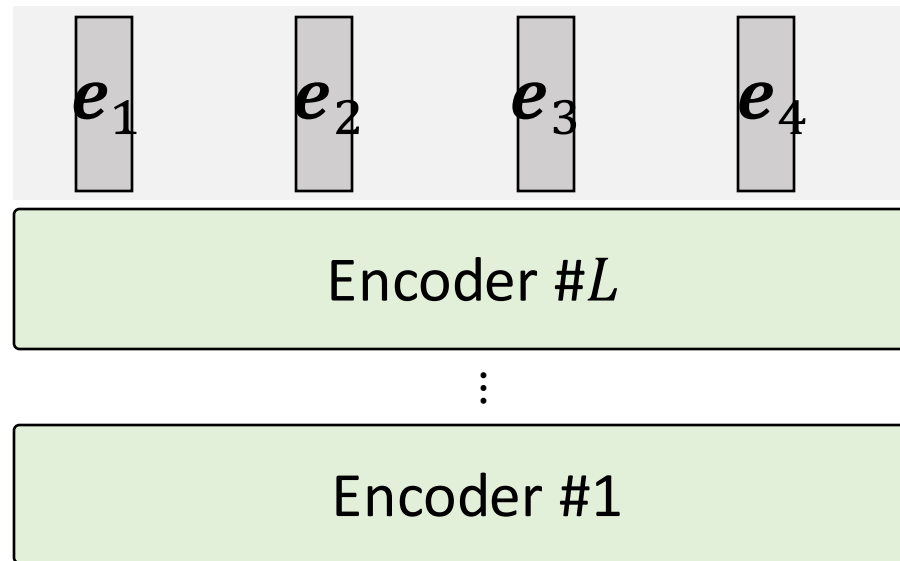
# Encoder-Decoder Transformer

**Examples**:

Attention is all you need, T5, BART.

**Good for**:

Machine translation, summarization. QA (when input/target are sufficiently different)
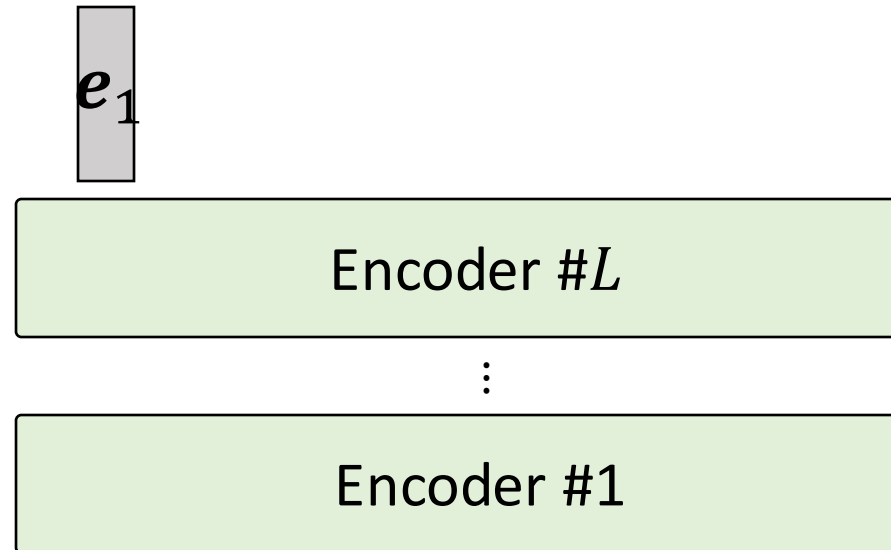
# Encoder-Decoder Transformer

**Examples**:

Attention is all you need, T5, BART.

