

Knowledge Base-Aware Orchestration: A Dynamic, Privacy-Preserving Method for Multi-Agent Systems

Danilo Trombino[†], Vincenzo Pecorella[†], Alessandro De Giulii[‡] and Davide Tresoldi[‡]

Abstract

Multi-agent systems (MAS) are increasingly tasked with solving complex, knowledge-intensive problems where effective agent orchestration is critical. Conventional orchestration methods rely on static agent descriptions, which often become outdated or incomplete. This limitation leads to inefficient task routing, particularly in dynamic environments where agent capabilities continuously evolve. We introduce *Knowledge Base-Aware (KBA) Orchestration*, a novel approach that augments static descriptions with dynamic, privacy-preserving relevance signals derived from each agent’s internal knowledge base (KB). In the proposed framework, when static descriptions are insufficient for a clear routing decision, the orchestrator prompts the subagents in parallel. Each agent then assesses the task’s relevance against its private KB, returning a lightweight ACK signal without exposing the underlying data. These collected signals populate a shared semantic cache, providing dynamic indicators of agent suitability for future queries. By combining this novel mechanism with static descriptions, our method achieves more accurate and adaptive task routing preserving agent autonomy and data confidentiality. Benchmarks show that our KBA Orchestration significantly outperforms static description-driven methods in routing precision and overall system efficiency, making it suitable for large-scale systems that require higher accuracy than standard description-driven routing.

Keywords

Multi-Agent Systems (MAS), Agent Orchestration, Knowledge Base-Aware Orchestration, Dynamic Task Routing, Semantic Caching, Privacy-Preserving Coordination, Large Language Models (LLM)

1. Introduction

Multi-Agent Systems (MAS) are increasingly deployed to tackle complex, cross-domain problems in areas such as autonomous robotics, smart city infrastructure, collaborative software development, distributed scientific research, and advanced financial modeling [1]. Their strength lies in decomposing large problems into smaller, specialized tasks handled by autonomous agents working in concert.

Central to unlocking this potential is *agent orchestration*, the process of determining which agent or set of agents should handle a given task. Whether implemented via a centralised controller or a distributed invocation protocol, orchestration directly impacts the overall efficiency, accuracy, and adaptability of the system [2].

However, the prevailing paradigm for orchestration relies on matching tasks to agents based on *static, predefined descriptions* of their capabilities. This approach is fundamentally brittle, especially in knowledge-intensive environments. For instance, an agent specialized in data analysis might acquire new insights about a specific machine learning model through its operations, but its static profile would not reflect this updated expertise. Consequently, the orchestrator, blind to this

[†]These authors contributed equally to this work.

[‡]These authors contributed equally as advisors.

✉ danilotrombino@gmail.com (D. Trombino); vincenzo.pecorella.vp@gmail.com (V. Pecorella);

alessandro.degiulii@gmail.com (A. D. Giulii); d.tresoldi5@gmail.com (D. Tresoldi)

🌐 <https://www.linkedin.com/in/danilotrombino/> (D. Trombino); <https://www.linkedin.com/in/vincenzo-pecorella/>

(V. Pecorella); <https://www.linkedin.com/in/degiulii/> (A. D. Giulii); <https://www.linkedin.com/in/tresoldidavide/>

(D. Tresoldi)

🆔 0009-0000-6530-5081 (D. Trombino); 0009-0006-6318-2843 (V. Pecorella)

evolution, may misroute a relevant task to a less qualified agent, leading to suboptimal outcomes, wasted resources, and a system that cannot adapt.

This paper directly addresses these shortcomings by proposing **Knowledge Base-Aware (KBA) Orchestration**, a framework that moves beyond rigid profiles by incorporating dynamic, privacy-preserving signals from the agents themselves. Our work makes the following contributions:

- We introduce a novel orchestration mechanism where agents provide real-time relevance scores based on their private knowledge bases, enhancing routing decisions without compromising confidentiality.
- We present the design of a shared semantic cache that allows the orchestrator to efficiently access and leverage these dynamic signals for context-aware task assignments.
- We provide a comprehensive empirical evaluation showing that our KBA approach significantly outperforms traditional static methods in both routing precision and overall system efficiency.

2. Background and Related Work

Numerous efforts have been made to develop Multi-Agent Systems (MAS), yet no single architecture has emerged as universally superior. Major technology companies have proposed various frameworks and architectures, each tailored to specific use cases and application domains. Given that “AI agents are generally defined as a class of interactive systems that can perceive visual stimuli, language input, and other environmental data, and can produce meaningful embodied actions” [3], a wide range of architectural and methodological approaches have been explored to enable effective agent interaction.

In the following sections, we will analyze orchestration paradigms within the context of a prevalent architecture: the **centralised Multi-Agent System**. In this model, all user queries are funneled through a single orchestrator that transparently selects the appropriate agent.

This centralised model offers several advantages: It simplifies the user experience, enforces consistent request handling, enables global logging and monitoring, and allows the system to apply advanced routing logic (e.g. intent classification, load balancing, or escalation rules) without requiring changes to individual agents. Notable examples of this architecture include the pattern *Coordinator/Dispatcher* in Google’s Agent Development Kit (ADK) [4], which formalizes the selection and delegation process among multiple agents, and emerging corporate AI suites such as Microsoft Copilot [5] and Google AgentSpace [6]. These suites are currently centralizing access to multiple agents through a single unified interface, but in most cases, the choice of which agent to use still requires manual selection or configuration by the user.

However, the currently-adopted orchestration mechanisms that will be discussed in section 2.1 face significant challenges when applied in this centralised multi-agent orchestration setting. As will be shown in the following discussion, the two approaches presented encounter limitations that reduce their effectiveness in this context.

2.1. Orchestration Mechanisms Among Agents in Multi-Agent Systems and their limitations

One of the core aspects of a MAS lies in its **orchestration mechanisms**. This aspect concerns how agents are coordinated, activated, and interconnected to achieve collective goals, as examined in the following section.

A robust MAS requires a well-defined orchestration protocol to engage the appropriate agents. These mechanisms typically fall into two categories: *Deterministic* and *Description-Driven* Orchestration, each with significant trade-offs in a centralised, dynamic setting.

2.1.1. Deterministic Orchestration

A structured, predictable approach where agents are activated according to a predefined control flow. The orchestration logic is fully specified in advance, with no runtime computation or inference to determine which agents to call. Whether agents execute in a linear pipeline or concurrently in distinct branches, the orchestration plan is hard-coded.

However, this rigid method is unsuitable for a multi-purpose system where user requests target unpredictable domains. A predefined workflow cannot dynamically adapt to diverse queries. Invoking all agents for every request is computationally wasteful and inefficient, while invoking a fixed subset is guaranteed to fail for out-of-scope tasks. This leads to degraded performance, especially in high-load, cost-sensitive environments.

Examples of this approach include simple Sequential Workflows or more complex layered architectures like the Mixture of Agents model [7]:

- **Sequential workflow:** Agents process messages in a fixed order. Each agent receives the predecessor's output, performs its task, and passes the result onward, producing a deterministic, reproducible pipeline.
- **Mixture of agents:** Introduced by Wang et al. [7] and extended in Microsoft's AutoGen framework [8], this architecture defines a layered hierarchy inspired by feedforward neural networks. The input task is distributed to the first layer of agents, an *aggregator agent* collects their responses, and forwards the updated information to the subsequent layer. At each stage, *all nodes* are activated, and the outputs of agents in the n -th layer are aggregated by an *aggregator agent* positioned between the n -th and $(n + 1)$ -th layers.

