

Single-Agent Generation Surpasses Multi-Agent Systems in Semantic Diversity

Anonymous ACL submission

Abstract

Multi-Agent Systems (MAS) are commonly used to improve reasoning diversity and robustness by simulating interactions among agents with distinct roles. However, prior work often entangles the contribution of the multi-agent architecture with that of prompt conditioning, making the source of observed diversity gains unclear. We address this confound with a controlled study on divergent thinking tasks, using identical prompt conditioning for MAS and single-agent settings. Under these matched conditions, single-agent setups consistently outperform multi-agent systems in semantic diversity. We attribute this gap to *information visibility*: parallel agents often converge on overlapping ideas, whereas a serial single-agent model can condition on its own generation history to avoid redundancy. We further find that a Multi-Output strategy, which prompts a single agent to produce multiple responses within a single inference pass, achieves the highest diversity without degrading logical validity. Together, these results suggest that costly inter-agent interactions may be unnecessary; instead, efficient information sharing and multi-output generation appear to be effective mechanisms for improving diversity, with implications for designing efficient agentic frameworks.

1 Introduction

Agentic frameworks built on Large Language Models (LLMs) have become a prominent research direction. Within such systems, diversity is widely regarded as a core mechanism for improving reasoning quality, broadening coverage, and increasing robustness (Wang et al., 2023; Yao et al., 2023a; Du et al., 2023; Liang et al., 2024). Multi-Agent Systems (MAS) have become a dominant paradigm for this purpose. MAS instantiate multiple agents with distinct roles and simulate their interactions to produce diverse reasoning trajectories. Although empirical studies support the effectiveness of this

approach (Li et al., 2023a; Liang et al., 2024; Estornell and Liu, 2024), the underlying source of these gains remains insufficiently examined.

In most MAS setups, agents share the same LLM backbone; accordingly, MAS is typically accompanied by carefully designed prompt conditioning that specifies distinct roles or perspectives for different agents (Guo et al., 2024; Li et al., 2023a; Park et al., 2023; Wang et al., 2024). In other words, MAS agents are explicitly seeded with a predefined set of diverse reasoning priors via carefully designed prompt engineering. However, to the best of our knowledge, existing MAS studies do not apply the same diversity-inducing prompt conditioning to their single-agent baselines when making comparisons. This raises a critical question: **are the observed diversity gains driven by the multi-agent architecture itself, or by the well-designed prompt conditioning?**

To investigate this issue, we design a controlled experiment in which single-agent models employ the same prompt conditioning strategies as their multi-agent counterparts. As part of this setup, we also strictly align output volume with MAS by adopting a *Multi-Output* strategy, where a single agent is prompted to generate multiple responses within a single inference pass. Under these controlled conditions, we find that single-agent setups exhibit higher generation diversity than multi-agent systems, contrary to prevailing intuition. Moreover, although the *Multi-Output* strategy is introduced to align output size under the same prompt conditioning, it also emerges as a significant factor in enhancing diversity.

Our results suggest that computationally expensive multi-agent interactions (Wang et al., 2025a; Zeng et al., 2025; Liu et al., 2024a; Wang et al., 2025b) may be unnecessary for achieving diversity, and that simpler, well-controlled prompting strategies can be effective.

2 Related Work

Diversity in Generation and Reasoning Diversity is fundamental for enhancing LLM reasoning capability, robustness, and creativity. However, standard greedy decoding often produces repetitive or generic outputs; consequently, stochastic sampling strategies such as nucleus sampling and top-k decoding are widely adopted to expand the generation space (Holtzman et al., 2020; Fan et al., 2018; Kool et al., 2019; Massarelli et al., 2020). Beyond variance introduced at decoding time, a growing body of work shows that aggregating diverse reasoning paths can substantially improve performance on complex tasks (Wang et al., 2023; Yao et al., 2023a; Besta et al., 2024; Long, 2023; Li et al., 2023b). Similarly, generating diverse candidate solutions for subsequent verification has been effective at reducing hallucinations and logical errors in mathematical and commonsense reasoning (Cobbe et al., 2021; Lightman et al., 2024; Weng et al., 2023; Dhuliawala et al., 2024). Beyond accuracy, diversity also improves robustness to prompt variation, because broader coverage reduces sensitivity to specific phrasing (Tam et al., 2023; Agrawal et al., 2023; Ippolito et al., 2019; Arora et al., 2023; Tevet and Berant, 2021).

Multi-Agent Systems for Eliciting Diversity To systematically induce diverse perspectives, recent research has converged on MAS as a promising paradigm. These frameworks orchestrate interactions among agents with distinct roles and memories to simulate complex social dynamics and collaborative problem-solving (Park et al., 2023; Wu et al., 2024a; Zhuge et al., 2023). A prevalent application is multi-agent debate, in which agents with conflicting viewpoints or assigned personas critique one another to elicit divergent thinking (Du et al., 2023; Liang et al., 2024; Chan et al., 2023; Xiong et al., 2024; Cohen et al., 2023). Theoretical and empirical studies have further analyzed how factors such as communication topology and cognitive synergy shape the diversity of outcomes (Estornell and Liu, 2024; Li et al., 2024; Wang et al., 2024; Zhang et al., 2024; Chen et al., 2024).

Despite these advances, a methodological ambiguity persists in the literature. Although studies often attribute performance gains to the multi-agent architecture itself, they typically compare MAS equipped with highly engineered, role-specific prompt conditioning against generic single-agent baselines. This confound makes it difficult to de-

termine whether the observed diversity arises from intrinsic multi-agent interactions or from prompt conditioning that could be equally effective within a single-agent setup.

Defining LLM agents. Recent literature diverges in how it defines “LLM agents.” One common view characterizes agents as stateful controllers that act and reason in external environments, often incorporating memory and planning (Park et al., 2023; Yao et al., 2023b; Shinn et al., 2023; Yang et al., 2024; Zhou et al., 2024a). Another line of work defines agents as prompt-conditioned conversational roles, where distinctions arise primarily from role specifications rather than architectural design (Li et al., 2023a; Wu et al., 2024b; Liu et al., 2024b; Zhou et al., 2024b; Schick et al., 2023; Du et al., 2023).

Our goal is to disentangle agent architecture from prompt conditioning, so we adopt the former definition. When this definition conflicts with prior usage, we provide clarifying notes.

Quantifying Diversity Lexical diversity is commonly measured using n-gram overlap ratios (e.g., Distinct-N) and pairwise similarity measures such as Self-BLEU (Li et al., 2016; Zhu et al., 2018; Zhang et al., 2018; Pillutla et al., 2021). For example, Liang et al. (2024) employed Self-BLEU to compare the output diversity of multi-agent and single-agent configurations.

However, lexical metrics often fail to capture deeper semantic differences. Accordingly, recent work has increasingly adopted semantic diversity measures based on embedding spaces and contrastive representations (Reimers and Gurevych, 2019; Zhang et al., 2020; Su et al., 2022). For example, Estornell and Liu (2024) showed that incorporating embedding-based measures into the “diversity pruning” step within a MAS significantly improved performance. Motivated by this line of work, we use embedding-based semantic diversity as our primary evaluation metric and retain Self-BLEU as a secondary measure for complementary analysis.

3 Task Formulation

Most benchmarks for evaluating LLMs, including Question Answering (QA) datasets and mathematical problem sets, are designed to elicit a single correct answer. Even when models reach that answer through different reasoning trajectories, a unique ground truth imposes an inherent ceiling on out-

put diversity. As a result, when multiple methods all saturate this ceiling, comparisons among them cannot reveal meaningful differences in their capacity for semantic variation. To properly evaluate semantic diversity across multiple valid answers, we therefore require tasks that naturally admit divergent yet coherent outputs.

