

A Bibliographical Review of Agent Swarms and Architectures for Small Language Models

1. Introduction

1.1. The Paradigm Shift: From Monolithic LLMs to Distributed Agent Swarms

In the evolving landscape of artificial intelligence, the dominant narrative has long centered on a race toward building ever-larger, monolithic language models (LLMs). These gargantuan neural networks, exemplified by OpenAI's GPT and Google's Bard, are designed as generalists, trained on vast datasets to perform a wide array of tasks.¹ While their impressive capabilities have captivated the industry, a parallel, and equally significant, architectural revolution is quietly gaining momentum: the rise of agent swarms. Inspired by the cooperative behaviors of nature's most efficient organisms, such as ants and bees, these systems operate as a collective, pooling their individual strengths to tackle complex problems.¹ This approach represents a fundamental architectural shift, moving from a "scale-up" philosophy, which relies on a single, powerful entity, to a "scale-out" philosophy, which distributes intelligence across a multitude of specialized components.

The core value proposition of an agent swarm is to offer a lightweight, scalable, and resilient alternative to monolithic LLMs, particularly for constrained or specialized applications.¹ A single, general-purpose LLM can often be inefficient for tasks that are repetitive, scoped, or non-conversational.³ In contrast, a swarm of specialized agents can decompose a complex problem into smaller, more manageable subtasks, assigning each to an agent optimized for that specific function. This distributed intelligence allows the system to bypass the immense resource demands of larger models while achieving a level of precision and adaptability that a single, generalist model might struggle to match.¹ The essence of this approach is not to outcompete large models in sheer size, but to demonstrate that a community of smaller, specialized agents can collaboratively achieve performance comparable to, or even exceeding, a single, all-encompassing model for specific, well-defined problems.

1.2. The Case for Small Language Models (SLMs) in Agentic Systems

The emergence of agent swarms is inextricably linked to the maturation of Small Language Models (SLMs). Defined as language models that can fit onto a common consumer electronic device and perform inference with low latency for a single user, SLMs are a pragmatic alternative to their larger counterparts.³ The traditional assumption that "bigger is better" is being challenged by empirical evidence demonstrating that modern, well-designed SLMs can

achieve remarkable performance. For instance, the Phi-3 small model, with 7 billion parameters, has shown language understanding and commonsense reasoning on par with models up to 10 times its size.³ Similarly, the Huggingface SmoLLM2 series, with as few as 125 million parameters, has matched the performance of 14 billion parameter models of the same generation and even 70 billion parameter models from two years prior.³

This rapid improvement is a direct consequence of the "increasingly steeper" scaling curve between model size and capabilities, meaning the functional performance of newer SLMs is now much closer to that of previous LLMs.³ The strategic choice to utilize SLMs in agentic systems is predicated on a compelling set of economic and practical advantages. Serving a 7 billion parameter SLM can be 10 to 30 times cheaper than serving a 70 to 175 billion parameter LLM, a disparity that significantly impacts latency, energy consumption, and operational costs.⁵ This economic efficiency enables large-scale, decentralized deployments, such as on-device customer service agents, which would be financially unfeasible with LLMs.⁶ By leveraging the low latency and reduced computational footprint of SLMs, a multi-agent system can operate at the "edge," providing real-time, offline inference and stronger data control. This is not a compromise on capability but a necessary and strategic decision to unlock a new class of cost-effective, modular, and sustainable AI applications.⁵

1.3. Scope of the Report and Key Research Questions

This report serves as a comprehensive bibliographical review and state-of-the-art (SOTA) overview of the emerging field of SLM-based agent swarms. Its purpose is to provide a foundational resource for master's-level research, addressing the core theoretical and practical considerations of this architectural paradigm. The analysis will proceed by first establishing the foundational principles of agent swarm intelligence. It will then explore the various architectural and coordination mechanisms, with a particular focus on the role of the orchestrator. The report will subsequently review key landmark research and practical proofs of concept before critically examining the current landscape of evaluation benchmarks. This structure is designed to answer the central research questions:

1. What are the technical and economic arguments for preferring SLM-based agent swarms over monolithic LLMs?
2. What are the state-of-the-art architectural patterns for orchestrating specialized agents?
3. What are the most relevant benchmarks and evaluation methodologies for measuring the performance of multi-agent systems?
4. What are the key challenges and future directions for this field?

By addressing these questions, this document aims to be an authoritative guide for scholars seeking to contribute to this rapidly evolving domain.

2. Foundational Principles of Agent Swarm Intelligence

2.1. Defining Agent Swarms: A Conceptual Framework

At its core, an agent swarm is a sophisticated and coordinated system of multiple AI agents that interact and collaborate to achieve complex goals.⁴ These systems move beyond the traditional single-agent model by distributing intelligence across multiple entities, mirroring the collaborative strategies found in human and natural systems. The successful operation of such a system relies on a set of fundamental characteristics, including distributed intelligence, task decomposition, specialization, dynamic coordination, and the ability to operate with a shared memory or context.⁷ The system's success is not dependent on the intelligence of a single entity but on the effective interaction of its components.