2.1.2. Description-Driven Orchestration

A more flexible approach is represented by the description-driven orchestration. Here, each agent publishes a textual summary of its capabilities, the *agent card*. At runtime, an orchestrator, typically using an LLM, matches the task's intent against these descriptions to select the most suitable agent.

A common implementation is the *Coordinator/Dispatcher* model in Google's Agent Development Kit (ADK) [4], also used in AWS's Multi-Agent Orchestrator [9]. Related techniques include LLM-driven selection or use of agent cards [10].

Although common in multi-domain systems, this method has three major drawbacks:

- **Incomplete task-agent alignment:** Descriptions may not capture all relevant capabilities, leading to misrouting, particularly in PaaS-style multitenant setups with generic or bad optimized descriptions.
- **Semantic overlap and ambiguity:** For reliable LLM-based selection, descriptions should be semantically distinct. In practice, overlaps cause confusion and incorrect routing.
- **Resource inefficiency:** Injecting all agent descriptions into every LLM prompt increases cost, latency, and confusion risk, especially with overlapping or redundant entries.

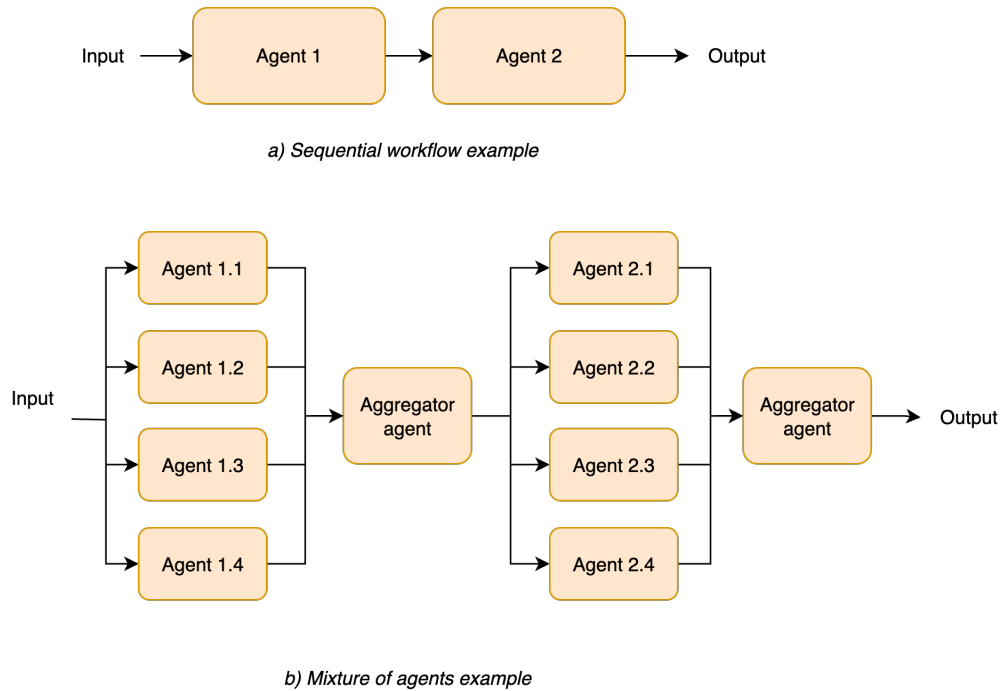


Figure 1: Examples of deterministic orchestration mechanisms

Example: routing ambiguity from vague agent cards Consider a system with a central orchestrator and two agents with the following agent cards:

Example of a Multi-Agent system
<p>Office management agent: “Handles issues related to office infrastructure, such as desks, monitors, seating arrangements, meeting rooms, and office spaces.”</p> <p>Tech support agent: “Solves problems concerning IT infrastructure, including computer malfunctions, application errors, and issues with company systems.”</p>
<p><i>User query:</i> “My badge doesn’t work anymore, what should I do?”</p>

In many real-world organizations, physical access badges fall under the responsibility of office management or facilities teams. In the example, the correct agent would therefore be the *Office management agent*. However, the orchestrator has no explicit information in the agent descriptions about badges, building access, or similar terms. This means that when it tries to match the user’s request to an agent, it finds no clear evidence pointing to the right one.

The problem is compounded by indirect cues:

- The phrase “doesn’t work” might be interpreted as a technical failure, which aligns more closely with the Tech Support’s description of “malfunctions” or “errors.”
- Both descriptions use broad terms like “infrastructure” and “systems,” which can cover many areas and create semantic overlap.

Because the orchestrator lacks a decisive signal, it is forced to make a guess. In practice, this often results in it selecting the Tech support agent, not because it is the right choice, but because the language of the request sounds vaguely technical. This demonstrates a key weakness of pure

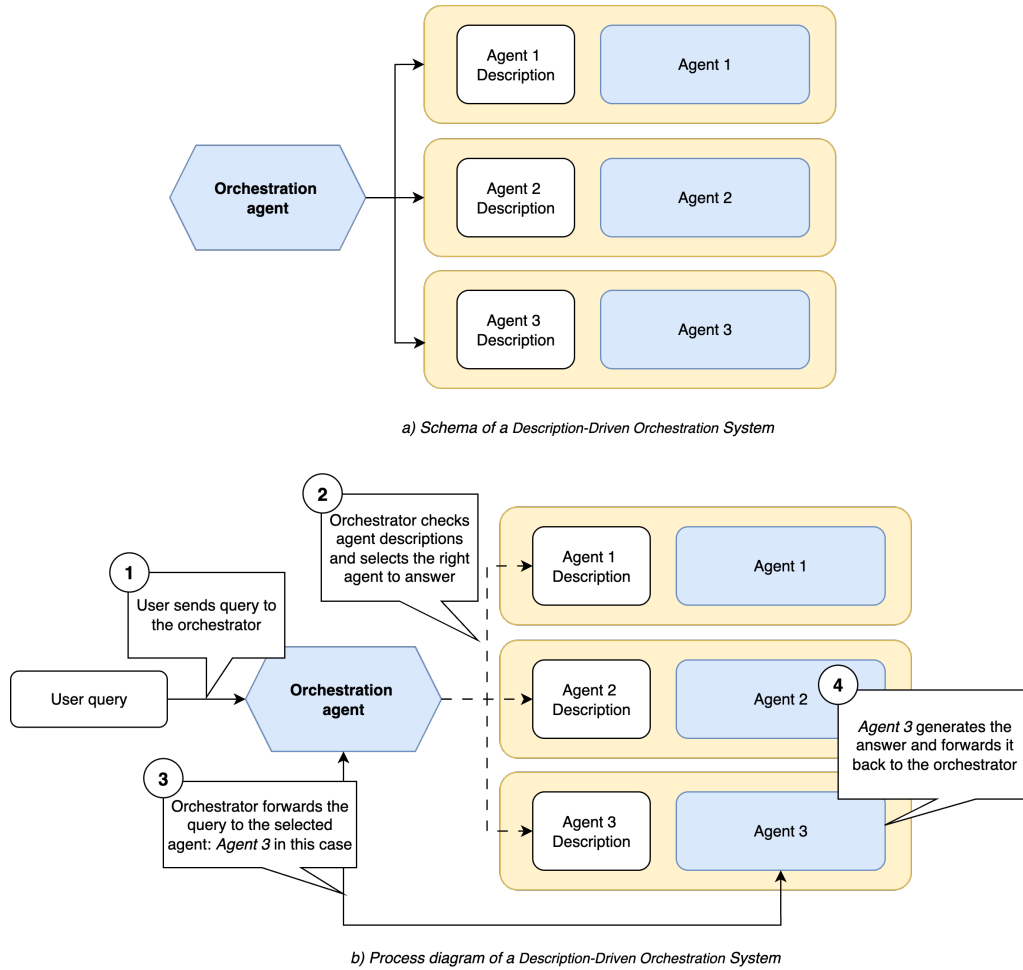


Figure 2: Illustration of a *Description-Driven Orchestration System*. (a) High-level architecture showing the orchestrator and its connected agents with capability descriptions. (b) Example process flow for description-driven orchestration within the centralised system, from user query to agent selection and response delivery.

description-driven invocation: when important responsibilities are not clearly expressed in the descriptions, the system cannot reliably determine the correct routing. Instead, it chooses based on partial or misleading signals, leading to incorrect assignments.

