CONSTELLATE

# How does ChatGPT work?

# Some clarifications

1. Conceptual questions vs.engineering questions
2. Human-engineered features vs. machine learning

Conceptual Questions
vs.
Engineering Questions

CONSTELLATE

Conceptual question:
What determines the meaning of a word?

**CONSTELLATE**

**Conceptual question:
What determines the meaning of a word?**

**The distributional hypothesis:**
words that are surrounded by similar words have similar meanings

Sentence 1: I have to make sure that the cat gets fed.

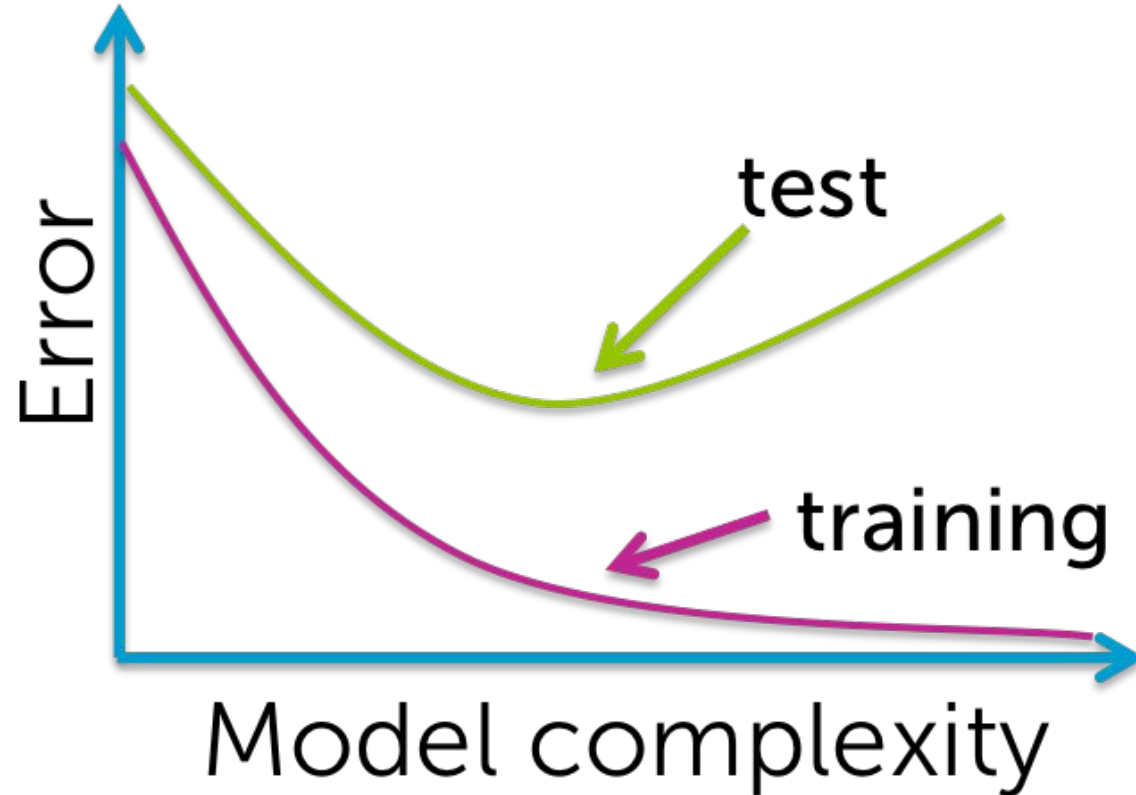Sentence 2: I forgot to make sure that the dog gets fed.