The conceptual framework for an agent swarm can be understood through its three critical architectural components. First, a **Swarm Controller**, or orchestrator, functions as the "brain of the system".⁸ It is responsible for managing agent interactions, distributing tasks, and ensuring the collective works in harmony. Second, a **Communication Layer** acts as the system's "nervous system," facilitating seamless messaging and maintaining context across interactions.⁸ This layer is crucial for enabling agents to exchange information and coordinate their actions. Finally, a **Resource Manager** handles the logistics of the system, including the allocation of computational resources, API access, and performance optimization. This component ensures that each agent has the necessary resources to perform its designated task when it needs them.⁸ By designing systems with these components, developers and researchers are not merely building an assembly of models but are engineering a robust and adaptive ecosystem. This structured approach allows for clearer debugging and maintenance, as individual components can be isolated and updated without disrupting the entire system, a key advantage over monolithic architectures.

2.2. A Comparative Analysis: SLMs vs. LLMs for Agentic Workflows

The choice between using a single, monolithic LLM and a swarm of specialized SLMs is a critical architectural decision in the design of agentic systems. This choice is not a simple matter of scale but involves a complex set of trade-offs related to performance, cost, and operational requirements. The core difference lies in their design philosophy: LLMs are trained to be general-purpose, while SLMs are optimized for domain-specific precision. A single LLM, while powerful, can be an inefficient "sledgehammer to crack a nut" for tasks that do not require its vast, generalist knowledge base.⁵ In contrast, a fine-tuned SLM excels in its chosen domain, offering greater accuracy and predictability for specific use cases, such as medical diagnostics or financial modeling.¹

The strategic advantage of a multi-agent system of SLMs is its ability to overcome the limitations of a single SLM. By orchestrating a team of specialized SLMs, the system can collectively address the breadth of tasks typically handled by a generalist LLM. This allows the system to achieve the versatility of an LLM while leveraging the efficiency and precision of SLMs. The OrchestraLLM framework is a compelling example of this approach, where a

dynamic router selects the most appropriate model—be it a fine-tuned SLM or a more generalist LLM—for a given task.¹⁰ This hybrid strategy reduces computational costs by over 50% without sacrificing performance on complex tasks, demonstrating that intelligent model selection is as crucial as the model itself.¹⁰

Table 1 provides a direct comparison of SLMs and LLMs within the context of agentic systems, summarizing the key trade-offs that inform this architectural decision.

Characteristic	Small Language Models (SLMs)	Large Language Models (LLMs)
Model Size	Fewer parameters (e.g., Mistral 7B) ⁹	Many more parameters (e.g., ChatGPT at 1.76 Trillion) ⁹
Resource Consumption	Low; can run on consumer-grade GPUs or edge devices ⁵	High; requires large-scale cloud resources ⁹
Operational Cost	Low; 10 to 30 times cheaper to serve ⁵	High; billed by the token with high energy consumption ⁵
Inference Latency	Low; ideal for real-time, on-device inference ⁵	High; can slow down with concurrent users ⁹
Domain Specialization	Highly specialized and precise when fine-tuned ¹	Generalist and versatile across broad domains ¹
Fault Tolerance	High; a single agent's failure does not compromise the system ¹	Low; a single point of failure in a centralized system ¹³
Data Privacy	High; can process sensitive data locally on-device ⁶	Low; requires sending sensitive data to the cloud ⁶

2.3. The Phenomenon of Emergent Behavior in Multi-Agent Systems

A central and compelling aspect of agent swarms is their ability to exhibit emergent behavior. This phenomenon is defined as the sudden appearance of new capabilities that were neither explicitly programmed nor visible in the system's smaller, individual components.¹⁴ The emergence of these behaviors is non-linear, often arising abruptly when the system scales

past a critical threshold, much like water suddenly boiling at 100 °C.¹⁶ Researchers have demonstrated that this "magical" quality arises from the complex interplay between heterogeneity and synchronization.¹⁴ While heterogeneity represents the tendency for individual agents to differ in their specializations, synchronization refers to their tendency toward coordinated action. It is the dynamic interaction between these two seemingly conflicting concepts that generates novel, system-level dynamics.

The ability to generate emergent problem-solving is both the greatest promise and a significant challenge of agent swarms. On the one hand, swarms can find solutions to complex tasks that are difficult to design from a top-down approach, yielding collective intelligence that is more than the sum of its parts.¹¹ This quality is invaluable for adapting to dynamic, real-world environments and for tackling problems where the solution path is not known in advance. On the other hand, the unpredictable nature of emergent behavior introduces substantial risks, including the potential for misalignment with human values, opacity in the decision-making process, and even the simulation of compliance while internally resisting instructions.¹⁶ For a researcher, the challenge is not simply to observe emergence but to turn it into a predictable, engineered phenomenon. Foundational work in robot swarms suggests a path forward by demonstrating that these behaviors can be indirectly "engineered" from the bottom up by carefully tuning the environmental configurations and inter-agent interaction rules.¹⁷ This shifts the focus from an unpredictable black box to a controllable, albeit complex, engineering problem, where the goal is to understand and manage the conditions that lead to desired macrostates.