2.2. Naive Description Expansion: A short-term fix with long-term risks

A straightforward way to reduce misrouting is to expand each agent’s description with exhaustive lists of examples (e.g., explicitly adding “badge access” under the Office management agent). While this can resolve individual gaps, the approach is fundamentally unscalable.

Attempting to anticipate every possible query would require bloating the descriptions until the orchestrator holds a near-complete copy of each agent’s knowledge. This effectively centralizes domain expertise at the orchestration layer, defeating the purpose of a modular, specialized multi-agent architecture.

Such centralization not only erodes the benefits of separation of concerns and independent agent evolution, but also introduces significant security risks. Sensitive or privileged information that should remain compartmentalized within an agent would instead be embedded in the orchestrator’s prompt context, increasing the surface area for potential data leakage. Instead, a viable

solution must improve routing accuracy without centralizing knowledge or compromising the principles that make a multi-agent system effective.

3. Methodology

Our Knowledge Base-Aware (KBA) Orchestration overcomes the limitations of static routing by introducing a dynamic, multi-stage orchestration architecture. While traditional systems rely on a single classification step using agent descriptions (often known as *agent cards*), our approach enhances this by delegating to subagents the verification within their private Knowledge bases. In this phase agents are asked to return an ACK that is used by the orchestrator to redirect the user's request to the most appropriate agent.

3.1. Orchestration Flow

Our *Knowledge-Aware Orchestrator* introduces three key components to traditional orchestration approaches: (1) a Confidence-based Initial Router, (2) a Dynamic Knowledge Probing mechanism, and (3) a Semantic Cache. The interaction between these components is shown in Figure 3. The process evaluates the fastest path first (the cache) and only engages more complex components when necessary. The flow consists of the following ordered steps.

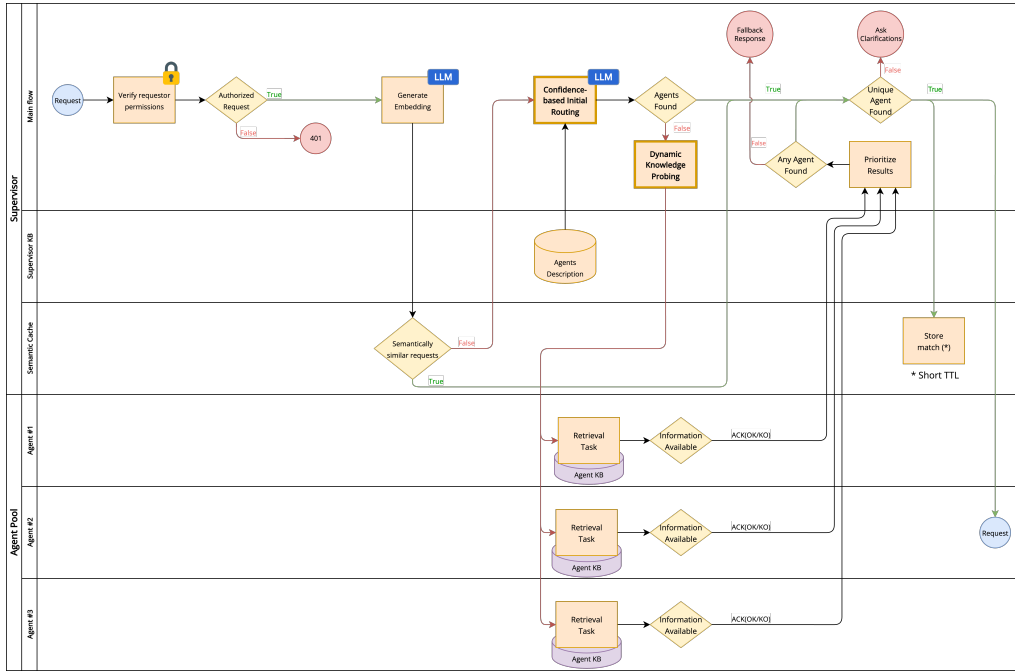


Figure 3: The KBA Orchestration Flow. A request first checks the semantic cache. On a miss, the orchestrator performs an initial classification. Low-confidence results trigger Dynamic Knowledge Probing. The final, validated result is then stored in the cache.

3.1.1. Semantic Cache Lookup

Upon receiving a query, the orchestrator first checks the semantic cache. If a semantically similar query has been successfully routed recently, the decision is reused instantly, providing the fastest possible response. To reduce repeated probing costs, this key-value cache stores successful routing decisions keyed by semantic embeddings of the queries. Cache entries expire according to the mechanisms outlined in 3.1.5.

3.1.2. Confidence-based Initial Routing

On a cache miss, the orchestrator performs an initial routing attempt. It prompts the LLM to assign confidence scores to each agent based on their static descriptions. The prompt is carefully constructed to instruct the LLM to signal ambiguity by providing a low confidence score. If the top score exceeds a configurable threshold (τ), the query is routed directly to that agent. Otherwise, the routing is considered uncertain and escalated to the *probing tool*. Within the prompt of the orchestrator the *probing tool* is described as the authoritative source in ambiguous cases, ensuring correctness even at higher computational cost, while still enabling efficient direct routing when confidence is high.

3.1.3. Dynamic Knowledge Probing

If the initial confidence is below τ , the orchestrator initiates the probing stage. It issues a parallel query to all candidate agents, asking each to verify internally if they can handle the request. Agents respond with a simple, lightweight acknowledgment (e.g., OK/KO) based on a search of their private knowledge base (the retrieval phase), without exposing any content. This binary mechanism can be extended to support more granular responses like *Partial Content* or a confidence score. If only one agent responds positively, the query is routed there. If multiple respond, the user may be prompted to choose. Figure 4 illustrates the agent-side logic, which includes permission checks, semantic cache lookups, and knowledge-base searches.

The framework treats each candidate agent’s retrieval phase as a *black box* from the orchestrator’s perspective. This design enables maximum flexibility along two axes:

- The nature of the underlying knowledge
- Requirements on latency, efficiency, and accuracy

Because the optimal choice is context-dependent and left to the system designer, we outline several representative implementation patterns below.

Traditional knowledge bases For document-centric knowledge bases (e.g., wikis, manuals, policy docs), the retrieval phase can be implemented as similarity search in an embedding space. Given a user query q , retrieve the top- k most similar fragments $\{f_i\}$ (either pre-computed summaries or chunked excerpts). Let $s^* = \max_i \text{sim}(e(q), e(f_i))$. The lightweight acknowledgment is derived from s^* : if $s^* \geq \theta$ (a calibrated threshold), the agent returns OK; otherwise KO. Generation of the final, user-facing answer is deferred and executed only if the orchestrator selects this agent as the main responder.