To this end, we construct a new dataset grounded in cognitive psychology and creativity studies, which have long examined the mechanisms of divergent thinking and open-ended problem solving. Drawing on this literature, we identify five representative task types that are known to elicit creative and associative reasoning:

1. **Impossible Situations Task** (Runco, 1999): Reasoning about implausible or paradoxical premises.
2. **Alternative Uses Task** (Guilford, 1967): Proposing unconventional uses for everyday objects.
3. **Improvement Task** (McCaffrey, 2012): Suggesting modifications to enhance common objects.
4. **Just Suppose Task** (Torrance, 1974): Exploring hypothetical or counterfactual scenarios.
5. **Bridge-the-Associative-Gap Task** (Gianotti et al., 2001): Connecting distant or conceptually unrelated ideas.

We use GPT-5 to generate 60 diverse, well-formed questions per task type, yielding a total of 300 open-ended questions. Each question is designed to admit multiple valid *solution perspectives*, enabling systematic evaluation of semantic diversity in both single-agent and multi-agent settings. We manually reviewed all questions to ensure clarity and conceptual breadth.

Perspectives Instead of Personas Prompt conditioning is often framed in terms of *personas* (e.g., "critic", "engineer", or "judge"), but we instead treat *perspectives* as the primary dimension of variation. Personas typically shape outputs through loosely specified roles that affect tone, background knowledge, or rhetorical style, introducing unstructured variability that is difficult to standardize or interpret. By contrast, we define perspectives as explicit, task-specific reasoning angles, such as emphasizing ethical tradeoffs, practical constraints, or long-term

consequences. This framing links variation to the task semantics rather than to superficial identity markers, allowing us to isolate the effects of reasoning diversity more precisely.

Diversity Metric: Vendi Score (Friedman and Dieng, 2023) To quantify semantic diversity, we adopt the *Vendi Score*, a reference-free and similarity-aware metric computed from a kernel similarity matrix K (e.g., cosine similarities between embedding vectors). Let $\{\lambda_i\}$ denote the eigenvalues of K , and define the normalized spectrum as $\tilde{\lambda}_i = \lambda_i / \sum_j \lambda_j$. The Vendi Score is defined as:

$$\text{VS}(K) = \exp\left(-\sum_i \tilde{\lambda}_i \log \tilde{\lambda}_i\right)$$

We compute semantic similarity among generated answers using three embedding models: Sentence-Transformers all-mpnet-base-v2 (768 dimensions), Qwen3-8B embedding (4096 dimensions), and OpenAI text-embedding-3-small (1536 dimensions).

Because the Vendi Score is sensitive to the number of items being compared, we control for this factor by fixing the number of *perspectives* per question. Each perspective corresponds to exactly one generated answer. Consequently, choosing $k \in \{2, 4, 8, 16\}$ jointly determines both the number of reasoning angles and the number of resulting outputs.

These perspectives are not manually written; instead, they are generated in advance via a dedicated prompting procedure (see Appendix A.2 for details). As shown in Figure 1, given a question and a target value of k , we use the same LLM as in answer generation to produce k diverse, task-relevant reasoning perspectives. This design ensures that variation across answers reflects semantically meaningful differences in perspective, while enabling scalable and consistent evaluation of diversity across experimental protocols.

As a complementary measure, we also compute Self-BLEU to evaluate lexical diversity. While our main analyses focus on the semantic variation captured by the Vendi Score, Self-BLEU provides a supplementary indicator of surface-level similarity across outputs.

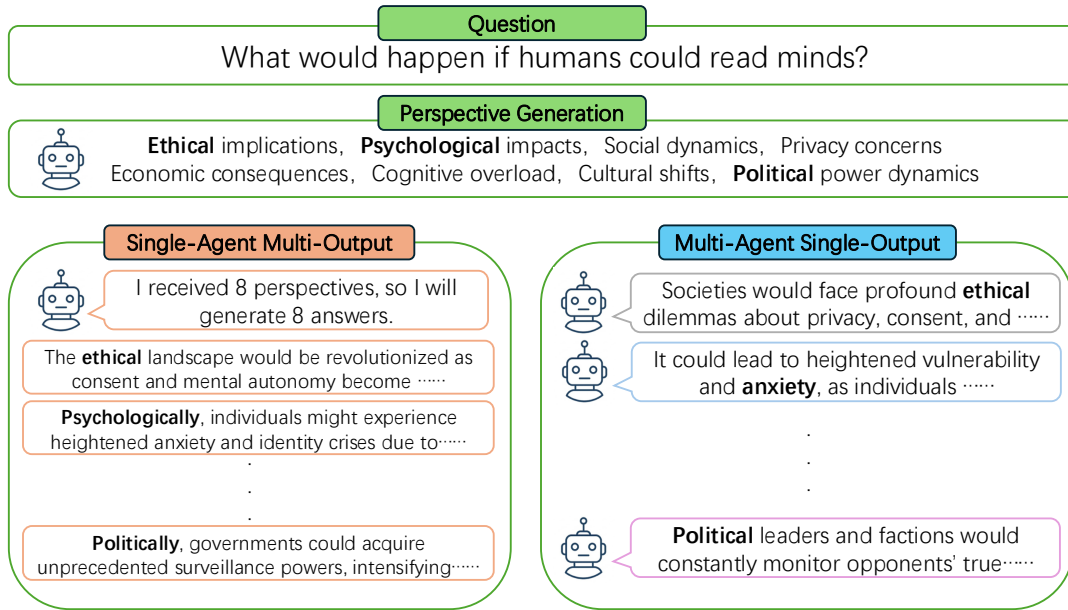


Figure 1: **Overview of the generation pipeline for single-agent and multi-agent settings.** Given an open-ended question, the system first generates k distinct reasoning perspectives (illustrated here with $k = 8$). For a fair comparison, both settings use the same prompt conditioning. In the *single-agent* setting, a single model instance receives all k perspectives simultaneously and is asked to generate one answer per perspective in a single inference pass. In the *multi-agent* setting, the k perspectives are distributed across k independent agents, and each agent is likewise asked to generate one answer per perspective.

4 Methodology

4.1 A Two-Factor Framework for Generation Protocols within generation round

To enable a rigorous comparison between the single-agent and multi-agent settings, we enforce identical prompt conditioning between them: both use the same prompt template P and the same set of k perspectives $A = \{a_1, \dots, a_k\}$. In addition, we strictly hold the total output volume constant at k answers in both settings.

To match the same output volume in a linear system, there are two strategies: the model can either switch perspectives iteratively to produce answers in a linear sequence, or ingest all perspectives at once and generate the corresponding answers in a single inference pass. Based on these operational differences, we can decouple the generation process along two orthogonal axes:

1. **Agent Architecture** (*Single-Agent vs. Multi-Agent*): This axis distinguishes the processing structure. The task is either handled by a single agent operating linearly or distributed across multiple independent agents operating in parallel. This categorization follows the

state-based definition of an agent in prior work (Park et al., 2023; Yao et al., 2023b; Shinn et al., 2023; Yang et al., 2024; Zhou et al., 2024a), which facilitates our controlled separation between agent architecture and prompt conditioning.

2. **Output Multiplicity** (*Single-Output vs. Multi-Output*): This axis specifies the decoding granularity per model invocation. *Single-Output* protocols generate exactly one answer per call, whereas *Multi-Output* protocols generate the full set of k answers in a single inference pass. In practice, we control the number of perspectives provided as input to the model, and pair this with the one answer each perspective instruction.

