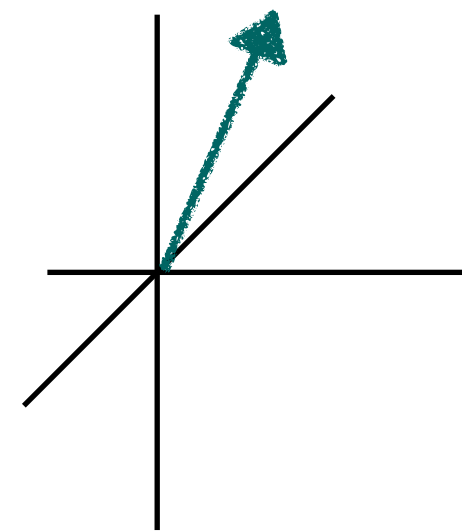
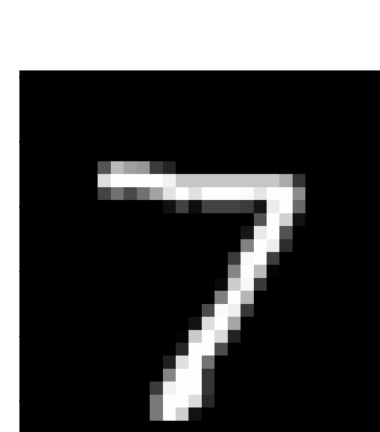
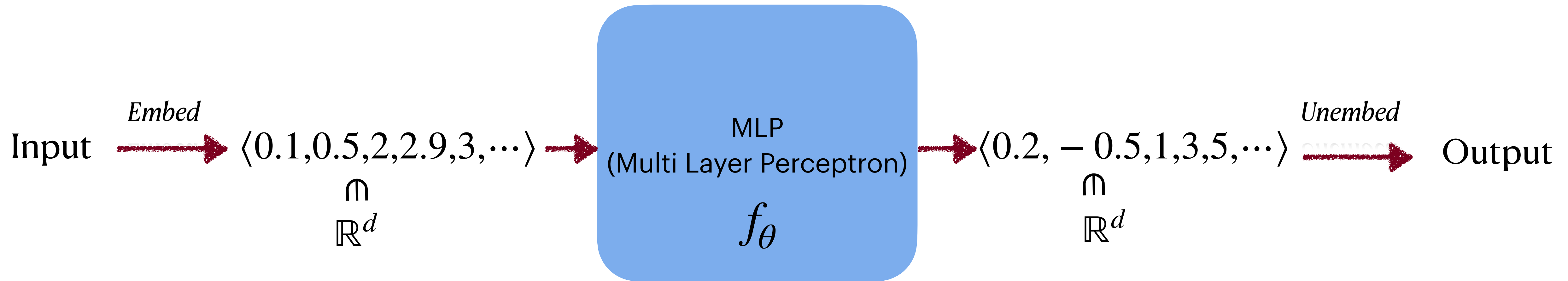


# **Analysis of reasoning models**

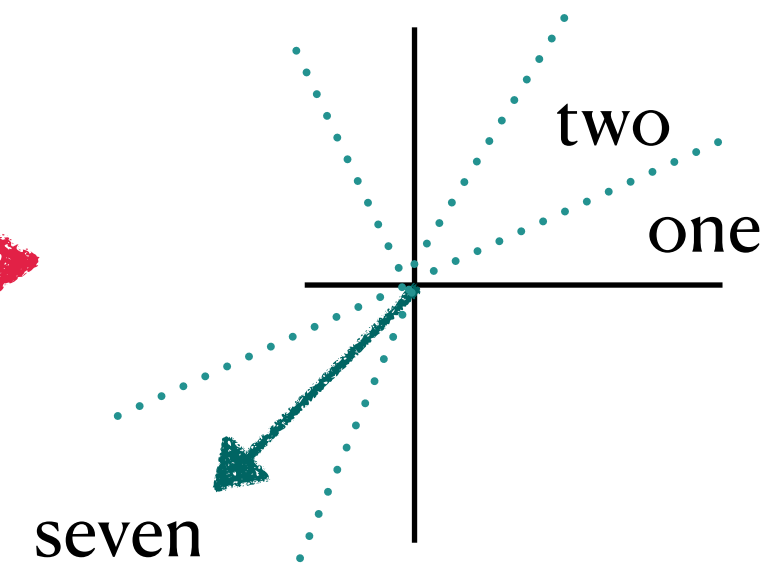
**JMM 2026, Washington DC**

Abhinav Chand 01/06/2026

# MLP

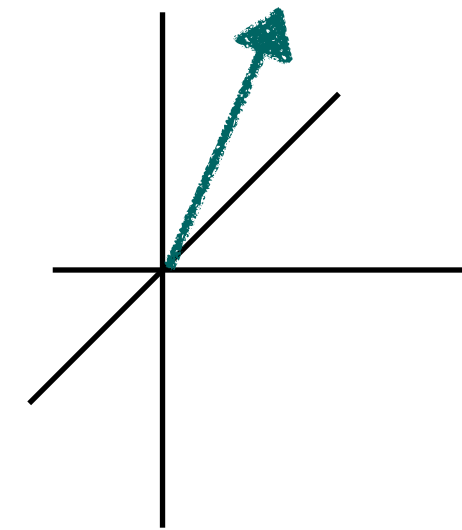


$f_\theta$

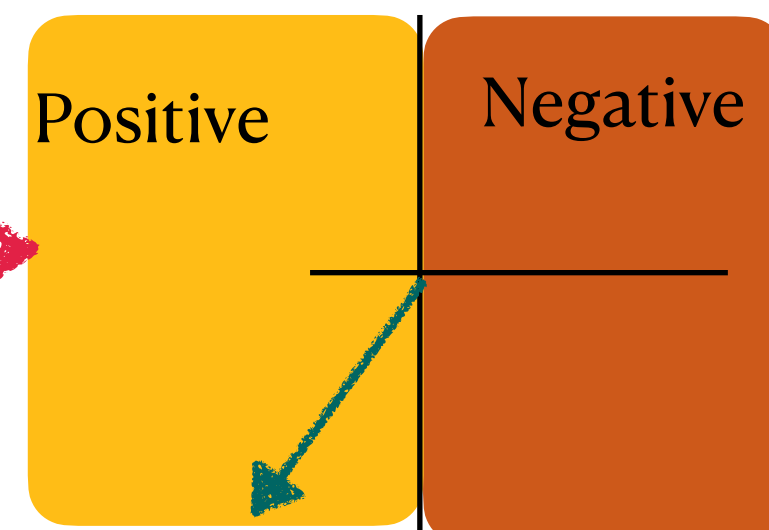


Seven

Amazing

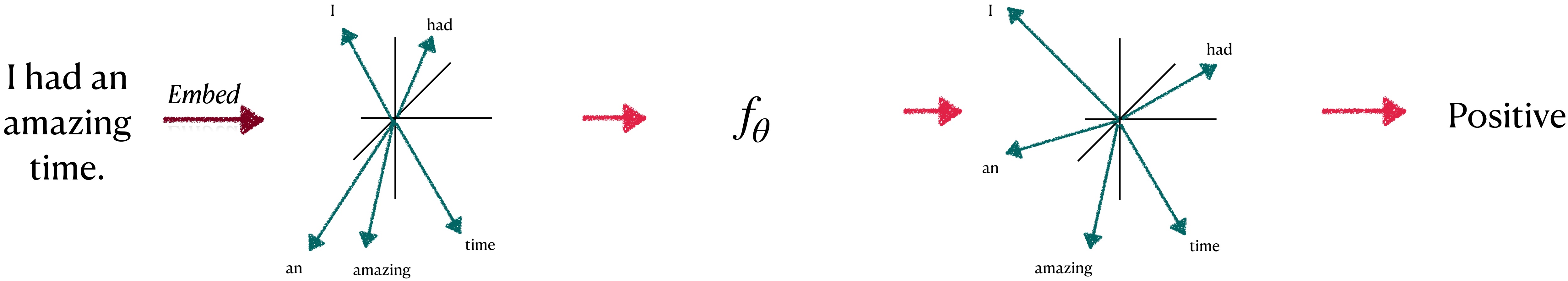
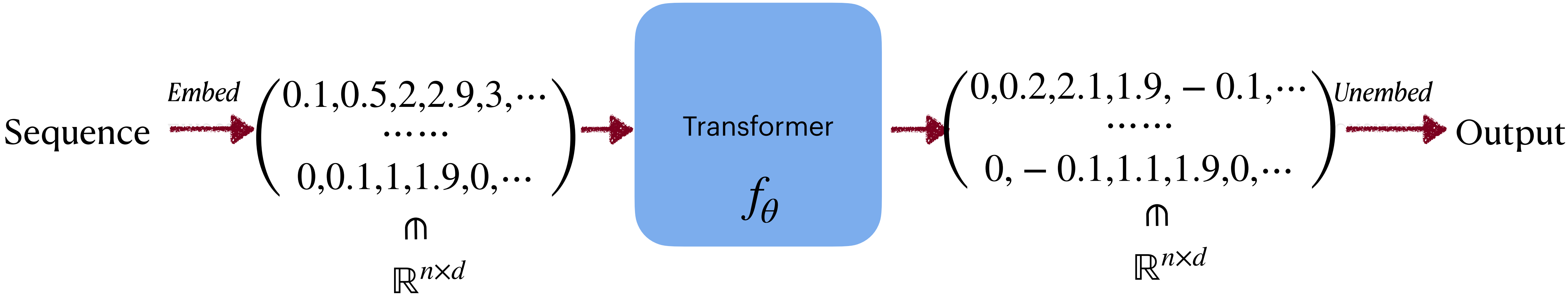


$f_\theta$



Positive

# Transformer

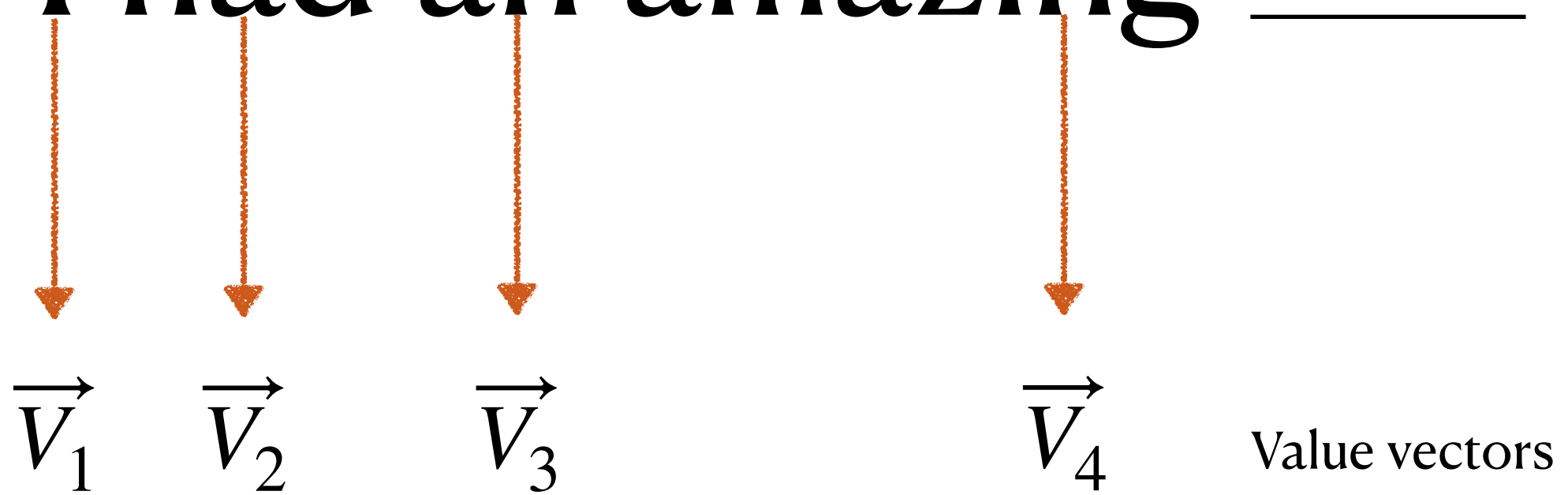


**Attention is all we need**

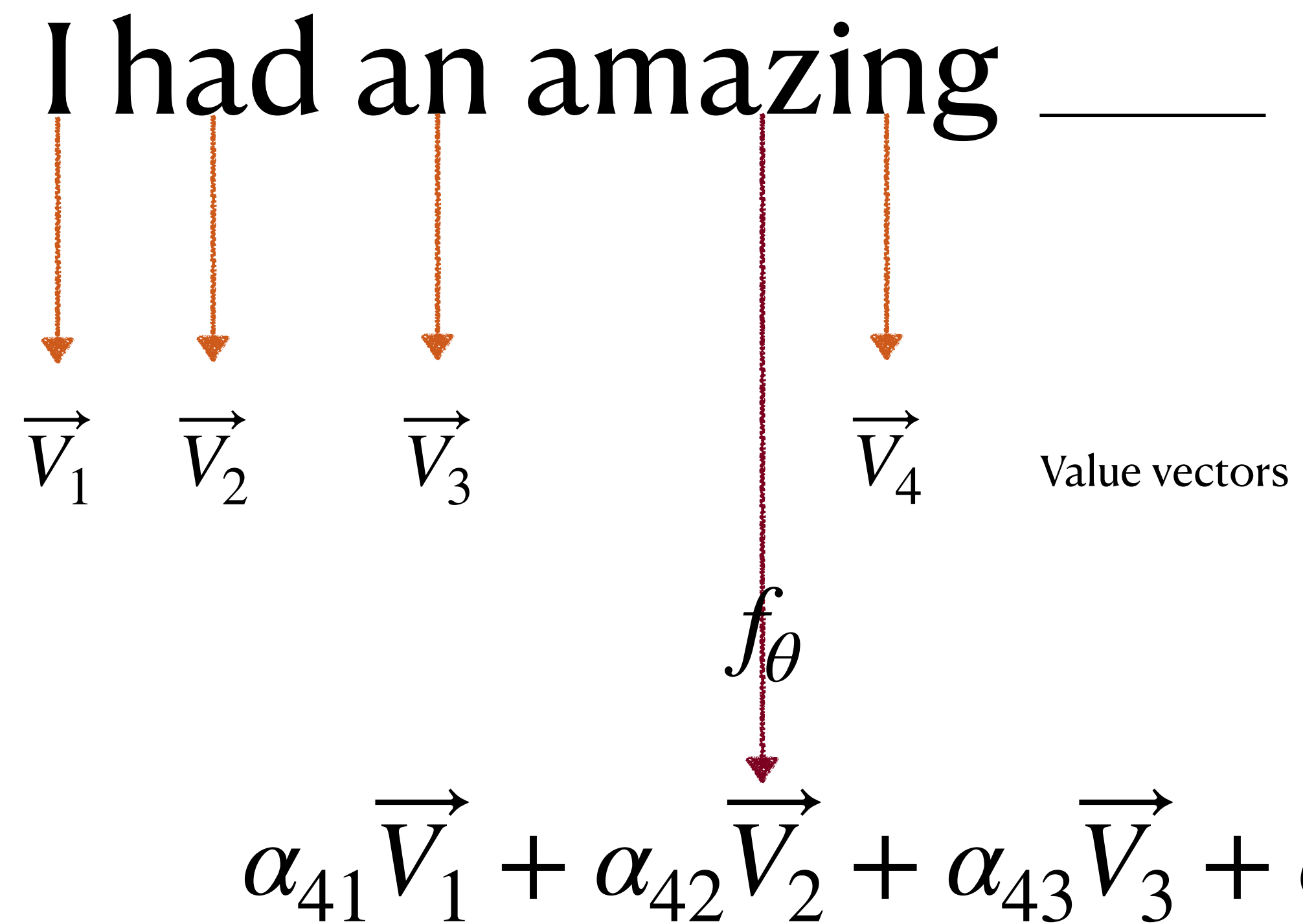
I had an amazing \_\_\_\_\_

# Attention is all we need

I had an amazing \_\_\_\_\_



# Attention is all we need



# Attention is all we need

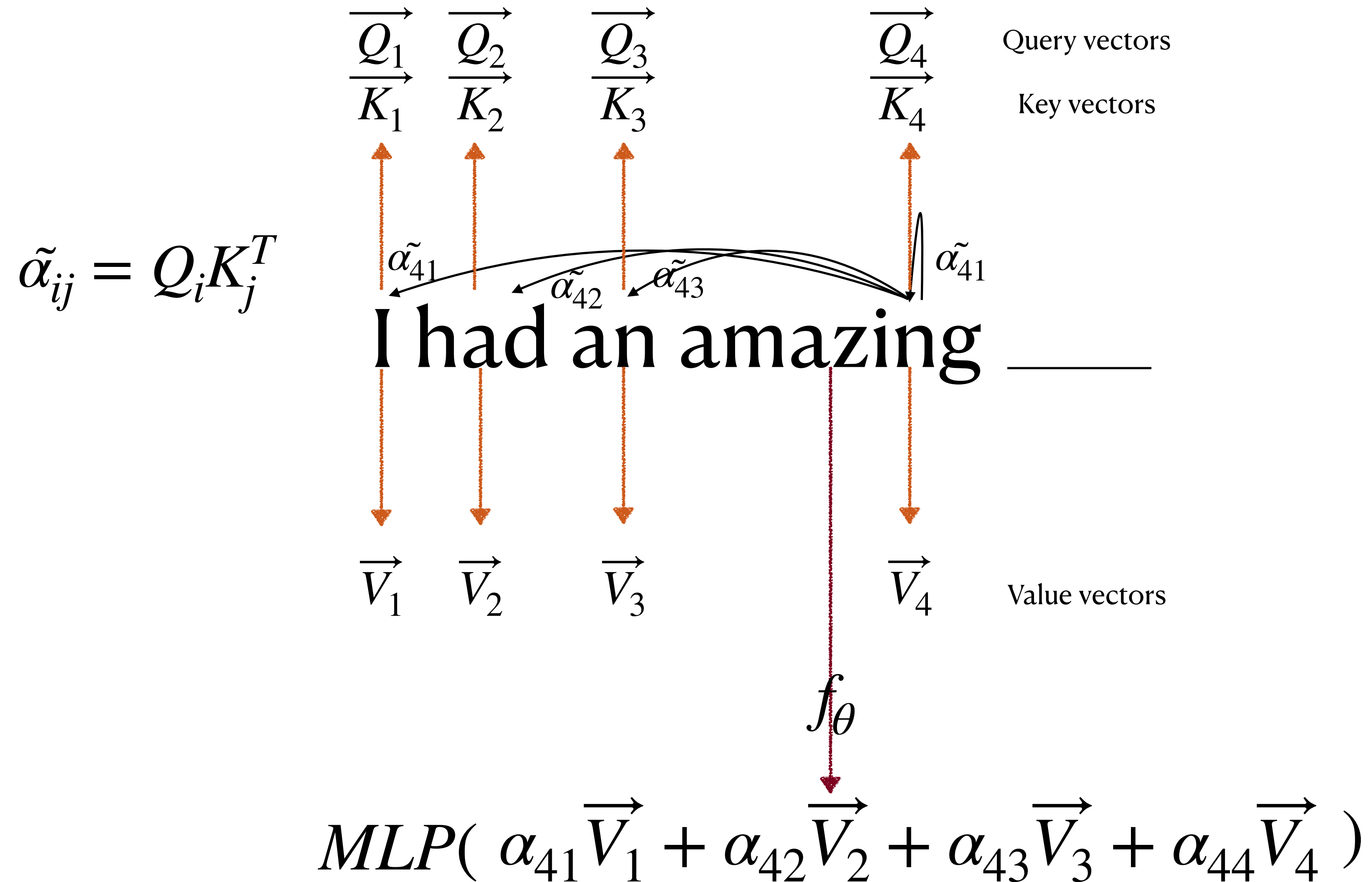
I had an amazing \_\_\_\_\_



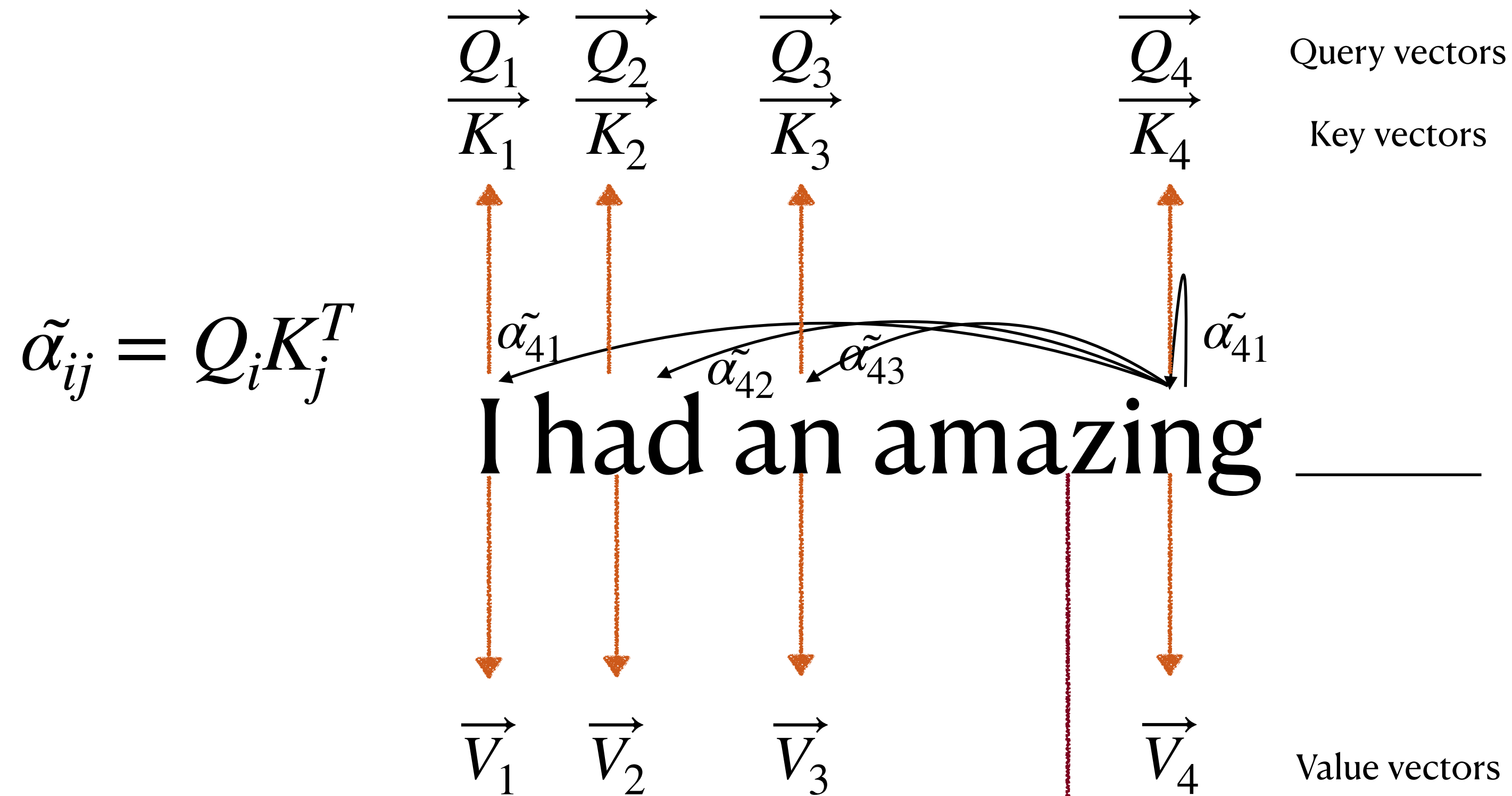
$f_\theta$

$$MLP( \alpha_{41} \vec{V}_1 + \alpha_{42} \vec{V}_2 + \alpha_{43} \vec{V}_3 + \alpha_{44} \vec{V}_4 )$$