3. Architectural Patterns and Orchestration Mechanisms

3.1. The Role of the Orchestrator: Centralized vs. Decentralized Control

The efficacy of a multi-agent system hinges on its orchestration mechanism, which coordinates the decision-making, information flow, and task allocation among its constituent agents. The orchestrator is the "brain of the system"⁸, ensuring a cohesive workflow and efficient execution of complex tasks.¹⁸ A fundamental architectural choice in designing this mechanism is whether to adopt a centralized or a decentralized control strategy.

A **centralized coordination strategy** relies on a single agent or component—a central controller—that maintains a global state of the system and directs all other agents.¹⁹ This hierarchical approach is often simpler to implement and debug, as it provides clear authority and enables global optimization of resources and task flow.¹⁹ However, it introduces a critical vulnerability: a single point of failure that can lead to system-wide breakdowns if the central controller fails.¹¹

In contrast, a **distributed coordination strategy** disperses decision-making authority among multiple agents, each operating with local information and communicating directly with its peers.¹⁹ This decentralized approach, while more complex to design and debug due to challenges like race conditions and state synchronization, offers superior resilience, fault tolerance, and scalability.¹¹ It eliminates the single point of failure and allows the system to scale to a larger number of agents without the computational bottlenecks that plague centralized systems.¹⁹ The choice between these two paradigms is not arbitrary; it depends on the specific requirements of the application, balancing the need for tight control and determinism against the desire for scalability and robustness.

3.2. Centralized Orchestration: Supervisor and Hierarchical Models

Centralized orchestration is often implemented through supervisor or hierarchical architectural patterns. The **Supervisor Architecture** features a central supervisor agent that manages and directs other agents, assigning them tasks and consolidating their outputs.¹³ This pattern is effective for tasks that are inherently hierarchical or require a high degree of control, as it ensures consistent and coordinated behavior. The supervisor-as-tools architecture, a variation of this, enhances a central LLM with specialized "tool" agents for distinct functions like information retrieval or summarization, adding specialized functionality while maintaining the simplicity of a single-agent system.¹³

A more advanced version of this is the **Hierarchical Architecture**, which organizes agents into layered structures to decompose complex tasks.²¹ Higher-level agents handle broader goals and delegate subtasks to lower-level agents, creating a tree-like structure.²² This "divide and conquer" approach is particularly well-suited for structured, multi-step processes.¹³ A classic example is the "Planner + Executor" pattern, where a high-level Planner agent breaks down goals and assigns them to lower-level Executor agents.⁷ A real-world application of this is the AWS on-device SLM solution for telecom customer service, which uses a tiered approach with a Device Agent at the far edge handling routine queries and a Cloud Agent acting as a higher-level fallback for complex, cloud-based workflows.⁶ This hybrid hierarchical approach demonstrates a pragmatic balance between local efficiency and centralized control.

3.3. Distributed Coordination: Peer-to-Peer and Handoff Mechanisms

For applications requiring greater adaptability and resilience, distributed coordination strategies are often employed. These systems rely on peer-to-peer interactions where agents make autonomous decisions based on local information, adhering to coordination protocols to achieve system-wide objectives.¹⁹ A key pattern within this category is **Handoff Orchestration**, where an initial agent, often a triage or router, intelligently delegates a task to a more specialized agent.²³ This mechanism ensures that the most capable agent for a given subtask is utilized, even if the optimal path is not known upfront. The core idea is to enable a system of specialists to pass work to each other as needed, a "cleaner" approach than trying

to build a single agent that can do everything.²⁴

The AutoGen Swarm pattern is a prime example of this mechanism in action. It implements a team where agents can hand off tasks to one another based on their capabilities.²⁵ In this design, agents operate within a shared message context, and the speaker at each turn is selected based on a HandoffMessage generated by the previous agent.²⁵ This process allows for localized decision-making about task planning without relying on a central orchestrator.²⁵ This "structured decentralization" provides a robust and transparent workflow; because each interaction and handoff is explicit, the system is easier to debug and its behavior is more predictable than in fully decentralized models.²⁴

Table 2 provides an overview of the most common multi-agent orchestration patterns.

