

MECHANISMS OF ATTENTION ATTENTION & CONTEXTUAL EMBEDDINGS

Large Language Models are built on Transformers. A Transformer is a specific network architecture fundamentally based on the mechanism of Attention.

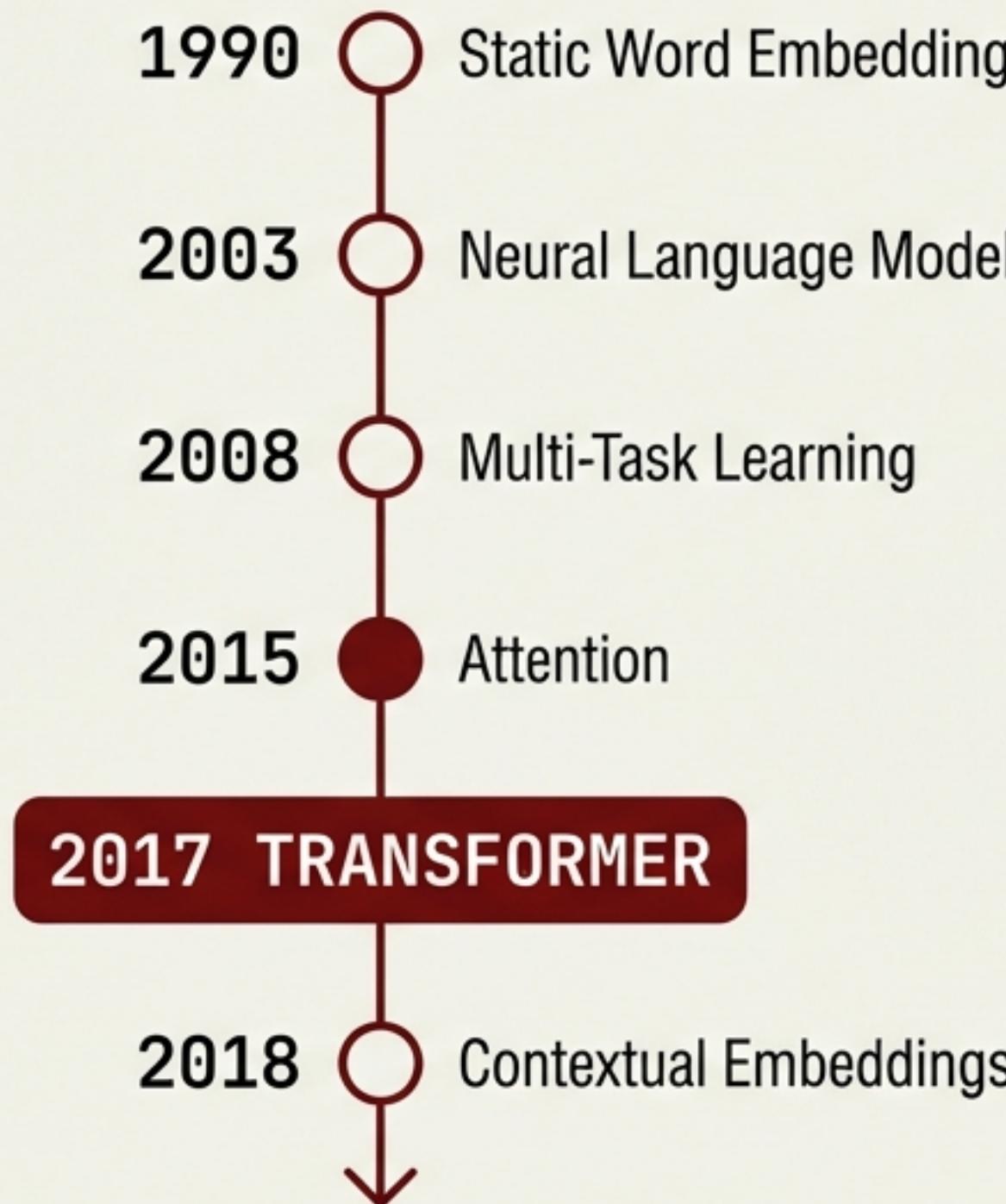
SEMINAL PAPER CITATION:

"Attention Is All You Need" (2017).