# Attention is all we need



# Attention is all we need



$$\langle \alpha_{41}, \alpha_{42}, \alpha_{43}, \alpha_{44} \rangle$$

$$= \text{softmax}(\langle \tilde{\alpha}_{41}, \tilde{\alpha}_{42}, \tilde{\alpha}_{43}, \tilde{\alpha}_{44} \rangle)$$

$$= \left\langle \dots, \frac{e^{\tilde{\alpha}_{4i}}}{\sum_j e^{\tilde{\alpha}_{4j}}}, \dots \right\rangle$$

$$MLP( \alpha_{41} \vec{V}_1 + \alpha_{42} \vec{V}_2 + \alpha_{43} \vec{V}_3 + \alpha_{44} \vec{V}_4 )$$

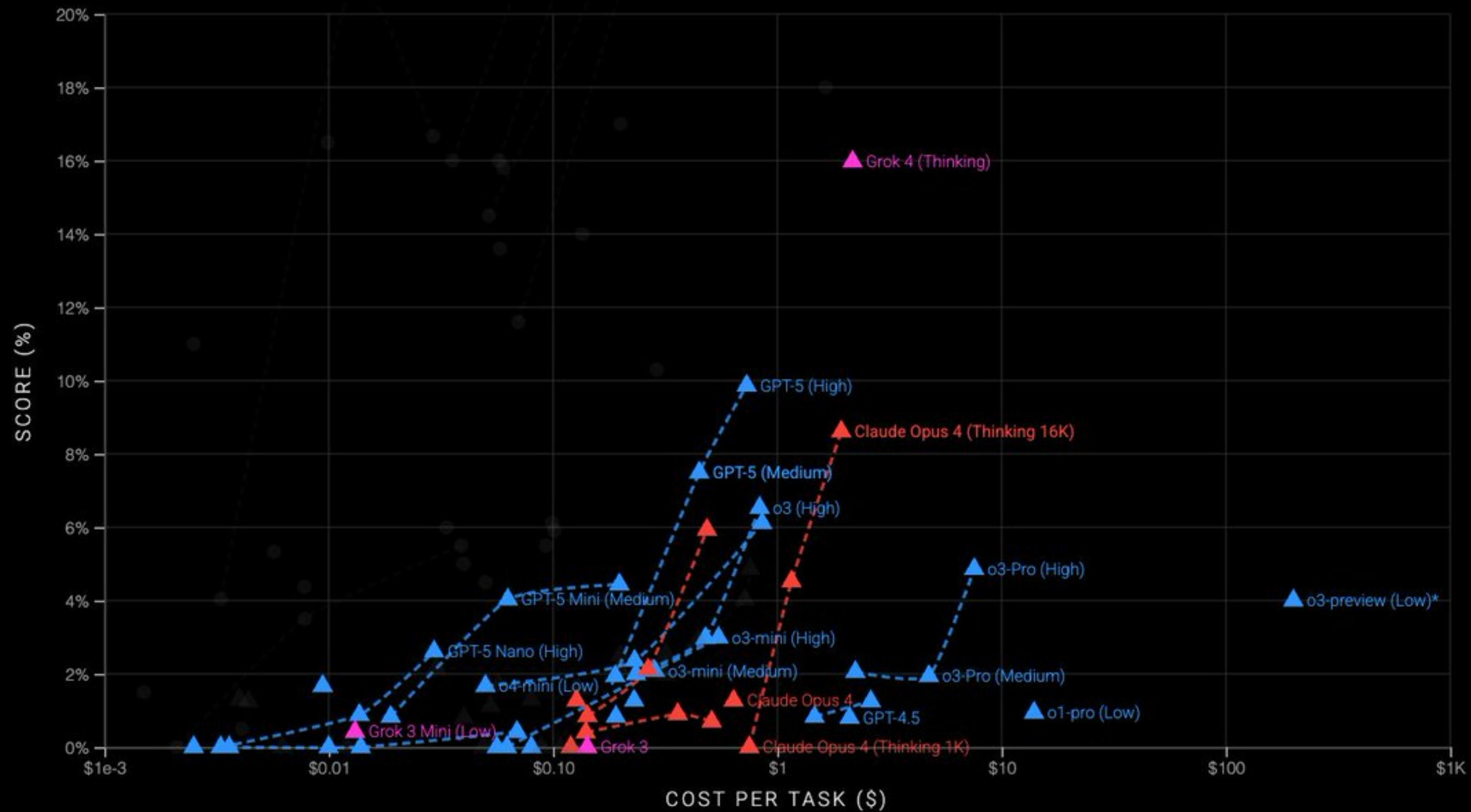
# ARC Prize



- Abstraction and Reasoning Corpus challenge
- ARC-AGI I (2019), ARC-AGI 2 (2025)

# August leaderboard

## ARC-AGI-2 LEADERBOARD



# Novel problem for deep learning

- Each puzzle is unique. Test puzzles not similar to training puzzles.
- Few shot learning

The image shows a user interface for a puzzle-solving task, divided into two main sections: 'EXAMPLES' and 'TEST'.

**EXAMPLES:** This section shows two examples of 18x18 grids. Each example consists of an 'Input' grid and an 'Output' grid. In the first example, the background is yellow, and the puzzle pieces are blue, red, and green. In the second example, the background is green, and the puzzle pieces are blue, red, yellow, and pink. Arrows indicate the transformation from input to output.

**TEST:** This section shows a 28x28 grid. The 'Input' grid has a light blue background with a few colored pieces (red, blue, green, yellow). The 'Output' grid is currently empty. Below the grids is a control panel with three steps:

1. Configure your output grid: A dropdown menu shows '28x28', with 'Resize' and 'Copy from input' buttons. 'Clear' and 'Reset' buttons are also present.
2. Edit your output grid cells: Includes 'Edit', 'Select', and 'Fill' buttons. A color palette below shows various colors, with a blue square selected and the number '8' next to it.
3. See if your output is correct: A 'Submit solution' button.

A vertical watermark 'TOGGLE ANIMATION' is visible on the right side of the interface.

# Hierarchical Reasoning Model



WTF... First they said no to Elon Musk.

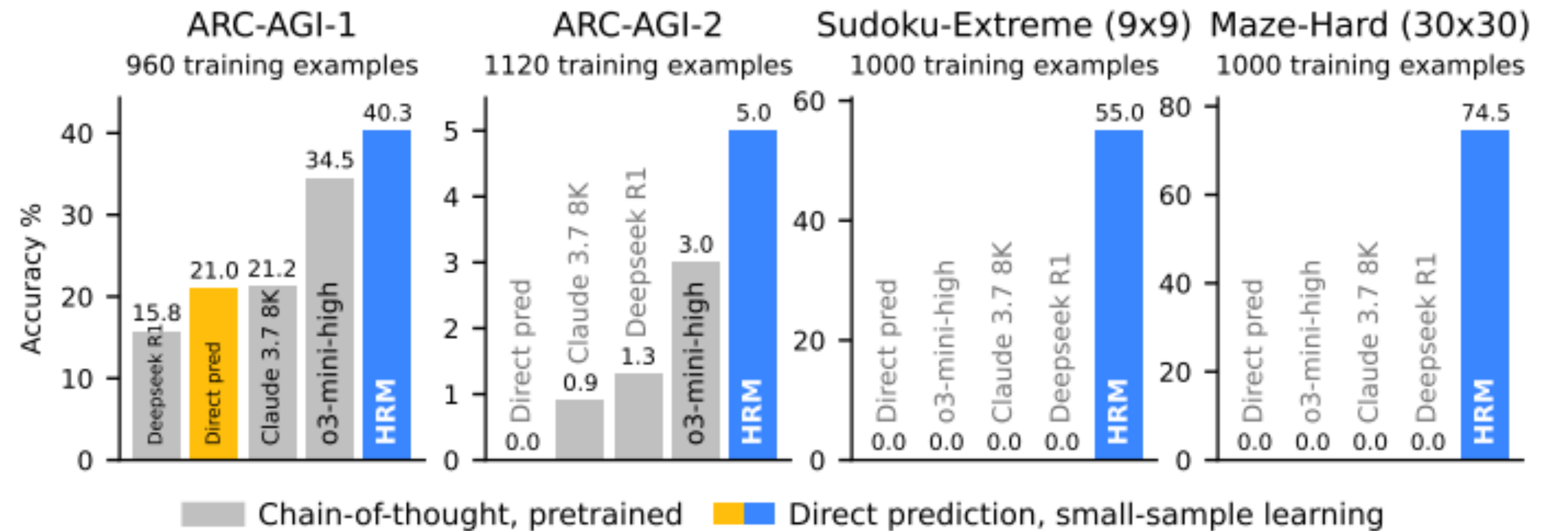
Then they built an AI that just shook the entire LLM industry.

- 27M parameters
- Trained on only 1000 examples
- Inspired by the human brain
- Beat Claude, o3-mini & Deepseek

You can't afford to ignore this: 🧵



8:47 AM · Aug 1, 2025

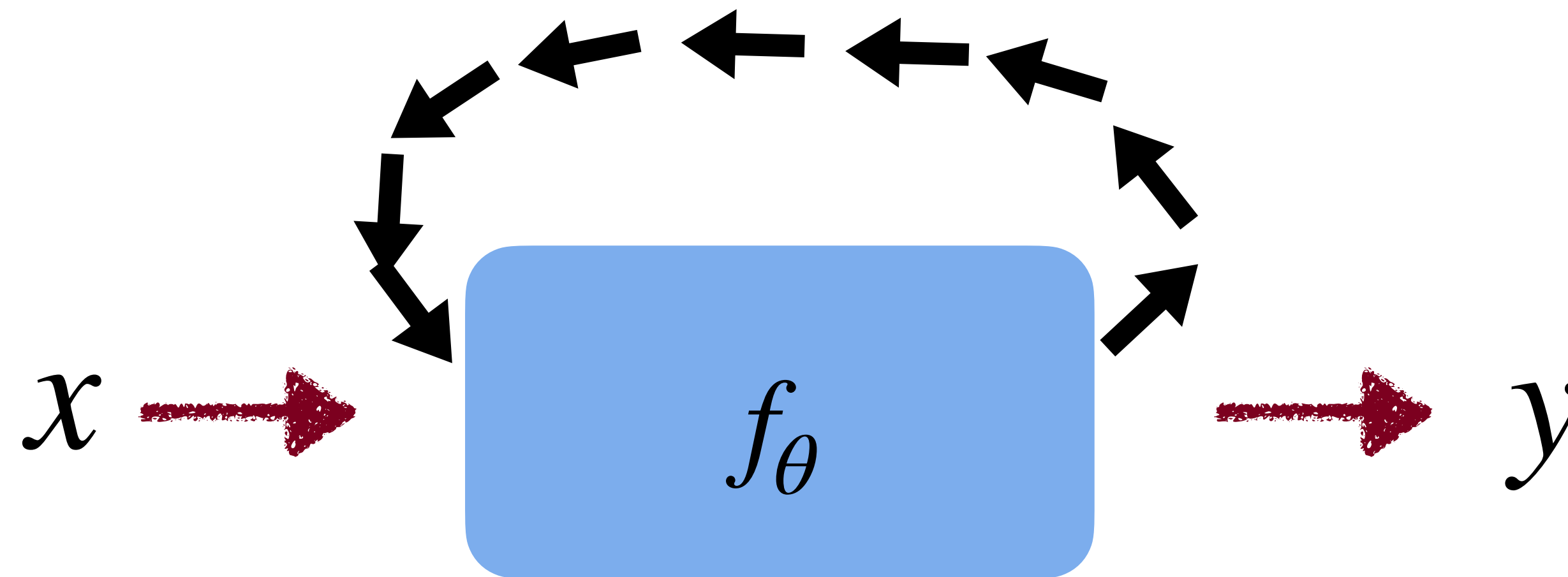


# Looped transformers



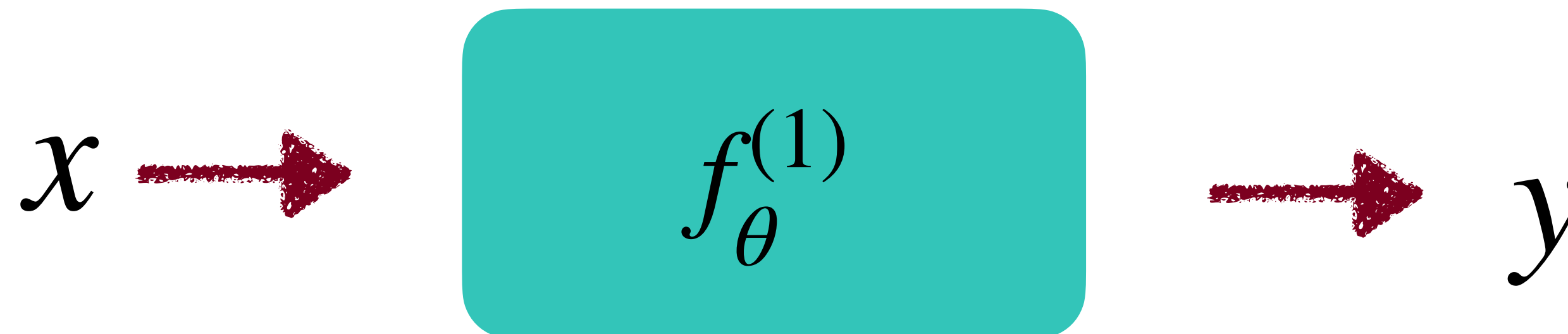
# Why looped transformers?

Looped



vs.

Vanilla



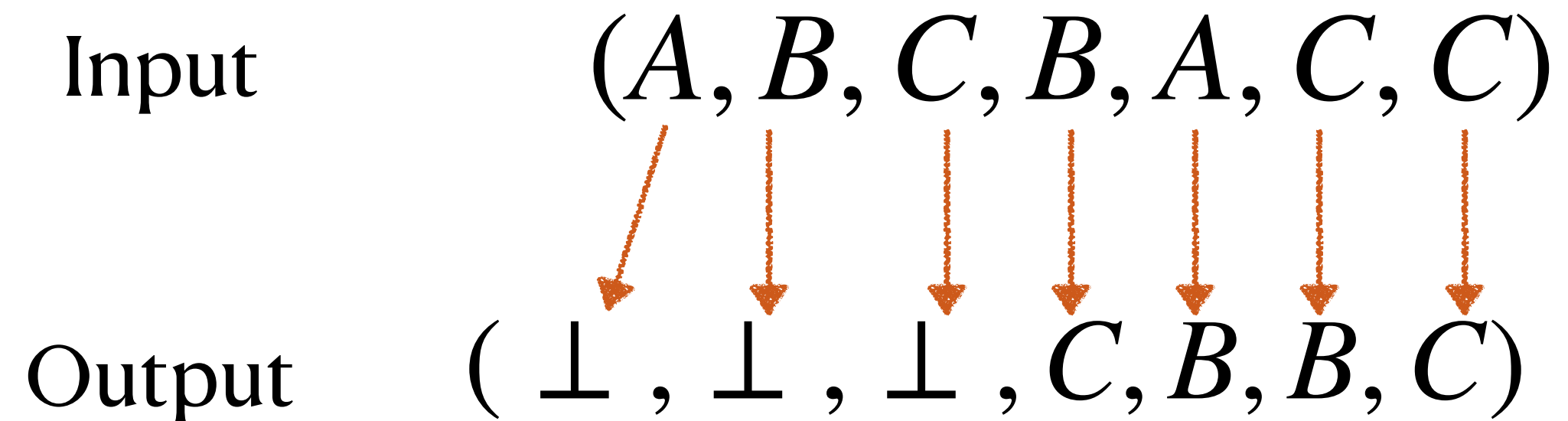
# Induction head task

$$\Sigma := \{A, B, C\}$$

Take a sequence of alphabets:

Input  $(A, B, C, B, A, C, C)$

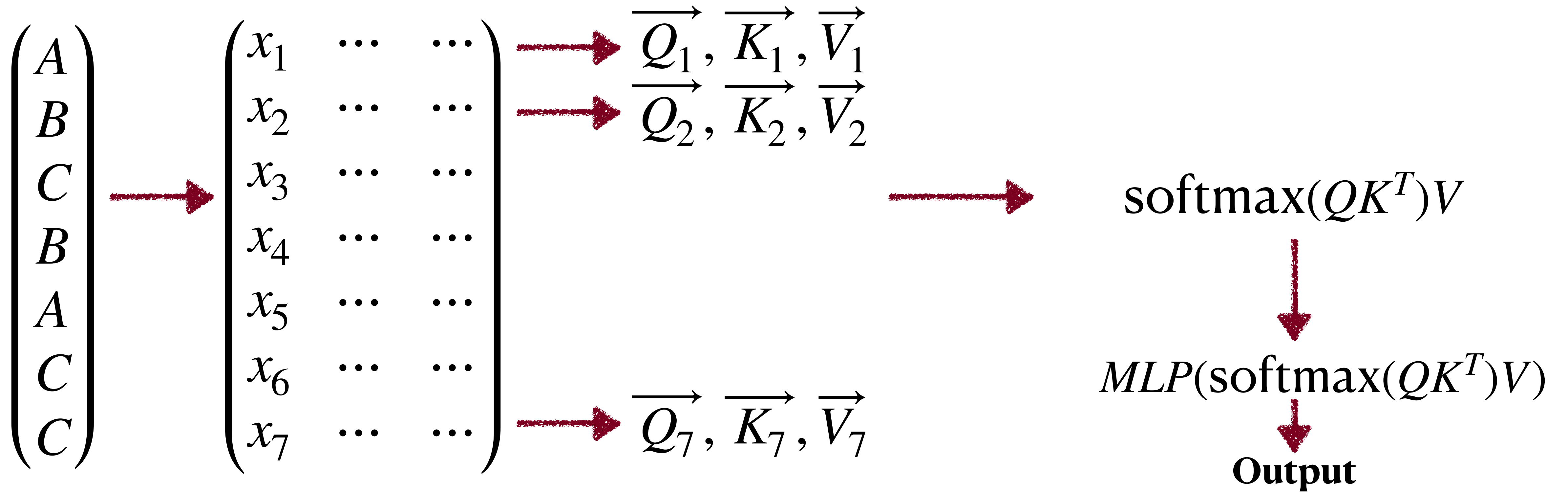
Output  $(\perp, \perp, \perp, C, B, B, C)$



- Look at previous occurrence.
- Output the token that follows the previous token.
- Output  $\perp$  if there is no previous occurrence.

