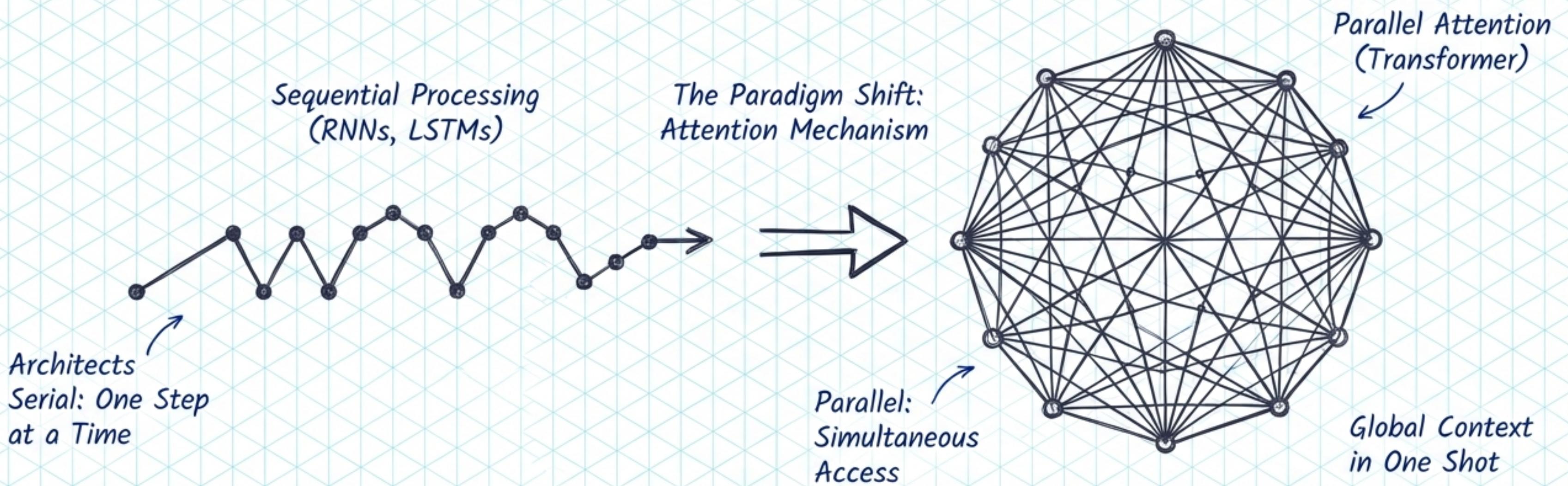


THE TRANSFORMER ARCHITECTURE

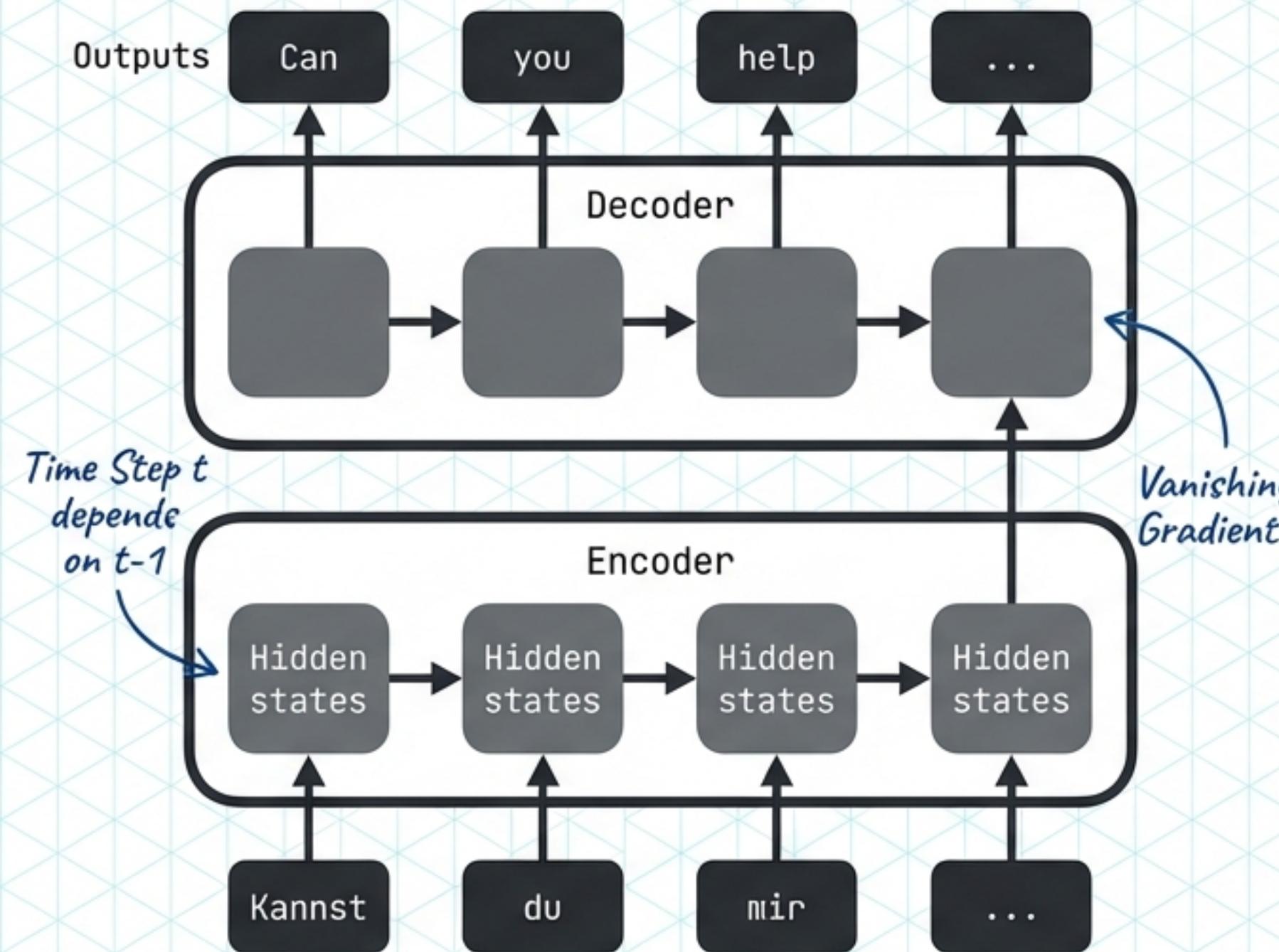
From Sequential Processing to Parallel Attention