Ashish Vaswani Noam Shazeer Niki Parmar
Jakob Uszkoreit Llion Jones Aidan N. Gomez
Lukasz Kaiser Illia Polosukhin

Google Brain | Google Research | U of Toronto

THE TIMELINE OF NLP ACCELERATION



For decades, language models relied on static definitions. The introduction of Attention Attention (**2015**) and **Transformers (2017)** marked the shift toward dynamic, context-aware understanding.

THE PROBLEM: STATIC EMBEDDINGS

Systems like **word2vec** created static maps. The embedding for a word was fixed. It did not reflect how meaning shifts based on the sentence.

Dictionary (Fixed)

"it": [0.2, -0.4, 0.9]

(Immutable Vector)

THE AMBIGUITY:

The chicken didn't cross the road because...



Does "it" mean chicken?

Does "it" mean road?

A static vector cannot know.

AMBIGUITY REQUIRES LOOKING AT NEIGHBORS

The identity of “it” is defined entirely by the context.

CASE A: Referring to the animal

The chicken didn’t cross the road because [it] was tired.

(Context defines ‘it’ as chicken)

CASE B: Referring to the street

The chicken didn’t cross the road because [it] was wide.

(Context defines “it” as road)

SOLUTION: CONTEXTUAL EMBEDDINGS

INTUITION

A word's meaning representation should differ across different contexts.

DEFINITION

Each word possesses a unique vector that expresses specific meaning derived from surrounding words.

How do we compute this?



ATTENTION

The mechanism that allows a word to "look" at its neighbors before defining itself.