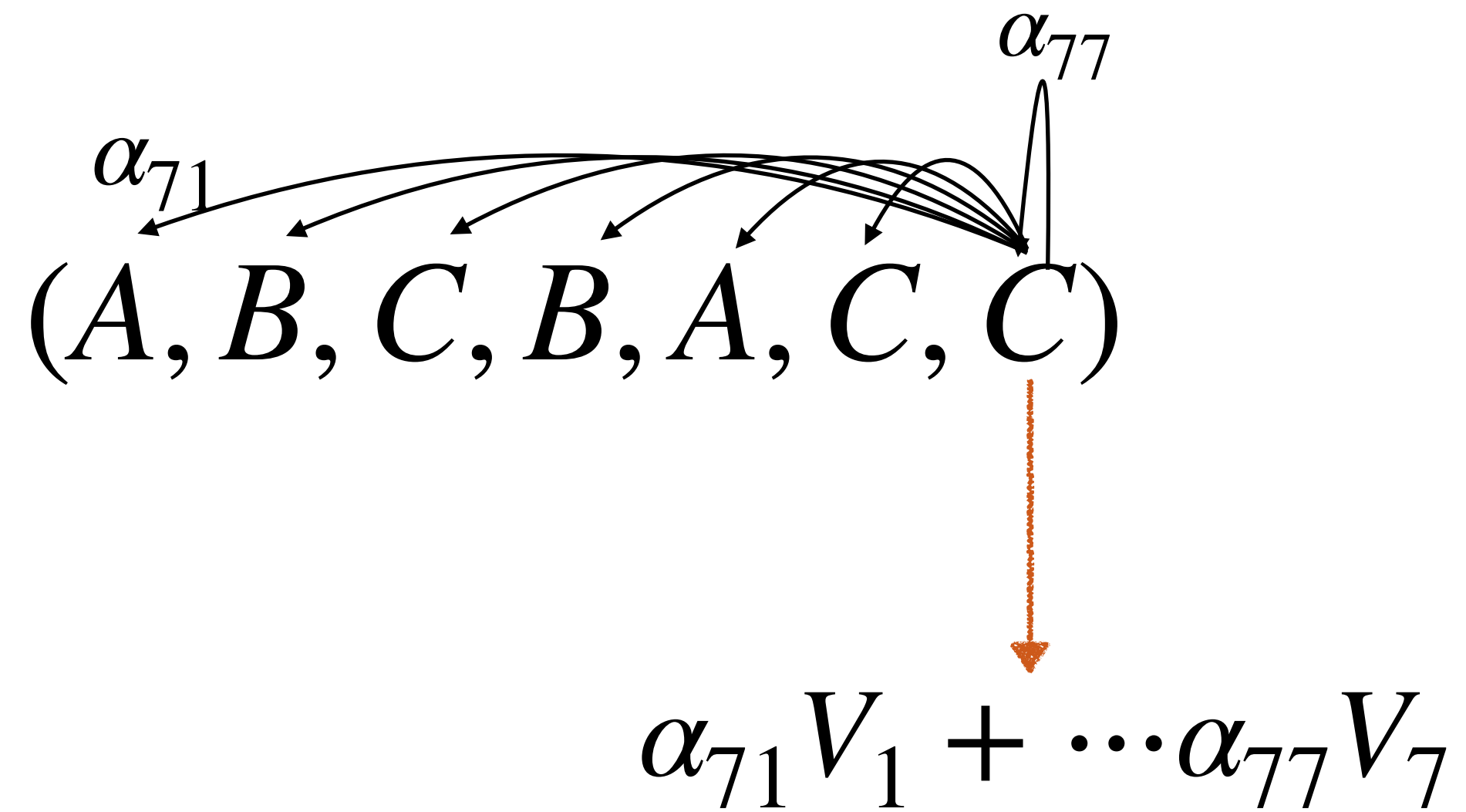
# Transformer architecture

$(A, B, C, B, A, C, C)$

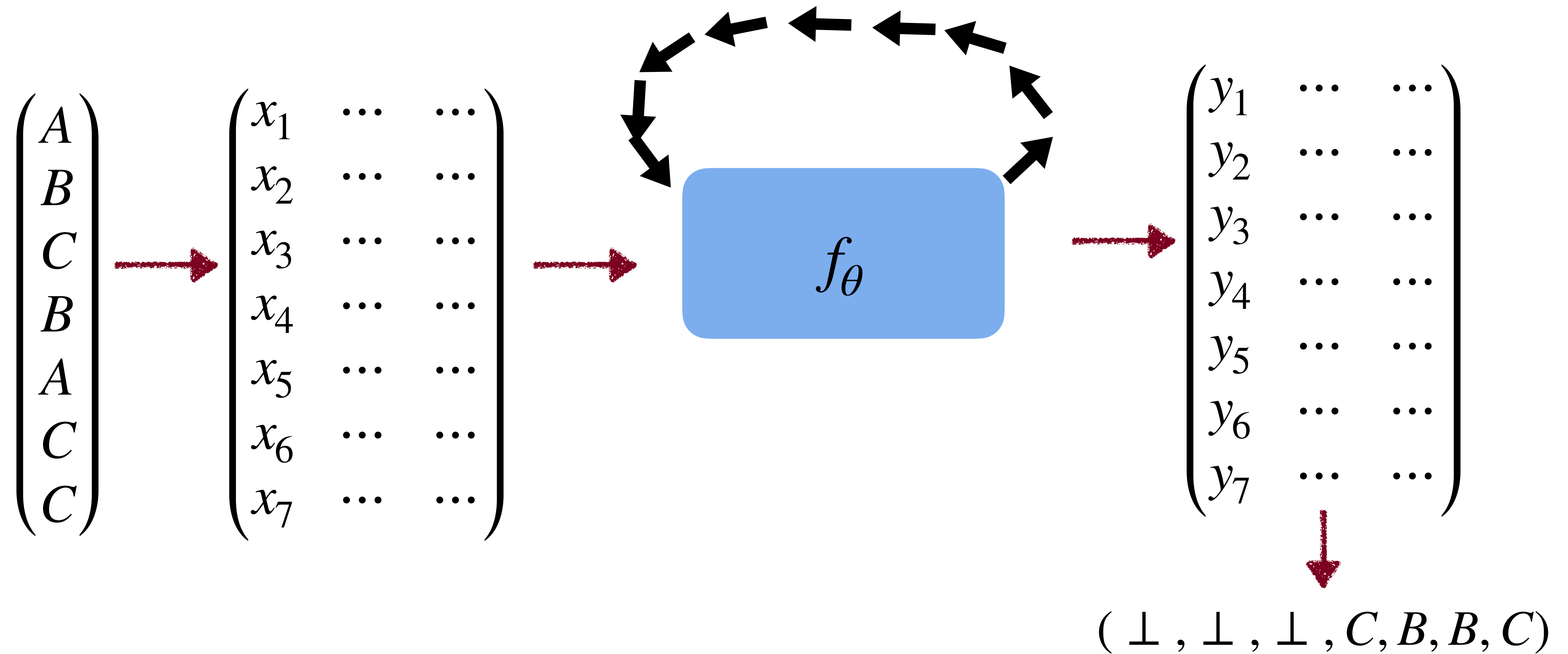


# Attention is all we need

$$\alpha_{ij} = Q_i K_j^T$$



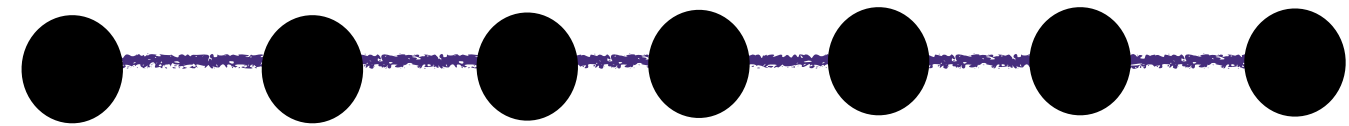
# Loop transformer architecture



# Theory first (mistake?)

## GNN solution

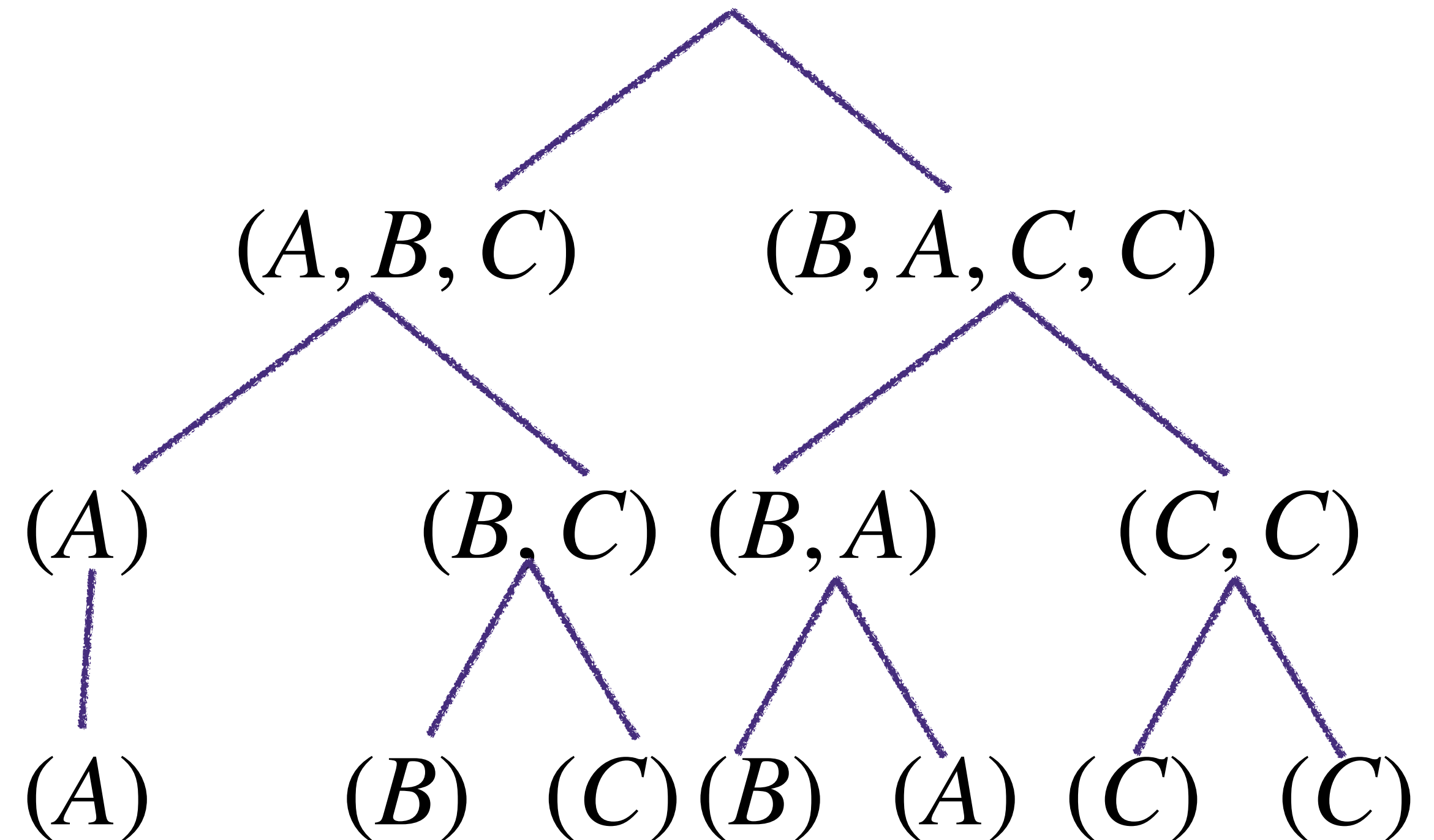
$(A, B, C, B, A, C, C)$



- Shift until you find a match.
- Path graph neural network

## Tree solution

$(A, B, C, B, A, C, C)$



# Experiment

$(A, B, C, B, A, C, C)$

↓ embed layer

$X$

↓

$$X = X + PE([1, 2, \dots, n])$$

$$PE : \mathbb{N} \rightarrow R^d$$

$$\text{softmax}(XW_Q(XW_K)^T)XW_V$$

Identity  $\rightarrow W_Q : R^d \rightarrow R^d$

Trainable  $\rightarrow W_K : R^d \rightarrow R^d$

Identity  $\rightarrow W_V : R^d \rightarrow R^d$

$$\text{Linear}(\text{softmax}(XW_K^T X^T)X)$$

↓

$z$

↓ softmax unembed layer

$(?, ?, ?, ?, ?, ?, ?)$

Add

# Recursive equation

$$X, \quad z_0 = \vec{0}$$

Transformer module:

$f_\theta$

$$z_1 = f_\theta(X + z_0)$$

$$z_2 = f_\theta(X + z_1)$$

... ..

$$z_k = f_\theta(X + z_{k-1})$$

$$Y = \text{Decoder}(z_k)$$

# Dataset

$(A, B, C, B, A, C, C, A, A, B)$

$\downarrow$  *embed layer*

$X$

$\downarrow$

$$X = X + PE([1, 2, \dots, n])$$

$\downarrow$

$$\text{softmax}(XW_Q(XW_K)^T)XW_V$$

$\downarrow$

$$\text{Linear}(\text{softmax}(XW_K^T X^T)X)$$

$\downarrow$

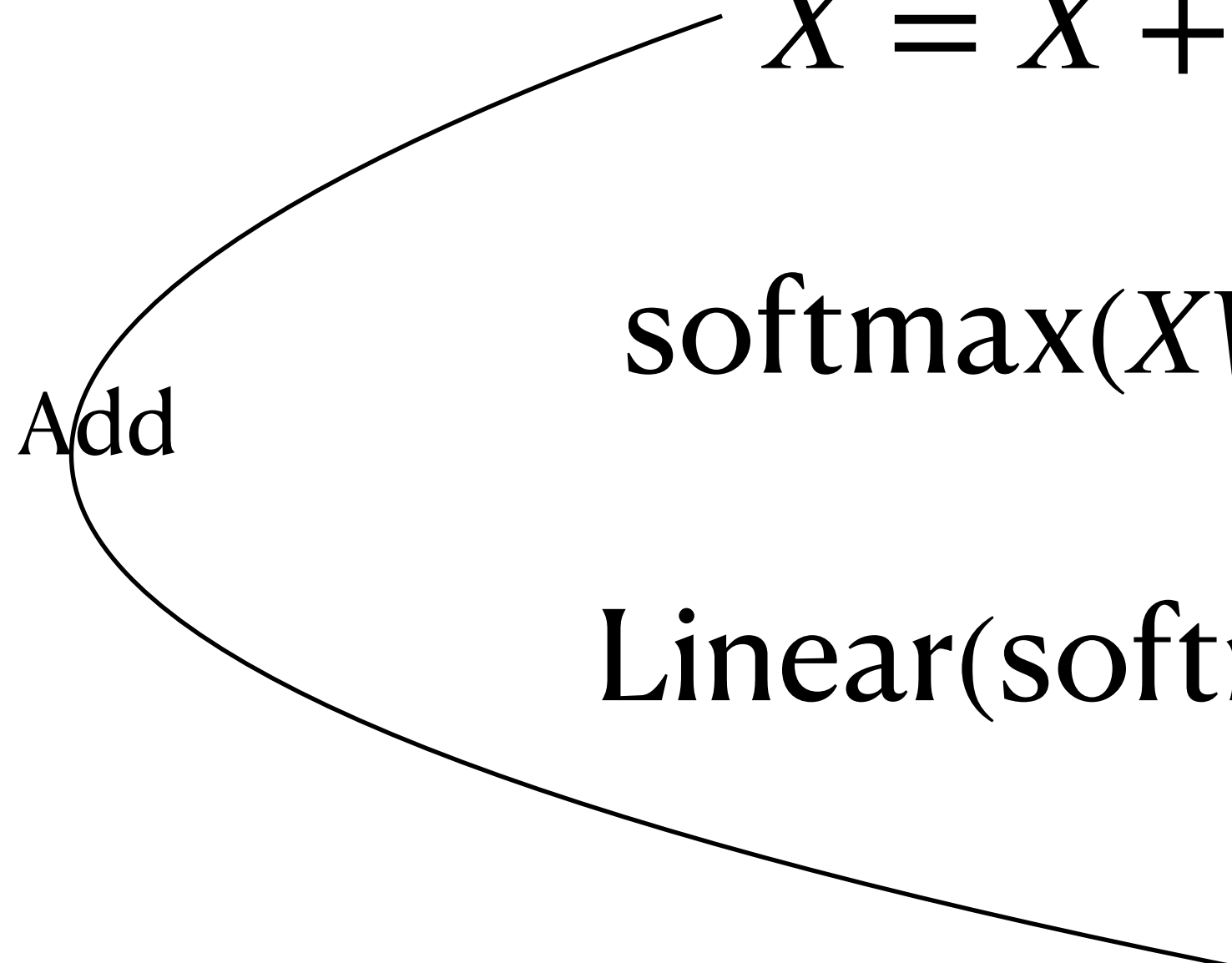
$z$

$\downarrow$  *softmax unembed layer*

$(?, ?, ?, ?, ?, ?, ?)$

- Sequence length = 10
- Alphabet size = 3
- Generated 1000 random sequences and corresponding induction head target
- 80% - 20% train-test split

Add



# Hyperparameters

$(A, B, C, B, A, C, C, A, A, B)$

$\downarrow$  *embed layer*

$X$

$\downarrow$

$$X = X + PE([1, 2, \dots, n])$$

$\downarrow$

$$\text{softmax}(XW_Q(XW_K)^T)XW_V$$

$\downarrow$

$$\text{Linear}(\text{softmax}(XW_K^T X^T)X)$$

$\downarrow$

$z$

$\downarrow$  *softmax unembed layer*

$(?, ?, ?, ?, ?, ?, ?)$

- Embedding dimension = 128
- Learning rate = 0.001
- Number of loops = 10
- Each embedding vector = random gaussian  $\sim N(0, I_d)$

