# Towards a Science of Collective AI: LLM-based Multi-Agent Systems Need a Transition from Blind Trial-and-Error to Rigorous Science

**Jingru Fan**[1,*], **Dewen Liu**[2,*], **Yufan Dang**[3], **Huatao Li**[1], **Yuheng Wang**[1], **Wei Liu**[4],
**Feiyu Duan**[2], **Xuanwen Ding**[2], **Shu Yao**[1], **Lin Wu**[1], **Ruijie Shi**[3], **Wai-Shing Leung**[3],
**Yuan Cheng**[1], **Zhongyu Wei**[2], **Cheng Yang**[5], **Chen Qian**[1,†], **Zhiyuan Liu**[3,†], **Maosong Sun**[3]

[1]Shanghai Jiao Tong University, [2]Fudan University [3]Tsinghua University,
[4]King's College London, [5]Peng Cheng Laboratory

{qianc, fanjingru510}@sjtu.edu.cn, dwliu23@m.fudan.edu.cn, liuzy@tsinghua.edu.cn

## Abstract

Recent advancements in Large Language Models (LLMs) have greatly extended the capabilities of Multi-Agent Systems (MAS), demonstrating significant effectiveness across a wide range of complex and open-ended domains. However, despite this rapid progress, the field still relies heavily on empirical trial-and-error. It lacks a unified and principled scientific framework necessary for systematic optimization and improvement. This bottleneck stems from the ambiguity of attribution: first, the absence of a structured taxonomy of factors leaves researchers restricted to unguided adjustments; second, the lack of a unified metric fails to distinguish genuine collaboration gain from mere resource accumulation. In this paper, we advocate for a transition to design science through an integrated framework. We advocate to establish the collaboration gain metric ($\Gamma$) as the scientific standard to isolate intrinsic gains from increased budgets. Leveraging $\Gamma$, we propose a factor attribution paradigm to systematically identify collaboration-driving factors. To support this, we construct a systematic MAS factor library, structuring the design space into control-level presets and information-level dynamics. Ultimately, this framework facilitates the transition from blind experimentation to rigorous science, paving the way towards a true science of Collective AI.

## 1. Introduction

With the enhanced reasoning capabilities of Large Language Models (LLMs), powerful autonomous agents capable of goal-oriented planning and task execution have emerged (Minaee et al., 2024; Zhou et al., 2025a; Wang et al., 2024b). By integrating internal memory and external tool-use, these agents can independently execute complex tasks, such as code generation and mathematical problem-solving (Xi et al., 2025; Schick et al., 2023). While these agents are inherently limited by their individual capability boundaries, making it difficult to maintain universal expertise across open-ended domains, inspired by the success of collective intelligence in biological swarms and human civilizations (detailed in Appendix A), Multi-Agent Systems (MAS) have emerged as a promising direction to break through these constraints (Chen et al., 2023; Li et al., 2024; Guo et al., 2024). By organizing diverse individuals to work collaboratively, MAS can transcend some limitations of individual agents, enabling agents to solve multi-turn, multi-function, and open-ended problems through collective intelligence (Tran et al., 2025; Maldonado et al., 2024). Consequently, the MAS paradigm has achieved significant success in complex and open fields such as scientific research (Ghareeb et al., 2025), software engineering (Qian et al., 2024), infrastructure management (Seff et al., 2023), healthcare services (Liu et al., 2025a), financial analysis (Xiao et al., 2025), and social sciences (Zhang et al., 2024).

However, despite the abundance of successful multi agent systems, the field's advancement remains reliant on engineering-driven blind trial-and-error, lacking scientific guidance for systematic improvement (Pan et al., 2025; Guo et al., 2024). Researchers thus default to optimizing performance via empirical adjustments to key factors—such as organizational structures (Li et al., 2024). This reliance on empirical adjustments essentially stems from the inability to conduct rigorous factor attribution. While such empirical optimization occasionally yields performance improvements, the true attribution of these gains remains obscured.

From the perspective of MAS construction, the ambiguity of attribution manifests as two critical challenges. Firstly, prior to initiating MAS optimization, researchers must select specific factors to verify performance gains; however, this selection is hindered by the absence of a structured taxonomy for the vast factor space (Maldonado et al., 2024). Without a standardized candidate factor set to anchor the process, researchers are restricted to unguided factor adjustments, unable to systematically examine potential collaboration-

---

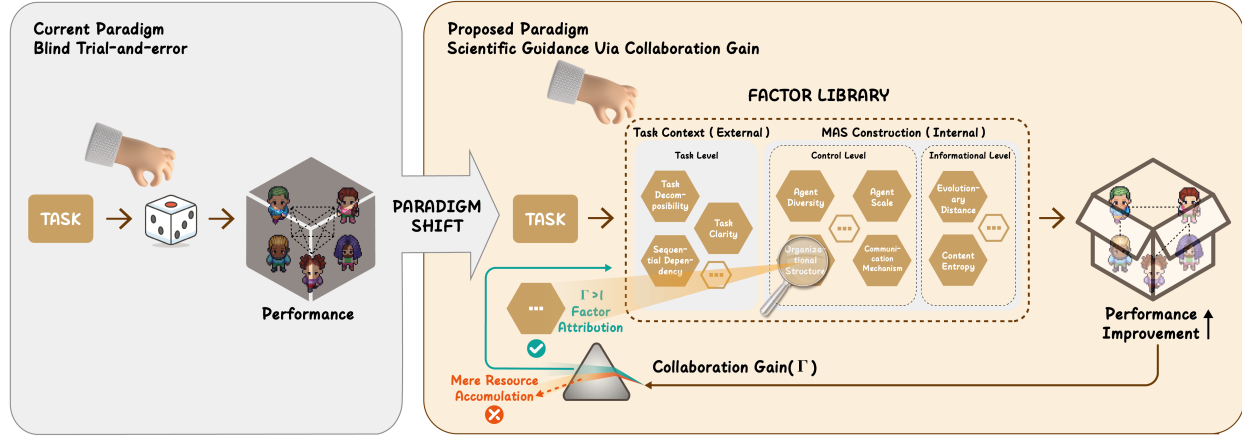*Equal contribution.
†Corresponding authors.

*Figure 1.* The Paradigm Shift: From Blind Trial-and-Error to a Science Guidance. Left (Current): An opaque black box where performance gains are stochastic and unattributable. Right (Proposed): A white-box paradigm. Researchers select factors from the library to construct the MAS; observed performance is then passed through the $\Gamma$ (the prism), which filters out mere resource accumulation to isolate genuine collaboration gain. This analytic step inherently executes factor attribution.

driving factors (Kim et al., 2025). Even more critically, in the dimension of value verification, the lack of a unified metric makes it impossible to distinguish whether a specific factor modification drives genuine collaboration gain or merely reflects mere resource accumulation(capacity growth without collaboration) (Zhang et al., 2025c). Current task-specific metrics inherently conflate these two qualitatively distinct effects: the intrinsic gain from collaboration (i.e., capability growth resulting from agent collaboration) and the performance improvement resulting merely from increased computational budget (Qian et al., 2025). Lacking both a structured framework for factors and a precise diagnostic signal for metrics, researchers cannot systematically identify the design factors that drive genuine collaboration (Renney et al., 2026). This limitation wastes substantial computational resources on ineffective designs and, more critically, precludes the establishment of scientific principles, preventing MAS from evolving into a rigorous science (Kapoor et al., 2025).

To address these challenges and guide the scientific construction of MAS, we propose three core recommendations that form an integrated paradigm, as illustrated in Figure 1. Recognizing that a scientifically valid evaluation standard is the prerequisite for distinguishing effective designs, we prioritize the proposal of the metric and attribution methodology, which then governs the systematic exploration of the factor space (Zhu et al., 2026):

1. **Build Collaboration Gain Metric** ($\Gamma$)**:** First and foremost, we call for establishing the collaboration gain metric ($\Gamma$) as a quantified feedback signal for system optimization. We define $\Gamma$ as the performance ratio of a MAS to a Single-Agent System (SAS) under equiv-

alent computational resource constraints. By isolating the genuine collaboration gain from mere resource accumulation (Du et al., 2025), this metric provides the essential quantitative feedback needed to identify valid optimization directions. It thus serves as a principled diagnostic signal that anchors the entire scientific framework.

2. **Factor Attribution Paradigm:** Leveraging this metric, we propose designing a factor attribution paradigm to identify genuine collaboration-driving factors. The core mechanism follows a two-step, sequential logic: we first examine whether modifying a factor leads to measurable performance improvement. Only when such improvement is observed do we apply the $\Gamma$ metric for scientific evaluation to verify the source of the gain—a candidate factor is confirmed as a genuine collaboration contributor if and only if $\Gamma > 1$, which indicates the improvement stems from genuine collaboration gain rather than mere resource accumulation (Malfa et al., 2025). This sequential validation can help transform factor attribution from trial-and-error into a more rigorous, evidence-based scientific process.

3. **Construct Factor Library:** To support this attribution, we advocate for a systematic MAS factor library—a systematically organized set of variables influencing system behavior. While the attribution paradigm provides the "method," this library provides the "scope" by systematizing the previously fragmented design space, providing a structured candidate factor set for the aforementioned attribution task. We decouple the factor space into the task context (external) and the MAS construction (internal) (Maldonado et al., 2024). The

latter is further organized into two essential levels: the control level, representing static architectural presets (e.g., organizational structure) (Dang et al., 2025), and the information level, characterizing dynamic execution mechanisms (e.g., content entropy) (Nikitin et al., 2024; Xia et al., 2025; Kuhn et al.). This structured framework provides a rigorous and reusable roadmap for systematic optimization.

By establishing the correlations between specific factors and collaboration gain via this integrated paradigm, we aim to replace blind trial-and-error with systematic scientific guidance (Panait & Luke, 2005), a shift that can help researchers move from less systematic approaches toward more proactive modulation of MAS factors to enhance collaborative efficacy. Since the transition toward a true science of Collective AI requires more than engineering-driven blind trial-and-error, it demands a fundamental rethinking of our research trajectory and a collective consensus on scientifically guided MAS construction (Zhu et al., 2026). We therefore advocate for the research community to embrace this integrated paradigm, providing a roadmap to transform MAS into a rigorous science and ensuring that the next generation of collective intelligence is built upon transparent scientific logic.

## 2. Problem Statement

While MAS have demonstrated outstanding performance across numerous tasks, the field remains heavily reliant on engineering-driven blind trial-and-error, unable to transition toward rigorous scientific construction (Li et al., 2024; Wu et al., 2025). From the perspective of MAS construction, this methodology hits a bottleneck stemming fundamentally from the ambiguity of attribution. This ambiguity obscures the connection between specific factors and performance gains, severely impeding the field's evolution from random trials to systematic scientific guidance. Specifically, this manifests as two critical challenges: the unguided search among a vast array of factors, and the inability of metrics to discern genuine collaboration gain from mere resource accumulation (Ma et al., 2025).

### 2.1. Absence of Structured Taxonomy for Factor Selection

The first barrier lies in the difficulty of factor selection, which currently resembles a blind search due to the absence of structured scientific guidance. Since MAS optimization is governed by a vast array of interacting systemic factors—such as agent scale (Qian et al., 2025; Wang et al., 2025c), organizational structure (Dang et al., 2025; Zhou et al., 2025b), and communication mechanism (Zou et al., 2025; Yang et al., 2025)—navigating this high-dimensional factor space has effectively become a blind search. This un-

guided search becomes computationally intractable as task complexity scales. Researchers face the risk of diminishing marginal returns, as coordination overheads—including communication and synchronization—often grow faster than collaborative benefits (Yang et al., 2026). Recent findings even indicate that when agent scale exceeds a specific threshold, system performance tends to stagnate or even degrade due to information overload (Qian et al., 2025).

Currently, research is trapped in a double bind: we attempt to solve increasingly difficult problems by raising system complexity, yet we lack the methodology to navigate that very complexity (Malfa et al., 2025). This necessitates a factor-based framework to transform this blind search into proactive, principled regulation.

### 2.2. Metrics Confounding Genuine Collaboration Gain and Resource Scaling

Even more critically, distinct from the factor selection challenge, valid attribution is prevented by the second gap: the inability of existing metrics to distinguish genuine collaboration from resource accumulation. The persistent reliance on blind trial-and-error is primarily driven by the fact that current evaluation methodologies tend to depend on metrics confined to task-specific performance (Sun et al., 2025; Schipper et al., 2025; Jing et al., 2025; White et al., 2025; Chan et al., 2024). While effective for benchmarking end-results, these metrics inherently conflate two qualitatively distinct effects: the intrinsic gain from collaboration (i.e., capability growth resulting from agent collaboration) and the performance improvement resulting merely from increased computational budget (Qian et al., 2025).

Consequently, it remains difficult to precisely determine whether a MAS architecture has truly achieved an expansion of the system's capability boundaries via collective intelligence, or has merely exploited a greater volume of resources—such as increased token consumption or agent scale—compared to a single-agent baseline (Du et al., 2025). This ambiguity obscures the causal link between systemic factors and performance outcomes, hindering the field's transition toward more predictive scientific guidance (Kuhn & Hacking, 1970).

