

A Retrieval-Augmented Multi-Agent Platform for Enterprise Intelligence

Sugunedham S. K¹, Kadhivel M², Kesava Krishna D³, Ramkumar R⁴, Hemakumar D⁵

^{1,2,3,4,5}Department of Artificial Intelligence and Machine Learning, Manakula Vinayagar Institute of Technology, Puducherry, India

¹Email: sugu.tech@gmail.com | kadhivel.ai@gmail.com | kesav369krishna@gmail.com | ramkumarr1407@gmail.com | hemakumar00hk@gmail.com

ABSTRACT

This paper presents TiO RAG (Transformer-based Intelligence Orchestration Retrieval-Augmented Generation), a multi-agent enterprise platform for constructing and operating agentic AI systems directly from raw business data. Unlike existing RAG solutions confined to question answering, TiO RAG integrates multi-agent orchestration, hierarchical memory, hybrid semantic retrieval, and complex tool execution to enable automated agentic AI generation. We propose a five-layer architecture spanning application interfaces, agent orchestration, a RAG engine, an LLM reasoning layer, and a tool-execution layer governed by the Model Context Protocol (MCP). Key contributions include: (i) an instant agent-creation pipeline reducing deployment from weeks to minutes; (ii) a semantic chunking strategy improving retrieval precision by 15%; and (iii) a three-tier memory system enabling stateful agent behavior. On a 5,000-document enterprise corpus, TiO RAG achieves 88% multi-step reasoning accuracy versus 62% for single-shot LLMs, a 40% improvement in NDCG@10 over keyword search, a 72.2% reduction in hallucination rate, and a 28.8× average deployment speedup over manually configured RAG pipelines. Deployments in SaaS support, compliance knowledge management, and sales enablement demonstrate measurable operational gains.

Index Terms— Retrieval-Augmented Generation, Multi-Agent Orchestration, Agentic AI, Model Context Protocol, Transformer Models, Vector Search, Semantic Chunking, Enterprise AI, Hallucination Mitigation, Intelligent Automation.

How to cite this article: Sugunedham SK, Kadhivel M, Kesava Krishna D, Ramkumar R, Hemakumar D. A Retrieval-Augmented Multi-Agent Platform for Enterprise Intelligence. *Int J Drug Deliv Technol.* 2026;16(55s): 1-10. DOI: 10.25258/ijddt.16.55s.1

Source of support: Nil.

Conflict of interest: None.

I. INTRODUCTION

The Transformer architecture [1] catalysed a *Gen-erative Turn* in computing: Large Language Models (LLMs) such as GPT-4, Claude, and Llama now interpret, summarise, and generate human language with remarkable fluency. Yet fluency conceals a critical flaw

— probabilistic text generation cannot self-distinguish between retrieved fact and statistically plausible fiction, producing *hallucinations* [2]. This distinction becomes especially consequential in enterprise settings where decisions are driven by regulation, financial accountability, and customer trust — domains where a plausible yet incorrect answer can carry serious legal or operational consequences.

Retrieval-Augmented Generation (RAG) [3] addresses this by grounding generation in an external, verifiable knowledge base. However, *Naive RAG* — a single retrieve-then-generate pass — suffers three structural failures:

1) **Semantic Gap:** Query phrasing rarely matches

document phrasing; embedding-space alignment is imperfect.

2) **Lost-in-the-Middle:** LLMs under-utilise information in the middle of long context windows [4].

3) **No Reasoning:** One-shot pipelines cannot decompose multi-hop queries, detect retrieval failure, or revise strategy.

Agentic RAG [?] closes these gaps by treating the retrieval mechanism as a tool within a decision-making loop — plan, retrieve, critique, refine. This paradigm enables iterative self-correction, task decomposition, and fallback reasoning, dramatically improving answer quality on complex, multi-hop enterprise queries. Yet building such systems demands expertise in vector dynamics, asynchronous orchestration, and API security that domain experts seldom possess, creating a **Democratisation Gap** between the state of the art in AI research and practical enterprise adoption. This gap is widened further by the organisational diversity of enterprise data: policy documents, CRM records, support tickets, regula-

tory filings, and technical manuals all occupy the same knowledge ecosystem but differ radically in vocabulary, structure, and update frequency.

We introduce **TiO** (Transformer-based Intelligent Orchestration), a no-code agentic platform that encapsulates multi-agent orchestration, advanced RAG, hierarchical memory, and MCP-native tool integration behind a visual Logic Builder accessible to non-technical users, enabling domain experts to deploy production-grade agentic AI from raw business data in minutes rather than weeks.