Add

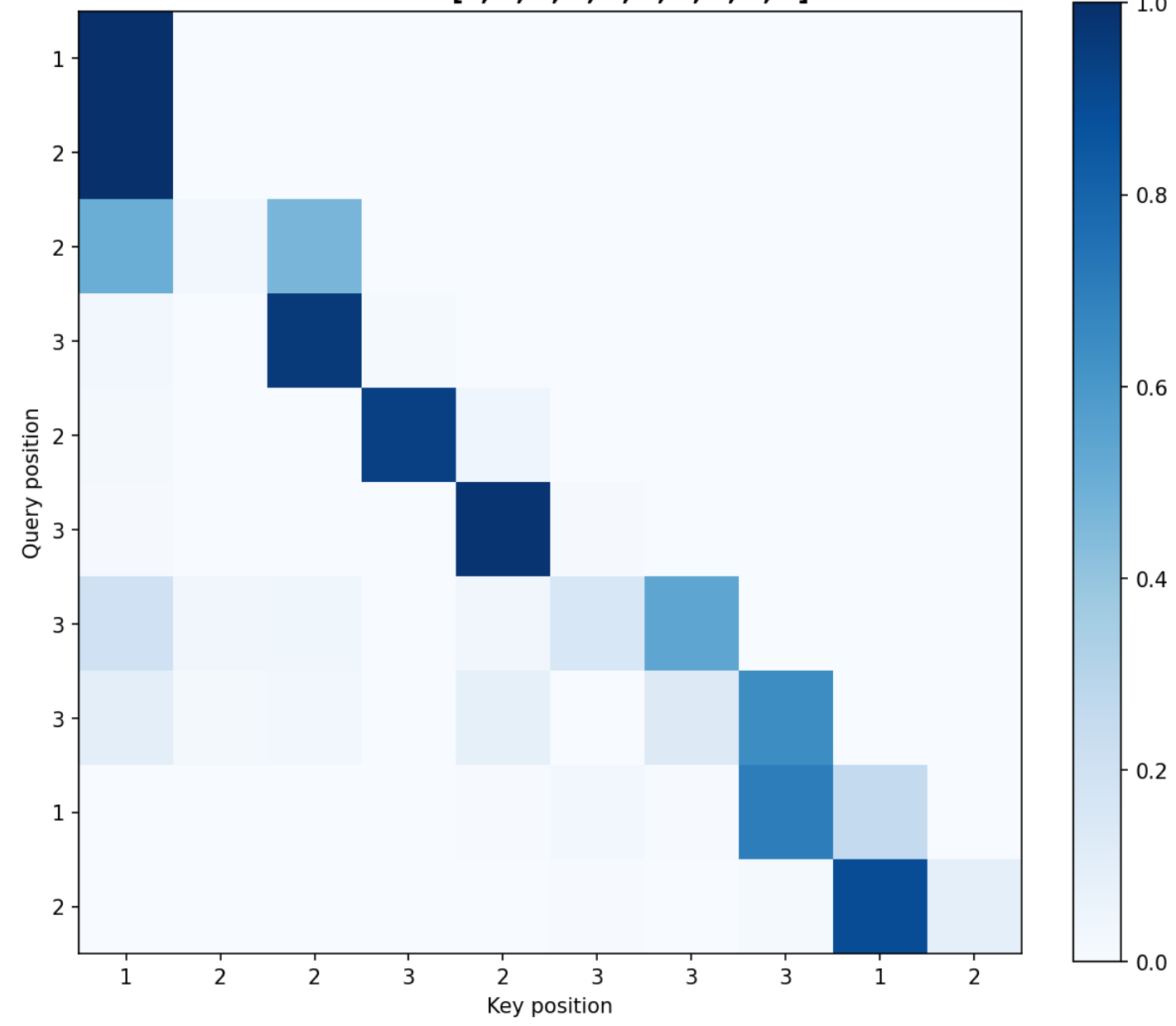




<b>Epoch</b>	<b>Loss</b>	<b>Train Acc</b>	<b>Test Acc</b>
<b>10</b>	0.6265	8.71%	9.60%
<b>20</b>	0.3384	27.90%	24.75%
<b>30</b>	0.1749	53.28%	48.48%
<b>40</b>	0.0889	79.55%	65.66%
<b>50</b>	0.0577	87.25%	73.23%
<b>60</b>	0.0429	86.74%	73.74%
<b>70</b>	0.0443	92.80%	80.81%
<b>80</b>	0.0343	85.48%	74.75%
<b>90</b>	0.0278	91.67%	76.77%
<b>100</b>	0.0235	95.33%	85.86%
<b>110</b>	0.0169	92.80%	82.83%
<b>120</b>	0.0101	95.71%	88.38%
<b>130</b>	0.0116	97.22%	93.43%
<b>140</b>	0.0091	95.71%	92.93%
<b>150</b>	0.0094	96.34%	90.40%
<b>160</b>	0.0064	98.86%	97.47%
<b>170</b>	0.0089	96.34%	95.96%
<b>180</b>	0.0046	99.49%	95.45%
<b>190</b>	0.0006	99.12%	95.96%
<b>200</b>	0.0005	99.87%	98.48%

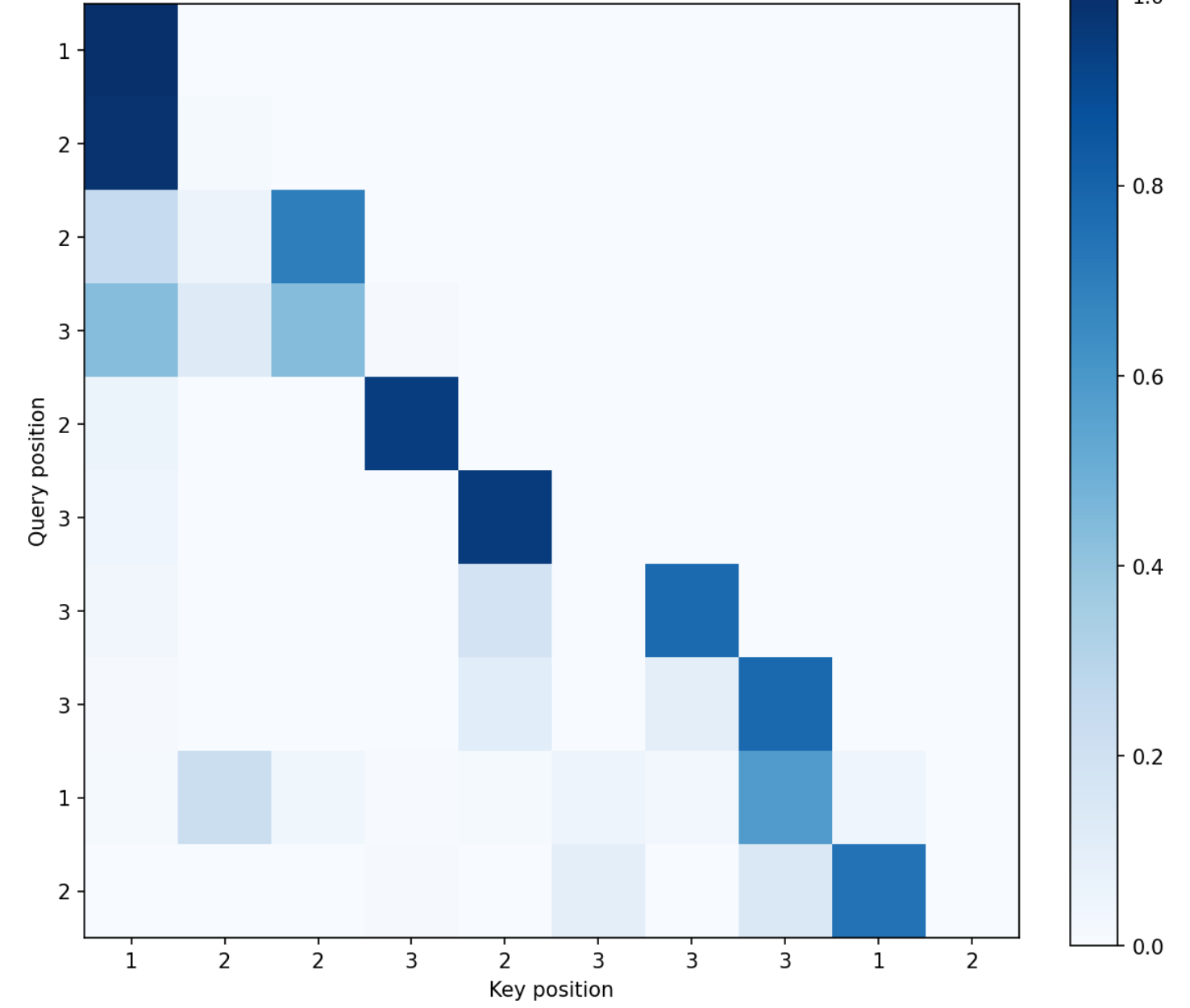
**Iteration 1**  
z: [0, 0, 0, 0, 0, 0, 0, 0, 0, 0]  
x+z: [1, 2, 2, 3, 2, 3, 3, 3, 1, 2]

**Input:** [1, 2, 2, 3, 2, 3, 3, 3, 1, 2] → **Target:** [0, 0, 2, 0, 3, 2, 3, 3, 2, 3]  
**Final Pred:** [0, 0, 2, 0, 3, 2, 3, 3, 2, 3]



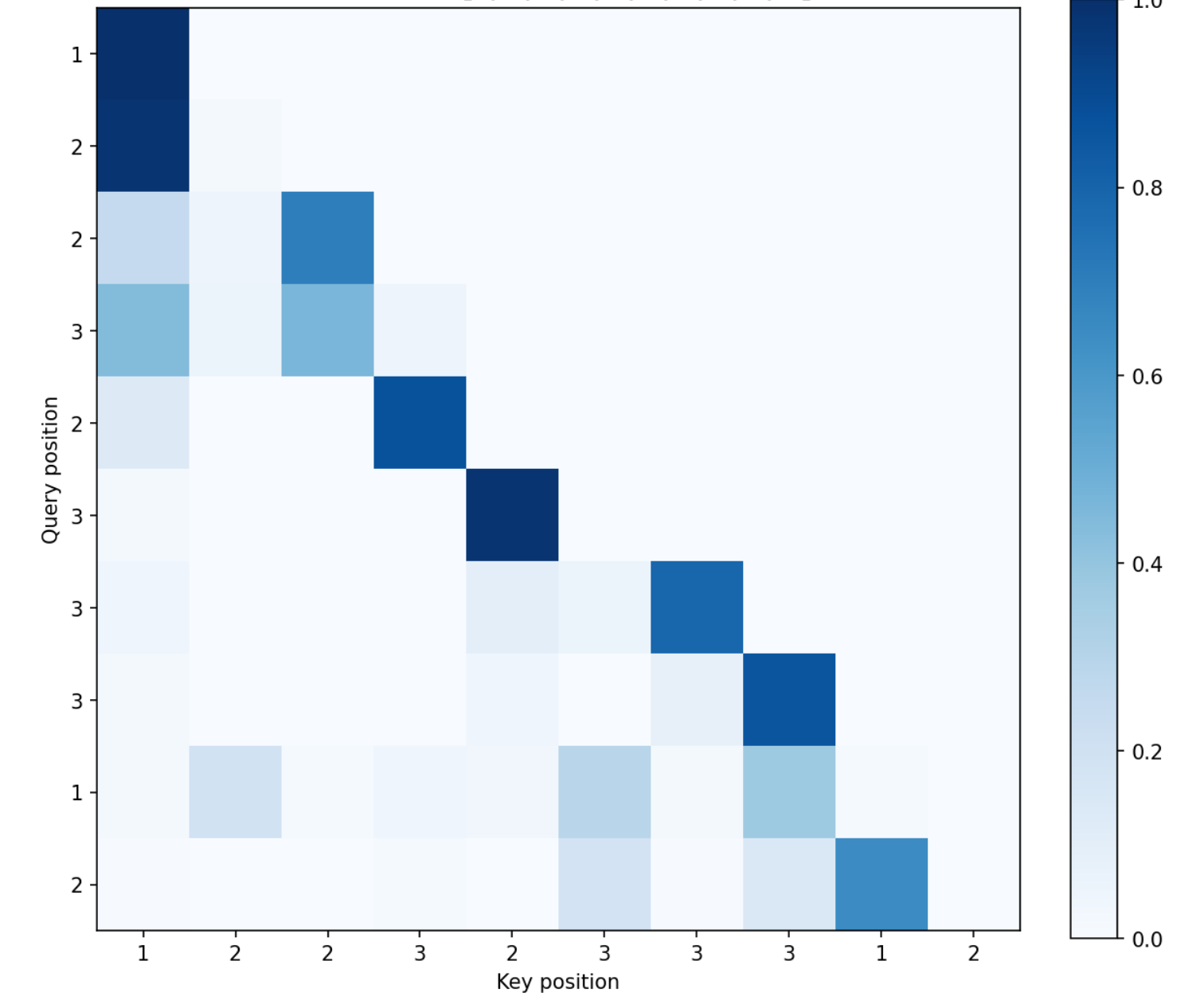
**Iteration 2**  
z: [0, 0, 0, 2, 3, 2, 0, 3, 3, 1]  
x+z: [1, 2, 2, 3, 2, 3, 3, 3, 1, 2]

**Input:** [1, 2, 2, 3, 2, 3, 3, 3, 1, 2] → **Target:** [0, 0, 2, 0, 3, 2, 3, 3, 2, 3]  
**Final Pred:** [0, 0, 2, 0, 3, 2, 3, 3, 2, 3]



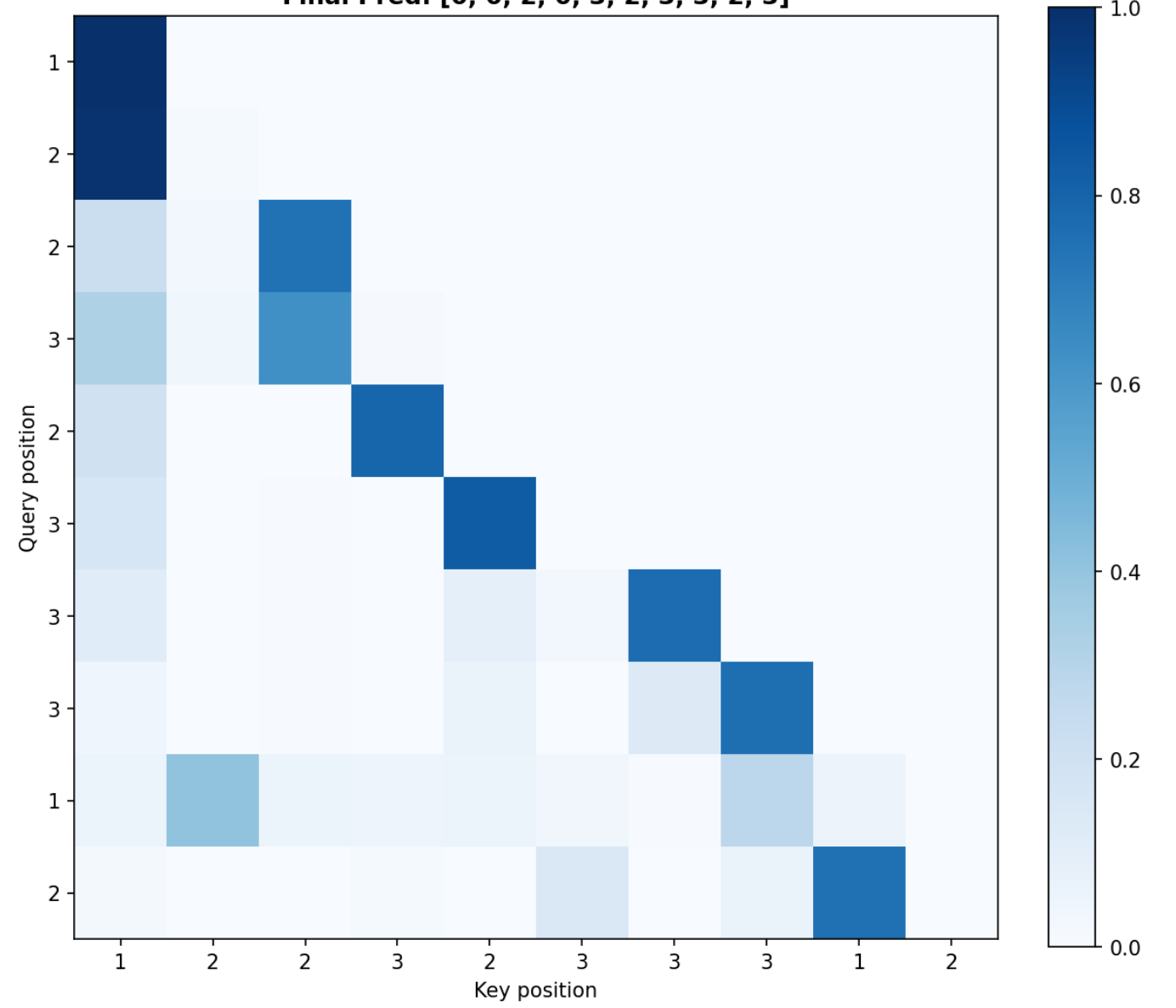
**Iteration 3**  
z: [0, 0, 0, 0, 3, 2, 3, 3, 0, 1]  
x+z: [1, 2, 2, 3, 2, 3, 3, 3, 1, 2]

**Input:** [1, 2, 2, 3, 2, 3, 3, 3, 1, 2] → **Target:** [0, 0, 2, 0, 3, 2, 3, 3, 2, 3]  
**Final Pred:** [0, 0, 2, 0, 3, 2, 3, 3, 2, 3]



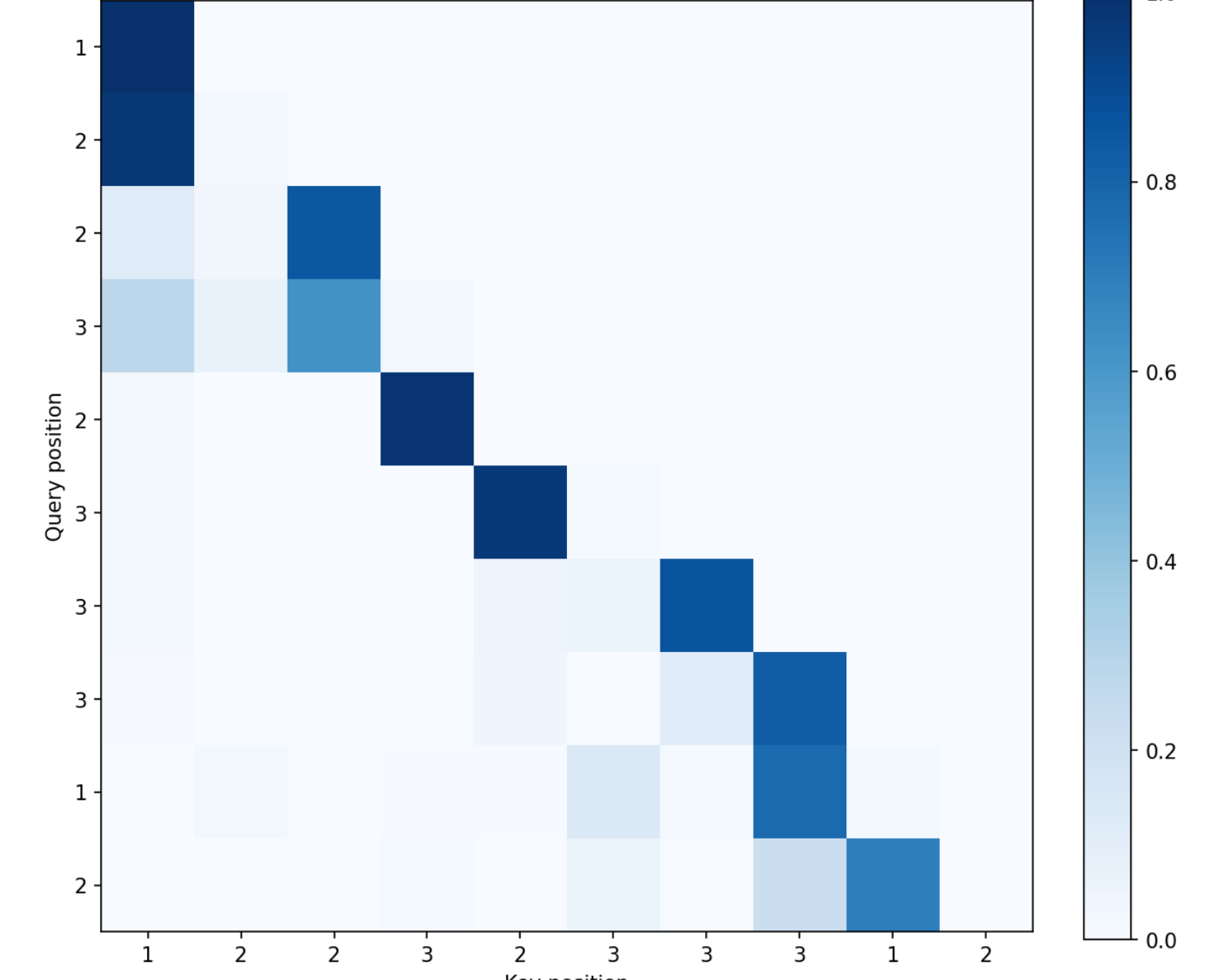
**Iteration 4**  
z: [0, 0, 2, 0, 3, 2, 3, 3, 1]  
x+z: [1, 2, 2, 3, 2, 3, 3, 3, 1, 2]

**Input:** [1, 2, 2, 3, 2, 3, 3, 3, 1, 2] → **Target:** [0, 0, 2, 0, 3, 2, 3, 3, 2, 3]  
**Final Pred:** [0, 0, 2, 0, 3, 2, 3, 3, 2, 3]



