# How does ChatGPT work?

CONSTELLATE

# 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
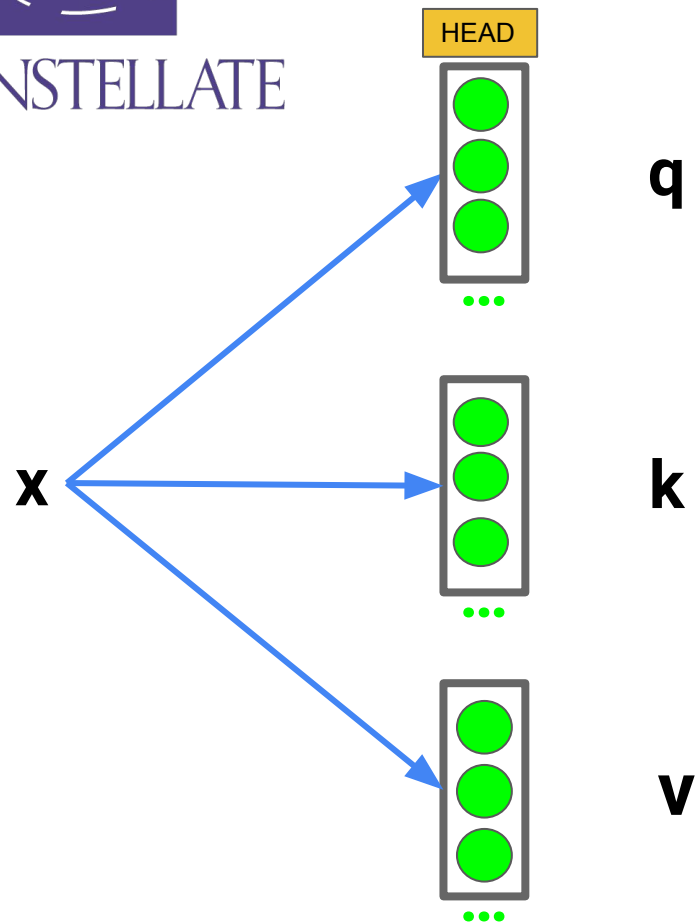  - Understand the power of matrix multiplication for parallel computation

A single self-attention head

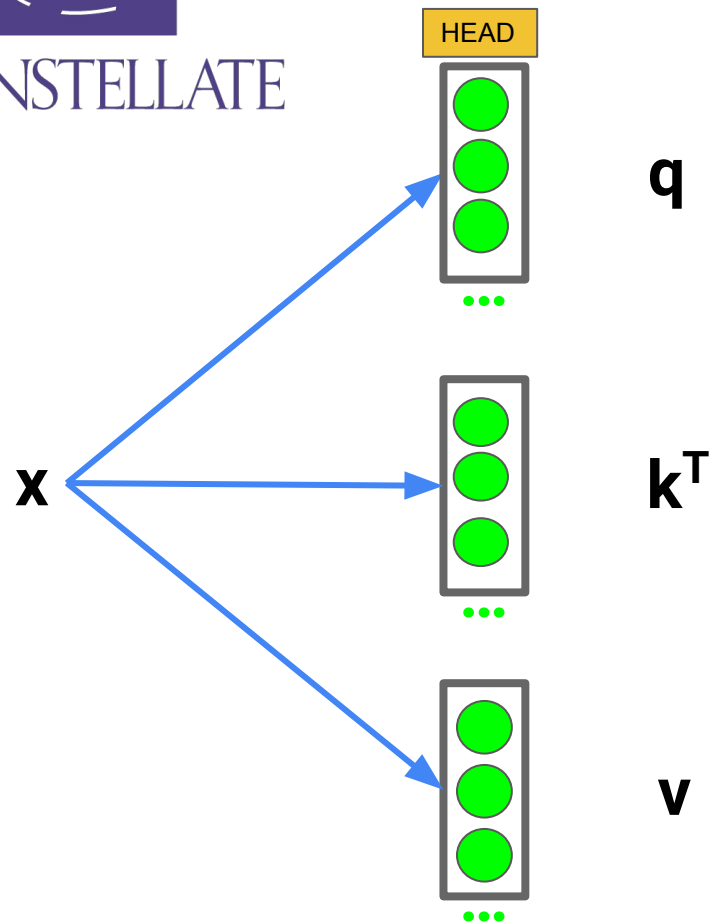# A visual of what we have learned so far