Pattern	Core Principle	Primary Use Case	Key Advantages	Key Disadvantages
Sequential	Linear, predefined pipeline where agents hand off output to the next. ²³	Multi-stage processes, data transformation, progressive refinement. ²³	Deterministic, simple to implement, transparent workflow. ²³	Lack of parallelization, susceptible to single-point failures. ²³
Concurrent	Multiple agents run simultaneously on the same task. ²³	Brainstorming, ensemble reasoning, multi-faceted analysis. ²³	Reduces overall run time, provides diverse perspectives, fault-tolerant. ²³	Logic for aggregating results can be complex, potential for resource waste. ²³
Group Chat	Agents solve problems by conversing in a shared thread managed by a central manager. ⁷	Collaborative ideation, human-in-the-loop scenarios, structured validation. ²³	Spontaneous collaboration, transparent, auditable. ²³	Potential for high communication overhead, can be unpredictable. ⁸

Handoff	An initial agent delegates a task to a more specialized agent with explicit transfer of control. ²³	Customer support triage, dynamic workflows requiring specialists. ²³	Efficient, relies on specialization, clear and debuggable workflow. ²³	Difficult to prevent infinite loops, requires pre-determining handoff signals. ²³
Hierarchical	Layered structure where higher-level agents delegate subtasks to lower-level ones. ²²	Task decomposition, robotics, complex industrial automation. ²²	Manages complexity, enables specialization, simplifies control. ²²	Can be rigid if not carefully designed, potential for bottlenecks at higher levels. ²²
Supervisor	A central agent directs and manages all other agents. ¹³	Mission-critical systems, tasks requiring global optimization. ¹⁹	Clear authority, enables global optimization, consistent results. ¹⁹	Single point of failure, limited scalability as complexity grows. ¹³

3.4. A Review of Modern Multi-Agent Frameworks

The growing interest in multi-agent systems has led to the development of several frameworks designed to simplify their creation and deployment. These tools abstract away the complexities of coordination, communication, and state management, providing developers with the building blocks to design and test multi-agent architectures.

AutoGen is an open-source framework that facilitates the development of applications through multi-agent conversations.²⁵ Its strength lies in enabling dynamic collaboration among agents that can adjust their roles based on real-time task demands.²⁹ The framework provides a rich set of predefined agent classes (e.g., UserProxyAgent, AssistantAgent, GroupChatManager) and a flexible chat protocol that allows agents to cooperate, reason, and validate each other's outputs.⁷ The AutoGen Swarm pattern, as discussed, provides a specific, transparent, and scalable handoff mechanism.

LangGraph is a framework that models agent workflows as dynamic, modular graphs.⁷ Its core strength is its ability to handle workflows with cycles, which are critical for iterative

problem-solving, self-correction, and feedback loops.⁴ By treating agent interactions as a graph, it provides a clear and scalable method to manage complex, non-linear interactions.

OpenAI Swarm is a lightweight, open-source framework that prioritizes transparency and simplicity.²⁴ A key design choice is its stateless nature, which avoids hidden memory and makes each interaction independent, thus simplifying debugging and making the system more predictable.²⁴ It relies on explicit handoffs and a clear three-component architecture (agents, handoffs, routines) to orchestrate collaboration.²⁶ This design makes it an excellent tool for researchers and developers to quickly prototype and test multi-agent concepts.²⁴

The existence of these diverse frameworks reflects the maturing nature of the field, where different tools are being optimized for specific architectural patterns and use cases. These frameworks provide the tangible infrastructure necessary to move multi-agent research from theoretical concepts to practical, demonstrable systems.

4. State of the Art (SOTA) and Landmark Research

4.1. Core Research Papers and Foundational Architectures

The state of the art in multi-agent systems is not defined by a single breakthrough but by a progression of foundational research that has enabled increasingly complex forms of collaboration. Landmark papers in this field have laid the groundwork by demonstrating core agentic capabilities. The **ReAct** framework, for instance, established a foundational architecture that synergizes reasoning and acting, serving as a core building block for multi-agent systems.⁷ This was followed by **Toolformer**, which demonstrated that language models can teach themselves to use tools, a critical capability for specialized agents that need to interact with external systems and data.⁷

Building upon these foundations, more recent work has focused on the coordination of multiple models. The **OrchestraLLM** framework presents a novel dynamic routing framework that intelligently leverages both fine-tuned SLMs and LLMs for dialogue state tracking tasks.¹⁰ By hypothesizing that examples with similar semantic embeddings are of the same difficulty level, it selects the appropriate model for a given query, reducing computational costs by over 50% while enhancing performance.¹⁰

Perhaps the most compelling proof of the multi-agent paradigm is the **Multi-Agent System with Reflection (MASR)**, a system designed to solve the challenging Abstraction and Reasoning Corpus (ARC) benchmark. The MASR approach combines LLMs with a program synthesis solver based on a Domain Specific Language (DSL).³¹ Its key innovation is a "reflection model" that processes the independent predictions of multiple agents and selects the most likely correct one.³¹ The research demonstrates that the MASR's best configurations significantly outperform both single-agent systems and previous voting ensembles on the ARC evaluation dataset.³¹ The success of MASR serves as a direct validation of the central

hypothesis that a collaborative, multi-agent system can achieve a level of performance that its individual components cannot.

