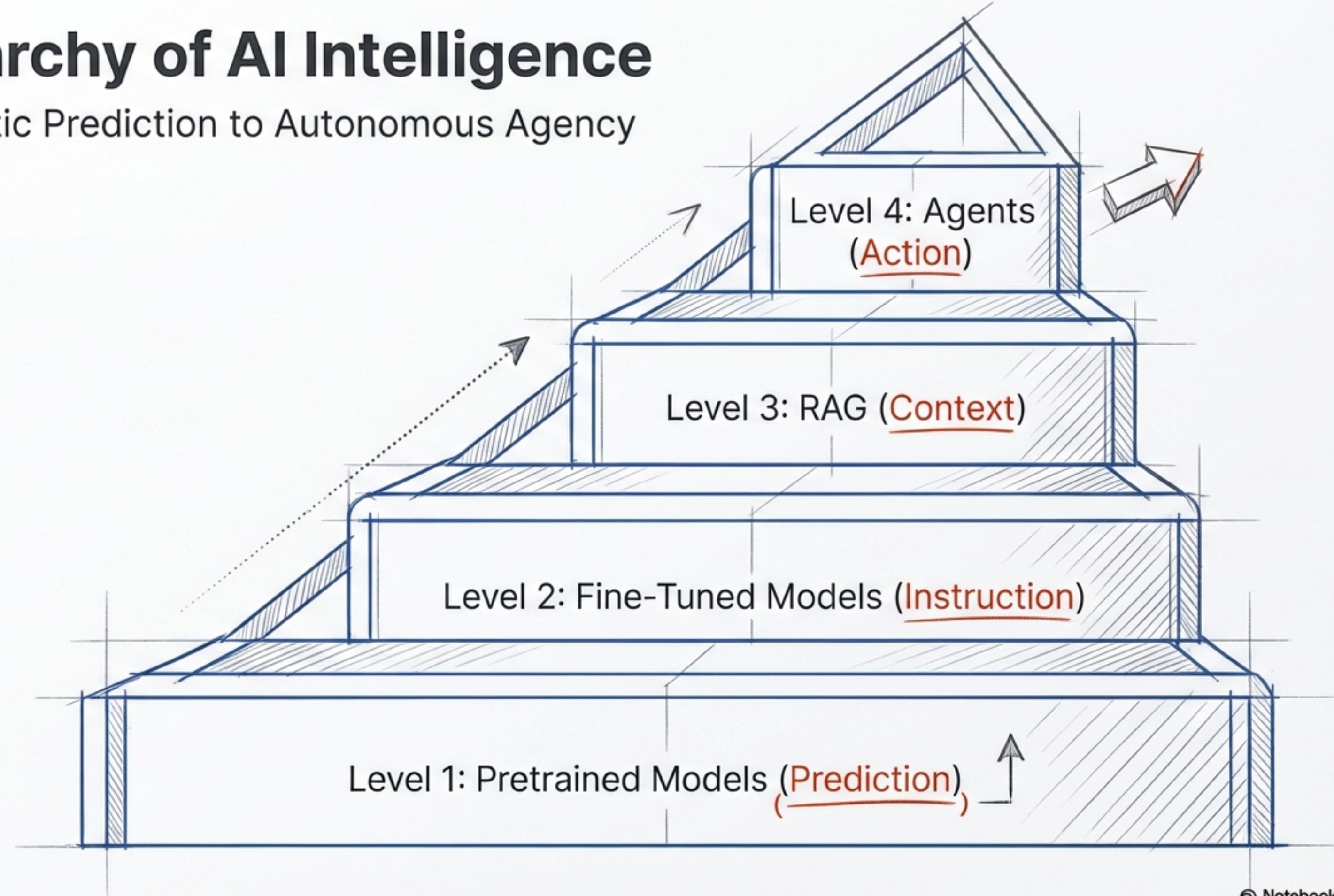


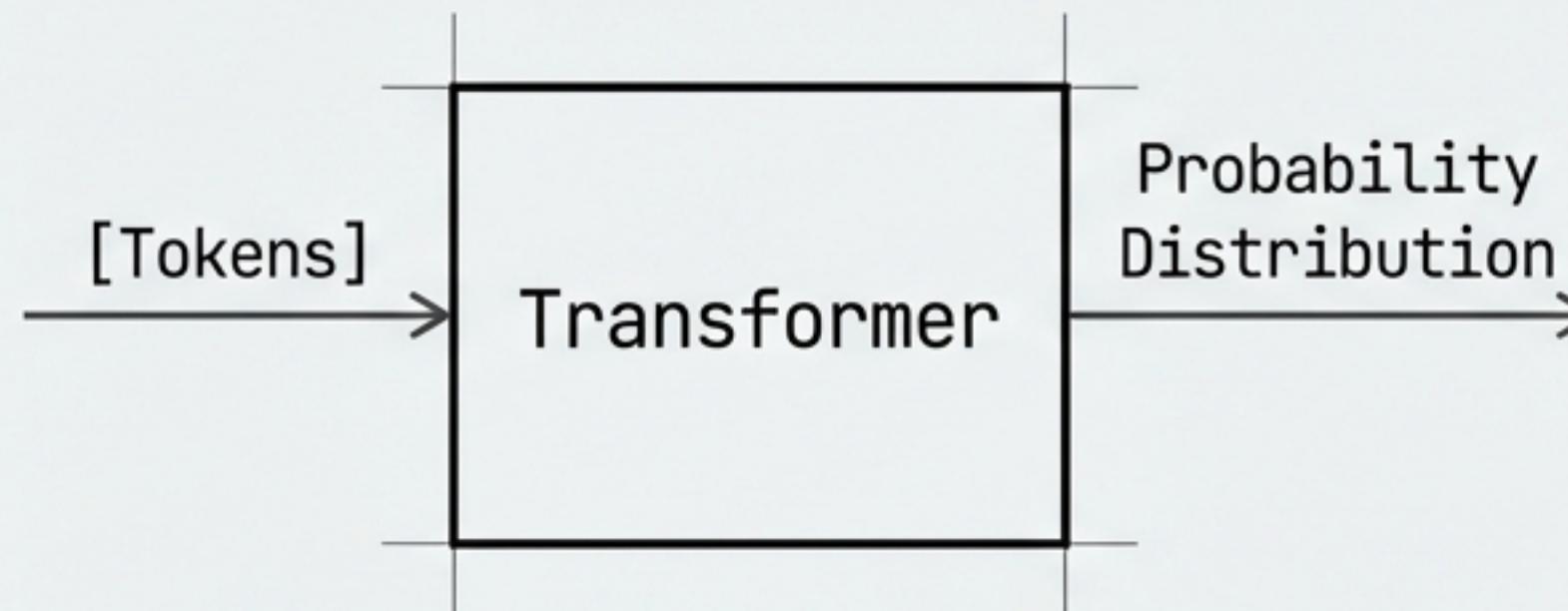
The Hierarchy of AI Intelligence

From Probabilistic Prediction to Autonomous Agency



Level 1: The Pretrained Model

The Well-Read Parrot



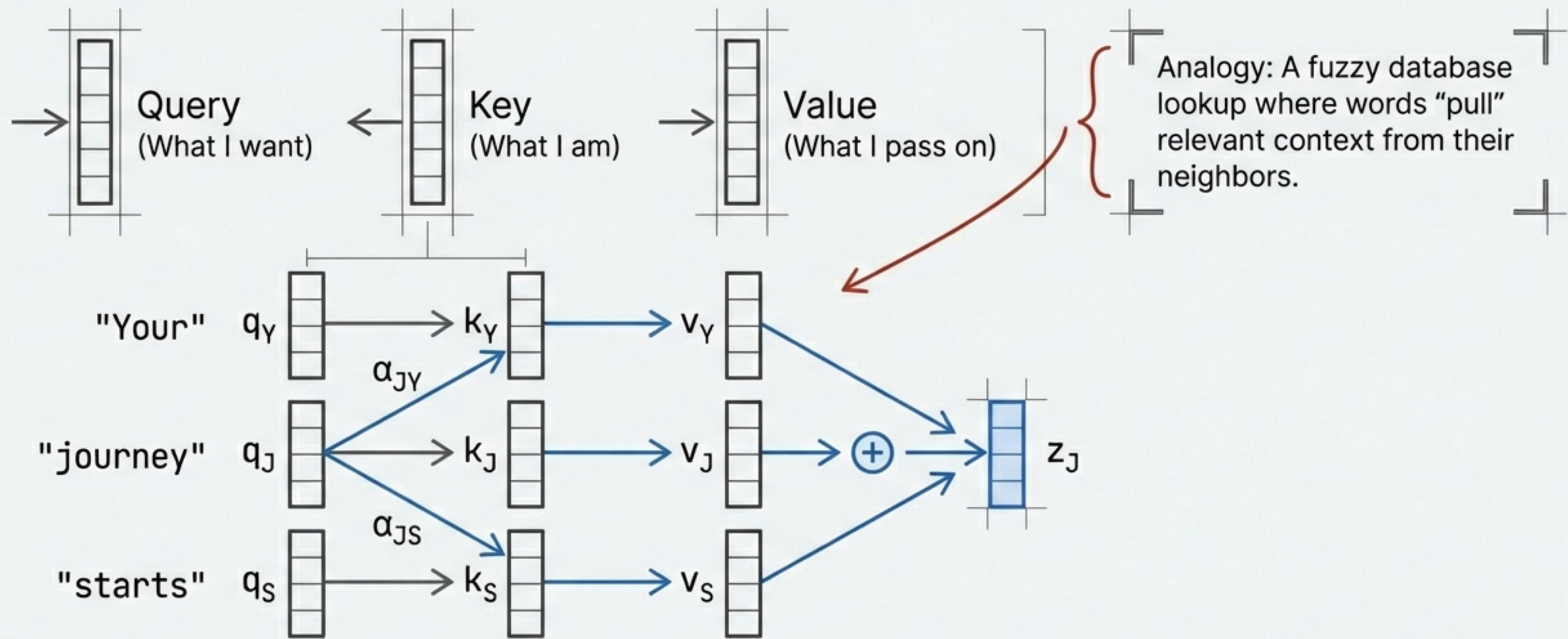
- Iteration
- ↓
1. [15496] Hello -> [257] a
 2. [15496, 257] Hello, I -> [716] am
 3. [15496, 257, 716] Hello, I am -> [257] a
- ...
- Append & Predict
- Append & Predict

Maximize $P(w_t | w_{1:t-1})$

Learns syntax & facts via compression.
No intent, only statistics.

The Engine: Self-Attention

Contextual Gravitation



The Math of Attention

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

Similarity Search
(Dot Product)

Normalization