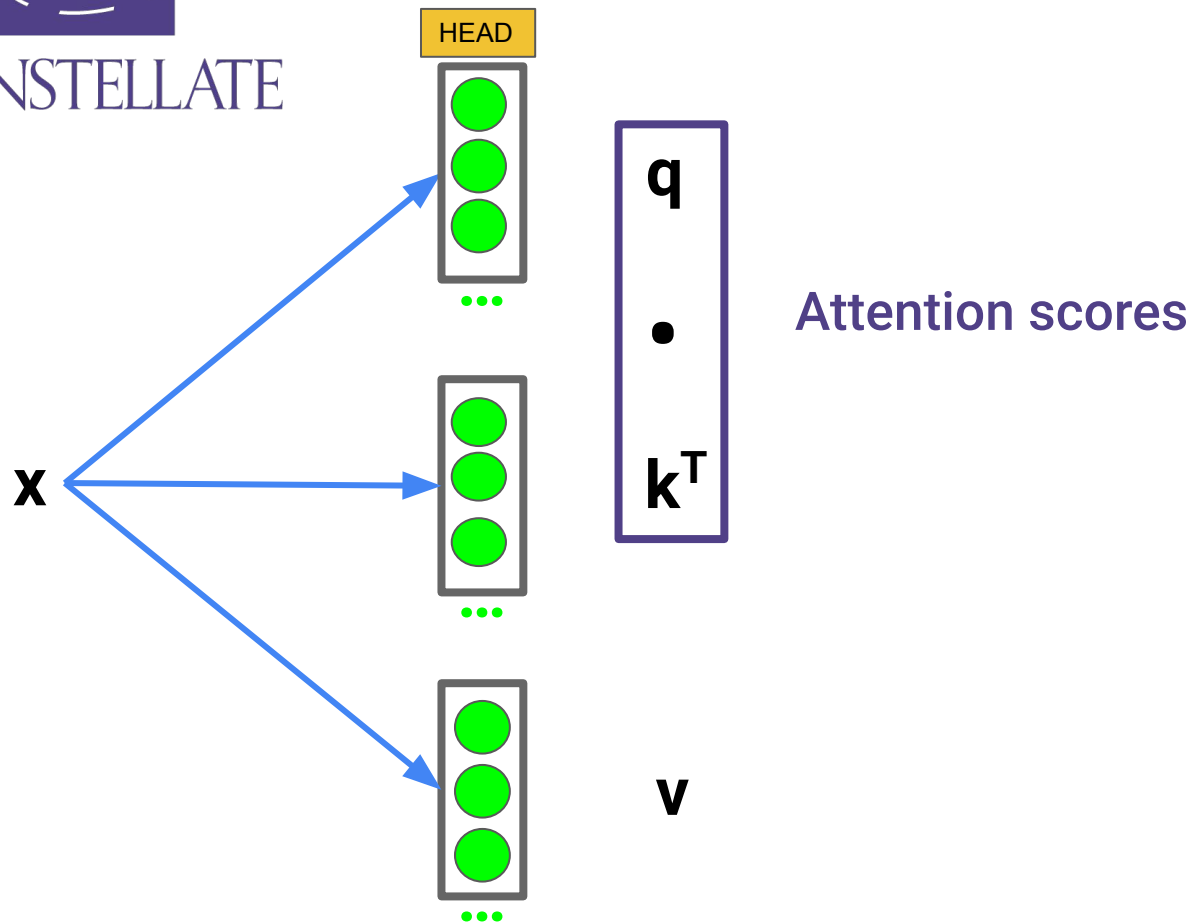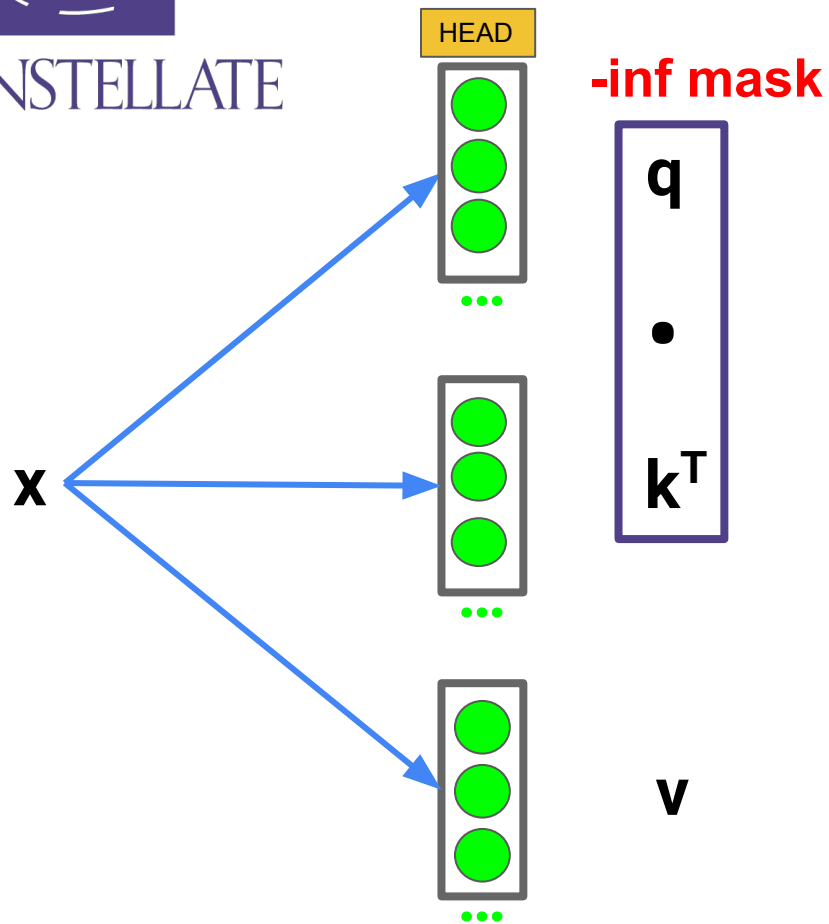A visual of what we have learned so far

**A visual of what we have learned so far**

HEAD

**-inf mask**

$q$

$\bullet$

$k^T$

$v$

A visual of what we have learned so far

A visual of what we have learned so far

Any questions?

# Multi-head attention mechanism

A visual of what we have learned so far

CONSTELLATE

Head size= 8

$$\begin{bmatrix} x_1^{t1} & x_2^{t1} & \dots & x_{32}^{t1} \\ x_1^{t2} & x_2^{t2} & \dots & x_{32}^{t2} \\ \dots & \dots & \dots & \dots \\ x_1^{t8} & x_2^{t8} & \dots & x_{32}^{t8} \end{bmatrix}$$

32 features

HEAD

-inf mask

$\mathbf{q}$

$\cdot$

$\mathbf{k^T}$

Softmax to each row

$\mathbf{w}$

$\cdot$

$\mathbf{v}$

# A single head



$$
\mathbf{q} \quad \text{8X8} \quad
\begin{bmatrix}
q_1^{t1} & q_2^{t1} & \cdots & q_8^{t1} \\
q_1^{t2} & q_2^{t2} & \cdots & q_8^{t2} \\
\cdots & \cdots & \cdots & \cdots \\
q_1^{t8} & q_2^{t8} & \cdots & q_8^{t8}
\end{bmatrix}
$$

$$
\mathbf{k}^{\mathsf{T}} \quad \text{8X8} \quad
\begin{bmatrix}
k_1^{t1} & k_1^{t2} & \cdots & k_1^{t8} \\
k_2^{t1} & k_2^{t2} & \cdots & k_2^{t8} \\
\cdots & \cdots & \cdots & \cdots \\
k_8^{t1} & k_8^{t2} & \cdots & k_8^{t8}
\end{bmatrix}
$$

$$
\mathbf{v} \quad \text{8X8} \quad
\begin{bmatrix}
v_1^{t1} & v_2^{t1} & \cdots & v_8^{t1} \\
v_1^{t2} & v_2^{t2} & \cdots & v_8^{t2} \\
\cdots & \cdots & \cdots & \cdots \\
v_1^{t8} & v_2^{t8} & \cdots & v_8^{t8}
\end{bmatrix}
$$