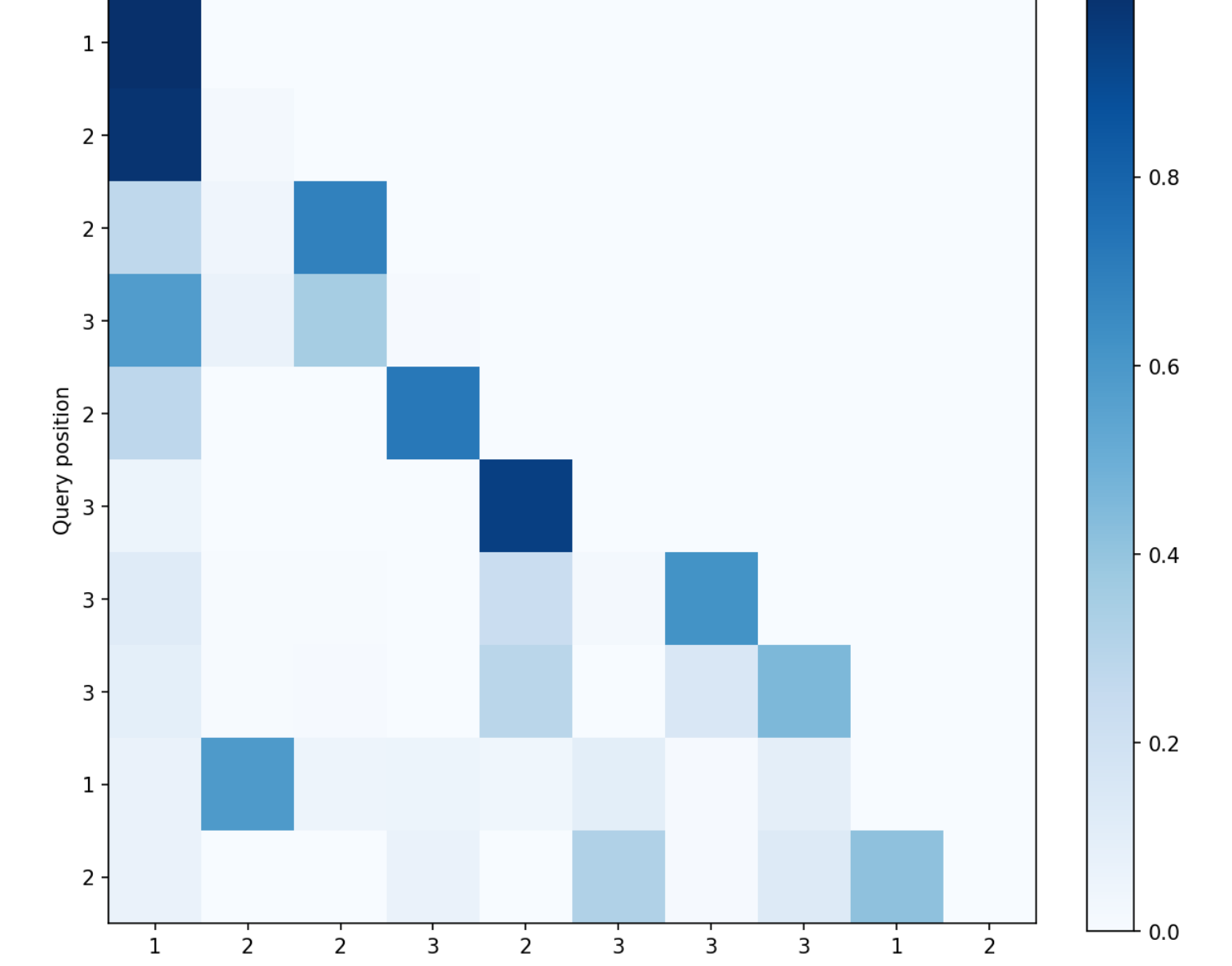
**Iteration 5**  
z: [0, 0, 2, 2, 3, 2, 3, 3, 2, 1]  
x+z: [1, 2, 2, 3, 2, 3, 3, 3, 1, 2]

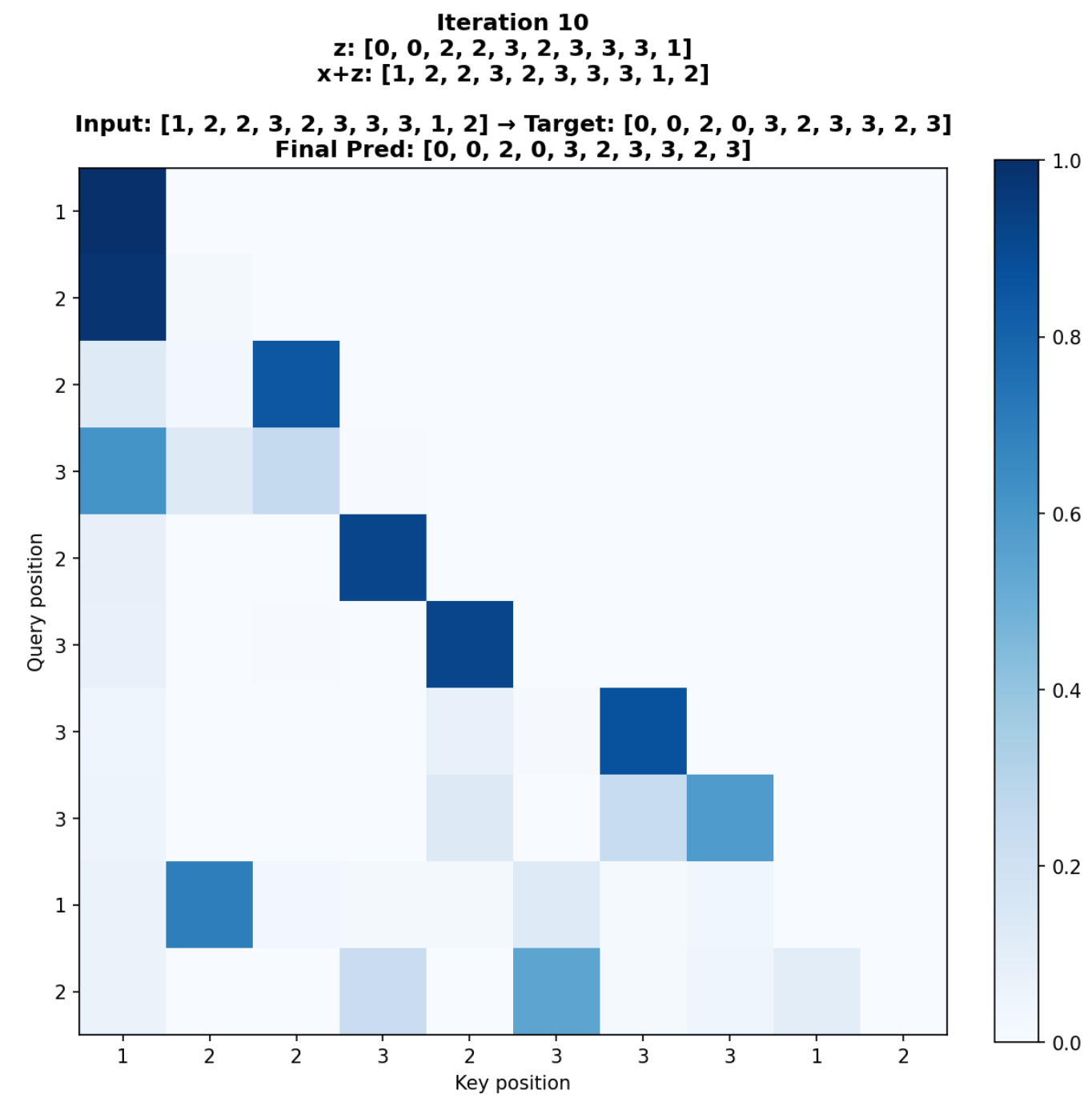
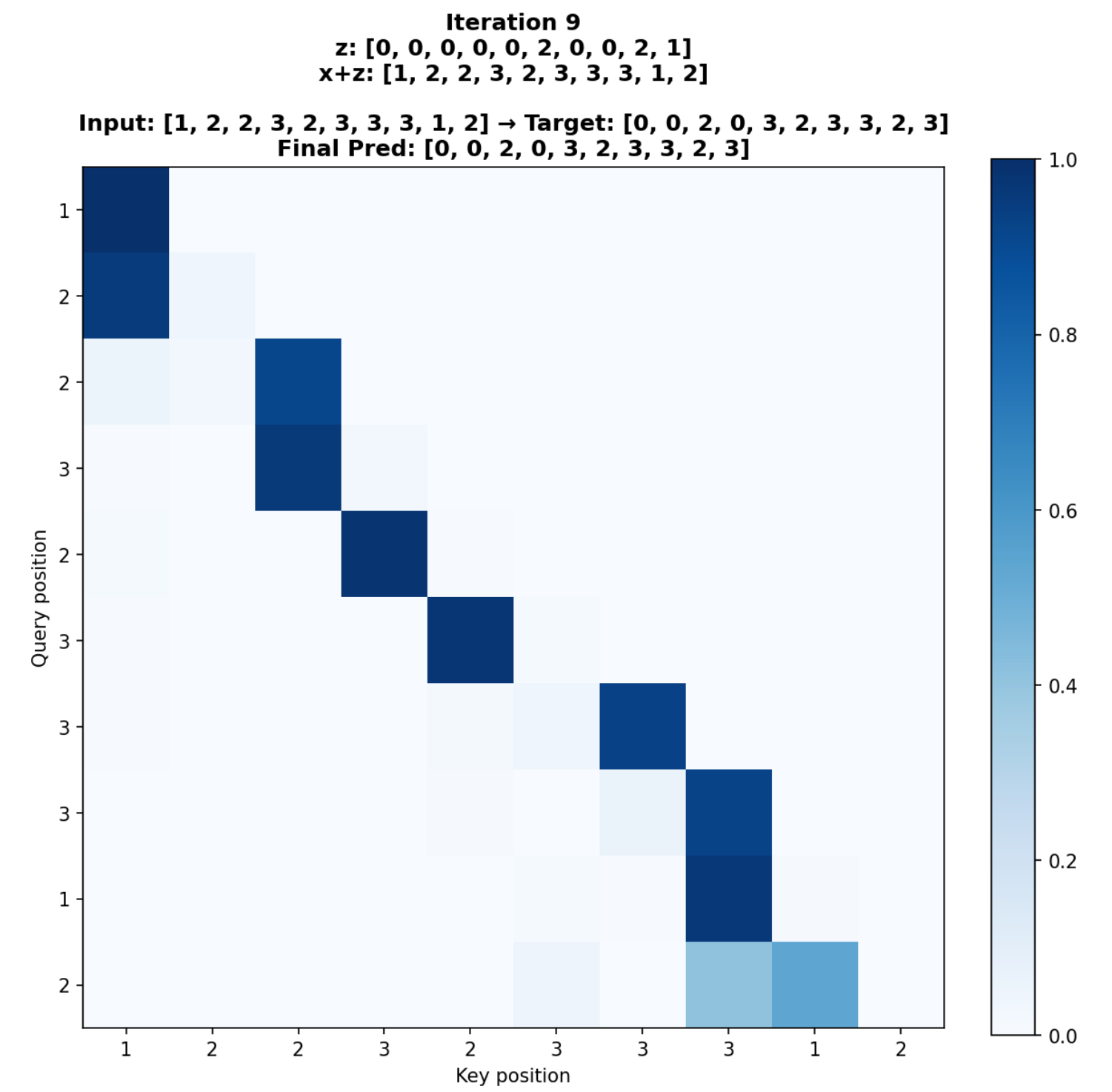
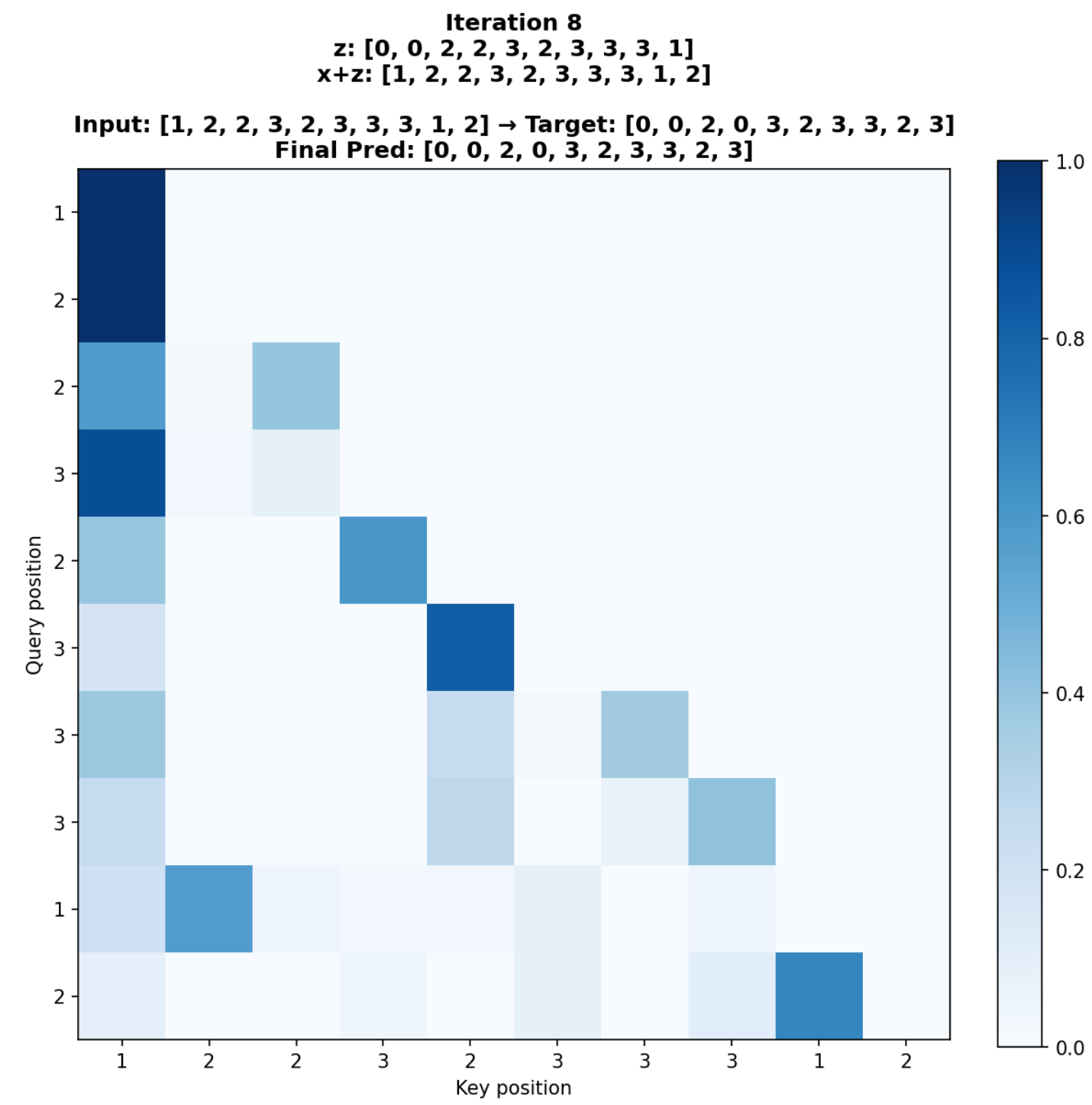
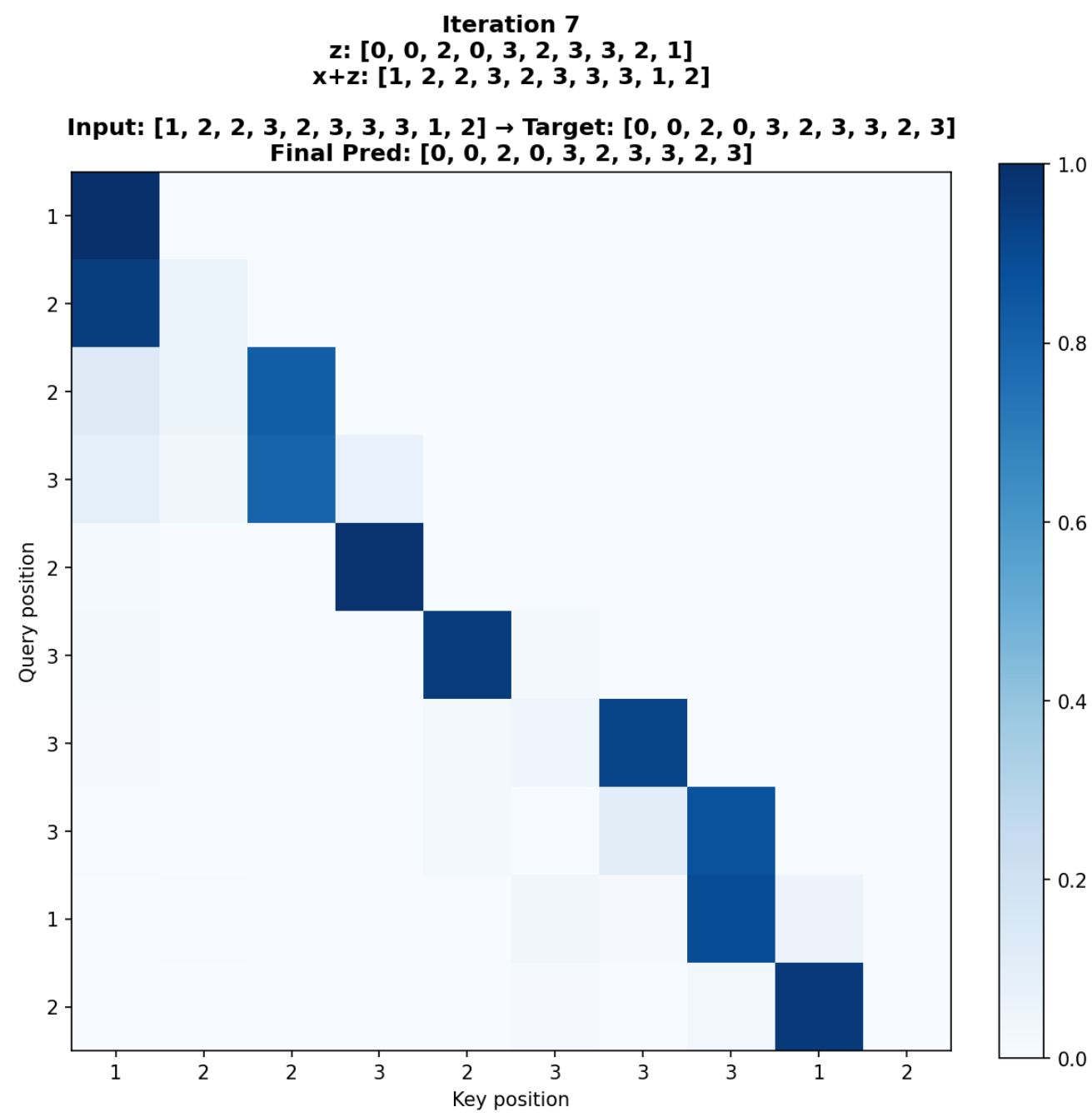
**Input:** [1, 2, 2, 3, 2, 3, 3, 3, 1, 2] → **Target:** [0, 0, 2, 0, 3, 2, 3, 3, 2, 3]  
**Final Pred:** [0, 0, 2, 0, 3, 2, 3, 3, 2, 3]



**Iteration 6**  
z: [0, 0, 2, 0, 3, 2, 3, 3, 3, 1]  
x+z: [1, 2, 2, 3, 2, 3, 3, 3, 1, 2]

**Input:** [1, 2, 2, 3, 2, 3, 3, 3, 1, 2] → **Target:** [0, 0, 2, 0, 3, 2, 3, 3, 2, 3]  
**Final Pred:** [0, 0, 2, 0, 3, 2, 3, 3, 2, 3]

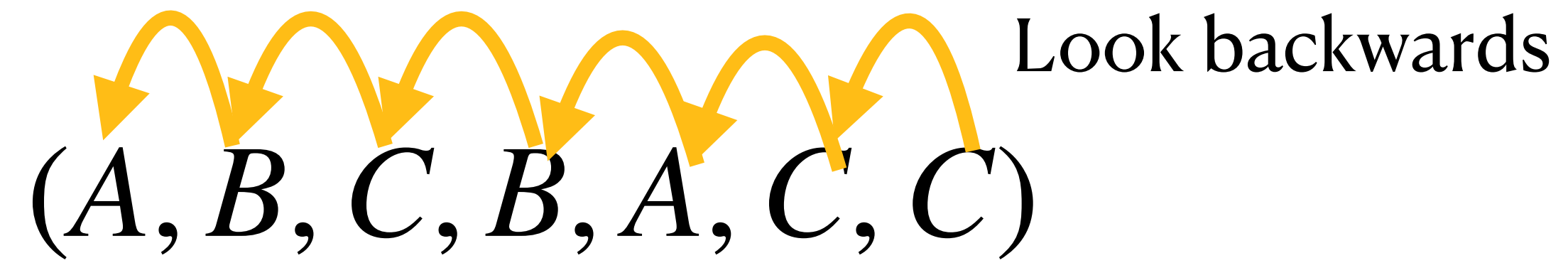




# Copy mechanism heuristics

*(A, B, C, B, A, C, C)*

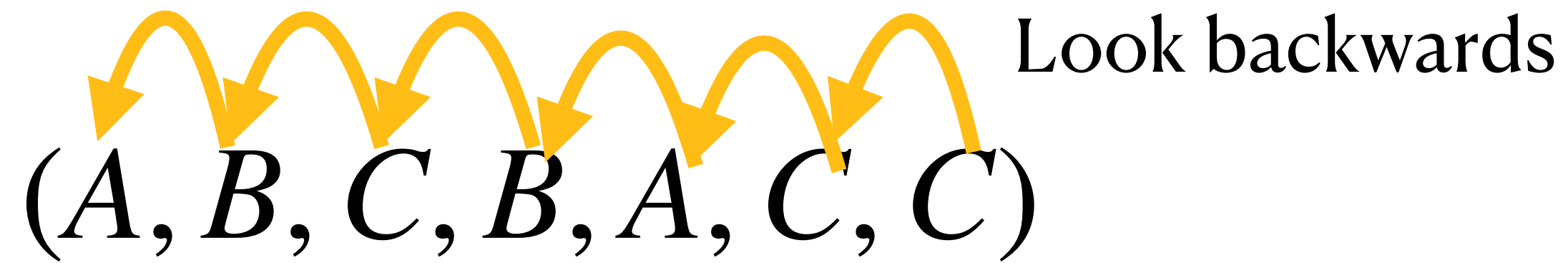
# Copy mechanism heuristics



# Copy mechanism heuristics

Look backwards

$(A, B, C, B, A, C, C)$



Copy the previous token



Previous token  $\begin{pmatrix} \perp, A, B, C, B, A, C \\ A, B, C, B, A, C, C \end{pmatrix}$

