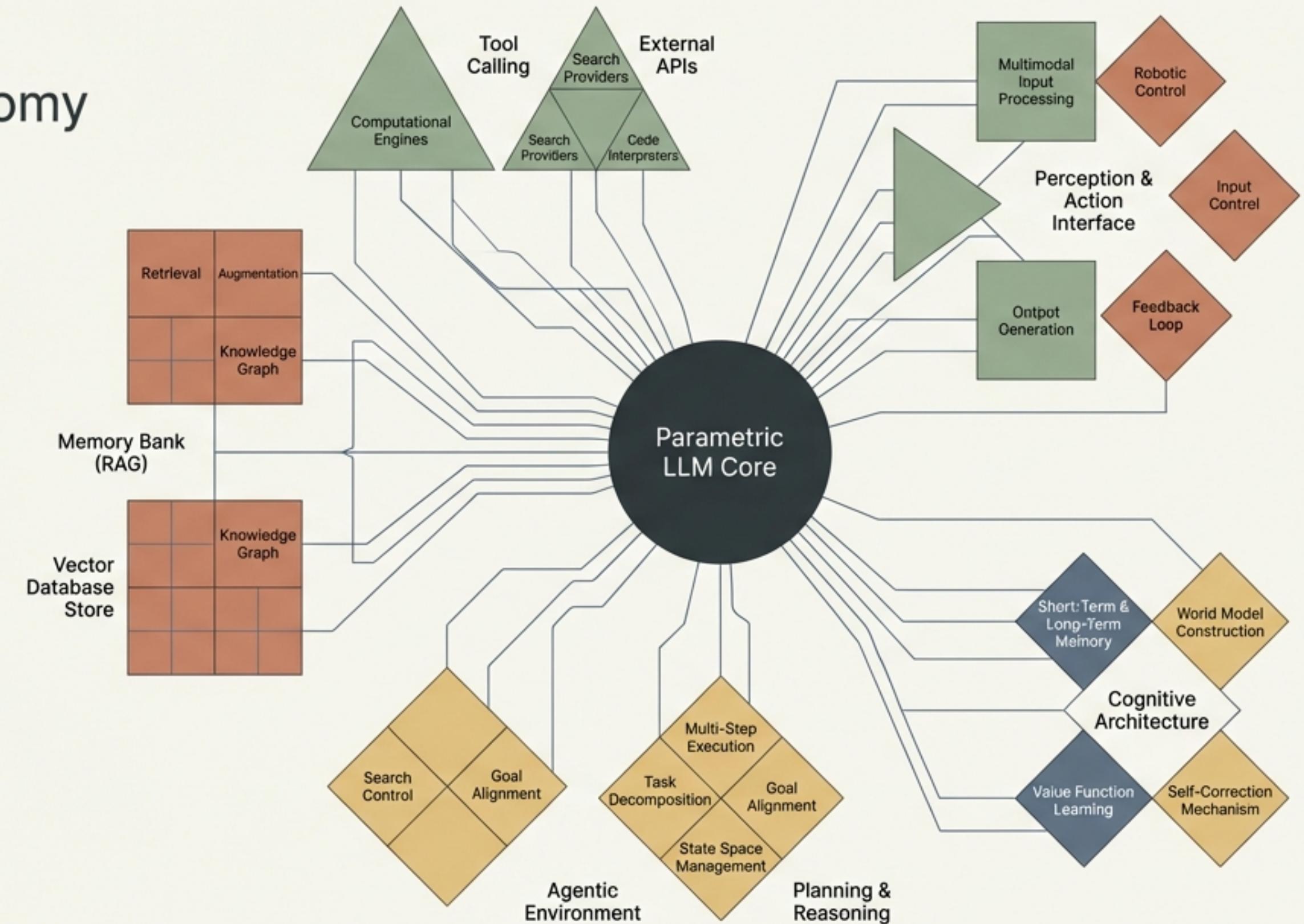


Augmented Large Language Models: The Evolution of Autonomy

A Comprehensive State-of-the-Art Treatment of RAG, Tool Calling, and Agents.



Inherent Constraints of Parametric Language Models

Knowledge Cutoff



$P_\theta(y | x)$ is conditioned only on training corpus D_{train} with temporal bound t_{max} .

Hallucination



Model assigns $P_\theta(y | x) > \epsilon$ to factually incorrect y due to distributional memorisation artifacts.

Computation Bottleneck



All reasoning must occur within fixed forward-pass depth L and hidden dimension d ; no external symbolic computation.

Static Knowledge



Parameters θ encode a snapshot; updating requires retraining with immense cost.

Groundedness Gap



No mechanism to verify generated claims against authoritative sources at inference time.

The Paradigm Shift to Augmented Large Language Models

Vanilla LLM:

$$P_{\theta}(\mathbf{y} \mid \mathbf{x}) = \prod P_{\theta}(y_t \mid \mathbf{y}_{<t}, \mathbf{x})$$

Augmented LLM (ALM):

$$P_S(\mathbf{y} \mid \mathbf{x}) = \prod P_{\theta}\left(y_t \mid \mathbf{y}_{<t}, \mathbf{x}, \bigoplus E_i(\phi_i(\mathbf{x}, \mathbf{y}_{<t}))\right)$$

$\phi_i(\cdot)$: Query formulation function constructing a request to external module E_i .

\bigoplus : Aggregation operator (concatenation, cross-attention fusion, interleaving, or gated mixture).

E_i : Non-parametric or external functional modules (retriever, tool API, code interpreter, memory buffer).