**Engineering question:
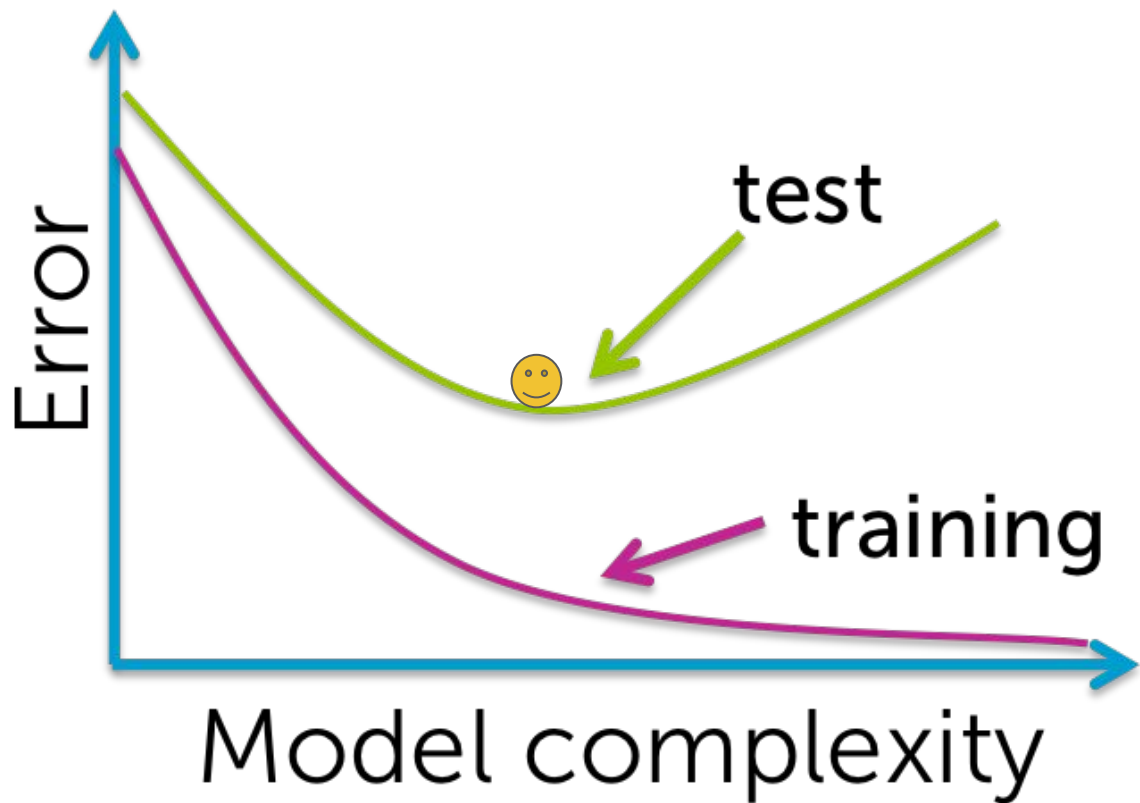Why the hidden layer has 300 neurons, not 200, not 500?**

# Human-engineered features
## vs.
## machine learning

# of bedrooms

# of bathrooms

Age

… … …

… … …

**Features**

# Skip gram

## Source Text

| The | quick | brown | fox jumps over the lazy dog. → |

**Training Samples**

(the, quick)
(the, brown)

| The | quick | brown | fox | jumps over the lazy dog. → |

(quick, the)
(quick, brown)
(quick, fox)

| The | quick | brown | fox | jumps | over the lazy dog. → |

(brown, the)
(brown, quick)
(brown, fox)
(brown, jumps)

| The | quick | brown | fox | jumps | over | the lazy dog. → |

(fox, quick)
(fox, brown)
(fox, jumps)
(fox, over)

Any questions?

**A brief review**

GPT: generative pre-trained transformer

# Transformer

a multi-layer neural network that relies on the parallel multi-head attention mechanism.

## Part 1: multi-layer neural network

**Part 2: multi-head attention mechanism**

A single self-attention head

# Learning objectives

- Concepts
  - Understand the temporal dimension of language
  - Understand some basic concepts in attention
    - query, key, value, self-attention
  - Understand a head as a representation of a certain relationship between words
- Computation
  - Understand the power of matrix multiplication for parallel computation

- **The temporal dimension of language**

- **A single self-attention head**

# The main task of LM

**Given a sequence of words, what is the most likely word that follows the sequence?**

The quick brown fox jumps over the lazy dog.

# The main task of LM

Given a sequence of words, what is the most likely word that follows the sequence?

$$w_1 w_2 w_3 w_4 \longrightarrow w_5$$

## The main task of LM

Given a sequence of words, what is the most likely word that follows the sequence?

$$w_1 w_2 w_3 w_4 \longrightarrow w_5$$

## Assumption

Every word in the sequence plays a role in determining what word follows this sequence.

$$w_1 w_2 w_3 w_4 \longrightarrow w_5$$

## Assumption

In a word sequence, every word is influenced by the past words and itself, but never by the future words.

$$w_1 w_2 w_3 w_4 w_5 w_6 w_7 w_8$$

$$w_1 w_2 w_3 w_4 w_5$$

$$W_1 W_2 W_3 W_4 \longrightarrow W_5$$

**Damian had a secret child, a girl, and had written in his will that his belongings will belong to \_\_\_\_**

# Interim summary

1. Language has a temporal dimension.

2. When making predictions about future words, we need to pay more attention to certain words in the given sequence.

# The task of ChatGPT

Given a sequence of words, generate the next word and continue this process.

