# Understanding Agent Scaling in LLM-Based Multi-Agent Systems via Diversity

**Yingxuan Yang** [1]  **Chengrui Qu** [2]  **Muning Wen** [1]  **Laixi Shi** [3]  **Ying Wen** [1]  **Weinan Zhang** [1]  **Adam Wierman** [2]
**Shangding Gu\*** [4]

## Abstract

LLM-based multi-agent systems (MAS) have emerged as a promising approach to tackle complex tasks that are difficult for individual LLMs. A natural strategy is to scale performance by increasing the number of agents; however, we find that such scaling exhibits strong diminishing returns in homogeneous settings, while introducing heterogeneity (e.g., different models, prompts, or tools) continues to yield substantial gains. This raises a fundamental question: *what limits scaling, and why does diversity help?* We present an information-theoretic framework showing that MAS performance is bounded by the intrinsic task uncertainty, not by agent count. We derive architecture-agnostic bounds demonstrating that improvements depend on how many *effective channels* the system accesses. Homogeneous agents saturate early because their outputs are strongly correlated, whereas heterogeneous agents contribute complementary evidence. We further introduce $K^*$, an *effective channel count* that quantifies the number of effective channels without ground-truth labels. Empirically, we show that heterogeneous configurations consistently outperform homogeneous scaling: 2 diverse agents can match or exceed the performance of 16 homogeneous agents. Our results provide principled guidelines for building efficient and robust MAS through diversity-aware design. Code and Dataset are available at the link: https://github.com/SafeRL-Lab/Agent-Scaling.

## 1. Introduction

Large language models (LLMs) have achieved remarkable performance across diverse tasks, including reasoning, cod-
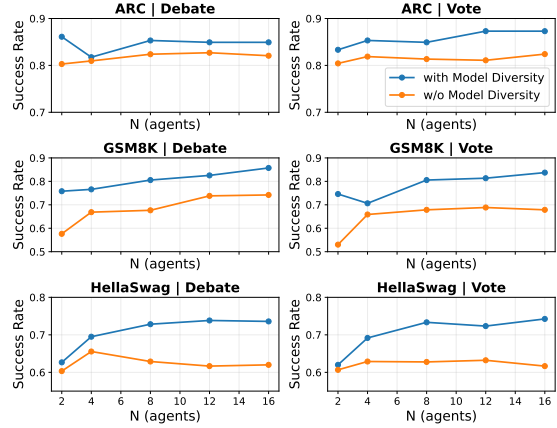


Figure 1: Effect of model diversity. We compare a mixture of three LLMs (Qwen-2.5-7B, Llama-3.1-8B, Mistral-7B) with the average of independent single-LLM runs.

ing, and open-domain question answering (Wei et al., 2022; Achiam et al., 2023). However, individual LLMs still struggle with complex problems that require multi-step reasoning, diverse perspectives, or complementary expertise (Huang et al., 2023). To address these limitations, LLM-based multi-agent systems (MAS) have emerged as a promising paradigm. By orchestrating multiple LLM agents through communication, coordination, or aggregation mechanisms, MAS can tackle challenges that are difficult for single models (Wu et al., 2024; Hong et al., 2024; Du et al., 2023). Recent studies have demonstrated that multi-agent collaboration can yield substantial improvements over single-agent baselines on tasks ranging from software engineering (Qian et al., 2024) to scientific reasoning (Guo et al., 2024).

Given the effectiveness of multi-agent systems, a natural question arises: *can we improve MAS performance simply by scaling the number of agents?* Intuitively, one might expect ensemble-style gains from aggregating more agent outputs (Li et al., 2024a; Wang et al., 2023). However, recent work (Kim et al., 2025) and our experiments reveal a more nuanced picture. As shown in Figure 2, scaling homogeneous agents (identical models, prompts, and configurations) exhibits strong diminishing returns: accuracy improves at small agent counts, but the marginal gain per additional agent, rapidly collapses toward zero. This suggests that simply adding more homogeneous agents (or allocating
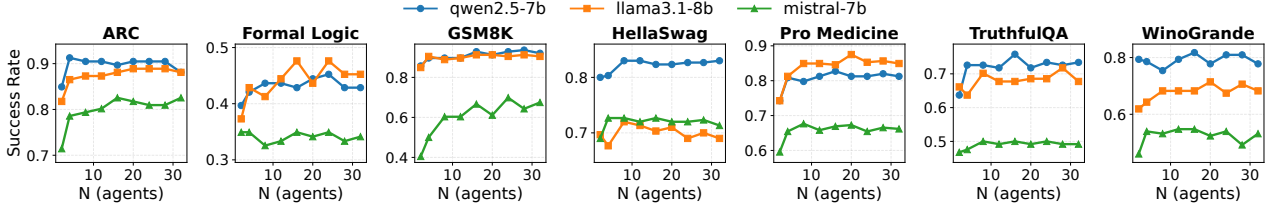
---

[1]Shanghai Jiao Tong University, [2]California Institute of Technology, [3] Johns Hopkins University, [4]UC Berkeley. * Correspondence to: <shangding.gu@berkeley.edu>.

*Preprint*

Figure 2: Scaling behavior of homogeneous multi-agent voting. Success rate versus agent count N on seven tasks for three base models. Performance improves with N but saturates, indicating clear diminishing marginal gains at larger agent counts.

more test-time compute) does not reliably introduce new *usable evidence* into the system, but may instead produce increasingly redundant trajectories.

In contrast, our experiments (Figure 1) show that introducing *diversity* yields sustained performance improvements. Here, diversity broadly refers to heterogeneity in agent configurations, such as backbone models, prompts or personas, and tool access, which empirically leads to more complementary, rather than redundant, information being introduced into the system. As a result, diverse systems can outperform homogeneous ones even with substantially fewer agent calls (Wang et al., 2024a; Zhang et al., 2024; Qian et al., 2025). Motivated by these observations, we ask: **what fundamentally limits scaling, and why does diversity help?**

We hypothesize that the primary bottleneck arises from *correlation among agent outputs*. Higher correlation induces greater redundancy, reducing the number of effective channels and leading to performance saturation (Chen et al., 2024; Choi et al., 2025). This intuition is illustrated in Figure 3, where heterogeneous agents provide complementary coverage and better information processing diversity compared to homogeneous systems (Yuen et al., 2025; Tang et al., 2025). To formalize this intuition, we develop an information-theoretic framework that characterizes MAS performance in terms of *effective channels*, the number of independent, non-redundant reasoning paths present in agent outputs, rather than raw agent count. For example, two agents that reason in nearly identical ways contribute only one effective channel, whereas two agents that follow genuinely different reasoning paths contribute two. Our analysis reveals that performance is bounded by intrinsic task uncertainty, and improvements depend on how many effective channels the system accesses.

Based on this framework, we introduce $K^*$, a label-free metric that quantifies effective channels without requiring ground-truth labels. Empirically, we demonstrate that heterogeneous configurations consistently outperform homogeneous scaling: with only **2 diverse agents**, we match or exceed the performance of **16 homogeneous agents**, achieving significant improvement across seven benchmarks.

Although diminishing returns in scaling have been observed empirically (Wang et al., 2024b; Kim et al., 2025), a unified theoretical framework explaining why and when this phenomenon occurs across different MAS workflows, such as voting (Wang et al., 2023), debate (Du et al., 2023; Khan et al., 2024), and centralized orchestration (Hong et al., 2024), remains lacking. Existing studies offer limited theoretical insight into how evidence accumulation is affected by agent redundancy as the system scales. To address this gap, we provide a unified information-theoretic explanation for diminishing returns in LLM-based MAS. Our contributions are summarized as follows:

- We derive architecture-independent performance bounds, demonstrating that MAS effectiveness is constrained by the intrinsic task uncertainty $H(Y|X)$, and that improvements arise from increasing the number of effective channels rather than scaling the agent count.

- We analyze representative MAS paradigms (vote, debate), showing that homogeneous configurations quickly saturate due to highly correlated evidence, whereas heterogeneity effectively reduces redundancy and expands the system's capacity for effective channels.

- We introduce $K^*$, an effective channel count that quantifies the number of non-redundant information sources in agent outputs. We empirically validate that $K^*$ tracks performance and provides principled guidelines for diversity-driven MAS design.

## 2. Related Works

**Information-Theoretic Analysis of LLM Reasoning.** Recent work has begun applying information theory to understand LLM behavior. Ton et al. (2024) quantify information gain at each chain-of-thought step, showing that effective reasoning requires each step to contribute new information. Gan et al. (2025) analyze cascading failures through information loss accumulation: when $I(t_\ell; r_\ell)$ grows superlinearly, conditional entropy increases rather than decreases. In multi-agent settings, Riedl (2025) use Time-Delayed Mutual Information to detect coordination vs. mere information sharing, and Chang (2025) track when agent dialogues converge vs. maintain distributed information. However, these works focus on *characterizing* information flow patterns rather than *explaining* why diversity constraints emerge or

*deriving* performance bounds. In contrast to prior work that primarily measures or characterizes information flow in LLM reasoning, we use information theory to *explain* diminishing returns via formal limits.

**LLM-based Multi-Agent Systems.** LLM–based MAS instantiate multiple interacting LLM agents to perform compound inference through communication, coordination, or aggregation mechanisms (Xi et al., 2023; Wang et al., 2024b; Guo et al., 2024). Existing designs span independent sampling and voting schemes related to self-consistency (Wang et al., 2023), decentralized debate and role-playing frameworks (Du et al., 2023; Khan et al., 2024; Li et al., 2024b;a; 2023), centralized orchestration frameworks such as Auto-Gen (Wu et al., 2024) and MetaGPT (Hong et al., 2024), as well as hybrid, evolving, or self-improving coordination strategies (Dang et al., 2025; Zhao et al., 2025). Cemri et al. (2025) identify systematic failure modes across multi-agent systems. Taken together, these findings suggest that MAS performance is influenced by multiple design factors, among which we specifically focus on the role of agent diversity.

**Empirical Studies of MAS Scaling and Diversity.** It has been shown that naïvely scaling the number of agents yields limited benefits when agent behaviors are homogeneous (Wang et al., 2024b; Chen et al., 2024), across majority voting (Qian et al., 2025), debate (Choi et al., 2025), and more general coordination mechanisms (Kim et al., 2025).

In contrast, a growing body of empirical evidence highlights the central role of diversity in MAS. Zhang et al. (2024) demonstrates that diversity leads to higher success rates in software engineering agents, while Wang et al. (2024a) finds that heterogeneous ensembles outperform homogeneous ones. Wu & Ito (2025) argue that preserving disagreement is preferable to enforcing early consensus. Related work shows that diversity benefits depend on task complexity (Tang et al., 2025) and that persona-based diversification has limitations (Samuel et al., 2024; Taillandier et al., 2025). However, these findings are restricted to specific diversity forms and narrow settings. We provide a unified theoretical and empirical analysis across multiple diversity types.

# 3. Problem Formulation

This section formalizes the notion of information flow in LLM-based multi-agent systems and establishes the theoretical foundations for understanding scaling behavior.

We first define the system setup, then introduce the key quantity that governs MAS performance: *usable evidence*. Finally, we derive upper bounds showing that achievable information gain is determined by agent diversity.

## 3.1. LLM-based Multi-Agent Systems

We begin by formally defining the class of systems we study.

**Definition 3.1** (LLM-based Multi-Agent System). An *LLM-based multi-agent system* consists of $N$ agents, each characterized by a *configuration* that specifies its backbone model, system prompt or persona, decoding strategy, and tool access. Given a task input $X$, the system executes a total of $n$ agent calls through a specified workflow (e.g., parallel voting, sequential debate) and aggregates the outputs to produce a final answer.