(Probabilities sum to 1)

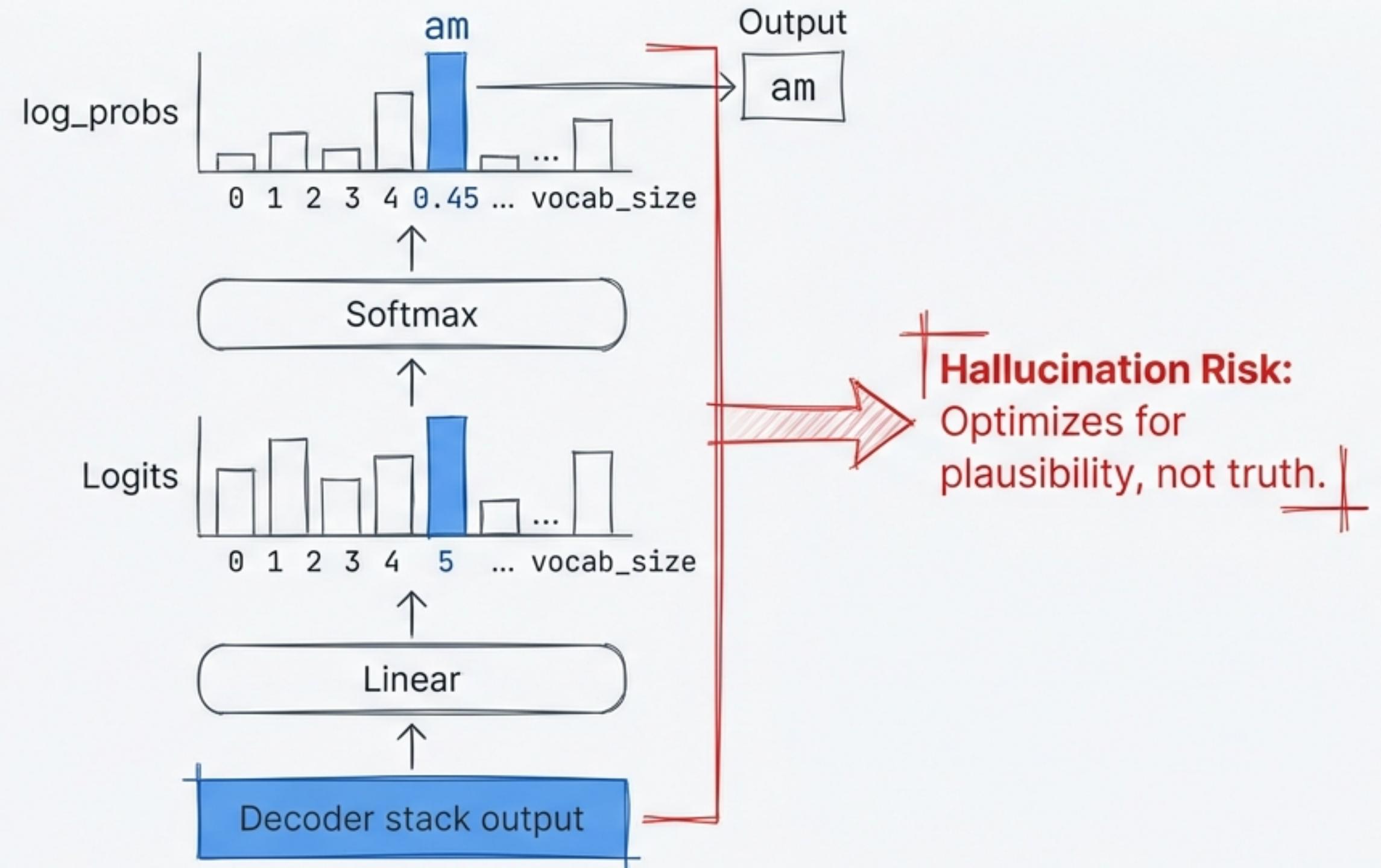
Scaling Factor
(Prevents gradient vanishing)

Weighted Information
Retrieval

The diagram illustrates the mathematical components of the Attention function. It starts with the formula $\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) * V$. Four arrows point from text labels to specific parts of the formula:

- An arrow points from "Similarity Search (Dot Product)" to the term QK^T .
- An arrow points from "Normalization (Probabilities sum to 1)" to the softmax function.
- An arrow points from "Scaling Factor (Prevents gradient vanishing)" to the term $\frac{1}{\sqrt{d_k}}$.
- An arrow points from "Weighted Information Retrieval" to the final result $* V$.