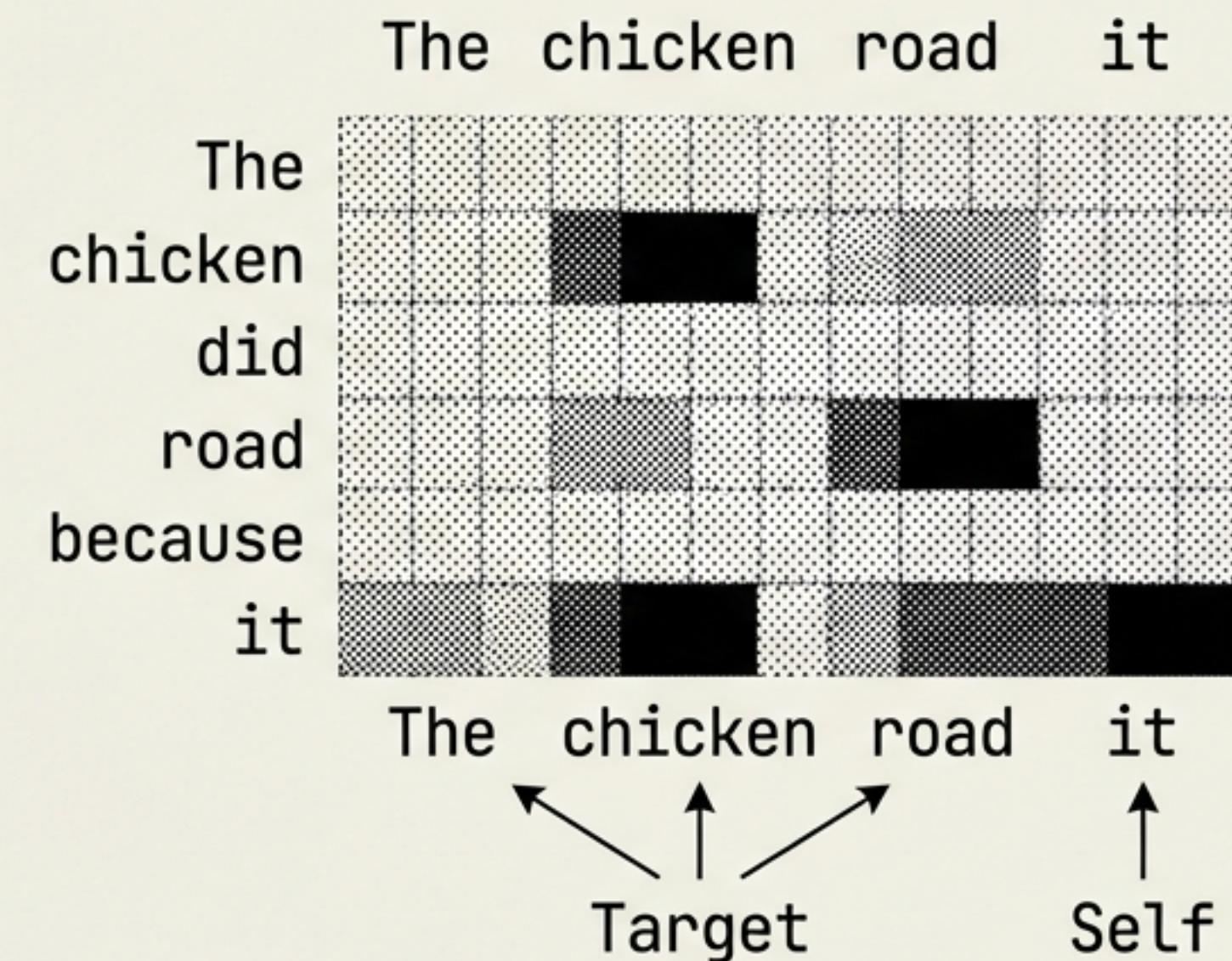
VISUALIZING ATTENTION

VISUALIZING ATTENTION

To build a contextual embedding, a word ‘attends to’ specific neighbors more than others. Selective integration of information.

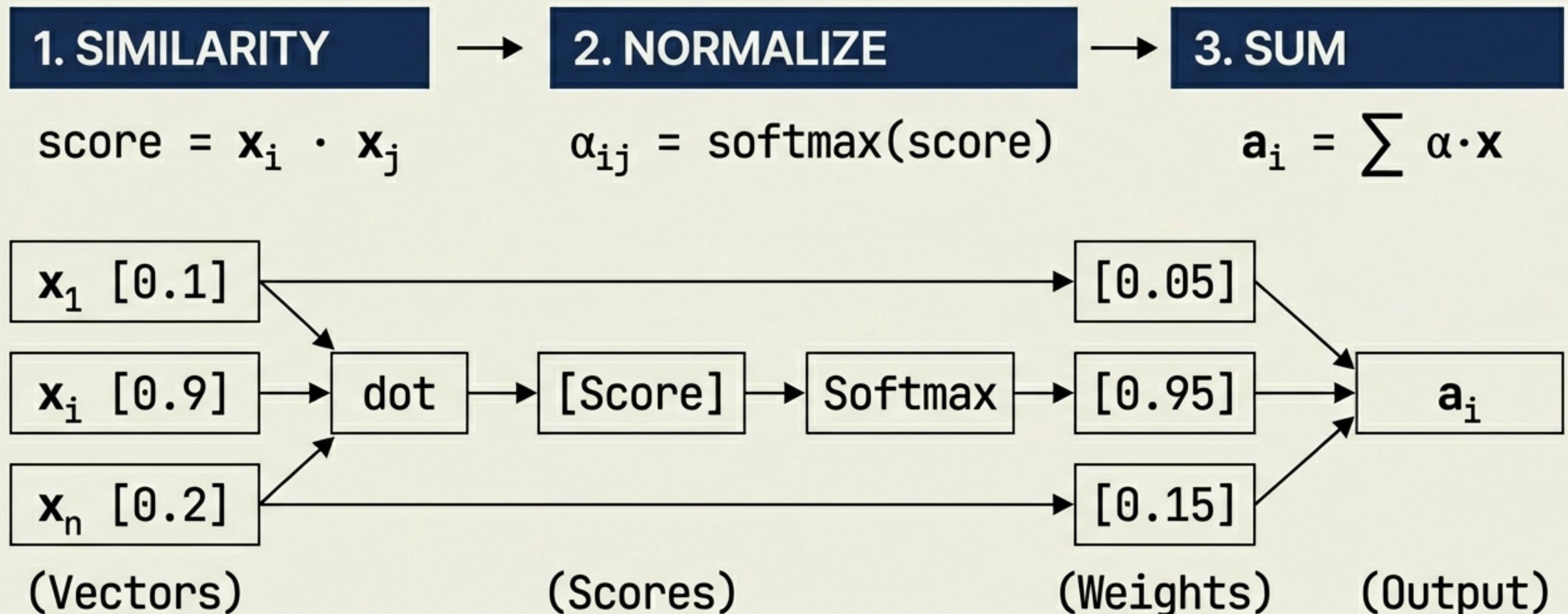
Look at the column for ‘it’ (far right). It is attending strongly (dark blocks) to ‘chicken’ and ‘road’.