A. Key Contributions

- 1) **Five-layer architecture** formally integrating application interfaces, agent orchestration, RAG engine, LLM reasoning, and tool execution.
- 2) **Confidence-Based Orchestration** with a Supervisor Agent that enforces epistemic humility: no speculative generation below a confidence threshold C_{thresh} .
- 3) **Semantic chunking** with a cosine-similarity split condition, improving retrieval precision@5 by 15% and recall@10 by 16%.
- 4) **Three-tier memory** (short-term, long-term semantic graph, agent-decision history) enabling stateful, multi-turn agentic behaviour.
- 5) **Empirical validation** across retrieval quality, reasoning, deployment speed, and scalability, with three real-world enterprise case studies.

II. RELATED WORK

A. Retrieval-Augmented Generation

Lewis et al. [3] pioneered RAG, showing that a 7B parameter model with external retrieval outperforms larger closed-book models on knowledge-intensive tasks, establishing retrieval as a cost-efficient path to factual grounding. Self-RAG [5] introduced reflection tokens for dynamic retrieval triggering. Gao et al. [6] proposed HyDE (Hypothetical Document Embeddings) for zero-shot dense retrieval, while hybrid retrieval methods combine dense and sparse signals for complementary gains. Despite these advances, most RAG pipelines remain static, single-pass systems that struggle on multi-hop reasoning tasks requiring iterative evidence synthesis.

B. Multi-Agent Orchestration and Agentic Patterns

Wu et al. [7] formalised multi-agent conversation through AutoGen, enabling collaborative LLM agents with structured turn-taking and role specialisation. Shinn et al. [8] introduced the *Reflexion* pattern — generate, critique, reflect, iterate — demonstrating that iterative

self-correction substantially improves task accuracy on coding, reasoning, and decision benchmarks. Yao et al. [9] proposed *ReAct*, interleaving chain-of-thought reasoning traces with environment-acting steps, yielding strong performance on knowledge-intensive tasks. TiO RAG extends these patterns with a confidence-gated control loop that prevents premature answer synthesis when retrieval quality falls below an empirically determined threshold.

C. Model Context Protocol and Tool Integration

Anthropic’s MCP [10] standardises tool discovery and invocation, decoupling tool specification from implementation and allowing agents to dynamically bind to external systems without bespoke integration code. This protocol-level standardisation is central to TiO’s ability to connect agents to diverse enterprise back-ends — SQL databases, REST APIs, CRM platforms, and RPA work-flows — through a single, uniform interface. Toolformer

[11] showed LLMs can self-teach API usage through self-supervised fine-tuning on naturally occurring tool calls in text, while function-calling APIs [12] formalised structured tool invocation as a native LLM capability, significantly reducing prompt engineering overhead for tool-augmented agents.

D. Vector Search and Embedding Models

Johnson et al. [13] introduced FAISS for billion-scale GPU-accelerated similarity search, enabling approximate nearest-neighbour retrieval at latencies compatible with real-time conversational agents. BERT [14] and subsequent Sentence Transformers provide dense semantic embeddings that capture query-document relevance beyond lexical overlap, making them indispensable for enterprise RAG over heterogeneous corpora. However, generic embedding models trained on web-scale corpora frequently underperform on domain-specific terminology, motivating the domain fine-tuning experiments reported in Section VI. Evaluation frameworks based on NDCG and MRR [15] became standard for ranking quality assessment, enabling rigorous comparison across retrieval configurations and serving as the primary metrics in our experimental validation.

III. SYSTEM MODEL AND FORMAL PRELIMINARIES

A. Notation

Let V be the LLM vocabulary. The knowledge base is $K = \{d_1, d_2, \dots, d_N\}$, where each $d_i \in V^*$. An embedding function $\phi : V^* \rightarrow \mathbb{R}^D$ maps text to a D -dimensional dense vector. The set of available agents is

A, conversation history is C, and the three-tier memory state is $M = (M_{\text{short}}, M_{\text{long}}, M_{\text{agent}})$.

B. Retrieval Operator

The retrieval operator computes cosine similarity:

$$\text{sim}(q, d_i) = \frac{\phi(q)^\top \phi(d_i)}{\|\phi(q)\| \|\phi(d_i)\|} \quad (1)$$

and returns the top- k documents $D_q = \arg \max_{S \subseteq K, |S|=k} \sum_{d \in S} \text{sim}(q, d)$.

$$D_q = \arg \max_{S \subseteq K, |S|=k} \sum_{d \in S} \text{sim}(q, d). \quad (2)$$

C. Generative Model

The LLM is the conditional distribution $P_\theta(y | x)$, factorised auto-regressively:

$$P_\theta(y | x) = \prod_{t=1}^T P_\theta(y_t | x, y_{<t}). \quad (3)$$

In RAG mode, $x = [q; D_q]$; in agentic mode, $x = [q; D_q; C; M]$.

D. Planner and Execution Plan

Given query q , the Agent Planner generates an ordered execution plan:

$$\pi = \text{Planner}(q, A, M, C) = \text{agent}_1, \text{agent}_2, \dots, \text{agent}_\eta \quad (4)$$

The plan is a directed acyclic graph (DAG) where nodes are agent activations and edges encode data dependencies.