## 3. Measuring Genuine Collaboration Gain: A Principled Metric to Guide Factor Attribution

To resolve this ambiguity of attribution and move beyond blind trial-and-error, we advocate for the construction of a diagnostic signal that quantifies a system's genuine collaboration gain as a rigorous feedback mechanism for optimization (Bo et al., 2024). While raw resource accumulation—such as expanding token budgets or agent popula-

tions—can indeed enhance performance, the diminishing marginal utility inherent in single-agent scaling suggests a fundamental ceiling for individual. Our proposed metric, therefore, serves as more than a performance benchmark; it functions as a rigorous filter for factor attribution. By decoupling genuine collaboration gain from resource-driven improvements, this metric facilitates a more systematic identification of factors that contribute to collective intelligence, thereby supporting a more informed approach to MAS design (Wang et al., 2025a).

## 3.1. Formal Definition and Theoretical Foundations

We contend that a key objective of MAS research is to capture the nonlinear leap facilitated by collective collaboration (Qian et al., 2025). This pursuit is fundamentally rooted in emergence theory, which posits that the macro-level behavior of a complex system should transcend the simple summation of its micro-level constituents (Fromm, 2005). In the context of MAS, this macro-level behavior manifests as the collaboration gain derived from inter-agent collaboration, while the micro-level summation corresponds to the cumulative performance of agents acting in isolation. While we do not deny the fundamental role of individual scaling and resource expansion in driving intelligence, we argue that the true scientific frontier lies in identifying whether—and to what extent—group-level organization can yield additional dividends beyond these scaling effects. To move toward a more rigorous science of Collective AI, we need to better disentangle these emergent effects from the contributions of mere resource expansion (Hoel et al., 2016).

To solve this, we propose the collaboration gain ($\Gamma$), a metric that quantifies the efficiency of a collective relative to a single-agent baseline under identical resource consumption (Tang et al., 2025). By benchmarking performance against this resource-equivalent reference, $\Gamma$ can serve as a probe for measuring genuine collaboration gain, helping to isolate the dividends of interaction from the impact of computational scale (Chen et al., 2024).

$$\Gamma = \frac{\Phi_M}{\Phi_S}{}^1 \qquad (1)$$

- $\Phi_M$ represents the collective performance of the MAS, specifically capturing the efficacy facilitated by its group-level organizational and collaborative mechanisms (Malfa et al., 2025).

- $\Phi_S$ denotes the non-collaborative baseline performance, representing the best achievable results of a

---

[1]We use the ratio here for its intuitive simplicity, but a more generalized form could be any function $f(\Phi_M, \Phi_S)$ that quantifies the relationship, potentially learned through regression or other methods to capture more complex, non-linear interactions.
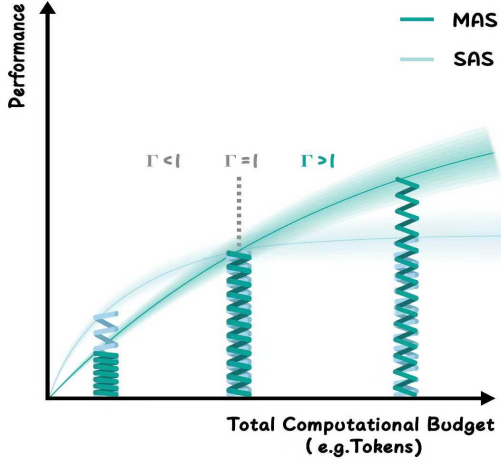
single agent when allocated a total resource budget equivalent to that of the MAS.

Under the null hypothesis of zero interaction effect, the system performance would simply equal the single-agent baseline ($\Phi_m = \Phi_s$), yielding $\Gamma = 1$ (Zhang et al., 2025c). Consequently, $\Gamma = 1$ serves as the theoretical floor for collaboration; any value above this threshold ($\Gamma > 1$) suggests that the system has achieved synergy, where "the whole is greater than the sum of its parts"(see Figure 2).

To operationalize this metric, we emphasize that the operational instantiation of both $\Phi.$ and the definition of resource equivalence are task-dependent (Zhu et al., 2025). First, the evaluation function $\Phi$ is adaptively selected to reflect the intrinsic goals of the task (e.g., accuracy, coverage, or efficiency) (Mohammadi et al., 2025). Second, the best achievable performance for $\Phi_S$ implies a saturated baseline tailored to the task's logical structure, ensuring the single agent utilizes the most effective non-collaborative strategies within the budget(detailed in Appendix B.1). Finally, the metrics for resource consumption—whether measured in token throughput, sampling steps, or tool-invocation counts—are calibrated to the specific bottlenecks of each task (detailed in Appendix B.2).

## 3.2. $\Gamma$-Driven Analysis of Factor Attribution

Building upon this foundation, $\Gamma$ functions as more than a performance metric; it serves as a diagnostic probe for isolating the actual collaborative yield of collective intelligence. However, simply calculating a number is insufficient for meaningful system design. Given the chaotic landscape of MAS design, attempting quantitative modeling (e.g., regression) on unverified factors introduces significant noise. Therefore, we establish a binary attribution framework to filter candidate factors into two distinct regimes based on their empirical validity.

### 3.2.1. CATEGORIZATION OF COLLABORATION GAIN

We use $\Gamma = 1$ as the definitive boundary of synergetic emergence—the equivalence point where architectural gain meets resource cost (Tang et al., 2025). This allows us to partition the factor space into two distinct factor categories:

**Class I: The Positive Factors** ($\Gamma > 1$)   Factors yielding $\Gamma > 1$ with statistical significance represent the successful realization of collaboration gain. This classification confirms that the specific design (e.g., a specific communication topology) enhances system performance beyond the single-agent limit under identical resource constraints (Kim et al., 2025). These factors are considered true collaboration drivers, as they demonstrate a proven ability to push the system beyond individual capability ceilings. They are the valid candidates for further quantitative refinement and tuning.

*Figure 2.* Conceptual Illustration of Collaboration Gain ($\Gamma$). The curves represent the performance of MAS and SAS under equivalent computational budgets to ensure comparability. When SAS performance equals or exceeds MAS, $\Gamma \lesssim 1$, indicating mere resource accumulation; conversely, $\Gamma > 1$ signifies the genuine collaboration beyond the single-agent.

**Class II: The Negative Factors** ($\Gamma \lesssim 1$)  Factors that fail to significantly outperform the baseline are classified into the negative set, indicating a state of attribution failure where the MAS collaboration lacks a clear advantage. This regime, characterized by $\Gamma \lesssim 1$, represents the absence of a synergetic dividend. In this state, any potential benefits of collaboration are either nullified by structural redundancy—where the protocol merely replicates individual reasoning—or actively suppressed by negative interference, such as coordination overhead and context fragmentation (Pan et al., 2025; Zhang et al., 2025c). Consequently, the negative factors serve as a critical pruning signal, compelling researchers to recognize that the current architecture fails to justify its added complexity and preventing the wasteful optimization of ineffective designs.

### 3.2.2. THE FACTOR ATTRIBUTION PROCESS

To operationalize this framework for factor attribution, we advocate for a sequence that begins by strictly preconditioning the experiment on a fixed computational budget (e.g., total token consumption) (Kim et al., 2025). This budget is utilized to construct a saturated single-agent baseline $\Phi_S$, ensuring it reflects the maximal non-collaborative capability within the allocated constraints, and to measure the MAS performance $\Phi_M$ under this identical resource allocation. We then utilize $\Gamma$ to measure the genuine collaboration gain, effectively decoupling structural advantages from mere resource expansion (Li et al., 2024). To verify the validity of the factor attribution against stochastic noise, we implement

a stability filtering stage; a factor is promoted to the positive set if it maintains a sustained $\Gamma > 1$ advantage. Finally, this process culminates in the final attribution, identifying whether the specific architectural design functions as a verified causal driver of the observed collaboration gain. This approach ensures that any MAS optimization is grounded in verified collaboration gain, preventing the wasteful fine-tuning of ineffective structures(detailed in Appendix D).

## 4. Systematizing the MAS Design Space: A Structured Factor Library

To facilitate effective factor attribution, we propose constructing a MAS factor library(see Figure 3). This library systematizes the previously fragmented design space of MAS factors, providing a structured candidate factor set for subsequent factor attribution.

This framework enables the factor attribution process to transcend reliance on blind trial-and-error, instead facilitating scientific exploration through a rigorous taxonomy comprising two primary domains: the task context(external) and the MAS construction(internal). Within the internal construction, we further stratify factors into two essential levels: the control level, representing static architectural presets (e.g., organizational structure), and the information level, characterizing dynamic execution mechanisms (e.g., content entropy) (Li et al., 2024; Maldonado et al., 2024). This classification provides a rigorous roadmap, ensuring that the construction of MAS is guided by a comprehensive and logical design space. For specific optimization strategies and literature-supported implementations corresponding to each factor, please refer to Appendix C.

### 4.1. Task Context (External) Factors

The attributes of a task define the external boundary conditions under which a MAS operates (Decker, 1995).They shape the logical structure of the problem space and constrain the bounds of the collaboration gain $\Gamma$. In particular, task attributes such as task decomposability, sequential dependency, and task clarity provide representative illustrations of how task context can affect the compatibility between problem formulation and collaborative agent architectures (Davis & Smith, 1983; Malone & Crowston, 1994; Decker & Lesser, 1993).

Empirical evidence suggests that mismatches between task structure and MAS design can neutralize or even degrade collaborative benefits (Kim et al., 2025). For example, highly sequential tasks with strong dependency chains often limit parallelism, reducing the effectiveness of distributed coordination, while tasks with low inherent clarity often necessitate explicit task decomposition and planning structures (e.g., in software engineering) (Qian et al., 2024).
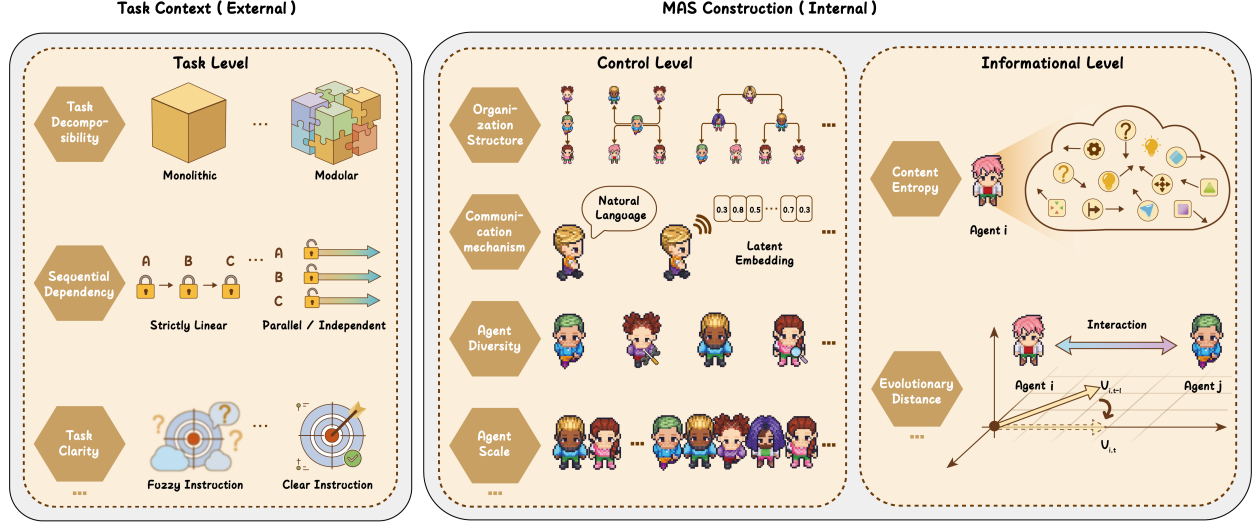
*Figure 3.* The MAS Factor Library Taxonomy. Factors are organized into task context(external) and MAS construction( internal), with the internal dimension spanning static control level presets and dynamic information level to guide rigorous system design.

These properties appear to shape the feedback behavior of $\Gamma$, suggesting that the effectiveness of collective cooperation is inherently influenced by the task's execution logic.

## 4.2. MAS Construction (Internal) Factors

We systematize the MAS factor space into two collaborative levels: control and information level. The control level encompasses static architectural presets, serving as the foundation that limits collaboration potential (Li et al., 2024; Guo et al., 2024; Maldonado et al., 2024). Complementarily, the information level characterizes dynamic execution mechanisms, revealing how this potential is activated and transformed into actual gains (Van Dyke Parunak & Brueckner, 2001; Nowak, 2006).