ATTENTION HEATMAP (Weights)



MECHANISM: ATTENTION AS A WEIGHTED SUM

Goal: Produce output vector $[a_i]$ from input vectors $[x_j]$



THE THREE ROLES OF A TOKEN

Raw vectors are limited. In a Transformer, we project vectors into three distinct roles, like a database lookup.



QUERY (Q)

"The Agent"

What I am currently looking for.

e.g. "I am 'it', who defines me?"



KEY (K)

"The Label"

How I identify myself to others.

e.g. "I am 'chicken', a noun."



VALUE (V)

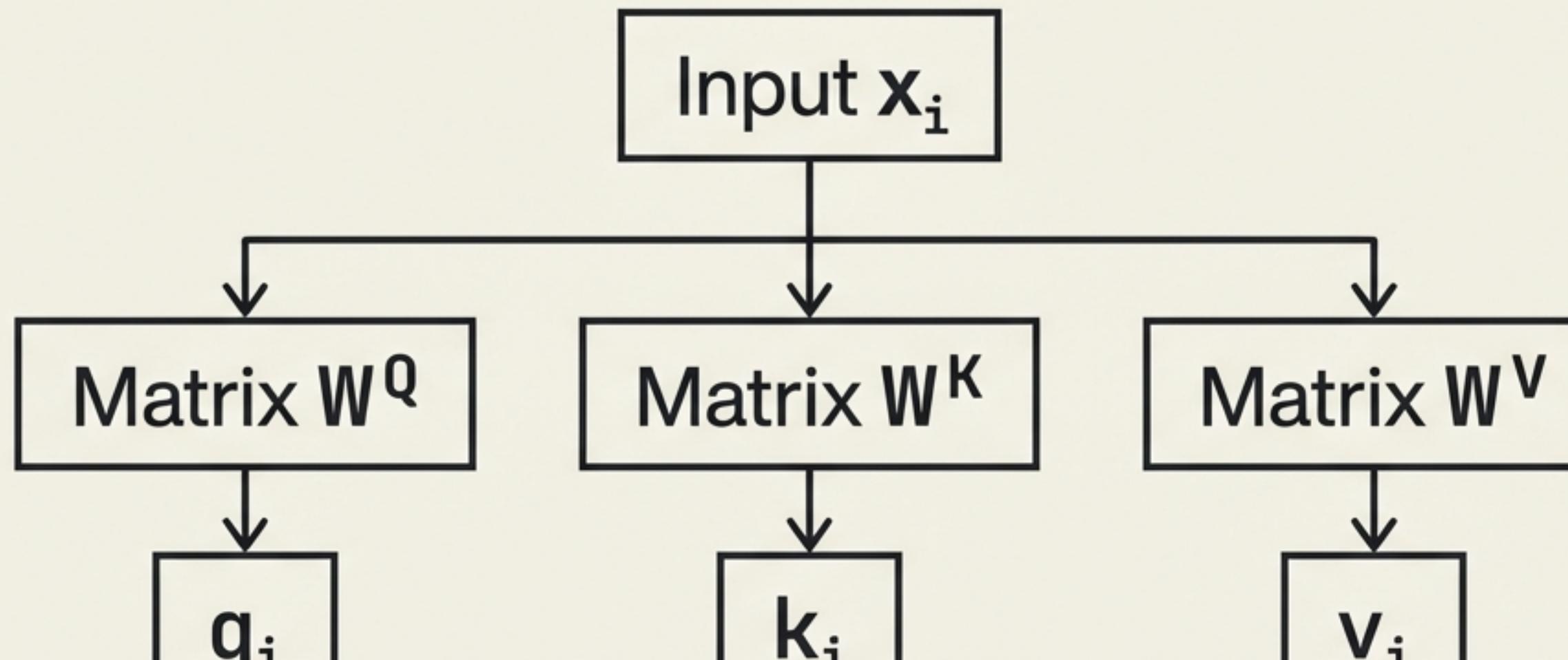
"Content"

What I pass on.

The info payload.