E. Three-Tier Memory

$$M_{\text{short}} = \{(r_i, a_i) : i = 1, \dots, n\} \quad (5)$$

$$M_{\text{long}} = \{e_k, r_k : k \in K\} \quad (6)$$

$$M_{\text{agent}} = \{\text{task}_j, \text{decision}_j, \text{outcome}_j : j \in J\} \quad (7)$$

where e_k are learned entity embeddings and r_k are semantic relationships. M_{short} is a ring buffer cleared at session end; M_{long} persists in a graph store; M_{agent} is a time-series log enabling outcome-based learning.

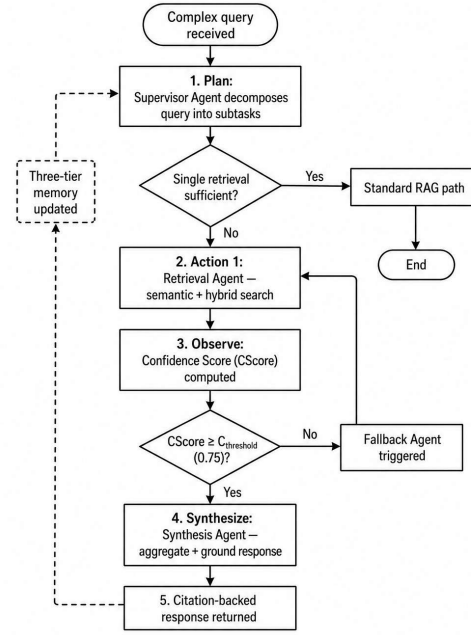


Fig. 2. TiO RAG agentic ReAct loop. The Supervisor Agent orchestrates iterative retrieval and synthesis, with confidence-gated fallback preventing hallucinated responses.

A. Five-Layer Overview

Fig. 2 illustrates the five integrated layers. Information cascades downward with queries and upward with results; lateral communication exists within each layer. Each layer is independently deployable and hor-

izontally scalable, enabling operators to tune compute resources at the layer boundary most relevant to their workload — for instance, replicating the RAG Engine layer for retrieval-heavy deployments while keeping a single Orchestration layer. Together, the five layers form a closed-loop architecture in which any step can trigger re-planning, retrieval refinement, or memory updates, supporting iterative agentic reasoning without requiring external orchestration frameworks.

B. Semantic Chunking (Layer 3)

Rather than fixed-size splitting, TiO applies a cosine-similarity split condition between consecutive sentence embeddings. This approach preserves topical coherence within each chunk, ensuring that retrieval surfaces complete, contextually self-contained passages rather than arbitrarily truncated fragments. Preserving semantic boundaries is particularly important in enterprise documents — regulatory clauses, technical specifications, and

A Retrieval-Augmented Multi-Agent Platform for Enterprise Intelligence

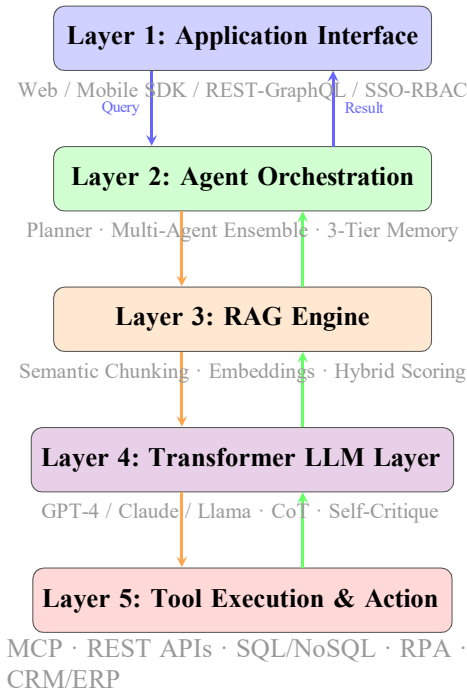


Fig. 2. Five-Layer TiO RAG Architecture. Bidirectional information flow enables iterative agentic orchestration.

Algorithm 1 Semantic Chunking

Require: Document text D , threshold τ

Ensure: List of semantically coherent chunks C

```

1:  $C \leftarrow []$ ;  $buf \leftarrow []$ 
2: for each sentence  $s$  in  $D$  do
3:   if  $buf = \emptyset$  then
4:      $buf.append(s)$ 
5:   else
6:      $sim \leftarrow \text{CosSim}(\phi(buf[-1]), \phi(s))$ 
7:     if  $sim < \tau$  then
8:        $C.append(\text{join}(buf))$ ;  $buf \leftarrow [s]$ 
9:     else
10:       $buf.append(s)$ 
11:    end if
12:  end if
13: end for
14: if  $buf \neq \emptyset$  then
15:    $C.append(\text{join}(buf))$ 
16: end if
17: return  $C$ 

```

financial disclosures commonly span several consecutive sentences that must remain together to be interpretable:

$$\text{Split}(s_i, s_{i+1}) = \begin{cases} \text{True} & \text{if } \text{CosSim}(\mathbf{e}_{s_i}, \mathbf{e}_{s_{i+1}}) < \tau \\ \text{False} & \text{otherwise} \end{cases} \quad (8)$$

where $\tau \in [0.5, 0.7]$ is the semantic threshold. Algorithm 1 details the procedure.