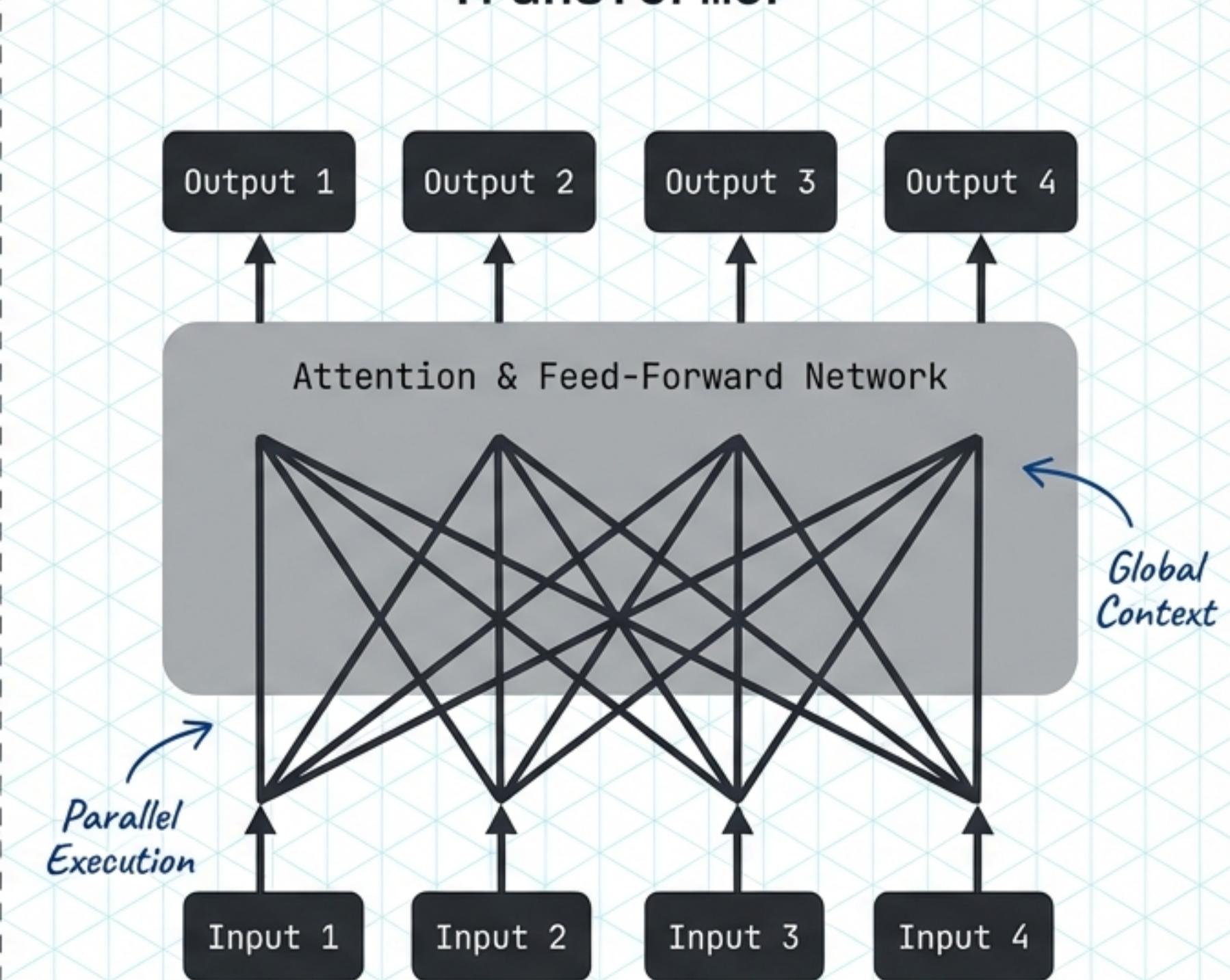
A Deep Dive for Engineering Graduate Students