$\mathbf{x}$

8X32

CONSTELLATE

q  8X8

•

k$^T$  8X8

w

x

8X32

v  8X8

# A single head



q    8X8

k$^T$    8X8

w    8X8

x    8X32

v    8X8

$$\begin{bmatrix} a_{v1}^{t1} & a_{v2}^{t1} & \dots & a_{v8}^{t1} \\ a_{v1}^{t2} & a_{v2}^{t2} & \dots & a_{v8}^{t2} \\ \dots & \dots & \dots & \dots \\ a_{v1}^{t8} & a_{v2}^{t8} & \dots & a_{v8}^{t8} \end{bmatrix}$$

8X8

A single head

CONSTELLATE

q    8X8

k$^T$    8X8

w    8X8

v    8X8

x

8X32

a$^{h1}$    8X8

A single head

x
8X32

q 8X8

·

k$^{T}$ 8X8

w 8X8

·

v 8X8

a$^{h1}$ 8X8

**x**

**HEAD 1**

$a^h_1$

8X32

8X8

# Multi-head

# Multi-head

$$\begin{bmatrix} a^{h1}_{t1_{v1}} & a^{h1}_{t1_{v2}} & \ldots & a^{h1}_{t1_{v8}} \\ a^{h1}_{t2_{v1}} & a^{h1}_{t2_{v2}} & \ldots & a^{h1}_{t2_{v8}} \\ \ldots & \ldots & \ldots & \ldots \\ a^{h1}_{t8_{v1}} & a^{h1}_{t8_{v2}} & \ldots & a^{h1}_{t8_{v8}} \end{bmatrix}$$

8X8

$$\begin{bmatrix} a^{h2}_{t1_{v1}} & a^{h2}_{t1_{v2}} & \ldots & a^{h2}_{t1_{v8}} \\ a^{h2}_{t2_{v1}} & a^{h2}_{t2_{v2}} & \ldots & a^{h2}_{t2_{v8}} \\ \ldots & \ldots & \ldots & \ldots \\ a^{h2}_{t8_{v1}} & a^{h2}_{t8_{v2}} & \ldots & a^{h2}_{t8_{v8}} \end{bmatrix}$$

8X8

$$\begin{bmatrix} a_{t1_{v1}}^{h1} & a_{t1_{v2}}^{h1} & \ldots & a_{t1_{v8}}^{h1} & a_{t1_{v1}}^{h2} & a_{t1_{v2}}^{h2} & \ldots & a_{t1_{v8}}^{h2} \\ a_{t2_{v1}}^{h1} & a_{t2_{v2}}^{h1} & \ldots & a_{t2_{v8}}^{h1} & a_{t2_{v1}}^{h2} & a_{t2_{v2}}^{h2} & \ldots & a_{t2_{v8}}^{h2} \\ \ldots & \ldots & \ldots & \ldots & \ldots & \ldots & \ldots & \ldots \\ a_{t8_{v1}}^{h1} & a_{t8_{v2}}^{h1} & \ldots & a_{t8_{v8}}^{h1} & a_{t8_{v1}}^{h2} & a_{t8_{v2}}^{h2} & \ldots & a_{t8_{v8}}^{h2} \end{bmatrix}$$

8X16

8X32

$$\begin{bmatrix} \begin{array}{ccc} a_{t1_{v1}}^{h1} & \cdots & a_{t1_{v8}}^{h1} \\ a_{t2_{v1}}^{h1} & \cdots & a_{t2_{v8}}^{h1} \\ \cdots & \cdots & \cdots \\ a_{t8_{v1}}^{h1} & \cdots & a_{t8_{v8}}^{h1} \end{array} \quad \begin{array}{ccc} a_{t1_{v1}}^{h2} & \cdots & a_{t1_{v8}}^{h2} \\ a_{t2_{v1}}^{h2} & \cdots & a_{t2_{v8}}^{h2} \\ \cdots & \cdots & \cdots \\ a_{t8_{v1}}^{h2} & \cdots & a_{t8_{v8}}^{h2} \end{array} \quad \begin{array}{ccc} a_{t1_{v1}}^{h3} & \cdots & a_{t1_{v8}}^{h3} \\ a_{t2_{v1}}^{h3} & \cdots & a_{t2_{v8}}^{h3} \\ \cdots & \cdots & \cdots \\ a_{t8_{v1}}^{h3} & \cdots & a_{t8_{v8}}^{h3} \end{array} \quad \begin{array}{ccc} a_{t1_{v1}}^{h4} & \cdots & a_{t1_{v8}}^{h4} \\ a_{t2_{v1}}^{h4} & \cdots & a_{t2_{v8}}^{h4} \\ \cdots & \cdots & \cdots \\ a_{t8_{v1}}^{h4} & \cdots & a_{t8_{v8}}^{h4} \end{array} \end{bmatrix}$$

**HEAD 1**  **HEAD 2**  **HEAD 3**  **HEAD 4**

Still remember the shape of the original input?