**Notation.** We distinguish the *number of agents* $N$ and the *number of agent calls* $n$. In single-round workflows such as majority voting, $n = N$. In multi-round workflows such as debate with $R$ rounds, $n = N \times R$. This distinction is important because our analysis focuses on how much information is extracted, regardless of which agent produces it. Agent configuration types are formally defined in Section 3.3.

## 3.2. Usable Evidence and Information Budget

Consider a task with input $X \in \mathcal{X}$ and ground-truth answer $Y \in \mathcal{Y}$. During inference, the MAS executes $n$ agent calls and produces a dialogue transcript:

$$Z_{1:n} = (Z_1, \ldots, Z_n), \tag{1}$$

where each output $Z_i$ may depend on the input $X$ and all preceding outputs $Z_{<i} = (Z_1, \ldots, Z_{i-1})$.

The central question is: *how much information about the answer $Y$ can the system extract from its agent calls?* We quantify this through the conditional mutual information:

$$\begin{aligned} I_{\text{MAS}}(n) &:= I(Z_{1:n}; Y \mid X) \\ &= H(Y \mid X) - H(Y \mid X, Z_{1:n}). \end{aligned} \tag{2}$$

This quantity, which we call *usable evidence*, measures the reduction in uncertainty about $Y$ achieved by observing the transcript, beyond what is already contained in the input $X$.

To understand how usable evidence accumulates, let $\Delta_i := I(Z_i; Y \mid X, Z_{<i})$ denote the *incremental contribution* of the $i$-th call, i.e., the new information it provides given all previous outputs. By the chain rule for mutual information:

$$I_{\text{MAS}}(n) = \sum_{i=1}^{n} \Delta_i. \tag{3}$$

This decomposition shows that MAS performance depends not on the total number of calls $n$, but on how much *non-redundant* evidence each call contributes. If agents produce highly correlated outputs, the incremental contributions $\Delta_i$ diminish rapidly, leading to saturation. As illustrated in Figure 3, heterogeneous agents provide complementary coverage and better diversity in information processing compared to homogeneous configurations.

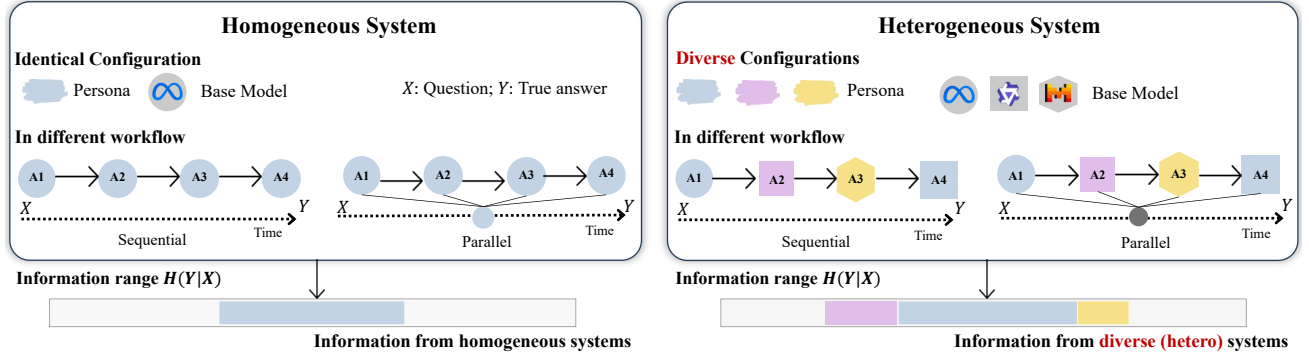The following theorem gives an upper bound on the achievable information:

Figure 3: A comparison between homogeneous and heterogeneous systems. In a homogeneous system, agents with identical configurations result in redundant behavior and limited information coverage. In contrast, heterogeneous agents, through diverse configurations (e.g., varying models or personas), provide complementary coverage and better diversity in the information processed, allowing for more effective problem-solving across different workflows.

**Theorem 3.2** (Finite Information Budget). *For any transcript $Z_{1:n}$,*

$$I_{\mathrm{MAS}}(n) \leq H(Y \mid X). \tag{4}$$

This bound states that no MAS can extract more information about $Y$ than the intrinsic task uncertainty $H(Y|X)$. The practical implication is that scaling benefits plateau once this ceiling is approached, and homogeneous systems may reach saturation much earlier than heterogeneous systems due to redundant evidence.

### 3.3. Agent Configuration Types

Agent diversity is operationalized through variations in backbone model, system prompt/persona, decoding strategy, and tool access. We index these choices by *configuration types*.

**Definition 3.3** (Agent Configuration Type). Each call $i \in \{1, \ldots, n\}$ is associated with a type $b(i) \in \mathcal{B}$. For each $b \in \mathcal{B}$, define the number of calls

$$m_b := \big|\{i \in \{1, \ldots, n\} : b(i) = b\}\big|, \quad \sum_{b \in \mathcal{B}} m_b = n. \tag{5}$$

### 3.4. Type-Dependent Ceilings Across workflows

We next state representative upper bounds showing that, across common MAS workflows, achievable information gain is controlled by the multiset of configuration types. All formal assumptions and proofs are deferred to Appendix A.

**Parallel interaction.** Let $I_b := I(Z^{(b)}; Y \mid X)$ denote the single-call information of type $b$ (see Appendix A.3 for the full derivation). Under a standard conditional-independence model for parallel sampling (Assumption A.6),

$$I_{\mathrm{MAS}}^{\mathrm{parallel}}(n) \leq H(Y \mid X) \wedge \sum_{b \in \mathcal{B}} m_b I_b. \tag{6}$$

**Sequential interaction.** Define the maximal per-step contribution of type $b$ by

$$I_b^{\mathrm{max}} := \sup_{z_{<i}} I(Z_i; Y \mid X, Z_{<i} = z_{<i}, b(i) = b). \tag{7}$$

Any sequential MAS satisfies

$$I_{\mathrm{MAS}}^{\mathrm{seq}}(n) \leq H(Y \mid X) \wedge \sum_{b \in \mathcal{B}} m_b I_b^{\mathrm{max}}. \tag{8}$$

Debate is a special case of sequential interaction and inherits the same ceiling.

**From ceilings to compute.** The bounds above depend on *structural properties* of the MAS, which types are instantiated and how they are composed, rather than on the raw call count $n$. Since these upper bounds do not depend on $n$, the raw call count is not the right quantity for characterizing MAS performance limits. This motivates us to identify a new quantity, the *effective channel count*, that more directly governs how much usable evidence a MAS can extract.

## 4. Why Diversity Matters

Section 3 establishes that MAS performance is bounded by intrinsic task uncertainty and that the upper bounds depend on configuration types rather than the raw call count $n$. This raises a natural question: what quantity *does* govern how much information a MAS actually extracts? In this section, we introduce the *effective channel count* $K$ to answer this question. We then show why homogeneous scaling often saturates (because $K$ stops growing), while heterogeneous designs can keep improving by increasing the amount of *complementary* evidence (larger $K$, and/or a higher evidence-coverage rate $\alpha$), which leads to a characteristic fast-then-slow gain curve.

## 4.1. Effective Channels: From Compute to Usable Evidence

An *effective channel* represents one independent source of task-relevant information in the MAS transcript. Intuitively, if two agents produce nearly identical reasoning, they contribute only one effective channel despite consuming two agent calls; if they reason along genuinely different paths, they contribute two. The *effective channel count* $K$ thus captures how many non-redundant information sources the system has, as opposed to the raw number of calls $n$. To formalize this, we introduce two interrelated concepts: the *complementarity rate* $\alpha$ and the *effective channel count* $K$.

**Definition 4.1** (Complementarity Rate). The *complementarity rate* $\alpha \in (0, 1)$ quantifies the probability that a new effective channel uncovers previously missing task-relevant evidence. Formally, $\alpha$ governs the rate at which additional channels reduce residual uncertainty about $Y$.

Intuitively, $\alpha$ reflects how "complementary" the information from different channels is. A high $\alpha$ indicates that each new channel is likely to provide fresh evidence, while a low $\alpha$ suggests substantial overlap with existing information.

**Definition 4.2** (Effective Channel Representation). An *effective channel representation* of the transcript $Z_{1:n}$ is a collection of $K$ channels:

$$\tilde{Z}_{1:K} = (\tilde{Z}^{(1)}, \ldots, \tilde{Z}^{(K)}) \quad \text{s.t.} \quad \tilde{Z}_{1:K} = \phi(Z_{1:n}), \quad (9)$$

for some (possibly lossy) aggregation map $\phi$, where $K$ is the *effective channel count*, representing the number of non-redundant information sources in the agent outputs.

$K$ and $\alpha$ are coupled: increasing diversity (larger $K$) is beneficial only if the new channels provide complementary evidence (captured by $\alpha$). The product $\alpha K$ thus serves as the fundamental quantity governing information recovery, as formalized in Theorem A.15.

Since $\tilde{Z}_{1:K}$ is a function of $Z_{1:n}$, the data processing inequality implies:

$$I(Z_{1:n}; Y \mid X) \geq I(\tilde{Z}_{1:K}; Y \mid X). \quad (10)$$

**Connecting $K$ and $\alpha$ to recoverable information.** To formalize the relationship between effective channels and information recovery, we introduce in Appendix A a minimal evidence-coverage model (Assumptions A.12 and A.13). Under this model, the information recovered from $K$ effective channels with complementarity rate $\alpha$ approaches the intrinsic task uncertainty at a geometric rate:

**Theorem 4.3** (Geometric Contraction with Effective Channels). *Under Assumptions A.12 and A.13, the residual un-*

*certainty after observing $K$ effective channels satisfies*

$$\boxed{\begin{aligned} & H(Y \mid X) - \mathbb{E}\Big[I(\tilde{Z}_{1:K}; Y \mid X)\Big] \\ & \leq (1-\alpha)^K H(Y \mid X) \leq e^{-\alpha K} H(Y \mid X). \end{aligned}} \quad (11)$$

*Equivalently, the* normalized residual *satisfies* $\mathbb{E}[H(Y \mid X, \tilde{Z}_{1:K})]/H(Y \mid X) \leq (1-\alpha)^K \leq e^{-\alpha K}$.

## 4.2. $K$ as the State Variable of MAS Scaling

The central question in MAS scaling is not whether $n$ increases, but whether $n$ induces growth in the *effective channel count* $K(n)$. This follows from Section 3: ceilings are fixed by intrinsic uncertainty $H(Y \mid X)$ and structural design (Section 3.4), while achievability improves with the number of non-redundant channels (Section 4.1).

**A direct heterog–homog advantage bound.** Consider two designs under the same compute budget $n$. Let $(K_{\text{homog}}, \alpha_{\text{homog}})$ and $(K_{\text{heterog}}, \alpha_{\text{heterog}})$ denote their effective channel counts and coverage rates in the evidence-coverage model. By Theorem A.15, each design admits a lower bound on recoverable information:

**Corollary 4.4** (Heterogeneity Advantage). *Under Assumptions A.12 and A.13, the lower bounds on recoverable information for the two designs are:*

$$\mathbb{E}\big[I_{\text{heterog}}\big] \geq H(Y \mid X)\big(1 - e^{-\alpha_{\text{heterog}}K_{\text{heterog}}}\big), \quad (12)$$

$$\mathbb{E}\big[I_{\text{homog}}\big] \geq H(Y \mid X)\big(1 - e^{-\alpha_{\text{homog}}K_{\text{homog}}}\big). \quad (13)$$

*When $\alpha_{\text{heterog}}K_{\text{heterog}} > \alpha_{\text{homog}}K_{\text{homog}}$, the heterogeneous design enjoys a strictly higher information-recovery guarantee: its lower bound on recoverable information, $H(Y \mid X)(1 - e^{-\alpha_{\text{heterog}}K_{\text{heterog}}})$, exceeds the corresponding homogeneous guarantee $H(Y \mid X)(1 - e^{-\alpha_{\text{homog}}K_{\text{homog}}})$.*

This is consistent with our empirical findings: as shown in Figure 1 and Table 1, heterogeneous configurations consistently recover more task-relevant information than homogeneous ones under matched compute. The corollary formalizes the intuition that heterogeneity helps by increasing $\alpha K$ through more non-redundant channels or higher complementarity.