THE SEQUENTIAL BOTTLENECK

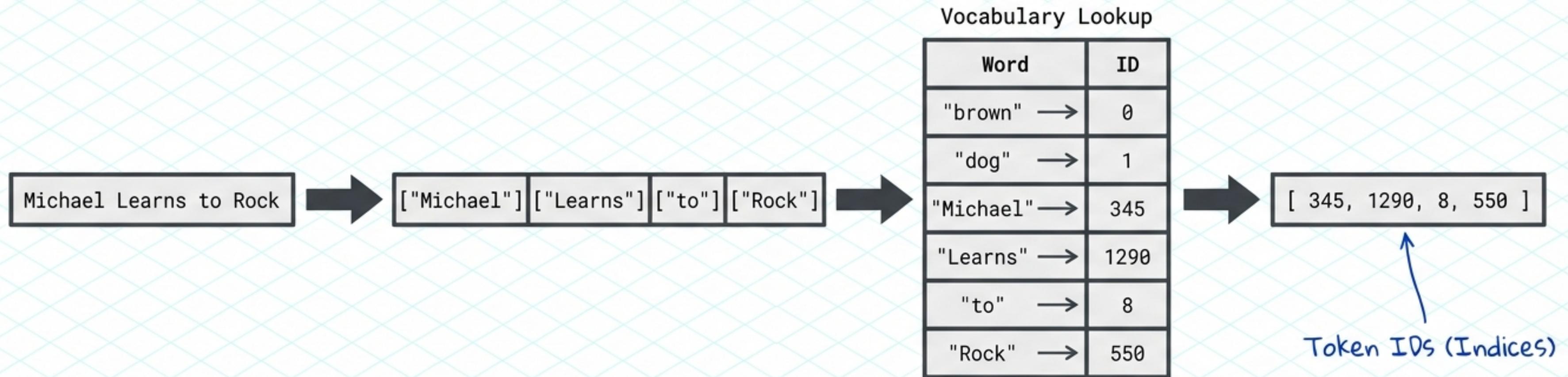
Recurrent Neural Network (RNN)