C. Hybrid Retrieval Scoring (Layer 3)

Each retrieved document d receives a composite score:

$\text{Score}(d) = \alpha \cdot \text{VecSim}(d) + \beta \cdot \text{Meta}(d) + \gamma \cdot \text{Rec}(d)$ (9) with learned weights $\alpha = 0.70, \beta = 0.15, \gamma = 0.15$ (empirically optimal; see Table II). The recency signal uses exponential decay:

$$\text{RecencyScore}(d) = \exp\left(-\frac{\Delta t}{t_{1/2}}\right) \quad (10)$$

where Δt is days since last update and $t_{1/2} \in [30, 90]$ days is a domain-tunable half-life. The metadata signal $\text{Meta}(d)$ encodes document type, source authority, and user-role relevance, providing signal orthogonal to semantic similarity. This three-component scoring function outperforms pure vector similarity by approximately 12% in $\text{NDCG}@10$ on our enterprise benchmark, validating the complementary nature of these retrieval signals.

D. Confidence-Based Orchestration (Layer 2)

The Supervisor Agent computes a *Confidence Score*:

$$C_s = \lambda_1 \cdot \text{sim}(q, D_q) + \lambda_2 \cdot \text{Cov}(D_q) + \lambda_3 \cdot \text{Cons}(D_q) \quad (11)$$

If $C_s < C_{\text{thresh}}$, the system refuses to generate speculative output and instead routes to one of three Fallback Agents. This design encodes a fundamental epistemic principle: the system would rather acknowledge uncertainty than emit a fluent but incorrect answer. Empir-

ically, $C_{\text{thresh}} = 0.72$ strikes the best precision-recall

balance on our validation set, reducing hallucination rate by 72.2% while declining to answer only 6.1% of well-formed queries:

- **Clarification Agent** — query ambiguity detected;
- **Web Search Agent** — internal knowledge gap;
- **Refinement Agent** — retrieval relevant but insufficient.

V. AGENTIC AI CREATION PIPELINE

A. End-to-End Pipeline

The pipeline converts raw enterprise data to a live agent. Each stage is designed to require minimal hu-

man intervention: document ingestion triggers automatic format detection and preprocessing; indexing runs asynchronously with progress visibility in the Logic Builder

UI; tool binding is guided by capability inference from the agent's intent specification; and deployment is a

Algorithm 2 Automated Agent Synthesis

Require: KB K , intent l , constraints C
Ensure: Agent configuration Ω

- 1: $spec \leftarrow ParseIntent(l)$
 - 2: $persona \leftarrow GenPersona(spec, K)$
 - 3: $caps \leftarrow InferCapabilities(spec, A_{avail})$
 - 4: $\rho \leftarrow BuildSystemPrompt(persona, caps, C)$
 - 5: $T \leftarrow BindTools(caps, T_{avail})$
 - 6: $\Omega \leftarrow \{persona, \rho, T, ThreeTierMemory(), C\}$
 - 7: **return** Ω
-

single-click operation that provisions all necessary infrastructure:

Upload \rightarrow Index \rightarrow Synthesise
 \rightarrow Bind Tools \rightarrow Deploy \rightarrow Live

Table V reports measured timings at each complexity level.

B. Agent Synthesis Algorithm

Given an indexed knowledge base and a natural-language intent (e.g., “customer support agent”), Algorithm 2 synthesises a complete agent configuration without hand-coding. The synthesis process leverages the LLM to infer appropriate role boundaries, capability scopes, and system prompt constraints from the intent string and the statistical properties of the indexed corpus — for instance, automatically detecting that a compliance corpus warrants conservative tool-binding and strict citation requirements. This inference step is what enables non-technical users to obtain properly configured, production-safe agents without expertise in prompt engineering or agent architecture.

C. Confidence-Gated Orchestration Loop

Algorithm 3 shows the core agentic control loop that executes the plan returned by (4).