Constructing the Augmented Context Window

$$\mathbf{C}_{\text{aug}}(\mathbf{x}, t) = [\mathbf{x} \parallel \mathbf{E}_{\text{ret}}(\phi_{\text{ret}}(\mathbf{x}, \mathbf{y}_{<t})) \parallel \mathbf{E}_{\text{tool}}(\phi_{\text{tool}}(\mathbf{x}, \mathbf{y}_{<t})) \parallel \mathbf{M}_{<t}]$$



\mathbf{x} : Base Input

\mathbf{E}_{ret} : Retrieved
Documents

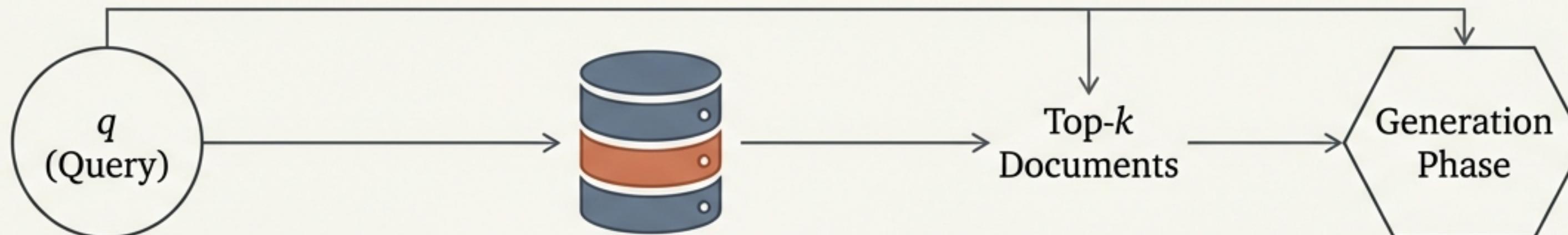
\mathbf{E}_{tool} : Tool Outputs

$\mathbf{M}_{<t}$: Memory Trace

Generation at each step becomes: $y_t \sim P_\theta(y_t | \mathbf{C}_{\text{aug}}(\mathbf{x}, t))$

Stage 1: Bridging the Groundedness Gap via Retrieval

Generation is conditioned dynamically on external documents retrieved at inference time.



Corpus $K = \{d_1, d_2, \dots, d_N\}$

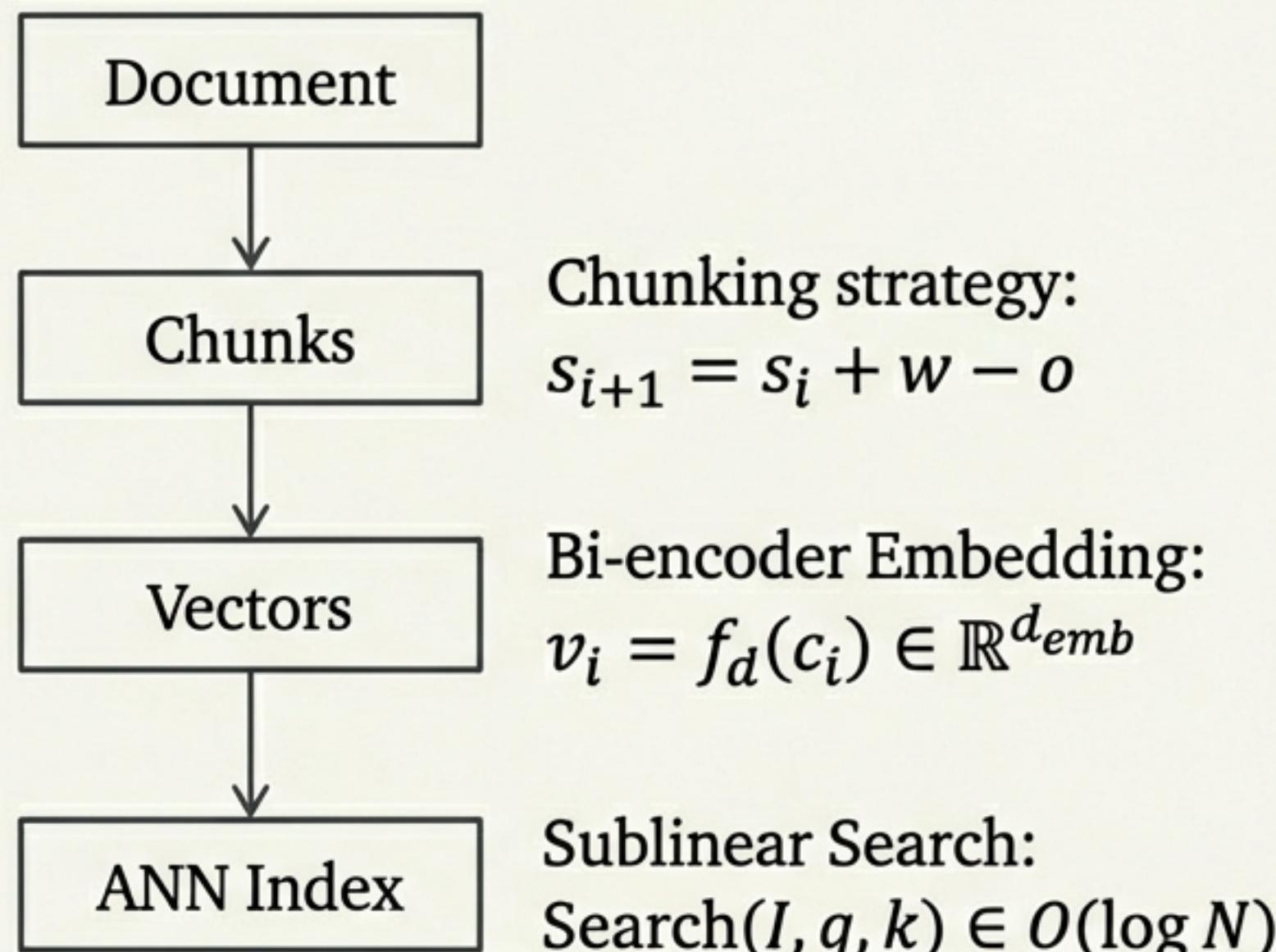
Top- k Marginalisation

$$P_{\text{RAG}}(y | q) \approx \sum \frac{\exp(\text{sim}(q, d)/\tau)}{\sum \exp(\text{sim}(q, d')/\tau)} \cdot P_\theta(y | q, d)$$