PROJECTIONS CREATE THE VECTORS

We use learned matrices to project input x into roles.



Query

Key

Value

$$q = x \cdot W^Q$$

$$k = x \cdot W^K$$

$$v = x \cdot W^V$$

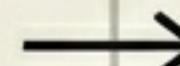
THE MECHANISM OF ATTENTION

THE LOGIC



1. COMPARE

Take the Query of the current token and Dot Product with Keys of context.



THE CORE EQUATION

$$\text{score}(x_i, x_j) = \frac{q_i \cdot k_j}{\sqrt{d_k}}$$



2. NORMALIZE

Scale by root dimension size, then Apply Softmax to get probabilities.



$$\alpha_{ij} = \text{softmax(score)}$$



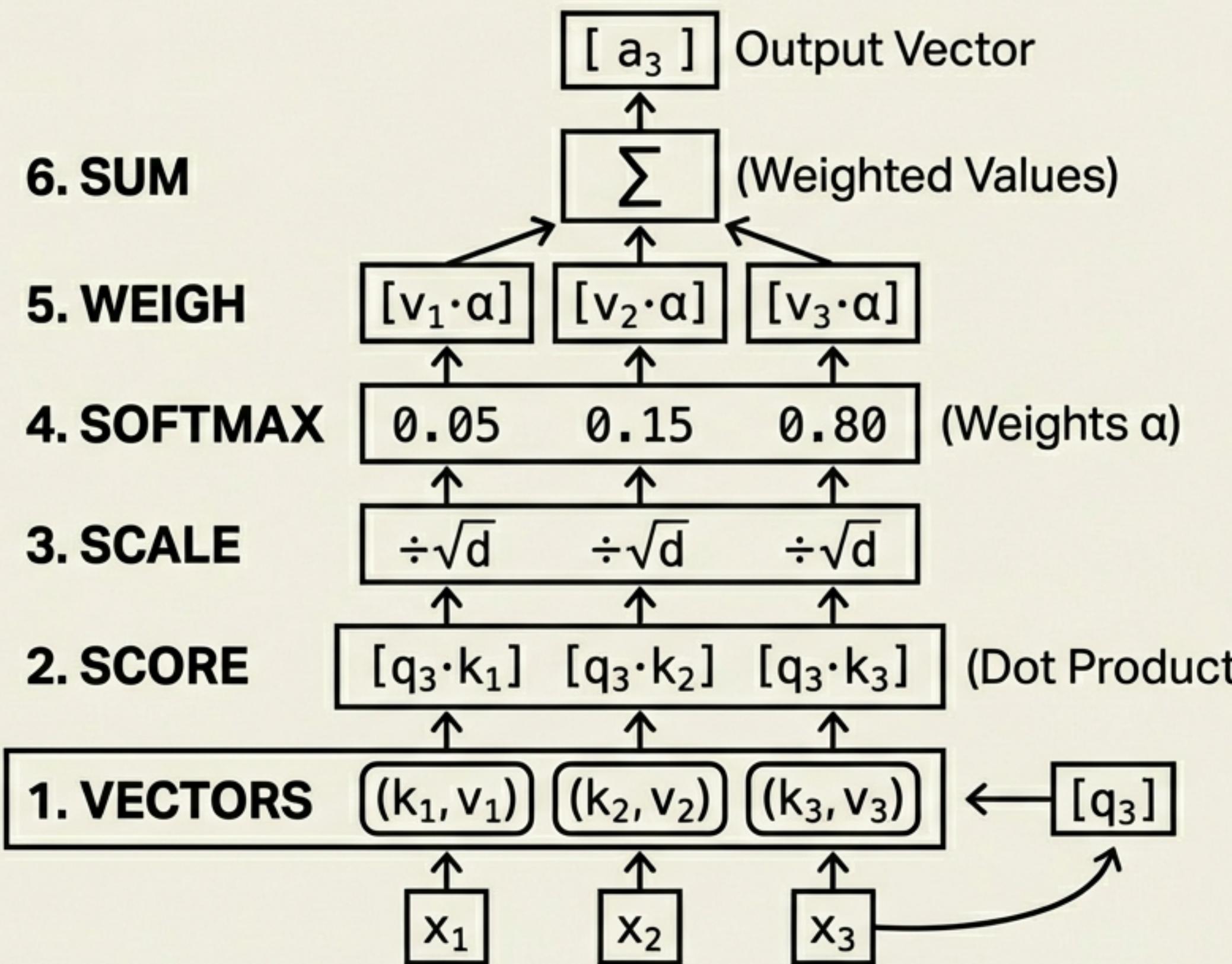
3. SUM VALUES

Multiply Value vectors by the weights and sum them.