The intersection of these dimensions yields four canonical protocols:

- **Single-Agent Single-Output** (SA-SO)
- **Single-Agent Multi-Output** (SA-MO)
- **Multi-Agent Single-Output** (MA-SO)
- **Multi-Agent Multi-Output** (MA-MO)

275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320

4.2 Protocol Definitions

We formally define the generation process for each protocol. Let G denote the generation function for each invocation of the LLM.

Single-Agent Single-Output (SA-SO). A single agent generates answers sequentially over k discrete steps. At step i , the agent adopts a distinct perspective a_i ; completing all k steps constitutes one full generation round. Critically, each step is conditioned on the set of answers generated in earlier steps. Considering that some MAS implementations use this protocol in practice, we discuss it from both single-agent and multi-agent perspectives in our subsequent analyses.

Let $R_{<i} = \{r_1, \dots, r_{i-1}\}$ denote the partial set of answers generated prior to step i . Then:

$$r_i = G(P \cup R_{<i}, a_i), \quad i = 1, \dots, k \quad (1)$$

$$R_{\text{SA-SO}} = \{r_1, \dots, r_k\} \quad (2)$$

Single-Agent Multi-Output (SA-MO). A single agent instance is invoked once to produce one corresponding answer for each perspective, distributing its decoding capacity internally across the k outputs; this setting unambiguously constitutes a single-agent system and serves as our definitive single-agent anchor throughout the paper.

$$R_{\text{SA-MO}} = G(P, A) = \{r_1, \dots, r_k\} \quad (3)$$

Multi-Agent Single-Output (MA-SO). Each perspective a_i is assigned to an independent agent instance that operates in parallel and produces a single answer in isolation:

$$r_i = G(P, a_i), \quad i = 1, \dots, k \quad (4)$$

$$R_{\text{MA-SO}} = \{r_1, \dots, r_k\} \quad (5)$$

Multi-Agent Multi-Output (MA-MO). For completeness, we define a hybrid protocol in which K parallel agents each produce m answers, keeping the total output size fixed at $k = K \cdot m$. The perspective set A is partitioned into K disjoint subsets A_1, \dots, A_K , each of size m . Each agent is assigned subset A_i and generates the corresponding answers using a multi-output strategy:

$$S_i = G(P, A_i) = \{r_{i,1}, \dots, r_{i,m}\} \quad (6)$$

$$R_{\text{MA-MO}} = \bigcup_{i=1}^K S_i = \{r_1, \dots, r_k\} \quad (7)$$

In our experiments, MA-MO is evaluated at $k = 8$ with $m \in \{2, 4\}$ and at $k = 16$ with $m \in \{2, 4, 8\}$.

4.3 Iterative Generation Rounds

In both single-agent and multi-agent settings, we adopt the canonical iterative formulation, in which each round conditions on the full set of outputs generated in all previous rounds; in the multi-agent setting, this conditioning serves as the communication mechanism. This design corresponds to a causal all-to-all (complete-graph) communication topology across rounds: outputs from every agent in earlier rounds are visible in later rounds, and more specialized interaction topologies can be viewed as pruned variants of this structure.

Let $R^{(t)}$ denote the set of outputs generated in round t . Each subsequent round builds on the full accumulated history by updating P at the beginning of the round:

$$P^{(t)} = P^{(t-1)} \cup R^{(t-1)}$$

$$R^{(0)} = \emptyset, \quad P^{(0)} = P$$

In our experiments, we fix the generation horizon to 4 rounds. Empirically, our metrics converges within this range, while additional rounds yield diminishing returns and can occasionally degrade performance. This threshold is sufficient to capture the essential iterative dynamics while maintaining computational efficiency.

5 Results and Discussions

5.1 Single-Agent vs. Multi-Agent under Single-Output

Across all evaluated chat models, embedding models, rounds, and values of k , we observe a consistent pattern: **SA-SO exhibits substantially higher diversity than MA-SO** (Figures 2, 5, and 8, and Table 1). A similar trend is observed under the Self-BLEU metric (Figure 11).

Under a strictly controlled experimental setup, this performance gap can be attributed to a structural asymmetry between the two paradigms. In the parallel MA-SO setting, agents generate answers independently at the current round. Although each agent has a distinct perspective, the lack of real-time visibility into peer outputs precludes mutual adjustment. With our manual inspection of a subset of experimental outputs, we found that agents remain unable to fully circumvent text degeneration (Holtzman et al., 2020; Li et al., 2016; Welleck

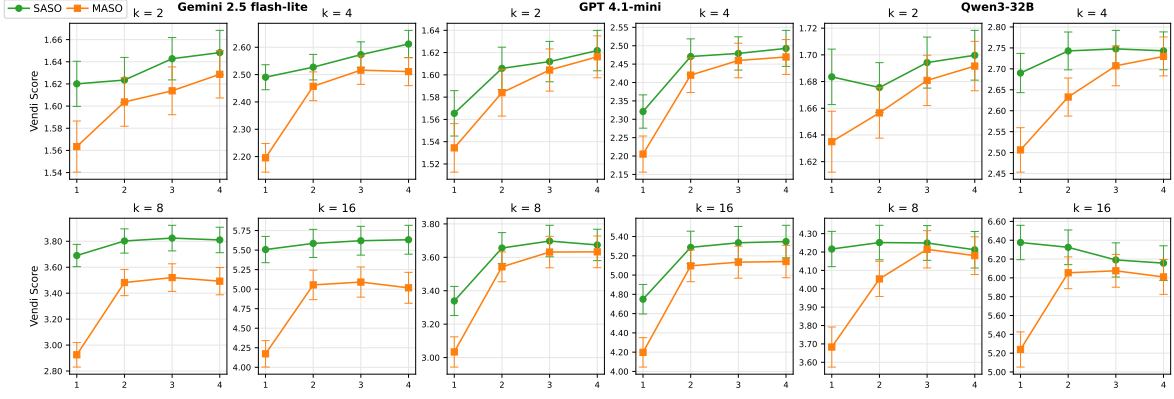


Figure 2: Comparison of Single-Agent vs. Multi-Agent architectures under Single-Output setting (mpnet-v2) The Single condition (green) represents SA-SO, and the Multi condition (orange) represents MA-SO. The x-axis indicates the generation round, and the y-axis reports the semantic diversity (Vendi score). Error bars denote 95% confidence intervals.