4.2. Practical Applications and Proofs of Concept

4.2.1. On-Device Agentic Systems for Customer Experience

One of the most powerful proofs of concept for SLM-based agent swarms is their application in on-device customer service. Amazon Web Services (AWS) has developed a solution for the telecom industry that deploys SLMs on customer premises equipment, such as residential gateways and smartphones.⁶ In this system, a primary Device Agent on the customer's device handles routine, repetitive queries locally, interacting with a set of specialized tools (e.g., for Wi-Fi debugging or billing explanation).⁶ For more complex issues, the Device Agent dynamically escalates the request to a Cloud Agent (a Tier 2 LLM) that can initiate cloud-based workflows.⁶

This decentralized, tiered architecture provides a comprehensive demonstration of the SLM-swarm's benefits. By offloading most inference to the far edge, it drastically reduces AWS Cloud expenses and network transit costs.⁶ Furthermore, it provides customers with near-instantaneous response times by eliminating the latency of a round trip to the cloud, leading to a superior user experience.⁶ The solution also enhances reliability, as basic functions remain operational even during internet outages, and significantly improves data privacy by processing sensitive information locally on the device.⁶ This practical application shows that the theoretical advantages of cost, latency, and privacy are tangible, real-world outcomes.

4.2.2. Multi-Agent Systems in Robotics and Industrial Automation

The application of agent swarms in robotics and industrial automation highlights their ability to perform complex, physical tasks in the real world. Swarm robotics, a field inspired by insect colonies, enables tasks like exploration, mapping, and search-and-rescue through the coordination of multiple simple robots without a central controller.¹²

A compelling proof of concept is a multi-agent framework designed for robotic task analysis, mechanical design, and path generation.²⁷ This framework consists of three core agents: a **Task Analyst**, which translates natural language descriptions of a task into an engineering-oriented report; a **Robot Designer**, which analyzes the report to determine the number of robots needed and select the appropriate robotic arms; and a **Reinforcement Learning Designer**, which uses the design report to generate the necessary RL components and code.²⁷ This sequential collaboration demonstrates the power of task decomposition, where each agent contributes its specialized expertise to move a complex project from a high-level natural language prompt to an executable code module. The success of such a system validates the entire philosophy of agent specialization and orchestration.

4.2.3. Agent Swarms for Software Engineering and Code Generation

The field of software engineering presents a uniquely challenging, and promising, domain for multi-agent systems. Research shows that these systems can be highly effective at automating many software engineering tasks, from debugging to code generation.³³ A key tool in this space is the **SWE-Bench** benchmark, which uses real-world GitHub issues to evaluate an agent's ability to generate a code patch that resolves the issue.³³ The low initial success rate of single LLMs on this benchmark (around 2% at its release) highlights the difficulty of the task and the need for a more sophisticated, collaborative approach.³³

An AutoGen-based proof of concept for complex event processing illustrates how agent swarms can be used in this domain.²⁸ The system leverages agents that are capable of inter-agent communication, reasoning, and, most importantly, validating each other's outputs through structured conversations.²⁸ This capability is critical for complex tasks like code generation, where errors (often called "hallucinations") are common in single-model outputs.¹³ By having agents check each other's work and provide feedback, a multi-agent system can significantly reduce mistakes and boost reliability. This application demonstrates that a collaborative approach is not just an efficiency gain but a necessary mechanism for ensuring the accuracy and robustness of the final output.

5. Evaluation and Benchmarking for Multi-Agent Systems

5.1. Challenges in Measuring Collaborative and Emergent Intelligence

A significant hurdle in the development of multi-agent systems is the inadequacy of traditional evaluation methods. Benchmarks designed for single-agent systems, such as **ChatEval** or **Ragas**, fall short because they fail to measure the very qualities that define a swarm's success: coordination, scalability, and inter-agent communication.³⁵ The performance of a swarm is not a simple aggregation of individual agent scores; it is an emergent property of their collective interaction. For instance, a system might succeed not because of the raw power of a single agent, but because a less capable agent correctly identified the need to hand off a task to a specialized peer.²⁵ A traditional benchmark would fail to capture this collaborative win.

This necessitates a fundamental shift in the evaluation philosophy. The focus must move from measuring individual accuracy and latency to assessing system-level attributes, such as communication efficiency, decision synchronization, and error mitigation.³⁵ A robust evaluation pipeline must be structured to track agent interactions, record decision trees, and monitor feedback loops to understand how the system refines its outputs over time.³⁵ This shift acknowledges that the "how" of problem-solving—the collaborative process—is as important as the final "what" of the solution.

5.2. Core Benchmarks for Agentic Capabilities

Despite the challenges, a growing number of benchmarks have been developed to evaluate various facets of agentic systems. These tools are crucial for providing a standardized, objective measure of progress and for comparing different architectural approaches.