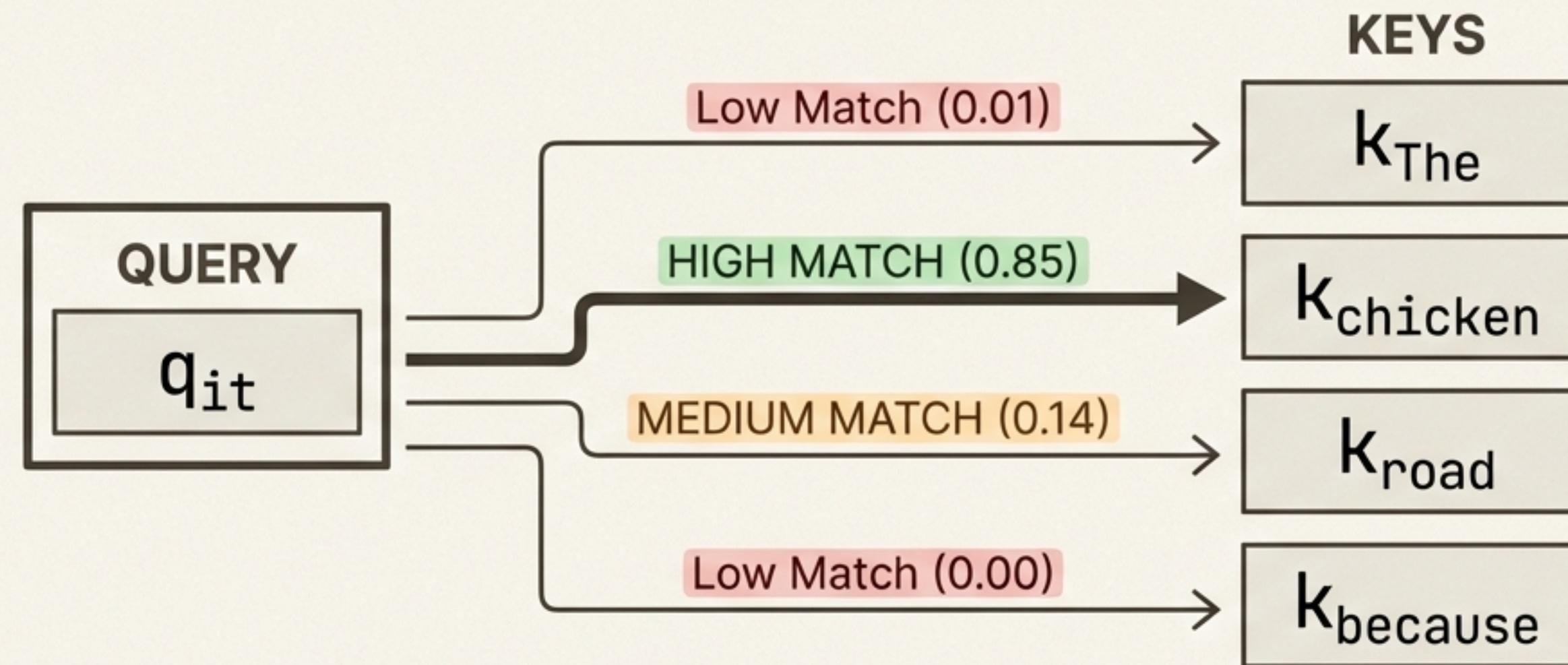
$$a_i = \sum (\alpha_{ij} \cdot v_j)$$

CALCULATING THE VALUE OF A TOKEN (a_3)



VISUALIZING THE COMPARISON (Query vs Keys)

Comparing the Query for 'it' against all other Keys.



Result: The final definition of 'it' will be composed mostly of the Value vector from 'chicken'.