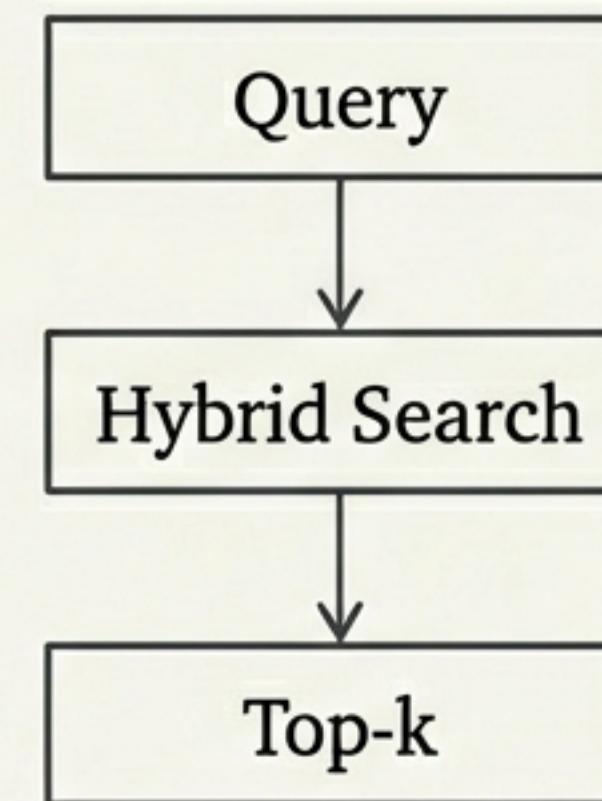
τ acts as the temperature parameter controlling retrieval distribution sharpness.

The Mechanics of the RAG Engine: Indexing and Retrieval

Offline Indexing Flow



Inference Retrieval Flow



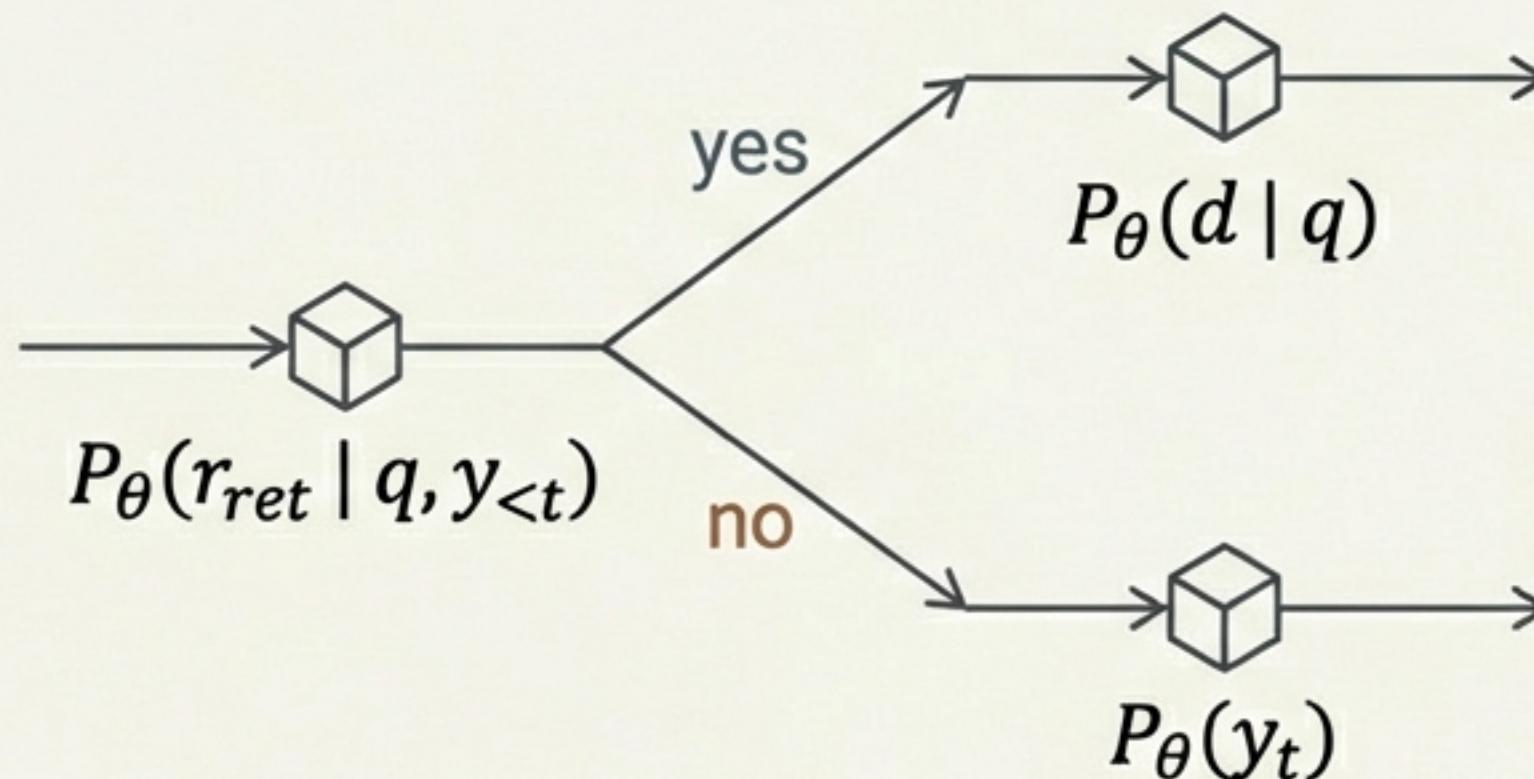
Hybrid Retrieval:

$$\begin{aligned} \text{score}_{\text{hybrid}}(q, d) &= \alpha \cdot \text{sim}_{\text{dense}}(q, d) \\ &+ (1 - \alpha) \cdot \text{BM25}_{\text{norm}}(q, d) \end{aligned}$$

Evolving RAG through Self-Reflection and Modularity



Self-RAG Conditional Logic:



$$P_{Self-RAG}(y, r | q) = \prod_{\theta} P_{\theta}(r_{ret} | q, y_{<t}) \cdot \begin{cases} P_{\theta}(d | q) \cdots & \text{if } r_{ret} = \text{yes} \\ P_{\theta}(y_t | q, y_{<t}) & \text{if } r_{ret} = \text{no} \end{cases}$$

Reflection Tokens:

?

r_{ret} : Retrieve? (yes/no)

\mathbb{S}

r_{sup} : Support Level?
(fully/partially/none)

✓

r_{rel} : Relevant? (yes/no)

★₅

r_{use} : Utility? (1-5 score)

Multi-Dimensional Evaluation of Retrieval and Generation

Retrieval

Context Precision

Context Recall

MRR

nDCG@k

Generation

Faithfulness (claims
grounded in context)

Answer Relevance

Hallucination Rate
(1 - Faithfulness)

End-to-End

Answer Correctness
(Semantic Similarity +
Factual Overlap)

RAGAS Framework

$RAGAS_score = \text{HarmonicMean}(Faithfulness, AnswerRelevance, ContextPrecision, ContextRecall)$

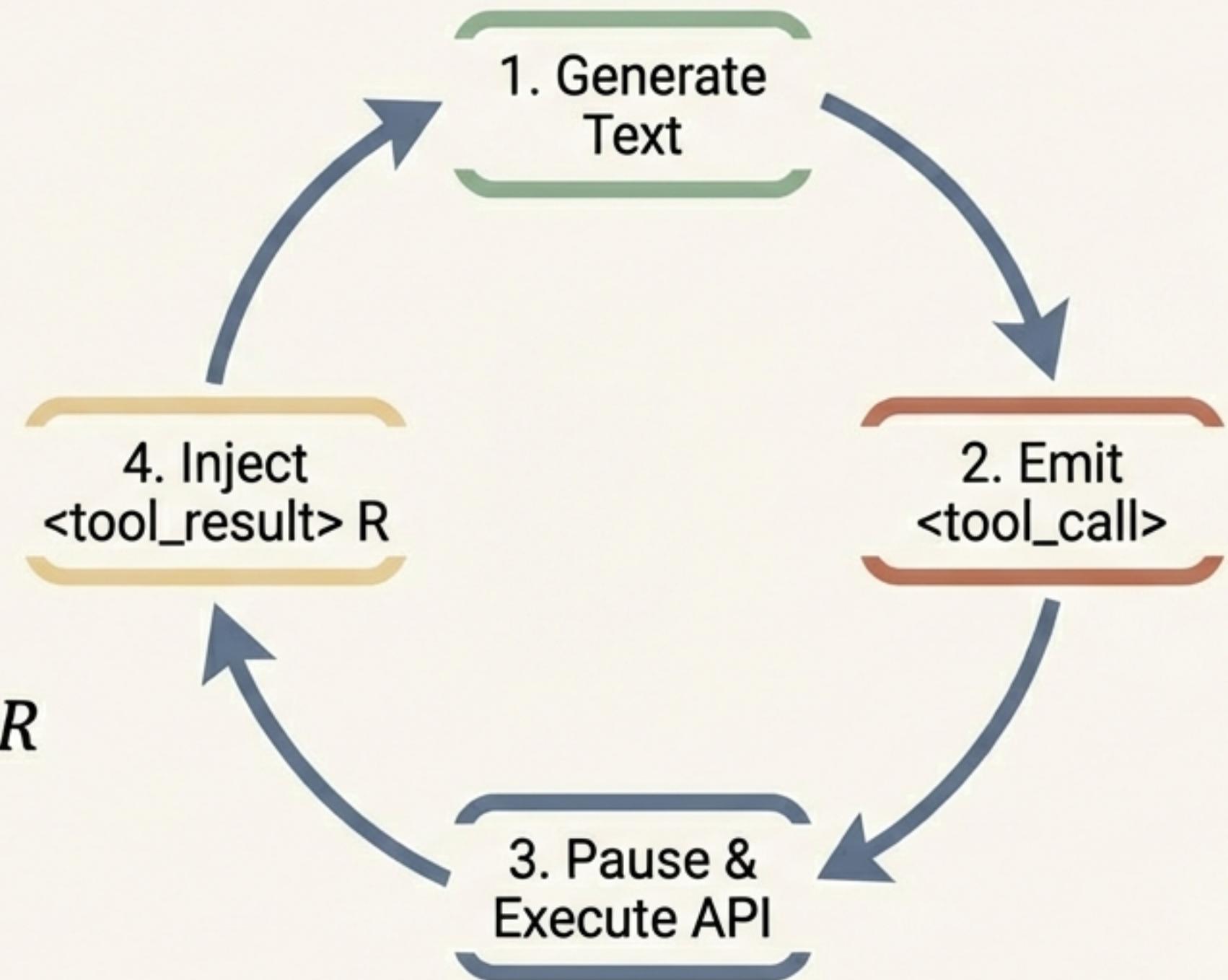
Stage 2: Enabling Computation via Tool Calling