**Good for**:

Machine translation, summarization. QA (when input/target are sufficiently different)

$e_1$ $e_2$ $e_3$ $e_4$

Encoder #$L$

⋮

Encoder #1

# Encoder-only Transformer

**Examples**:

BERT, RoBERTa, DeBERTa, X-BERT

**Good for**:

Classification, sequence tagging, sentiment analysis
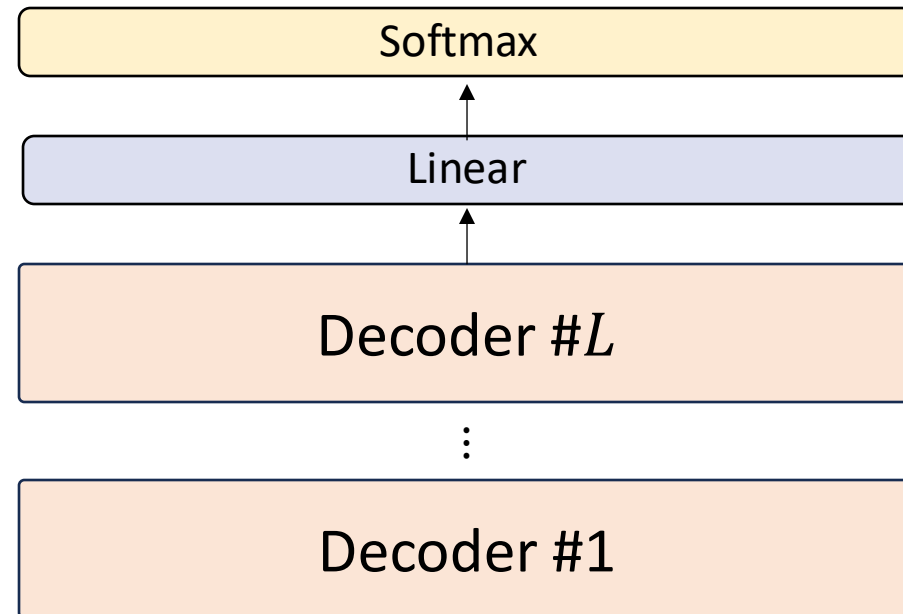(Understand text, but not generate them)