### 4.2.1. CONTROL LEVEL

Control level encompasses static architectural factors, serving as the foundation that presets collaboration potential limits (Li et al., 2024). Its core logic formally maps unstructured requirements into a structured system design, defining operational boundaries via specific structural factors. This level integrates four dimensions: organizational structure, communication mechanisms, agent diversity, and agent scale (Guo et al., 2024; Maldonado et al., 2024). Improper design creates structural bottlenecks, resulting in poor collaboration gain.

**Organizational structure** constitutes the architectural configuration and connection topology of the MAS, establishing the fundamental pathways available for agent interaction (Li et al., 2024). We identify this as a critical factor because the

potential of even high-performing agents can be constrained by a restrictive topology, which may impose bottlenecks that limit information propagation and collaborative efficiency. (Qian et al., 2025).

Research in this domain is undergoing a transition from preset, static hierarchies to adaptive, dynamic architectures. For example, agentic supernets allow systems to replace fixed workflows with flexible topologies, achieving higher accuracy at only a fraction of traditional inference costs (Zhang et al., 2025c). A well-constructed topology is designed to optimize the density of connections and minimize structural redundancy, thereby providing the architectural foundation necessary for the collaboration gain to manifest as tangible performance improvements.

**Communication mechanism** governs internal collaboration and environmental feedback, encompassing modalities ranging from explicit natural language to implicit latent space representations (Li et al., 2023). Particularly within text-based interactions, unconstrained protocols often lead to a "communication explosion", where redundant dialogues cause exponential resource scaling and diminishing marginal returns (Kim et al., 2025; Wang et al., 2024b; Tran et al., 2025). We argue that excessive and unstructured redundancy can suppress $\Gamma$ by drowning critical task signals in communicative noise, thereby preventing the system from exceeding single-agent performance under the same budget.

To address this issue, recent research focuses on communication-efficient design, optimizing interaction structures and protocols to reduce redundancy while preserving coordination capacity (Chen et al., 2025). For example,

one-shot pruning on spatial-temporal message graphs filters redundant transmissions, reducing token usage by 28.1%–72.8% while maintaining system robustness (Zhang et al., 2025b). Through such structured coordination control, the mechanism lowers coordination costs and seeks to create the necessary conditions for converting structural potential into a verified collaboration gain ($\Gamma > 1$).

**Agent diversity** characterizes the degree of functional heterogeneity within a MAS, representing a tunable factor of differentiation. Formally, let $\mathcal{I} = \{i_1, \ldots, i_N\}$ be the set of agents, where each agent $i_k$ is defined by the tuple $\mathcal{A}_k = \langle P_k, M_k, T_k, R_k \rangle$, denoting base parameters, memory mechanisms, toolsets, and role configurations, respectively (Wang et al., 2024b; Xi et al., 2025; Schick et al., 2023). Instead of a fixed state, diversity is a range that can be smoothly adjusted; by changing how much these functional parts overlap or differ, we can control whether the system stays uniform or becomes highly specialized.

We argue that functional heterogeneity acts as a compensatory mechanism for individual cognitive bias, expanding the collective solution space while enabling cross-verification (Jiang et al., 2023). By carefully adjusting differences in toolsets ($T_k$), roles ($R_k$), memory structures ($M_k$), or reasoning modalities, one can strategically amplify complementary strengths among agents, thereby boosting collective performance (Ye et al., 2025). By synthesizing specialized perspectives, diversity provides a potential pathway to realize a positive collaboration gain ($\Gamma > 1$).

**Agent scale** defines the total number of agents, which forms the basis for interaction complexity. As scale increases, the expanded solution space and interaction paths provide a structural foundation for collaborative emergence. Specifically, coordinating over 1,000 agents facilitates a logistic growth pattern, allowing collaborative emergence to appear much earlier than in traditional neural scaling (Qian et al., 2025). For instance, the OASIS system supports millions of agents, demonstrating how larger scales can simulate stronger collective dynamics (Zhang et al., 2024). Similarly, in financial trading TradingAgents achieves better performance through multi-role collaboration (Xiao et al., 2025). In this framework, agent scale serves as a tunable parameter that opens a pathway for $\Gamma$ to manifest as significant collective intelligence through phase transitions in performance.

### 4.2.2. INFORMATION LEVEL

Information level characterizes the dynamic activation of potential into actual gains during execution. While control factors define the architectural boundaries, information factors—emerging in real time—serve as proxies for the system's internal cognitive trajectory. We include this level to provide mechanistic transparency into the collaborative process, enabling quantitative tracking of how static con-

figurations manifest as dynamic collaboration gain $\Gamma$. We analyze two dimensions: *intra-agent* uses content entropy to measure solution certainty and convergence, while *inter-agent* uses evolutionary distance to track semantic flow and interaction contributions.

**Content entropy** measures the certainty of the solution space by capturing intent distribution within the semantic feature space (Xia et al., 2025). The instantaneous content entropy at time $t$ is defined as:

$$H_t = -\sum_i p(x_i | C_t) \log p(x_i | C_t) \qquad (2)$$

where $H_t$ is the entropy at time $t$, $i$ indexes the set of possible discrete content types generated by agents, and $p(x_i \mid C_t)$ is the conditional probability of content $x_i$ given system state $C_t$.

In the MAS framework, content entropy $H_t$ serves as a quantitative proxy for the system's informational state transition. A monotonic decrease in $H_t$ characterizes the convergence trajectory from stochastic exploration toward a stabilized collective state (Navajas et al., 2018). Conversely, persistent high entropy signifies a coordination deficit, where the lack of effective constraints leaves the system in a state of decision-making divergence (Guo, 2025). However, since a reduction in entropy only describes agreement in a formal sense, it does not necessarily equate to a high-quality solution (Yang et al., 2025). For instance, in cases of "contextual breakdown," the system may exhibit pseudo-convergence because agents prematurely ignore key information. Therefore, a simple decrease in entropy cannot distinguish between valid consensus and degenerative information loss, requiring a deeper analysis of the interaction context, with more detailed examples provided in Appendix D.

**Evolutionary distance** characterizes the dynamic "work" of the system, quantifying the intensity of semantic displacement during the interaction process. This factor defines the degree to which interaction behaviors induce changes in the internal states of agents from an information-theoretic perspective (Jaques et al., 2019):

$$D_t = \sum_{i=1}^{N} \left(1 - \frac{\mathbf{v}_{i,t} \cdot \mathbf{v}_{i,t-1}}{\|\mathbf{v}_{i,t}\| \|\mathbf{v}_{i,t-1}\|}\right) \qquad (3)$$

where $D_t$ denotes the evolutionary distance of the system at time $t$, characterizing the overall update intensity of semantic content relative to the previous round. Here, $N$ represents the total number of agents in the MAS, and $i$ is the agent index. The term $\frac{\mathbf{v}_{i,t} \cdot \mathbf{v}_{i,t-1}}{\|\mathbf{v}_{i,t}\| \|\mathbf{v}_{i,t-1}\|}$ measures the cosine similarity between the state vectors(e.g., semantic state embeddings) of the $i$-th agent at consecutive time steps; subtracting this value from 1 gives the cosine distance, which reflects the change in the agent's semantic state.

In an ideal execution, an optimal evolutionary distance suggests the emergence of meaningful new information (Cover, 1999). However, abnormal fluctuations in this factor require a two-sided analysis: a change that is too small typically indicates that the system is stuck in redundant repetition. Conversely, an excessive change may not be a positive signal, as it often suggests a contextual breakdown, where agents may lose their connection to previous context, leading to outputs completely decoupled from the preceding logic. Therefore, only when evolutionary distance and content entropy maintain a balanced relationship can the system achieve effective performance growth while staying logically consistent (March, 1991). An empirical analysis with detailed case studies is provided in Appendix D.

## 5. Alternative Views

This paper proposes steering MAS design onto a structured, scientific track. While founded on robust logic, this approach may encounter skepticism. This section addresses three primary counterarguments.

### 5.1. Operational Complexity and Resource Costs

Critics argue that our framework is too difficult to implement in practice. They point out that analyzing so many different factors creates an overwhelming amount of work, and that measuring the specific collaboration gain requires experimental setups that are far too strict (Dorri et al., 2018). From this perspective, $\Gamma$ is seen as an unnecessary burden; critics prefer simpler, more direct metrics like absolute accuracy, which are much easier to track (Windl et al., 2022).

However, we argue that prioritizing short-term convenience over understanding why a system actually works creates a significant blind spot. Simple metrics fail to reveal whether performance stems from mere resource scaling or from genuine architectural synergy. Instead of continuing with blind trial-and-error, the field should move toward adopting collaboration gain as a precise feedback signal to evaluate the true benefits of specific design factors (Fang & Casadevall, 2011). While this rigorous evaluation requires more controlled experimentation, it provides the necessary transparency to distinguish between superficial improvements and fundamental architectural breakthroughs. This shift is essential for transforming MAS construction from an empirical craft into a predictable and reproducible science (Sharma et al., 2025).

### 5.2. Contradiction between Holism and Reductionism

Critics from systems theory argue that breaking a MAS down into individual factors is too simplistic. They worry that by focusing on separate parts, we lose sight of "emergence"—the way a system's collective behavior becomes greater than the sum of its parts (Mazzocchi, 2012; Pigliucci, 2014).

In response, we argue that systematic analysis is a necessary step toward understanding collective power. Just as we study individual genes to comprehend the complexity of life, we must identify specific design factors to see how they interact to create emergence (Fang & Casadevall, 2011). Rather than ignoring the "whole", we advocate for a perspective that provides a clearer roadmap to see how these individual "ingredients" combine to produce collective intelligence. By identifying the role of each factor, we move from simply observing that a system works to understanding how its internal components drive that success (Carrodano, 2025).

### 5.3. Ambiguities in Transitioning from Correlation to Causal Discovery

Skeptics suggest that identifying design factors yields only statistical correlations rather than true causal mechanisms, risking a static repository that fails to generalize to new environments (Bai, 2022; Shu & Ye, 2023).

In response, we argue that factor attribution is a necessary foundation for deeper understanding. We call upon the community to adopt $\Gamma$ as a standard feedback signal to reflect the true efficacy of collective collaboration. By using this signal, we can systematically identify which design factors actually drive performance gains. This process provides the essential data needed for future causal modeling—a methodology that maps the direct link between specific design "actions" and their resulting "effects". Establishing this foundation will pave the way toward a more transparent and designable framework, shifting the field from merely observing that a system works to predicting how to make it work, facilitating that collective intelligence becomes a reproducible outcome (Schölkopf, 2022).

## 6. Conclusion

Challenging the prevailing reliance on engineering-driven blind trial-and-error, this paper advocates for a paradigm shift to resolve the fundamental ambiguity of attribution in MAS. We propose an integrated framework centered on three core strategic recommendations. First, we advocate for adopting the collaboration gain metric as a scientific evaluation tool to decouple genuine collaboration gain from mere resource accumulation. Building upon this metric, we propose a factor attribution paradigm to systematically identify the true causal drivers of performance. Finally, we call for constructing a systematic MAS factor library to provide a structured design space for these attribution tasks. Collectively, these recommendations aim to transition the construction of MAS from empirical practices toward systematic scientific guidance. By ensuring that perfor-

mance improvements are theoretically traceable, we offer a roadmap for the field to evolve into a rigorous science, facilitating the reliable engineering of collective intelligence.

# References

Bai, H. The epistemology of machine learning. *Filosofija. Sociologija*, 2022.

Ballerini, M., Cabibbo, N., Candelier, R., Cavagna, A., Cisbani, E., Giardina, I., Lecomte, V., Orlandi, A., Parisi, G., Procaccini, A., et al. Interaction ruling animal collective behavior depends on topological rather than metric distance: Evidence from a field study. *Proceedings of the national academy of sciences(PNAS)*, 2008.

Bigeard, A., Nashold, L., Krishnan, R., and Wu, S. Finance agent benchmark: Benchmarking llms on real-world financial research tasks. *arXiv:2508.00828*, 2025.

Bo, X., Zhang, Z., Dai, Q., Feng, X., Wang, L., Li, R., Chen, X., and Wen, J.-R. Reflective multi-agent collaboration based on large language models. In *Conference on Neural Information Processing Systems(NeurIPS)*, 2024.

Bonabeau, E., Dorigo, M., and Theraulaz, G. *Swarm intelligence: from natural to artificial systems*. Oxford university press, 1999.

Bran, A. M., Cox, S., Schilter, O., Baldassari, C., White, A., and Schwaller, P. Augmenting large language models with chemistry tools. 2023.

Carrodano, C. General framework of nonlinear factor interactions using bayesian networks for risk analysis applied to road safety and public health. *Scientific Reports*, 2025.

Cavagna, A., Cimarelli, A., Giardina, I., Parisi, G., Santagati, R., Stefanini, F., and Viale, M. Scale-free correlations in starling flocks. *Proceedings of the National Academy of Sciences(PNAS)*, 2010.

Chan, C.-M., Chen, W., Su, Y., Yu, J., Xue, W., Zhang, S., Fu, J., and Liu, Z. Chateval: Towards better LLM-based evaluators through multi-agent debate. In *International Conference on Learning Representations(ICLR)*, 2024.