Vocabulary Augmentation

$$V_{aug} = V \cup V_{tool}$$

Conditional Loop

$$y_t \sim \begin{cases} P_\theta(y_t | y_{<t}, x, R) & \text{if } y_{t-1} \notin V_{tool} \\ \text{EXECUTE(parse_tool_call}(y_{<t})) \rightarrow R & \text{if } y_{t-1} = </\text{tool_call}> \end{cases}$$



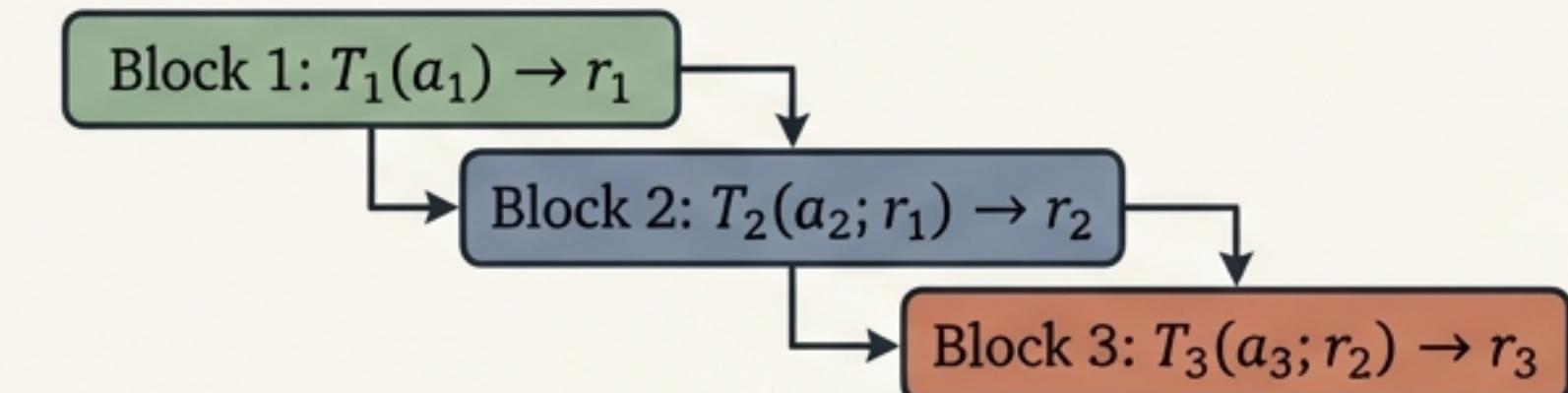
Structuring and Optimising Tool Execution

Tool Specification Tuple:

$$T = (\text{name}, \text{description}, \text{parameters}, \text{returns})$$

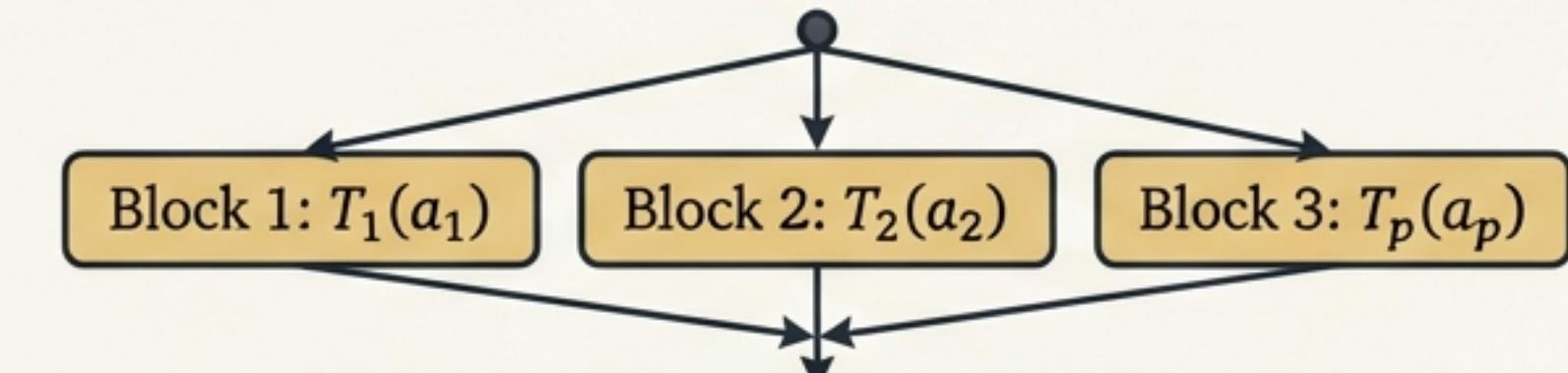
- name
- description
- parameters
- returns

Sequential Execution



$$y = M_\theta(x, r_1 = T_1(a_1), r_2 = T_2(a_2; r_1), \dots)$$

Parallel Execution



$$\{r_1, \dots, r_p\} = \text{PARALLEL_EXEC}(T_1(a_1), \dots, T_p(a_p))$$

Training Paradigm: Tool-use is learned via Supervised Fine-Tuning (\mathcal{L}_{SFT}) strictly on model-generated tokens, or DPO optimising for correct tool trajectories.