# The task of ChatGPT

$$w_1 w_2 \ldots w_8 \rightarrow w_9$$

$$w_2 w_3 \ldots w_9 \rightarrow w_{10}$$ **Context length:8**

$$w_3 w_4 \ldots w_{10} \rightarrow w_{11}$$

**We are using the past words to predict the future words.**

**The question becomes:**
How do we encode the relationship between the words in a given sequence?

A word is represented by a feature vector

$$\begin{bmatrix} w_1^{f_1} & w_1^{f_2} & \ldots & w_1^{f_{32}} \end{bmatrix}$$

$$1 \times 32$$

# If we have a word sequence of len 8

$$\begin{bmatrix} w_1^{f_1} & w_1^{f_2} & \ldots & w_1^{f_{32}} \\ w_2^{f_1} & w_2^{f_2} & \ldots & w_2^{f_{32}} \\ \ldots & \ldots & \ldots & \ldots \\ w_8^{f_1} & w_8^{f_2} & \ldots & w_8^{f_{32}} \end{bmatrix}$$

$$8 \times 32$$

**How do we encode the relationship between the words in a given sequence?**

**We use attention heads!**

CONSTELLATE

**For each word token**

**head** → **query**

**head** → **key**

**head** → **value**

# Self-attention head

**For each word token**

**head**

**query** What is the word looking for?

**key** What information does the word contain?

Self-attention head

HEAD

x

q

k

**Head size**

HEAD

# neurons
= head size

**q**

# neurons
= head size

**k**

x

# Query

$$
\begin{array}{l}
t1 \\
t2 \\
\ldots \\
t8
\end{array}
\begin{bmatrix}
w_1^{f_1} & w_1^{f_2} & \ldots & w_1^{f_{32}} \\
w_2^{f_1} & w_2^{f_2} & \ldots & w_2^{f_{32}} \\
\ldots & \ldots & \ldots & \ldots \\
w_8^{f_1} & w_8^{f_2} & \ldots & w_8^{f_{32}}
\end{bmatrix}
$$

**Word sequence**

$$
\begin{bmatrix}
q_1^{t1} & q_2^{t1} & \ldots & q_{32}^{t1} \\
q_1^{t2} & q_2^{t2} & \ldots & q_{32}^{t2} \\
\ldots & \ldots & \ldots & \ldots \\
q_1^{t8} & q_2^{t8} & \ldots & q_{32}^{t8}
\end{bmatrix}
$$

**Query matrix**

# Key

$$\begin{array}{c} t1 \\ t2 \\ \ldots \\ t8 \end{array} \begin{bmatrix} w_1^{f_1} & w_1^{f_2} & \ldots & w_1^{f_{32}} \\ w_2^{f_1} & w_2^{f_2} & \ldots & w_2^{f_{32}} \\ \ldots & \ldots & \ldots & \ldots \\ w_8^{f_1} & w_8^{f_2} & \ldots & w_8^{f_{32}} \end{bmatrix}$$

**Word sequence**

$$\begin{bmatrix} k_1^{t1} & k_2^{t1} & \ldots & k_{32}^{t1} \\ k_1^{t2} & k_2^{t2} & \ldots & k_{32}^{t2} \\ \ldots & \ldots & \ldots & \ldots \\ k_1^{t8} & k_2^{t8} & \ldots & k_{32}^{t8} \end{bmatrix}$$

**Key matrix**

$$
\begin{array}{c}
t1 \\
t2 \\
\ldots \\
t8
\end{array}
\begin{bmatrix}
w_1^{f_1} & w_1^{f_2} & \ldots & w_1^{f_{32}} \\
w_2^{f_1} & w_2^{f_2} & \ldots & w_2^{f_{32}} \\
\ldots & \ldots & \ldots & \ldots \\
w_8^{f_1} & w_8^{f_2} & \ldots & w_8^{f_{32}}
\end{bmatrix}
$$

**Word sequence**

$$
\begin{bmatrix}
k_1^{t1} & k_2^{t2} & \ldots & k_8^{t8} \\
k_2^{t1} & k_2^{t2} & \ldots & k_2^{t8} \\
\ldots & \ldots & \ldots & \ldots \\
k_{32}^{t1} & k_{32}^{t2} & \ldots & k_{32}^{t8}
\end{bmatrix}
$$