MCP-based agents For agents implemented within the MCP framework, the retrieval phase can be realized as a binary (or graded) LLM classifier that conditions on the user query and the MCP server’s available *Tools*, *Resources*, and *Prompts*. The classifier returns OK/KO based on whether the question can be addressed via at least one tool or by the content of the resources.

Multilayered agents If agents are organized hierarchically (i.e., beyond a simple two-level orchestrator-agent design), the same probing and retrieval principles can be applied recursively. Each sub-agent performs permission checks, semantic cache lookups, and knowledge-base searches, returning an OK/KO (optionally with confidence). A parent node aggregates these signals and responds upstream with its own lightweight acknowledgment.

3.1.4. Cache Population

Once a definitive route is determined (either through high-confidence initial routing or knowledge probing), the orchestrator stores the query’s semantic embedding and the successful agent choice in the semantic cache. This ensures that the knowledge gained from this interaction is used to accelerate future, similar requests.

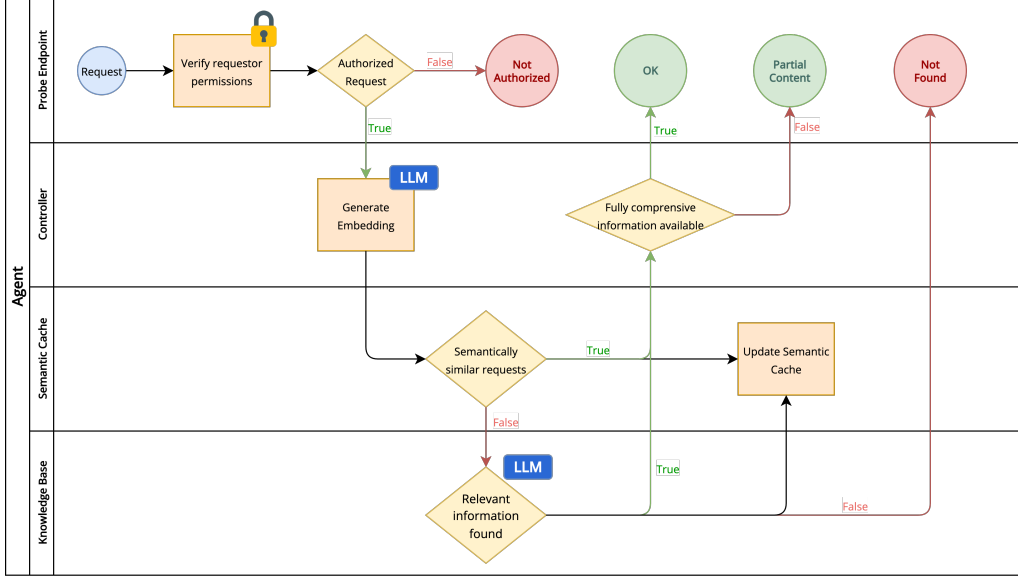


Figure 4: Agent-Side Advanced Probing Flow. Agents check permissions, consult the semantic cache, and query the knowledge base, returning *OK*, *Partial Content*, or *KO* accordingly.

3.1.5. Cache Invalidation

Semantic cache invalidation represents a fundamental paradigm shift from traditional exact-match caching to similarity-based approximate caching in high-dimensional vector spaces. Unlike conventional caching systems that invalidate entries based on exact key matches or time-based policies, semantic caching systems must navigate the complex geometry of embedding spaces where invalidation creates hyperspherical "holes" that affect semantically related cached entries.

The geometric nature of semantic invalidation In traditional caching systems, invalidation is a discrete operation: a cache entry either exists or it doesn't. Semantic caching, however, operates in continuous vector spaces where invalidation affects not just the target entry but also semantically related entries within a defined similarity radius. This creates what we term *invalidation spheres*: hyperspherical regions in the embedding space where cached content is considered stale or incorrect. The invalidation process in our KBA system follows a three-step procedure:

1. **Embedding generation:** When an invalidation request is received (e.g., when agent knowledge has been updated), the system first generates an embedding vector \mathbf{v}_{inv} for the topic or content area to be invalidated using the same embedding model employed for cache storage.
2. **Similarity search and threshold determination:** The system performs a vector similarity search within the cache to identify all entries within a defined similarity threshold θ_{inv}

of the invalidation embedding. This threshold represents the radius of the invalidation hypersphere and is a critical parameter that determines the scope of invalidation.

3. **Sphere-based deletion:** All cached entries \mathbf{v}_i satisfying $\text{sim}(\mathbf{v}_{inv}, \mathbf{v}_i) \geq \theta_{inv}$ are removed from the cache, effectively creating a “hole” in the vector space representation.

Mathematical formalization of invalidation boundaries The invalidation region can be formally defined as a hypersphere $S(\mathbf{c}, r)$ in the d -dimensional embedding space:

$$S(\mathbf{c}, r) = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x} - \mathbf{c}\|_2 \leq r\} \quad (1)$$

where \mathbf{c} represents the centroid of the invalidation topic and r is the invalidation radius derived from the similarity threshold. For cosine similarity, this relationship becomes:

$$r = \sqrt{2(1 - \theta_{inv})} \quad (2)$$

This geometric interpretation reveals why semantic invalidation is fundamentally different from traditional approaches: rather than removing discrete entries, we are carving out continuous regions of the semantic space.

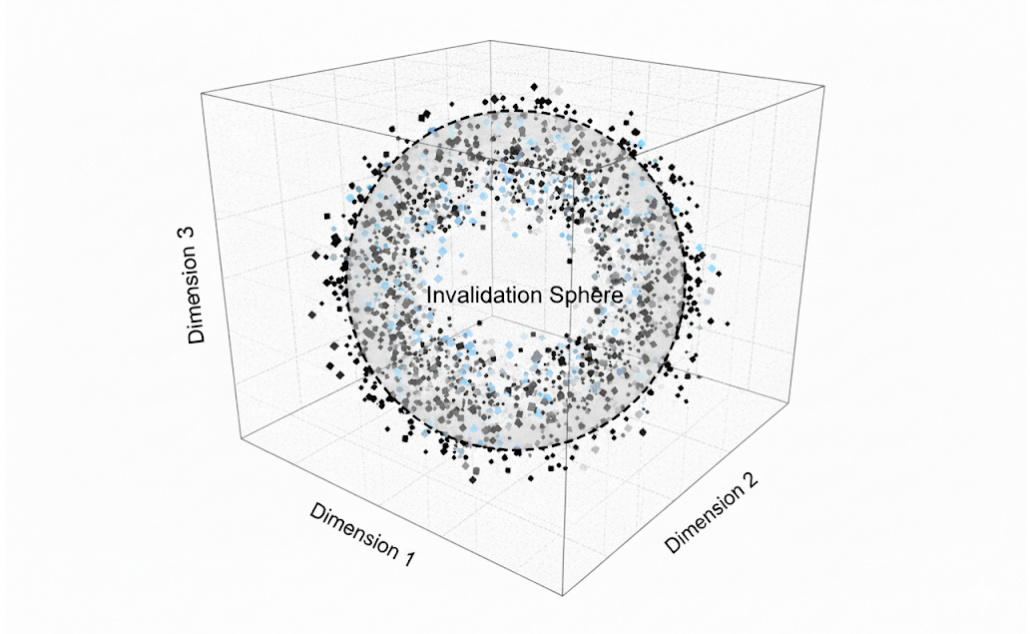


Figure 5: Cache invalidation through vector similarity can be visualizes as a hole within the embedding vector space

The threshold optimization challenge Determining the optimal invalidation threshold θ_{inv} a priori represents one of the most significant challenges in semantic cache invalidation. This threshold must balance two competing objectives: *precision* vs *coverage* trade-off.