**Fast-then-slow scaling: the $1 - e^{-\alpha K}$ shape.** Corollary A.16 implies that recoverable information grows at least as

$$\mathbb{E}\big[I(\tilde{Z}_{1:K}; Y \mid X)\big] \geq H(Y \mid X)\big(1 - e^{-\alpha K}\big). \quad (14)$$

The shape of (14) directly predicts diminishing returns: the marginal gain from one additional effective channel satisfies

$$\big(1 - e^{-\alpha(K+1)}\big) - \big(1 - e^{-\alpha K}\big) = (1 - e^{-\alpha})e^{-\alpha K}, \quad (15)$$

Table 1: Effect of persona diversity. $\Delta$ denotes improvement from heterogeneity. All agents share the same base model pool (Qwen-2.5-7B, Llama-3.1-8B, and Mistral-7B); only persona assignments differ between Homog and Heterog.

| Dataset | Single Agent | N | Vote | | | Debate | | | Dataset | Single Agent | N | Vote | | | Debate | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | Homog | Heterog | $\Delta$ | Homog | Heterog | $\Delta$ | | | | Homog | Heterog | $\Delta$ | Homog | Heterog | $\Delta$ |
| GSM8K | 50.8 | 2 | 86.5 | 87.3 | +0.8 | 76.2 | 75.4 | -0.8 | ARC | 77.8 | 2 | 78.6 | 81.8 | +3.2 | 84.9 | 87.3 | +2.4 |
| | | 4 | 84.9 | 88.1 | +3.2 | 73.8 | 79.4 | +5.6 | | | 4 | 79.4 | 85.7 | +6.3 | 79.4 | 84.1 | +4.8 |
| | | 8 | 90.5 | 93.7 | +3.2 | 75.4 | 85.7 | +10.3 | | | 8 | 84.1 | 86.5 | +2.4 | 84.9 | 85.7 | +0.8 |
| | | 12 | 86.5 | 90.5 | +4.0 | 77.8 | 87.3 | +9.5 | | | 12 | 85.7 | 89.7 | +4.0 | 82.5 | 87.3 | +4.8 |
| | | 16 | 89.7 | 92.1 | +2.4 | 83.3 | 88.1 | +4.8 | | | 16 | 84.9 | 88.9 | +4.0 | 84.9 | 84.9 | 0.0 |
| Formal Logic | 32.0 | 2 | 45.2 | 48.4 | +3.2 | 34.1 | 38.9 | +4.8 | Truthful QA | 71.8 | 2 | 74.2 | 77.4 | +3.2 | 71.0 | 77.4 | +6.4 |
| | | 4 | 47.6 | 52.4 | +4.8 | 42.9 | 53.2 | +10.3 | | | 4 | 75.0 | 75.8 | +0.8 | 71.8 | 79.8 | +8.0 |
| | | 8 | 47.6 | 55.6 | +7.9 | 49.2 | 53.2 | +4.0 | | | 8 | 76.6 | 79.0 | +2.4 | 76.6 | 78.2 | +1.6 |
| | | 12 | 48.4 | 57.9 | +9.5 | 48.4 | 54.8 | +6.4 | | | 12 | 75.0 | 79.0 | +4.0 | 73.4 | 79.8 | +6.4 |
| | | 16 | 50.0 | 54.0 | +4.0 | 43.6 | 51.6 | +8.0 | | | 16 | 78.2 | 81.5 | +3.3 | 75.0 | 84.7 | +9.7 |
| HellaSwag | 66.1 | 2 | 62.3 | 73.7 | +11.4 | 50.3 | 75.0 | +24.7 | Wino grande | 57.1 | 2 | 51.6 | 60.3 | +8.7 | 58.7 | 50.0 | -8.7 |
| | | 4 | 68.7 | 75.3 | +6.6 | 66.0 | 73.0 | +7.0 | | | 4 | 54.0 | 69.1 | +15.1 | 53.2 | 62.7 | +9.5 |
| | | 8 | 70.0 | 79.0 | +9.0 | 69.7 | 76.0 | +6.3 | | | 8 | 57.9 | 69.1 | +11.2 | 61.9 | 69.1 | +7.2 |
| | | 12 | 72.3 | 79.0 | +6.7 | 69.3 | 78.3 | +9.0 | | | 12 | 58.7 | 70.6 | +11.9 | 62.7 | 70.6 | +7.9 |
| | | 16 | 72.0 | 79.9 | +7.9 | 70.3 | 76.4 | +6.1 | | | 16 | 60.3 | 69.8 | +9.5 | 57.9 | 64.3 | +6.4 |
| Pro Medicine | 68.6 | 2 | 78.3 | 78.7 | +0.4 | 76.8 | 71.3 | -5.5 | Average | 60.6 | 2 | 68.1 | 72.5 | **+4.4** | 64.6 | 67.9 | **+3.3** |
| | | 4 | 80.5 | 81.6 | +1.1 | 76.8 | 76.5 | -0.3 | | | 4 | 69.9 | 76.1 | **+6.2** | 66.3 | 72.7 | **+6.4** |
| | | 8 | 81.3 | 83.5 | +2.2 | 81.6 | 82.7 | +1.1 | | | 8 | 72.6 | 79.6 | **+7.0** | 71.3 | 75.8 | **+4.5** |
| | | 12 | 80.2 | 82.7 | +2.5 | 81.3 | 83.8 | +2.5 | | | 12 | 72.4 | 81.0 | **+8.6** | 70.8 | 77.4 | **+6.6** |
| | | 16 | 80.5 | 81.8 | +1.3 | 80.5 | 83.3 | +2.8 | | | 16 | 73.6 | 81.1 | **+7.5** | 70.8 | 76.2 | **+5.4** |

which is largest at small $K$ and decays exponentially thereafter. This yields a clean explanation for the empirically observed *fast-then-slow* improvement pattern as the number of agents $n$ increases: early gains occur when $K(n)$ is still growing, while later gains diminish once $K(n)$ saturates.

### 4.3. Measuring Effective Channels Without Labels: $K^*$

The effective channel count $K$ cannot be computed directly at inference time because it depends on the unknown ground-truth $Y$. We therefore introduce $K^*$, a *label-free proxy* that estimates the number of effective channels from agent outputs in embedding space: $K^*$ is large when outputs are diverse and approaches 1 when outputs are similar.

**Definition.** Let $\mathrm{Emb}(\cdot)$ be an embedding model. Given outputs $\{Z_i\}_{i=1}^n$, define normalized embeddings

$$\hat{\mathbf{z}}_i \;:=\; \frac{\mathrm{Emb}(Z_i)}{\|\mathrm{Emb}(Z_i)\|_2} \in \mathbb{R}^d, \qquad (16)$$

and the cosine-similarity Gram matrix $G \in \mathbb{R}^{n \times n}$:

$$G_{ij} \;:=\; \langle \hat{\mathbf{z}}_i, \hat{\mathbf{z}}_j \rangle. \qquad (17)$$

Trace-normalize to obtain $\rho := G/\mathrm{Tr}(G)$ with $\mathrm{Tr}(\rho) = 1$, and let $\{\lambda_j\}_{j=1}^n$ be the eigenvalues of $\rho$. We define the entropy effective rank

$$K^* := 2^{H(\rho)}, \quad \text{where} \quad H(\rho) = -\sum_{j=1}^n \lambda_j \log_2 \lambda_j. \quad (18)$$

**Interpretation.** $K^*$ counts how many "independent directions" the agent outputs span in embedding space. When all agents produce nearly identical outputs (e.g., paraphrases of the same reasoning), their embeddings are collinear and $K^* \approx 1$: the system effectively has a single information channel. When agents produce genuinely different outputs whose embeddings point in different directions with roughly equal magnitude, $K^*$ grows toward $n$: each agent contributes a distinct channel. For example, if four agents all solve a math problem using the same algebraic approach with minor wording differences, their outputs cluster in one direction and $K^* \approx 1$. If instead the agents employ genuinely different strategies (e.g., algebraic manipulation, geometric reasoning, and numerical estimation), $K^*$ will be notably larger than 1, reflecting that the system draws on multiple independent lines of evidence. Formally, $1 \le K^* \le n$. $K^*$ reaches its maximum $n$ when the normalized Gram matrix $\rho$ has a uniform spectrum (all eigenvalues equal to $1/n$), which occurs when outputs are orthogonal and carry equal energy. Proofs of these properties are given in Appendix A.

## 5. Experiments

This section validates 3 core claims: (i) scaling *homogeneous* MAS exhibits diminishing returns, (ii) *heterogeneity* consistently outperforms pure scaling under matched compute, and (iii) performance gains are governed by the *number of effective channels* rather than the raw agent count.

### 5.1. Experimental Setup

**Tasks.** We consider a diverse set of reasoning and knowledge benchmarks, including GSM8K (Cobbe et al., 2021), ARC (Clark et al., 2018), Formal Logic (Hendrycks et al., 2021a;b), TruthfulQA (Lin et al., 2022), Hel-
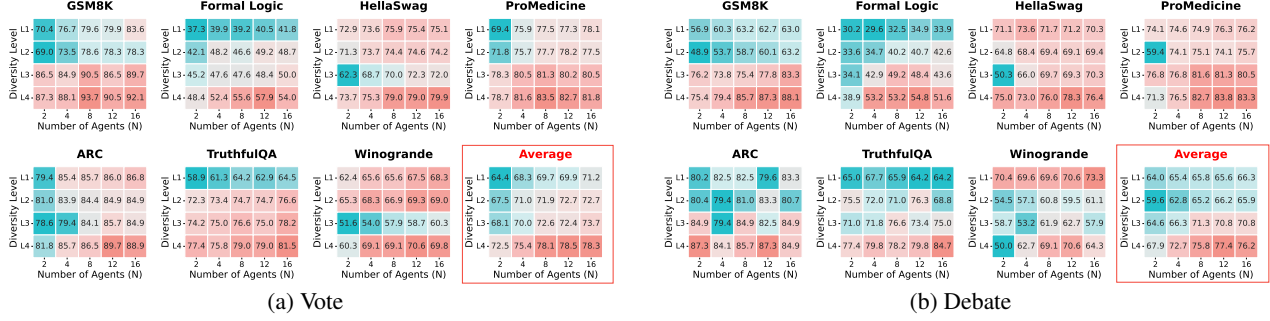
Figure 4: Diversity analysis of the **Vote** and **Debate** mechanisms across all datasets. Each subfigure corresponds to one dataset and visualizes absolute performance as a $4 \times 5$ heatmap, where rows represent progressively enriched diversity layers (L1–L4) and columns denote the number of agents $N$. Colors indicate success rate values: cyan (lowest) to red (highest).

laSwag (Zellers et al., 2019), WinoGrande (ai2, 2019), and Pro Medicine (Hendrycks et al., 2021a;b). These tasks span arithmetic reasoning, formal deduction, commonsense reasoning, and domain knowledge, covering both deterministic and ambiguous settings.

**Models.** Agents are instantiated using three open-source LLMs: Qwen-2.5-7B (Qwen Team, 2024), Llama-3.1-8B (Grattafiori et al., 2024), and Mistral-7B (Jiang et al., 2023). In the *single-model* setting, all agents within a MAS share the same base model; in the *MIX* setting, agents within a single MAS can use different base models, enabling model-level heterogeneity.