# A visual of what we have learned so far



CONSTELLATE

$$\begin{bmatrix} x_1^{t1} & x_2^{t1} & \ldots & x_{32}^{t1} \\ x_1^{t2} & x_2^{t2} & \ldots & x_{32}^{t2} \\ \ldots & \ldots & \ldots & \ldots \\ x_1^{t8} & x_2^{t8} & \ldots & x_{32}^{t8} \end{bmatrix}$$

8X32

Head size= 8

HEAD

-inf mask

$\mathbf{q}$

$\cdot$

$\mathbf{k^T}$

Softmax to each row

$\mathbf{w}$

$\cdot$

$\mathbf{v}$

8X32

$$\left[\begin{array}{ccc} a^{h1}_{t1_{v1}} & \cdots & a^{h1}_{t1_{v8}} \\ a^{h1}_{t2_{v1}} & \cdots & a^{h1}_{t2_{v8}} \\ \cdots & \cdots & \cdots \\ a^{h1}_{t8_{v1}} & \cdots & a^{h1}_{t8_{v8}} \end{array}\right. \quad \begin{array}{ccc} a^{h2}_{t1_{v1}} & \cdots & a^{h2}_{t1_{v8}} \\ a^{h2}_{t2_{v1}} & \cdots & a^{h2}_{t2_{v8}} \\ \cdots & \cdots & \cdots \\ a^{h2}_{t8_{v1}} & \cdots & a^{h2}_{t8_{v8}} \end{array} \quad \begin{array}{ccc} a^{h3}_{t1_{v1}} & \cdots & a^{h3}_{t1_{v8}} \\ a^{h3}_{t2_{v1}} & \cdots & a^{h3}_{t2_{v8}} \\ \cdots & \cdots & \cdots \\ a^{h3}_{t8_{v1}} & \cdots & a^{h3}_{t8_{v8}} \end{array} \quad \left.\begin{array}{ccc} a^{h4}_{t1_{v1}} & \cdots & a^{h4}_{t1_{v8}} \\ a^{h4}_{t2_{v1}} & \cdots & a^{h4}_{t2_{v8}} \\ \cdots & \cdots & \cdots \\ a^{h4}_{t8_{v1}} & \cdots & a^{h4}_{t8_{v8}} \end{array}\right]$$