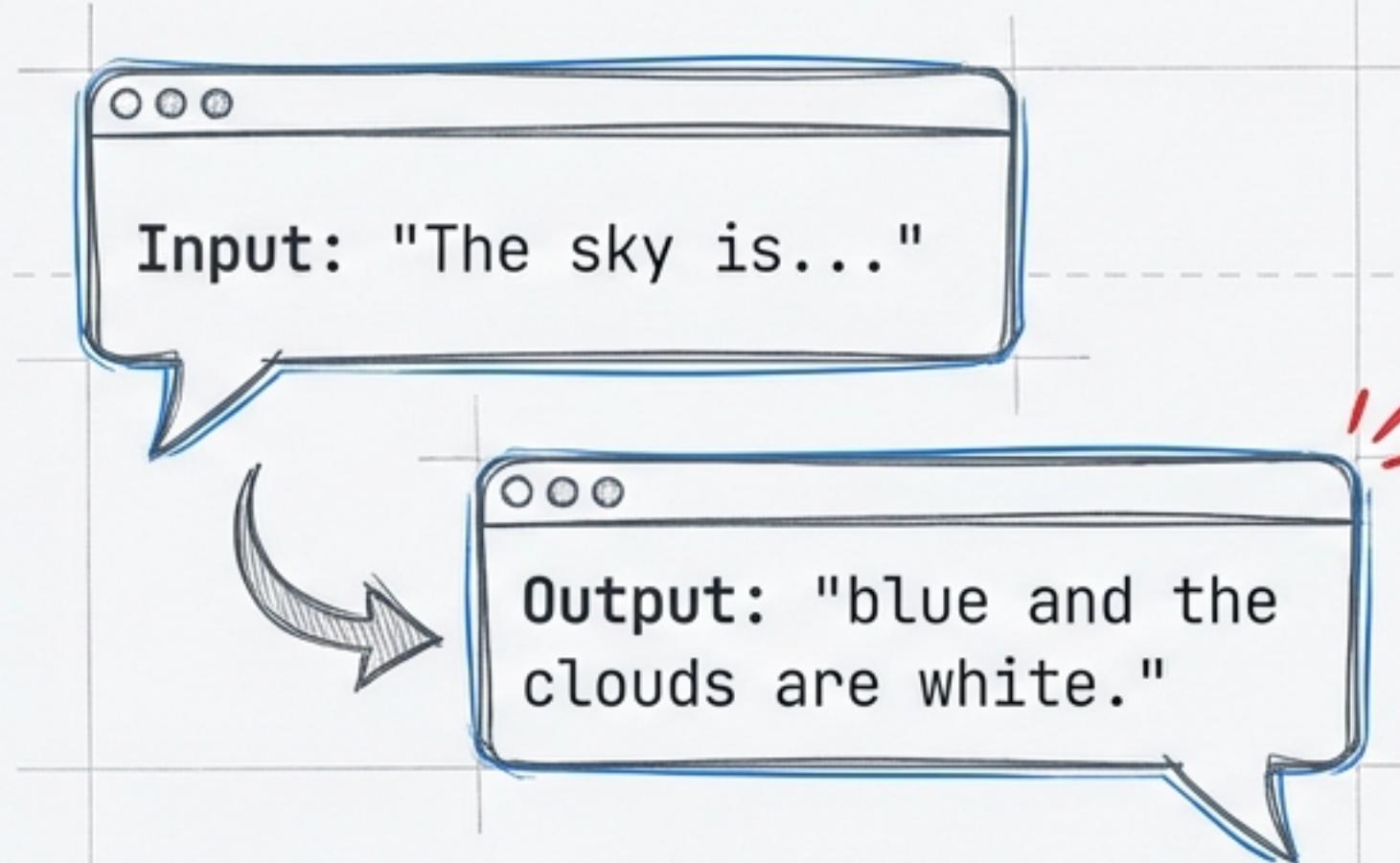
The Limitation: Probabilistic, Not Factual



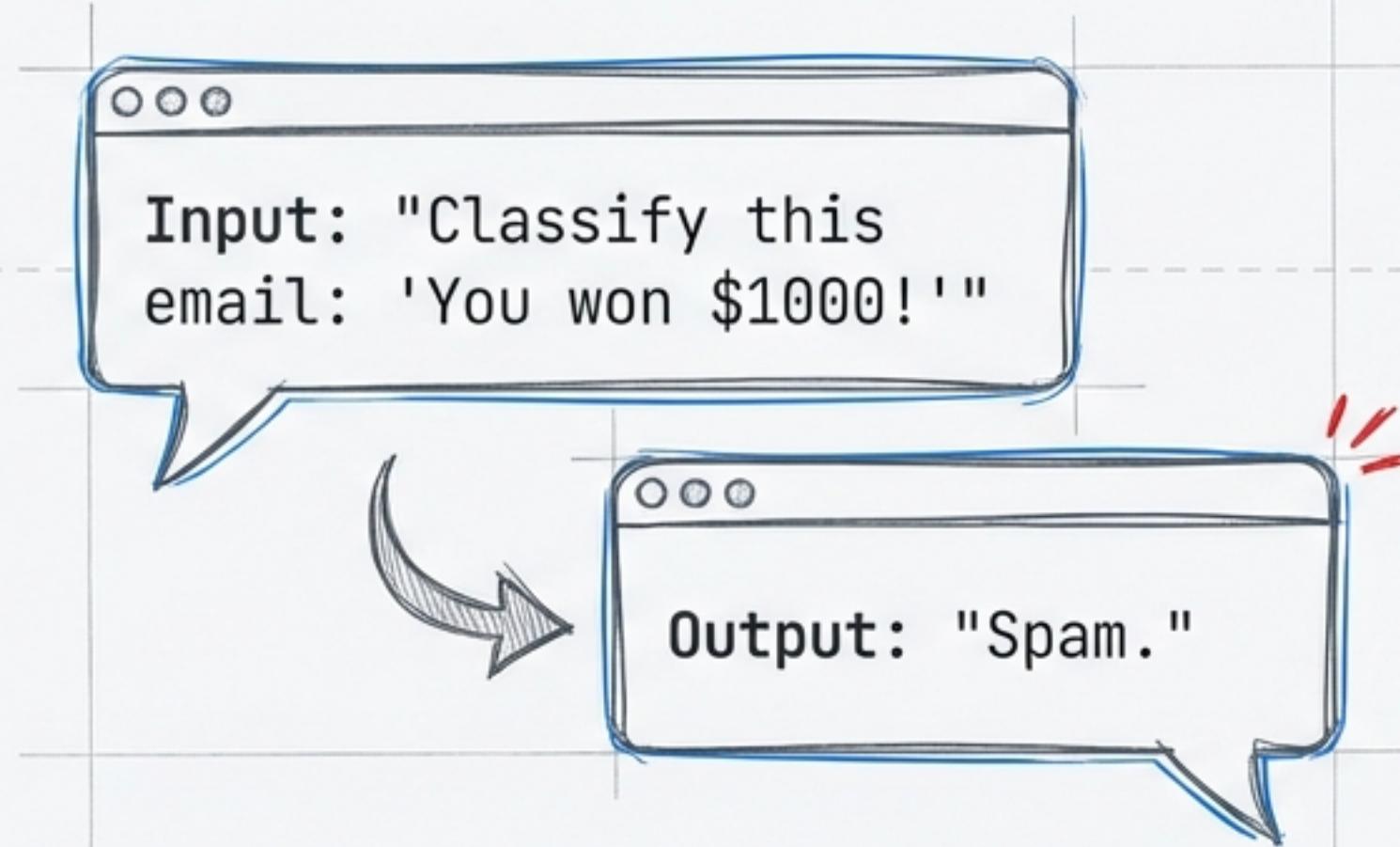
Level 2: The Specialist (Fine-Tuning)