Chen, J., Saha, S., and Bansal, M. Reconcile: Round-table conference improves reasoning via consensus among diverse llms. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, 2024.

Chen, W., Su, Y., Zuo, J., Yang, C., Yuan, C., Chan, C.-M., Yu, H., Lu, Y., Hung, Y.-H., Qian, C., et al. Agentverse: Facilitating multi-agent collaboration and exploring emergent behaviors. In *International Conference on Learning Representations(ICLR)*, 2023.

Chen, W., Yuan, J., Qian, C., Yang, C., Liu, Z., and Sun, M. Optima: Optimizing effectiveness and efficiency for llm-based multi-agent system. 2025.

Couzin, I. D., Krause, J., Franks, N. R., and Levin, S. A. Effective leadership and decision-making in animal groups on the move. *Nature*, 2005.

Cover, T. M. *Elements of information theory*. 1999.

Dang, Y., Qian, C., Luo, X., Fan, J., Xie, Z., Shi, R., Chen, W., Yang, C., Che, X., Tian, Y., Xiong, X., Han, L., Liu, Z., and Sun, M. Multi-agent collaboration via evolving orchestration. In *Conference on Neural Information Processing Systems(NeurIPS)*, 2025.

Davis, R. and Smith, R. G. Negotiation as a metaphor for distributed problem solving. *Artificial intelligence*, 1983.

Decker, K. and Lesser, V. Quantitative modeling of complex computational task environments. In *AAAI*, 1993.

Decker, K. S. *Environment centered analysis and design of coordination mechanisms*. PhD thesis, 1995.

Di Caro, G. and Dorigo, M. Antnet: Distributed stigmergetic control for communications networks. *Journal of Artificial Intelligence Research*, 1998.

Di Caro, G., Ducatelle, F., and Gambardella, L. M. Anthocnet: an adaptive nature-inspired algorithm for routing in mobile ad hoc networks. *European transactions on telecommunications*, 2005.

Dorigo, M., Maniezzo, V., and Colorni, A. Ant system: optimization by a colony of cooperating agents. *IEEE transactions on systems, man, and cybernetics, part b (cybernetics)*, 1996.

Dorri, A., Kanhere, S. S., and Jurdak, R. Multi-agent systems: A survey. *IEEE Access*, 2018.

Du, Z., Qian, C., Liu, W., Xie, Z., Wang, Y., Qiu, R., Dang, Y., Chen, W., Yang, C., Tian, Y., Xiong, X., and Han, L. Multi-agent collaboration via cross-team orchestration. 2025.

Eberhart, R. C., Shi, Y., and Kennedy, J. *Swarm intelligence*. Elsevier, 2001.

Edmondson, A. Psychological safety and learning behavior in work teams. *Administrative science quarterly*, 1999.

Fang, F. C. and Casadevall, A. The evolution of molecular biology. *Nature Reviews Molecular Cell Biology*, 2011.

Franks, N. R. Army ants: A collective intelligence. 1989.

Fromm, J. Types and forms of emergence. 2005.

Galton, F. Vox populi. 1907.

Ghareeb, A. E., Chang, B., Mitchener, L., Yiu, A., Szostkiewicz, C. J., Laurent, J. M., Razzak, M. T., White, A. D., Hinks, M. M., and Rodriques, S. G. Robin: A multi-agent system for automating scientific discovery. *arXiv:2505.13400*, 2025.

Gong, R., Huang, Q., Ma, X., Noda, Y., Durante, Z., Zheng, Z., Terzopoulos, D., Fei-Fei, L., Gao, J., and Vo, H. Mindagent: Emergent gaming interaction. In *Findings of the Association for Computational Linguistics(NAACL)*, 2024.

Guo, S., Shariatmadari, A. H., Xiong, G., Huang, A., Kim, M., Williams, C. M., Bekiranov, S., and Zhang, A. Ideabench: Benchmarking large language models for research idea generation. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 2*, 2025.

Guo, T., Chen, X., Wang, Y., Chang, R., Pei, S., Chawla, N. V., Wiest, O., and Zhang, X. Large language model based multi-agents: A survey of progress and challenges. 2024.

Guo, X. Measuring reasoning utility in llms via conditional entropy reduction. 2025.

Hoel, E. P., Albantakis, L., Marshall, W., and Tononi, G. Can the macro beat the micro? integrated information across spatiotemporal scales. *Neuroscience of Consciousness*, 2016.

Hoff, M. P. Mountain gorillas: Three decades of research at karisoke, 2002.

Hoff, N., Wood, R., and Nagpal, R. Distributed colony-level algorithm switching for robot swarm foraging. 2013.

Hong, L. and Page, S. E. Groups of diverse problem solvers can outperform groups of high-ability problem solvers. *Proceedings of the National Academy of Sciences*, 2004.

Hong, S., Zhuge, M., Chen, J., Zheng, X., Cheng, Y., Wang, J., Zhang, C., Wang, Z., Yau, S. K. S., Lin, Z., et al. Metagpt: Meta programming for a multi-agent collaborative framework. 2023.

Hu, X.-m., Zhang, J., Xiao, J., and Li, Y. Protein folding in hydrophobic-polar lattice model: a flexible ant-colony optimization approach. *Protein and Peptide Letters*, 2008.

Huot, F., Amplayo, R. K., Palomaki, J., Jakobovits, A. S., Clark, E., and Lapata, M. Agents' room: Narrative generation through multi-step collaboration. In *The Thirteenth International Conference on Learning Representations(ICLR)*.

Isaeva, V. Self-organization in biological systems. *Biology Bulletin*, 2012.

Jaques, N., Lazaridou, A., Hughes, E., Çaglar Gülçehre, Ortega, P. A., Strouse, D., Leibo, J. Z., and de Freitas, N. Social influence as intrinsic motivation for multi-agent deep reinforcement learning. In *ICML*, 2019.

Jiang, D., Ren, X., and Lin, B. Y. Llm-blender: Ensembling large language models with pairwise ranking and generative fusion. 2023.

Jing, L., Huang, Z., Wang, X., Yao, W., Yu, W., Ma, K., Zhang, H., Du, X., and Yu, D. Dsbench: How far are data science agents from becoming data science experts? In *International Conference on Learning Representations(ICLR)*, 2025.

Kapoor, S., Stroebl, B., Siegel, Z. S., Nadgir, N., and Narayanan, A. AI agents that matter. *Transactions on Machine Learning Research(TMLR)*, 2025.

Karaboga, D. and ASLAN, S. A discrete artificial bee colony algorithm for detecting transcription factor binding sites in dna sequences. *Genetics and Molecular Research*, 2016.

Karaboga, D. and Basturk, B. A powerful and efficient algorithm for numerical function optimization: artificial bee colony (abc) algorithm. *Journal of global optimization*, 2007.

Kennedy, J. and Eberhart, R. Particle swarm optimization. In *Proceedings of ICNN'95-international conference on neural networks*, 1995.

Kernbach, S., Thenius, R., Kernbach, O., and Schmickl, T. Re-embodiment of honeybee aggregation behavior in an artificial micro-robotic system. *Adaptive Behavior*, 2009.

Kim, Y., Gu, K., Park, C., Park, C., Schmidgall, S., Heydari, A. A., Yan, Y., Zhang, Z., Zhuang, Y., Malhotra, M., et al. Towards a science of scaling agent systems. *arXiv:2512.08296*, 2025.

Kuhn, L., Gal, Y., and Farquhar, S. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation.

Kuhn, T. S. and Hacking, I. The structure of scientific revolutions. 1970.

Labrou, Y. and Finin, T. W. Semantics and conversations for an agent communication language. 1997.

Labrou, Y., Finin, T., and Peng, Y. Agent communication languages: The current landscape. *IEEE Intelligent Systems and Their Applications*, 1999.

LD, D. B. and Krishna, P. V. Honey bee behavior inspired load balancing of tasks in cloud computing environments. *Applied soft computing*, 2013.

Li, G., Hammoud, H., Itani, H., Khizbullin, D., and Ghanem, B. Camel: Communicative agents for" mind" exploration of large language model society. *Conference on Neural Information Processing Systems(NeurIPS)*, 2023.

Li, X., Wang, S., Zeng, S., Wu, Y., and Yang, Y. A survey on llm-based multi-agent systems: workflow, infrastructure, and challenges. *Vicinagearth*, 2024.

Lin, B. Y., Zhou, W., Shen, M., Zhou, P., Bhagavatula, C., Choi, Y., and Ren, X. Commongen: A constrained text generation challenge for generative commonsense reasoning. In *Findings of the Association for Computational Linguistics (ACL)*, 2020.

Liu, F., Niu, Y., Zhang, Q., Wang, K., Dong, Z., Wong, I. N., Cheng, L., Li, T., Duan, L., Li, K., et al. A foundational architecture for ai agents in healthcare. *Cell Reports Medicine*, 2025a.

Liu, F., Yang, Z.-R., Liu, C., SONG, T., Gao, X., and Liu, H. MM-agent: LLM as agents for real-world mathematical modeling problem. In *Conference on Neural Information Processing Systems(NeurIPS)*, 2025b.

Liu, Y. and Niu, B. A novel pso model based on simulating human social communication behavior. *Discrete Dynamics in Nature and Society*, 2012.

Lorenz, J., Rauhut, H., Schweitzer, F., and Helbing, D. How social influence can undermine the wisdom of crowd effect. *Proceedings of the national academy of sciences(PNAS)*, 2011.

Lumer, E. D. and Faieta, B. Diversity and adaptation in populations of clustering ants. 1994.

Ma, M., Zhang, J., Yang, F., Kang, Y., Lin, Q., Yang, T., Rajmohan, S., and Zhang, D. Dover: Intervention-driven auto debugging for llm multi-agent systems. *arXiv:2512.06749*, 2025.

Maldonado, D., Cruz, E., Torres, J. A., Cruz, P. J., and Benitez, S. d. P. G. Multi-agent systems: A survey about its components, framework and workflow. *IEEE Access*, 2024.

Malfa, E. L., Malfa, G. L., Marro, S., Zhang, J. M., Black, E., Luck, M., Torr, P., and Wooldridge, M. J. Large language models miss the multi-agent mark. In *Conference on Neural Information Processing Systems(NeurIPS)*, 2025.

Malone, T. W. and Crowston, K. The interdisciplinary study of coordination. *ACM Computing Surveys (CSUR)*, 1994.

March, J. G. Exploration and exploitation in organizational learning. *Organization science*, 1991.

Mazzocchi, F. Complexity and the reductionism–holism debate in systems biology. *Wiley Interdisciplinary Reviews: Systems Biology and Medicine*, 2012.

Mech, L. D. Alpha status, dominance, and division of labor in wolf packs. *Canadian journal of zoology*, 1999.

Minaee, S., Mikolov, T., Nikzad, N., Chenaghlu, M., Socher, R., Amatriain, X., and Gao, J. Large language models: A survey. *arXiv:2402.06196*, 2024.

Mirjalili, S., Mirjalili, S. M., and Lewis, A. Grey wolf optimizer. *Advances in engineering software*, 2014.

Mohammadi, M., Li, Y., Lo, J., and Yip, W. Evaluation and benchmarking of llm agents: A survey. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (SIGKDD)*, 2025.

Navajas, J., Niella, T., Garbulsky, G., Bahrami, B., and Sigman, M. Aggregated knowledge from a small number of debates outperforms the wisdom of large crowds. *Nature Human Behaviour*, 2018.

Nikitin, A., Kossen, J., Gal, Y., and Marttinen, P. Kernel language entropy: Fine-grained uncertainty quantification for llms from semantic similarities. *Conference on Neural Information Processing Systems(NeurIPS)*, 2024.

Nowak, M. A. *Evolutionary dynamics: exploring the equations of life*. Harvard university press, 2006.

Pan, M. Z., Cemri, M., Agrawal, L. A., Yang, S., Chopra, B., Tiwari, R., Keutzer, K., Parameswaran, A., Ramchandran, K., Klein, D., et al. Why do multiagent systems fail? 2025.

Panait, L. and Luke, S. Cooperative multi-agent learning: The state of the art. *Autonomous agents and multi-agent systems*, 2005.

Park, J. S., O'Brien, J. C., Cai, C. J., Morris, M. R., Liang, P., and Bernstein, M. S. Generative agents: Interactive simulacra of human behavior. 2023.

Pigliucci, M. Between holism and reductionism: a philosophical primer on emergence. *Biological Journal of the Linnean Society*, 2014.

Qian, C., Liu, W., Liu, H., Chen, N., Dang, Y., Li, J., Yang, C., Chen, W., Su, Y., Cong, X., Xu, J., Li, D., Liu, Z., and Sun, M. Chatdev: Communicative agents for software development. In *Annual Meeting of the Association for Computational Linguistics(ACL)*, 2024.