**Key matrix**

# Query · Key$^\mathsf{T}$

$$\mathbf{q^{t1}} \quad \mathbf{q^{t2}} \quad \mathbf{...} \quad \mathbf{q^{t8}} \quad \begin{bmatrix} q_1^{t1} & q_2^{t1} & \cdots & q_{32}^{t1} \\ q_1^{t2} & q_2^{t2} & \cdots & q_{32}^{t2} \\ \cdots & \cdots & \cdots & \cdots \\ q_1^{t8} & q_2^{t8} & \cdots & q_{32}^{t8} \end{bmatrix}$$

$$\mathbf{k^{t1}} \quad \mathbf{k^{t2}} \quad \mathbf{...} \quad \mathbf{k^{t8}}$$

$$\cdot \quad \begin{bmatrix} k_1^{t1} & k_2^{t2} & \cdots & k_8^{t8} \\ k_2^{t1} & k_2^{t2} & \cdots & k_2^{t8} \\ \cdots & \cdots & \cdots & \cdots \\ k_{32}^{t1} & k_{32}^{t2} & \cdots & k_{32}^{t8} \end{bmatrix}$$

**Query matrix**                     **Key matrix**

$$\begin{bmatrix} \mathbf{q}^{t1} \cdot \mathbf{k}^{t1}, & \mathbf{q}^{t1} \cdot \mathbf{k}^{t2}, & \ldots, & \mathbf{q}^{t1} \cdot \mathbf{k}^{t8} \\ \mathbf{q}^{t2} \cdot \mathbf{k}^{t1}, & \mathbf{q}^{t2} \cdot \mathbf{k}^{t2}, & \ldots, & \mathbf{q}^{t2} \cdot \mathbf{k}^{t8} \\ \ldots & \ldots & \ldots & \ldots \\ \mathbf{q}^{t8} \cdot \mathbf{k}^{t1}, & \mathbf{q}^{t8} \cdot \mathbf{k}^{t2}, & \ldots, & \mathbf{q}^{t8} \cdot \mathbf{k}^{t8} \end{bmatrix}$$

CONSTELLATE

**But, what about the temporal dimension of language?**

$$\begin{bmatrix} \mathbf{q}^{t1} \cdot \mathbf{k}^{t1}, & \mathbf{q}^{t1} \cdot \mathbf{k}^{t2}, & \ldots, & \mathbf{q}^{t1} \cdot \mathbf{k}^{t8} \\ \mathbf{q}^{t2} \cdot \mathbf{k}^{t1}, & \mathbf{q}^{t2} \cdot \mathbf{k}^{t2}, & \ldots, & \mathbf{q}^{t2} \cdot \mathbf{k}^{t8} \\ \ldots & \ldots & \ldots & \ldots \\ \mathbf{q}^{t8} \cdot \mathbf{k}^{t1}, & \mathbf{q}^{t8} \cdot \mathbf{k}^{t2}, & \ldots, & \mathbf{q}^{t8} \cdot \mathbf{k}^{t8} \end{bmatrix}$$

CONSTELLATE

$$\begin{bmatrix} \mathbf{q}^{t1} \cdot \mathbf{k}^{t1}, & \mathbf{q}^{t1} \cdot \mathbf{k}^{t2}, & \ldots, & \mathbf{q}^{t1} \cdot \mathbf{k}^{t8} \\ \mathbf{q}^{t2} \cdot \mathbf{k}^{t1}, & \mathbf{q}^{t2} \cdot \mathbf{k}^{t2}, & \ldots, & \mathbf{q}^{t2} \cdot \mathbf{k}^{t8} \\ \ldots & \ldots & \ldots & \ldots \\ \mathbf{q}^{t8} \cdot \mathbf{k}^{t1}, & \mathbf{q}^{t8} \cdot \mathbf{k}^{t2}, & \ldots, & \mathbf{q}^{t8} \cdot \mathbf{k}^{t8} \end{bmatrix}$$

# Attention score

**-inf**

$$
\begin{bmatrix}
\mathbf{q}^{t1} \cdot \mathbf{k}^{t1}, & \mathbf{q}^{t1} \cdot \mathbf{k}^{t2}, & \dots, & \mathbf{q}^{t1} \cdot \mathbf{k}^{t8} \\
\mathbf{q}^{t2} \cdot \mathbf{k}^{t1}, & \mathbf{q}^{t2} \cdot \mathbf{k}^{t2}, & \dots, & \mathbf{q}^{t2} \cdot \mathbf{k}^{t8} \\
\dots & \dots & \dots & \dots \\
\mathbf{q}^{t8} \cdot \mathbf{k}^{t1}, & \mathbf{q}^{t8} \cdot \mathbf{k}^{t2}, & \dots, & \mathbf{q}^{t8} \cdot \mathbf{k}^{t8}
\end{bmatrix}
$$