From Completion to Instruction

Pretrained (Base Model)



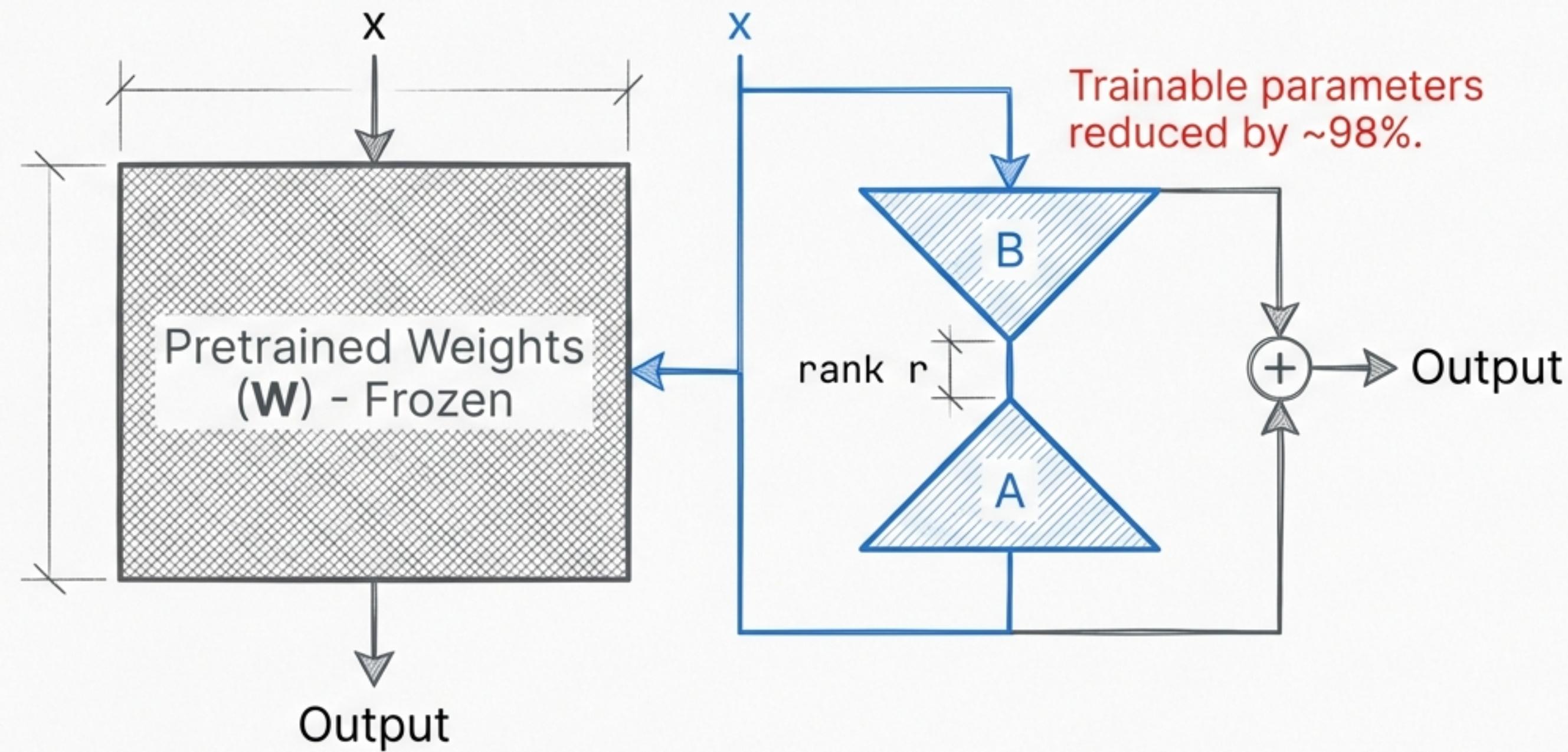
Fine-Tuned (Instruction Model)



Supervised Fine-Tuning (SFT) constrains the distribution to specific tasks.

The Efficiency Solution: LoRA

Low-Rank Adaptation



The LoRA Equation

$$W_{\text{updated}} = W + \Delta W \approx W + (B * A)$$

$$h = W_0 * X + B * A * X$$

Inter Bold
Frozen



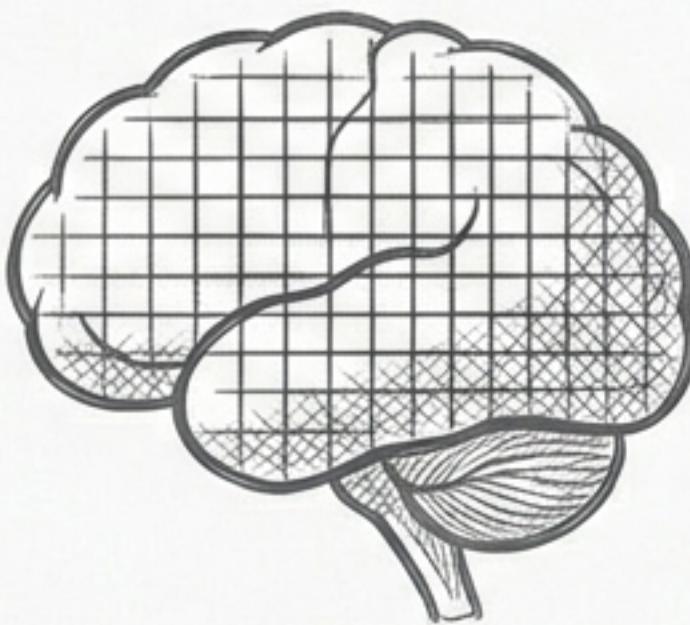
Inter Bold
Trainable Adapter



A is $d \times r$, B is $r \times d$, where $r \ll d$

Level 3: The Researcher (RAG)