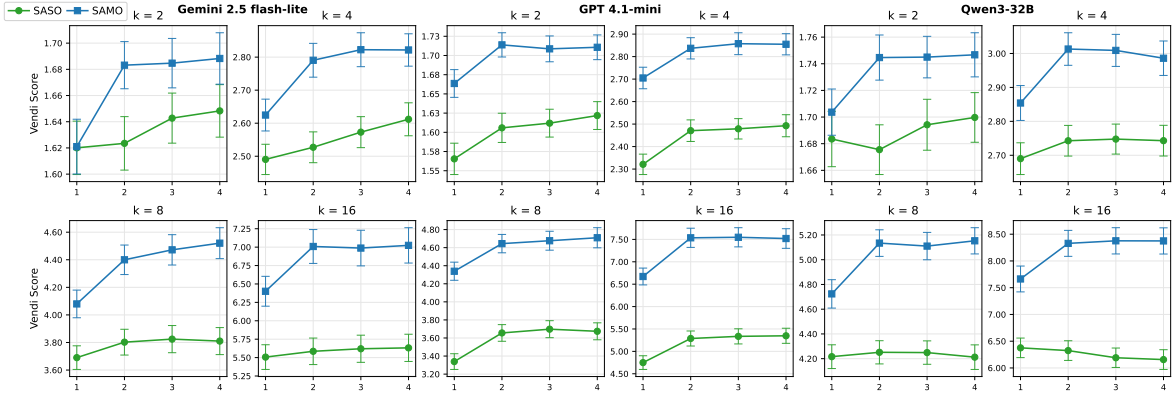


Figure 3: Comparison of Single-Output vs. Multi-Output strategies under Single-Agent architecture (mpnet-v2) The Single-output condition (green) corresponds to SA-SO, while the Multi-output condition (blue) corresponds to SA-MO. The x-axis indicates the generation round, and the y-axis reports the semantic diversity (Vendi score). Error bars denote 95% confidence intervals.

et al., 2020), often converging on similar outputs despite their assigned perspective. In contrast, the linear structure of SA-SO enables the model to condition every output on the preceding history, avoiding redundancy in real time. Thus, single-agent generation achieves higher diversity, even when operating over the identical set of perspectives. We refer to this difference as the degree of *information visibility*, defined as the extent to which an agent’s access to information produced within the system is restricted during a single inference pass.

This *information visibility* is also supported by the marked increase in MA-SO diversity from the first to the second generation round, when agents first gain access to the interaction information. The increases suggests that higher *information visibility* can substantially alleviate text degeneration. MA-SO’s *information visibility* disadvantage is particu-

larly pronounced in the first generation round, and this gap is also reflected in our experimental results, where SA-SO typically shows a substantial lead in round 1.

Taken together, our controlled comparisons indicate that agent multiplicity does not inherently increase diversity or even reduce it. MA-SO does improve over rounds as interaction information accumulates, but SA-SO maintains higher *information visibility* throughout, and correspondingly remains more diverse across the vast majority of subsequent steps. Even if we set aside our agent definition and analyze SA-SO as a form of MAS, our conclusions remain unchanged: across both MAS formulations, higher *information visibility* enables the system to generate more diverse outputs.

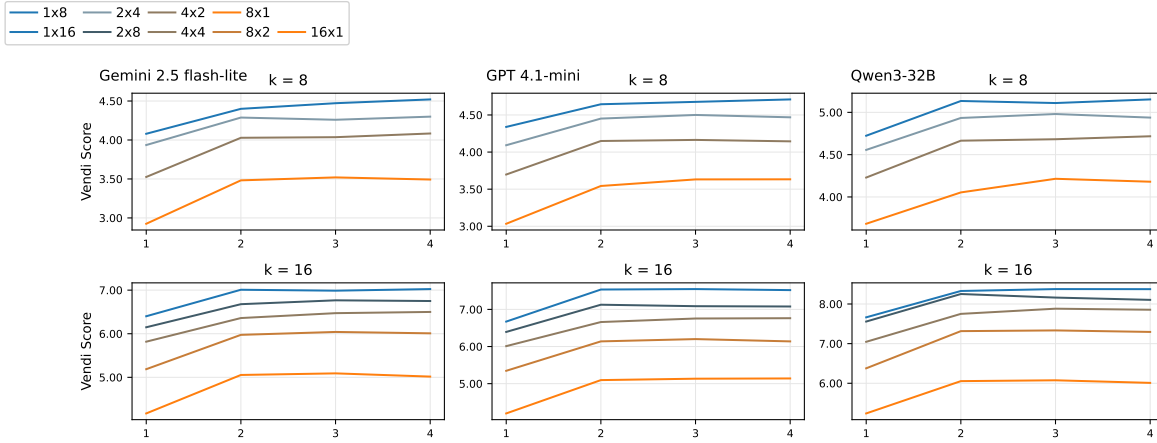


Figure 4: **Exploring multi-agent multi-output (MA-MO) configurations (mpnet-v2)** This figure accompanies the analysis in *Additional Derivation: Exploring MA-MO*. Each configuration is denoted as $n \times m$, where n is the number of agents and each agent generates m outputs, with the total number of answers fixed at $k = n \cdot m$. The special cases $1 \times k$ and $k \times 1$ correspond to the SA-MO and MA-SO settings, respectively. The x-axis indicates the generation round, and the y-axis reports semantic diversity quantified by the Vendi score.

5.2 Single-Output vs. Multi-Output under Single-Agent

Under a fixed single-agent architecture, we observe a consistent pattern across all conditions: **SA-MO yields substantially higher diversity than SA-SO** (Figures 3, 6, and 9, and Table 1). Moreover, this gap widens as k increases. A parallel trend is observed under the Self-BLEU metric (Figure 12).

Although both protocols condition later outputs on earlier content, they implement this conditioning differently. In SA-SO, earlier answers are explicitly reintroduced in the prompt across successive invocations, whereas in SA-MO they are generated within a single auto-regressive decoding stream. Auto-regressive models apply causal self-attention uniformly over preceding tokens (Radford et al., 2019; Brown et al., 2020; Touvron et al., 2023), so the placement of earlier content within a continuous sequence can affect how it is attended to; differences in positional encoding and the presence of sequence boundaries may further modulate how prior context is utilized.

While our controlled experiments isolate protocol-level differences, they can't identify the precise mechanism underlying the observed diversity gap. Beyond the model architectural considerations discussed above, we propose a data-centric hypothesis: Multi-Output instructions may mirror the human-authored enumerations that are pervasive in pretraining corpora, where list items are implicitly expected to be non-redundant and to contribute dis-

tinct aspects in accordance with the Gricean Maxim of Quantity (Grice, 1975). When asked to produce multiple answers in a single inference pass, the model may generalize this learned convention and actively vary content across the requested outputs, consistent with evidence that pretrained LLMs encode non-trivial discourse-level regularities (Koto et al., 2021; Shen et al., 2021).

Regardless of the underlying mechanism, our results support a methodological conclusion: generating multiple outputs per invocation is a simple yet effective choice for improving output diversity in agentic framework. When we set aside our agent definition and analyze SA-SO as a form of MAS, we can find that under identical prompt conditioning, the single-agent setup outperforms the multi-agent setup in diversity generation, suggested that the long-held intuition that splitting a model into multiple interacting agents inherently yields greater diversity may be misguided.

5.3 Additional Derivation: Exploring MA-MO

We now turn to the derived setting that follows naturally from our two-factor analysis: MA-MO.

As shown in Figures 4, 7, and 10, and Table 2, a consistent pattern emerges: when the total output count is held constant, for all settings, **MA-MO exhibits a largely monotonic decrease in diversity as the agent count increases**. A parallel trend is observed under the Self-BLEU metric (Figure 13).

Increasing the number of agents necessarily re-

duces the number of multi-output generations available to each agent, which in turn lowers the system’s *information visibility*. These results further corroborate our earlier conclusion that *information visibility* is a plausible driver of the observed diversity differences.

5.4 Post-hoc Quality Audit: Does Multi-Output Harm Validity?

Motivation. Although our prompts explicitly require logical consistency and close relevance to the input question, it is not clear *a priori* whether MO harms answer quality. In particular, if MO amplifies hallucinations or invalid reasoning, any gains in diversity may come at the cost of reliability. To quantify this potential trade-off, we perform an LLM-as-judge quality audit on the SA-MO and MA-SO outputs.