# Attention score

$$
\begin{bmatrix}
\mathbf{q}^{t1} \cdot \mathbf{k}^{t1} & \text{-inf} & \ldots & \ldots & \text{-inf} \\
\mathbf{q}^{t2} \cdot \mathbf{k}^{t1} & \mathbf{q}^{t2} \cdot \mathbf{k}^{t2} & \text{-inf} & \ldots & \text{-inf} \\
\ldots & \ldots & \ldots & \ldots & \text{-inf} \\
\mathbf{q}^{t8} \cdot \mathbf{k}^{t1} & \mathbf{q}^{t8} \cdot \mathbf{k}^{t2} & \mathbf{q}^{t8} \cdot \mathbf{k}^{t3} & \ldots & \mathbf{q}^{t8} \cdot \mathbf{k}^{t8}
\end{bmatrix}
$$

# Attention scores to weights

$$\begin{bmatrix} \mathbf{q}^{t1} \cdot \mathbf{k}^{t1} & \text{-inf} & \ldots & \ldots & \text{-inf} \\ \mathbf{q}^{t2} \cdot \mathbf{k}^{t1} & \mathbf{q}^{t2} \cdot \mathbf{k}^{t2} & \text{-inf} & \ldots & \text{-inf} \\ \ldots & \ldots & \ldots & \ldots & \text{-inf} \\ \mathbf{q}^{t8} \cdot \mathbf{k}^{t1} & \mathbf{q}^{t8} \cdot \mathbf{k}^{t2} & \mathbf{q}^{t8} \cdot \mathbf{k}^{t3} & \ldots & \mathbf{q}^{t8} \cdot \mathbf{k}^{t8} \end{bmatrix}$$

**softmax**

**softmax**

**softmax**

**softmax**

# Attention scores to weights

$$
\begin{bmatrix}
P_{\mathbf{q}^{t1} \cdot \mathbf{k}^{t1}} & 0 & \ldots & \ldots & 0 \\
P_{\mathbf{q}^{t2} \cdot \mathbf{k}^{t1}} & P_{\mathbf{q}^{t2} \cdot \mathbf{k}^{t2}} & 0 & \ldots & 0 \\
\ldots & \ldots & \ldots & \ldots & 0 \\
P_{\mathbf{q}^{t8} \cdot \mathbf{k}^{t1}} & P_{\mathbf{q}^{t8} \cdot \mathbf{k}^{t2}} & P_{\mathbf{q}^{t8} \cdot \mathbf{k}^{t3}} & \ldots & P_{\mathbf{q}^{t8} \cdot \mathbf{k}^{t8}}
\end{bmatrix}
$$

sum up to 1

sum up to 1

sum up to 1

sum up to 1

# Attention scores to weights

$$\begin{bmatrix} 1 & 0 & 0 & \ldots & 0 \\ 0.22 & 0.78 & 0 & \ldots & 0 \\ \ldots & \ldots & \ldots & \ldots & \ldots \\ 0.02 & 0.12 & 0.24 & \ldots & 0.45 \end{bmatrix}$$

sum up to 1

sum up to 1

sum up to 1

sum up to 1

**weight matrix**

A visual of what we have learned so far

HEAD

x

q

k$^T$

A visual of what we have learned so far

HEAD

x

q

•

k$^T$

# A visual of what we have learned so far



HEAD

x

-inf mask

q

•

k$^T$

Softmax to each row

A visual of what we have learned so far

CONSTELLATE

**For each word token**

**head** → **query**

**head** → **key**

**head** → **value**

$$\begin{bmatrix} v_1^{t1} & v_2^{t1} & \dots & v_8^{t1} \\ v_1^{t2} & v_2^{t2} & \dots & v_8^{t2} \\ \dots & \dots & \dots & \dots \\ v_1^{t8} & v_2^{t8} & \dots & v_8^{t8} \end{bmatrix}$$