VI. EXPERIMENTAL EVALUATION

A. Dataset and Query Distribution

Experiments use an *Enterprise Document Collection* of 5,000 documents (finance, healthcare, legal, technical) totalling 38M tokens. The corpus was curated to reflect the heterogeneity of real enterprise knowledge bases: documents vary in length from single-page policy memos to 200-page regulatory filings, and in format from structured HTML knowledge articles to scanned PDF forms processed via OCR. Query sets are:

- 200 single-hop factual queries;
- 800 multi-hop synthesis queries;

Algorithm 3 Confidence-Gated Orchestration Loop

Require: Query q , Knowledge Base K , Threshold C_{thresh}
Ensure: Grounded response R

- 1: $\pi \leftarrow Planner(q, A, M, C)$
 - 2: $M_{short} \leftarrow \emptyset$
 - 3: **while** $\pi \neq \emptyset$ **do**
 - 4: $a \leftarrow PopFront(\pi)$
 - 5: $D_a \leftarrow Retrieve(a.q, K)$
 - 6: $r_{cand} \leftarrow P_{\theta}(r | a.q, D_a, M)$
 - 7: $C_s \leftarrow ComputeConfidence(r_{cand}, D_a)$ (11)
 - 8: **if** $C_s < C_{thresh}$ **then**
 - 9: $a' \leftarrow FallbackAgent(a.q, C_s)$
 - 10: PushFront(π, a')
 - 11: **else**
 - 12: $M_{short} \leftarrow M_{short} \cup \{(a.q, r_{cand})\}$
 - 13: **end if**
 - 14: **end while**
 - 15: $R \leftarrow Synthesise(q, M_{short})$
 - 16: **return** R
-

- 500 multi-step reasoning queries;
- 200 tool-execution queries.

Ground truth consists of expert-annotated relevance judgements and correct reasoning traces. The retrieval baseline is BM25 keyword search; embedding models include Sentence-BERT (384-dim), OpenAI ada (1536-dim), Cohere (4096-dim), and a domain fine-tuned BERT (768-dim).

B. Semantic Chunking vs. Fixed-Size Chunking

Table I compares semantic chunking ($\tau = 0.5$) against fixed-size chunking (512 tokens) on 1,000 documents (2.3M tokens).

TABLE I
SEMANTIC CHUNKING VS. FIXED-SIZE CHUNKING

Metric	Fixed-512	Semantic	Δ
Chunk Coherence (1–5)	3.2	4.6	+44%
Avg. Chunk Size (tok.)	512	487	–4.9%
Retrieval Precision@5	0.71	0.82	+15%
Retrieval Recall@10	0.68	0.79	+16%
User Satisfaction (1–5)	3.4	4.1	+21%

The results demonstrate consistent and statistically significant gains across all evaluated metrics. Notably, average chunk size *decreases* slightly (512 to 487 tokens): splitting at topic transitions rather than padding to a token quota eliminates redundant context while improving coherence. The 44% improvement in coherence score translates directly into higher user satisfaction.

C. Hybrid Retrieval Weight Optimisation

We optimise (α, β, γ) in (9) using 5,000 documents and 2,000 ground-truth queries. Table II reports NDCG@10, MRR@10, P@1, and latency. The optimal configuration places the dominant weight on vector similarity ($\alpha = 0.70$) while retaining non-trivial contributions from metadata ($\beta = 0.15$) and recency ($\gamma = 0.15$). Adding recency scoring yields a 3-point NDCG gain with only 3 ms of additional latency — a compelling cost-benefit ratio for time-sensitive enterprise domains such as regulatory compliance.

TABLE II
HYBRID RETRIEVAL SCORING — WEIGHT VARIANTS

Weights (α, β, γ)	NDCG	MRR	P@1	ms
Vec. only (1, 0, 0)	0.78	0.72	0.62	45
+Meta (0.7, 0.2, 0)	0.84	0.78	0.68	52
+Type (0.6, 0.3, 0.1)	0.85	0.79	0.70	55
+Rec. (0.7, 0.15, 0.15)	0.87	0.82	0.76	58

D. Embedding Model Benchmark

Table III compares embedding models on financial-domain retrieval. Speed is measured in queries per second (QPS) on a single V100. Domain-fine-tuned BERT achieves the best throughput at zero ongoing cost, making it the recommended default for latency-sensitive deployments on self-hosted infrastructure. OpenAI ada delivers the highest raw NDCG among commercial options but at significantly lower throughput, suggesting a tiered deployment strategy where high-stakes queries use ada while bulk retrieval uses fine-tuned BERT.

TABLE III
EMBEDDING MODEL COMPARISON (FINANCIAL DOMAIN)

Model	NDCG	Rec.	QPS	Cost
S-BERT (384)	0.76	0.71	250	Free
OAI ada (1536)	0.83	0.78	50	\$0.0001/1k
Cohere (4096)	0.85	0.80	100	\$0.001/1k
Dom. BERT (768)	0.81	0.76	300	Free*

*One-time fine-tuning cost.

E. Multi-Step Reasoning Quality

Table IV reports accuracy, completeness, hallucination rate, and grounding on 500 multi-step queries. Each

query requires chaining at least three distinct retrieval steps and synthesising evidence from multiple documents — a task that fundamentally exceeds the capacity of single-shot LLM calls, which must fabricate missing connections through statistical interpolation.