Stage 3: LLM Agents as Self-Directing Cognitive Engines

An autonomous system

$$A = (M_\theta, P, T, E, S)$$

Action Mapping Equation:

$$\begin{aligned} a_t &= \pi_\theta(o_{\leq t}, m_t, g) = \\ &= M_\theta(\text{PROMPT}(o_{\leq t}, m_t, g, T)) \end{aligned}$$

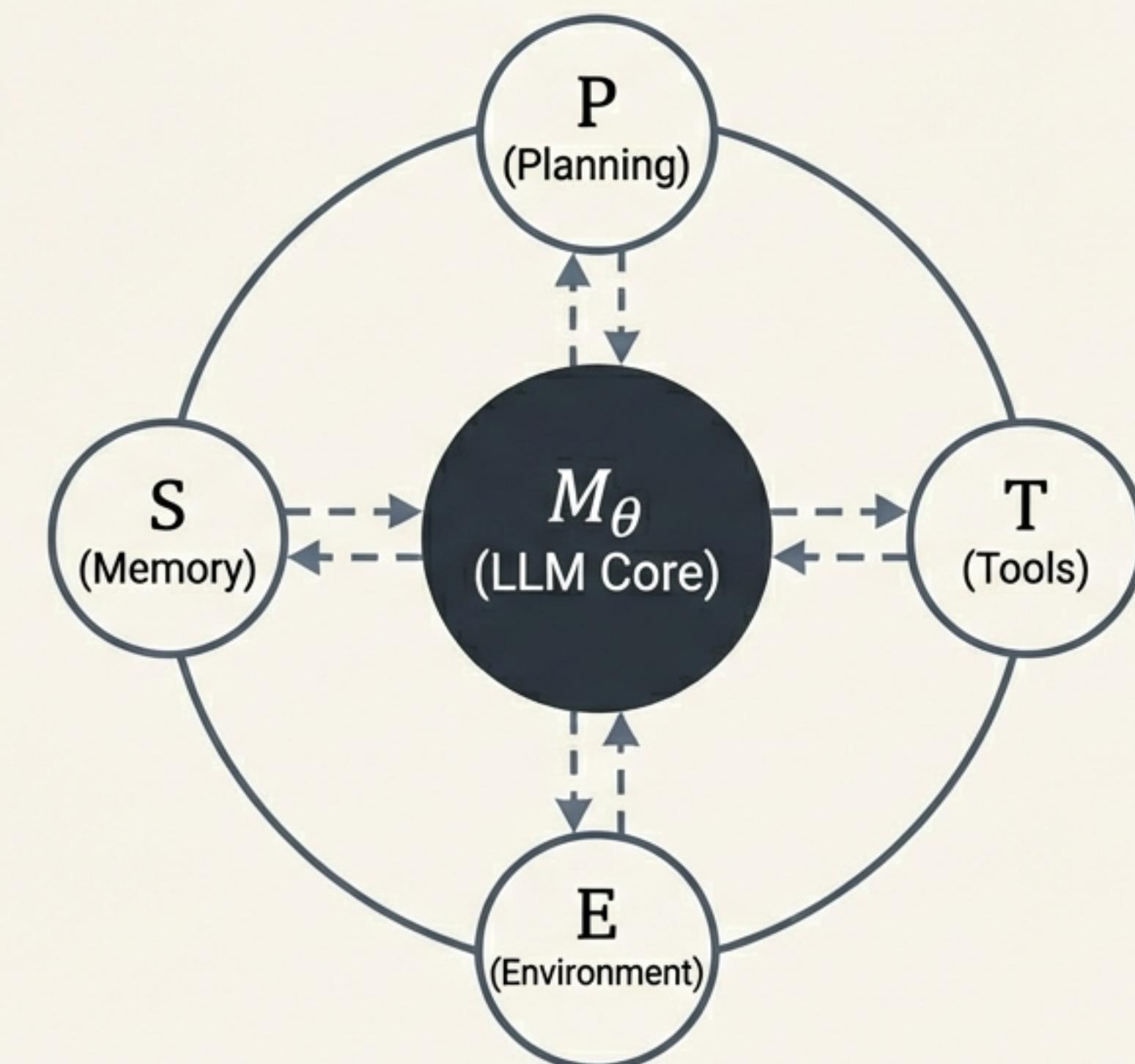
$o_{\leq t}$: observation history

m_t : memory state

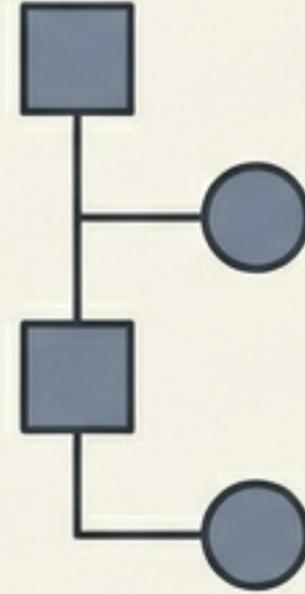
g : goal

a_t : action (think, respond, tool_call, delegate)