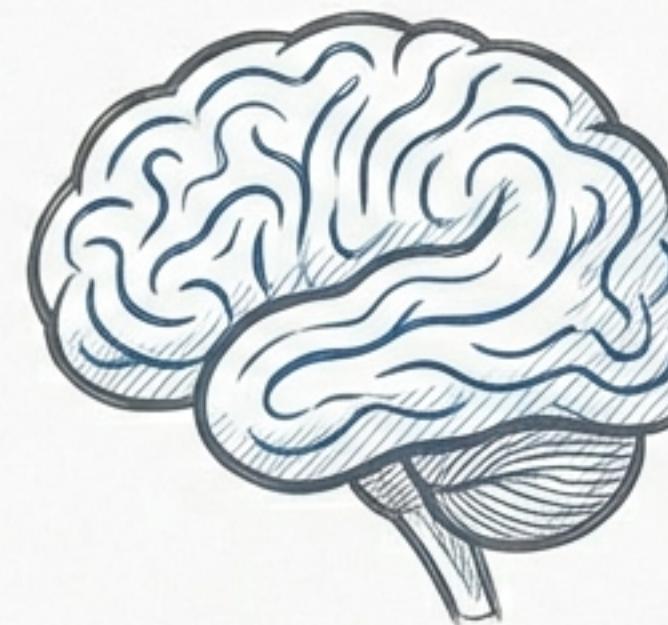
Overcoming the Knowledge Cutoff



Limited by
training data

Parametric Memory (Weights) ↗

Frozen, Expensive, Hallucinations



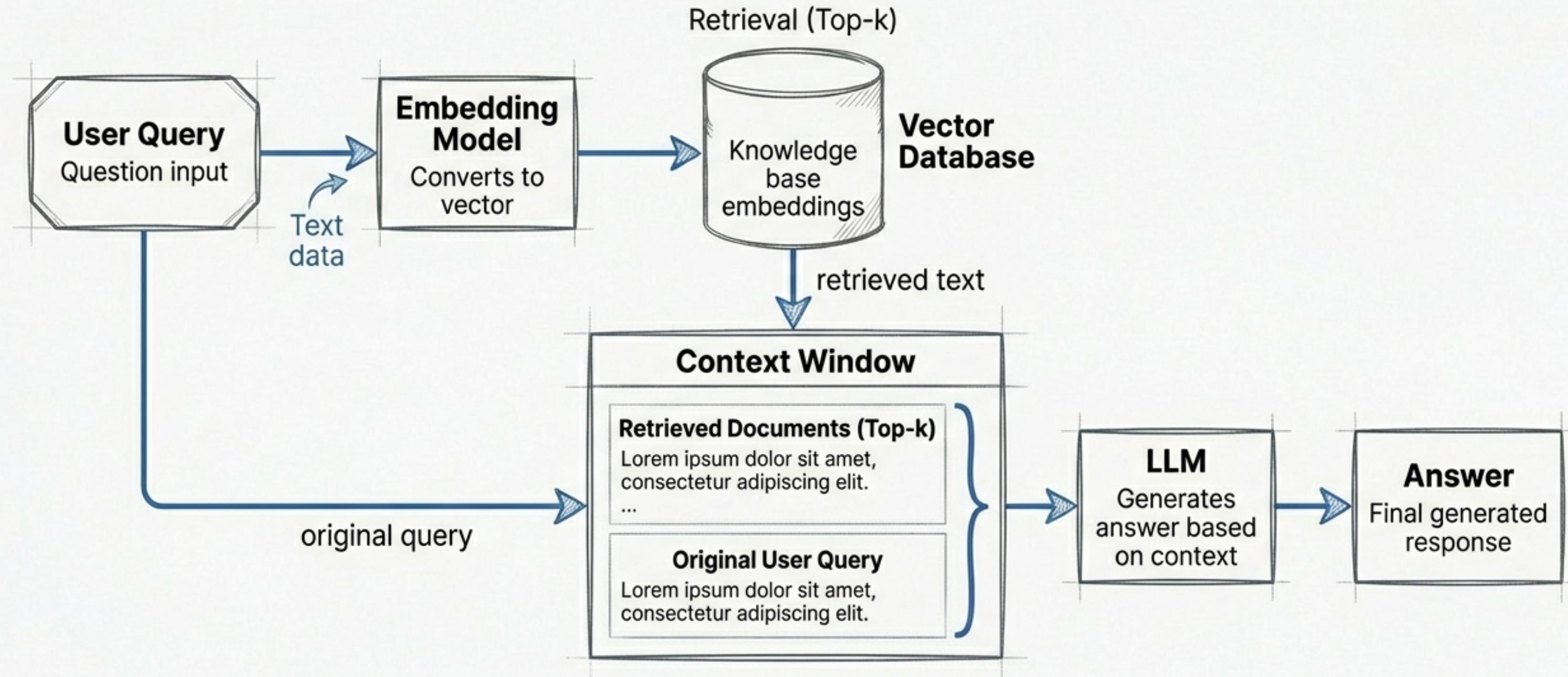
Non-Parametric Memory (Context)