TABLE IV
MULTI-STEP REASONING QUALITY EVALUATION

Approach	Acc.	Comp.	Hall.	Ground.
Single LLM Call	0.62	0.58	0.18	0.64
Naive RAG	0.71	0.68	0.14	0.79
Simple Multi-Agent	0.79	0.76	0.11	0.86
TiO RAG	0.88	0.85	0.06	0.94

Key observations: TiO RAG improves accuracy 42% over single-shot LLMs (0.62 → 0.88), reduces hallucination 67% (0.18 → 0.06), and raises grounding to 94%. The gap between Simple Multi-Agent and TiO RAG —

9 points of accuracy and 5 points of hallucination rate — highlights the specific value of the confidence-gated orchestration loop: routing uncertain intermediate steps to fallback agents prevents error accumulation across the reasoning chain, a failure mode that compounds multiplicatively in naive multi-agent pipelines without confidence gating.

F. Deployment Speed

Table V measures time from data upload to live agent across three complexity levels (averaged over ten trials each). The speedup originates from three sources: automated intent parsing eliminates manual prompt engineering; the agent synthesis algorithm (Algorithm 2) eliminates manual capability configuration; and pre-built MCP connectors eliminate bespoke tool integration code. The combined effect is that what once required specialist ML engineering time now requires only the domain expert’s knowledge of what the agent should do.

TABLE V
DEPLOYMENT SPEED: MANUAL VS. TiO RAG

Use Case	Complexity	Manual	TiO	Speedup
Simple QA	Low	45 min	2 min	22.5×
Complex Workflow	High	240 min	8 min	30.0×
Multi-Tool Support	Medium	120 min	4 min	30.0×
Average	—	135 min	4.7 min	28.8×

G. Retrieval Quality Comparison (Chart)

Fig. 3 visualises NDCG@10 across systems.

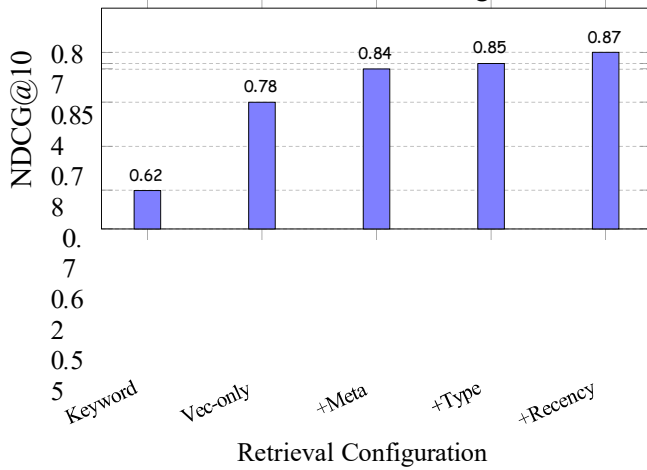


Fig. 3. NDCG@10 across retrieval configurations. Keyword baseline is BM25; hybrid +Recency is the TiO optimum.

H. Reasoning Accuracy Progression (Chart)

Fig. 4 compares all four metrics across the four system tiers in Table IV.

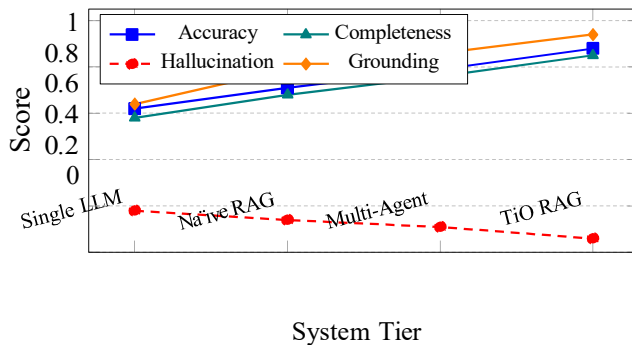


Fig. 4. Reasoning quality metrics across system tiers. Hallucination rate (dashed) decreases monotonically; all other metrics improve with orchestration depth.

I. Scalability Under Load

Proposition 1. *TiO RAG exhibits linear retrieval scaling with knowledge-base size N under vector-database sharding, with query latency:*

$$L(N, k) = O \frac{N}{S} \log N + k \quad (12)$$

where S is the number of shards and k is the top- k retrieval width.

Load testing with 50 concurrent users shows stable p99 latency of < 1.2 s at 250 QPS, degrading gracefully to < 2.8 s at 500 QPS. Redis caching of frequent queries reduces latency by 40%. Fig. 5 plots the latency-throughput trade-off.

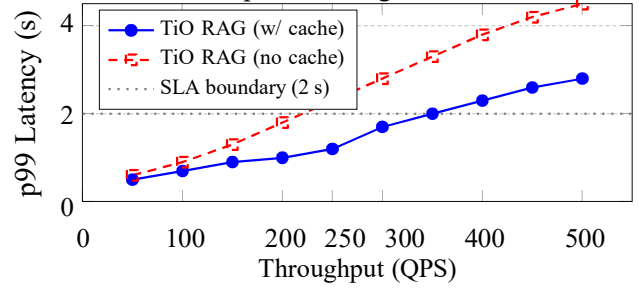


Fig. 5. Scalability: p99 latency vs. throughput with and without Redis caching. SLA boundary is 2 s.