**MAS Workflows.** We consider two representative collaboration mechanisms (Choi et al., 2025): **Vote**, where agents independently generate answers and a majority decision is taken after 1 round, and **Debate**, where agents interact sequentially for 4 rounds before producing a final answer. For each mechanism, we vary the number of agents $N \in \{2, 4, 8, 12, 16\}$. Compute budgets are matched by fixing the total number of agent calls.

**Diversity Configurations.** We organize agent heterogeneity into four progressively enriched layers to isolate the contribution of each diversity source:

- **L1: No Diversity.** All agents share the same base model and the same default system prompt (no persona). This serves as the homogeneous baseline. Results are averaged over the three single-model runs.
- **L2: Persona Diversity Only.** All agents share the same base model, but each agent receives a distinct persona prompt (e.g., "You are an expert mathematician" vs. "You are a careful logician"). Results are averaged over the three single-model runs.
- **L3: Model Diversity Only.** Agents are drawn from different base models (Qwen, Llama, Mistral) but all use the same default system prompt.
- **L4: Full Diversity.** Agents differ in both base model and persona prompt, combining model-level and prompt-level heterogeneity.

Table 2: Efficiency gains from diversity. Number of agents needed to match L1 (N=16) baseline. Higher diversity achieves equivalent performance with fewer agents.

| Method | Config | Agents to Match L1 (N=16) | Accuracy at that N | Peak Accuracy (any N) |
|---|---|---|---|---|
| Vote | L1 | 16 (baseline) | 65.34 | 65.49 |
| | L2 | 8 | 65.44 | 66.01 |
| | L3 | 4 | 67.29 | 71.54 |
| | L4 | 2 | 67.71 | **76.86** |
| Debate | L1 | 16 (baseline) | 65.48 | 65.48 |
| | L2 | 12 | 66.08 | 66.08 |
| | L3 | 4 | 66.26 | 71.33 |
| | L4 | 2 | 67.90 | **77.43** |

This controlled design allows us to isolate and compare the contributions of model diversity and persona diversity.

### 5.2. Finding 1: Scaling Homogeneous MAS Exhibits Diminishing Returns

We first examine whether increasing the number of agents improves performance in homogeneous settings. Figure 2 shows success rates and marginal gains for both voting- and debate-based MAS across multiple tasks and base models.

Across all settings, we observe a consistent pattern: accuracy improves only at small agent counts, after which marginal gains $\Delta \text{Success}/\Delta N$ rapidly collapse toward zero. In several cases, performance even degrades as $N$ increases.

As predicted by our theoretical framework (Theorem A.15), this saturation occurs because homogeneous agents produce highly correlated outputs, so additional calls fail to increase the effective channel count $K$. In other words, allocating more test-time computation via homogeneous scaling does not reliably inject new usable evidence into the system.

### 5.3. Finding 2: Diversity Consistently Beats Scale

We compare homogeneous scaling with heterogeneous designs under matched compute in Table 1, which reports the performance of Vote and Debate mechanisms across all tasks and agent counts. In nearly all cases, heterogeneous
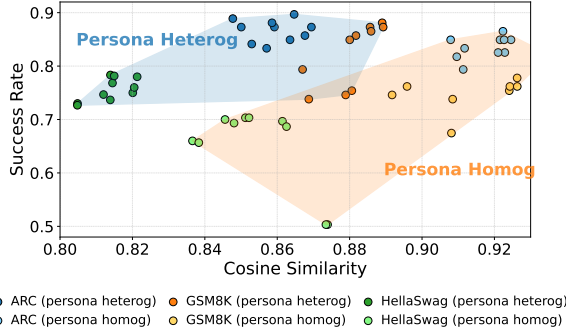
Figure 5: Correlation between cosine similarity and success rate. Homogeneous settings show higher similarity but lower performance; heterogeneous personas preserve diversity and improve accuracy.
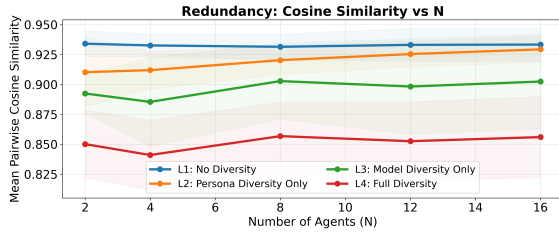


Figure 6: Mean pairwise cosine similarity vs. agent count. Higher diversity (L1→L4) consistently reduce redundancy.

configurations significantly outperform homogeneous ones, with gains increasing as $N$ grows. Figure 4 provides a detailed view of this effect. Enriching diversity from L1 to L4 yields consistent performance improvements for both Vote and Debate. Notably, model diversity (L3) and persona diversity (L2) each deliver non-trivial gains, while their combination (L4) consistently performs best.

Table 2 shows the minimum number of heterogeneous agents required to outperform homogeneous configurations. For both Vote and Debate, L4 (full diversity) with just **2 agents** surpasses the performance of L1 (no diversity) with **16 agents**. This represents an $8\times$ reduction in agent count for equivalent or better accuracy. This result directly reflects the theory: by Corollary 4.4, the heterogeneous design achieves a higher $\alpha K$ product, so fewer agents suffice to reach the same information-recovery level.

We also compare heterogeneous model mixtures against independent single-model runs. Figure 1 demonstrates that a mixture of three LLMs outperforms the average performance of the individual models, confirming that the improvements stem from complementary effective channels rather than simple averaging.

### 5.4. Finding 3: Performance Gains Are Governed by the Number of Effective Channels

Our theory predicts that homogeneous agents produce highly correlated outputs, contributing few effective chan-

Table 3: Relation between K* and Accuracy on ARC.

| Method | Config | Performance | | Channels | | Answer-Cond. | |
|---|---|---|---|---|---|---|---|
| | | Acc. | $\Delta$Acc | $K^*$ | $\Delta K^*$ | $K_c^*$ | $K_w^*$ |
| Debate | L1 | 81.6% | – | 1.197 | – | 1.184 | 1.177 |
| | L2 | 81.0% | -0.7 | 1.348 | +0.152 | 1.315 | 1.234 |
| | L3 | 83.3% | +1.7 | 1.246 | +0.049 | 1.220 | 1.160 |
| | L4 | **85.9%** | **+4.2** | **1.517** | **+0.320** | **1.472** | 1.288 |
| Vote | L1 | 81.3% | – | 1.201 | – | 1.183 | 1.173 |
| | L2 | 81.5% | +0.2 | 1.349 | +0.149 | 1.318 | 1.222 |
| | L3 | 83.8% | +2.5 | 1.245 | +0.044 | 1.223 | 1.161 |
| | L4 | **87.5%** | **+6.1** | **1.521** | **+0.321** | **1.484** | 1.297 |

nels and leading to saturation. We now verify this empirically, proceeding from a simple redundancy proxy (pairwise cosine similarity) to the effective channel measure ($K^*$).

#### 5.4.1. HIGH OUTPUT SIMILARITY HINDERS PERFORMANCE

A key reason homogeneous scaling saturates is that additional agent calls increasingly produce *correlated* outputs, yielding limited *new* evidence. To quantify this redundancy, we embed each agent output (the full reasoning trace) using NV-Embed-v2 (Lee et al., 2025) and compute the *mean pairwise cosine similarity*: for $n$ agent outputs with normalized embeddings $\hat{\mathbf{z}}_1, \ldots, \hat{\mathbf{z}}_n$, this is $\bar{\rho} = \frac{2}{n(n-1)} \sum_{i<j} \langle \hat{\mathbf{z}}_i, \hat{\mathbf{z}}_j \rangle$. While $\bar{\rho}$ is not an information-theoretic quantity, it provides a consistent proxy for output overlap: higher $\bar{\rho}$ indicates that agents explore fewer non-redundant directions, which constrains the growth of effective channels.
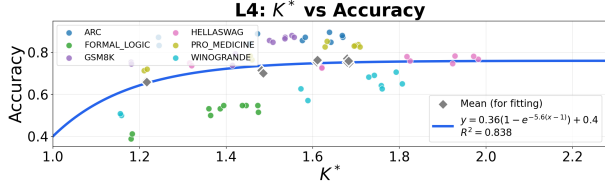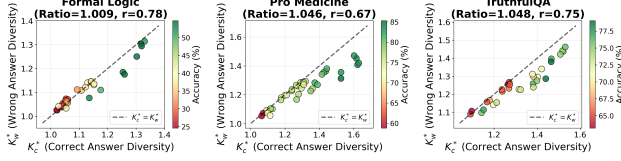
Figure 5 shows that homogeneous persona settings produce higher similarity yet do not translate this additional compute into higher success rates, whereas heterogeneous personas maintain lower similarity and achieve stronger performance. Moreover, Figure 6 reveals a systematic scaling trend: for every diversity layer, redundancy increases with agent count $N$, implying that larger homogeneous ensembles mainly amplify existing trajectories rather than introducing qualitatively new evidence. Crucially, redundancy decreases monotonically from L1 to L4, consistent with our hypothesis that heterogeneity mitigates output correlation and thus enlarges the number of effective channels.

While these results confirm a qualitative relationship between output diversity and performance, pairwise cosine similarity is a coarse measure. To obtain a more precise and theoretically grounded characterization, we next turn to the effective channel count $K^*$ introduced in Section 4.3.

#### 5.4.2. DIVERSE CHANNELS IMPROVES PERFORMANCE

We compute $K^*$ by embedding each agent output with NV-Embed-v2 (Lee et al., 2025), forming the cosine-similarity matrix $G$, trace-normalizing it to $\rho$ with $\text{Tr}(\rho) = 1$, and defining $K^*$ as the entropy effective rank of $\rho$ (Eq. 18).

Figure 7: Correlation between $K^*$ and accuracy.



Figure 8: Decomposition of $K^*$ on three tasks. Points below the diagonal correspond to configurations where correct answer diversity dominates.

**Diversity increases $K^*$.** As shown in Table 3, $K^*$ consistently increases with diversity level from L1 to L4 under both Vote and Debate mechanisms, validating $K^*$ as a robust indicator of system diversity without ground-truth labels.

**Higher $K^*$ leads to better performance.** The increase in $K^*$ is accompanied by higher accuracy in most cases (Table 3). Figure 7 further confirms this positive correlation, depicting a strong linear relationship between $K^*$ and task accuracy across configurations. Moreover, consistent with Theorem A.15, the marginal improvement in accuracy diminishes as $K^*$ grows, reflecting the geometric decay $(1 - \alpha)^K$ predicted by our theory. We observe a minor anomaly in L2 under Debate, where $K^*$ increases but accuracy slightly decreases; we investigate this through the decomposition of $K^*$ below.

**Mechanistic Decomposition: $K_c^*$ vs. $K_w^*$.** To determine if the growth in $K^*$ represents useful evidence or merely increased noise, we decompose it into $K_c^*$ (correct reasoning diversity) and $K_w^*$ (incorrect reasoning diversity). Let $\hat{y}_i$ represent the final answer of agent $i$, and $Y$ be the ground-truth label. We define:

$$\mathcal{I}_c = \{i : \hat{y}_i = Y\}, \qquad \mathcal{I}_w = \{i : \hat{y}_i \neq Y\}$$

Here, $\mathcal{I}_c$ is the set of correct agents, and $\mathcal{I}_w$ is the set of incorrect agents. We then compute the effective number of channels for each set:

$$K_c^* = K^*(Z_c), \qquad K_w^* = K^*(Z_w)$$

where $Z_c$ and $Z_w$ are the sub-matrices of the original data matrix $Z$, corresponding to correct and incorrect agents.