$$
\begin{bmatrix}
P_{\mathbf{q}^{t1} \cdot \mathbf{k}^{t1}} & 0 & \dots & \dots & 0 \\
P_{\mathbf{q}^{t2} \cdot \mathbf{k}^{t1}} & P_{\mathbf{q}^{t2} \cdot \mathbf{k}^{t2}} & 0 & \dots & 0 \\
\dots & \dots & \dots & \dots & 0 \\
P_{\mathbf{q}^{t8} \cdot \mathbf{k}^{t1}} & P_{\mathbf{q}^{t8} \cdot \mathbf{k}^{t2}} & P_{\mathbf{q}^{t8} \cdot \mathbf{k}^{t3}} & \dots & P_{\mathbf{q}^{t8} \cdot \mathbf{k}^{t8}}
\end{bmatrix}
\bullet
\begin{bmatrix}
v_1^{t1} & v_2^{t1} & \dots & v_8^{t1} \\
v_1^{t2} & v_2^{t2} & \dots & v_8^{t2} \\
\dots & \dots & \dots & \dots \\
v_1^{t8} & v_2^{t8} & \dots & v_8^{t8}
\end{bmatrix}
$$

$$\begin{bmatrix} P_{\mathbf{q}^{t1}\cdot\mathbf{k}^{t1}} & 0 & \ldots & \ldots & 0 \\ P_{\mathbf{q}^{t2}\cdot\mathbf{k}^{t1}} & P_{\mathbf{q}^{t2}\cdot\mathbf{k}^{t2}} & 0 & \ldots & 0 \\ \ldots & \ldots & \ldots & \ldots & 0 \\ P_{\mathbf{q}^{t8}\cdot\mathbf{k}^{t1}} & P_{\mathbf{q}^{t8}\cdot\mathbf{k}^{t2}} & P_{\mathbf{q}^{t8}\cdot\mathbf{k}^{t3}} & \ldots & P_{\mathbf{q}^{t8}\cdot\mathbf{k}^{t8}} \end{bmatrix} \cdot \begin{bmatrix} v_1^{t1} & v_2^{t1} & \ldots & v_8^{t1} \\ v_1^{t2} & v_2^{t2} & \ldots & v_8^{t2} \\ \ldots & \ldots & \ldots & \ldots \\ v_1^{t8} & v_2^{t8} & \ldots & v_8^{t8} \end{bmatrix}$$

$$\begin{bmatrix} P_{\mathbf{q}^{t1}\cdot\mathbf{k}^{t1}} & 0 & \ldots & \ldots & 0 \\ P_{\mathbf{q}^{t2}\cdot\mathbf{k}^{t1}} & P_{\mathbf{q}^{t2}\cdot\mathbf{k}^{t2}} & 0 & \ldots & 0 \\ \ldots & \ldots & \ldots & \ldots & 0 \\ P_{\mathbf{q}^{t8}\cdot\mathbf{k}^{t1}} & P_{\mathbf{q}^{t8}\cdot\mathbf{k}^{t2}} & P_{\mathbf{q}^{t8}\cdot\mathbf{k}^{t3}} & \ldots & P_{\mathbf{q}^{t8}\cdot\mathbf{k}^{t8}} \end{bmatrix} \cdot \begin{bmatrix} v_1^{t1} & v_2^{t1} & \ldots & v_8^{t1} \\ v_1^{t2} & v_2^{t2} & \ldots & v_8^{t2} \\ \ldots & \ldots & \ldots & \ldots \\ v_1^{t8} & v_2^{t8} & \ldots & v_8^{t8} \end{bmatrix}$$

$$= \begin{bmatrix} P_{\mathbf{q}^{t2}\cdot\mathbf{k}^{t1}} \cdot v_1^{t1} + P_{\mathbf{q}^{t2}\cdot\mathbf{k}^{t2}} \cdot v_1^{t2} & P_{\mathbf{q}^{t2}\cdot\mathbf{k}^{t1}} \cdot v_2^{t1} + P_{\mathbf{q}^{t2}\cdot\mathbf{k}^{t2}} \cdot v_2^{t2} & \ldots & P_{\mathbf{q}^{t2}\cdot\mathbf{k}^{t1}} \cdot v_8^{t1} + P_{\mathbf{q}^{t2}\cdot\mathbf{k}^{t2}} \cdot v_8^{t2} \end{bmatrix}$$

# References

Jurafsky, Daniel, and James H. Martin. (2023). *Speech and language processing: An introduction to natural language processing, computational linguistics, and speech recognition.*

Karpathy, Andrej. (2023), GitHub repository, https://github.com/karpathy/nanoGPT

Karpathy Andrej. (2023). Let's build GPT: from scratch, in code, spelled out. [Andrej Karpathy]. YouTube. Retrieved September 5, 2023 from https://www.youtube.com/watch?v=kCc8FmEb1nY