In open-ended generation, however, hallucination is difficult to define operationally, particularly when no clear reference is available. As a result, separating hallucinated content from answers that are merely unverifiable or intentionally creative can be ill-posed (van Deemter, 2024; Ji et al., 2023). We therefore take a conservative approach and restrict our audit to **strict logical errors**, such as explicit self-contradictions, mutually incompatible causal claims, or inconsistent definitions—*independent of factuality or stylistic fluency* (Lin et al., 2022).

LLM-based screening with targeted human verification. We use GPT-4o as a conservative logical auditor to screen approximately **210K** answers from SA-MO and MA-SO, and it flags only **74** cases (**0.03%**) as containing logical errors, of which **81.08%** come from MA-SO. Given this extremely low base rate, exhaustively annotating all outputs would be inefficient. We therefore form a 400-item evaluation set by including all 74 LLM-flagged answers and adding 326 randomly sampled unflagged ones. Two annotators independently label the set: Annotator 1 marks **63** answers as erroneous, whereas Annotator 2 marks **57**. Among the 74 LLM-flagged answers, Annotator 1 and Annotator 2 confirm errors in **43.97%** and **52.59%** of cases, respectively. Overall, only **32** items receive consensus error labels from both human annotators.

These results underscore a known challenge in evaluating open-ended outputs: even under a narrow validity criterion, human judgments can diverge because annotators apply different interpre-

tive thresholds. Nonetheless, GPT-4o aligns well on clear-cut cases, flagging **87.5%** (28 out of 32) of the errors on which both annotators agree. This indicates that although fine-grained agreement remains difficult to obtain, LLMs can still function as conservative auditors for *unambiguous* logical flaws.

Overall, we find no evidence that MO instructions measurably degrade logical validity. Given the low incidence of strict logical errors and the LLM’s high recall on cases where annotators agree, we conclude that MO’s diversity gains are not accompanied by an increase in obvious logical defects.

6 Conclusion

Increasing diversity with an agentic framework is widely viewed as a practical way to improve the reasoning and generative capabilities of LLMs. Multi-agent systems are often seen as a natural mechanism for achieving this diversity, since multiple agents can approximate human-like deliberation and multi-perspective reasoning.

In this work, we re-examined this assumption via a controlled study that isolates the effect of agent architecture on generation diversity. Under matched experimental conditions, single-agent configurations consistently and substantially outperformed their multi-agent counterparts across all settings.

Our analysis attributed the relative disadvantage of MAS to two structural factors: the lower **information visibility** induced by multi agent architecture, and the fact that **multi-output** delivers larger gains under comparable conditions, thereby widening the gap between single-agent and the traditional multi-agent settings.

These findings challenge the intuition that Multi-Agent Architecture inherently increases diversity. Unlike human groups, LLM agents may not be suitable for inter-agent discussion architecture. Instead, a single model that internally generates multiple distinct answers could effectively yield higher semantic diversity.

Beyond the empirical results, our study pointed to a broader implication: **multi-output generation itself is a key driver of diversity**. This finding points to a simple and efficient way to increase generation diversity, offering a new direction for designing high-performance agentic frameworks.

608 Limitations

609 While our framework provided a controlled com-
610 parison between single-agent and multi-agent
611 paradigms, several limitations remained to be ad-
612 dressed.

613 First, our experiments showed that Multi-Output
614 increases diversity, but the mechanism underlying
615 this effect remains hypothetical. We also do not
616 characterize the upper limit of its generative capac-
617 ity; there may exist a critical point beyond which
618 generating additional outputs yields diminishing
619 returns or even degrades performance. We leave
620 both questions for future investigation.

621 Second, for consistency and comparability, we
622 **fixed the temperature parameter to 1.0** across
623 all models. However, MAS could potentially lever-
624 age **heterogeneous temperature configurations**
625 to induce greater diversity across agents.

626 Third, in terms of **Quality Audit**, the incidence
627 of detected errors is extremely low, and the open-
628 ended, reference-free nature of our dataset makes
629 quality judgments inherently subjective. To ac-
630 commodate this property, we employed unconven-
631 tional evaluation procedures. Future work could
632 develop more principled **evaluation frameworks**
633 specifically tailored to multi-answer generation and
634 diversity-oriented reasoning.

635 References

636 Sweta Agrawal, Chunting Zhou, Mike Lewis, Luke
637 Zettlemoyer, and Ghazvininejad Marjan. 2023. In-
638 context learning with diverse examples. In *Proceed-*
639 *ings of EMNLP*.

640 Simran Arora, Avanika Narayan, Mayee F Chen, Laurel
641 Orr, Neel Gu, Christopher Re, Aaron Sidford, and
642 Frederic Sala. 2023. Ask me anything: A simple
643 strategy for prompting language models. In *Proceed-*
644 *ings of ICLR*.

645 Maciej Besta, Nils Blach, Ales Kubicek, Robert Ger-
646 stenberger, Lukas Gianinazzi, Joanna Gajda, and 1
647 others. 2024. Graph of thoughts: Solving elaborate
648 problems with large language models. In *Proceed-*
649 *ings of AAAI*.

650 Tom B Brown, Benjamin Mann, Nick Ryder, Melanie
651 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind
652 Neelakantan, and 1 others. 2020. Language models
653 are few-shot learners. In *Proceedings of NeurIPS*,
654 volume 33, pages 1877–1901.

655 Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu,
656 Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan Liu.
657 2023. Chateval: Towards better llm-based evaluat-
658 ors through multi-agent debate. In *Proceedings of*
659 *EMNLP*.

Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang,
Chenfei Yuan, Chi-Min Chan, and 1 others. 2024.
Agentverse: Facilitating multi-agent collaboration
and exploring emergent behaviors. In *Proceedings of*
ICLR.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,
Mark Chen, Heewoo Jun, Lukasz Kaiser, and 1 others.
2021. Training verifiers to solve math word problems.
arXiv preprint arXiv:2110.14168.

Roi Cohen, May Hamri, Mor Geva, and Amir Globerson.
2023. Lm vs lm: Detecting factual errors via
cross examination. In *Proceedings of EMNLP*.

Shehzaad Dhuliawala, Mojtaba Komeili, Jing Xu,
Roberta Raileanu, Xian Li, Asli Asghari, and Jason
Weston. 2024. Chain-of-verification reduces hallu-
cination in large language models. In *Findings of*
ACL.

Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenen-
baum, and Igor Mordatch. 2023. Improving factual-
ity and reasoning in language models through multia-
gent debate. In *Proceedings of ICML*.

Andrew Estornell and Yang Liu. 2024. Multi-llm de-
bate: Framework, principals, and interventions. In
Proceedings of NeurIPS.

Angela Fan, Mike Lewis, and Yann Dauphin. 2018.
Hierarchical neural story generation. In *Proceedings*
of ACL.

Dan Friedman and Adji Bousso Dieng. 2023. The vendi
score: A diversity evaluation metric for machine
learning. In *Proceedings of ICML*.

Lorena R R Gianotti, Christine Mohr, Diego Pizzagalli,
Dietrich Lehmann, and Peter Brugger. 2001. ASSO-
ciative processing and paranormal belief. *Psychiatry*
and Clinical Neurosciences, 55(6):595–603. Intro-
duces the Bridge-the-Associative-Gap (BAG) task.

H. Paul Grice. 1975. Logic and conversation. In *Syntax*
and Semantics, Vol. 3: Speech Acts. Academic Press.

Joy Paul Guilford. 1967. *The Nature of Human Intelli-*
gence. McGraw-Hill, New York.

Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang,
Shichao Pei, Nitesh V Chawla, Olaf Wiest, and Xi-
angliang Zhang. 2024. Large language model based
multi-agents: A survey of progress and challenges.
In *Proceedings of IJCAI*.

Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and
Yejin Choi. 2020. The curious case of neural text
degeneration. In *Proceedings of ICLR*.

Daphne Ippolito, Reno Kriz, Joao Sedoc, Maria Kri-
vokapic, and Chris Callison-Burch. 2019. Compari-
son of diverse decoding methods for natural language
generation. In *Proceedings of ACL*.

711	Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. <i>ACM Computing Surveys</i> , 55(12):1–38.	763
712		764
713		765
714		766
715		767
716	Wouter Kool, Herke van Hoof, and Max Welling. 2019. Stochastic beams and where to find them: The gumbel-top- k trick for sampling sequences. In <i>Proceedings of ICLR</i> .	768
717		769
718		770
719		771
720	Fajri Koto, Jey Han Lau, and Timothy Baldwin. 2021. Discourse probing of pretrained language models. In <i>Proceedings of NAACL</i> .	772
721		773
722		774
723	Guohao Li, Hasan Abed Al Kader Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023a. CAMEL: Communicative agents for “mind” exploration of large language model society. In <i>Proceedings of NeurIPS</i> .	775
724		776
725		777
726		778
727		779
728	Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2016. A diversity-promoting objective function for neural conversation models. In <i>Proceedings of NAACL</i> .	780
729		781
730		782
731		783
732	Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, Bei Chen, Jian-Guang Lou, and Weizhu Chen. 2023b. On the advance of diversity-based reasoning with large language models. In <i>Proceedings of ICLR</i> .	784
733		785
734		786
735		787
736	Yunxuan Li, Yibing Du, Jiageng Zhang, Yiqun Li, and Xinyu Hu. 2024. Improving multi-agent debate with sparse communication topology. In <i>Findings of EMNLP</i> .	788
737		789
738		790
739		791
740	Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, and 1 others. 2024. Encouraging divergent thinking in large language models through multi-agent debate. In <i>Proceedings of ICLR</i> .	792
741		793
742		794
743		795
744		796
745	Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, and 1 others. 2024. Let’s verify step by step. In <i>Proceedings of ICLR</i> .	797
746		798
747		799
748		800
749	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. TruthfulQA: Measuring how models mimic human falsehoods. In <i>Proceedings of ACL</i> , pages 3214–3252.	801
750		802
751		803
752		804
753	Tongxuan Liu, Xingyu Wang, Weizhe Huang, Wenjiang Xu, Yuting Zeng, Lei Jiang, Hailong Yang, and Jing Li. 2024a. Groupdebate: Enhancing the efficiency of multi-agent debate using group discussion. <i>arXiv preprint arXiv:2409.14051</i> .	805
754		806
755		807
756		808
757		
758	Xiao Liu, Hao Yu, Hanchen Zhang, and 1 others. 2024b. Agentbench: Evaluating llms as agents. In <i>Proceedings of ICLR</i> .	809
759		810
760		811
761		812
762		
	Luca Massarelli, Fabio Petroni, Aleksandra Piktus, Myle Ott, Tim Rocktäschel, Vassilis Plachouras, Fabrizio Silvestri, and Sebastian Riedel. 2020. How decoding strategies affect the verifiability of generated text. In <i>Findings of EMNLP</i> .	813
		814
		815
	Tony McCaffrey. 2012. Innovation relies on the obscure: A key to overcoming the classic problem of functional fixedness. <i>Psychological Science</i> , 23(3):215–218.	
	Joon Sung Park, Joseph C O’Brien, Carrie J Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In <i>Proceedings of UIST</i> .	
	Krishna Pillutla, Swabha Swayamdipta, Rowan Zellers, John Thickstun, Sean Welleck, Yejin Choi, and Zaid Harchaoui. 2021. Mauve: Measuring the gap between neural text and human text using divergence frontiers. In <i>Proceedings of NeurIPS</i> .	
	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. <i>OpenAI Technical Report</i> .	
	Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. In <i>Proceedings of EMNLP</i> .	
	Mark A Runco. 1999. Divergent thinking. In Mark A Runco and Steven R Pritzker, editors, <i>Encyclopedia of Creativity</i> , volume 1, pages 577–582. Academic Press, San Diego, CA.	
	Timo Schick, Jane Dwivedi-Yu, Roberto Dessi, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools. In <i>Proceedings of NeurIPS</i> .	
	Aili Shen, Meladel Mistica, Bahar Salehi, Hang Li, Timothy Baldwin, and Jianzhong Qi. 2021. Evaluating document coherence modeling. <i>Transactions of the Association for Computational Linguistics</i> .	
	Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: Language agents with verbal reinforcement learning. In <i>Proceedings of NeurIPS</i> .	
	Yixuan Su, Tian Lan, Yan Wang, Dani Yogatama, Lingpeng Kong, and Nigel Collier. 2022. A contrastive framework for neural text generation. In <i>Proceedings of NeurIPS</i> .	
	Derek Tam, Sachit Menon, Mohit Bansal, Shashank Srivastava, and Colin Raffel. 2023. Evaluating the robustness of large language models to prompt variations. In <i>Proceedings of EMNLP</i> .	
	Guy Tevet and Jonathan Berant. 2021. Evaluating the evaluation of diversity in natural language generation. In <i>Proceedings of EMNLP</i> .	