Qian, C., Xie, Z., Wang, Y., Liu, W., Zhu, K., Xia, H., Dang, Y., Du, Z., Chen, W., Yang, C., Liu, Z., and Sun, M. Scaling large language model-based multi-agent collaboration. In *International Conference on Learning Representations(ICLR)*, 2025.

Renney, H., Nethercott, M. N., Renney, N., and Hayes, P. Llm-enabled multi-agent systems: Empirical evaluation and insights into emerging design patterns & paradigms, 2026.

Reynolds, C. W. Flocks, herds and schools: A distributed behavioral model. In *Proceedings of the 14th annual conference on Computer graphics and interactive techniques*, 1987.

Schaller, G. B. *The Serengeti lion: a study of predator-prey relations*. 2009.

Schick, T., Dwivedi-Yu, J., Dessì, R., Raileanu, R., Lomeli, M., Hambro, E., Zettlemoyer, L., Cancedda, N., and Scialom, T. Toolformer: Language models can teach themselves to use tools. *Conference on Neural Information Processing Systems(NeurIPS)*, 2023.

Schipper, O., Zhang, Y., Du, Y., Pechenizkiy, M., and Fang, M. Pillagerbench: Benchmarking llm-based agents in competitive minecraft team environments. 2025.

Schölkopf, B. Causality for machine learning. 2022.

Seeley, T. D. *The wisdom of the hive: the social physiology of honey bee colonies*. 2009.

Seeley, T. D. *Honeybee democracy*. 2011.

Seeley, T. D., Visscher, P. K., and Passino, K. M. Group decision making in honey bee swarms: When 10,000 bees go house hunting, how do they cooperatively choose their new nesting site? *American scientist*, 2006.

Seff, A., Cera, B., Chen, D., Ng, M., Zhou, A., Nayakanti, N., Refaat, K. S., Al-Rfou, R., and Sapp, B. Motionlm: Multi-agent motion forecasting as language modeling. In *International Conference on Computer Vision (ICCV)*, 2023.

Sharma, G., Munteanu, V., Ghiasi, N. M., Banerjee, J., Varma, S., Foschini, L., Ellrott, K., Mutlu, O., Ciorb, D., and Ophoff, R. A. A decentralized future for the open-science databases. 2025.

Sherman, P. W., Jarvis, J. U., and Alexander, R. D. *The biology of the naked mole-rat*. Princeton University Press, 2017.

Shinn, N., Cassano, F., Gopinath, A., Narasimhan, K. R., and Yao, S. Reflexion: language agents with verbal reinforcement learning. In *Conference on Neural Information Processing Systems(NeurIPS)*, 2023.

Shu, X. and Ye, Y. Knowledge discovery: Methods from data mining and machine learning. *Social Science Research*, 2023.

Soleimani-Pouri, M., Rezvanian, A., and Meybodi, M. R. Finding a maximum clique using ant colony optimization and particle swarm optimization in social networks. 2013.

Sumpter, D. J. Collective animal behavior. In *Collective animal behavior*. 2010.

Sun, H., Zhang, S., Niu, L., Ren, L., Xu, H., Fu, H., Zhao, F., Yuan, C., and Wang, X. Collab-overcooked: Benchmarking and evaluating large language models as collaborative agents. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2025.

Tang, B., Liang, H., Jiang, K., and Dong, X. On the importance of task complexity in evaluating LLM-based multi-agent systems. In *NeurIPS 2025 Workshop on Scaling Environments for Agents*, 2025.

Tetlock, P. E. The impact of accountability on judgment and choice: Toward a social contingency model. In *Advances in experimental social psychology*. 1992.

Tran, K.-T., Dao, D., Nguyen, M.-D., Pham, Q.-V., O'Sullivan, B., and Nguyen, H. D. Multi-agent collaboration mechanisms: A survey of llms. *arXiv:2501.06322*, 2025.

Van Dyke Parunak, H. and Brueckner, S. Entropy and self-organization in multi-agent systems. In *Proceedings of the fifth international conference on Autonomous agents*, 2001.

Wang, G., Xie, Y., Jiang, Y., Mandlekar, A., Xiao, C., Zhu, Y., Fan, L., and Anandkumar, A. Voyager: An open-ended embodied agent with large language models. *Transactions on Machine Learning Research(TMLR)*, 2024a.

Wang, H., Zhao, S., Wang, J., Qiang, Z., Qin, B., and Liu, T. Beyond frameworks: Unpacking collaboration strategies in multi-agent systems, 2025a.

Wang, J., WANG, J., Athiwaratkun, B., Zhang, C., and Zou, J. Mixture-of-agents enhances large language model capabilities, 2025b.

Wang, L., Ma, C., Feng, X., Zhang, Z., Yang, H., Zhang, J., Chen, Z.-Y., Tang, J., Chen, X., Lin, Y., Zhao, W. X., Wei, Z., and Wen, J. A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 2024b.

Wang, Q., Wang, T., Tang, Z., Li, Q., Chen, N., Liang, J., and He, B. Megaagent: A large-scale autonomous llm-based multi-agent system without predefined sops. In *ACL (Findings)*, 2025c.

Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E. H., Le, Q. V., and Zhou, D. Chain-of-thought prompting elicits reasoning in large language models. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2022.

White, I., Nottingham, K., Maniar, A., Robinson, M., Lillemark, H., Maheshwari, M., Qin, L., and Ammanabrolu, P. Collaborating action by action: A multi-agent llm framework for embodied reasoning. 2025.

Windl, M., Feger, S. S., Zijlstra, L., Schmidt, A., and Wozniak, P. W. 'it is not always discovery time': Four pragmatic approaches in designing ai systems. In *Conference on Human Factors in Computing Systems(CHI)*, 2022.

Woolley, A. W., Chabris, C. F., Pentland, A., Hashmi, N., and Malone, T. W. Evidence for a collective intelligence factor in the performance of human groups. *science*, 2010.

Wu, Q., Bansal, G., Zhang, J., Wu, Y., Li, B., Zhu, E., Jiang, L., Zhang, X., Zhang, S., Liu, J., Awadallah, A. H., White, R. W., Burger, D., and Wang, C. Autogen: Enabling nextgen LLM applications via multi-agent conversations. In *Conference on Language Modeling (COLM)*, 2024.

Wu, Y., Li, D., Chen, Y., Jiang, R., Zou, H. P., Huang, W.-C., Li, Y., Fang, L., Wang, Z., and Yu, P. S. Multi-agent autonomous driving systems with large language models: A survey of recent advances. 2025.

Xi, Z., Chen, W., Guo, X., He, W., Ding, Y., Hong, B., Zhang, M., Wang, J., Jin, S., Zhou, E., et al. The rise and potential of large language model based agents: A survey. *Science China Information Sciences*, 2025.

Xia, Z., Xu, J., Zhang, Y., and Liu, H. A survey of uncertainty estimation methods on large language models. 2025.

Xiao, Y., Sun, E., Luo, D., and Wang, W. Tradingagents: Multi-agents LLM financial trading framework. In *The First MARW: Multi-Agent AI in the Real World Workshop at AAAI 2025*, 2025.

Yang, X., Li, L., Zhou, H., Zhu, T., Qu, X., Fan, Y., Wei, Q., Ye, R., Kang, L., Qin, Y., Kou, Z., Liu, D., Li, Q., Ding, N., Chen, S., and Shao, J. Toward efficient agents: Memory, tool learning, and planning. 2026.

Yang, Y., Ma, Y., Feng, H., Cheng, Y., and Han, Z. Minimizing hallucinations and communication costs: Adversarial debate and voting mechanisms in llm-based multi-agents. *Applied Sciences*, 2025.

Yao, S., Zhao, J., Yu, D., Du, N., Shafran, I., Narasimhan, K. R., and Cao, Y. React: Synergizing reasoning and acting in language models. 2022.

Ye, R., Liu, X., Wu, Q., Pang, X., Yin, Z., Bai, L., and Chen, S. X-mas: Towards building multi-agent systems with heterogeneous llms. 2025.

Zhang, G., Niu, L., Fang, J., Wang, K., BAI, L., and Wang, X. Multi-agent architecture search via agentic supernet. In *International Conference on Machine Learning(ICML)*, 2025a.

Zhang, G., Yue, Y., Li, Z., Yun, S., Wan, G., Wang, K., Cheng, D., Yu, J. X., and Chen, T. Cut the crap: An economical communication pipeline for llm-based multi-agent systems. In *International Conference on Learning Representations(ICLR)*, 2025b.

Zhang, H., Cui, Z., Chen, J., Wang, X., Zhang, Q., Wang, Z., Wu, D., and Hu, S. Stop overvaluing multi-agent debate – we must rethink evaluation and embrace model heterogeneity. 2025c.

Zhang, Z., Yin, Z., and Shao, J. Genss: A generalized and scalable llm-based agents social network simulator. In *NeurIPS 2024 Workshop on Open-World Agents*, 2024.

Zhou, C., Li, Q., Li, C., Yu, J., Liu, Y., Wang, G., Zhang, K., Ji, C., Yan, Q., He, L., Peng, H., Li, J., Wu, J., Liu, Z., Xie, P., Xiong, C., Pei, J., Yu, P. S., and Sun, L. A comprehensive survey on pretrained foundation models: a history from bert to chatgpt. *Int. J. Mach. Learn. Cybern.*, 2025a.

Zhou, H., Wan, X., Sun, R., Palangi, H., Iqbal, S., Vulić, I., Korhonen, A., and Arık, S. Ö. Multi-agent design: Optimizing agents with better prompts and topologies. *arXiv:2502.02533*, 2025b.

Zhu, C., Hong, S., Wu, J., Chawla, K., Tang, C., Yin, Y., Wolfe, N., Babinsky, E., and Liu, D. Raffles: Reasoning-based attribution of faults for llm systems, 2026.

Zhu, K., Du, H., Hong, Z., Yang, X., Guo, S., Wang, Z., Wang, Z., Qian, C., Tang, R., Ji, H., and You, J. Multiagentbench: Evaluating the collaboration and competition of llm agents. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, 2025.

Zou, J., Yang, X., Qiu, R., Li, G., Tieu, K., Lu, P., Shen, K., Tong, H., Choi, Y., He, J., et al. Latent collaboration in multi-agent systems. *arXiv:2511.20639*, 2025.

# Appendix

# A. Theoretical Background of Collective Intelligence

### A.1. Collective Intelligence in Biological Systems

The concept of *Swarm Intelligence* in artificial intelligence (Eberhart et al., 2001) originates from the observation of cooperative behaviors in biological populations (Bonabeau et al., 1999). Cooperative swarm behaviors of organisms can typically be categorized into three paradigms.

- **Decentralized Swarm Intelligence.** In this paradigm, there is no leader or centralized controller within the system. Instead, the complexity and functionality of swarm behavior arise entirely from local interactions among individuals. The macroscopic performance of the collective far exceeds the sum of individual capabilities. For example, army ants exhibit highly organized "raiding" formations during foraging and can dynamically construct living bridges made of tens of thousands of ants to cross obstacles (Franks, 1989). Another iconic example is the murmuration of starlings, where thousands of birds perform synchronous, fluid, and instantaneous aerial maneuvers to evade predators. Studies show that such large-scale coordination occurs without a leader; each bird only needs to align its speed and direction with those of its 6–7 nearest neighbors (Ballerini et al., 2008). These local rules propagate through the network, allowing perturbations to travel across the entire flock (Cavagna et al., 2010), resulting in sophisticated collective behaviors that far surpass the intelligence of individual birds.

- **Centralized Organization.** In this paradigm, one or a few core individuals exert decisive, top-down influence on the structure and function of the collective through their physiological state or behavior. A representative case is the naked mole-rat colony, in which a single breeding female (the queen) dominates reproduction. She releases specific chemical signals and engages in aggressive behaviors that physiologically suppress reproduction in all other females, forcing them to serve as workers for the colony (Sherman et al., 2017). Similarly, gorilla groups are led by a dominant male known as the "Silverback." Although he does not micromanage individual activities, he holds ultimate authority over decisions such as when the group moves, rests, or changes routes, and he mediates internal conflicts while protecting the group from external threats (Hoff, 2002).

- **Hybrid Organization.** This paradigm combines centralized and decentralized elements. While the group has a core individual exerting special influence, many key decisions are still made collectively in a distributed manner. The European honeybee colony exemplifies this hybrid structure. The queen serves as the sole reproductive individual and maintains social cohesion through pheromones—a hallmark of centralization (Seeley, 2009). However, when the colony must find a new nest site, the decision-making process becomes highly democratic. Scout bees "debate" candidate locations using the symbolic waggle dance. The final decision is not made by the queen, but through consensus once the strength of dance signals for one location reaches a quorum (Seeley, 2011).

Although individual organisms generally possess limited intelligence, evolution has endowed them with an extraordinary diversity of communication mechanisms, as illustrated in the Figure 4. These mechanisms span chemical and physical modalities and even forms of social learning analogous to human behavior, transmitting rich information either explicitly or implicitly.