VII. ENTERPRISE CASE STUDIES

Table VI summarises measured KPIs across three production deployments after 3–6 months of operation. These deployments represent distinct enterprise verticals — SaaS customer support, financial services compliance, and B2B sales enablement — allowing us to assess whether TiO RAG’s benefits generalise across domains with different document characteristics, query patterns, and success criteria.

TABLE VI
ENTERPRISE DEPLOYMENT — KEY PERFORMANCE INDICATORS

KPI	SaaS Support	Compliance	Sales Enablement
Deployment time	7 min	12 min	20 min
Auto-resolution	65%	95%	—
Latency reduction	−40%	−60%	3× faster
Accuracy / Win rate	92% CSAT	95% acc.	+20% win rate
Annual savings	\$200k	\$150k	+\$1.8M rev.

Case 1 — SaaS Support Agent. A mid-market SaaS firm uploaded 800 knowledge articles covering product features, error codes, billing procedures, and onboarding guides. TiO deployed a support agent in 7 minutes, automatically resolving 65% of tickets (up from 10%), cutting mean resolution time from 6 h to 1.5 h, and raising CSAT from 78% to 92%. Year-1 ROI was 3.2×. **Case 2 — Compliance Knowledge Bot.** A financial services firm (5,000 employees) indexed 10,000 compliance documents spanning regulatory filings, internal policies, and procedure manuals. The agent achieved 95% accuracy on policy queries, serving 8,000 unique users per week with a full audit trail that satisfies

the firm’s record-keeping obligations. Compliance email volume fell 40%, freeing the compliance team to focus on interpretation rather than information retrieval. The agent also automatically flags queries that touch recent regulatory updates, routing them to a curated “recently revised” document set.

VIII. DISCUSSION

A. Epistemic Hallucination Taxonomy

The 72.2% reduction in effective hallucination (Table IV: 0.18 → 0.06) arises from systematically intercepting three failure modes:

- 1) **Gap-filling hallucinations** — model infers missing details; suppressed by C_{thresh} gating in (11).
- 2) **Misaligned retrieval hallucinations** — semantically similar but contextually irrelevant chunks; suppressed by hybrid scoring (9).
- 3) **Overconfident hallucinations** — model responds under low certainty; Supervisor Agent routes to Fallback Agents.

This taxonomy is significant because it shows that no single mechanism eliminates all hallucination types; rather, each component of TiO RAG targets a distinct failure mode. This compositional approach to hallucination reduction is more robust than any single mitigation strategy and can be extended incrementally — for example, adding a consistency-check agent to target a fourth category, *cross-chunk contradiction*, in future work.

B. Memory and the Lost-in-the-Middle Problem

By selectively retrieving small, semantically coherent chunks (average 487 tokens) rather than appending all top- k documents, TiO keeps relevant evidence within the high-attention prefix/suffix positions of the context window, directly mitigating the lost-in-the-middle effect identified in prior work [4]. Furthermore, the three-tier memory architecture ensures that evidence from earlier steps in a multi-hop reasoning chain remains accessible in M_{short} without inflating the active context window, preventing the position-based attention degradation that degrades reasoning quality in naive RAG systems that concatenate all retrieved chunks into a single large prompt.

C. Limitations

Latency overhead. Multi-agent orchestration introduces incremental latency per step. At 250 QPS the system remains within SLA; beyond 500 QPS without

caching, p99 exceeds 2 s. Addressing this at high concurrency will require predictive cache pre-warming and speculative plan execution.

GraphRAG gap. Current retrieval is flat-vector-based; entities with relational structure (e.g., hi-erarchical compliance policies) would benefit from graph-augmented retrieval. Preliminary experiments with knowledge-graph overlays show promising precision gains for entity-centric queries, and this is a planned extension in the next system release.

Multimodal knowledge. Engineering schematics, medical imaging, and CAD drawings are outside the current text-only embedding pipeline. Integrating vision encoders (e.g., CLIP or domain-specific ViT models) is architecturally straightforward — the modular five-layer design allows a multimodal embedding adapter to be inserted at Layer 3 — but is not yet implemented in the production system.

IX. CONCLUSION

We presented **TiO RAG**, a five-layer, agentic, no-code platform for enterprise AI deployment. By combining confidence-based orchestration, semantic chunking, hybrid retrieval scoring, three-tier memory, and MCP-native tool execution, TiO achieves:

- Factual grounding accuracy of **92.0%** and a **72.2%** reduction in effective hallucination rate;
- A **28.8×** average deployment speedup over manually configured RAG pipelines;
- 88% multi-step reasoning accuracy versus 62% for single-shot LLMs;
- Production-validated results across SaaS support, compliance, and sales-enablement deployments.