**The Empirical Boundary.** Figure 8 suggests an empirical boundary in the $(K_c^*, K_w^*)$ plane: high-accuracy configurations concentrate in the region where $K_c^* > K_w^*$ (below

the diagonal line). The intuition is as follows: when multiple agents arrive at the correct answer through genuinely *different* reasoning paths ($K_c^*$ is high), the correct answer receives support from independent evidence sources, making it more robust under aggregation. Conversely, when incorrect answers are also diverse ($K_w^*$ is high), the error "votes" are spread across many competing alternatives, which can dilute the correct signal. Thus, $K_c^* > K_w^*$ indicates that correct reasoning benefits from diverse support while incorrect reasoning remains fragmented, and this favorable signal-to-noise ratio is a prerequisite for robust MAS performance.

### 5.5. Design Guidelines for LLM-based MAS

Our analysis of effective channels yields several data-driven design guidelines for MAS development:

- **Match diversity to task type.** $K^*$ predicts accuracy strongly on reasoning tasks but weakly on knowledge-heavy tasks. For tasks requiring complex multi-step reasoning (e.g., *GSM8K*, *ARC*), investing in diversity yields significant performance gains. In contrast, for tasks dominated by factual retrieval (e.g., *Winogrande*), the diversity investment should be more conservative.
- **Ensure correct-path dominance.** Systems with high $K_c^*/K_w^*$ achieve substantially higher accuracy. In practice, this means that when introducing diversity, one should focus on increasing the diversity of *correct* reasoning paths, for example by using personas that encourage different valid problem-solving strategies (e.g., algebraic vs. geometric approaches in math tasks), rather than indiscriminately adding diversity that may also amplify incorrect reasoning (e.g., random temperature increases that introduce more errors).
- **Right-size agent count.** Homogeneous systems plateau at $N \approx 4$, while heterogeneous systems continue to benefit from scaling up to $N \approx 8$. Beyond this point, adding more agents results in diminishing returns and wasted compute resources. Thus, it is important to find a balance in agent count to avoid inefficiency.

## 6. Conclusion

This paper shows that simply increasing agent count in multi-agent systems results in diminishing returns, both for homogeneous and heterogeneous configurations. However, heterogeneity improves performance by introducing more diverse, non-redundant information, delaying saturation. We introduce $K^*$, a label-free measure of effective channels, which reveals that performance gains are driven by the balance between correct-path diversity and redundancy. These results suggest that the challenge in multi-agent scaling lies in the effective allocation of diverse information channels rather than just raw computational power.

# Impact Statement

This work establishes an information-theoretic framework for understanding scaling behavior in LLM-based multi-agent systems. We discuss the scope, limitations, and implications of our contributions below.

**Theoretical Contributions and Scope.** Our framework provides architecture-agnostic upper bounds and lower bounds showing that MAS performance is fundamentally limited by diversity. The geometric contraction result (Theorem A.15) offers a principled explanation for the empirically observed "fast-then-slow" pattern. However, our theoretical analysis relies on idealized assumptions: the evidence-bits model (Assumption A.12) assumes perfect sufficiency and conditional independence of latent evidence, and the coverage model (Assumption A.13) assumes uniform and independent coverage probabilities. Real-world agents may exhibit more complex dependency structures.

**Limitations of $K^*$.** While $K^*$ provides a practical label-free proxy for effective channels, it measures *semantic* diversity in embedding space rather than *task-relevant* information diversity. As shown in Section 5.4.2, the decomposition into $K_c^*$ and $K_w^*$ reveals that not all diversity is beneficial, only diversity among correct reasoning paths reliably improves performance. Furthermore, $K^*$ depends on the choice of embedding model, and its correlation with accuracy varies across task types (stronger for reasoning tasks, weaker for knowledge-intensive tasks). Developing task-adaptive diversity metrics remains an open problem.

**Empirical Scope.** Our experiments focus on 7B-8B scale open-weight models across seven benchmarks. Whether the diversity-over-scale principle extends to larger models, closed-source APIs, or more complex agentic workflows (e.g., tool use, long-horizon planning) requires further investigation. Additionally, our analysis considers vote and debate mechanisms; other coordination protocols may exhibit different scaling behaviors.

# References

Winogrande: An adversarial winograd schema challenge at scale. 2019.

Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

Cemri, M., Pan, M. Z., Yang, S., Agrawal, L. A., Chopra, B., Tiwari, R., Keutzer, K., Parameswaran, A., Klein, D., Ramchandran, K., Zaharia, M., Gonzalez, J. E., and Stoica, I. Why do multi-agent LLM systems fail? In *The Thirty-ninth Annual Conference on Neural Informa-tion Processing Systems Datasets and Benchmarks Track*, 2025.

Chang, E. Y. *Multi-LLM Agent Collaborative Intelligence: The Path to Artificial General Intelligence*. Edward Y. Chang, 2025.

Chen, L., Davis, J. Q., Hanin, B., Bailis, P., Stoica, I., Zaharia, M., and Zou, J. Are more llm calls all you need? towards scaling laws of compound inference systems. *arXiv preprint arXiv:2403.02419*, 2024.

Choi, H. K., Zhu, X., and Li, S. Debate or vote: Which yields better decisions in multi-agent large language models? *arXiv preprint arXiv:2508.17536*, 2025.

Clark, P., Cowhey, I., Etzioni, O., Khot, T., Sabharwal, A., Schoenick, C., and Tafjord, O. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv:1803.05457v1*, 2018.

Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., Hesse, C., and Schulman, J. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.

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 *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025.

Du, Y., Li, S., Torralba, A., Tenenbaum, J. B., and Mordatch, I. Improving factuality and reasoning in language models through multiagent debate. *arXiv preprint arXiv:2305.14325*, 2023.

Gan, Z., Liao, Y., and Liu, Y. Rethinking external slow-thinking: From snowball errors to probability of correct reasoning. *arXiv preprint arXiv:2501.15602*, 2025.

Grattafiori, A. et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

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. *arXiv preprint arXiv:2402.01680*, 2024.

Hendrycks, D., Burns, C., Basart, S., Critch, A., Li, J., Song, D., and Steinhardt, J. Aligning ai with shared human values. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021a.

Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., and Steinhardt, J. Measuring massive multitask language understanding. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021b.

Hong, S., Zhuge, M., Chen, J., Zheng, X., Cheng, Y., Zhang, C., Wang, J., Wang, Z., Yau, S. K. S., Lin, Z., Zhou, L., Ran, C., Xiao, L., Wu, C., and Schmidhuber, J. Metagpt: Meta programming for a multi-agent collaborative framework. *arXiv preprint arXiv:2308.00352*, 2024.

Huang, J., Gu, X., Chen, L., Han, J., and Kraska, T. Large language models cannot self-correct reasoning yet. *arXiv preprint arXiv:2310.01798*, 2023.

Jiang, A. Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D. S., de las Casas, D., Bressand, F., Lengyel, G., Lample, G., Saulnier, L., Lavaud, L. R., Lachaux, M.-A., Stock, P., Scao, T. L., Lavril, T., Wang, T., Lacroix, T., and Sayed, W. E. Mistral 7b, 2023.

Khan, A., Hughes, J., Valentine, D., Ruis, L., Sachan, K., Radhakrishnan, A., Grefenstette, E., Bowman, S. R., Rocktäschel, T., and Perez, E. Debating with more persuasive LLMs leads to more truthful answers. In *Proceedings of the 41st International Conference on Machine Learning*, ICML'24, 2024.

Kim, Y., Gu, K., Park, C., Park, C., Schmidgall, S., Heydari, A. A., Yan, Y., Zhang, Z., Zhuang, Y., Malhotra, M., Liang, P. P., Park, H. W., Yang, Y., Xu, X., Du, Y., Patel, S., Althoff, T., McDuff, D., and Liu, X. Towards a science of scaling agent systems. *arXiv preprint arXiv:2512.08296*, 2025.

Lee, C., Roy, R., Xu, M., Raiman, J., Shoeybi, M., Catanzaro, B., and Ping, W. Nv-embed: Improved techniques for training llms as generalist embedding models, 2025.

Li, G., Al Kader Hammoud, H. A., Itani, H., Khizbullin, D., and Ghanem, B. Camel: Communicative agents for "mind" exploration of large language model society. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, NeurIPS '23, 2023.

Li, J., Zhang, Q., Yu, Y., Fu, Q., and Ye, D. More agents is all you need. *arXiv preprint arXiv:2402.05120*, 2024a.

Li, Y., Du, Y., Zhang, J., Hou, L., Grabowski, P., Li, Y., and Ie, E. Improving multi-agent debate with sparse communication topology. In Al-Onaizan, Y., Bansal, M., and Chen, Y.-N. (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 7281–7294, Miami, Florida, USA, November 2024b. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-emnlp.427.

Lin, S., Hilton, J., and Evans, O. Truthfulqa: Measuring how models mimic human falsehoods, 2022.

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 *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 15174–15186. Association for Computational Linguistics, 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 *The Thirteenth International Conference on Learning Representations*, 2025.

Qwen Team. Qwen2.5: A party of foundation models. *arXiv preprint arXiv:2412.15115*, 2024.

Riedl, C. Emergent coordination in multi-agent language models. *arXiv preprint arXiv:2510.05174*, 2025.

Samuel, V., Zou, H. P., Zhou, Y., Chaudhari, S., Kalyan, A., Rajpurohit, T., Deshpande, A., Narasimhan, K., and Murahari, V. Personagym: Evaluating persona agents and llms. *arXiv preprint arXiv:2407.18416*, 2024.

Taillandier, P., Zucker, J. D., Grignard, A., Gaudou, B., Huynh, N. Q., and Drogoul, A. Integrating llm in agent-based social simulation: Opportunities and challenges. *arXiv preprint arXiv:2507.19364*, 2025.

Tang, B., Liang, H., Jiang, K., and Dong, X. On the importance of task complexity in evaluating llm-based multi-agent systems. *arXiv preprint arXiv:2510.04311*, 2025.

Ton, J.-F., Taufiq, M. F., and Liu, Y. Understanding chain-of-thought in llms through information theory. *arXiv preprint arXiv:2411.11984*, 2024.

Wang, J., Wang, J., Athiwaratkun, B., Zhang, C., and Zou, J. Mixture-of-agents enhances large language model capabilities. *arXiv preprint arXiv:2406.04692*, 2024a.

Wang, L., Ma, C., Feng, X., Zhang, Z., Yang, H., Zhang, J., Chen, Z., 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*, 18 (6):186345, 2024b. doi: 10.1007/s11704-024-40231-1.

Wang, X., Wei, J., Schuurmans, D., Le, Q., Chi, E., Narang, S., Chowdhery, A., and Zhou, D. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2023.

Wei, J., Wang, X., Schuurmans, D., Maeda, M., Chi, E., Xia, S., Le, Q., and Zhou, D. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837, 2022.

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 next-gen LLM applications via multi-agent conversations. In *First Conference on Language Modeling*, 2024.

Wu, Z. and Ito, T. The hidden strength of disagreement: Unraveling the consensus-diversity tradeoff in adaptive multi-agent systems. *arXiv preprint arXiv:2502.16565*, 2025.

Xi, Z., Chen, W., Guo, X., He, W., Ding, Y., Hong, B., Zhang, M., Wang, J., Jin, S., Zhou, E., Zheng, R., Fan, X., Wang, X., Xiong, L., Zhou, Y., Wang, W., Jiang, C., Zou, Y., Liu, X., Yin, Z., Dou, S., Weng, R., Cheng, W., Zhang, Q., Qin, W., Zheng, Y., Qiu, X., Huang, X., and Gui, T. The rise and potential of large language model based agents: A survey. *arXiv preprint arXiv:2309.07864*, 2023.