In decentralized systems, numerous individuals following simple rules can self-organize into complex and adaptive macroscopic patterns (Isaeva, 2012), giving rise to emergent intelligent behavior. Such phenomena have inspired a wide range of computational algorithms (Sumpter, 2010). For instance, implicit communication protocols in nature often rely on stigmergy: individuals release chemical cues that modify the environment, and others interpret these environmental changes to coordinate indirectly. This mechanism inspired multi-agent control systems and ant colony optimization algorithms. The Table 1 summarizes representative algorithms inspired by cooperative behaviors observed in biological groups.

### A.2. Collective Intelligence in Human Societies

**Collective intelligence is equally evident in human societies.** At the macro level, millennia of social development and evolution are themselves products of collective intelligence. Most achievements of civilization did not spring from a single genius; they were established gradually through long-term, large-scale collaboration and social production among countless individuals. In both history and current practice, aggregates often outperform individuals. For example, Galton observed in 1907 that the average of hundreds of fairgoers' estimates of an ox's weight was almost exactly correct (Galton, 1907).
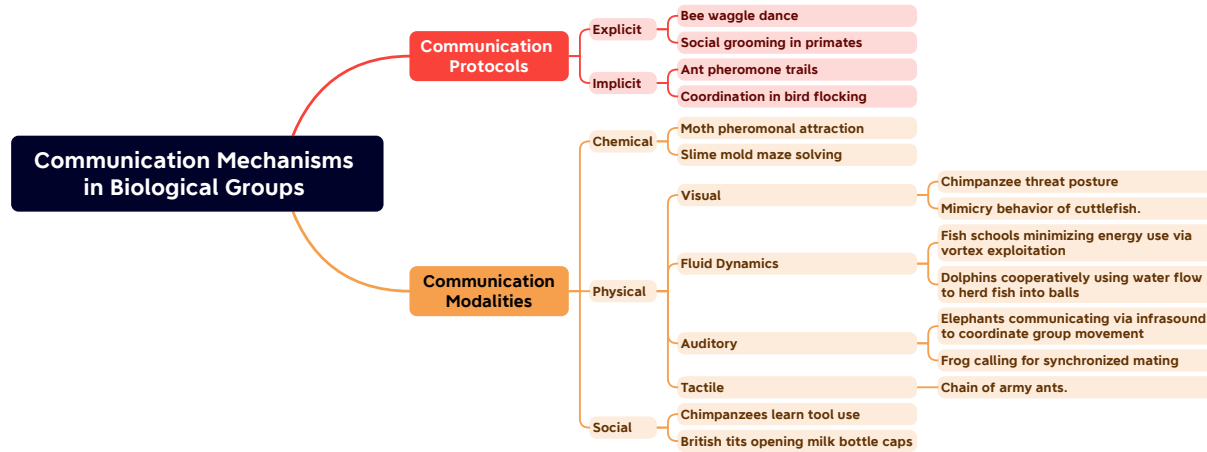
*Figure 4.* The taxonomy of communication mechanisms in biological groups.

*Table 1.* Representative algorithms inspired by behaviors in biological groups

| Application Domain | Representative Algorithm | Biological Inspiration |
|---|---|---|
| Search and Optimization | Ant Colony Optimization (ACO) (Dorigo et al., 1996) <br> Particle Swarm Optimization (PSO) (Kennedy & Eberhart, 1995) <br> Artificial Bee Colony (ABC) (Karaboga & Basturk, 2007) <br> Grey Wolf Optimizer (GWO) (Mirjalili et al., 2014) | Ants finding the shortest path via pheromone trails <br> Social search behavior of bird flocks and fish schools <br> Intelligent foraging and information sharing in honeybee colonies <br> Social hierarchy and cooperative hunting in grey wolves |
| Robotics Control | Boids Algorithm (Reynolds, 1987) <br> BEECLUST Algorithm (Kernbach et al., 2009) <br> Virtual Pheromone Algorithm (Hoff et al., 2013) | Collective flight of bird flocks <br> Aggregation behavior of young honeybees in warm hive regions <br> Pheromone trail following in ant foraging |
| Bioinformatics | ACO for Protein Folding (Hu et al., 2008) <br> ABC for Motif Discovery (Karaboga & ASLAN, 2016) <br> Ant-Based Clustering (Lumer & Faieta, 1994) | Ants finding shortest paths via pheromone trails <br> Intelligent foraging behavior of bee colonies <br> Ant clustering of corpses and larvae |
| Computer Networks | AntNet (Di Caro & Dorigo, 1998) <br> AntHocNet (Di Caro et al., 2005) <br> Honeybee Behavior Load Balancing (LD & Krishna, 2013) | Indirect communication in ant foraging <br> Ant-based path discovery and maintenance <br> Foraging and information-sharing behavior of honeybees |
| Social Networks | PSO as a Social Model (Liu & Niu, 2012) <br> ACO for Community Detection (Soleimani-Pouri et al., 2013) | Collective motion and information sharing in bird and fish groups <br> Path construction via pheromone trails in ants |

In the internet era, large-scale online collaboration makes the manifestation of human collective intelligence even more salient. Networked communication allows thousands of people to interact, share, and cooperate rapidly, producing knowledge and making decisions together. To synthesize research on collective intelligence in human societies, this section organizes the literature around five core dimensions (see Figure 5):

- **Group goals and task types.** Group goals set the direction for collective action, while task types shape *how—and how readily—*collective intelligence can manifest. Clear, shared goals align individual efforts; vague or contested goals weaken focus and coordination. Task characteristics further determine fit: highly complex tasks that require complementary knowledge and skills suit group collaboration; very simple tasks or those reliant on rare individual insight may suit individuals.

  Decomposability matters: decomposable tasks benefit from division of labor and parallelism; non-decomposable tasks require tight global coordination. From organizational behavior and management, task complexity and variety create scope for diversity, independence, and aggregation (Hong & Page, 2004; Woolley et al., 2010). Evidence from organizational studies shows that strong teams are defined not by size or hierarchy but by psychological safety and shared meaning (Edmondson, 1999); members feel safe to speak up and drive work from the bottom up.

- **Group structure.** Structure determines who participates and how people connect and interact, shaping both production and quality of collective intelligence. Small groups often rely on dense communication; large groups reduce coordination costs through layered or networked designs that range from strict hierarchies to loose, market-like forms. Functional,

*Figure 5.* A taxonomy of mechanisms enabling collective intelligence in human societies.

cognitive, and identity diversity expand the knowledge base, foster innovation, and sharpen collective judgment. Decentralization is a key condition for effective bottom-up collaboration: distributed power and knowledge let individuals act on local information, reducing central bottlenecks and increasing adaptability. By contrast, when structures are closed, power is concentrated and communication is restricted; dissent is suppressed and decision quality declines (Tetlock, 1992).

- **Information diffusion and communication.** Communication mechanisms determine how knowledge flows and integrates across members. Effective communication aggregates dispersed information into collective insight, enabling the group to outperform individuals. Technology greatly widens the scope and depth of interaction. Open collaboration platforms let tens of thousands share documents, exchange views, and correct errors in real time, achieving high coordination without central control.

  Design still matters: overload, redundant diffusion, and herding weaken performance. For example, if discussions revolve only around known information and ignore heterogeneous sources, collective judgment can drift (Lorenz et al., 2011). Good mechanism design provides both *channels* and *filters*: it ensures diverse views enter the public space and are transformed into group wisdom through aggregation.

- **Coordination and decision-making.** Coordination and decision rules determine how dispersed actions and opinions translate into collective action. Coordination can be centralized, led by core individuals or institutions, suitable for urgent or highly coupled tasks, or decentralized, relying on self-organization and local feedback (Malone & Crowston, 1994). Decentralized coordination is common in open communities, prediction markets, and open-source software. It emphasizes consensus, voting, market signals, or algorithmic aggregation to turn diverse opinions into group decisions. Effective mechanisms must aggregate heterogeneous information while preventing authority capture, balancing broad deliberation with timely closure and avoiding domination by a few voices.

- **Incentives and norms.** Incentives motivate participation; norms ensure orderly and trustworthy cooperation. Incentives include money, reputation, achievement, interest, and mission. Many open projects rely on non-monetary incentives, especially reputation, to sustain large-scale voluntary effort, whereas prediction markets or bounty platforms rely on financial rewards. Norms align individual behavior with group goals through formal rules, shared values, and informal conventions. Together they shape participation depth, psychological safety, and collective identity, strengthening bottom-up collaborative dynamics. Poorly designed incentives and norms can intensify conformity pressure and superficial consensus, undermining independence and diversity.

16

The following cases illustrate these five dimensions in practice (Table 2).

*Table 2.* Comparative analysis of human collective intelligence across five dimensions.

| Case | Group Goals & Task Types | Group Structure | Information Diffusion & Communication | Coordination & Decision-Making | Incentives & Norms |
|---|---|---|---|---|---|
| Wikipedia | Build and continuously update an open encyclopedia; tasks modularized by article. | Decentralized volunteer community; a small number of administrators for governance. | Public revision history and talk pages; transparent, traceable edits. | Consensus first, supplemented by voting and administrator arbitration. | Intrinsic motivation and reputation; adherence to policies and patrolling guidelines. |
| Linux open source community | Develop and maintain the OS kernel; modular subsystems evolve in parallel. | Peer based layered structure; maintainers steward subsystems and integration. | Public mailing lists and code hosting records for discussion and review. | "Propose, review, and merge" workflow; technical consensus; maintainers decide when needed. | Mix of intrinsic motivation and corporate support; coding standards and code review ensure quality. |
| Waze navigation crowdsourcing | Real time aggregation of traffic and map updates; data crowdsourcing and rapid iteration. | Mass users plus tiered volunteer editors co maintain the map. | In app data collection and community forums; real time distribution of traffic information. | Algorithmic route coordination combined with permissioned human moderation. | Gamified points and reputation; editor guidelines and community rules. |
| Prediction markets | Aggregate probabilistic judgments on event outcomes via trading contracts. | Open, decentralized trading community; no fixed hierarchy. | Prices immediately reflect information and are visible to all participants. | Market clearing yields collective forecasts; independent trades drive rapid calibration. | Monetary incentives and market rules; anti manipulation safeguards and compliance. |
| Foldit citizen science game | Gamified search for optimal protein folds. | Global players with self organization; team collaboration in parallel. | In game real time feedback; community forums share strategies and solutions. | Leaderboards and team merging to combine superior solutions. | Badges and rankings plus scientific mission; anti cheating norms. |

## A.3. Collective Intelligence in LLM-based Multi-Agent Systems

Unlike biological collectives that rely on low-dimensional stimulus–response mechanisms and scale effects, or human societies that depend on institutions and cultural norms, LLM-based collective intelligence is characterized by natural language as the core medium of coordination and language agents—autonomous entities capable of task understanding and reasoning—as the fundamental units. Through linguistic consensus, such systems achieve internal self-organization and collaboration, marking a new stage of explicit symbolization and semanticization in the generation of intelligence (Qian et al., 2024; Park et al., 2023; Wu et al., 2024). The following discussion examines its features from three perspectives: the characteristics of interactive agents, their interaction patterns, and their organizational architectures.

- **Characteristics of Interactive Agents.** Language agents, powered by large-scale language models, integrate multifaceted abilities such as language understanding, knowledge reasoning, and task execution. They exhibit a high degree of autonomy and adaptability, enabling goal-directed reasoning and self-optimization within complex semantic environments (Yao et al., 2022; Shinn et al., 2023). Autonomy manifests in their capacities for intent generation, independent decision-making, and self-planning (Bran et al., 2023). A language agent can independently interpret task descriptions expressed in natural language, decompose objectives, generate reasoning chains, and form internal representations of execution paths. Adaptability, on the other hand, is reflected in their ability to perform linguistic reflection and strategic adjustment based on environmental feedback and task dynamics, thereby establishing a closed loop of self-improvement (Wang et al., 2024a; Park et al., 2023). Notably, interactive agents' cognitive profiles mirror human-like planning and behavioral adjustment through language-based reasoning. This synergy of autonomy and adaptability defines their intelligent individuality, enabling task comprehension and dynamic optimization—the essential building blocks of LLM-based collective intelligence.

- **Interaction Patterns of Interactive Agents.** Just as humans rely on natural language for communication and collaboration, interactive agents also use natural language as the central medium of coordination, forming the key mechanism of LLM-based collective intelligence. Unlike traditional MAS that communicate through fixed protocols (e.g., FIPA-ACL or KQML) and struggle to handle dynamic contexts or task shifts (Labrou & Finin, 1997; Labrou et al., 1999), language agents can dynamically generate, interpret, and adjust semantic structures during interaction, achieving semantic coordination through continual linguistic negotiation (Wu et al., 2024; Li et al., 2023). At the collective level, natural language interaction serves not only as a channel for information exchange but also as a driving force for knowledge sharing and group evolution (Park et al., 2023; Hong et al., 2023). Through continuous semantic negotiation and feedback, interactive agents can gradually form collective memory and shared norms under decentralized conditions, thereby giving rise to stable structures of collaboration. Similar to the cultural role of language

in human societies, language-agent communities can achieve a transition from individual understanding to collective cognition through the accumulation of semantic consensus. The interactions among agents—characterized by semantic explicitness, dynamic coordination, and co-construction of knowledge—endow collective intelligence with enhanced flexibility, adaptability, and sociality, forming a language-centered paradigm of "symbolic collective intelligence."