Cognitive Architecture



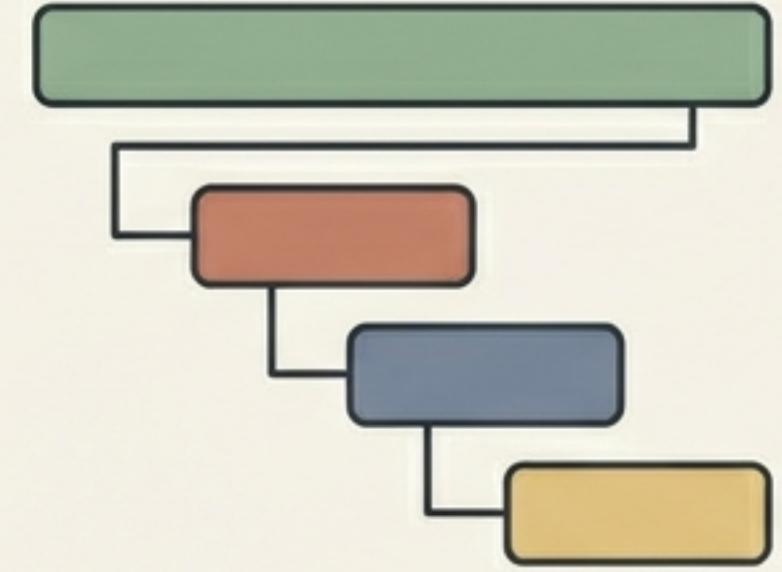
Architectures for Autonomous Agent Reasoning



1. ReAct

Interleaves reasoning and actions.

Trajectory =
 $(t_1, a_1, o_1, t_2, a_2 \dots)$



2. Plan-and-Execute

Separates planning phase from execution.

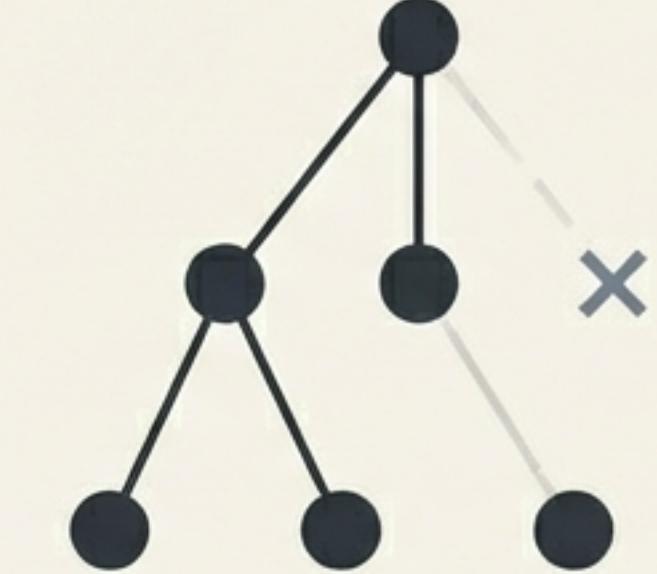
$$G = \text{PLANNER}(M_\theta, q)$$



3. Reflexion

Explicit self-evaluation loop.

$$\begin{aligned} S_{\text{memory}} &\leftarrow S_{\text{memory}} \\ &\cup \{\text{Reflection}_t\} \end{aligned}$$



4. Tree of Thoughts (ToT)

Explores multiple paths via search.