Dynamic, Cheap, Grounded

↗ Augments with
real-time data

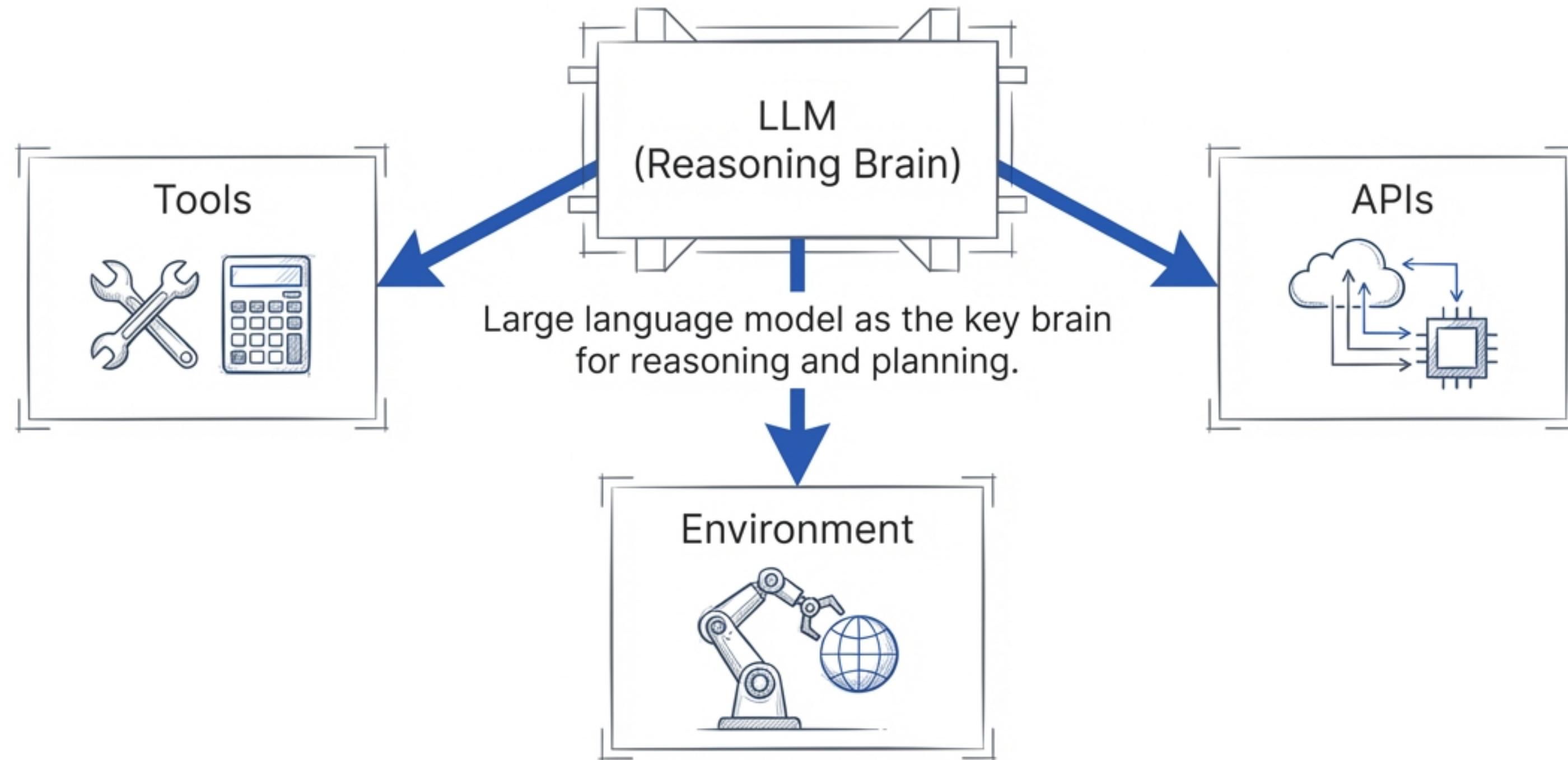
Retrieval Augmented Generation ↙

RAG Architecture: The 'Open Book' Exam

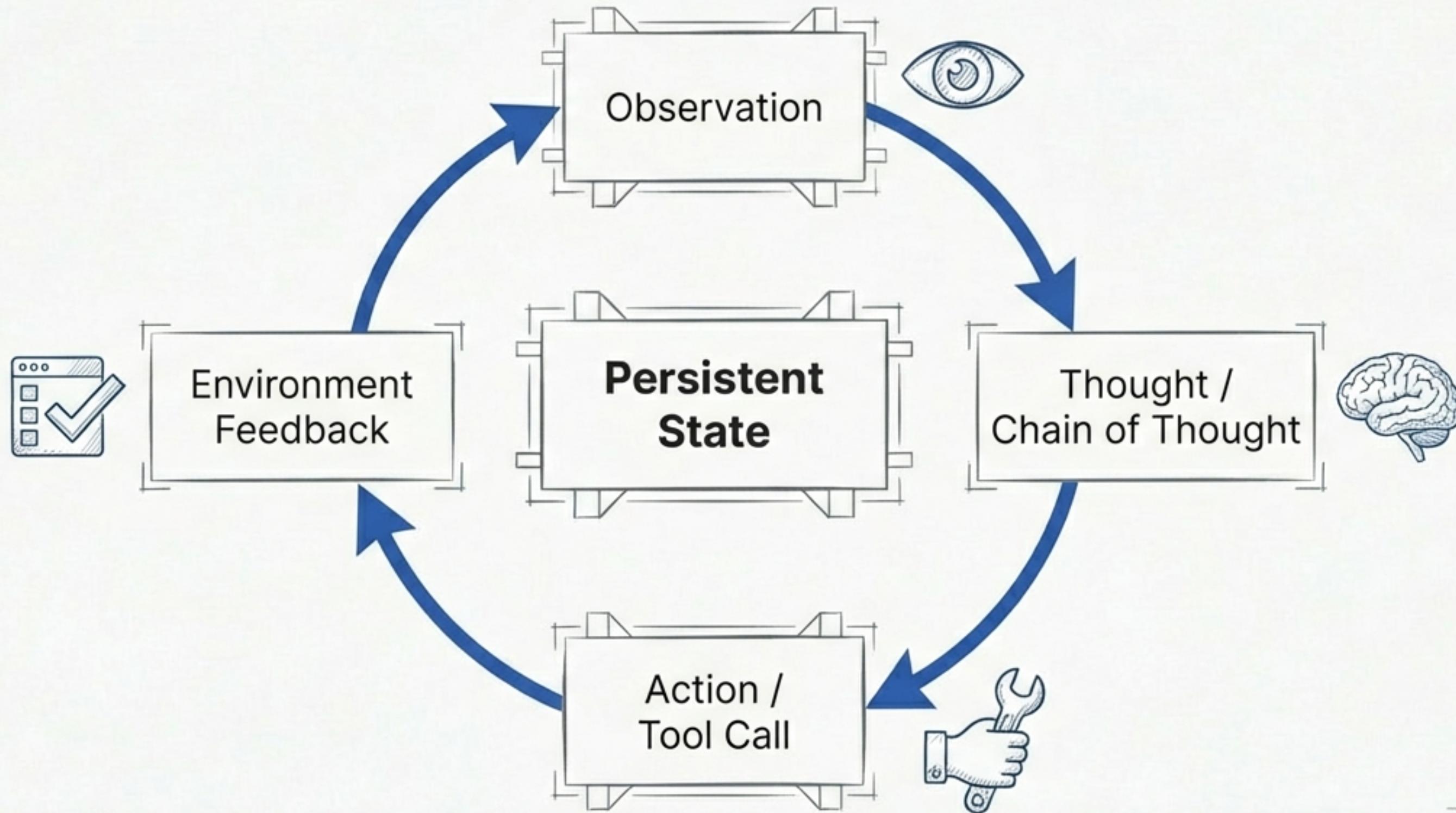


Level 4: The Agent (The Actor)

From Passive Text to Active Doing

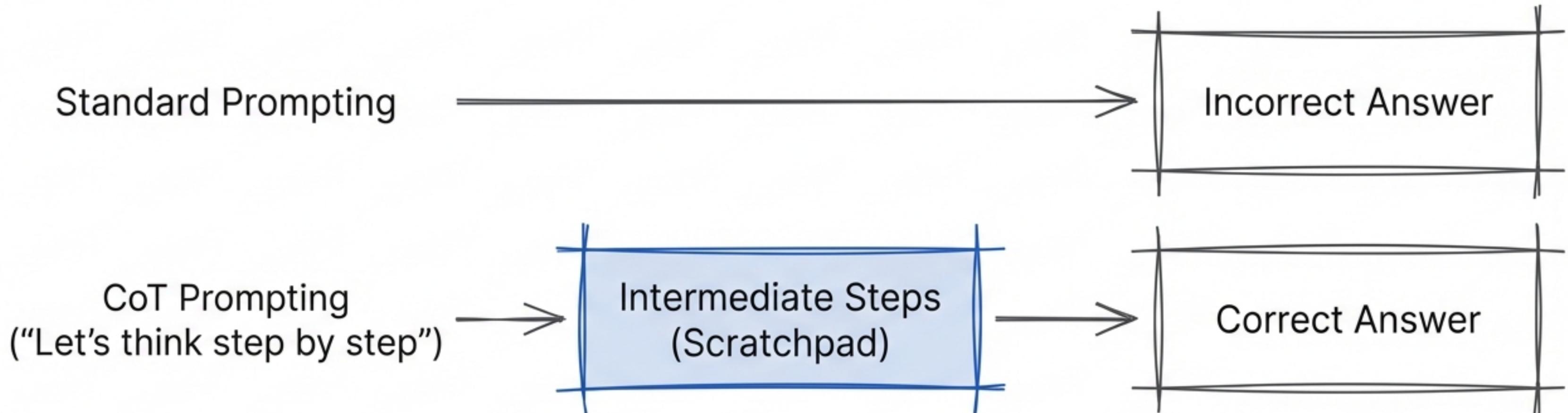


The Agentic Loop



Unlocking Reasoning: Chain of Thought (CoT)

System 1 vs. System 2 Thinking



Generating tokens is computation.
More tokens = more compute time = better reasoning.

CoT in Action

I have 3 apples. My dad has 2 more apples than me. How many do we have?

~~Zero Shot~~

~~Answer: 5~~

CoT

Step 1: I have 3 apples.

Step 2: Dad has $3 + 2 = 5$ apples.