Yuen, S., Medina, F. G., Su, T., Du, Y., and Sobey, A. J. Intrinsic memory agents: Heterogeneous multi-agent llm systems through structured contextual memory. *arXiv preprint arXiv:2508.08997*, 2025.

Zellers, R., Holtzman, A., Bisk, Y., Farhadi, A., and Choi, Y. Hellaswag: Can a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 2019.

Zhang, K., Yao, W., Liu, Z., Feng, Y., Liu, Z., Murthy, R., Lan, T., Li, L., Lou, R., Xu, J., et al. Diversity empowers intelligence: Integrating expertise of software engineering agents. *arXiv preprint arXiv:2408.07060*, 2024.

Zhao, W., Yuksekgonul, M., Wu, S., and Zou, J. Sirius: Self-improving multi-agent systems via bootstrapped reasoning. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025.

# A. Proofs and Technical Details

This appendix provides full proofs and technical details omitted from the main text. Section A.1 reviews standard information-theoretic identities. Section A.2 proves Theorem 3.2. Sections A.3–A.4 derive upper bounds for common MAS workflows. Section A.5 presents the evidence-bits coverage model and proves Theorem A.15 and Corollary A.16 from the main text. Section A.6 proves basic properties of the effective channel count $K^*$.

## A.1. Information-Theoretic Preliminaries

We recall standard definitions and lemmas from information theory.

**Definition A.1** (Conditional Mutual Information). For random variables $A, B, C$,

$$I(A; B \mid C) = H(A \mid C) - H(A \mid B, C) = H(B \mid C) - H(B \mid A, C). \tag{19}$$

**Lemma A.2** (Chain Rule for Mutual Information). *For random variables $X, Y_1, \ldots, Y_n, Z$,*

$$I(Y_1, \ldots, Y_n; X \mid Z) = \sum_{i=1}^{n} I(Y_i; X \mid Z, Y_{<i}). \tag{20}$$

**Lemma A.3** (Data Processing Inequality). *If $A \to B \to C$ forms a Markov chain, then $I(A; C) \leq I(A; B)$.*

**Lemma A.4** (Incremental Information as Entropy Difference). *Let $\Delta_i := I(Z_i; Y \mid X, Z_{<i})$. Then*

$$\Delta_i = H(Y \mid X, Z_{<i}) - H(Y \mid X, Z_{\leq i}), \tag{21}$$

*where $Z_{\leq i} = (Z_1, \ldots, Z_i)$.*

*Proof.* By the definition of conditional mutual information,

$$\Delta_i = I(Z_i; Y \mid X, Z_{<i}) \tag{22}$$
$$= H(Y \mid X, Z_{<i}) - H(Y \mid X, Z_{<i}, Z_i) \tag{23}$$
$$= H(Y \mid X, Z_{<i}) - H(Y \mid X, Z_{\leq i}). \qquad \square$$

## A.2. Finite Information Budget (Upper Bound)

For completeness, the total information an MAS can extract is always upper-bounded by the intrinsic task uncertainty.

**Theorem A.5** (Finite Information Budget). *For any MAS transcript $Z_{1:n}$,*

$$I(Z_{1:n}; Y \mid X) \leq H(Y \mid X). \tag{24}$$

*Moreover, writing $I(Z_{1:n}; Y \mid X) = \sum_{i=1}^{n} \Delta_i$ with $\Delta_i := I(Z_i; Y \mid X, Z_{<i})$, we have $\sum_{i=1}^{n} \Delta_i \leq H(Y \mid X)$ and $\Delta_i \to 0$ as $i \to \infty$.*

*Proof.* By the definition of conditional mutual information,

$$I(Z_{1:n}; Y \mid X) = H(Y \mid X) - H(Y \mid X, Z_{1:n}) \leq H(Y \mid X), \tag{25}$$

since conditional entropy is nonnegative. By Lemma A.2,

$$I(Z_{1:n}; Y \mid X) = \sum_{i=1}^{n} I(Z_i; Y \mid X, Z_{<i}) = \sum_{i=1}^{n} \Delta_i. \tag{26}$$

Because $\Delta_i \geq 0$ and the partial sums are uniformly bounded by $H(Y \mid X)$, we must have $\Delta_i \to 0$ as $i \to \infty$. $\qquad \square$

## A.3. Parallel Voting: Assumptions and Upper Bounds

This section derives the parallel-voting upper bounds used in the main text. The key message is that repeated sampling from the same configuration produces redundant evidence.

A.3.1. CONDITIONAL INDEPENDENCE FOR PARALLEL SAMPLING

**Assumption A.6** (Conditional Independence for Parallel Sampling (All Types)). Consider a parallel MAS with agent configuration types $b(i) \in \mathcal{B}$. There exist channels $\{K_b(\cdot \mid x, y)\}_{b \in \mathcal{B}}$ such that, for every $i$,

$$P(Z_i = z \mid X = x, Y = y, Z_{<i}) = K_{b(i)}(z \mid x, y), \tag{27}$$

and the outputs are mutually independent conditioned on $(X, Y)$:

$$P(Z_{1:n} \mid X, Y) = \prod_{i=1}^{n} P(Z_i \mid X, Y, b(i)). \tag{28}$$

Define the single-call information for type $b$:

$$I_b := I(Z^{(b)}; Y \mid X), \tag{29}$$

where $Z^{(b)}$ denotes one output from type $b$ in isolation.

A.3.2. A REDUNDANCY IDENTITY

**Lemma A.7** (Three-Way Mutual Information Decomposition). *For any random variables $A, B, C, D$,*

$$I(A; B \mid C, D) = I(A; B \mid C) + I(A; D \mid B, C) - I(A; D \mid C). \tag{30}$$

*Proof.* Apply chain rule in two ways:

$$I(A; B, D \mid C) = I(A; B \mid C) + I(A; D \mid B, C), \tag{31}$$
$$I(A; B, D \mid C) = I(A; D \mid C) + I(A; B \mid C, D). \tag{32}$$

Equating and rearranging yields the claim. $\square$

**Corollary A.8** (Incremental Gain under Parallel Sampling). *With $A = Z_i$, $B = Y$, $C = X$, $D = Z_{<i}$,*

$$I(Z_i; Y \mid X, Z_{<i}) = I(Z_i; Y \mid X) + I(Z_i; Z_{<i} \mid X, Y) - I(Z_i; Z_{<i} \mid X). \tag{33}$$

*Under Assumption A.6, $I(Z_i; Z_{<i} \mid X, Y) = 0$, hence*

$$I(Z_i; Y \mid X, Z_{<i}) = I(Z_i; Y \mid X) - I(Z_i; Z_{<i} \mid X) \leq I(Z_i; Y \mid X). \tag{34}$$

This formalizes redundancy: previous outputs can only reduce the new information.

**Implication: redundancy controls early saturation.** Upper bounds identify *what* limits the total information gain. To explain *when* saturation occurs, consider the incremental contribution $\Delta_i := I(Z_i; Y \mid X, Z_{<i})$. Eq. (34) provides an explicit decomposition:

$$\Delta_i = I(Z_i; Y \mid X) - I(Z_i; Z_{<i} \mid X), \tag{35}$$

where the *redundancy term* $I(Z_i; Z_{<i} \mid X)$ quantifies how much the $i$-th output overlaps with previous outputs. Thus, early saturation arises when repeated calls increase $I(Z_i; Z_{<i} \mid X)$, leaving little additional evidence to accumulate. Homogeneous agents typically induce large redundancy due to similar reasoning trajectories, while heterogeneity mitigates overlap and sustains $\Delta_i$.

Since $I(Z_i; Y \mid X, Z_{<i}) \geq 0$, the identity also implies $I(Z_i; Z_{<i} \mid X) \leq I(Z_i; Y \mid X)$ under Assumption A.6.

A.3.3. HOMOGENEOUS PARALLEL BOUND

**Proposition A.9** (Homogeneous Parallel Upper Bound). *Assume $m$ parallel samples from a single type $b$ under Assumption A.6. Then*

$$I(Z_{1:m}; Y \mid X) \leq H(Y \mid X) \wedge mI_b. \tag{36}$$

*Proof.* By chain rule,

$$I(Z_{1:m}; Y \mid X) = \sum_{i=1}^{m} I(Z_i; Y \mid X, Z_{<i}). \tag{37}$$

Using Eq. (34) and $I(Z_i; Y \mid X) = I_b$ for all $i$,

$$I(Z_{1:m}; Y \mid X) \le \sum_{i=1}^{m} I_b = m I_b. \tag{38}$$

The finite budget further implies $I(Z_{1:m}; Y \mid X) \le H(Y \mid X)$. Combining yields Eq. (36). □

### A.3.4. HETEROGENEOUS PARALLEL BOUND

**Theorem A.10** (Heterogeneous Parallel Upper Bound). *Consider parallel voting with configuration types $\mathcal{B}$. Let type $b$ be sampled $m_b$ times, with total $n = \sum_{b \in \mathcal{B}} m_b$. Then*

$$I(Z_{1:n}; Y \mid X) \ \le \ H(Y \mid X) \ \wedge \ \sum_{b \in \mathcal{B}} m_b \, I_b. \tag{39}$$

*Proof.* Apply the chain rule:

$$I(Z_{1:n}; Y \mid X) = \sum_{i=1}^{n} I(Z_i; Y \mid X, Z_{<i}). \tag{40}$$

By Eq. (34), each term is bounded by $I(Z_i; Y \mid X) = I_{b(i)}$. Summing over steps grouped by type gives $\sum_{b \in \mathcal{B}} m_b \, I_b$. The finite budget gives the minimum with $H(Y \mid X)$. □

## A.4. Sequential Pipelines and Debate: Upper Bounds

In sequential settings, each output conditions on the interaction history. This invalidates conditional independence, but the chain rule remains valid.

### A.4.1. MAXIMAL PER-STEP CONTRIBUTION

Define the maximal incremental contribution for agent configuration type $b$:

$$I_b^{\max} := \sup_{z_{<i}} I(Z_i; Y \mid X, Z_{<i} = z_{<i}, b(i) = b). \tag{41}$$

**Proposition A.11** (Sequential Pipeline Upper Bound). *For any sequential MAS with $n$ steps,*

$$I(Z_{1:n}; Y \mid X) \ \le \ H(Y \mid X) \ \wedge \ \sum_{i=1}^{n} I_{b(i)}^{\max}. \tag{42}$$

*Proof.* By chain rule,

$$I(Z_{1:n}; Y \mid X) = \sum_{i=1}^{n} I(Z_i; Y \mid X, Z_{<i}). \tag{43}$$

For each $i$, by definition of $I_{b(i)}^{\max}$ we have $I(Z_i; Y \mid X, Z_{<i}) \le I_{b(i)}^{\max}$. Summing yields the stated bound, and the finite budget gives the minimum with $H(Y \mid X)$. □

**Debate.** Two-agent debate is a special case of sequential interaction and inherits the same ceiling $H(Y \mid X)$. This formalizes why debate cannot systematically improve over voting if agents remain redundant.

## A.5. Lower Bound via Independent Evidence-Bits Coverage

This section formalizes the "effective channels" view used in the main text. It proves Theorem A.15 (geometric contraction of the residual uncertainty) and Corollary A.16 (the saturated lower bound in expectation), which together imply a characteristic rapid-then-saturating improvement curve $1 - e^{-\alpha K}$ emphasized in Eq. (14) of the main text.

### A.5.1. EVIDENCE BITS MODEL

**Assumption A.12** (Independent Evidence Bits). There exist latent variables $U = (U_1, \ldots, U_M)$ such that:

1. (**Sufficiency**) $H(Y \mid X, U) = 0$.

2. (**Conditional independence**) $U_1, \ldots, U_M$ are independent conditioned on $X$.

3. (**Matching uncertainty scale**) $H(U \mid X) = H(Y \mid X)$.