816	Ellis Paul Torrance. 1974. <i>Torrance Tests of Creative Thinking: Norms-Technical Manual</i> . Scholastic Testing Service, Bensenville, IL.	869
817		870
818		871
819	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, and 1 others. 2023. LLaMA: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> .	872
820		873
821		874
822		875
823		876
824		877
825	Kees van Deemter. 2024. The pitfalls of defining hallucination. <i>Computational Linguistics</i> , 50(2):807–816.	878
826		879
827	Qian Wang, Zhenheng Tang, Zichen Jiang, Nuo Chen, Tianyu Wang, and Bingsheng He. 2025a. Agenttaxo: Dissecting and benchmarking token distribution of llm multi-agent systems. In <i>ICLR 2025 Workshop on Foundation Models in the Wild</i> .	880
828		881
829		882
830		883
831		884
832	Qian Wang, Tianyu Wang, Zhenheng Tang, Qinbin Li, Nuo Chen, Jingsheng Liang, and Bingsheng He. 2025b. Megaagent: A large-scale autonomous llm-based multi-agent system without predefined sops. In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> .	885
833		886
834		887
835		888
836		889
837		890
838	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. Self-consistency improves chain of thought reasoning in language models. In <i>Proceedings of ICLR</i> .	891
839		892
840		893
841		894
842		895
843	Zhenhailong Wang, Shaoguang Mao, Wenshan Wu, Tao Ge, Furu Wei, and Heng Ji. 2024. Unleashing the emergent cognitive synergy in large language models via multi-role co-operation. In <i>Proceedings of ICLR</i> .	896
844		897
845		898
846		899
847	Sean Welleck, Ilya Kulikov, Stephen Roller, Emily Dinan, Kyunghyun Cho, and Jason Weston. 2020. Neural text generation with unlikelihood training. In <i>Proceedings of ICLR</i> .	900
848		901
849		902
850		903
851	Yixuan Weng, Minjun Zhu, Fei Xia, Bin Li, Shizhu He, Kang Liu, and Jun Zhao. 2023. Large language models are better reasoners with self-verification. <i>arXiv preprint arXiv:2212.09561</i> .	904
852		905
853		906
854		907
855	Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Peng, Toby Walsh, and Ahmed Awadallah. 2024a. Autogen: Enabling next-gen llm applications. In <i>Proceedings of ICLR</i> .	908
856		909
857		910
858		911
859	Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, Ahmed Hassan Awadallah, Ryen W. White, Doug Burger, and Chi Wang. 2024b. Autogen: Enabling next-gen llm applications via multi-agent conversations. In <i>Proceedings of COLM</i> .	912
860		913
861		914
862		915
863		916
864		917
865		918
866	Kai Xiong, Xiao Ding, Yixin Cao, Ting Liu, and Bing Qin. 2024. Examining inter-consistency of large language models collaboration. In <i>Proceedings of ICLR</i> .	919
867		920
868		921
	John Yang, Carlos E. Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan, and Ofir Press. 2024. Swe-agent: Agent-computer interfaces enable automated software engineering. In <i>Proceedings of NeurIPS</i> .	
	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik Narasimhan. 2023a. Tree of thoughts: Deliberate problem solving with large language models. In <i>Proceedings of NeurIPS</i> .	
	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023b. React: Synergizing reasoning and acting in language models. In <i>Proceedings of ICLR</i> .	
	Yuting Zeng, Weizhe Huang, Lei Jiang, Tongxuan Liu, XiTai Jin, Chen Tianying Tiana, Jing Li, and Xiaohua Xu. 2025. S ² -mad: Breaking the token barrier to enhance multi-agent debate efficiency. In <i>Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> .	
	Ceyao Zhang, Kaijie Yang, Xujian Hu, Zeyuan Wang, Guang Li, Yuhui Sun, and Cheng Zhang. 2024. Proagent: Building proactive cooperative ai with large language models. In <i>Proceedings of ICLR</i> .	
	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert. In <i>Proceedings of ICLR</i> .	
	Yizhe Zhang, Michel Galley, Jianfeng Gao, Zhe Gan, Xiujun Li, Chris Brockett, and Bill Dolan. 2018. Generating informative and diverse conversational responses via adversarial information maximization. In <i>Proceedings of NeurIPS</i> .	
	Andy Zhou, Kai Yan, Michal Shlapentokh-Rothman, Haohan Wang, and Yu-Xiong Wang. 2024a. Language agent tree search unifies reasoning, acting, and planning in language models. In <i>Proceedings of ICML</i> .	
	Shuyan Zhou, Frank F. Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan Bisk, Daniel Fried, Uri Alon, and Graham Neubig. 2024b. Webarena: A realistic web environment for building autonomous agents. In <i>Proceedings of ICLR</i> .	
	Yaoming Zhu, Sidi Lu, Lei Zheng, Jiaxian Guo, Weinan Zhang, Jun Wang, and Yong Yu. 2018. Taxygen: A benchmarking platform for text generation models. In <i>Proceedings of SIGIR</i> .	
	Mingchen Zhuge, Haozhe Liu, Francesco Faccio, David Rutherford, Gabriele Santacatterina, and 1 others. 2023. Mindstorms in natural language-based society of mind. In <i>Proceedings of NeurIPS</i> .	

922	A Appendix				
923	A.1 Question Set				
924	We constructed our question set by adapting classic				
925	creativity and divergent-thinking tasks from cog-				
926	nitve psychology. Each category contained 60				
927	prompts in the full dataset; here we list representa-				
928	tive examples.				
929	1. Alternative Uses Task Generate unusual or				
930	creative uses for common objects:				
931	• What are alternative uses for a paperclip?				
932	• What are alternative uses for a brick?				
933	• What are alternative uses for a cardboard box?				
934	• What are alternative uses for a mirror?				
935	• What are alternative uses for a backpack?				
936	2. Improvement Task Suggest ways to improve				
937	the design or function of everyday items:				
938	• How could you improve a bicycle to make it				
939	safer?				
940	• How could you improve a smartphone to last				
941	longer?				
942	• How could you improve a refrigerator to waste				
943	less food?				
944	• How could you improve a chair to be more				
945	ergonomic?				
946	• How could you improve a pen to write in all				
947	conditions?				
948	3. Just Suppose Task Explore counterfactual				
949	“what if” scenarios:				
950	• Just suppose humans could breathe underwa-				
951	ter—what would change?				
952	• Just suppose everyone could teleport in-				
953	stantly—what would change?				
954	• Just suppose language barriers disap-				
955	peared—what would change?				
956	• Just suppose no one could lie—what would				
957	change?				
958	• Just suppose the sun never set—what would				
959	change?				
		4. Impossible Situations / Consequences Task			960
		Imagine outcomes of implausible or paradoxical			961
		conditions:			962
		• What would happen if gravity stopped work-			963
		ing for one day?			964
		• What would happen if humans never needed			965
		sleep?			966
		• What would happen if money lost all value			967
		overnight?			968
		• What would happen if humans could fly like			969
		birds?			970
		• What would happen if humans could live un-			971
		derwater?			972
		5. Bridge-the-Associative-Gap Task Form con-			973
		nections between unrelated or distant concepts:			974
		• What is the connection between a cloud and a			975
		pillow?			976
		• What is the connection between a ladder and			977
		a song?			978
		• What is the connection between a violin and			979
		the wind?			980
		• What is the connection between a shoe and a			981
		story?			982
		• What is the connection between a river and a			983
		road?			984
		Dataset Summary. Each task category included			985
		60 questions (300 in total). All prompts were veri-			986
		fied for clarity and open-endedness, ensuring mul-			987
		tiiple valid yet distinct answers could be generated			988
		across perspectives.			989
		A.2 Prompt Design			990
		To ensure consistent prompting across experi-			991
		mental settings, we used three structured tem-			992
		plates corresponding to different stages of gener-			993
		ation: the Perspective Generation Prompt,			994
		Answer Generation System Prompt, and Answer			995
		Generation User Prompt (Iteration). The			996
		prompt concatenation was implemented using			997
		LangChain, and Pydantic was used to produce			998
		structured outputs.			999
		All models were decoded with a fixed tempera-			1000
		ture of 1.0. All other decoding and sampling param-			1001
		eters were left at their respective model defaults.			1002

(1) Perspective Generation Prompt

MISSION: You are a creative thinker agent specialized in generating diverse and varied perspectives toward the given question. Your core expertise is in exploring multiple perspectives, approaches, and solutions to any given question.

Your mission is to generate exactly $\{k\}$ diverse perspectives that maximize variety and avoid redundancy.

CORE INSTRUCTIONS:

1. **Maximize Diversity:** Generate exactly $\{k\}$ answer’s perspectives that are as different as possible from each other and from any previous answers.

2. **Explore Different Dimensions:** Consider various aspects like:

- Different approaches or methodologies
- Various perspectives or viewpoints
- Different scales or levels of analysis
- Alternative frameworks or paradigms
- Contrasting assumptions or premises

3. **Concise answers:** Generate perspective with just one phrase or a single word.

Question: $\{Q\}$

Generate $\{k\}$ diverse perspectives for analyzing this question.

(2) Answer Generation System Prompt

MISSION: You are a creative yet analytical thinker agent. Thinking about the problem from the following perspective:

Perspective: $\{A_i\}$

CORE INSTRUCTIONS and WORKFLOW

1. **Keep the direction firmly anchored in logic:** Every idea must clearly address the core intent of the question — stay focused and avoid drifting off-topic.

2. **Maximize Creativity:** Generate one answer for EACH perspective. Within that relevance, be as imaginative as possible — propose unconventional, cross-disciplinary, or thought-provoking ideas.

3. **Output neatly:** Express each idea concisely in one sentence.