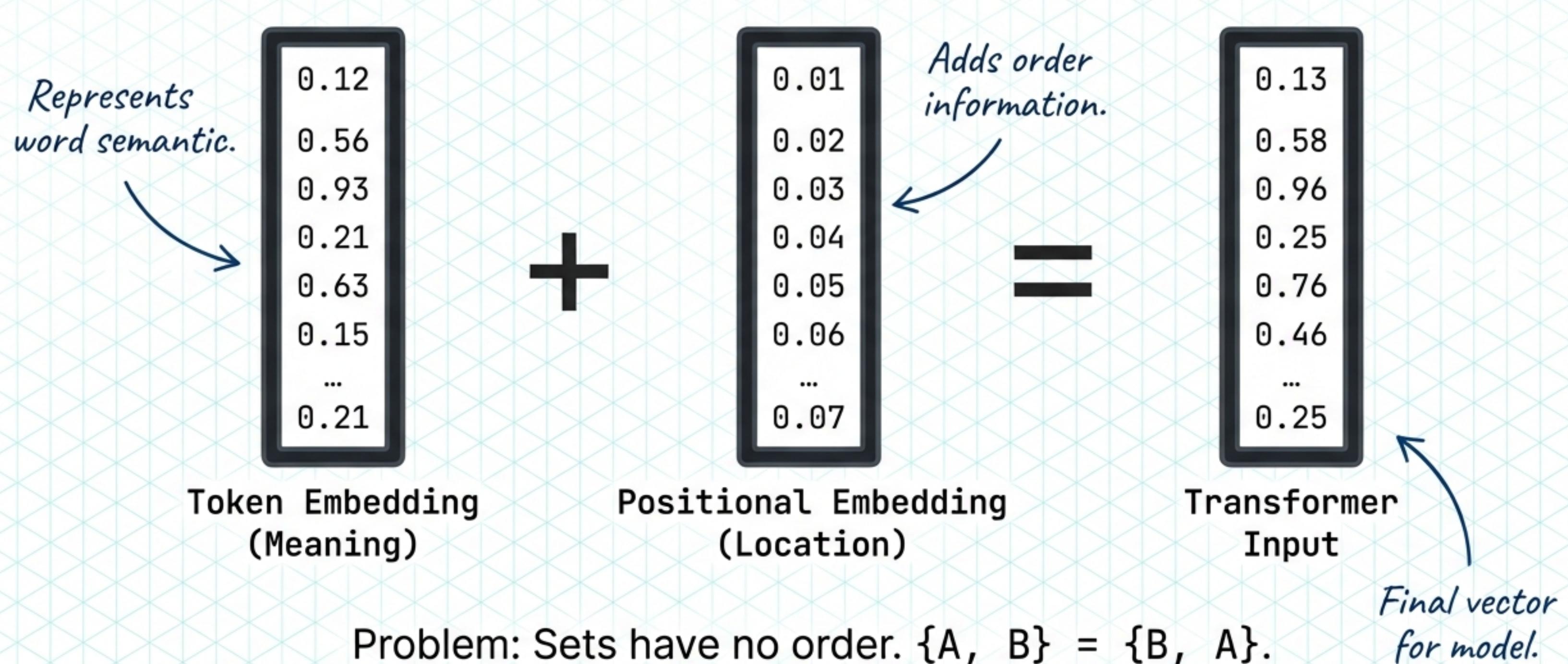
Transformer



INPUT PIPELINE: TOKENIZATION



EMBEDDINGS & POSITIONAL ENCODING



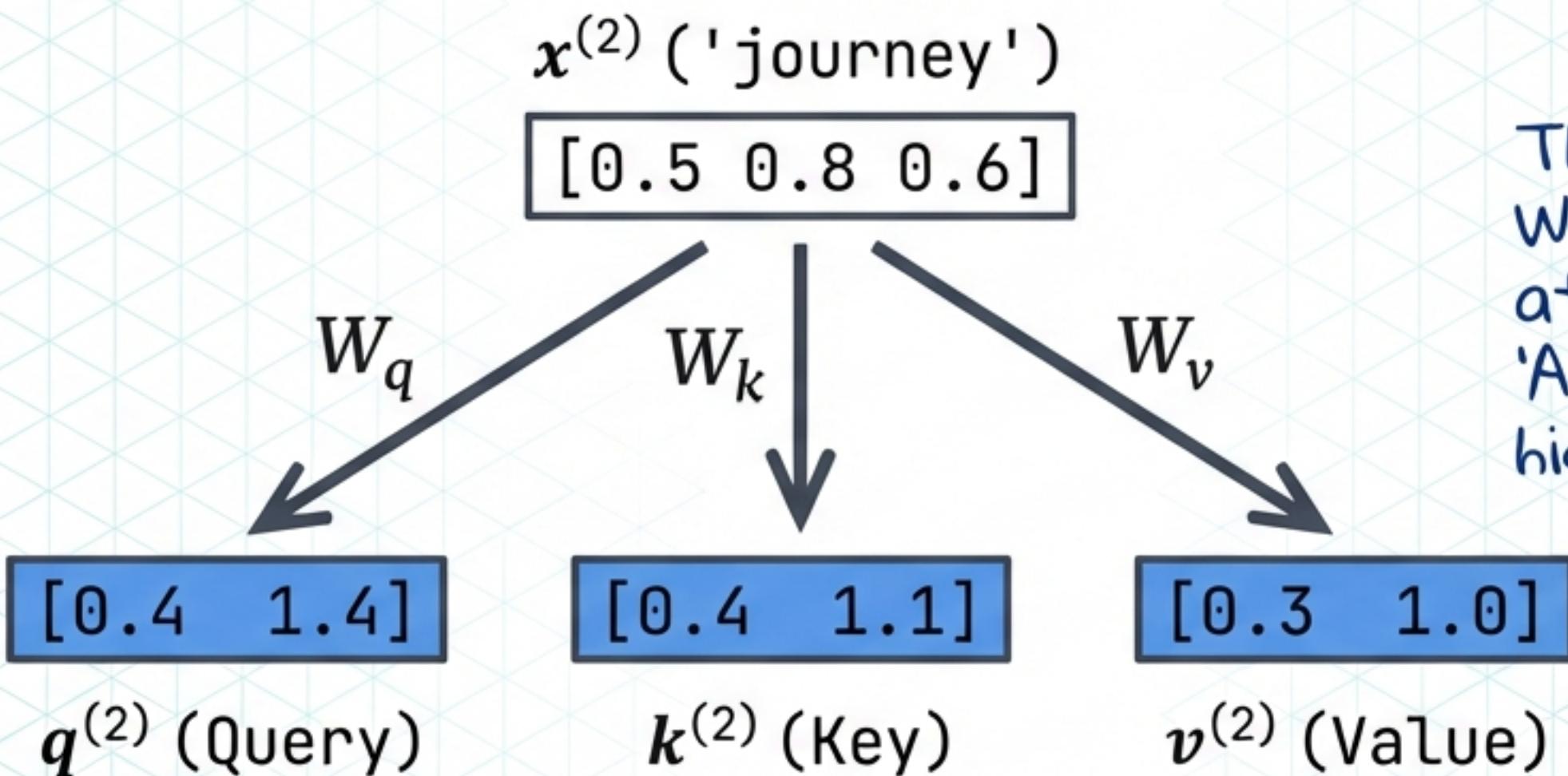
CORE MECHANISM: SELF-ATTENTION

Solving ambiguity through context.

The **animal** didn't cross the street because **it** was too tired.

The model asks:
What does "it" attend to?