**HEAD 1**      **HEAD 2**      **HEAD 3**      **HEAD 4**

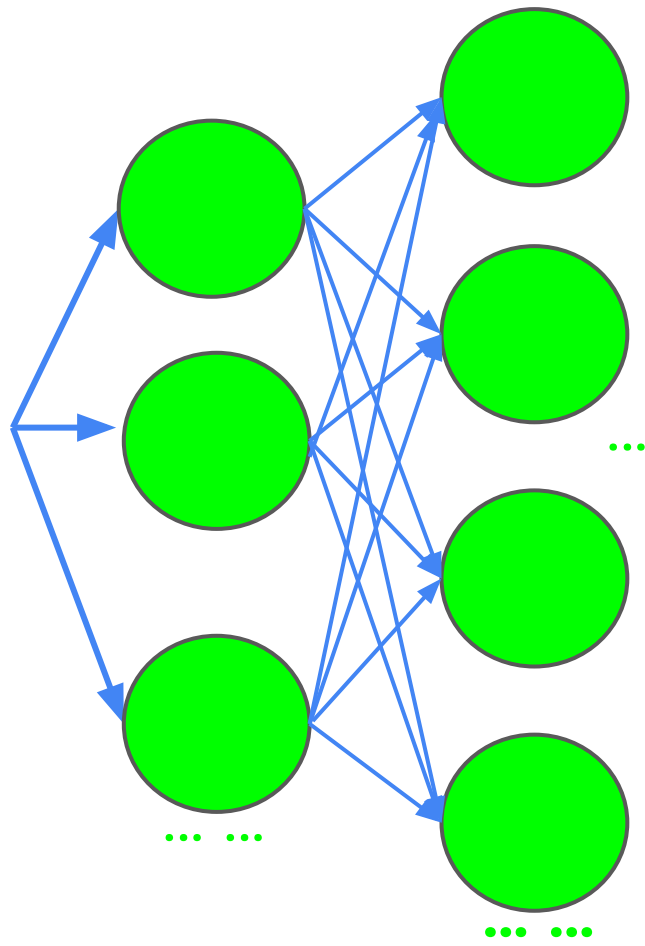**Feedforward**

$$\begin{bmatrix} a_{t1_{v1}}^{h1} & \dots & a_{t1_{v8}}^{h1} & a_{t1_{v1}}^{h2} & \dots & a_{t1_{v8}}^{h2} & a_{t1_{v1}}^{h3} & \dots & a_{t1_{v8}}^{h3} & a_{t1_{v1}}^{h4} & \dots & a_{t1_{v8}}^{h4} \\ a_{t2_{v1}}^{h1} & \dots & a_{t2_{v8}}^{h1} & a_{t2_{v1}}^{h2} & \dots & a_{t2_{v8}}^{h2} & a_{t2_{v1}}^{h3} & \dots & a_{t2_{v8}}^{h3} & a_{t2_{v1}}^{h4} & \dots & a_{t2_{v8}}^{h4} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ a_{t8_{v1}}^{h1} & \dots & a_{t8_{v8}}^{h1} & a_{t8_{v1}}^{h2} & \dots & a_{t8_{v8}}^{h2} & a_{t8_{v1}}^{h3} & \dots & a_{t8_{v8}}^{h3} & a_{t8_{v1}}^{h4} & \dots & a_{t8_{v8}}^{h4} \end{bmatrix}$$
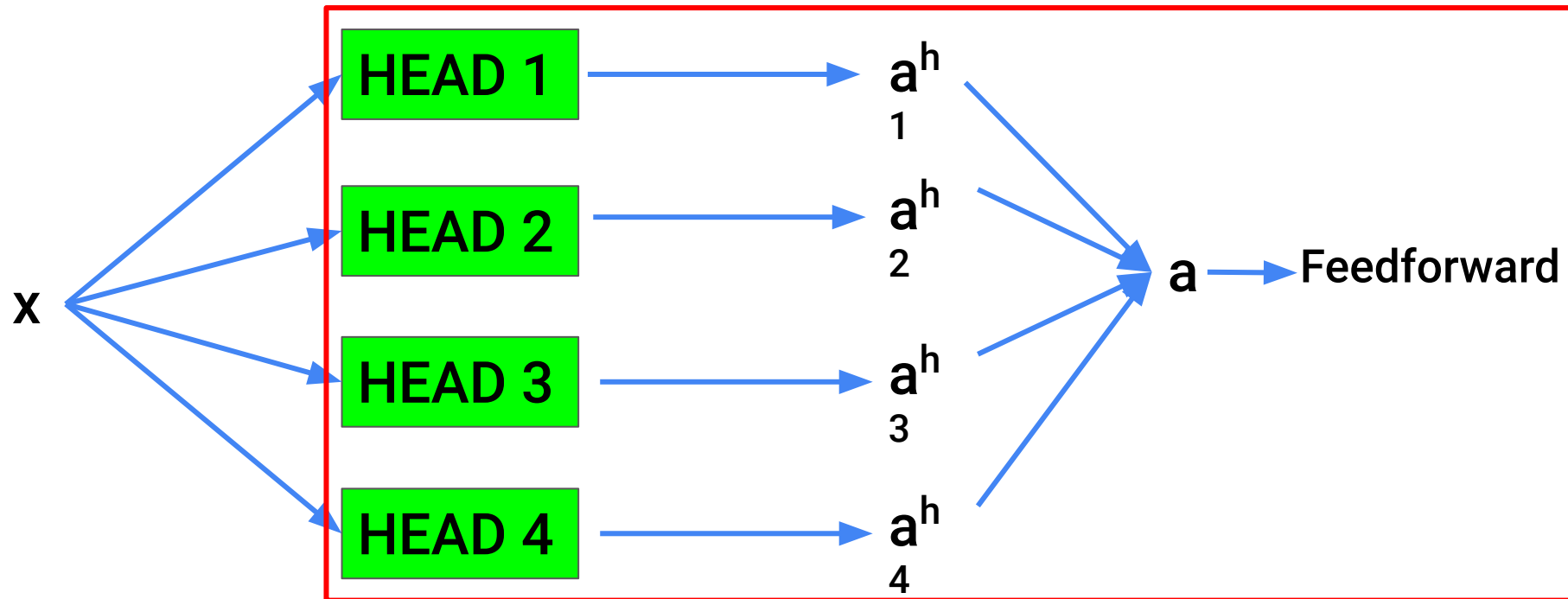
8X32

# Attention block

$x \rightarrow$ Block 1 $\rightarrow a^{b1} \rightarrow$ Block 2 $\rightarrow a^{b2} \rightarrow \ldots\ldots \rightarrow a^{b96}$

Any questions?

## Motivation

**The order of the words in a sequence matters.**

**Even though she did not win the award, she was satisfied.**

**Even though she did win the award, she was not satisfied.**

$$\begin{bmatrix} x_1^{t1} & x_2^{t1} & \dots & x_{32}^{t1} \\ x_1^{t2} & x_2^{t2} & \dots & x_{32}^{t2} \\ \dots & \dots & \dots & \dots \\ x_1^{t8} & x_2^{t8} & \dots & x_{32}^{t8} \end{bmatrix} + \begin{bmatrix} p_1^{t1} & p_2^{t1} & \dots & p_{32}^{t1} \\ p_1^{t2} & p_2^{t2} & \dots & p_{32}^{t2} \\ \dots & \dots & \dots & \dots \\ p_1^{t8} & p_2^{t8} & \dots & p_{32}^{t8} \end{bmatrix}$$

CONSTELLATE

$$\begin{bmatrix} x_1^{t1} & x_2^{t1} & \dots & x_{32}^{t1} \\ x_1^{t2} & x_2^{t2} & \dots & x_{32}^{t2} \\ \dots & \dots & \dots & \dots \\ x_1^{t8} & x_2^{t8} & \dots & x_{32}^{t8} \end{bmatrix} + \begin{bmatrix} p_1^{t1} & p_2^{t1} & \dots & p_{32}^{t1} \\ p_1^{t2} & p_2^{t2} & \dots & p_{32}^{t2} \\ \dots & \dots & \dots & \dots \\ p_1^{t8} & p_2^{t8} & \dots & p_{32}^{t8} \end{bmatrix}$$

$$\begin{bmatrix} x_1^{t1} & x_2^{t1} & \ldots & x_{32}^{t1} \\ x_1^{t2} & x_2^{t2} & \ldots & x_{32}^{t2} \\ \ldots & \ldots & \ldots & \ldots \\ x_1^{t8} & x_2^{t8} & \ldots & x_{32}^{t8} \end{bmatrix} + \begin{bmatrix} 0 & 0 & \ldots & 0 \\ 1 & 1 & \ldots & 1 \\ \ldots & \ldots & \ldots & \ldots \\ 7 & 7 & \ldots & 7 \end{bmatrix}$$

# Positional encoding

$$
\begin{array}{cccc}
\text{i=0} & \text{i=1} & \text{......} & \text{i=32}
\end{array}
$$

$$
\begin{bmatrix} x_1^{t1} & x_2^{t1} & \ldots & x_{32}^{t1} \end{bmatrix}
$$

$$
PE_{(pos,\, 2i)} = sin\left(\frac{pos}{10000^{\frac{2i}{d}}}\right) \qquad \text{d=32}
$$