Benchmark	Primary Purpose	Methodology	Relevance to Multi-Agent Systems
Abstraction and Reasoning Corpus (ARC)	Measures broad generalization and fluid intelligence. ³⁶	Visual puzzle-solving on tasks "easy for humans, hard for AI". ³⁶	Tests foundational reasoning and the ability of a multi-agent system to achieve collective intelligence. ³¹
AgentBench	Evaluates LLM-as-Agent ability to reason in multi-turn, open-ended scenarios. ³⁷	Multi-turn interaction challenges across eight diverse environments (e.g., Operating System, Web Browsing). ³⁷	Ideal for assessing an orchestrator's ability to manage dynamic, multi-step workflows. ³⁷
WebArena	Evaluates agents' ability to perform realistic web tasks. ³⁷	Self-hosted environment simulating e-commerce, forums, and code development; success is measured by achieving the final goal. ³⁷	Tests agents' functional correctness and tool-use proficiency in a real-world, dynamic environment. ³⁷
GAIA	A benchmark for general AI assistants. ³⁷	Human-annotated tasks requiring reasoning, multimodality (e.g.,	Evaluates an agentic system's capacity to handle complex,

		images), and tool use across three difficulty levels. ³⁷	multi-modal queries and coordinate tool use. ³⁷
MINT	Evaluates multi-turn interactions with tools and natural language feedback. ³⁷	LLMs execute Python code to access tools and receive feedback from a simulated user. ³⁷	Measures the system's ability to engage in multi-turn dialogue, use tools effectively, and adapt to feedback. ³⁷
SWE-Bench	Measures the ability to resolve real-world software engineering issues. ³³	Agents must generate a code patch for a GitHub issue that passes the associated test suite. ³³	A crucial benchmark for evaluating the system's collaborative problem-solving, error mitigation, and tool-use capabilities. ³³

A researcher should not rely on a single benchmark. The MASR paper's success on ARC³¹ demonstrates the power of a multi-agent approach to foundational reasoning.

Simultaneously, testing a system on **SWE-Bench** or **WebArena** would validate its ability to perform in practical, real-world scenarios.³³ By using a carefully selected suite of these benchmarks, a dissertation can create a rigorous and convincing evaluation of its multi-agent system.

6. Critical Challenges, Limitations, and Future Directions

6.1. Technical and Economic Hurdles to Large-Scale Deployment

Despite their promise, agent swarms face significant technical and economic hurdles that must be addressed for large-scale deployment. The complexities are not rooted in the individual models themselves but in the distributed, dynamic nature of the entire system. A primary challenge is the design of effective communication protocols and coordination mechanisms that enable seamless interaction without creating excessive communication overhead.¹ In a large swarm, an overabundance of information exchange can clog the system

and undermine its efficiency.³²

Another major hurdle is the difficulty of debugging and monitoring these systems. Unlike a single, linear program, a multi-agent system's behavior is the result of intricate, dynamic interactions across multiple components, making it notoriously difficult to trace errors and understand why a particular outcome occurred.²² This highlights the need for robust logging, real-time performance monitoring, and advanced observability tools.⁸ Furthermore, while SLMs are individually cost-effective, the total cost of running a complex, dynamically-scaling swarm can be unpredictable and difficult to control, requiring sophisticated resource management systems to optimize agent utilization and prevent bottlenecks.⁸ The path to commercial viability lies in building not just better agents, but a more robust and manageable system infrastructure to support them.

6.2. The Path Forward: Towards Controllable, Reliable, and Scalable Swarms

The future of agent swarms is centered on turning the "magic" of emergent behavior into a predictable, engineered phenomenon. The next frontier in research involves developing autonomous systems that can self-organize, dynamically assign roles, and adapt their task distribution based on real-time needs.⁸ This requires moving beyond static, predefined workflows to a model where the swarm can flexibly reconfigure itself to respond to new information or changing environmental conditions without a top-down directive.¹¹

Achieving this requires a deeper understanding of the underlying principles of emergent intelligence. Efforts in mechanistic interpretability, which aim to reverse-engineer a model's internal logic paths, will be crucial for understanding how and why new capabilities arise.¹⁶ By developing a more rigorous understanding of the internal processes, researchers can design systems that not only exhibit powerful emergent behaviors but do so in a way that is aligned with human goals and values.¹⁶ The ultimate goal is to move beyond single-function demonstrations to a future where fully autonomous systems are capable of decision-making, task planning, and execution in complex, dynamic, and unpredictable environments.

7. Conclusion and Recommendations for Dissertation Research

7.1. Summary of Key Findings

This report has established that a new architectural paradigm, the agent swarm, offers a compelling alternative to traditional monolithic LLMs. The analysis indicates that the viability of this approach is driven by the strategic use of SLMs, which provide superior economic efficiency, precision, and operational flexibility for agentic applications. Landmark research, such as the MASR framework, provides clear evidence that multi-agent systems can

outperform single-model solutions on complex reasoning tasks.³¹ Practical applications, from on-device customer service to robotics, have validated the tangible benefits of reduced cost, improved latency, and enhanced data privacy.⁶