"Animal" gets the highest weight.

THE MECHANICS: Q, K, V PROJECTIONS



The model asks:
What does 'it'
attend to?
'Animal' gets the
highest weight.

Crucial: W matrices are learned.
Q, K, V are derived per token.

ANALOGY 1: THE CASTING DIRECTOR



QUERY (Q)

The Search:
“I need an action
hero.”

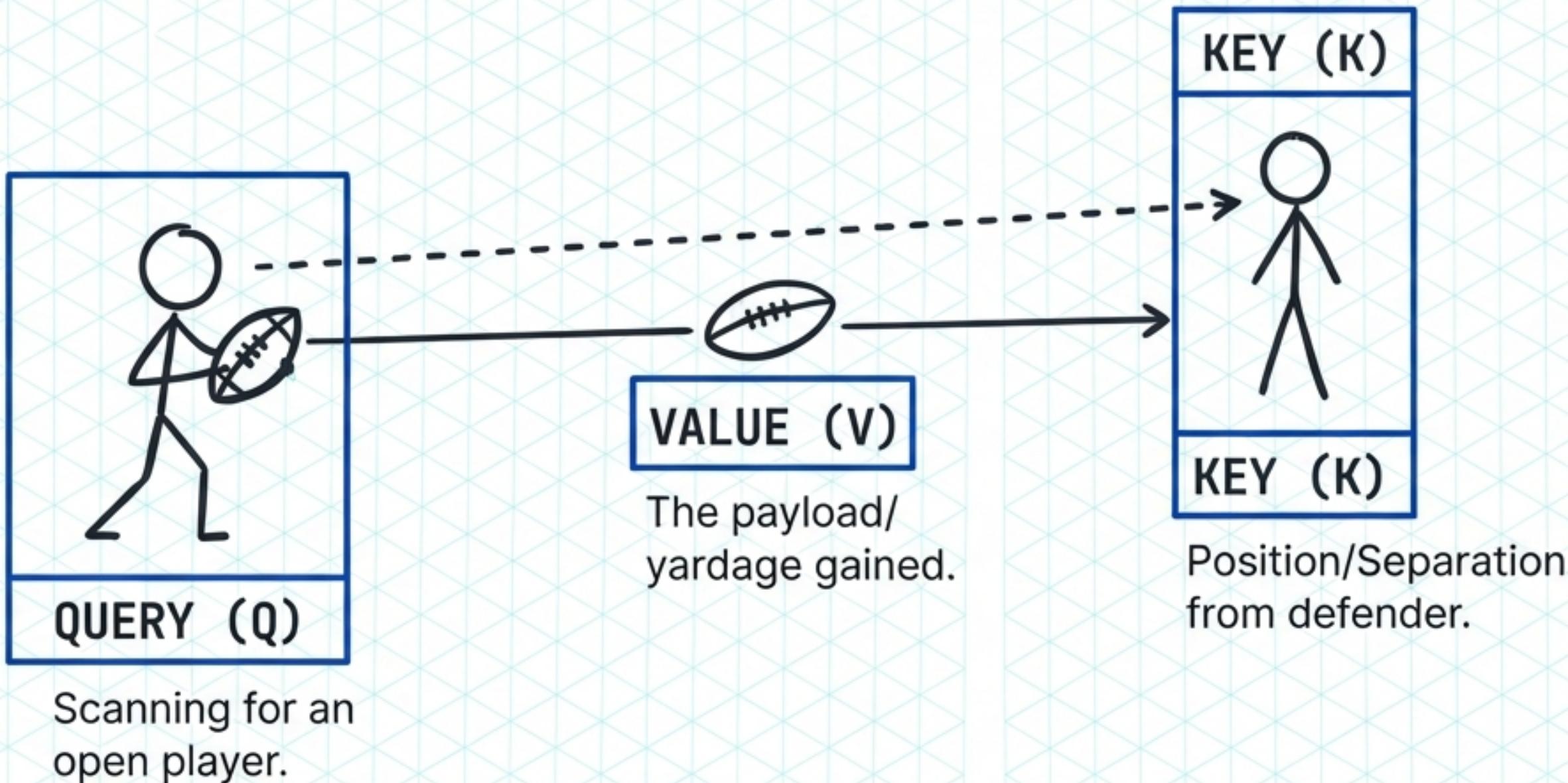
KEY (K)