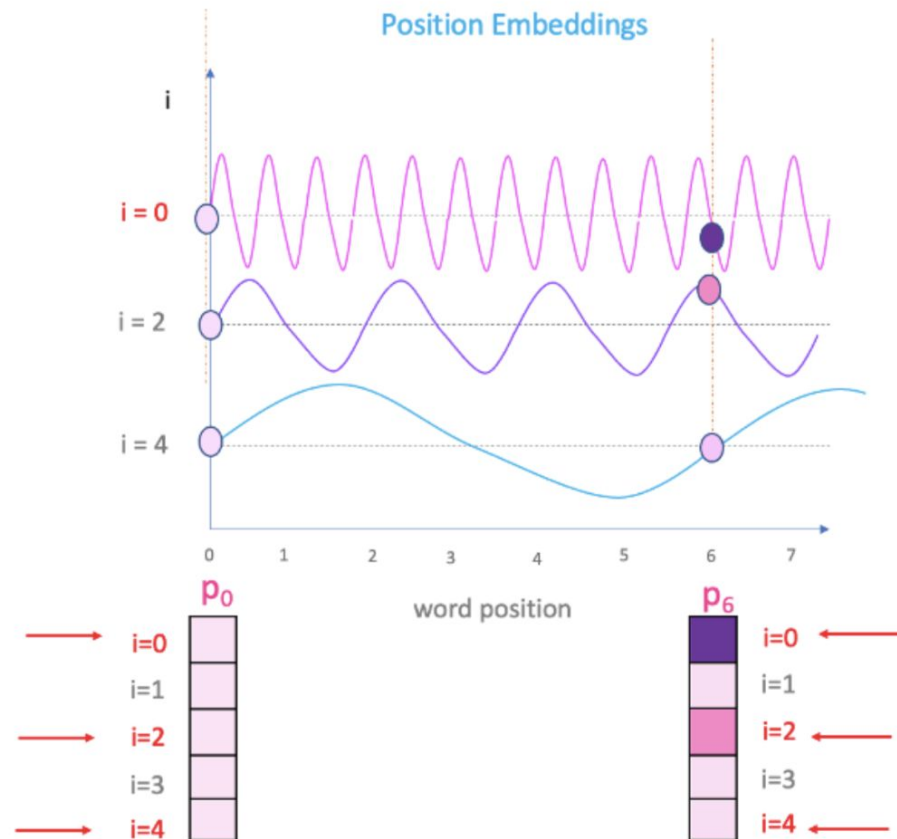$$
PE_{(pos,\, 2i+1)} = cos\left(\frac{pos}{10000^{\frac{2i}{d}}}\right)
$$

# Positional encoding

$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{\frac{2i}{d}}}\right)$$



Position Embeddings

Even though she did **not** win the award, she was satisfied.

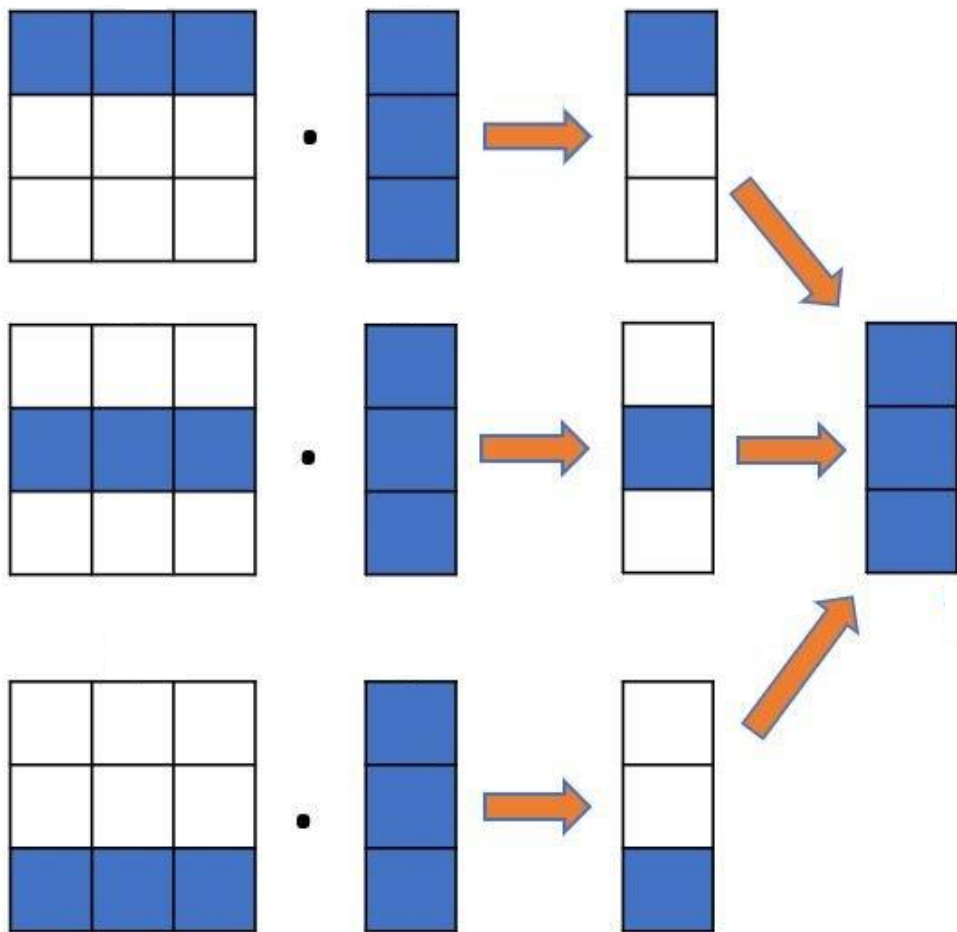Even though she did win the award, she was **not** satisfied.

Any questions?

# Calculate all observations at one time: Matrix multiplication

$$\mathbf{x}^{(1)} \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \ldots & x_n^{(1)} \\ x_1^{(2)} & x_2^{(2)} & \ldots & x_n^{(2)} \\ \ldots & \ldots & \ldots & \ldots \\ x_1^{(m)} & x_2^{(m)} & \ldots & x_n^{(m)} \end{bmatrix} \bullet \begin{matrix} \mathbf{w^c_1} & \mathbf{w^c_2} & & \mathbf{w^{cK}} \\ \end{matrix} \begin{bmatrix} w_1^{c1} & w_1^{c2} & \ldots & w_1^{cK} \\ w_2^{c1} & w_2^{c2} & \ldots & w_2^{cK} \\ \ldots & \ldots & \ldots & \ldots \\ w_n^{c1} & w_n^{c2} & \ldots & w_n^{cK} \end{bmatrix} + \begin{bmatrix} b^{c1} & b^{c2} & \ldots & b^{cK} \end{bmatrix}$$