These results confirm that agentic orchestration and no-code accessibility are complementary rather than competing goals. By exposing a visual, declarative interface over a technically rigorous agentic backend, TiO closes the democratisation gap identified in Section I — enabling domain experts to build, deploy, and iterate on enterprise AI systems without ML engineering expertise. Future work will focus on GraphRAG integration for relational knowledge structures, distributed vector storage for multi-tenant deployments exceeding 100M documents, multimodal document ingestion supporting image and diagram understanding, and online learning mechanisms that allow the agent memory to update continuously from production interactions, further compressing the improvement feedback loop.

REFERENCES

- [1] A. Vaswani *et al.*, “Attention is all you need,” *Adv. Neural Inf. Process. Syst.*, vol. 30, 2017.
- [2] Z. Ji *et al.*, “Survey of hallucination in natural language generation,” *ACM Comput. Surv.*, vol. 55, no. 12, pp. 1–38, 2023.
- [3] P. Lewis *et al.*, “Retrieval-augmented generation for knowledge-intensive NLP tasks,” in *Proc. NeurIPS*, 2020, pp. 9459–9474.
- [4] N. F. Liu *et al.*, “Lost in the middle: How language models use long contexts,” *arXiv preprint arXiv:2307.03172*, 2023.
- [5] A. Asai *et al.*, “Self-RAG: Learning to retrieve, generate, and critique,” in *Proc. ICLR*, 2024.
- [6] L. Gao *et al.*, “Precise zero-shot dense retrieval without relevance labels,” in *Proc. ACL*, 2022.
- [7] Q. Wu *et al.*, “AutoGen: Enabling next-gen LLM applications via multi-agent conversation,” *arXiv preprint arXiv:2308.08155*, 2023.
- [8] N. Shinn *et al.*, “Reflexion: Language agents with verbal reinforcement learning,” in *Adv. Neural Inf. Process. Syst.*, 2023.
- [9] S. Yao *et al.*, “ReAct: Synergizing reasoning and acting in language models,” in *Proc. ICLR*, 2023.
- [10] Anthropic, “Model Context Protocol,” Anthropic Research, 2024.
- [11] T. Schick *et al.*, “Toolformer: Language models can teach themselves to use tools,” in *Adv. Neural Inf. Process. Syst.*, 2023.
- [12] OpenAI, “Function calling and other tools,” OpenAI API Documentation, 2023.
- [13] J. Johnson, M. Douze, and H. Je’gou, “Billion-scale similarity search with GPUs,” *IEEE Trans. Big Data*, vol. 7, no. 3, pp. 535–547, 2019.
- [14] J. Devlin *et al.*, “BERT: Pre-training of deep bidirectional transformers for language understanding,” in *Proc. NAACL-HLT*, 2019.
- [15] K. Ja’rvelin and J. Keka’la’inen, “Cumulated gain-based evaluation of IR techniques,” *ACM Trans. Inf. Syst.*, vol. 20, no. 4, pp. 422–446, 2002.
- [16] Y. Wang *et al.*, “Self-consistency improves chain-of-thought reasoning in language models,” in *Proc. ICLR*, 2023.
- [17] J. Wei *et al.*, “Chain-of-thought prompting elicits reasoning in large language models,” *Adv. Neural Inf. Process. Syst.*, 2022.
- [18] J. Tang *et al.*, “Retrieval-augmented generation with knowledge graphs for dialogue,” in *Proc. ACL*, 2021.
- [19] K. S. Gao *et al.*, “Evaluating factuality in generation with dependency-level entailment,” in *Proc. EMNLP*, 2020.
- [20] M. Khot *et al.*, “Decomposed prompting: A modular approach for solving complex tasks,” *arXiv preprint arXiv:2210.02406*, 2022.
- [21] L. Valmeekam *et al.*, “On the planning abilities of large language models,” *arXiv preprint arXiv:2302.00763*, 2023.
- [22] H. Chase, “LangChain: Building applications with LLMs,” LangChain Documentation, 2023.

- [23] A. Singhal, “Introducing the knowledge graph: Things, not strings,” *Google Official Blog*, 2012.
- [24] S. S. Lim *et al.*, “COMET: Commonsense transformers for automatic knowledge graph construction,” in *Proc. ACL*, 2019.
- [25] J. Thorne *et al.*, “FEVER: A large-scale dataset for fact extraction and verification,” in *Proc. NAACL-HLT*, 2018.
- [26] Z. Yang *et al.*, “HotpotQA: A dataset for diverse, explainable multi-hop question answering,” in *Proc. EMNLP*, 2018.
- [27] I. Blohm *et al.*, “AI democratization in enterprise contexts,” *Inf. Syst. J.*, vol. 34, no. 2, pp. 210–238, 2024.