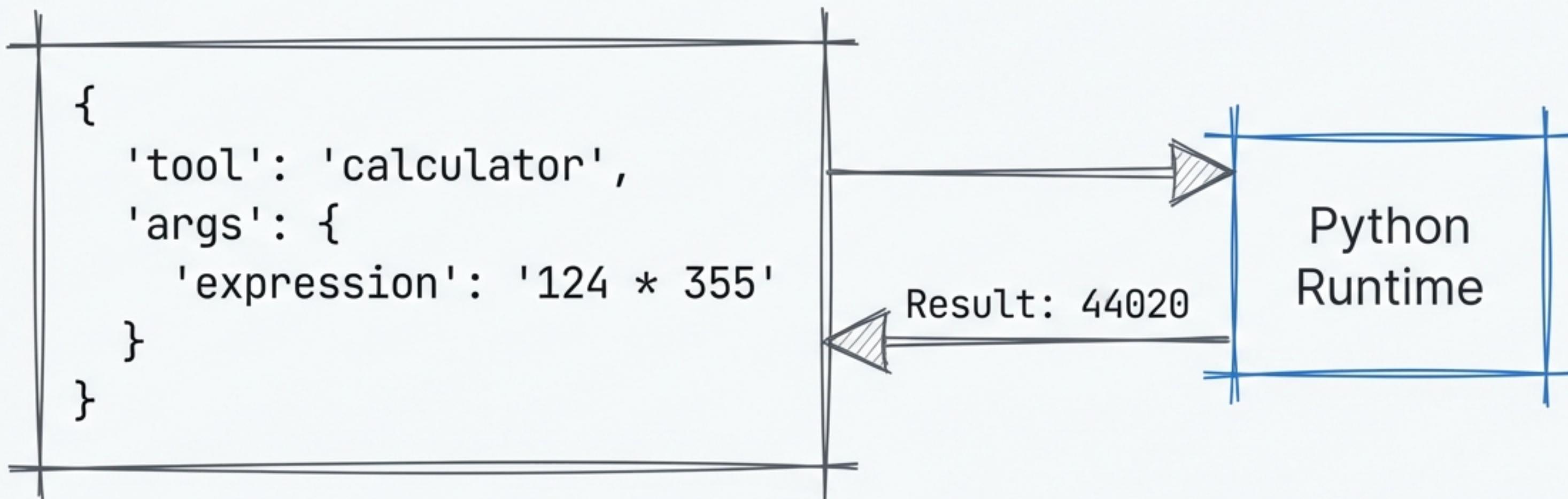
Step 3: Total = $3 + 5 = 8$.

Answer: 8



The Hands of the Agent: Tool Use

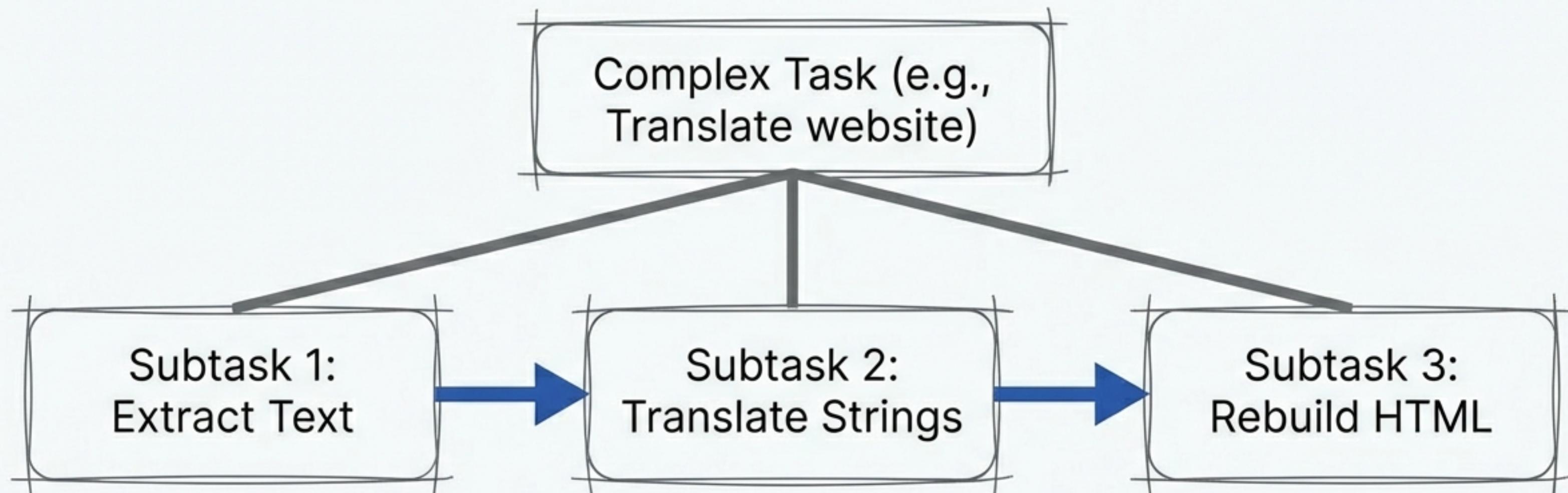
Model Controller Protocol (MCP)



The model outputs structured data to control external software.

Planning & Decomposition

Least-to-Most Prompting



Decompose complex problems into sequential dependencies.

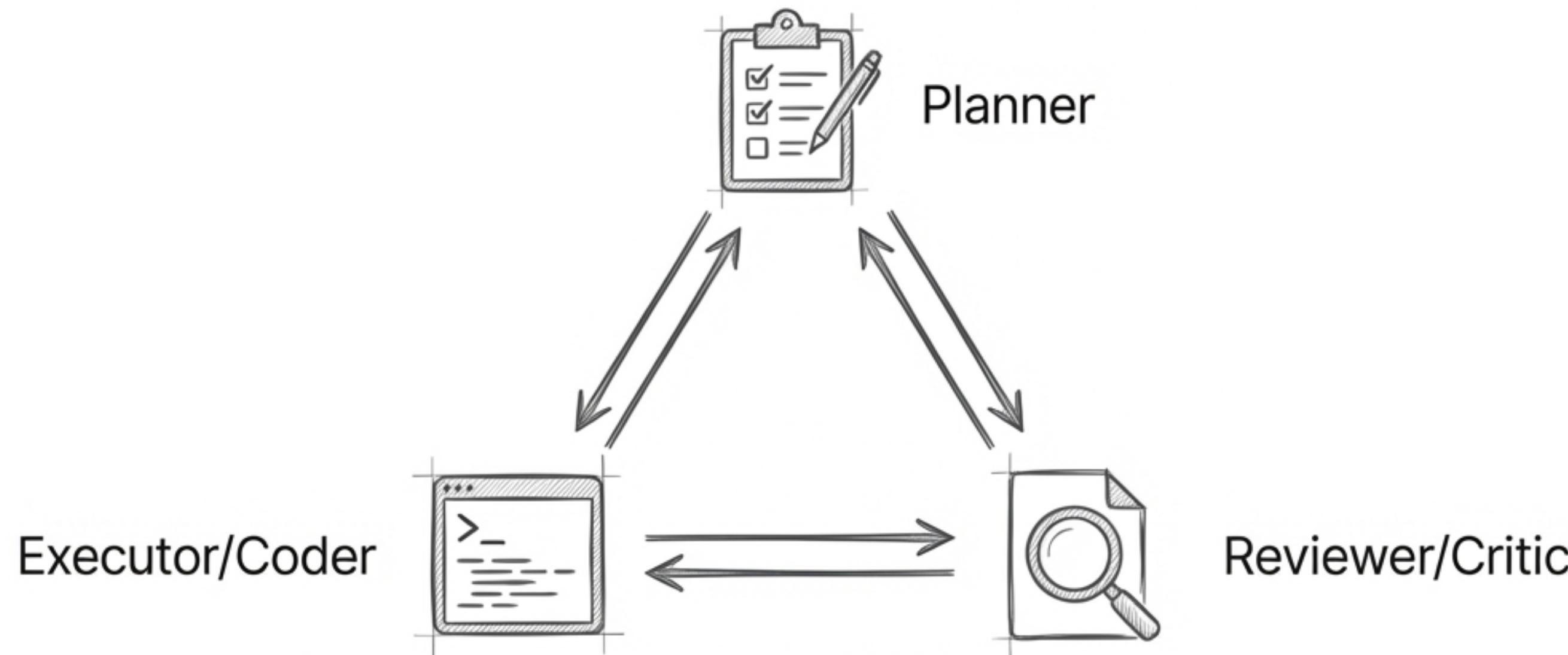
ReAct: Reasoning + Acting

```
> I need to find the CEO of Company X.  
> Search('Company X CEO')  
> Result: Jane Doe  
> Now I need to find her age.  
> Search('Jane Doe age')
```