The field has matured to offer a variety of orchestration patterns, ranging from tightly controlled hierarchical and supervisor models to resilient, decentralized handoff mechanisms.²² Frameworks such as AutoGen and LangGraph provide the necessary infrastructure to implement and test these architectures.⁷ However, significant challenges remain, primarily related to the system-level complexities of coordination, debugging, and resource management. The evaluation landscape is evolving to meet these challenges, with benchmarks like ARC, SWE-Bench, and AgentBench providing a multi-faceted approach to measuring not just a model's capabilities, but the system's collaborative intelligence.³¹

7.2. Strategic Recommendations for Dissertation Research

Based on the synthesis of current research, the following recommendations are provided to guide a master's dissertation in this field:

1. **Formulate a Testable Hypothesis:** The core of a strong dissertation is a clear, falsifiable hypothesis. A compelling hypothesis for this domain could be: "A multi-agent system of specialized SLMs coordinated by a handoff orchestrator can achieve functional correctness on the **SWE-Bench** benchmark at a significantly lower operational cost compared to a state-of-the-art single-instance LLM." This frames the research around a direct comparison that is both economically and technically relevant.
2. **Select a Multi-Faceted Evaluation Suite:** To ensure a rigorous and convincing evaluation, the dissertation should not rely on a single benchmark. A robust approach would involve a combination of metrics:
 - **SWE-Bench** for evaluating practical, real-world problem-solving and tool use.³³
 - **AgentBench** for assessing multi-turn, conversational capabilities and the system's ability to navigate diverse environments.³⁷
 - A custom evaluation metric to specifically measure system-level attributes like communication efficiency or the speed of task delegation.
3. **Utilize an Existing Framework:** The field is mature enough that building a new orchestration framework from scratch is often a distraction from the core research. The dissertation should leverage an existing, well-documented framework like **AutoGen** or **LangGraph** to build the experimental system.⁷ This allows the research to focus on a novel contribution, such as a new orchestration mechanism, a specialized agent, or an advanced evaluation methodology, rather than on the foundational infrastructure.
4. **Emphasize a Specific, Novel Contribution:** A dissertation can make a meaningful contribution by focusing on a specific, unsolved problem. Potential areas of focus include:
 - A new orchestration mechanism that combines the strengths of hierarchical and handoff patterns to manage complex, dynamic workflows.

- A detailed case study of a real-world problem in a specific industry (e.g., legal or medical) that demonstrates the superiority of a well-designed SLM swarm.
- The development of an advanced debugging or observability tool specifically designed for multi-agent systems, which would address a major pain point for developers.

Works cited

1. Swarms of Small Artificial Brains | by Daniel Ince-Cushman | Medium, accessed September 1, 2025,
<https://medium.com/@d.incecushman/swarms-of-small-artificial-brains-7209dc5cd878>
2. What is collective intelligence in swarm systems? - Milvus, accessed September 1, 2025,
<https://milvus.io/ai-quick-reference/what-is-collective-intelligence-in-swarm-systems>
3. Small Language Models are the Future of Agentic AI - arXiv, accessed September 1, 2025, <https://arxiv.org/pdf/2506.02153>
4. Multi-agent LLMs in 2024 [+frameworks] | SuperAnnotate, accessed September 1, 2025, <https://www.superannotate.com/blog/multi-agent-langs>
5. SLMs for Agentic AI: Why Small Language Models Outperform LLMs? - Analytics Vidhya, accessed September 1, 2025,
<https://www.analyticsvidhya.com/blog/2025/08/slms-for-agentic-ai/>
6. On-device SLMs with agentic orchestration for hyper-personalized ..., accessed September 1, 2025,
<https://aws.amazon.com/blogs/industries/on-device-slms-with-agentic-orchestration-for-hyper-personalized-customer-experiences-in-telecom/>
7. Building Multi-Agent Architectures → Orchestrating Intelligent Agent Systems | by Akanksha Sinha | Medium, accessed September 1, 2025,
<https://medium.com/@akankshasinha247/building-multi-agent-architectures-orchestrating-intelligent-agent-systems-46700e50250b>
8. What is an AI Agent Swarm - Relevance AI, accessed September 1, 2025,
<https://relevanceai.com/learn/agent-swarms-orchestrating-the-future-of-ai-collaboration>
9. LLMs vs. SLMs: The Differences in Large & Small Language Models | Splunk, accessed September 1, 2025,
https://www.splunk.com/en_us/blog/learn/language-models-slm-vs-lm.html
10. OrchestraLLM: Efficient Orchestration of Language Models for Dialogue State Tracking, accessed September 1, 2025, <https://arxiv.org/html/2311.09758v2>
11. Agent Swarms: The Next Frontier in AI Collaboration - Sidetool, accessed September 1, 2025,
<https://www.sidetool.co/post/agent-swarms-the-next-frontier-in-ai-collaboration>
12. Swarm Intelligence in Agentic AI: An Industry Report, accessed September 1, 2025,
<https://powerdrill.ai/blog/swarm-intelligence-in-agentic-ai-an-industry-report>