Assumption A.12(iii) calibrates the latent "evidence bits" to exactly match the intrinsic task uncertainty. In particular, recovering all evidence bits eliminates residual uncertainty about $Y$.

### A.5.2. FRACTIONAL COVERAGE BY EFFECTIVE CHANNELS

**Assumption A.13** (Fractional Evidence Coverage). Let $\tilde{Z}_{1:K}$ denote $K$ effective channels extracted from an MAS transcript. For each evidence bit $U_j$ and each channel $k \in \{1, \ldots, K\}$, define a Bernoulli indicator $C_{j,k} \in \{0,1\}$: $C_{j,k} = 1$ means channel $k$ reveals $U_j$. Assume:

1. $\mathbb{P}(C_{j,k} = 1) = \alpha$ for some fixed $\alpha \in (0,1)$.

2. For each fixed $j$, $\{C_{j,k}\}_{k=1}^{K}$ are independent.

3. If $\exists k$ such that $C_{j,k} = 1$, then $H(U_j \mid X, \tilde{Z}_{1:K}) = 0$.

Assumption A.13 is a minimal complementarity model: each *new effective channel* has a constant probability $\alpha$ of covering any remaining evidence bit, independently across channels.

### A.5.3. RESIDUAL CONTRACTION AND SATURATED LOWER BOUND

**Lemma A.14** (Expected Geometric Decay of Residual Uncertainty). *Under Assumptions A.12 and A.13,*

$$\mathbb{E}\big[H(Y \mid X, \tilde{Z}_{1:K})\big] \leq (1-\alpha)^K H(Y \mid X). \tag{44}$$

*Proof.* By Assumption A.12(i), $Y$ is a function of $(X, U)$, hence

$$H(Y \mid X, \tilde{Z}_{1:K}) \leq H(U \mid X, \tilde{Z}_{1:K}). \tag{45}$$

Subadditivity of conditional entropy yields

$$H(U \mid X, \tilde{Z}_{1:K}) \leq \sum_{j=1}^{M} H(U_j \mid X, \tilde{Z}_{1:K}). \tag{46}$$

Fix $j$. If $U_j$ is revealed by at least one effective channel, then by Assumption A.13(iii), $H(U_j \mid X, \tilde{Z}_{1:K}) = 0$; otherwise, $H(U_j \mid X, \tilde{Z}_{1:K}) \leq H(U_j \mid X)$. The probability that $U_j$ is not revealed by any of the $K$ channels is $(1-\alpha)^K$ by Assumption A.13(ii). Therefore,

$$\mathbb{E}\big[H(U_j \mid X, \tilde{Z}_{1:K})\big] \leq (1-\alpha)^K H(U_j \mid X). \tag{47}$$

Summing over $j$ gives

$$\mathbb{E}\big[H(U \mid X, \tilde{Z}_{1:K})\big] \leq (1-\alpha)^K \sum_{j=1}^{M} H(U_j \mid X). \tag{48}$$

Finally, by Assumption A.12(ii), $H(U \mid X) = \sum_{j=1}^{M} H(U_j \mid X)$, and by (iii) $H(U \mid X) = H(Y \mid X)$. Combining completes the proof. $\square$

**Theorem A.15** (Geometric Contraction with Effective Channels). *Under Assumptions A.12 and A.13,*

$$H(Y \mid X) - \mathbb{E}\big[I(\tilde{Z}_{1:K}; Y \mid X)\big] \;=\; \mathbb{E}\big[H(Y \mid X, \tilde{Z}_{1:K})\big] \;\leq\; (1-\alpha)^K H(Y \mid X). \tag{49}$$

*Consequently,*

$$\boxed{\begin{aligned} &H(Y \mid X) - \mathbb{E}\Big[I(\tilde{Z}_{1:K}; Y \mid X)\Big] \\ &\leq\; (1-\alpha)^K H(Y \mid X) \;\leq\; e^{-\alpha K} H(Y \mid X). \end{aligned}} \tag{50}$$

*Equivalently, the normalized residual satisfies* $\mathbb{E}[H(Y \mid X, \tilde{Z}_{1:K})]/H(Y \mid X) \leq (1-\alpha)^K \leq e^{-\alpha K}$.

*Proof.* By definition, $I(\tilde{Z}_{1:K}; Y \mid X) = H(Y \mid X) - H(Y \mid X, \tilde{Z}_{1:K})$. Taking expectations yields

$$H(Y \mid X) - \mathbb{E}[I(\tilde{Z}_{1:K}; Y \mid X)] = \mathbb{E}[H(Y \mid X, \tilde{Z}_{1:K})]. \tag{51}$$

Apply Lemma A.14 to obtain (49). The exponential form follows from $(1-\alpha)^K \leq e^{-\alpha K}$. $\qquad\square$

**Corollary A.16** (Saturated Lower Bound (in Expectation)). *Under the same assumptions,*

$$\boxed{\mathbb{E}\big[I(\tilde{Z}_{1:K}; Y \mid X)\big] \;\geq\; H(Y \mid X)\Big(1 - (1-\alpha)^K\Big) \;\geq\; H(Y \mid X)\Big(1 - e^{-\alpha K}\Big).} \tag{52}$$

*Proof.* Rearrange the identity $\mathbb{E}[I(\tilde{Z}_{1:K}; Y \mid X)] = H(Y \mid X) - \mathbb{E}[H(Y \mid X, \tilde{Z}_{1:K})]$ and apply Lemma A.14. The exponential form follows from $(1-\alpha)^K \leq e^{-\alpha K}$. $\qquad\square$

### A.5.4. HETEROGENEITY ADVANTAGE AS AN $\alpha K$ COMPARISON

This subsection provides a formal underpinning for the main-text comparison (Corollary 4.4): heterogeneity improves expected recoverable information whenever it increases the effective evidence term $\alpha K$.

**Lemma A.17** (Monotonicity in $\alpha K$). *Define* $f(t) := 1 - e^{-t}$ *for* $t \geq 0$. *Then for any* $t_1, t_2 \geq 0$, $t_2 > t_1$ *implies* $f(t_2) > f(t_1)$.

*Proof.* $f'(t) = e^{-t} > 0$ for all $t \geq 0$, hence $f$ is strictly increasing. $\qquad\square$

**Corollary A.18** (Heterogeneity Advantage from Corollary A.16). *Consider two designs summarized by* $(K_{\mathrm{homog}}, \alpha_{\mathrm{homog}})$ *and* $(K_{\mathrm{heterog}}, \alpha_{\mathrm{heterog}})$ *under Assumptions A.12–A.13. By Corollary A.16, the lower bounds on recoverable information for the two designs are:*

$$\mathbb{E}[I_{\mathrm{heterog}}] \geq H(Y \mid X)\big(1 - e^{-\alpha_{\mathrm{heterog}} K_{\mathrm{heterog}}}\big), \tag{53}$$

$$\mathbb{E}[I_{\mathrm{homog}}] \geq H(Y \mid X)\big(1 - e^{-\alpha_{\mathrm{homog}} K_{\mathrm{homog}}}\big). \tag{54}$$

*When* $\alpha_{\mathrm{heterog}} K_{\mathrm{heterog}} > \alpha_{\mathrm{homog}} K_{\mathrm{homog}}$, *the heterogeneous design enjoys a strictly higher information-recovery guarantee, since by Lemma A.17 the function* $1 - e^{-t}$ *is strictly increasing in* $t$.

*Proof.* Apply Corollary A.16 to each design to obtain (53) and (54). Since $\alpha_{\mathrm{heterog}} K_{\mathrm{heterog}} > \alpha_{\mathrm{homog}} K_{\mathrm{homog}}$ and $f(t) = 1 - e^{-t}$ is strictly increasing (Lemma A.17), the lower bound for the heterogeneous design is strictly larger than that for the homogeneous design. $\qquad\square$

### A.6. Properties of the Effective Channel Count $K^*$

This section proves basic properties of the label-free proxy $K^*$ used in the main text (Section 4.3). We restate the definition for completeness.

**Setup.** Given $n$ outputs, let $\hat{\mathbf{z}}_i := \mathrm{Emb}(Z_i)/\|\mathrm{Emb}(Z_i)\|_2 \in \mathbb{R}^d$ be the normalized embeddings, and let $M \in \mathbb{R}^{n \times d}$ be the embedding matrix whose $i$-th row is $\hat{\mathbf{z}}_i^\top$. Define the cosine-similarity Gram matrix $G_{ij} := \langle \hat{\mathbf{z}}_i, \hat{\mathbf{z}}_j \rangle$ (equivalently $G = MM^\top$) and its trace-normalization

$$\rho := \frac{G}{\mathrm{Tr}(G)}. \tag{55}$$

Let $\{\lambda_j\}_{j=1}^n$ be the eigenvalues of $\rho$. The von Neumann entropy is

$$H(\rho) := -\sum_{j=1}^n \lambda_j \log_2 \lambda_j, \tag{56}$$

and the effective channel count is $K^* := 2^{H(\rho)}$.

**Proposition A.19** (Basic Properties of $K^*$). *For any nonzero embedding matrix $M$,*

1. $1 \le K^* \le n$.

2. $K^* = 1$ iff $\rho$ is rank-1 (all embeddings are collinear up to scaling).

3. $K^* = n$ iff $\rho = \frac{1}{n} I_n$ (embeddings are orthogonal with equal norm).

4. $K^*$ is continuous in $M$ (when $\mathrm{Tr}(G) > 0$) and invariant to permutation of outputs.

*Proof.* **(i) Bounds.** Entropy satisfies $0 \le H(\rho) \le \log_2 n$, hence $1 = 2^0 \le 2^{H(\rho)} \le 2^{\log_2 n} = n$.

**(ii) $K^* = 1$.** $H(\rho) = 0$ iff the spectrum is $(1, 0, \ldots, 0)$, which holds iff $\rho$ is rank-1. This corresponds to all rows of $M$ being collinear, i.e., all embeddings are identical up to scaling.

**(iii) $K^* = n$.** $H(\rho) = \log_2 n$ iff $\lambda_j = 1/n$ for all $j$, which occurs when $\rho = \frac{1}{n} I_n$. This corresponds to $G$ being proportional to the identity, i.e., embeddings are orthogonal with equal norm.

**(iv) Continuity and permutation invariance.** The map $M \mapsto G = MM^\top$ is continuous. Normalization by $\mathrm{Tr}(G)$ is continuous when $\mathrm{Tr}(G) > 0$. Eigenvalues of a symmetric matrix vary continuously with the entries, and entropy is continuous on the simplex. Permutation of outputs corresponds to $G \mapsto PGP^\top$ for a permutation matrix $P$, which preserves eigenvalues. $\square$

## B. Supplementary Experiments

### B.1. Closed-Source Model Experiments

We extend our analysis to closed-source models (gpt-4.1-mini, gpt-5-mini) on the Formal Logic benchmark to test whether the heterogeneity advantage generalizes across model families. Table 4 compares closed- and open-source models under homogeneous (Base) and heterogeneous (Heterog) configurations at $N = 2$ and $N = 16$.

**Key findings.** The results confirm that the heterogeneity advantage generalizes to closed-source models, while revealing that its magnitude and scaling behavior vary across model families.