- **Organizational Architectures of Interactive Agents.** Analogous to biological collectives, groups of interactive agents display diverse and flexible organizational forms. In nature, wolf packs and lion prides rely on dominant individuals for centralized command and coordination, representing typical centralized architectures (Mech, 1999; Schaller, 2009); flocks of birds and schools of fish, conversely, exhibit self-organized global order through local perception and dynamic feedback (Couzin et al., 2005; Reynolds, 1987); bee colonies—while centered around the queen—combine symbolic communication mechanisms (e.g., pheromones, waggle dances) with distributed decision-making (Seeley et al., 2006), demonstrating a hybrid structure where centralization and decentralization coexist. Similarly, the organizational structures of language-agent systems can be broadly categorized into centralized, decentralized, and hybrid types. In a centralized architecture, a core agent acts as a "central planner," responsible for unified global task management; in a decentralized architecture, agents interact on an equal footing, and group behavior emerges spontaneously from individual interactions; hybrid architectures combine centralized global planning with decentralized local autonomy, maintaining task coherence while preserving agent independence. Selecting or composing an optimal architecture necessitates a strategic trade-off between global control, individual autonomy, and coordination efficiency, tailored specifically to the task's scale and environmental complexity.

Building upon this foundation, an intriguing question arises: when multiple language-model-based agents form a collective, their interaction patterns become almost indistinguishable from the communication modes of human societies in both form and function. They likewise rely on language as the medium of coordination, achieving shared understanding through semantic negotiation. At the organizational level, they exhibit structural differentiation similar to that observed in human or biological collectives — such as hierarchical role systems, decentralized collaboration networks, and hybrid dynamic organizational forms. This structural and interactive resemblance leads to a profound hypothesis: could collectives of language-model agents also give rise to a form of "collective intelligence" analogous to that found in biological or human societies?

## B. Technical Details of Collaboration Gain

### B.1. Task-Specific Baseline Implementation

We emphasize that the operational instantiation of the baseline $\Phi_S$ is not a fixed constant but a task-dependent "cognitive ceiling." It must be adaptively selected to reflect the intrinsic goals of the task and calibrated to specific resource bottlenecks. Building on these principles, we categorize the aggregation logic of $\Phi_S$ based on the nature of task outputs:

- **Accumulative Tasks (e.g., multi-target tracking):** For tasks where performance is linearly additive, the baseline $\Phi_S$ is defined as the sum of independent contributions derived from dividing the total resource budget among independent single-agent instances. If $\Gamma > 1$, it indicates that collaboration gains—such as dynamic division of labor preventing redundant search paths—have successfully offset the coordination costs, yielding a result superior to simple parallel processing.

- **Coverage Tasks (e.g., mathematical problem sets):** For tasks measured by task set pass rate, the baseline $\Phi_S$ aligns with "search breadth under equivalent attempts". Given a total budget of $N$ tokens, $\Phi_S$ is calculated as the union of problems solved by a single agent through independent sampling (e.g., self-consistency with $k$ paths) under that budget. This decouples random gains from "brute-force sampling," ensuring that the $\Gamma$ metric measures the collaborative penetration of complex logic rather than mere probability superposition.

- **Single-Solution Tasks (e.g., code generation, complex planning):** For tasks delivering a single, cohesive output, the baseline $\Phi_S$ represents the "individual limit after deep deliberation." The single agent is granted computational support equivalent to the total MAS budget to execute advanced reasoning strategies (e.g., deeper self-reflection or extended Chain-of-Thought (Wei et al., 2022)). This captures whether collective intelligence truly achieves a non-linear leap over the single-agent capability ceiling via interaction, rather than simply thinking longer.

**B.2. Dimensions of Resource Equivalence**

To isolate genuine collaboration gain, the MAS and the SAS baseline must be calibrated to the same resource consumption level. Depending on the task bottleneck, "Resource Equivalence" refers to one or more of the following dimensions:

- **Model Capability Alignment:** The single-agent baseline must utilize the most capable model base within the MAS. This ensures that the evaluation excludes base-model dividends or illusory gains caused by simply upgrading the underlying foundation model.

- **Resource-Matched Evaluation:** The single agent must exhaust a total inference cost (e.g., token length, sampling steps, or tool-invocation counts) equal to that of the entire MAS. This explicitly measures resource consumption to exclude redundant gains derived merely from resource expansion.

## C. The MAS Factor Library

This appendix provides a literature review and specific implementation for the factor library proposed in section 4. We advocate for a scope that exclusively prioritizes collective-level factors, while deliberately excluding single-agent enhancements—such as memory or tool-use—to ensure a focused diagnostic of collective intelligence.

*Table 3.* Classification of Optimization Factors for MAS: A Factor Library Based on Prior Work. This table summarizes key optimization factors and their associated high-level strategies from recent MAS studies, offering a structured view of how various approaches impact collective behavior and system performance.

| Category | Factor | Reference & Title | Specific Optimization Strategy |
|---|---|---|---|
| **Task Context (External)** | **Task Attributes** | (Liu et al., 2025b) *MM-Agent* | **Sequential Dependency:** Decompose open-ended tasks into four distinct sequential stages. |
| | | (Bigeard et al., 2025; Kim et al., 2025) *Finance Agent* | **High Task Decomposability:** Enable sub-tasks to be executed in parallel through a centralized coordinator. |
| | | (Tang et al., 2025) *Depth-Width Writing (DW$^2$)* | **Task Complexity:** Model tasks by Depth (sequential reasoning steps) and Width (capability diversity). Prioritize MAS for high-depth tasks. |
| | | (Huot et al.) *Agents' Room* | **Task Openness:** Open-ended tasks (e.g., fiction writing) inherently favor MAS through specialized sub-task decomposition. |
| **MAS Construction (Internal)** | **Agent Scale** | (Qian et al., 2025) *MacNet* | **Topological Scaling via DAGs:** Structure massive agent populations into Directed Acyclic Graphs (DAGs). |
| | | (Wang et al., 2025c) *MegaAgent* | **SOP-Free Dynamic Scaling:** Dynamically generate agent populations to match real-time task complexity. |
| | **Agent Diversity** | (Ye et al., 2025) *X-MAS* | **Heterogeneous Model Integration:** Leverage model heterogeneity by assigning specialized LLMs to agent roles based on domain-specific capabilities. |
| | | (Jiang et al., 2023) *LLM-Blender* | **Ensemble-Based Generative Fusion:** Aggregate diverse model outputs via pairwise ranking and generative fusion to synthesize a superior response. |
| | | (Wang et al., 2025b) *MoA* | **Layered Collaborative Refinement:** Orchestrate diverse LLMs in a multi-layer structure where agents iteratively refine responses by attending to collective outputs from the previous layer. |
| | **Organizational Structure** | (Dang et al., 2025) *Puppeteer* | **RL-Driven Dynamic Orchestration:** Dynamically sequence agent interactions via an RL-trained orchestrator that adapts to evolving task states without static workflows. |
| | | (Chan et al., 2024) *ChatEval* | **Multi-Agent Debate Topology:** Organize agents into a debate structure to synthesize diverse perspectives through autonomous multi-round argumentation. |

| Category | Factor | Reference & Title | Specific Optimization Strategy |
|---|---|---|---|
| | | (Gong et al., 2024) *MindAgent* | **Centralized Task Coordination:** Dynamically schedule tasks and dispatch instructions to coordinate concurrent multi-agent and human interactions. |
| | | (Zhou et al., 2025b) *MASS* | **Interleaved Design Optimization:** Iteratively optimize local prompts, workflow topologies, and global prompts in a multi-stage search process. |
| | | (Zhang et al., 2025a) *MaAS* | **Supernet-Based Architecture Search:** Dynamically sample query-dependent architectures from a probabilistic supernet. |
| | **Communication Mechanism** | (Zhang et al., 2025b) *AgentPrune* | **One-Shot Graph Pruning:** Use one-shot pruning on the spatial-temporal message-passing graph to eliminate redundant interactions. |
| | | (Yang et al., 2025) *Entropy-Debate* | **Entropy Compression & Adversarial Debate:** Enhance communication efficiency via entropy-compression and employ adversarial debates and voting mechanisms for cross-verification. |
| | | (Chen et al., 2025) *Optima* | **Token-Efficient Training:** Leverage SFT/DPO with cost-aware rewards to balance task performance against token consumption in agent dialogues. |
| | | (Zou et al., 2025) *LatentMAS* | **Training-Free Latent Collaboration:** Establish a shared latent working memory to transfer latent thoughts (last-layer hidden embeddings) directly. |

## D. Validation of the Factor Attribution Paradigm

This appendix validates the factor attribution paradigm proposed in the main text, specifically demonstrating the utility of collaboration gain ($\Gamma$) as a diagnostic probe. As an exploratory pre-experiment to assess the feasibility of $\Gamma$ as a probe for control-level factors, we systematically varied specific factors in the factor library (organization structure, agent diversity, agent scale) and examined the corresponding changes in $\Gamma$. Our observations in this specific setting suggest that this metric can help distinguish between genuine collaboration gain ($\Gamma > 1$) and resource redundancy ($\Gamma < 1$).

### D.1. Validation I: Diagnostic Sensitivity of $\Gamma$ to Factors

To validate the efficacy of $\Gamma$ as an effective probe for control-level factors, a preliminary sensitivity analysis was first conducted, with a specific focus on organizational structure. We measured $\Gamma$ across three distinct datasets (SRDD (Qian et al., 2024), IdeaBench (Guo et al., 2025), and CommonGen(Lin et al., 2020)).



*Figure 6.* Visual Representation of Agent Topological Structures. The upper row illustrates the chain topology, characterized by a sequential, unidirectional data flow among five distinct agents. The lower row depicts a tree topology (specifically a branching structure), where the workflow diverges after the third agent.

We held agent scale (5 Agents) and model backbone (Qwen3-30B-A3B-Instruct-2507) constant while varying only the organizational-structure configuration between a sequential chain and a branching tree. For each benchmark, we randomly sampled 30 distinct tasks to evaluate the system's performance, reporting the averaged metrics.

For baseline construction, the SAS serves as a resource-equivalent reference, enhanced with extended CoT within the same total token budget as the MAS. This ensures that $\Gamma > 1$ identifies genuine structural dividends rather than advantages derived from allowing a single agent to "think longer" under the same resource allocation.
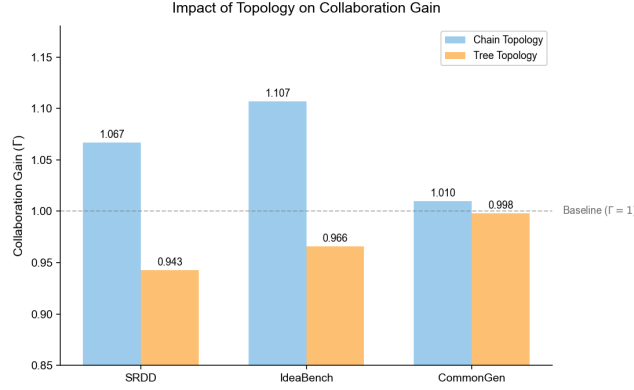


*Figure 7.* Impact of Topological Structure on Collaboration Gain ($\Gamma$)

The results in Figure 7 suggest the metric's potential sensitivity in distinguishing topological effectiveness within the tested domains. Specifically, observations of $\Gamma > 1$ in the chain topology suggest that this structure facilitated information accumulation in these tasks; however, the modest magnitude of these gains indicates that this topology is not a universal remedy, as it remains vulnerable to collapse under naive scaling strategies (detailed in the D.2.4 analysis). In contrast, the metric reflected the performance penalty in tree topologies ($\Gamma < 1$) for these specific scenarios, indicating where the overhead of coordinating branching paths might negate the potential benefits of parallelization.

In this specific experimental setting, the metric $\Gamma$ accurately reflected the performance disparity between the two topologies. It captured the performance drop observed in the tree topology ($\Gamma < 1$) and the gain in the chain topology ($\Gamma > 1$). Within the context of chain and tree topologies evaluated here, this demonstrates the potential of $\Gamma$ as a diagnostic signal to detect effective versus reduced collaboration outcomes within a given context, serving as a basis for the subsequent factor analysis.

### D.2. Validation II: Factor Attribution Mechanism: A detailed case study

Having established the general sensitivity of the collaboration gain metric ($\Gamma$), we present a detailed case study of the "Navigation Buddy" task from the SRDD dataset (Qian et al., 2024).Our objective is not to propose a state-of-the-art multi-agent framework, but to utilize this case as a controlled environment to illustrate how the collaboration gain metric might assist in factor attribution and inform the evolutionary construction of MAS. We analyze three phases through step-wise factor modification.