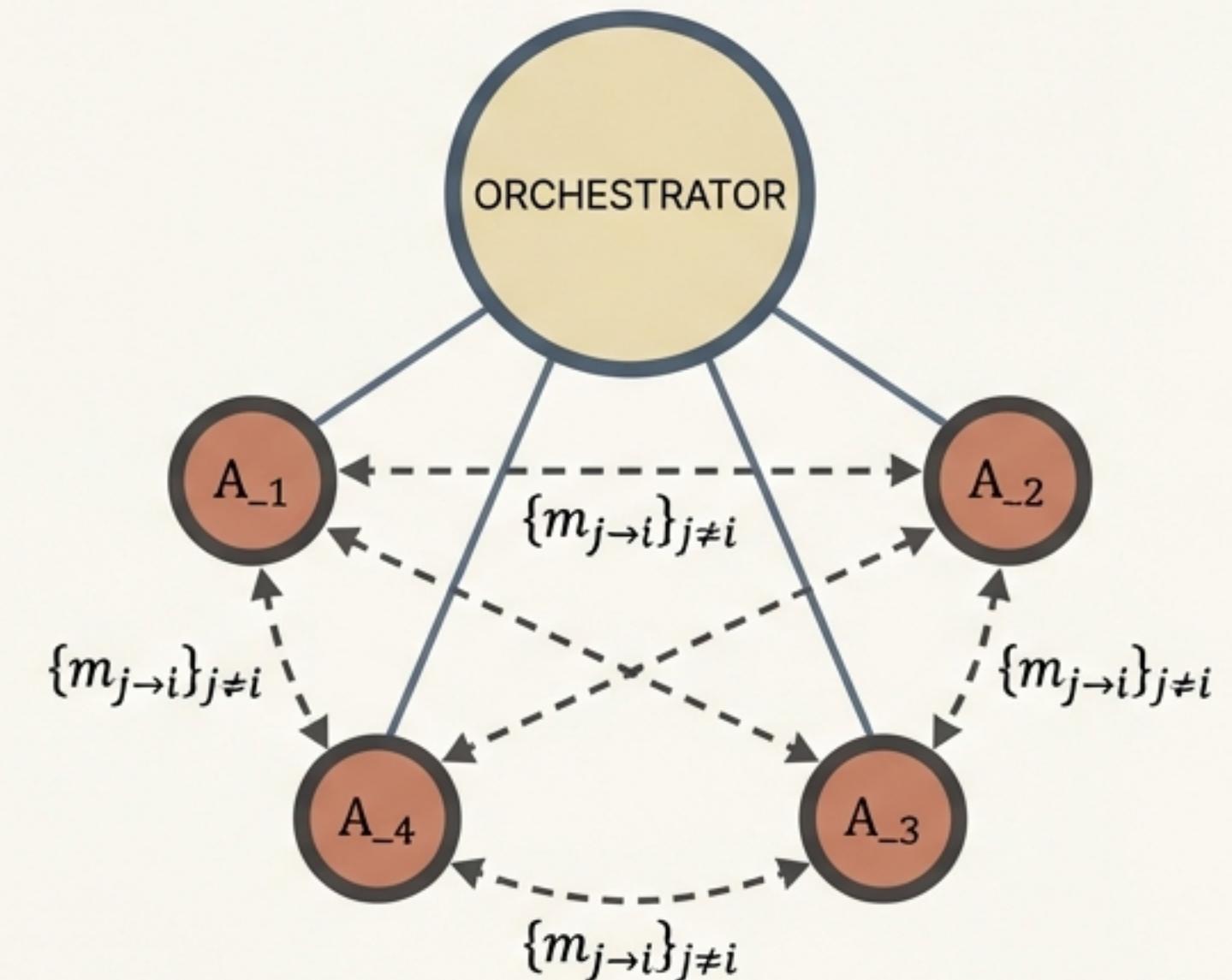
Prunes where
 $V(s) < \delta$

Sustaining Context and Scaling via Multi-Agent Systems

Memory Taxonomy

Working Memory	Context window (Per-session)
Episodic Memory	Action-observation traces (Cross-session)
Semantic Memory	Vector database of facts (Permanent)
Procedural Memory	Fine-tuned weights/prompts (Permanent)

Multi-Agent Coordination



$$y = \text{ORCHESTRATOR}\left(\{r_i\}_{i=1}^n\right)$$

where $r_i = A_i(x_i, \text{shared_state}, \{m_{j \rightarrow i}\}_{j \neq i})$

Quantifying Autonomous Performance and Efficiency

Task Success Rate

$$\frac{\text{Number of goals achieved}}{\text{Total number of goals}}$$

Average Steps to Completion

$$E[\text{Length of valid trajectory } \tau_{\text{success}}]$$

Tool Call Accuracy

$$\frac{\text{Validly formatted tool executions}}{\text{Total tool calls emitted}}$$

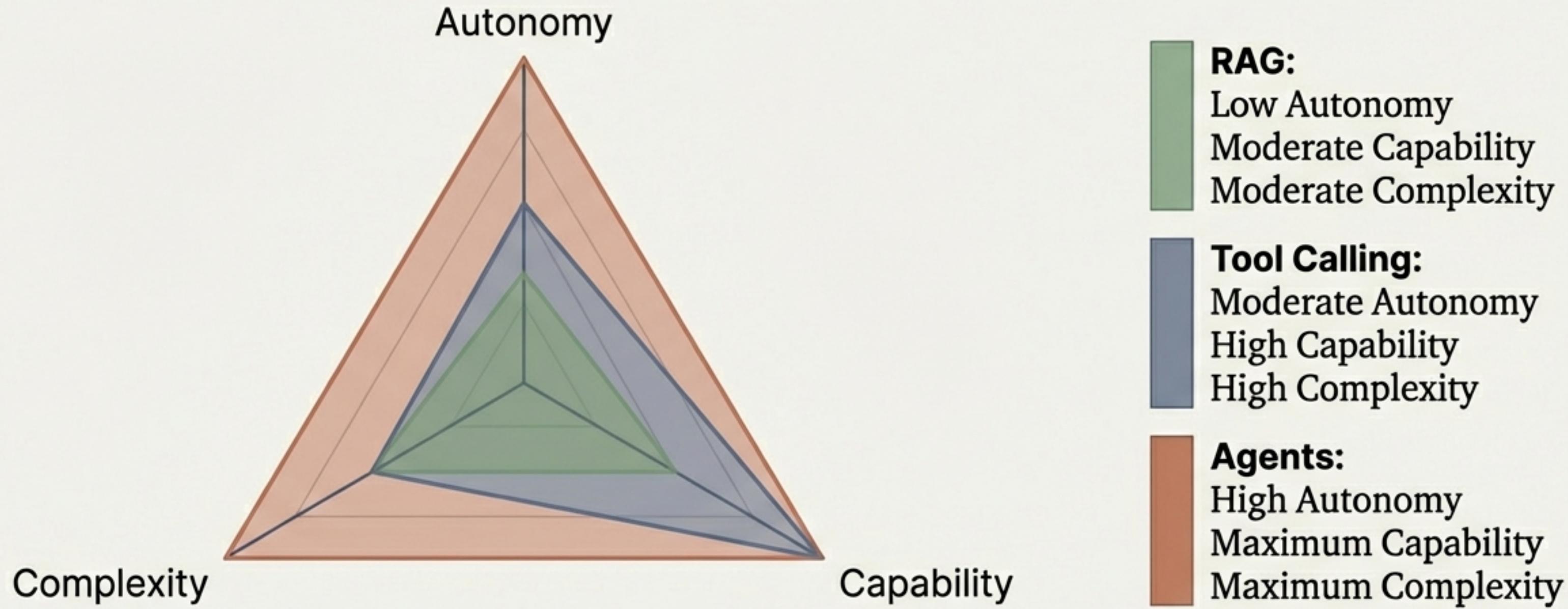
Reflection Improvement Rate

$$\frac{(\text{score}_{\text{trial}_{t+1}} - \text{score}_{\text{trial}_t})}{\text{score}_{\text{trial}_t}}$$

Cost Efficiency

$$\frac{\text{task_success_rate}}{\text{total_tokens_consumed}}$$

The Unified View of Augmented Architectures



The progression from Retrieval to Agents represents exponential gains in autonomy and capability, traded against engineering complexity. The optimal operating point is dictated by the specific requirements for task accuracy, inference latency, compute cost, and the required degree of deterministic trust.