The Attributes:
“I do stunts &
combat.”

VALUE (V)

The Content:
The actual
performance.

ANALOGY 2: THE FOOTBALL PASS



Key Insight: The QB (Q) selects the Receiver (K) based on fit, and delivers the Ball (V) as the actual payload.

ANALOGY 3: THE VOTER



QUERY (Q)

Voter Interests:
“Economy & Education”.



Match!



KEY (K)

Platform: “Lower
Taxes, New Schools”.



VALUE (V)

Policies Enacted.

Key Insight: The Voter (Q) matches their interests with a Platform (K) to elect the Policies (V) that will be enacted.

THE ALGORITHM: SCALED DOT-PRODUCT ATTENTION

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$

Convert to Probabilities (0 to 1)

Similarity Score (Alignment)

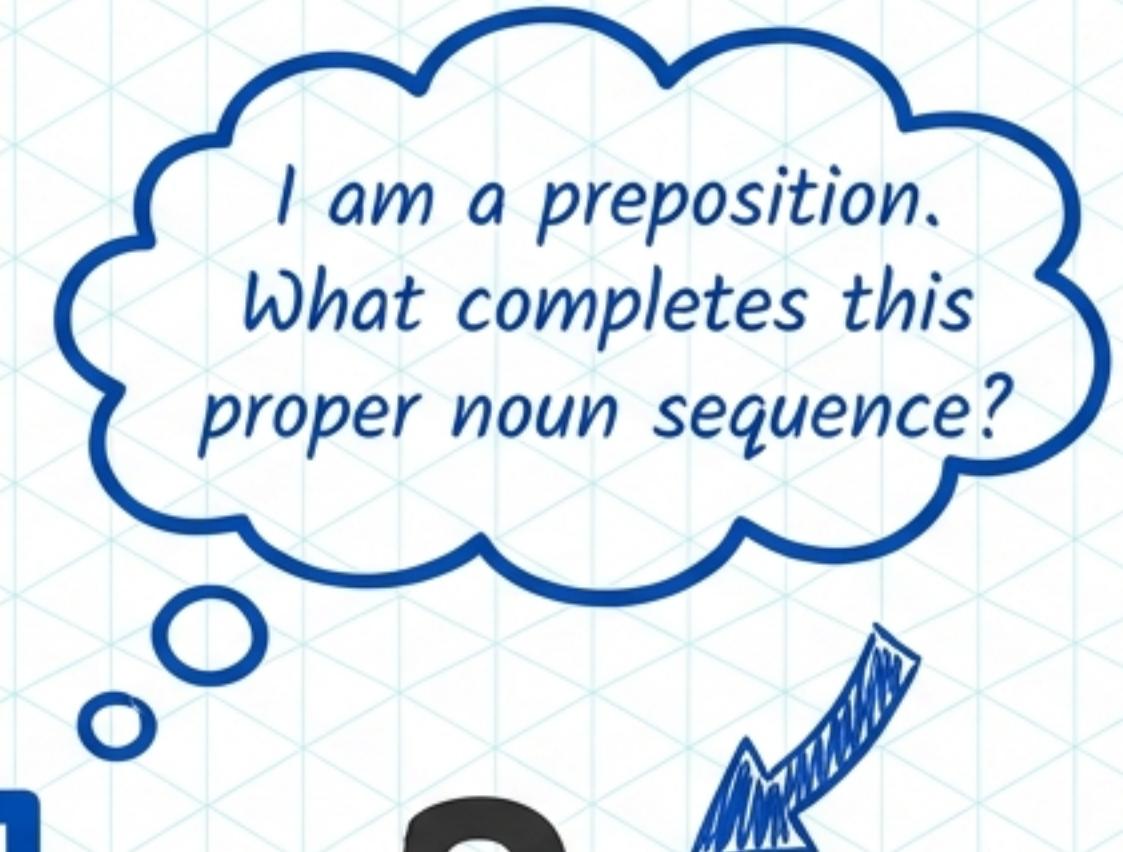
Scaling Factor:
Prevents vanishing gradients in Softmax