- **Heterogeneity consistently improves over homogeneous baselines.** All five models exhibit positive $\Delta_{\mathrm{Het}}$ in at least one interaction mechanism, confirming that the advantage is not specific to open-source settings.
- **Models with weaker homogeneous baselines benefit more from heterogeneity.** gpt-5-mini achieves near-zero accuracy under homogeneous settings (0–6%) but reaches 35–55% with heterogeneous prompting ($\Delta_{\mathrm{Het}}$ of +39–41%). In contrast, gpt-4.1-mini and the open-source models, which already achieve 25–46% under homogeneous settings, show more modest gains (+3–10%).
- **Scaling trends diverge across model families.** Open-source models exhibit positive scaling ($\Delta_N > 0$) under both configurations. gpt-4.1-mini, however, shows *negative* scaling in debate: accuracy drops from 50.79% to 42.86% ($\Delta_N = -7.93$) even under heterogeneous settings, indicating that adding more agents can hurt when the base model is already strong. gpt-5-mini shows the opposite pattern: under heterogeneous settings it benefits substantially from more agents ($\Delta_N = +19.84$ for Debate), whereas its homogeneous scaling remains near-flat.

Table 4: Closed-source vs. Open-source on Formal Logic (%). $\Delta_{\text{Het}}$: average Heterog gain over Base. $\Delta_N$: accuracy change from $N=2$ to $N=16$ under the Heterog configuration.

| Model | Method | Base | | Heterog | | $\Delta_{\text{Het}}$ | $\Delta_N$ |
|---|---|---|---|---|---|---|---|
| | | N=2 | N=16 | N=2 | N=16 | | |
| *Closed-source* | | | | | | | |
| gpt-4.1-mini | vote | 46.83 | 48.41 | 55.56 | 52.38 | +6.35 | −3.18 |
| | debate | 46.03 | 39.68 | 50.79 | 42.86 | +3.97 | −7.93 |
| gpt-5-mini | vote | 0.00 | 6.35 | 35.71 | 49.21 | **+39.29** | +13.50 |
| | debate | 0.79 | 6.35 | 34.92 | 54.76 | **+41.27** | +19.84 |
| *Open-source* | | | | | | | |
| Qwen-2.5-7B | vote | 38.10 | 44.44 | 45.24 | 50.00 | +6.35 | +4.76 |
| | debate | 30.95 | 34.13 | 25.40 | 38.10 | −0.79 | +12.70 |
| Llama-3.1-8B | vote | 45.24 | 42.86 | 44.44 | 54.76 | +5.55 | +10.32 |
| | debate | 24.60 | 31.75 | 35.71 | 39.68 | +9.52 | +3.97 |
| Mistral-7B | vote | 35.71 | 42.06 | 34.92 | 41.27 | −0.79 | +6.35 |
| | debate | 34.92 | 44.44 | 39.68 | 46.83 | +3.58 | +7.15 |

## B.2. Robustness to Embedding Model Choice

A potential concern is whether our effective channel metric $K^*$ depends critically on the choice of embedding model. To address this, we recompute $K^*$ using a different embedding model, gte-Qwen2-1.5B-instruct (1536 dimensions), and compare the results against our primary model NV-Embed-v2 (4096 dimensions). We conduct this comparison across seven datasets (ARC, Formal Logic, GSM8K, HellaSwag, Pro Medicine, TruthfulQA, WinoGrande), varying agent counts $N \in \{2, 4, 8, 12, 16\}$ and interaction mechanisms (Vote and Debate).

Since embedding dimensionality affects absolute $K^*$ values, direct comparison of raw values across models is not meaningful. Instead, we assess robustness by measuring whether the two embeddings agree on *relative ordering*: within each (configuration type, dataset) pair, do both embeddings rank different (method, $N$) combinations consistently? Across all matched pairs, we observe an average Spearman correlation of $\rho = 0.91$, with over 95% of pairs showing $\rho > 0.5$. This indicates that both embeddings consistently identify which experimental settings produce more diverse outputs, even though their absolute scales differ.

Furthermore, both embeddings yield $K^*$ metrics that positively correlate with task accuracy (NV-Embed-v2: $r = 0.40$; gte-Qwen2: $r = 0.23$), confirming that our core finding, diversity predicts performance, is not an artifact of a particular embedding choice. We use NV-Embed-v2 in the main experiments as it achieves stronger predictive power.

## B.3. Is $K^*$ More Than a Proxy for Scale and Configuration?

Since $K^*$ is computed from agent outputs whose diversity naturally varies with agent count $N$ and configuration type, a key question is whether $K^*$ captures information *beyond* these design variables, or merely serves as a redundant proxy for them. To disentangle this, we fit a baseline regression that predicts task accuracy from $N$ and configuration labels alone, then measure the incremental variance explained ($\Delta R^2$) when $K^*$ or its components are added.

Table 5: **Incremental Explanatory Power of Effective Channels.** The baseline model using only agent count ($N$) and configuration labels explains little variance ($R^2 = 0.062$). Adding $K^*$ substantially improves fit ($\Delta R^2 = +0.147$), and conditioning on answer correctness ($K_c^*$) yields the largest gain ($\Delta R^2 = +0.331$), while $K_w^*$ contributes negligibly.

| Model | $R^2$ | Adj. $R^2$ | $\Delta R^2$ | AIC |
|---|---|---|---|---|
| Baseline ($N$ + Config) | 0.062 | 0.044 | – | 1806.6 |
| Baseline + $K^*$ | 0.209 | 0.190 | +0.147 | 1771.1 |
| Baseline + $K_c^*$ | 0.393 | 0.378 | +0.331 | 1713.0 |
| Baseline + $K_c^* + K_w^*$ | 0.396 | 0.379 | +0.334 | 1713.8 |
| Baseline + $K_c^*/K_w^*$ | 0.325 | 0.309 | +0.263 | 1736.4 |

Table 5 reveals three findings. First, the baseline model with only $N$ and configuration labels achieves $R^2 = 0.062$, confirming that scale and configuration alone are poor predictors of MAS performance. Second, adding $K^*$ raises $\Delta R^2$ by $+0.147$, demonstrating that it captures structural information about output diversity that is not reducible to agent count or configuration choice. Third, and most importantly, replacing $K^*$ with its correctness-conditioned component $K_c^*$ more than doubles the incremental gain ($\Delta R^2 = +0.331$), while further adding $K_w^*$ yields negligible improvement ($\Delta R^2$: $+0.331 \to +0.334$). This asymmetry directly supports our central thesis: what drives MAS performance is not output diversity in general, but specifically the diversity of *correct* reasoning paths. Increasing the number of distinct ways agents arrive at the right answer is far more predictive than total channel count or the diversity of incorrect responses.

### B.4. Sanity Checks: Are $K^*$–Performance Relations Accidental?

We further test whether the observed relationship between effective channels and performance could arise by chance. To this end, we conduct permutation-based randomization tests that preserve the marginal distribution of accuracy while destroying any structural association with $K^*$.

Table 6: **Permutation Sanity Check (1000 shuffles).** Observed correlations between effective-channel metrics and accuracy lie far outside the null distribution, confirming that the relationship is not due to chance.

| Metric | Observed $r$ | $z$-score | $p$ |
|---|---|---|---|
| $K^*$ | 0.388 | 5.87 | <0.001 |
| $K_c^*$ | 0.535 | 7.75 | <0.001 |
| $K_c^*/K_w^*$ | 0.503 | 7.23 | <0.001 |

As shown in Table 6, all effective-channel metrics exhibit $z$-scores well above 5 under permutation testing, with $p < 10^{-3}$. This rules out the possibility that the observed correlations arise from random alignment or dataset-specific artifacts. Notably, $K_c^*$ again yields the strongest signal, reinforcing the interpretation that correct-path diversity is the dominant driver of multi-agent performance.

### B.5. Case Study: Heterogeneity Effects Across Models and Workflows

Table 7 reports a comprehensive ablation study on the Formal Logic benchmark, varying base models, agent counts ($N = 2$–16), and interaction mechanisms. Across nearly all settings, heterogeneous configurations outperform homogeneous ones, often by substantial margins. Importantly, these gains do not arise from scaling alone. For example, in both Vote and Debate, increasing $N$ beyond moderate values frequently yields diminishing or unstable returns in homogeneous settings, while heterogeneous systems maintain consistent improvements. This pattern holds across all three base models and their mixture, indicating that the benefit of heterogeneity is robust to model choice and interaction protocol.

Table 8 isolates the effect of model mixing by comparing a heterogeneous mixture (MIX) against the best-performing single model under the same agent count. At $N \geq 4$, MIX consistently outperforms the strongest individual model by large margins, reaching up to $+14.28\%$ absolute accuracy at $N = 8$.

Crucially, these gains cannot be explained by model selection alone. Even when the best single model is used with heterogeneous prompting, the MIX configuration achieves higher performance, demonstrating genuine synergy across models rather than simple averaging or dominance effects.

Table 7: Model Ablation on Formal Logic: Impact of Heterogeneity from $N = 2$ to $N = 16$

| Base Model | Agents ($N$) | Vote (Round 0) | | | Debate (Final) | | |
|---|---|---|---|---|---|---|---|
| | | Homog | Heterog | $\Delta_{H-M}$ | Homog | Heterog | $\Delta_{H-M}$ |
| Qwen-2.5-7B | 2 | 38.10% | 45.24% | +7.14% | 30.95% | 25.40% | -5.55% |
| | 4 | 42.06% | 53.97% | +11.91% | 30.16% | 34.92% | +4.76% |
| | 8 | 43.65% | 50.00% | +6.35% | 28.57% | 38.10% | +9.53% |
| | 12 | 44.44% | 52.38% | +7.94% | 31.75% | 35.71% | +3.96% |
| | 16 | 44.44% | 50.00% | +5.56% | 34.13% | 38.10% | +3.97% |
| Llama-3.1-8B | 2 | 45.24% | 44.44% | -0.80% | 24.60% | 35.71% | +11.11% |
| | 4 | 42.86% | 53.97% | +11.11% | 23.02% | 24.60% | +1.58% |
| | 8 | 41.27% | 52.38% | +11.11% | 27.78% | 35.71% | +7.93% |
| | 12 | 43.65% | 53.97% | +10.32% | 30.95% | 38.89% | +7.94% |
| | 16 | 42.86% | 54.76% | +11.90% | 31.75% | 39.68% | +7.93% |
| Mistral-7B | 2 | 35.71% | 34.92% | -0.79% | 34.92% | 39.68% | +4.76% |
| | 4 | 34.92% | 36.51% | +1.59% | 35.71% | 44.44% | +8.73% |
| | 8 | 32.54% | 37.30% | +4.76% | 40.48% | 38.89% | -1.59% |
| | 12 | 38.89% | 38.10% | -0.79% | 42.06% | 42.86% | +0.80% |
| | 16 | 42.06% | 41.27% | -0.79% | 44.44% | 46.83% | +2.39% |
| MIX | 2 | 45.24% | 48.41% | +3.17% | 34.13% | 38.89% | +4.76% |
| | 4 | 47.62% | 52.38% | +4.76% | 42.86% | 53.17% | +10.31% |
| | 8 | 47.62% | 55.56% | +7.94% | 49.21% | 53.17% | +3.96% |
| | 12 | 48.41% | 57.94% | +9.53% | 48.41% | 54.76% | +6.35% |
| | 16 | 50.00% | 53.97% | +3.97% | 43.65% | 51.59% | +7.94% |

Table 8: Formal Logic: Synergy of Model Mixing (MIX vs. Best Single Model)

| Agents ($N$) | Best Single (Heterog) | MIX (Heterog) | $\Delta_{\text{MIX vs. Best}}$ | MIX (Homog) | $\Delta_{\text{H-M (MIX)}}$ |
|---|---|---|---|---|---|
| 2 | 39.68% | 38.89% | -0.79% | 34.13% | +4.76% |
| 4 | 44.44% | 53.17% | **+8.73%** | 42.86% | +10.31% |
| 8 | 38.89% | 53.17% | **+14.28%** | 49.21% | +3.96% |
| 12 | 42.86% | 54.76% | **+11.90%** | 48.41% | +6.35% |
| 16 | 46.83% | 51.59% | **+4.76%** | 43.65% | +7.94% |