#### D.2.1. EXPERIMENTAL SETUP AND CONTROLS

**Experimental Control and Budget.** To eliminate computational volume as a confounding factor, we enforce a strict token budget ($B_{max} = 20,000$) for both systems. The MAS distributes this budget across all agents and interaction rounds, while the SAS is allocated the full equivalent budget. Crucially, the SAS is explicitly enhanced with CoT strategies to ensure a "strong baseline," ensuring that performance gaps are attributed to architectural intelligence rather than computational quantity.

**Evaluation Metrics.** We adopt objective metrics from ChatDev (Qian et al., 2024), focusing on completeness, executability, and consistency. We define software quality ($Q$) as the product of these three dimensions ($Q$ = completeness $\times$ executability $\times$ consistency). This composite metric is utilized to calculate the collaboration gain ($\Gamma$), providing a unified standard for assessing software reliability and adherence to requirements.

**Model Consistency.** To ensure fairness, all employed models belong to the same 30B parameter class. We utilize

*Qwen3-30B-A3B-Instruct-2507* as the foundation model for the SAS and all MAS roles in phase I.

D.2.2. PHASE I: PROBING THE "AGENT DIVERSITY(ROLE DIVERSITY)" FACTOR: $\Gamma > 1$

**Factor Configuration Table**

*Table 4.* MAS Construction (Internal) based on section 4 Framework

| Factor | Configuration |
|---|---|
| Organization Structure | Chain |
| Communication Mechanism | Natural Language |
| Agent Diversity | Role Diversity (Heterogeneous Roles: CEO for Strategy, CTO for Architecture, Programmer for Coding). |
| Agent Scale | 3 Agents (1 CEO, 1 CTO, 1 Programmer) |

**Qualitative Analysis**

We conducted a qualitative analysis of the multi-agent collaboration by examining the logged interaction records. Due to space constraints, we present only selected outputs that highlight the collaborative advantage of the MAS. Each retained output is annotated to indicate its strategic significance.

---

**User - CEO: Core CEO Output (Selected)**

Define operational mechanics, workflow, and quality constraints for Navigation Buddy, a real-time navigation system.
Ensure developers implement robust, shortcut-free code; maintain production reliability and security.
Analysis: Centralizes strategy and standards, reducing risk of fragmented or inconsistent implementation.

---

**CEO - CTO: Core CTO Output (Selected)**

*Listing 1.* Project File Structure

```
1  {
2    "file_structure": [
3      "navigation_buddy/__init__.py",
4      "navigation_buddy/config.py",
5      "navigation_buddy/core/__init__.py",
6      "navigation_buddy/core/position.py",
7      "navigation_buddy/core/route_planner.py",
8      "navigation_buddy/core/traffic_data.py",
9      "navigation_buddy/integrations/gps_provider.py",
10     "navigation_buddy/integrations/map_provider.py",
11     "navigation_buddy/services/__init__.py",
12     "navigation_buddy/services/directions_service.py",
13     "navigation_buddy/services/realtime_updates.py",
14     "navigation_buddy/utils/__init__.py",
15     "navigation_buddy/utils/logging.py",
16     "navigation_buddy/utils/time_utils.py",
17     "navigation_buddy/main.py"
18   ]
19  }
```

Analysis: Transform an open-ended problem into a constrained design space.

---

**CTO - Programmer: Core Programmer Output (Selected)**

*Listing 2.* Input Validation & Logic Orchestration

```
1      def get_directions(self, start, end):
2          # CONSTRAINT CHECK: Input Validation
3          if not self.location_service.validate_location(start):
4              raise ValueError("Invalid start location")
```

```
5        if not self.location_service.validate_location(end):
6            raise ValueError("Invalid end location")
7        # Dynamic Resolution
8        if start == "Current Location":
9            start = self.location_service.get_current_location()
10       # Core Business Logic
11       route = self.routing_engine.compute_route(start, end)
12       return {
13           "start": start,
14           "end": end,
15           "route": route,
16           "steps": self._format_steps(route)
17       }
```

Analysis: This reflects CTO-defined architectural contracts, reducing coupling and error propagation paths.

### Single Agent System (SAS): Evidence Output (Selected)

*Listing 3.* Stubbed & Linear Logic

```
1  def get_current_location():
2      # Returns constant instead of implementation
3      return "XYZ"
4
5  def navigate_to(location):
6      # Linear script ignoring required event loops
7      route = calculate_route(current, location)
8      print(route)
```

We compare the generated code from the multi-agent system against a resource-equivalent single-agent baseline. Despite CoT prompting, the SAS showed consistency issues in long-generation tasks, forgetting planning-phase constraints and using placeholders as the context window filled.

### Quantitative Evaluation

We utilize the composite $Q$ score as the primary basis for calculating the collaboration gain ($\Gamma$) proposed in Section 3. We compare the MAS against the resource-equivalent SAS. In this selected case study, both systems can successful execute; however, disparities persisted in completeness and consistency.

*Table 5.* Quantitative comparison of code generation performance

| Setting | Comp. | Exec. | Cons. | $Q$ | $\Gamma$ |
|---|---|---|---|---|---|
| Single Agent (Baseline) | 0.35 | 1 | 0.76 | 0.27 | / |
| Multi-Agent System | **0.42** | 1 | **0.81** | **0.34** | **1.26** |

The MAS achieved a collaboration gain ($\Gamma > 1$), demonstrating better performance than the resource-equivalent single-agent baseline.

D.2.3. PHASE II: OPTIMIZING VIA "AGENT DIVERSITY (MODEL)" FACTOR ($\Gamma \uparrow$)

Building upon the findings from the phase I experiment, where functional orthogonality via role diversity resulted in a positive collaboration gain, we posit that the *agent diversity* factor can be further optimized through Model Heterogeneity.

### Factor Configuration Update

In phase II, we introduce model heterogeneity by switching the programmer agent to *Qwen3-Coder-30B-A3B-Instruct*. This setup maintains a constant parameter scale while isolating the impact of specialized coding expertise.

### Qualitative Analysis

Rather than increasing surface complexity, the change manifested as the emergence of new engineering capability classes that were absent in the homogeneous configuration.

```
1   class RealTimeUpdates:
2       """Manages real-time navigation updates"""
3       def start_monitoring(self):
4           if not self.is_monitoring:
5               self.is_monitoring = True
6               # Autonomous decision: Threaded execution for non-blocking I/O
7               self.monitoring_thread = threading.Thread(
8                   target=self._monitor_loop
9               )
10              self.monitoring_thread.daemon = True
11              self.monitoring_thread.start()
12      def _monitor_loop(self):
13          while self.is_monitoring:
14              try:
15                  self._process_updates()
16                  time.sleep(30)
17              except Exception as e:
18                  print(f"Error in monitoring loop: {e}")
```

Analysis: An autonomous concurrency model that was not specified at the architectural level.

The output demonstrates enhanced structural completeness, exemplified by features such as threaded concurrency and observability decorators.

**Quantitative Evaluation**

To rigorously quantify the effect of the enhanced agent diversity factor, we revisit the collaboration gain ($\Gamma$) metric based on the composite software quality score ($Q$). Table 6 reports the comparative results. The data exhibits a stepwise improvement in system performance as the agent diversity factor is progressively enriched.

*Table 6.* Quantitative comparison of code generation performance

| Setting | Comp. | Exec. | Cons. | $Q$ | $\Gamma$ |
|---|---|---|---|---|---|
| Single Agent (Baseline) | 0.35 | 1 | 0.76 | 0.27 | / |
| Multi-Agent System | **0.60** | **1** | **0.83** | **0.50** | **1.85** |

The increase in $\Gamma$ corresponds to the introduction of model heterogeneity. In this case, this factor appeared to facilitate a capability transition from "functional correctness" to "engineering abstraction" (e.g., the autonomous implementation of concurrency patterns), suggesting that strategic model-role alignment can contribute to system performance.

D.2.4. PHASE III: BOUNDARY TESTING THE "AGENT SCALE" FACTOR ($\Gamma < 1$)

Following the validation of role heterogeneity (phase I) and model heterogeneity (phase II), we examine a deeper scalability question: Does change the factor of agent scale produce better collaboration gain? To find and verify this factor hypothesis, we assume that expanding the implementation layer—by stacking multiple Programmer agents in sequence—should improve output quality through iterative refinement, redundancy, and error correction.

**Factor Configuration Update**

To isolate the effect of agent scale, the experimental configuration was updated to a sequential implementation chain expanded to 5 agents. To ensure valid collaboration within this sequential topology, Programmer 2 and 3 were explicitly instructed to ingest the predecessor's code and perform incremental optimization (e.g., bug fixing and logic hardening).

**Qualitative Analysis**

Contrary to expectations, the qualitative analysis indicates a decline in the multi-agent system's performance in this sequential setup. In this sequential setup, we observed a structural regression in the multi-agent system's output. Following the initial complexity construction by Programmer 1, the subsequent agent did not maintain the multi-file architecture. Instead, the final output was simplified into a single, self-contained script, prioritizing executability over the original structural integrity.

```python
# main.py (Programmer 2)
class NavigationBuddy:
    def __init__(self):
        self.current_location = (0.0, 0.0)
        self.route = [
            "Head east for 2.5 km...",
            "Turn right onto Main Street."
        ]
def main():
    lat_input = float(input("Latitude: "))
```

**Analysis:** The agent fails to preserve the multi-file architectural model. To maintain executability, it collapses the system into a monolithic script.

## Quantitative Evaluation

*Table 7.* Quantitative comparison of code generation performance

| Setting | Comp. | Exec. | Cons. | $Q$ | $\Gamma$ |
|---|---|---|---|---|---|
| Single Agent (Baseline) | 0.35 | 1 | 0.76 | 0.27 | / |
| Multi-Agent System | **0.23** | 1 | **0.74** | **0.17** | **0.63** |

Collaboration gain collapses to $\Gamma < 1$, indicating that, within this specific experimental design, naive scale expansion did not yield cumulative intelligence.

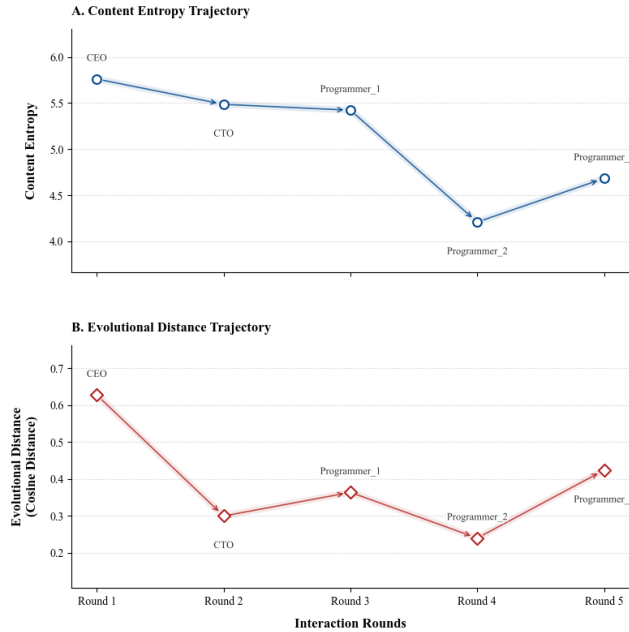## Process Dynamics: Diagnostic Analysis via Information Flow



*Figure 8.* Evolutionary Dynamics of Information Level

To interpret the collapse in collaboration gain ($\Gamma < 1$), we conducted a diagnostic analysis combining the detailed interaction logs with the information flow metrics shown in Figure 8.

In the initial phase, the system exhibits a stable transmission pattern: content entropy shows a gradual decline, while evolutionary distance increases (particularly from CTO to Programmer_1). This trajectory corresponds to a healthy "semantic expansion," where abstract requirements are effectively translated into concrete code implementations without

losing focus. However, a critical anomaly emerges during the transition from Programmer_1 to Programmer_2. Here, we observe a precipitous drop in both content entropy and evolutionary distance. Cross-referencing this with the interaction logs suggests a "fracture" in the information flow caused by context overload. The extensive multi-file architecture generated by Programmer_1 appears to have exceeded the effective processing window of the subsequent agent, causing Programmer_2 to discard structural details and focus on a narrowed subset of the data. Finally, the subsequent rise in metrics at Programmer_3 indicates that the agent generated new content, but this expansion was grounded in the already broken, incomplete context provided by Programmer_2. This fluctuation implies that the observed loss of collaboration gain likely stems from a bottleneck in sustaining complex context across extended sequential chains.

### D.3. Conclusion

The analysis presented above, preliminarily illustrates that collaboration gain serves as a principled diagnostic signal for factor attribution within MAS. By controlling for resource consumption, $\Gamma$ enables the identification of specific factors—such as agent diversity—that either facilitate or impede system performance. Combined with information-level metrics, $\Gamma$ provides feedback for detecting internal bottlenecks, offering guidance for evaluating and refining MAS beyond empirical observation.