

Attention Matrix Properties

Setup: 4x4 Attention Matrix

Consider a sentence with 4 tokens: **[Token₁, Token₂, Token₃, Token₄]**

The attention matrix **A** has shape (4, 4):

	K ₁	K ₂	K ₃	K ₄
Q ₁	A _{1,1}	A _{1,2}	A _{1,3}	A _{1,4}
Q ₂	A _{2,1}	A _{2,2}	A _{2,3}	A _{2,4}
Q ₃	A _{3,1}	A _{3,2}	A _{3,3}	A _{3,4}
Q ₄	A _{4,1}	A _{4,2}	A _{4,3}	A _{4,4}

Question 1: What does element A₂₃ represent?

Answer: A₂₃ represents the **attention weight** that **Token₂** (query) pays to **Token₃** (key).

Interpretation:

- **Row index (2)**: The token that is "paying attention" (the query position)
- **Column index (3)**: The token that is "being attended to" (the key position)
- **A₂₃**: How much Token₂ focuses on Token₃ when computing its representation

Reading the Matrix:

- **Rows**: Each row shows where one token attends to ALL other tokens
- **Columns**: Each column shows which tokens attend to a specific token
- **A₂₃**: Element at row 2, column 3

Question 2: If $A_{23} = 0.8$, what does it mean?

Answer: Token₂ is **strongly attending** to Token₃.

Specific Meaning:

1. **High attention weight:** 0.8 out of maximum 1.0 is very high
2. **Token₂'s representation** will be heavily influenced by Token₃'s value vector
3. **80% of the "focus"** from Token₂ goes to Token₃

Context (Remember: Rows Sum to 1)

Since softmax ensures row sums equal 1:

$$A_{21} + A_{22} + A_{23} + A_{24} = 1.0$$

If $A_{23} = 0.8$, then:

$$A_{21} + A_{22} + A_{24} = 0.2$$

This means:

- Token₂ pays **80% attention** to Token₃
- Token₂ pays only **20% combined** to Token₁, Token₂ (itself), and Token₄
- Token₃ is **dominant** in determining Token₂'s output representation

Example

Sentence: "The cat chased mice"

- Token₁ = "The"
- Token₂ = "cat"
- Token₃ = "chased"
- Token₄ = "mice"

If $A_{23} = 0.8$:

- "cat" strongly attends to "chased"
- When computing the representation for "cat", the model heavily incorporates information from "chased"

- This captures the subject-verb relationship
- The output for "cat" will strongly reflect what action it's performing

Summary

Element	Query (Row)	Key (Column)	Meaning
A_{23}	Token_2	Token_3	How much Token_2 attends to Token_3
$A_{23} = 0.8$	"cat"	"chased"	80% of "cat"'s representation comes from "chased"

Key insight: High attention weight (0.8) means strong token interaction - the query token's output will be dominated by the key token's value.