MULTI-HEAD ATTENTION

Concept & Equation

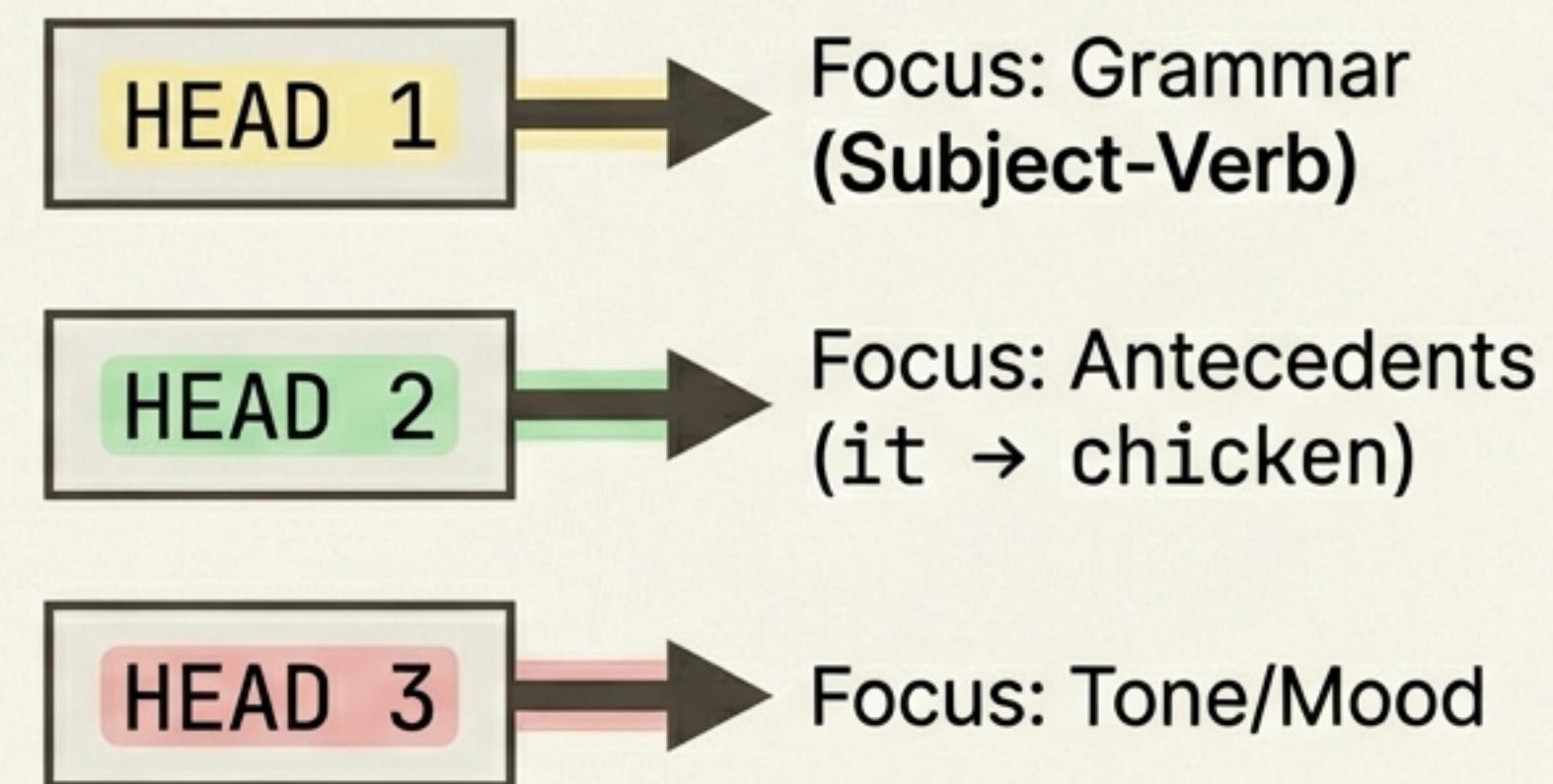
Transformers run multiple attention ‘heads’ in parallel.

$$\text{MultiHead}(x) = \text{Concat}(\mathbf{h}_1, \dots, \mathbf{h}_h) \cdot W^0$$