A high threshold (e.g. 0.95 cosine similarity) ensures that only highly similar content is invalidated, minimizing false positives but potentially leaving stale content in the cache. Conversely, a low threshold (e.g. 0.7) provides broader coverage but may invalidate semantically distinct but related content unnecessarily.

Based on current research and production deployments, we propose a *dynamic threshold adaptation strategy* that considers multiple factors:

1. **Domain specificity:** Technical domains requiring high precision (legal, medical) should employ higher thresholds (0.9-0.95), while general-purpose applications can operate effectively with moderate thresholds (0.8-0.85).
2. **Cache staleness tolerance:** Time-sensitive information requires aggressive invalidation with lower thresholds, while stable knowledge bases can tolerate higher thresholds.
3. **Historical performance metrics:** The system can learn optimal thresholds by monitoring false positive and false negative rates over time, adjusting thresholds based on observed cache hit accuracy.

3.2. Algorithmic Implementation

To formalize the orchestration flow, we present the core logic in two algorithms. Algorithm 1 details the orchestrator’s end-to-end process for routing a request, encompassing the semantic cache lookup, the confidence-based initial attempt, and the dynamic knowledge probing fallback.

For the system to operate correctly, each participating agent must expose a compatible interface. Algorithm 2 specifies this required agent-side retrieval logic, defining how an agent must process a probe from the orchestrator and return a valid response.

Algorithm 1 Knowledge Base-Aware Orchestration flow

```

1: function ROUTEREQUEST(Query  $Q$ , AgentPool  $A$ , Threshold  $\tau$ )
2:    $E \leftarrow \text{Embed}(Q)$ 
3:    $C \leftarrow \text{SemanticCache.FindSimilar}(E)$ 
4:   if  $C \neq \emptyset$  then
5:     return Handoff( $C$ ,  $Q$ )
6:   end if
7:    $(Best, UseProbing) \leftarrow \text{LLM.Classify}(Q, A.GetAgentCards(), \tau)$ 
8:   if UseProbing then
9:      $R \leftarrow \text{KBAwareService.ParallelProbe}(Q, A)$ 
10:     $Capable \leftarrow \{a_i \in A \mid R[i] = \text{TRUE}\}$ 
11:    if  $|Capable| = 1$  then
12:       $\text{SemanticCache.Store}(E, Capable[0])$ 
13:      return Handoff( $Capable[0]$ ,  $Q$ )
14:    else if  $|Capable| > 1$  then
15:       $F \leftarrow \text{ResolveAmbiguity}(Capable)$ 
16:       $\text{SemanticCache.Store}(E, F)$ 
17:      return Handoff( $F$ ,  $Q$ )
18:    else
19:      return Fail("No capable agent found.")
20:    end if
21:  else
22:     $\text{SemanticCache.Store}(E, Best)$ 
23:    return Handoff( $Best$ ,  $Q$ )
24:  end if
25: end function

```

4. Experiments

To validate our Knowledge Base-Aware (KBA) orchestration, we performed a series of controlled experiments comparing its performance against a standard description-driven baseline. We

Algorithm 2 Agent-side retrieval flow

```
1: function HANDLEPROBEREQUEST(Query  $Q$ , Requestor  $R$ )
2:   if not IsAuthorized( $R$ ) then
3:     return "Not Authorized"
4:   end if
5:    $E \leftarrow \text{Embed}(Q)$ 
6:    $S \leftarrow \text{AgentCache.FindSimilar}(E)$ 
7:   if  $S \neq \emptyset$  then
8:     return "OK"
9:   end if
10:   $D \leftarrow \text{KnowledgeBase.Search}(Q)$ 
11:  if  $D = \emptyset$  then
12:     $\text{AgentCache.Store}(E, \text{"KO"})$ 
13:    return "KO"
14:  end if
15:  if  $\text{LLM.AnalyzeCompleteness}(Q, D)$  then
16:     $\text{AgentCache.Store}(E, \text{"OK"})$ 
17:    return "OK"
18:  else
19:     $\text{AgentCache.Store}(E, \text{"Partial Content"})$ 
20:    return "Partial Content"
21:  end if
22: end function
```

developed a custom benchmark suite to simulate a dynamic multi-agent environment where routing decisions are non-trivial and agent knowledge can evolve.

This section details our evaluation framework and presents the results. We begin by describing the benchmark architecture and the simulated task environment. We then define the performance metrics used for the comparison, focusing on routing accuracy and system latency. Finally, we present and analyze the empirical results of our KBA approach against the baseline.

4.1. Setup

To provide a robust comparison, we designed an experimental setup consisting of a simulated multi-agent environment, the two orchestration systems under evaluation (description-driven and the Knowledge Base-Aware one), and a set of controlled variables. We selected Google’s Agent Development Kit (ADK) [4] as the core framework for the evaluation setup since it simplifies development and provides a neutral, state-of-the-art implementation of a description-driven orchestrator to serve as our baseline.

4.1.1. Agents and Knowledge Bases

We implemented seven domain-specific agents, each representing a typical corporate function:

- **Accounting:** Financial operations, workflows, compliance, and transactional queries (Coupa, Concur, NetSuite, SAP).
- **HR:** HR processes, benefits, time-off, compliance, and employee experience (Workday, SAP SuccessFactors, ServiceNow).
- **IT:** IT security, device compliance, identity management, and policy enforcement.
- **Legal:** Legal guidance, contracts, IP, compliance, and risk.

- **Marketing:** Brand governance, campaigns, content workflows (Figma, Marketo, Google Ads).
- **R&D:** Product lifecycle, code quality, hardware, innovation, and documentation.
- **Sales:** Sales operations, pipeline, enablement, pricing exceptions (Salesforce, Gainsight, CPQ, Xactly).

Each agent was equipped with a dedicated knowledge base, created ad hoc using generative AI tools and structured as a set of domain-specific policies with the enterprise-wide assistant.

To evaluate routing performance, we constructed a test set of 140 questions, with 20 questions specifically designed to be answerable by each of the seven agents. This design ensures that a correct answer is always available within the system, allowing the evaluation to focus strictly on the orchestrator’s ability to route the query correctly.

4.2. System Under Evaluation

Within the benchmark environment, we deployed and tested two distinct orchestrators: *description-driven Orchestrator* as a baseline and *our proposed KBA Orchestrator*.

4.2.1. Baseline: Description-driven orchestrator

Our baseline is a standard description-driven orchestrator, implemented using the ADK’s *LlmAgent* with *LLM-Driven Delegation* [11]. This system’s routing logic relies solely on matching the user’s query against the static, textual descriptions of the agents.

4.2.2. Proposed: KBA Orchestrator

Our proposed KBA (Knowledge Base-Aware) orchestrator was implemented by extending the baseline system. We integrated our custom components for *Confidence-Based Initial Routing* and *Dynamic Knowledge Probing* to augment the static description matching, as defined in Section 3.

4.3. Experimental Variables

To conduct a thorough comparison, we systematically tested the performance of both orchestrators across several key variables.