$e_1$

| Encoder #$L$ |
| :---: |

⋮

| Encoder #1 |
| :---: |

# Decoder-only Transformer

**Examples**:
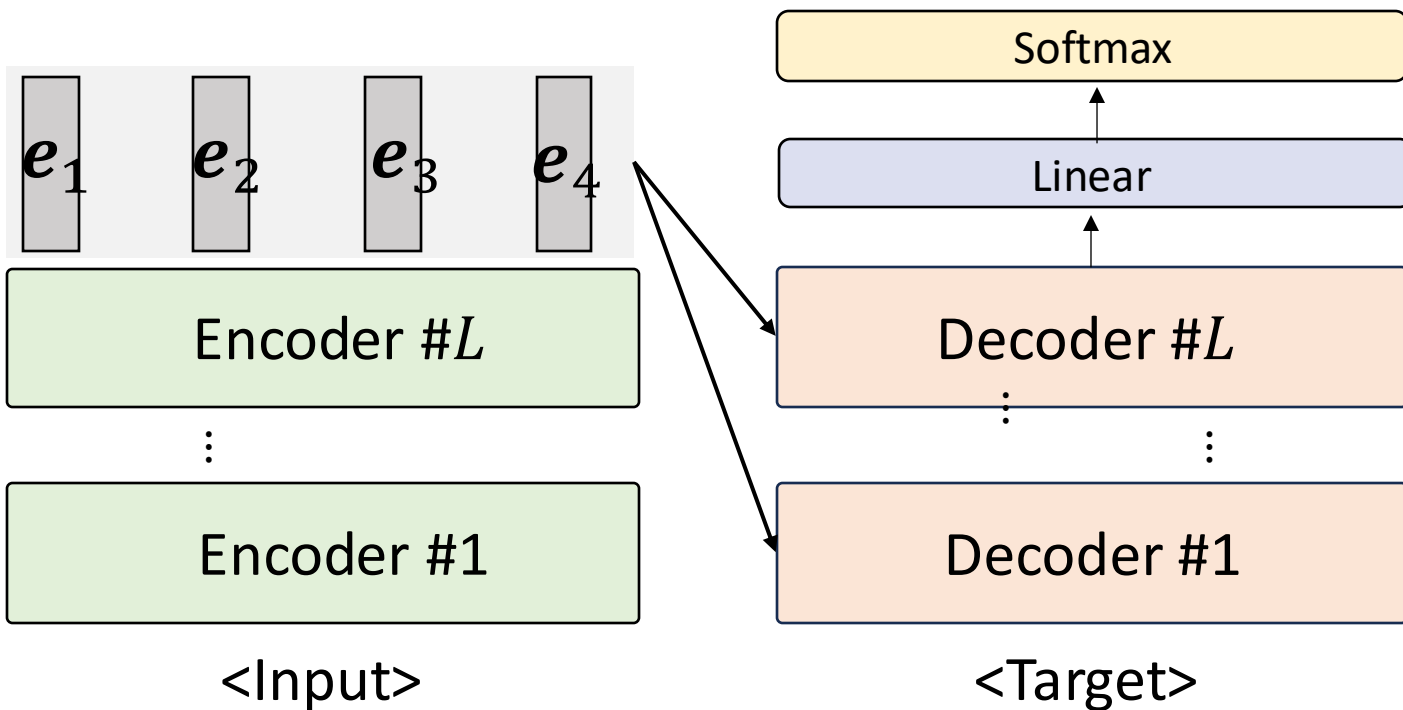
GPT-X (OpenAI), PaLM (Google), LLaMA (Meta) BLOOM (BigScience)
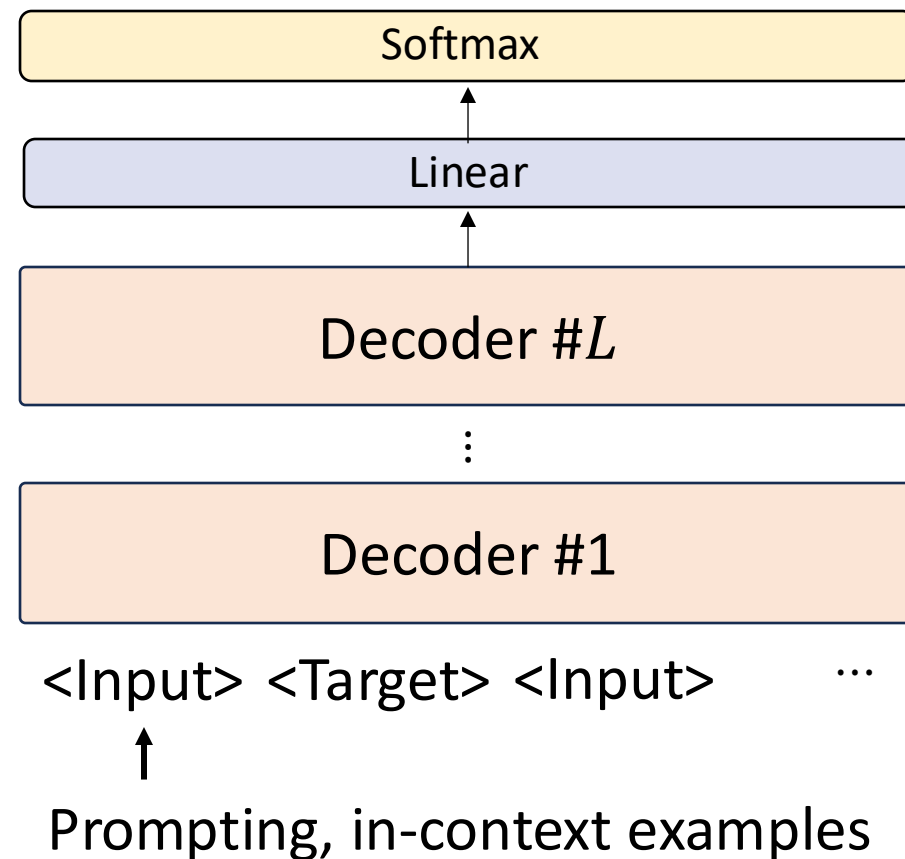
**Good for**:

Text generation, multi-round conversation
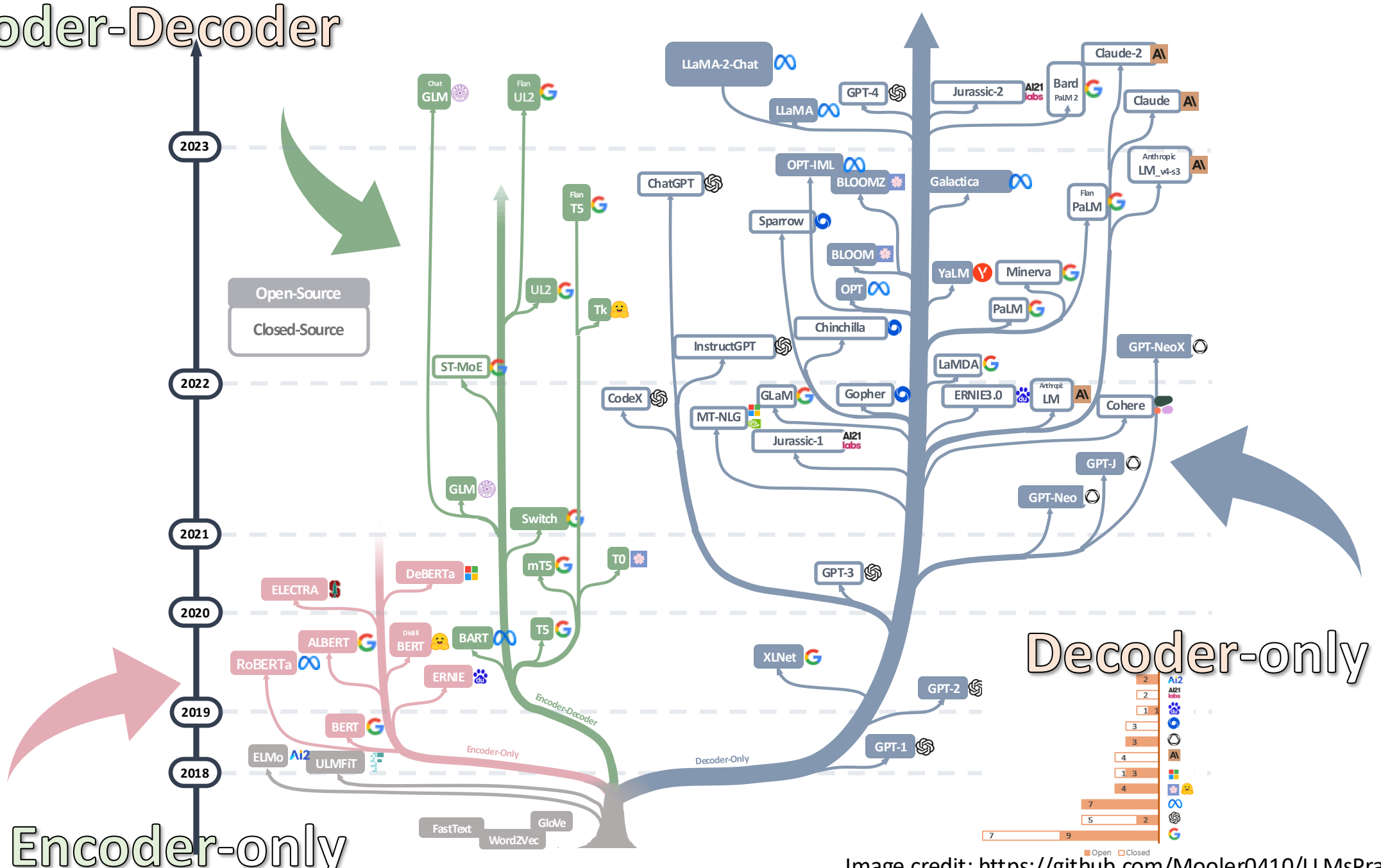
# Encoder-Decoder Transformer

# Decoder-only Transformer

$e_1$ $e_2$ $e_3$ $e_4$

Encoder #$L$

Encoder #1

<Input>

Softmax

Linear

Decoder #$L$

Decoder #1

<Target>

Softmax

Linear

Decoder #$L$

Decoder #1

<Input> <Target> <Input> ⋯

Prompting, in-context examples

**Different** parameters for encoder/decoder

**Shared** parameters

Image credit: https://github.com/Mooler0410/LLMsPracticalGuide

LLM Evolutionary Tree