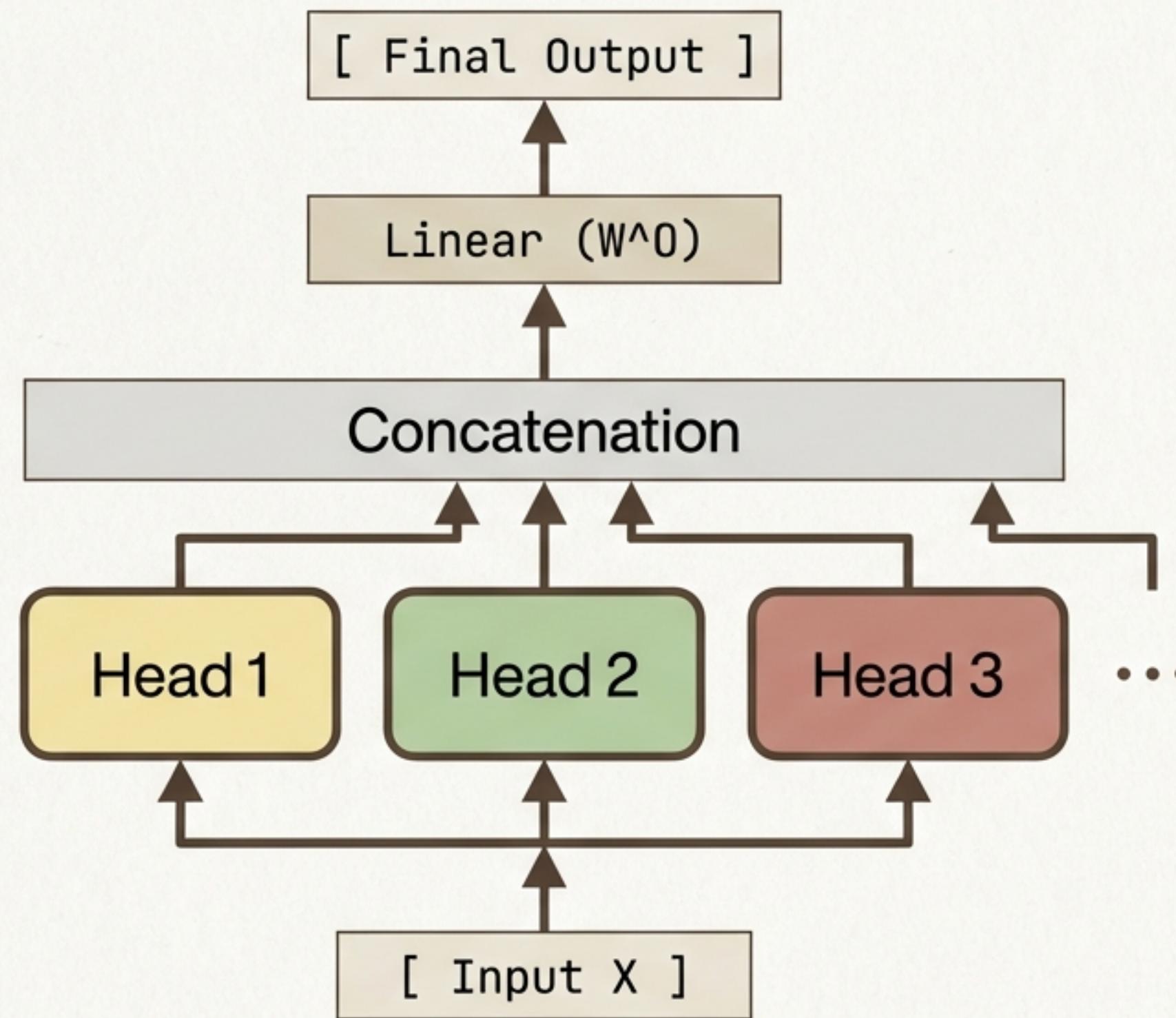
↗
Output Weight Matrix

Intuition & Examples

Different heads look for different relationships.



THE MULTI-HEAD ARCHITECTURE



SUMMARY: RESOLVING CONTEXT

1. ENRICHMENT

Attention enriches the representation of a token by incorporating information from its neighbors.

2. MECHANISM

It projects words into Queries, Keys, and Values to calculate relevance scores.

RESULT

3. RESULT

The final embedding for “it” is no longer static—it is a composite of the word itself plus the relevant parts of “chicken” or “road.”