# Copy mechanism heuristics

Look backwards  
 $(A, B, C, B, A, C, C)$



Copy the previous token  
↓



Previous token  $\begin{pmatrix} \perp, A, B, C, B, A, C \\ A, B, C, B, A, C, C \end{pmatrix}$

# Copy mechanism heuristics

Look backwards  
(A, B, C, B, A, C, C)

Copy the previous token

Previous token  $\left( \begin{array}{l} \perp, A, B, C, B, A, C \\ A, B, C, B, A, C, C \end{array} \right)$

Key

Query

Look for matching token

# Copy mechanism heuristics

Look backwards  
(A, B, C, B, A, C, C)

Copy the previous token

Previous token  $\left( \begin{array}{l} \perp, A, B, C, B, A, C \\ A, B, C, B, A, C, C \end{array} \right)$   $\begin{array}{l} \longrightarrow \text{Key} \\ \longrightarrow \text{Query} \end{array}$

Look for matching token

$\left( \begin{array}{l} A, B, C, B, A, C, C \\ \dots\dots\dots B \end{array} \right)$

# Copy mechanism heuristics

Look backwards  
(A, B, C, B, A, C, C)

Copy the previous token

Previous token

( ⊥ , A, B, C, B, A, C )  
( A, B, C, B, A, C, C )

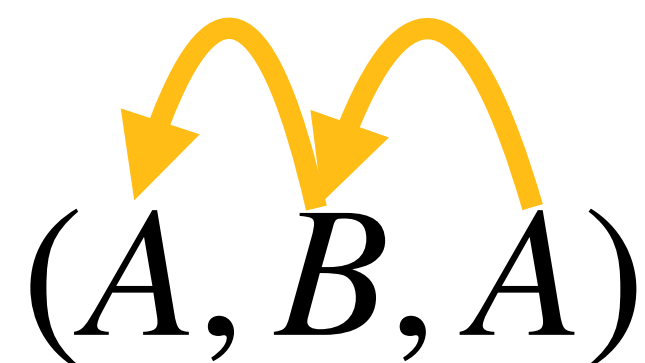
Key  
Query

Look for matching token

( A, B, C, B, A, C, C )  
( ..... B )

# Heuristic vectorial explanation For preceding position

$(A, B, A)$

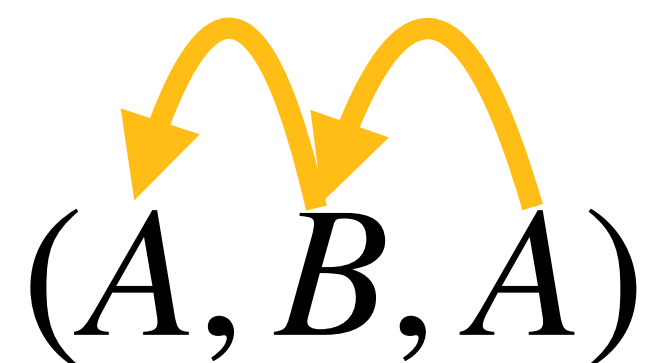


Previous token  $\begin{pmatrix} A, B, A \\ \dots B \end{pmatrix}$

$$\begin{matrix} A \\ \begin{bmatrix} 1 \\ 0 \\ \cos(0) \\ \sin(0) \end{bmatrix} \end{matrix}, \begin{matrix} B \\ \begin{bmatrix} 0 \\ 1 \\ \cos(\frac{\pi}{4}) \\ \sin(\frac{\pi}{4}) \end{bmatrix} \end{matrix}, \begin{matrix} A \\ \begin{bmatrix} 1 \\ 0 \\ \cos(\frac{\pi}{2}) \\ \sin(\frac{\pi}{2}) \end{bmatrix} \end{matrix}$$

# Heuristic vectorial explanation For preceding position

$(A, B, A)$



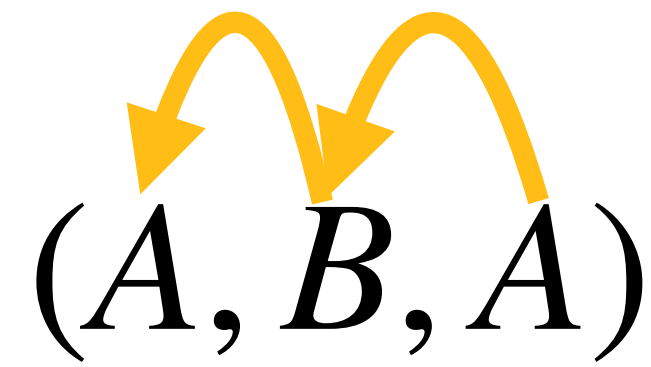
Previous token  $\begin{pmatrix} A, B, A \\ \dots & B \end{pmatrix}$

token information  $\begin{bmatrix} A \\ 1 \\ 0 \end{bmatrix}$ ,  $\begin{bmatrix} B \\ 0 \\ 1 \end{bmatrix}$ ,  $\begin{bmatrix} A \\ 1 \\ 0 \end{bmatrix}$

position information  $\begin{bmatrix} \cos(0) \\ \sin(0) \end{bmatrix}$ ,  $\begin{bmatrix} \cos(\frac{\pi}{4}) \\ \sin(\frac{\pi}{4}) \end{bmatrix}$ ,  $\begin{bmatrix} \cos(\frac{\pi}{2}) \\ \sin(\frac{\pi}{2}) \end{bmatrix}$

# Heuristic vectorial explanation For preceding position

$(A, B, A)$



$$\vec{V} = \text{token information} \begin{bmatrix} A \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} B \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} A \\ 1 \\ 0 \end{bmatrix}$$

$$\vec{Q} = \vec{K} = \text{position information} \begin{bmatrix} \cos(0) \\ \sin(0) \end{bmatrix}, \begin{bmatrix} \cos(\frac{\pi}{4}) \\ \sin(\frac{\pi}{4}) \end{bmatrix}, \begin{bmatrix} \cos(\frac{\pi}{2}) \\ \sin(\frac{\pi}{2}) \end{bmatrix}$$



Previous token  $\begin{pmatrix} A, B, A \\ \dots & B \end{pmatrix}$

# Heuristic vectorial explanation For preceding position

$(A, B, A)$

$\vec{V} =$  token information  
 $\vec{Q} = \vec{K} =$  position information

$$\begin{matrix} A & B & A \\ \begin{bmatrix} 1 \\ 0 \\ \cos(0) \\ \sin(0) \end{bmatrix}, & \begin{bmatrix} 0 \\ 1 \\ \cos(\frac{\pi}{4}) \\ \sin(\frac{\pi}{4}) \end{bmatrix}, & \begin{bmatrix} 1 \\ 0 \\ \cos(\frac{\pi}{2}) \\ \sin(\frac{\pi}{2}) \end{bmatrix} \end{matrix}$$



Previous token  $\begin{pmatrix} A, B, A \\ \dots B \end{pmatrix}$

# Heuristic vectorial explanation For preceding position

$(A, B, A)$

$\vec{V} =$  token information  
 $\vec{Q} = \vec{K} =$  position information

$$\begin{matrix} A & B & A \\ \begin{bmatrix} 1 \\ 0 \\ \cos(0) \\ \sin(0) \end{bmatrix}, & \begin{bmatrix} 0 \\ 1 \\ \cos(\frac{\pi}{4}) \\ \sin(\frac{\pi}{4}) \end{bmatrix}, & \begin{bmatrix} 1 \\ 0 \\ \cos(\frac{\pi}{2}) \\ \sin(\frac{\pi}{2}) \end{bmatrix} \end{matrix}$$

Previous token  $\begin{pmatrix} A, B, A \\ \dots B \end{pmatrix}$

$$\begin{matrix} \begin{bmatrix} 0 \\ 0 \\ \cos(0) \\ \sin(0) \end{bmatrix}, & \begin{bmatrix} 1 \\ 0 \\ \cos(\frac{\pi}{4}) \\ \sin(\frac{\pi}{4}) \end{bmatrix}, & \begin{bmatrix} 0 \\ 1 \\ \cos(\frac{\pi}{2}) \\ \sin(\frac{\pi}{2}) \end{bmatrix} \end{matrix}$$

# Heuristic vectorial explanation for matching tokens



$$\begin{matrix} A \\ \begin{bmatrix} 1 \\ 0 \\ \cos(0) \\ \sin(0) \end{bmatrix} \end{matrix}, \begin{matrix} B \\ \begin{bmatrix} 0 \\ 1 \\ \cos(\frac{\pi}{4}) \\ \sin(\frac{\pi}{4}) \end{bmatrix} \end{matrix}, \begin{matrix} A \\ \begin{bmatrix} 1 \\ 0 \\ \cos(\frac{\pi}{2}) \\ \sin(\frac{\pi}{2}) \end{bmatrix} \end{matrix}$$

# Heuristic vectorial explanation for matching tokens



$$\begin{array}{l} \text{token information} \\ \text{position information} \end{array} \begin{array}{c} A \\ \left[ \begin{array}{c} 1 \\ 0 \\ \cos(0) \\ \sin(0) \end{array} \right] \end{array}, \begin{array}{c} B \\ \left[ \begin{array}{c} 0 \\ 1 \\ \cos(\frac{\pi}{4}) \\ \sin(\frac{\pi}{4}) \end{array} \right] \end{array}, \begin{array}{c} A \\ \left[ \begin{array}{c} 1 \\ 0 \\ \cos(\frac{\pi}{2}) \\ \sin(\frac{\pi}{2}) \end{array} \right] \end{array}$$

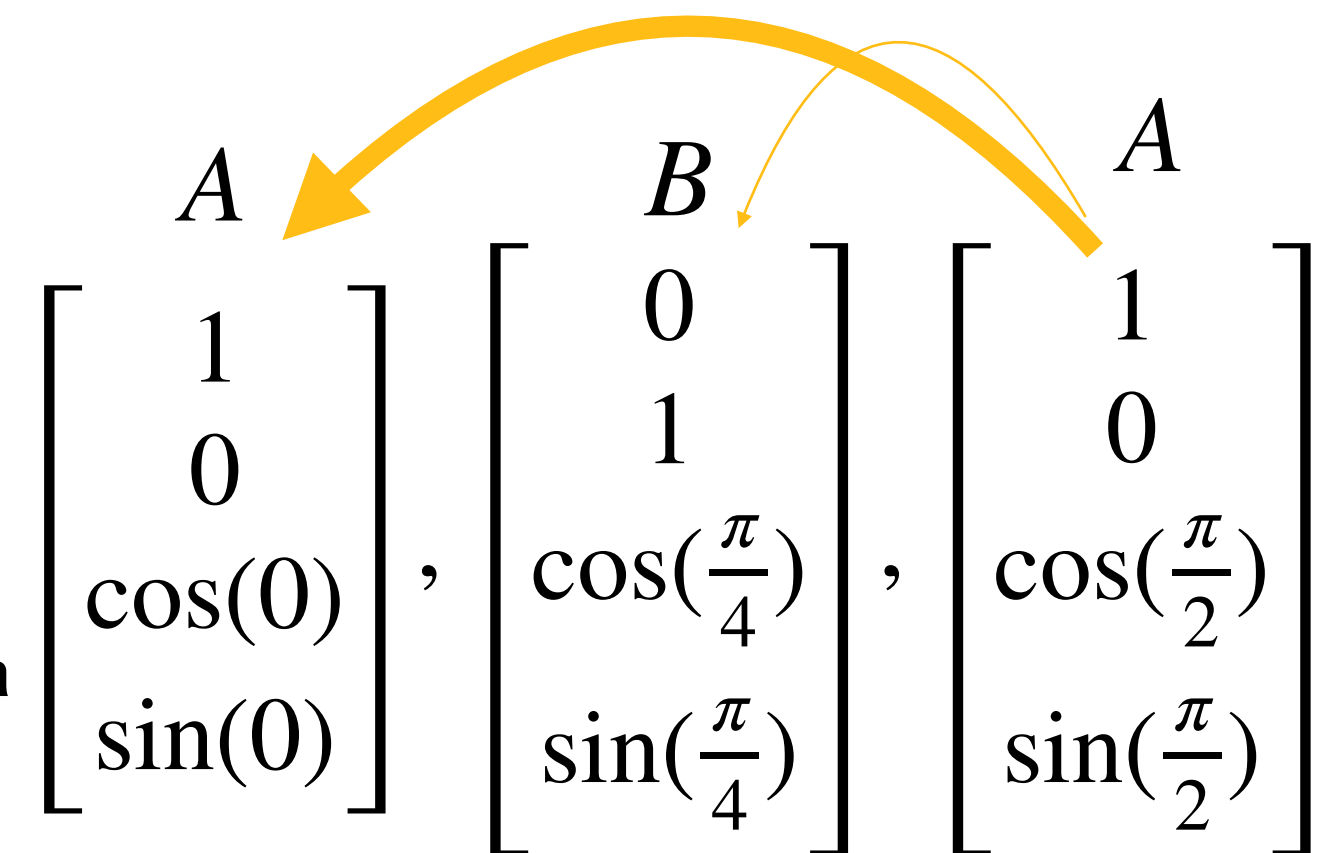
# Heuristic vectorial explanation for matching tokens

$(A, C, A, B, A)$



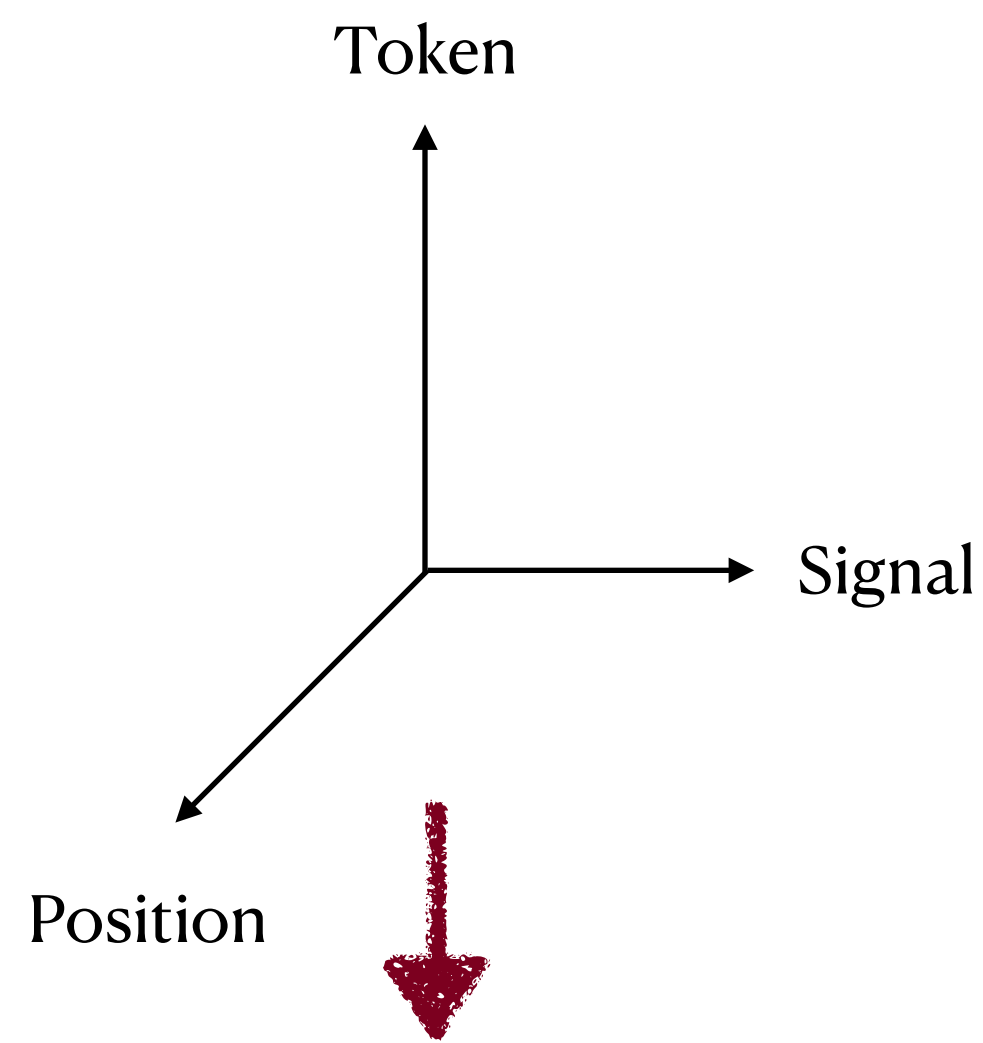
$\vec{Q} = \vec{K} =$  token information

position information

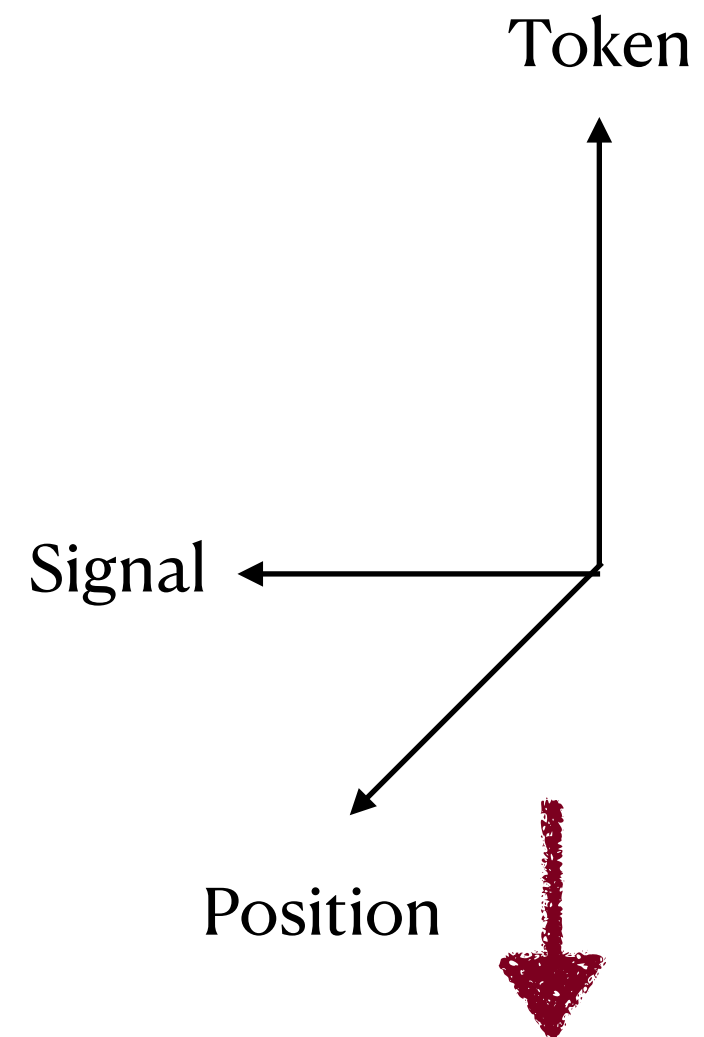
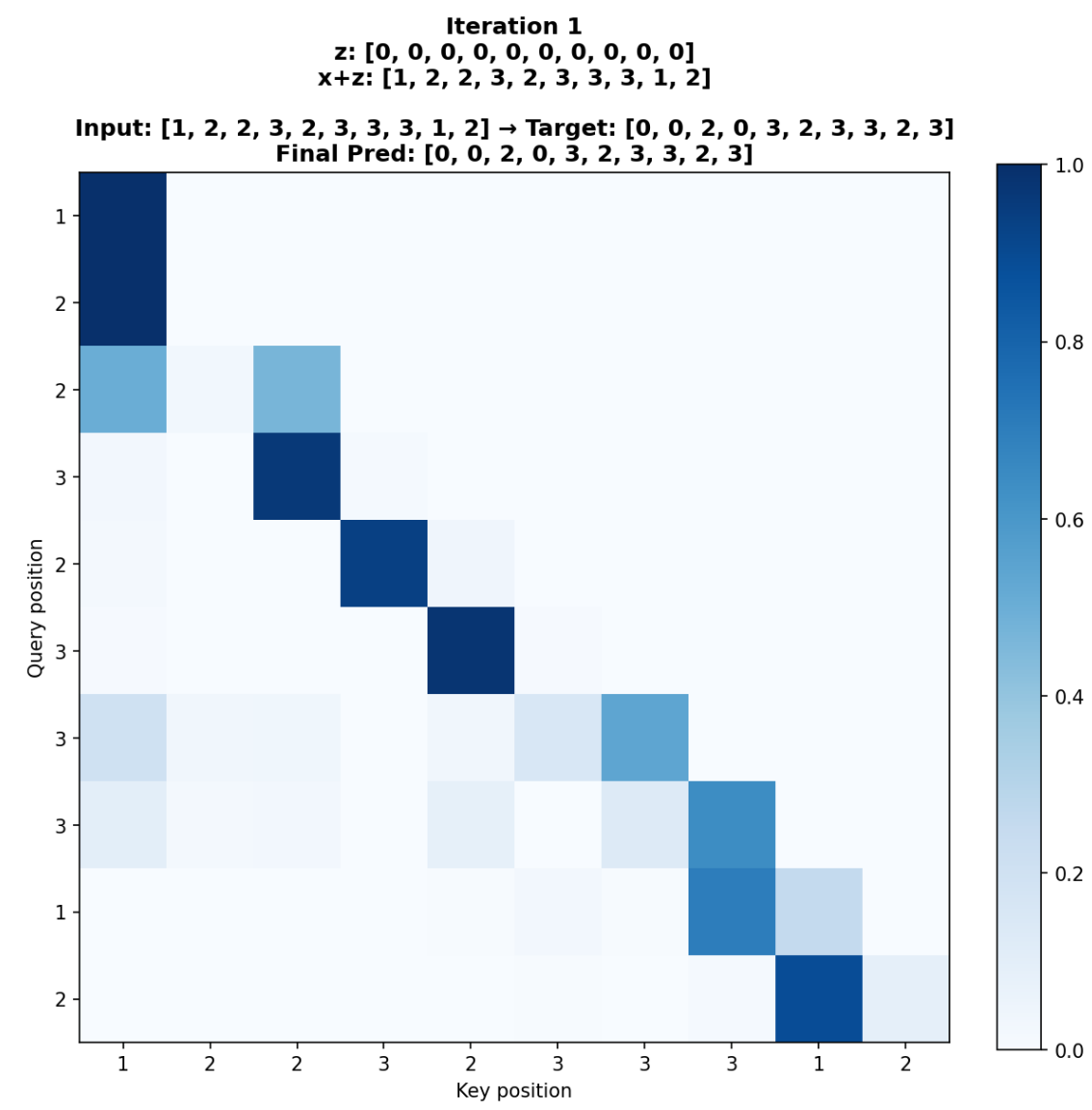

$$\begin{bmatrix} 1 \\ 0 \\ \cos(0) \\ \sin(0) \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ \cos(\frac{\pi}{4}) \\ \sin(\frac{\pi}{4}) \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ \cos(\frac{\pi}{2}) \\ \sin(\frac{\pi}{2}) \end{bmatrix}$$