Reasoning Trace
Inter Regular

Environment
Interaction
Inter Regular

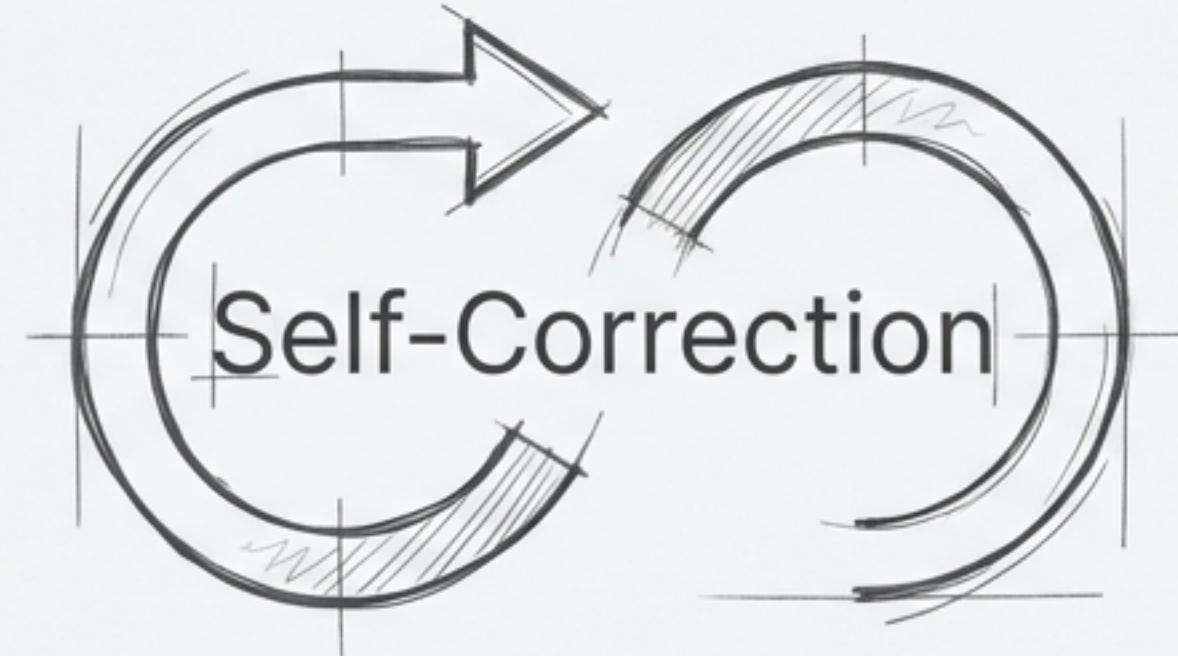
Multi-Agent Collaboration



Specialization allows for complex task solving
beyond a single context window.

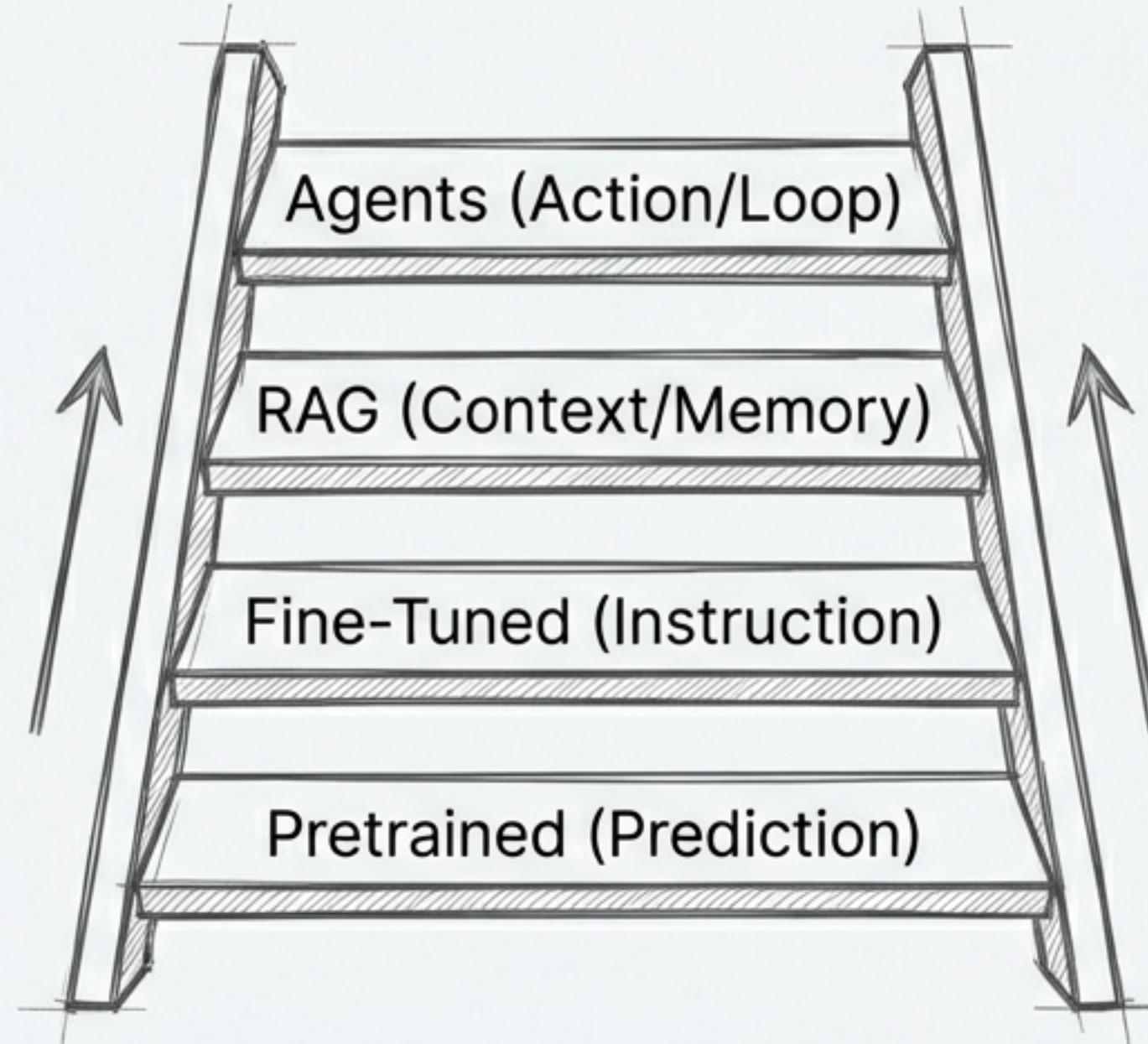
The Frontier: Self-Correction

The Challenge of Intrinsic Verification



- Models often hallucinate corrections.
- Risk of changing correct answers to incorrect ones.
- Future: External verifiers and Reward Models.

Summary: The Evolution of Cognition



Paradigm Shift: “What comes next?” → “**What should I do?**”