Weighted Aggregation of Information

EXAMPLE: THE QUERY

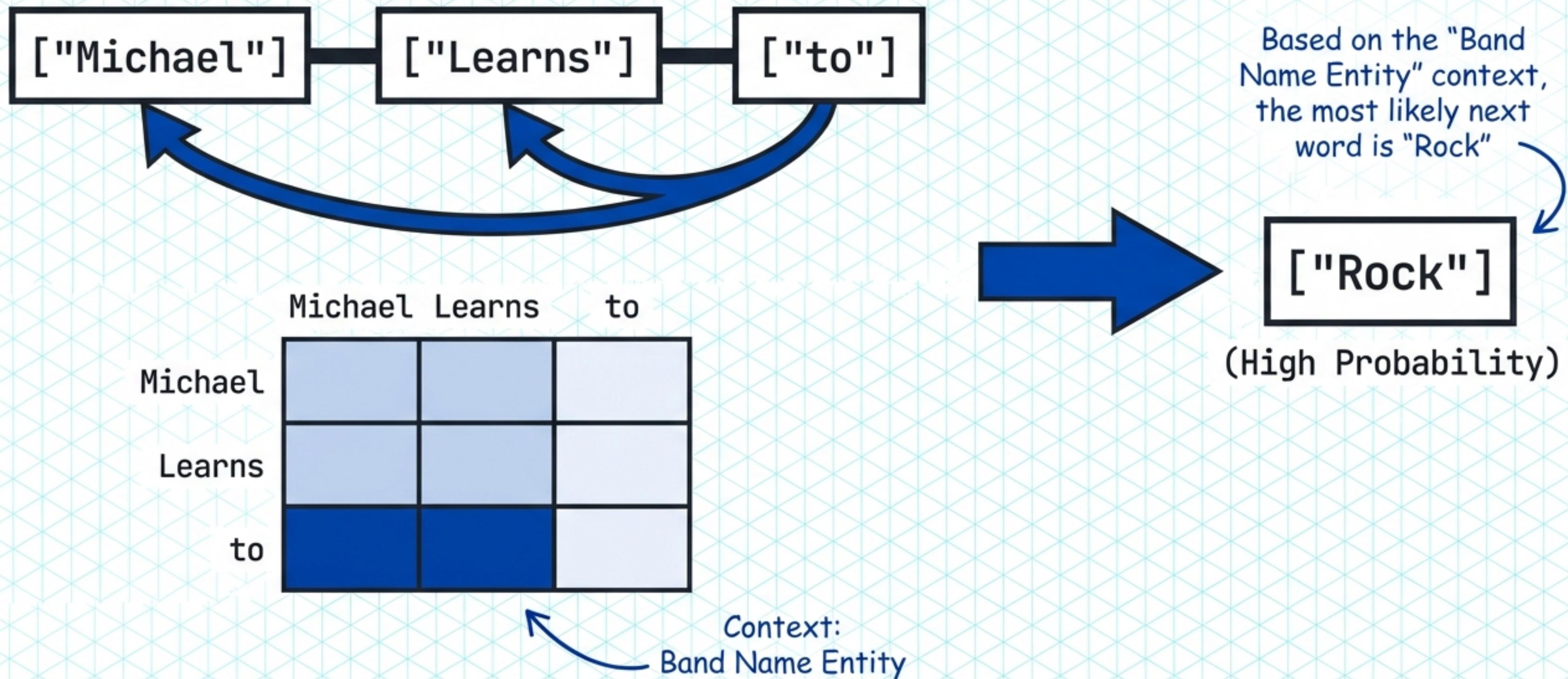
["Michael"] ["Learns"] **["to"]**

Q_{to}



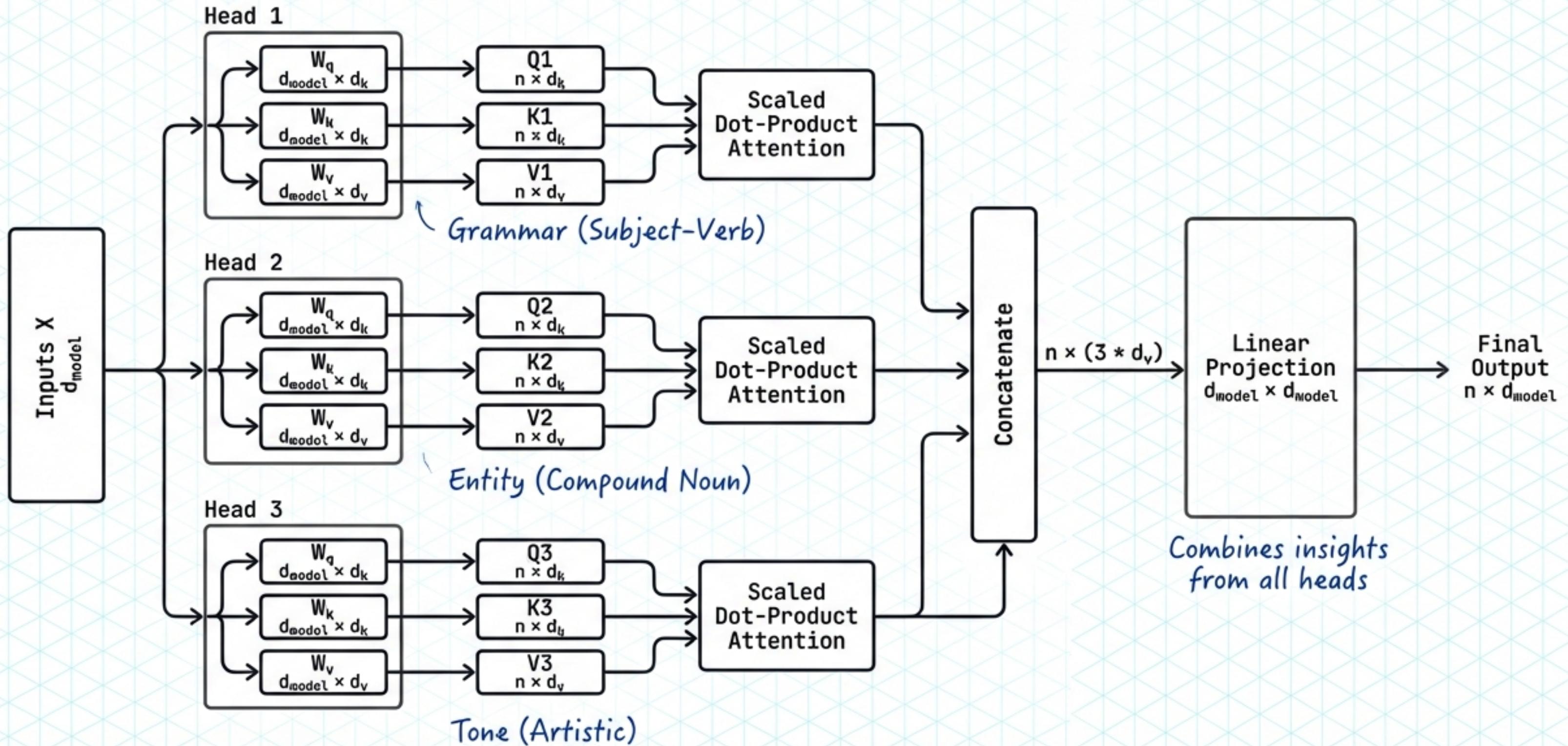
?

EXAMPLE: ATTENTION & PREDICTION



MULTI-HEAD ATTENTION

Capturing different types of relationships simultaneously.



MASKED (CAUSAL) ATTENTION

Preventing the model from seeing the future.

| | Michael | Learns | to | Rock |
|---------|---------|--------|------|------|
| Michael | 1.0 | | | |
| Learns | 0.55 | 0.45 | | |
| to | 0.30 | 0.35 | 0.35 | |
| Rock | 0.20 | 0.25 | 0.25 | 0.30 |

During training,
“Michael” cannot
attend to “Learns”.
Future is unknown.

Masked (Causal) Attention: Upper triangle zeros out future positions.

THE TRANSFORMER BLOCK

The Fundamental Building Block of GPT.