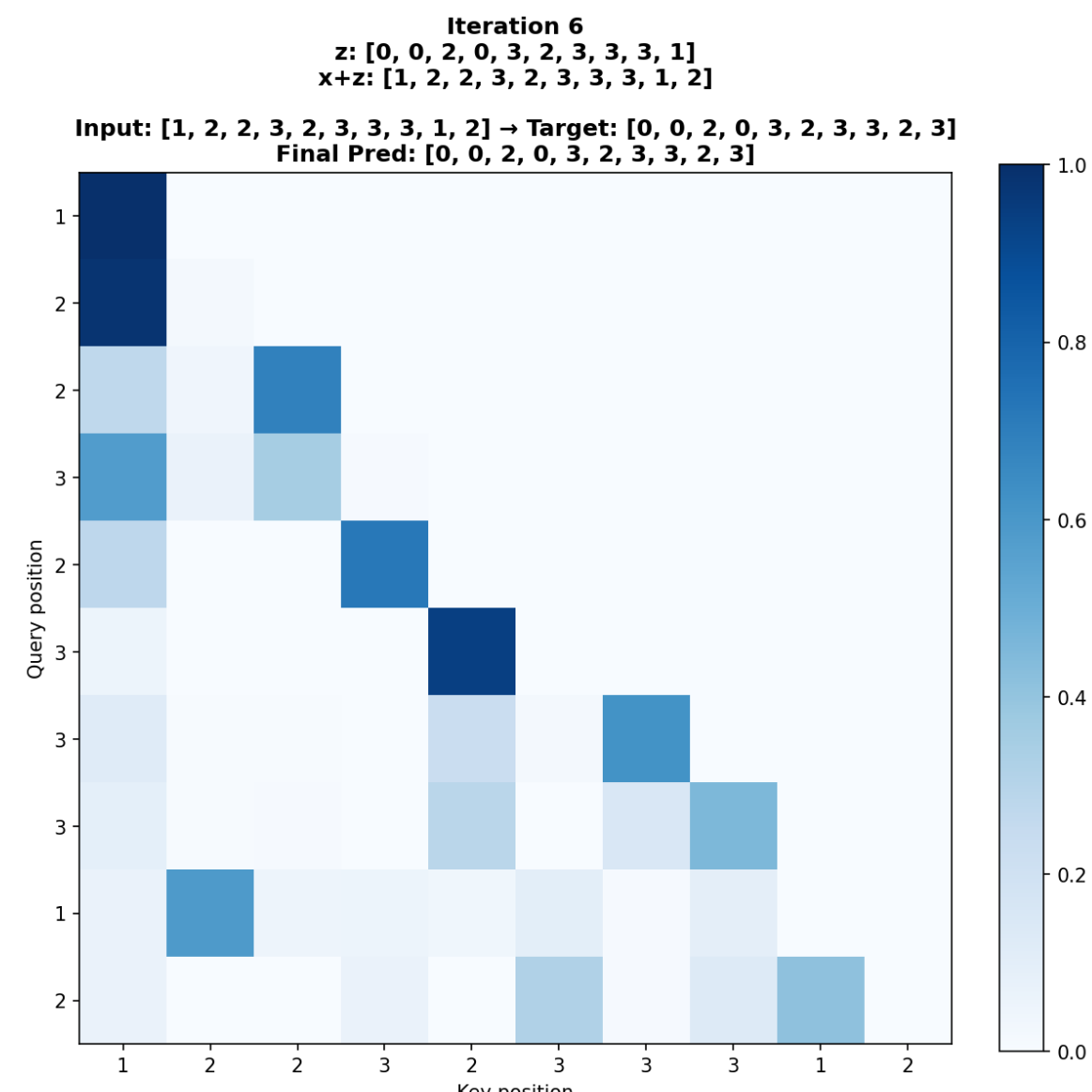
# Transformation of the latent $z$



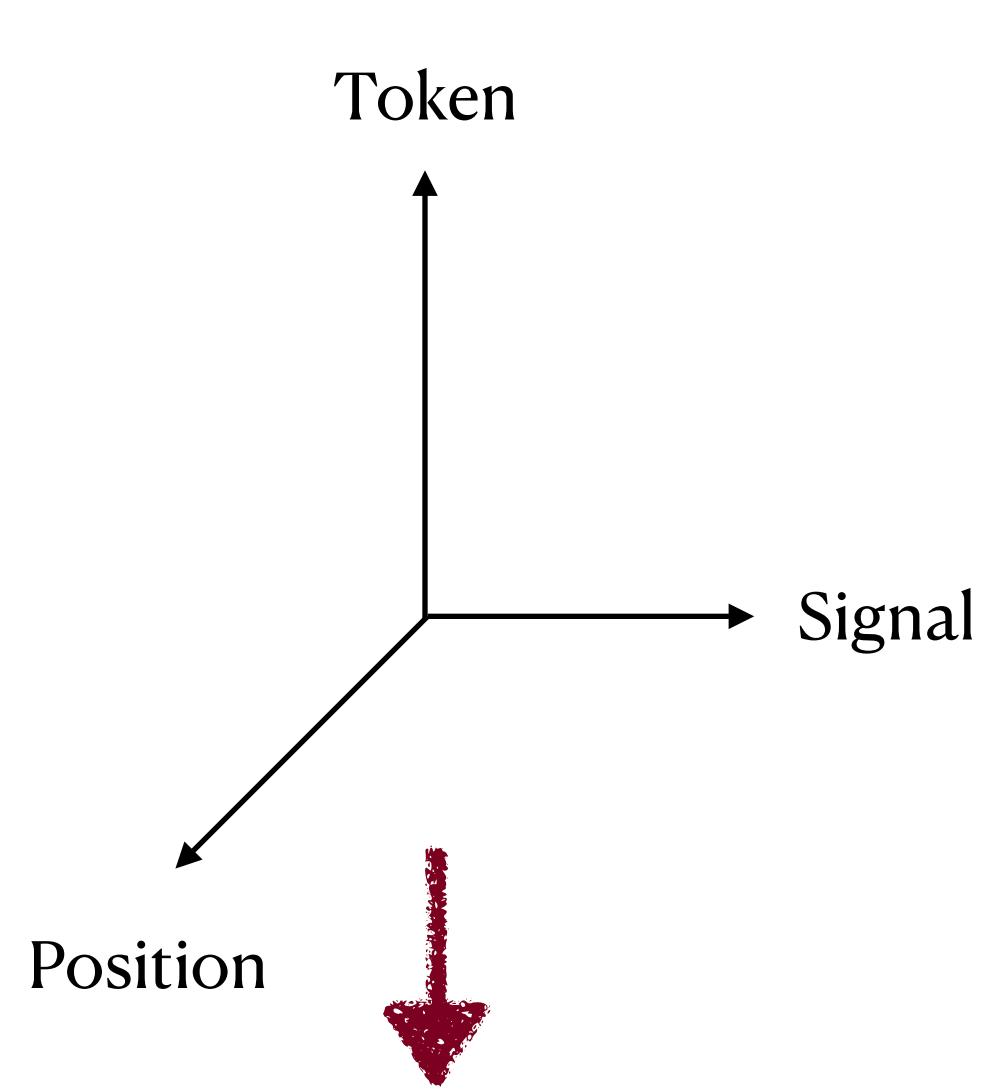
Look at previous token and copy it



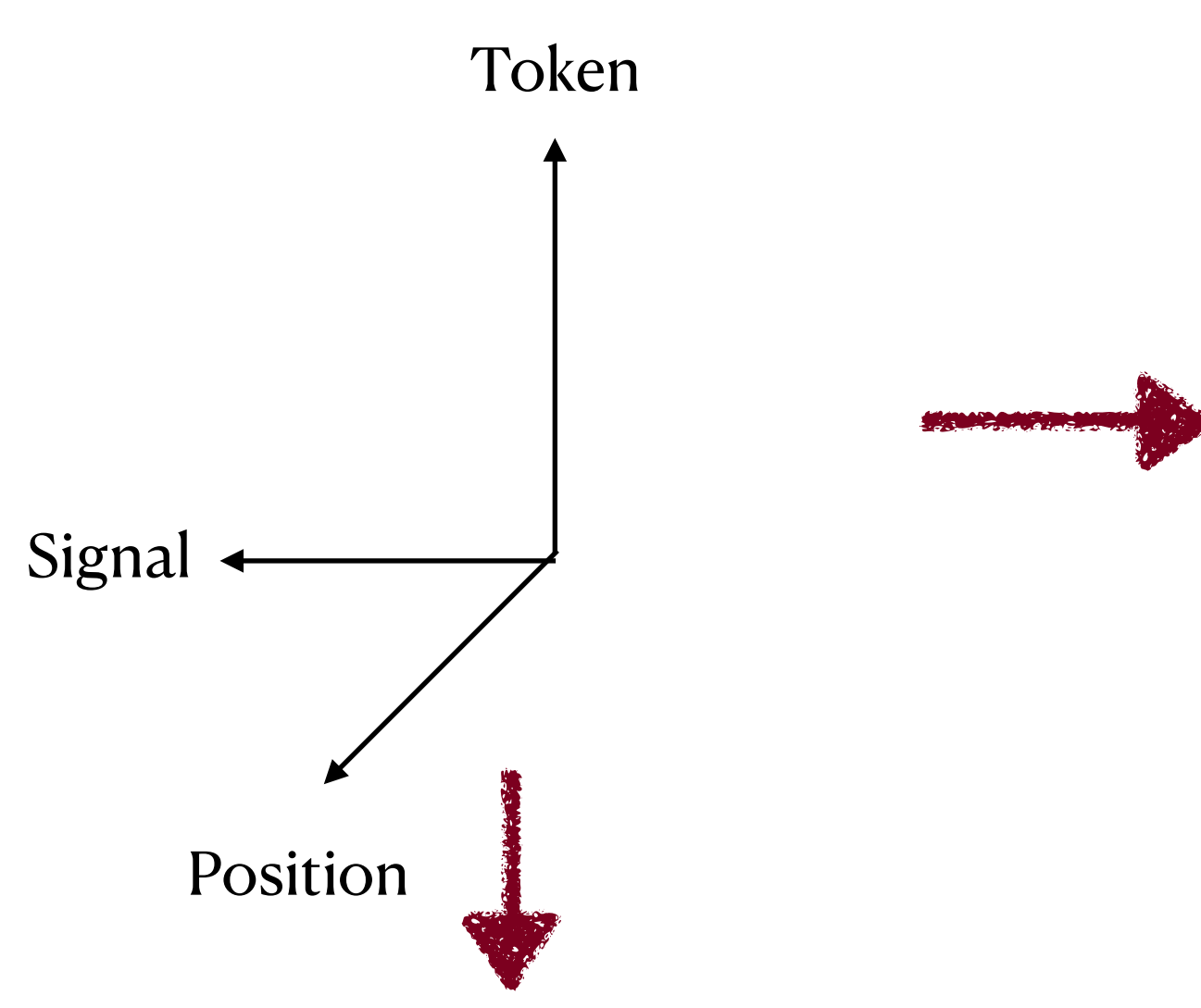
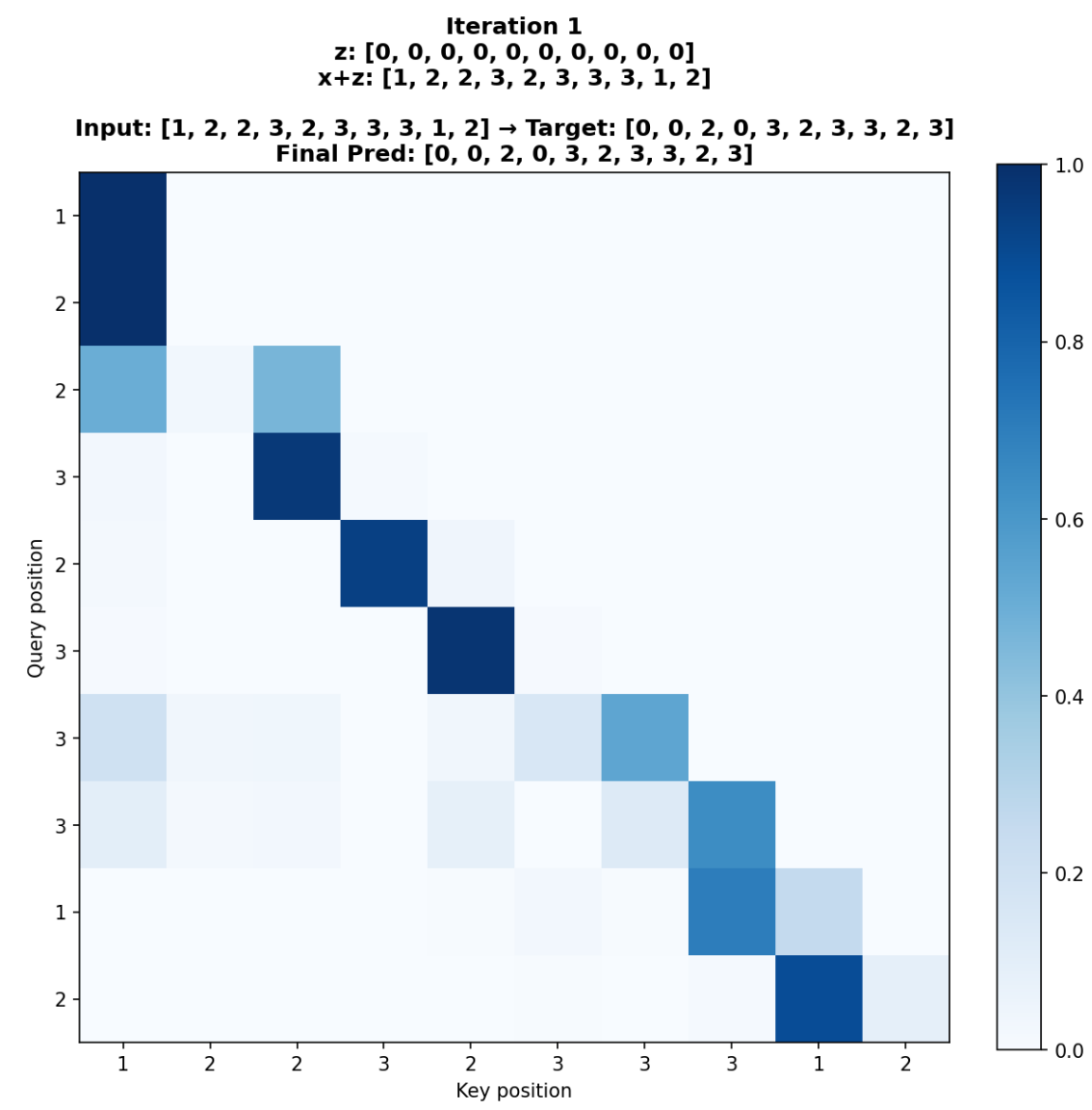
Look at previous occurrence (induction)