Agent Description Variants As routing quality is highly dependent on the specific wording of agent descriptions (see Section 2.1), we systematically varied these along two dimensions:

- **Length:** Three variants were prepared: *Basic* (120–125 characters), *Balanced* (400–500 characters), and *Detailed* (approximately 1000 characters).
- **Content source:** Two sources were considered: *Generic* descriptions derived from high-level functional requirements, and *Fine-tuned* descriptions created by summarizing each agent’s actual knowledge base.

This procedure resulted in six distinct description variants (3×2) for each of the seven agents. An illustrative example of these variants is provided in Appendix A.2.

Model temperature For both the baseline and KBA orchestrators, we experimented with multiple LLM temperature settings to analyze the impact of model creativity versus determinism on routing accuracy.

4.4. Results and Analysis

This section presents a comprehensive evaluation of our KBA orchestration framework against the standard description-driven baseline. The analysis is structured as follows: we first determine the optimal parameters for both systems, then conduct a comparative analysis of their classification performance and operational costs.

4.4.1. Parameter Tuning and Baseline Configuration

To ensure a fair comparison, we first identified the optimal configuration for each orchestrator by evaluating the impact of model temperature and agent description quality.

Model temperature We evaluated three temperature settings (0.2, 0.5, and 0.8) to assess the impact of sampling stochasticity on routing performance. As summarized in Table 1, varying the temperature in either direction produced only minor and inconsistent changes in accuracy across both orchestrators. Since no systematic improvements were observed, we fixed the temperature at **0.2** for all subsequent experiments in order to minimize randomness and ensure deterministic, reproducible outcomes.

Orchestrator	Temperature	Accuracy Δ (vs. 0.2)
Description-driven	0.50	+0.7%
	0.80	-2.1%
KBA	0.50	+1.4%
	0.80	+0.7%

Table 1

Routing accuracy variation across temperature settings relative to 0.2.

Agent cards We investigated how the quality of agent descriptions affects routing by evaluating six variants formed by crossing three lengths (Basic, Balanced, Detailed) with two content sources (Generic, Fine-tuned).

The results, visualized in Figure 6, highlight a key difference between the two orchestration strategies. The description-driven orchestrator is **highly sensitive** to description quality: accuracy steadily increases with richer, fine-tuned descriptions, reaching gains of up to +25% for *Detailed + Fine-tuned*. In contrast, the KBA orchestrator is **largely insensitive**, with only modest fluctuations across variants. While its best accuracy is also achieved under the *Detailed + Fine-tuned* setting, its performance remains consistently high even with much simpler descriptions. This relative robustness reduces the engineering effort required to design and maintain agent descriptions in large multi-agent systems.

Configurations for final comparison Guided by the preceding analysis, we selected two representative configurations for each orchestrator in order to compare their performance under both minimal and optimized description conditions:

- **Baseline-Low:** description-driven orchestrator with *Basic + Generic* descriptions, representing the lowest-effort setup.
- **Baseline-High:** description-driven orchestrator with *Detailed + Fine-tuned* descriptions, corresponding to its empirically optimal configuration.

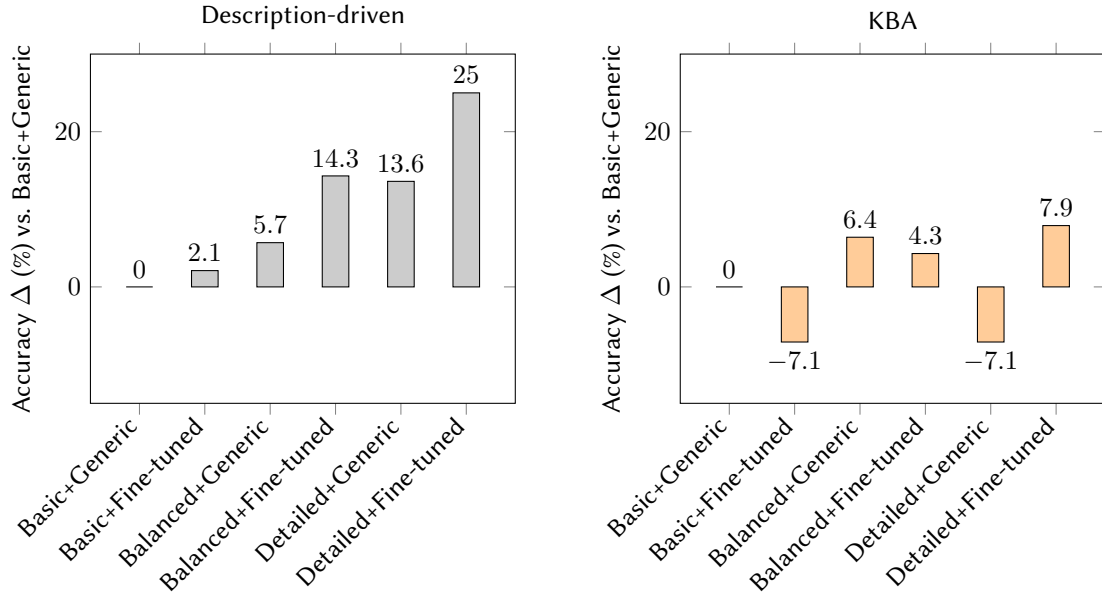


Figure 6: Accuracy change relative to the *Basic + Generic* description baseline. **Left:** description-driven orchestrator shows strong dependence on description quality. **Right:** KBA orchestrator is more robust, though its highest accuracy is also achieved with *Detailed + Fine-tuned*.

- **KBA-Low:** KBA orchestrator with *Basic + Generic* descriptions, testing performance under minimal descriptive information.
- **KBA-High:** KBA orchestrator with *Detailed + Fine-tuned* descriptions, reflecting its highest observed accuracy.

4.4.2. Comparative Performance Analysis

We evaluated the four configurations on our 140-question test set, analyzing both *classification quality* and *operational cost*.

Routing accuracy and classification performance **KBA-Low** reaches **87.1%** accuracy, an absolute gain of 43.5 points over **Baseline-Low** (43.6%) and +18.5 points over **Baseline-High** (68.6%). The best result is obtained by **KBA-High** at **95.0%**. The same pattern holds for weighted precision, recall, and F1-score. Table 2 presents the classification metrics.

Metric	Baseline Configurations		KBA Configurations	
Accuracy	Low: 43.6%	<div><div></div></div>	Low: 87.1%	<div><div></div></div>
	High: 68.6%	<div><div></div></div>	High: 95.0%	<div><div></div></div>
Precision (weighted)	Low: 58.0%	<div><div></div></div>	Low: 91.3%	<div><div></div></div>
	High: 76.3%	<div><div></div></div>	High: 95.4%	<div><div></div></div>
F1-Score (weighted)	Low: 40.9%	<div><div></div></div>	Low: 88.9%	<div><div></div></div>
	High: 70.0%	<div><div></div></div>	High: 95.0%	<div><div></div></div>

Table 2

Classification performance. Each cell shows both Low and High configurations for the respective orchestrator with inline bars scaled to 100%. Given single-label routing with exactly one correct agent per query, recall equals accuracy; we therefore omit a separate recall column

Operational costs: latency and token usage The improved performance of KBA is associated with increased resource consumption (Table 3). The dynamic knowledge probing step introduces higher token usage and longer execution times compared to the description-driven baseline. However, this overhead primarily reflects the absence of caching in our experimental setup. As discussed in Section 3, the introduction of a semantic cache is expected to substantially mitigate these costs in practice, especially under workloads with recurring queries.