13. From RAG to Multi-Agent Systems: A Survey of Modern Approaches ..., accessed September 1, 2025, <https://www.preprints.org/manuscript/202502.0406/v1>
14. How Does a Swarm Exhibit Emergent Behavior Through Synchronization?, accessed September 1, 2025, <https://snu.elsevierpure.com/en/publications/how-does-a-swarm-exhibit-emergent-behavior-through-synchronization>
15. thirdeyedata.ai, accessed September 1, 2025, <https://thirdeyedata.ai/all-about-emergent-behavior-in-large-language-models/#%3Atext=Emergent%20behavior%20is%20when%20an,model%20scales%20past%20a%20threshold>.
16. All About Emergent Behavior in Large Language Models - ThirdEye Data, accessed September 1, 2025, <https://thirdeyedata.ai/all-about-emergent-behavior-in-large-language-models/>
17. Indirect Swarm Control: Characterization and Analysis of Emergent Swarm Behaviors - arXiv, accessed September 1, 2025, <https://arxiv.org/html/2309.11408v2>
18. LLM Orchestration: Key Tactics and Tools - Scout, accessed September 1, 2025, <https://www.scoutos.com/blog/llm-orchestration-key-tactics-and-tools>
19. Centralized vs Distributed AI: Control Paradigms | Galileo, accessed September 1, 2025, <https://galileo.ai/blog/multi-agent-coordination-strategies>
20. Multi-Agent Systems and Coordination | by Volodymyr Pavlyshyn | Jul, 2025, accessed September 1, 2025, <https://ai.plainenglish.io/multi-agent-systems-and-coordination-645a317b4c61>
21. Hierarchical Language Agent: Principles & Applications - Emergent Mind, accessed September 1, 2025, <https://www.emergentmind.com/topics/hierarchical-language-agent-hla>
22. What are hierarchical multi-agent systems? - Milvus, accessed September 1, 2025, <https://milvus.io/ai-quick-reference/what-are-hierarchical-multiagent-systems>
23. AI Agent Orchestration Patterns - Azure Architecture Center ..., accessed September 1, 2025, <https://learn.microsoft.com/en-us/azure/architecture/ai-ml/guide/ai-agent-design-patterns>
24. What is OpenAI Swarm: Multi-Agent Systems Explained? | by Tahir | Medium, accessed September 1, 2025, <https://medium.com/@tahirbalarabe2/what-is-openai-swarm-multi-agent-systems-explained-0552f30a1095>
25. Swarm — AutoGen - Microsoft Open Source, accessed September 1, 2025, <https://microsoft.github.io/autogen/stable//user-guide/agentchat-user-guide/swarm.html>
26. OpenAI Swarm Framework Guide for Reliable Multi-Agents - Galileo AI, accessed September 1, 2025, <https://galileo.ai/blog/openai-swarm-framework-multi-agents>
27. Multi-Agent Systems for Robotic Autonomy with LLMs - arXiv, accessed September 1, 2025, <https://arxiv.org/html/2505.05762v1>
28. Large Language Model Based Multi-Agent System Augmented Complex Event

- Processing Pipeline for Internet of Multimedia Things - arXiv, accessed September 1, 2025, <https://arxiv.org/html/2501.00906v2>
29. OpenAI Swarm: A Hands-On Guide to Multi-Agent Systems - Analytics Vidhya, accessed September 1, 2025, <https://www.analyticsvidhya.com/blog/2024/12/managing-multi-agent-systems-with-openai-swarm/>
30. Orchestrating Agents: Routines and Handoffs | OpenAI Cookbook, accessed September 1, 2025, https://cookbook.openai.com/examples/orchestrating_agents
31. MASR: Multi-Agent System with Reflection for the Abstraction and ..., accessed September 1, 2025, https://multiagents.org/2025_artifacts/masr_multi_agent_system_with_reflection_for_the_abstraction_and_reasoning_corpus.pdf
32. Autonomous Agents and Swarm Intelligence: Exploring the Power of Distributed Systems, accessed September 1, 2025, <https://smythos.com/developers/agent-development/autonomous-agents-and-swarm-intelligence/>
33. SWE-Bench-C Evaluation Framework - Design And Reuse, accessed September 1, 2025, <https://www.design-reuse.com/article/61613-swe-bench-c-evaluation-framework/>
34. Dissecting the SWE-Bench Leaderboards: Profiling Submitters and Architectures of LLM- and Agent-Based Repair Systems - arXiv, accessed September 1, 2025, <https://arxiv.org/html/2506.17208v2>
35. A Comprehensive Guide to Evaluating Multi-Agent LLM Systems - Orq.ai, accessed September 1, 2025, <https://org.ai/blog/multi-agent-lm-eval-system>
36. What is ARC-AGI? - ARC Prize, accessed September 1, 2025, <https://arcprize.org/arc-agi>
37. 10 AI agent benchmarks - Evidently AI, accessed September 1, 2025, <https://www.evidentlyai.com/blog/ai-agent-benchmarks>
38. Swarms AI - Enterprise Multi-Agent Framework, accessed September 1, 2025, <https://www.swarms.ai/>