$$= \begin{bmatrix} \mathbf{w}^{c1}\mathbf{x}^{(1)} + b^{c1} & \mathbf{w}^{c2}\mathbf{x}^{(1)} + b^{c2} & \ldots & \mathbf{w}^{cK}\mathbf{x}^{(1)} + b^{cK} \\ \mathbf{w}^{c1}\mathbf{x}^{(2)} + b^{c1} & \mathbf{w}^{c2}\mathbf{x}^{(2)} + b^{c2} & \ldots & \mathbf{w}^{cK}\mathbf{x}^{(2)} + b^{cK} \\ \ldots & \ldots & \ldots & \ldots \\ \mathbf{w}^{c1}\mathbf{x}^{(m)} + b^{c1} & \mathbf{w}^{c2}\mathbf{x}^{(m)} + b^{c2} & \ldots & \mathbf{w}^{cK}\mathbf{x}^{(m)} + b^{cK} \end{bmatrix}$$

# Calculate all observations at one time: Matrix multiplication

$$
\begin{array}{c}
\mathbf{x}^{(1)} \\
\mathbf{x}^{(2)} \\
\\
\mathbf{x}^{(m)}
\end{array}
\begin{bmatrix}
x_1^{(1)} & x_2^{(1)} & \dots & x_n^{(1)} \\
x_1^{(2)} & x_2^{(2)} & \dots & x_n^{(2)} \\
\dots & \dots & \dots & \dots \\
x_1^{(m)} & x_2^{(m)} & \dots & x_n^{(m)}
\end{bmatrix}
\bullet
\begin{matrix}
\mathbf{w^c} & \mathbf{w^c} & & \mathbf{w^{cK}} \\
\mathbf{1} & \mathbf{2} & &
\end{matrix}
\begin{bmatrix}
w_1^{c1} & w_1^{c2} & \dots & w_1^{cK} \\
w_2^{c1} & w_2^{c2} & \dots & w_2^{cK} \\
\dots & \dots & \dots & \dots \\
w_n^{c1} & w_n^{c2} & \dots & w_n^{cK}
\end{bmatrix}
+
\begin{bmatrix}
b^{c1} & b^{c2} & \dots & b^{cK}
\end{bmatrix}
$$

$$
=
\begin{bmatrix}
\mathbf{w}^{c1}\mathbf{x}^{(1)} + b^{c1} & \mathbf{w}^{c2}\mathbf{x}^{(1)} + b^{c2} & \dots & \mathbf{w}^{cK}\mathbf{x}^{(1)} + b^{cK} \\
\mathbf{w}^{c1}\mathbf{x}^{(2)} + b^{c1} & \mathbf{w}^{c2}\mathbf{x}^{(2)} + b^{c2} & \dots & \mathbf{w}^{cK}\mathbf{x}^{(2)} + b^{cK} \\
\dots & \dots & \dots & \dots \\
\mathbf{w}^{c1}\mathbf{x}^{(m)} + b^{c1} & \mathbf{w}^{c2}\mathbf{x}^{(m)} + b^{c2} & \dots & \mathbf{w}^{cK}\mathbf{x}^{(m)} + b^{cK}
\end{bmatrix}
$$

# Calculate all observations at one time: Matrix multiplication

$$
\begin{array}{c}
\mathbf{x}^{(1)} \\
\mathbf{x}^{(2)} \\
\\
\mathbf{x}^{(m)}
\end{array}
\begin{bmatrix}
x_1^{(1)} & x_2^{(1)} & \ldots & x_n^{(1)} \\
x_1^{(2)} & x_2^{(2)} & \ldots & x_n^{(2)} \\
\ldots & \ldots & \ldots & \ldots \\
x_1^{(m)} & x_2^{(m)} & \ldots & x_n^{(m)}
\end{bmatrix}
\bullet
\begin{matrix}
\mathbf{w^c} & \mathbf{w^c} & & \mathbf{w^{cK}}
\end{matrix}
\begin{bmatrix}
\mathbf{1} w_1^{c1} & \mathbf{2} w_1^{c2} & \ldots & w_1^{cK} \\
w_2^{c1} & w_2^{c2} & \ldots & w_2^{cK} \\
\ldots & \ldots & \ldots & \ldots \\
w_n^{c1} & w_n^{c2} & \ldots & w_n^{cK}
\end{bmatrix}
+
\begin{bmatrix}
b^{c1} & b^{c2} & \ldots & b^{cK}
\end{bmatrix}
$$

$$
=
\begin{bmatrix}
\mathbf{w}^{c1}\mathbf{x}^{(1)} + b^{c1} & \mathbf{w}^{c2}\mathbf{x}^{(1)} + b^{c2} & \ldots & \mathbf{w}^{cK}\mathbf{x}^{(1)} + b^{cK} \\
\mathbf{w}^{c1}\mathbf{x}^{(2)} + b^{c1} & \mathbf{w}^{c2}\mathbf{x}^{(2)} + b^{c2} & \ldots & \mathbf{w}^{cK}\mathbf{x}^{(2)} + b^{cK} \\
\ldots & \ldots & \ldots & \ldots \\
\mathbf{w}^{c1}\mathbf{x}^{(m)} + b^{c1} & \mathbf{w}^{c2}\mathbf{x}^{(m)} + b^{c2} & \ldots & \mathbf{w}^{cK}\mathbf{x}^{(m)} + b^{cK}
\end{bmatrix}
$$