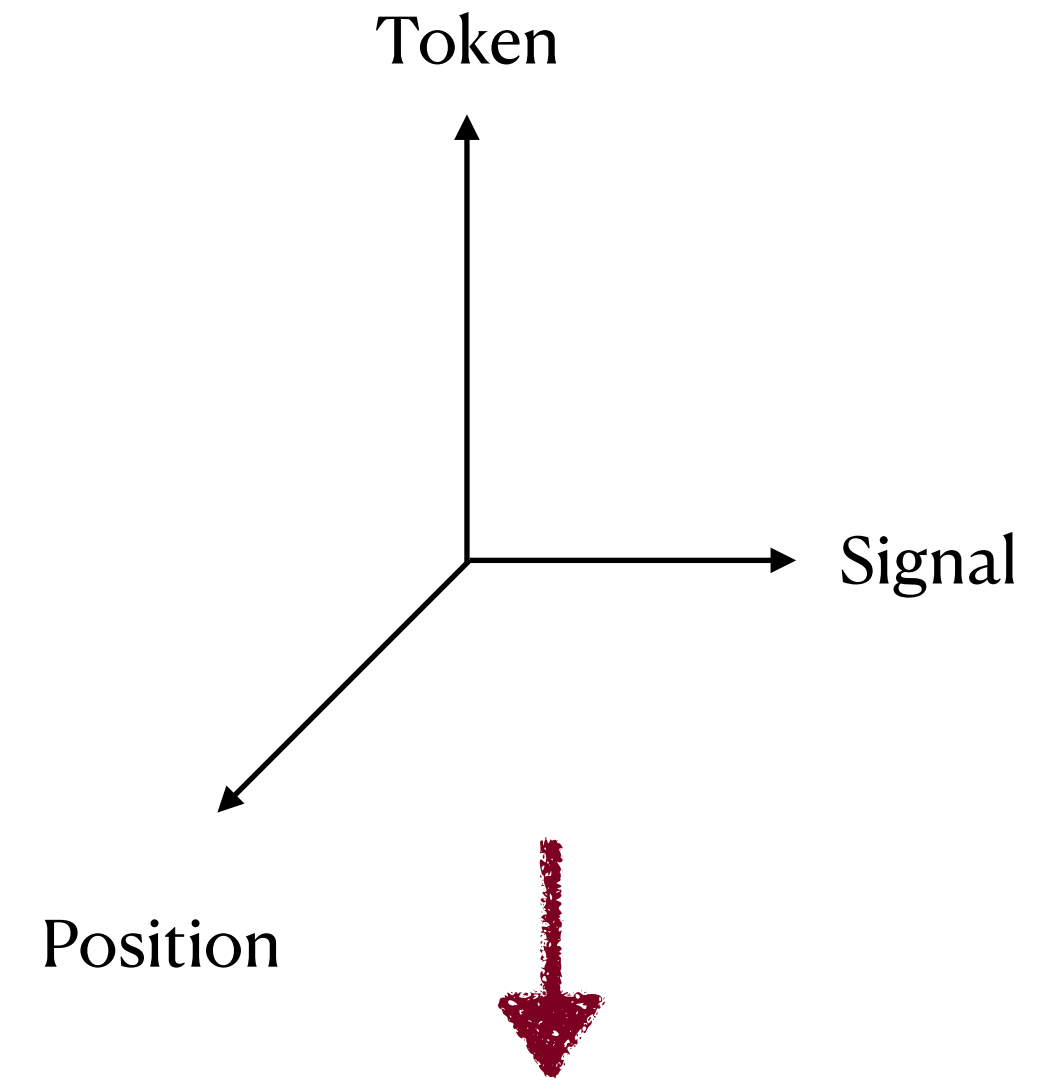
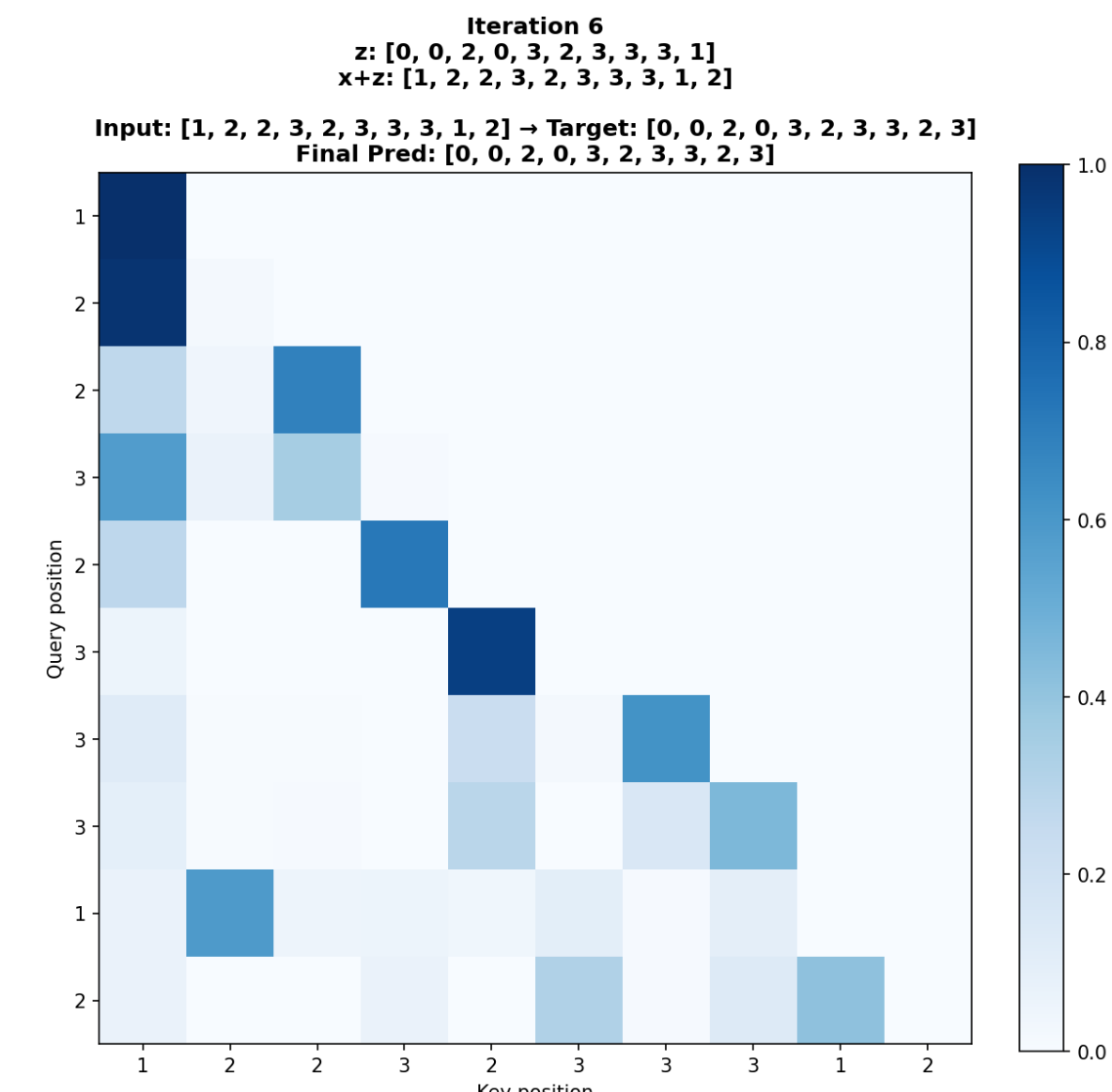
# Transformation of the latent $z$



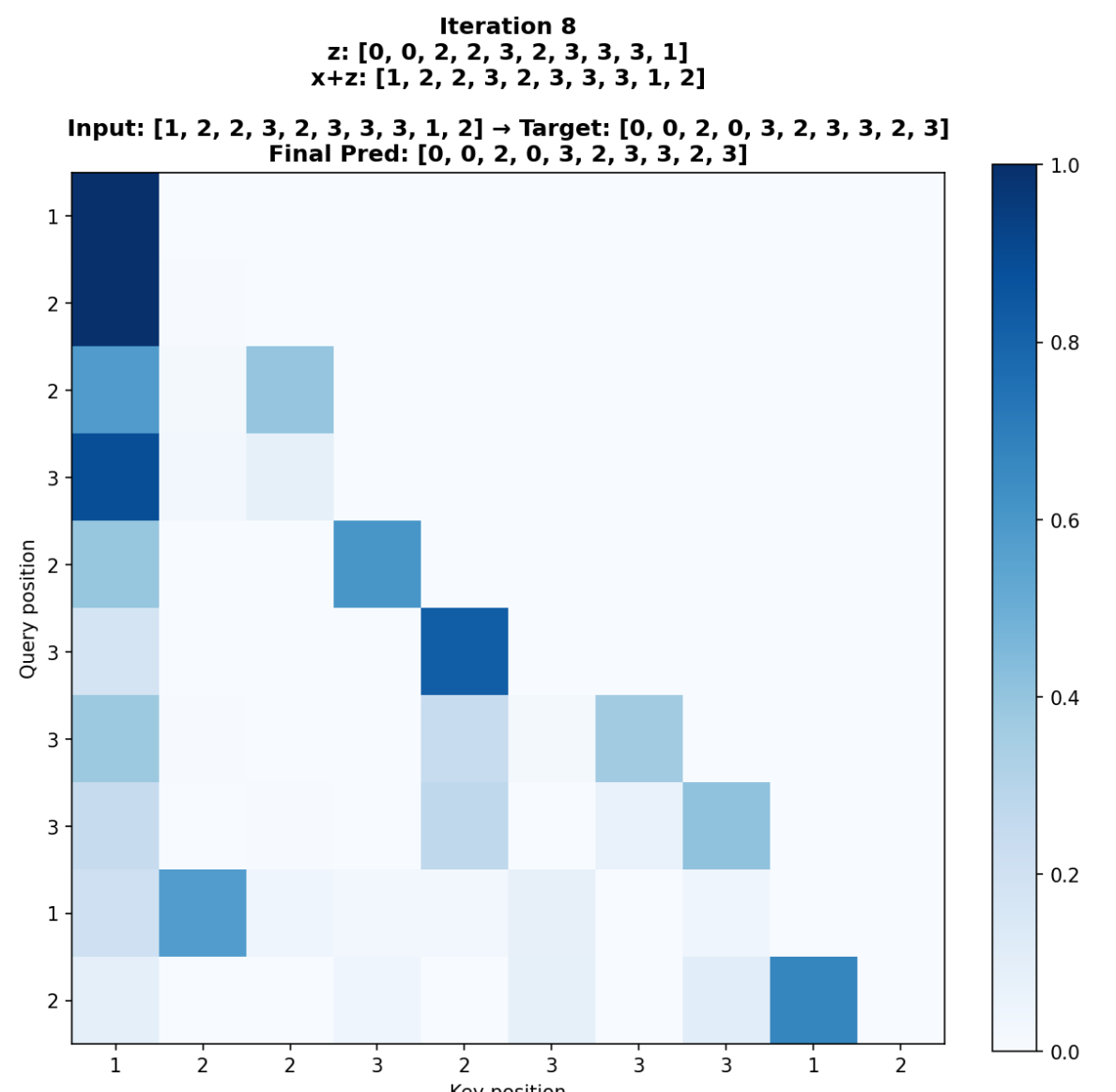
Look at previous token and copy it



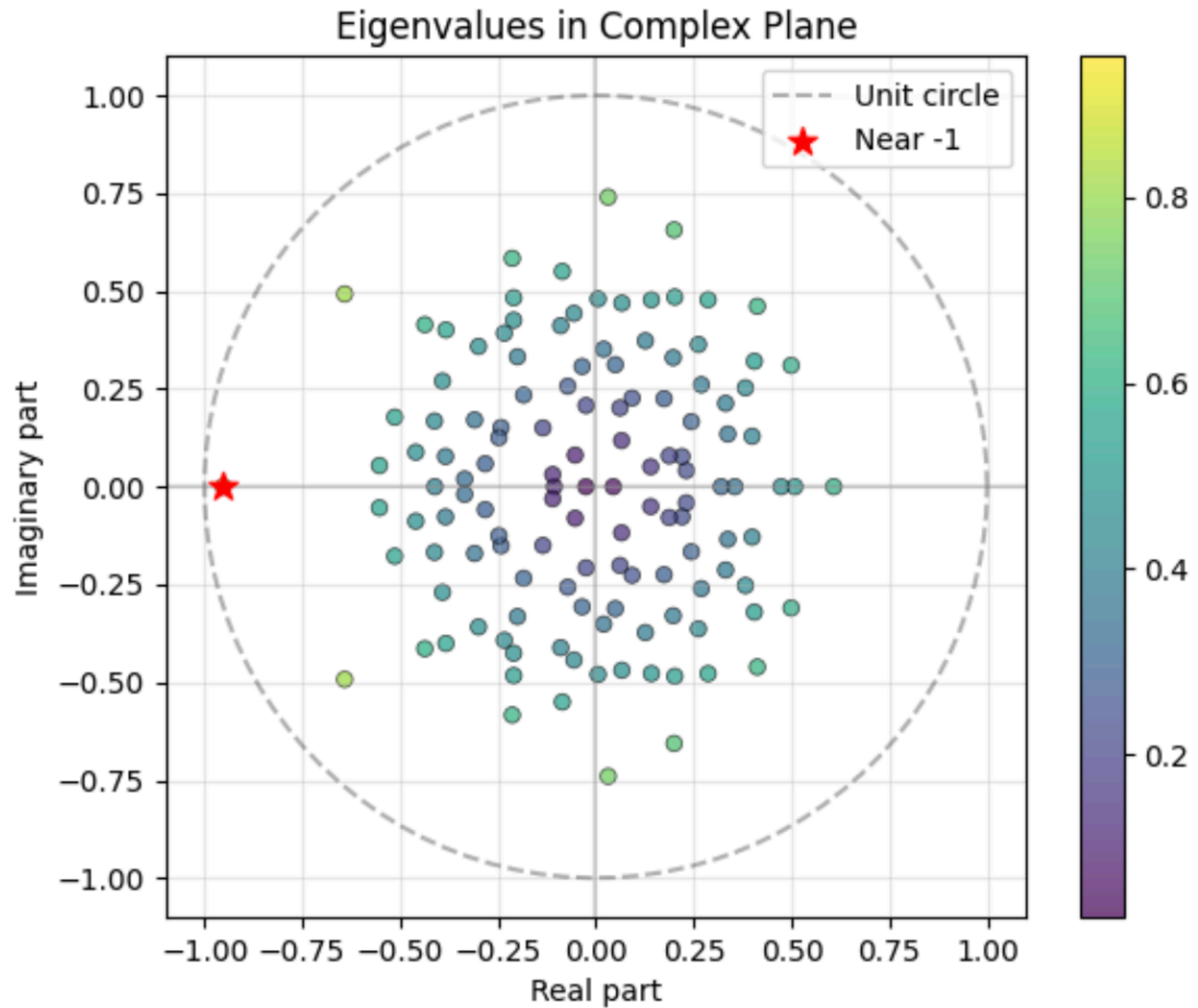
Look at previous occurrence (induction)



Look at previous token and copy it

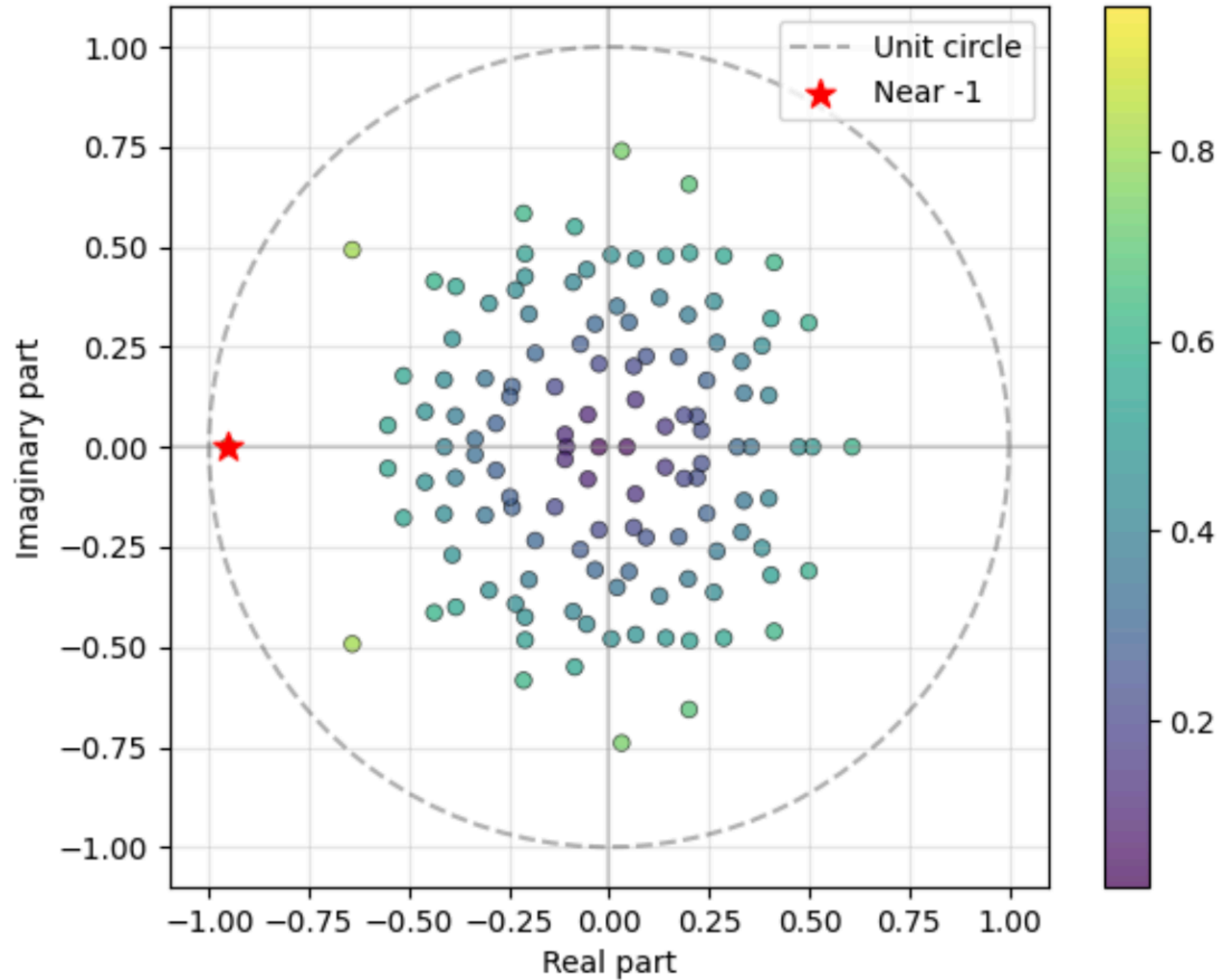


# Feed forward layer

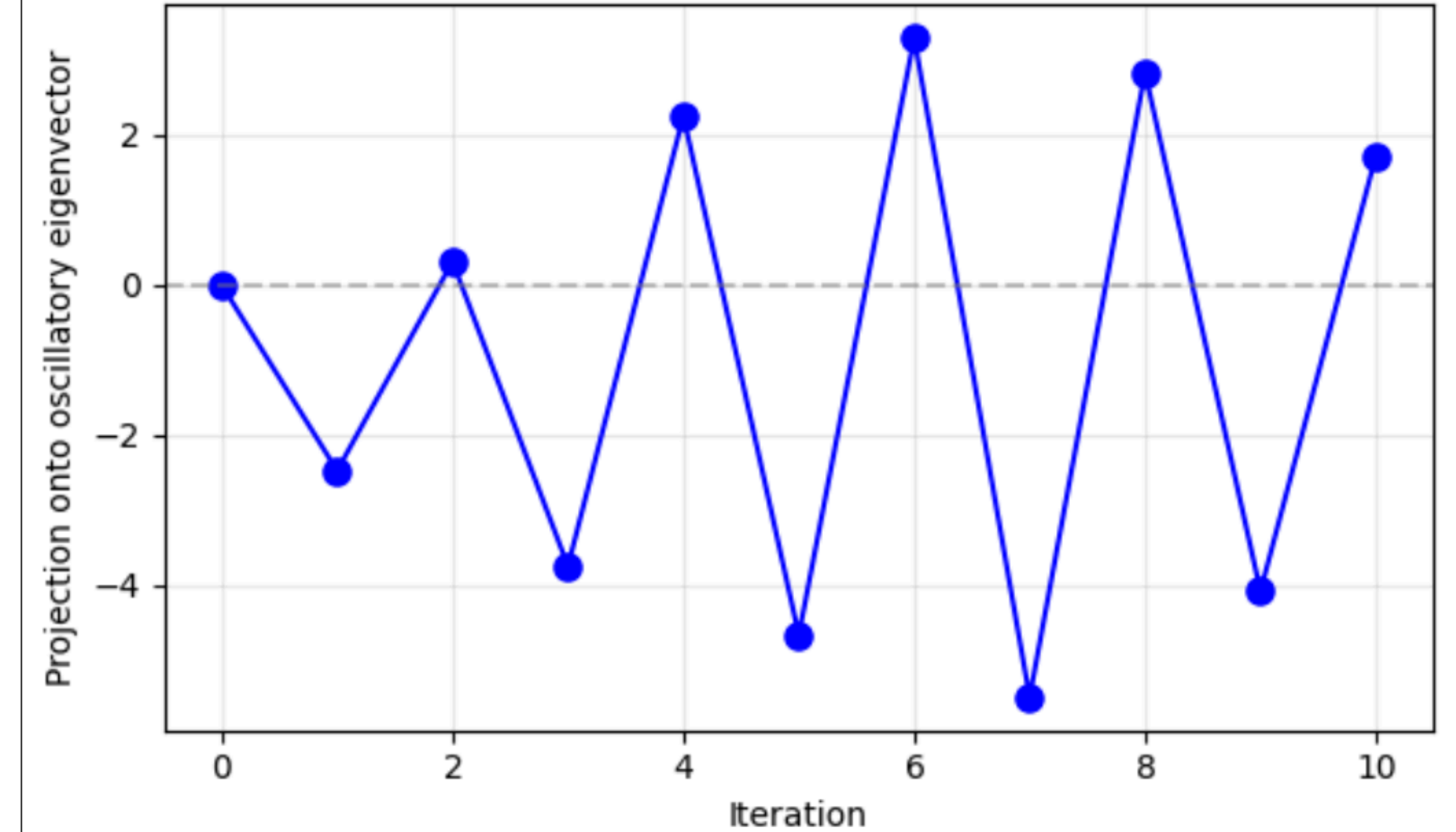


# Feed forward layer

Eigenvalues in Complex Plane



z projection onto eigenvector ( $\lambda = -0.952 + 0.000j$ )



# Conclusion

- Loop transformers come up with a clever way to solve the induction head task.

# Ongoing work

Conjecture: Looped transformers generalize better and are more sample efficient than vanilla transformers.

Idea: Look at loss landscapes of

$$\sigma(W_2\sigma(W_1X + b_1) + b_2) \quad \text{vs.} \quad \sigma(W\sigma(WX + b) + b)$$

Vanilla                      Looped

- Need to make the role of the latent vector more rigorous.
- Formalize some ARC-AGI tasks and demonstrate how loop transformers solve them.

← **Post**



**François Chollet**

@fchollet



The Transformer architecture is fundamentally a parallel processor of context, but reasoning is a sequential, iterative process.

To solve complex problems, a model needs a "scratchpad" not just in its output CoT, but in its internal state. A differentiable way to loop, branch, and backtrack until the model finds a solution that works.

11:49 AM · Dec 23, 2025 · **129.2K** Views



144



230



1.9K



875



**Read 144 replies**

**Thank you**