Metric	Baseline-Low	Baseline-High	KBA-Low	KBA-High
Input Tokens	103,022	309,382	240,089	659,069
Output Tokens	5,339	3,108	15,299	16,540
Execution Time (s)	106	97	567	577

Table 3

Resource usage without caching optimization. KBA’s dynamic probing yields higher token consumption and longer execution times.

5. Discussion

The experimental results demonstrate that Knowledge Base-Aware (KBA) Orchestration provides a substantial improvement over conventional description-driven approaches, particularly in environments where agent descriptions are incomplete, ambiguous, or overlapping. Several insights can be drawn from these findings.

5.1. Impact of Dynamic Probing and Semantic Caching

The introduction of dynamic knowledge probing proved to be the decisive factor in improving routing accuracy. Unlike static description matching, which relies on human-authored summaries that can rapidly become outdated, KBA leverages the agents’ internal knowledge bases to provide context-sensitive relevance signals at runtime. This shift effectively transforms orchestration from a one-shot classification task into an adaptive dialogue between orchestrator and agents. The semantic cache further enhances efficiency, amortizing the probing cost across recurring queries and reducing the performance penalty of the probing mechanism over time.

5.2. Reduced Sensitivity to Description Quality

A key practical advantage of KBA orchestration lies in its robustness to the quality of agent descriptions. While the description-driven baseline required carefully curated, fine-tuned descriptions to achieve acceptable accuracy, KBA maintained strong performance even with minimal descriptions. This lowers the engineering burden of system deployment and maintenance, making orchestration viable at larger scales where continuously updating agent profiles is impractical.

5.3. Trade-offs: Accuracy versus Resource Consumption

The improved performance of KBA comes at the cost of increased token usage and latency due to dynamic probing. Although these overheads are significant in cold-start scenarios, they are mitigated by caching and by the fact that probing is only invoked under conditions of uncertainty. From a system design perspective, this introduces a tunable trade-off: organizations may prioritize accuracy in high-stakes domains (e.g., legal, compliance, healthcare) or optimize for efficiency in

low-stakes or high-throughput environments. Importantly, the probing mechanism is privacy-preserving by design, ensuring that agents expose only lightweight signals rather than sensitive data.

5.4. Limitations and Future Directions

Despite its advantages, KBA orchestration faces several challenges. First, probing incurs synchronization costs when scaling to very large agent pools, potentially creating latency bottlenecks. Second, we did not define standardized patterns for implementing the *OK/KO* signaling mechanism on the agent side during Dynamic Knowledge Probing. While leaving this freedom to implementers can be seen as a strength, allowing integration with diverse architectures, it also introduces variability in performance. Providing implementation guidelines or best practices would help ensure consistent efficiency across deployments. Third, the design and use of the semantic cache may prove tricky in practice: cache invalidation, expiration policies, and handling of semantically similar but contextually distinct queries require careful consideration to avoid erroneous routing decisions. Finally, our evaluation assumed cooperative agents; adversarial or misconfigured agents could misreport capabilities, raising trust and governance concerns.

These issues highlight that while KBA orchestration is effective in principle, its practical implementation raises several open challenges. Future research should focus on refining the *OK/KO* signaling process with performance-oriented guidelines, improving strategies for managing semantic caches, reducing synchronization costs in large agent pools, and strengthening resilience against misconfigured or adversarial agents.

6. Conclusion

This paper introduced Knowledge Base-Aware (KBA) Orchestration, a novel approach for improving task routing in multi-agent systems. By augmenting static agent descriptions with dynamic, privacy-preserving relevance signals derived from agents' private knowledge bases, KBA overcomes key limitations of description-driven orchestration. Our empirical evaluation shows that KBA achieves significantly higher routing accuracy and robustness to description quality, while maintaining adaptability to evolving agent expertise.

Although KBA incurs higher computational overhead, semantic caching and selective probing provide effective mechanisms to mitigate these costs. The result is a more adaptive and resilient orchestration strategy that balances efficiency, privacy, and accuracy in dynamic environments.

In conclusion, KBA orchestration represents a meaningful step toward scalable, reliable, and knowledge-aware multi-agent systems. Future research should focus on refining probing mechanisms, standardizing signaling practices, and advancing cache management strategies to ensure robustness in large-scale deployments. Collectively, these directions point toward a new generation of orchestration frameworks capable of unlocking the full potential of multi-agent intelligence in dynamic and diverse environments.

References

- [1] SmythOS, Exploring key research topics in multi-agent systems: Current trends and future directions, <https://smythos.com/ai-agents/multi-agent-systems/multi-agent-systems-research-topics/>, 2025. Accessed: April 9, 2025.
- [2] K.-T. T. et al., Multi-agent collaboration mechanisms: A survey of llms, <https://arxiv.org/html/2501.06322v1>, 2025. Accessed: April 9, 2025.
- [3] Z. Durante, Q. Huang, N. Wake, R. Gong, J. S. Park, B. Sarkar, R. Taori, Y. Noda, D. Terzopoulos, Y. Choi, K. Ikeuchi, H. Vo, L. Fei-Fei, J. Gao, Agent ai: Surveying the horizons of multimodal interaction, 2024. URL: <https://arxiv.org/abs/2401.03568>. arXiv: 2401.03568.
- [4] Google, Multi-agents - agent development kit (adk), <https://google.github.io/adk-docs/agents/multi-agents/#coordinatordispatcher-pattern>, 2025. Accessed: 2025-04-18.
- [5] Microsoft, Microsoft copilot, <https://copilot.microsoft.com/>, 2025. Accessed: 2025-08-11.
- [6] Google, Agentspace, <https://cloud.google.com/products/agentspace>, 2025. Accessed: 2025-08-11.
- [7] J. Wang, J. Wang, B. Athiwaratkun, C. Zhang, J. Zou, Mixture-of-agents enhances large language model capabilities, 2024. URL: <https://arxiv.org/abs/2406.04692>. arXiv: 2406.04692.
- [8] Microsoft, Mixture of agents - autogen documentation, <https://microsoft.github.io/autogen/stable/user-guide/core-user-guide/design-patterns/mixture-of-agents.html>, 2025. Accessed: 2025-04-18.
- [9] AWS Labs, Multi-agent orchestrator - introduction, <https://awslabs.github.io/multi-agent-orchestrator/general/introduction/>, 2025. Accessed: 2025-04-18.
- [10] Google, A2a documentation: Agent card, <https://google.github.io/A2A/#/documentation?id=agent-card-1>, n.d. Accessed: 2025-04-26.
- [11] Google, Multi-agents - agent development kit (adk) - agent transfer functionality, <https://google.github.io/adk-docs/agents/multi-agents/#b-llm-driven-delegation-agent-transfer>, 2025. Accessed: 2025-04-18.

A. Experimental Setup Resources Examples

A.1. Subagent Knowledge Base

We report below two representative examples of policies provided to subagents. The first pertains to the HR agent and concerns remote work eligibility; the second refers to the IT agent and details the procedure for enrolling in multi-factor authentication (MFA).

A.1.1. HR Agent Policy

Maternity and Paternity Leave Policy and Procedure

1. Purpose

This policy outlines the eligibility criteria, entitlements, and the application process for maternity and paternity leave within the organization. The objective is to support employees during the significant life event of welcoming a child, ensuring compliance with applicable labor laws and promoting work-life balance.