# Calculate all observations at one time: Matrix multiplication

$$
\begin{array}{c} \mathbf{x}^{(1)} \\ \mathbf{x}^{(2)} \\ \\ \mathbf{x}^{(m)} \end{array}
\begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \ldots & x_n^{(1)} \\ x_1^{(2)} & x_2^{(2)} & \ldots & x_n^{(2)} \\ \ldots & \ldots & \ldots & \ldots \\ x_1^{(m)} & x_2^{(m)} & \ldots & x_n^{(m)} \end{bmatrix}
\bullet
\begin{array}{ccc} \mathbf{w^c} & \mathbf{w^c} & \mathbf{w^{cK}} \\ \mathbf{1} & \mathbf{2} & \end{array}
\begin{bmatrix} w_1^{c1} & w_1^{c2} & \ldots & w_1^{cK} \\ w_2^{c1} & w_2^{c2} & \ldots & w_2^{cK} \\ \ldots & \ldots & \ldots & \ldots \\ w_n^{c1} & w_n^{c2} & \ldots & w_n^{cK} \end{bmatrix}
+ \begin{bmatrix} b^{c1} & b^{c2} & \ldots & b^{cK} \end{bmatrix}
$$

$$
= \begin{bmatrix}
\mathbf{w}^{c1}\mathbf{x}^{(1)} + b^{c1} & \mathbf{w}^{c2}\mathbf{x}^{(1)} + b^{c2} & \ldots & \mathbf{w}^{cK}\mathbf{x}^{(1)} + b^{cK} \\
\mathbf{w}^{c1}\mathbf{x}^{(2)} + b^{c1} & \mathbf{w}^{c2}\mathbf{x}^{(2)} + b^{c2} & \ldots & \mathbf{w}^{cK}\mathbf{x}^{(2)} + b^{cK} \\
\ldots & \ldots & \ldots & \ldots \\
\mathbf{w}^{c1}\mathbf{x}^{(m)} + b^{c1} & \mathbf{w}^{c2}\mathbf{x}^{(m)} + b^{c2} & \ldots & \mathbf{w}^{cK}\mathbf{x}^{(m)} + b^{cK}
\end{bmatrix}
$$

# Calculate all observations at one time: Matrix multiplication

$$
\begin{matrix}
\mathbf{x}^{(1)} \\
\mathbf{x}^{(2)} \\
\\
\mathbf{x}^{(m)}
\end{matrix}
\begin{bmatrix}
x_1^{(1)} & x_2^{(1)} & \ldots & x_n^{(1)} \\
x_1^{(2)} & x_2^{(2)} & \ldots & x_n^{(2)} \\
\ldots & \ldots & \ldots & \ldots \\
x_1^{(m)} & x_2^{(m)} & \ldots & x_n^{(m)}
\end{bmatrix}
\bullet
\begin{matrix}
\mathbf{w^c} & \mathbf{w^c} & & \mathbf{w^{cK}} \\
\mathbf{1} & \mathbf{2} & & \\
\end{matrix}
\begin{bmatrix}
w_1^{c1} & w_1^{c2} & \ldots & w_1^{cK} \\
w_2^{c1} & w_2^{c2} & \ldots & w_2^{cK} \\
\ldots & \ldots & \ldots & \ldots \\
w_n^{c1} & w_n^{c2} & \ldots & w_n^{cK}
\end{bmatrix}
+ \begin{bmatrix} b^{c1} & b^{c2} & \ldots & b^{cK} \end{bmatrix}
$$

$$
= \begin{bmatrix}
\mathbf{w}^{c1}\mathbf{x}^{(1)} + b^{c1} & \mathbf{w}^{c2}\mathbf{x}^{(1)} + b^{c2} & \ldots & \mathbf{w}^{cK}\mathbf{x}^{(1)} + b^{cK} \\
\mathbf{w}^{c1}\mathbf{x}^{(2)} + b^{c1} & \mathbf{w}^{c2}\mathbf{x}^{(2)} + b^{c2} & \ldots & \mathbf{w}^{cK}\mathbf{x}^{(2)} + b^{cK} \\
\ldots & \ldots & \ldots & \ldots \\
\mathbf{w}^{c1}\mathbf{x}^{(m)} + b^{c1} & \mathbf{w}^{c2}\mathbf{x}^{(m)} + b^{c2} & \ldots & \mathbf{w}^{cK}\mathbf{x}^{(m)} + b^{cK}
\end{bmatrix}
$$

CONSTELLATE

$$\mathbf{x}^{(1)} \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \ldots & x_n^{(1)} \\ x_1^{(2)} & x_2^{(2)} & \ldots & x_n^{(2)} \\ \ldots & \ldots & \ldots & \ldots \\ x_1^{(m)} & x_2^{(m)} & \ldots & x_n^{(m)} \end{bmatrix} \bullet \begin{bmatrix} w_1^{c1} & w_1^{c2} & \ldots & w_1^{cK} \\ w_2^{c1} & w_2^{c2} & \ldots & w_2^{cK} \\ \ldots & \ldots & \ldots & \ldots \\ w_n^{c1} & w_n^{c2} & \ldots & w_n^{cK} \end{bmatrix} + \begin{bmatrix} b^{c1} & b^{c2} & \ldots & b^{cK} \end{bmatrix}$$

with labels: $\mathbf{x}^{(1)}$, $\mathbf{x}^{(2)}$, $\mathbf{x}^{(m)}$ on the first matrix; $\mathbf{w^c}_1$, $\mathbf{w^c}_2$, $\mathbf{w^{cK}}$ on the second matrix.

$$= \begin{bmatrix} \mathbf{w}^{c1}\mathbf{x}^{(1)} + b^{c1} & \mathbf{w}^{c2}\mathbf{x}^{(1)} + b^{c2} & \ldots & \mathbf{w}^{cK}\mathbf{x}^{(1)} + b^{cK} \\ \mathbf{w}^{c1}\mathbf{x}^{(2)} + b^{c1} & \mathbf{w}^{c2}\mathbf{x}^{(2)} + b^{c2} & \ldots & \mathbf{w}^{cK}\mathbf{x}^{(2)} + b^{cK} \\ \ldots & \ldots & \ldots & \ldots \\ \mathbf{w}^{c1}\mathbf{x}^{(m)} + b^{c1} & \mathbf{w}^{c2}\mathbf{x}^{(m)} + b^{c2} & \ldots & \mathbf{w}^{cK}\mathbf{x}^{(m)} + b^{cK} \end{bmatrix}$$

# 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

Vaswani, Ashish et al. (2017). Attention is all you need. In *Advances in Neural Information Processing Systems*. 5998–6008.