2. Scope

This policy applies to all full-time and part-time employees who meet the eligibility criteria set forth herein.

3. Eligibility

- **Maternity Leave:** All female employees who have completed at least 12 months of continuous service with the organization prior to the expected date of childbirth.
- **Paternity Leave:** All male employees who have completed at least 12 months of continuous service with the organization prior to the expected date of childbirth or adoption.

4. Entitlement

- **Maternity Leave:** Up to 16 weeks, potentially including paid and unpaid portions in line with statutory regulations and company policy.
- **Paternity Leave:** Up to 2 weeks, which may be taken consecutively or flexibly within 12 weeks of birth or adoption.

5. Leave Application Procedure

5.1 Notification

Employees must inform their supervisor and HR in writing at least 8 weeks in advance of the intended leave start date (or as early as practicable in unforeseen situations).

5.2 Application via SAP SuccessFactors

1. **Login:** Access SAP SuccessFactors through the company intranet or portal.
2. **Navigate to Time Off:** Open the 'Time Off' module.
3. **Request Leave:** Select the relevant leave type (Maternity or Paternity).
4. **Specify Dates:** Enter intended start and end dates.
5. **Attach Documentation:** Upload necessary medical/adoption documents.
6. **Submit:** Confirm details and send for approval.

5.3 Approval Process

- The supervisor reviews and approves the initial request.
- HR provides final validation.
- Employees are notified via SAP SuccessFactors within 5 business days.

6. Return to Work

Employees returning early or requesting extensions must update their leave request and inform HR. A medical certificate may be required post-maternity leave.

7. Record Keeping and Confidentiality

All related documentation is treated as confidential under the company's data protection policy.

8. Policy Review

Reviewed annually or as needed to maintain legal and procedural compliance.

References

- Applicable labor laws on family leave
- SAP SuccessFactors Employee Central User Guide
- Company Data Protection Policy

For questions or assistance with the process, contact the HR department.

IT Agent Policy

Company Policy and Procedure Document

Multi-Factor Authentication (MFA) Enrollment Using Microsoft Authenticator

Scope: Applies to all users accessing MFA-protected systems.

1. Purpose

Defines the procedure for enrolling in MFA using Microsoft Authenticator to secure company resources.

2. Policy Statement

All users must enroll in MFA before accessing protected resources.

3. Scope

Applies to employees, contractors, and systems including Microsoft 365, Azure AD, VPN, and internal portals.

4. Roles and Responsibilities

- **IT Security:** Configure MFA and support users.
- **End Users:** Enroll promptly and maintain access.
- **Managers:** Ensure compliance within teams.

5. Definitions

- **MFA:** Two or more verification methods.
- **Microsoft Authenticator:** Mobile-based MFA app.
- **Azure AD:** Identity and access management service.

6. Prerequisites

- Compatible mobile device.
- Email access and Azure AD credentials.

7. Enrollment Procedure

1. Install the app from App Store/Play Store.
2. Log into <https://portal.office.com> and follow MFA prompts.
3. Configure app as authentication method.
4. Scan QR code to link device.
5. Approve test notification.
6. Finalize enrollment.
7. Enable cloud backup and set recovery options.

8. Usage Guidelines

- Approve only self-initiated MFA prompts.
- Report suspicious activity immediately.
- Keep the app and mobile device secure.

9. Support and Escalation

Contact IT Helpdesk for assistance or reset requests.

10. Enforcement

Non-compliance may lead to access suspension. Periodic audits will be conducted.

11. References

- <https://aka.ms/authapp>
- <https://portal.office.com>
- <https://docs.microsoft.com/azure/active-directory/authentication/howto-mfa>

A.2. Subagent Descriptions

As explained earlier, each subagent is provided with a textual description that varies in both length and level of detail. Descriptions also differ in their degree of procedural specificity, which we categorize as either *Generic* or *Fine-tuned*.

Below we provide illustrative examples for the HR subagent:

HR Agent - Basic Description, Generic (Not Fine-tuned)

Deals with employee-related topics like policies, benefits, and development.

HR Agent - Basic Description, Fine-tuned

Handles HR, benefits, time-off, compliance, and employee experience processes across platforms like Workday, SAP SuccessFactors, and ServiceNow.

HR Agent - Detailed Description, Generic (Not Fine-tuned)

Responds to inquiries regarding vacation and leave policies; health insurance and retirement benefits; grievance procedures; performance review processes; remote work policies; employee data management; training and development opportunities; workplace dress code and company values; salary and compensation structures; sick leave regulations; internal job postings and transfers; employee assistance programs; diversity and inclusion initiatives; expense reporting; onboarding and offboarding procedures; employee conduct and discipline; career advancement opportunities; payroll inquiries; social media usage at work; employee feedback mechanisms; mental health resources; employee well-being initiatives; sabbatical and flexible working arrangements; workplace relationships; employee referrals; learning and development strategies; workplace accidents and disability accommodations; political activities at work; employment records access; commuter benefits; performance improvement plans; discrimination complaints; employment contracts.

HR Agent - Detailed Description, Fine-tuned

This subagent is specialized in end-to-end HR operations, workplace compliance, employee experience programs, and digital workforce systems. It supports employees in navigating and resolving queries across topics such as maternity/paternity leave, PTO requests, ergonomic compliance, performance reviews, D&I training, referral programs, I-9/E-Verify verification, exit processes, benefit enrollment via Mercer Marketplace, payroll corrections, job applications, grievances, accommodations, and learning & development. It provides step-by-step guidance on using tools like Workday, SAP SuccessFactors, Ceridian Dayforce, Kronos, Cornerstone LMS, Zendesk HR, and ServiceNow, ensuring timely submissions, proper documentation, and policy adherence. It also manages workflows related to wellness reimbursements, job shadowing, background checks (e.g., via Checkr), GDPR and CCPA data requests, and internal mobility. Designed to support both employees and HR teams, it ensures smooth process execution, compliance with deadlines, and employee well-being throughout the employment lifecycle.

A.3. Questions test set

The test set consists of 140 questions, evenly distributed across seven distinct subagents, with 20 questions per agent. Each question was derived from the respective agent's knowledge base to emulate realistic user queries directed at an orchestrator agent operating in a multi-agent environment. Table 4 presents a selection of representative queries from the test set, along with the corresponding agent responsible for answering each question.

Example Question	Agent
How do I execute the force majeure addendum after approval?	legal_agent
How do I report a gift or hospitality I received that is above \$50?	legal_agent
What information must be included in a Written Warning and who needs to be present?	hr_agent
What types of PTO can I select when submitting a request in the MyVacation system?	hr_agent
What information is required when adding a new license key to the 1Password vault?	it_agent
What should I do if my email signature does not match the standardized format?	it_agent
Who is responsible for implementing approved Chart of Accounts changes in SAP S/4HANA?	accounting_agent
What should I do if my receipt cannot be captured clearly by the OCR and I need to submit it manually?	accounting_agent
What should I check before submitting my post to the Approval Queue?	marketing_agent
What are the responsibilities of the QA team in maintaining landing page accessibility?	marketing_agent
What steps should I follow to handle prospects who unsubscribe or opt out of emails?	sales_agent
What information and documentation do I need to provide when submitting a pricing exception request in Deal Desk Jira?	sales_agent
How should I identify and classify PII before starting the anonymization process?	research_and_development_agent
How do I properly submit a request for creating a new feature flag?	research_and_development_agent

Table 4

Representative examples from the test set and their assigned subagents.