(3) Answer Generation User Prompt (Iteration)

QUESTION: “ $\{Q\}$ ”

PREVIOUS ANSWERS:
 $\{R\}$

INSTRUCTIONS:

Carefully analyze the list of PREVIOUS ANSWERS provided and draw inspiration from them.

Ensure your answers contain no logical flaws.

Generate new answers that introduce novel ideas while avoiding redundancy with PREVIOUS ANSWERS.

Make sure every idea is clearly and logically connected to the question.

Ensure the answers array contains exactly $\{m\}$ items, in the same order as your assigned perspectives, with one sentence per item.

Algorithm 1

Require: Question Q , perspectives set A , number of agents K

Ensure: Answer set R

- 1: Divide A into K subsets $\{P_1, \dots, P_K\}$
- 2: $R \leftarrow \emptyset$
- 3: **for** $t = 0 \rightarrow 3$ **do** ▷ 1 Initial + 3 Iterations
- 4: **for all** agents $i = 1 \rightarrow K$ **in parallel do**
- 5: $R_i \leftarrow G(Q, A_i, R)$
- 6: **end for**
- 7: $R^{(t)} \leftarrow \bigcup_{i=1}^K R_i$
- 8: $R \leftarrow R \cup R^{(t)}$
- 9: **end for**
- 10: **return** R

Explanation. SA-MO, MA-SO and MA-MO settings follow the same generation algorithm and share identical prompt templates described above. The only difference lies in the number of agents K : in SA-MO, we set $K = 1$, meaning that a single model instance sequentially generates all answers. In contrast, in MA-SO and MA-MO, $K > 1$, and each agent independently operates on a distinct subset of perspectives in parallel. This ensures that any observed differences in output diversity or quality arise solely from the collaborative dynamics rather than prompt or structural variations.

1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020

Algorithm 2 About SA-SO

Require: Question Q , perspectives set $A = \{P_1, \dots, P_k\}$, number of perspectives k

Ensure: Answer set R

- 1: $R \leftarrow \emptyset$
- 2: **for** $t = 0 \rightarrow 3$ **do** ▷ 1 Initial + 3 Iterations
- 3: **for** $i = 1 \rightarrow k$ **do** ▷ iterate over A
- 4: $R_i \leftarrow G(Q, P_i, R)$
- 5: $R \leftarrow R \cup R_i$
- 6: **end for**
- 7: **end for**
- 8: **return** R

Explanation. The SA-SO algorithm also use the identical prompt templates described above.

1021
1022

A.3 Additional Results Table and Figure

1023

Model	k	Setting	Round 1				Round 2				Round 3				Round 4						
			MPNet		Qwen		MPNet		Qwen		MPNet		Qwen		MPNet		Qwen				
			Value	Δ (%)	Value	Δ (%)	Value	Δ (%)	Value	Δ (%)	Value	Δ (%)	Value	Δ (%)	Value	Δ (%)	Value	Δ (%)			
GPT 4.1-mini	2	SAMO	1.66 (0.03)	6.27	1.70 (0.01)	1.64 (0.02)	4.09	1.71 (0.02)	1.73 (0.01)	1.66 (0.01)	1.71 (0.02)	1.72 (0.01)	1.67 (0.01)	1.71 (0.02)	1.72 (0.01)	1.67 (0.01)	1.71 (0.02)	1.72 (0.01)			
		SASO	1.57 (0.03)	-0.98	1.58 (0.03)	0.01	1.57 (0.02)	-0.01	1.57 (0.02)	-0.01	1.57 (0.02)	-0.01	1.57 (0.02)	-0.01	1.57 (0.02)	-0.01	1.57 (0.02)	-0.01	1.57 (0.02)		
		MASO	1.53 (0.04)	-0.34	1.54 (0.03)	0.01	1.54 (0.03)	0.00	1.54 (0.03)	0.00	1.54 (0.03)	0.00	1.54 (0.03)	0.00	1.54 (0.03)	0.00	1.54 (0.03)	0.00	1.54 (0.03)		
	4	SAMO	2.70 (0.18)	1.85	2.85 (0.10)	2.54 (0.09)	2.84 (0.17)	2.90 (0.08)	2.86 (0.18)	2.91 (0.08)	2.86 (0.18)	2.91 (0.08)	2.86 (0.18)	2.91 (0.08)	2.86 (0.18)	2.91 (0.08)	2.86 (0.18)	2.91 (0.08)	2.86 (0.18)		
		SASO	2.32 (0.16)	-0.58	2.38 (0.14)	0.26	2.33 (0.17)	-0.05	2.32 (0.17)	-0.05	2.32 (0.17)	-0.05	2.32 (0.17)	-0.05	2.32 (0.17)	-0.05	2.32 (0.17)	-0.05	2.32 (0.17)		
		MASO	2.21 (0.19)	-0.99	2.34 (0.14)	0.53	2.23 (0.12)	-0.31	2.42 (0.17)	2.05	2.49 (0.11)	1.43	2.40 (0.09)	-0.61	2.46 (0.17)	0.77	2.48 (0.11)	0.58	2.42 (0.08)	-0.33	
	8	SAMO	4.34 (0.78)	1.67	4.47 (0.44)	3.88 (0.36)	4.64 (0.80)	4.79 (0.41)	4.79 (0.41)	4.68 (0.85)	4.76 (0.42)	4.73 (0.35)	4.71 (0.96)	4.73 (0.35)	4.71 (0.96)	4.73 (0.35)	4.71 (0.96)	4.73 (0.35)	4.71 (0.96)		
		SASO	3.34 (0.59)	-0.98	3.40 (0.47)	0.16	3.33 (0.28)	-0.17	3.36 (0.60)	0.03	3.31 (0.46)	-0.05	3.34 (0.28)	-0.05	3.36 (0.60)	0.03	3.31 (0.46)	-0.05	3.34 (0.28)	-0.05	
		MASO	3.03 (0.63)	-0.31	3.28 (0.52)	0.25	3.04 (0.39)	-0.27	3.54 (0.62)	0.50	3.63 (0.43)	0.09	3.61 (0.47)	-0.02	3.56 (0.49)	-0.05	3.48 (0.29)	-0.12	3.56 (0.49)	-0.05	
	16	SAMO	6.67 (2.55)	2.00	7.00 (1.48)	5.62 (1.26)	7.24 (1.32)	42.52	7.23 (1.32)	42.52	7.23 (1.32)	42.52	7.23 (1.32)	42.52	7.23 (1.32)	42.52	7.23 (1.32)	42.52	7.23 (1.32)		
		SASO	4.25 (1.70)	-0.48	4.51 (1.39)	0.16	4.64 (0.82)	21.03	5.29 (2.03)	42.52	5.31 (1.37)	45.58	4.97 (0.85)	25.47	5.34 (2.06)	41.46	5.20 (1.29)	46.23	4.93 (0.88)	27.14	
		MASO	4.20 (1.69)	-0.11	4.63 (1.50)	0.35	4.11 (1.04)	-0.18	5.10 (1.96)	3.62	5.23 (1.36)	-0.42	4.88 (0.90)	-0.81	5.13 (2.00)	3.78	5.07 (1.32)	-0.44	4.79 (0.82)	-0.76	
Gemini 2.5 Flash-lite	2	SAMO	1.62 (0.03)	0.06	1.67 (0.02)	0.14	1.60 (0.02)	-0.62	1.62 (0.03)	3.68	1.67 (0.02)	2.33	1.62 (0.02)	0.58	1.64 (0.03)	2.55	1.67 (0.02)	2.68	1.65 (0.02)	0.98	
		SASO	1.56 (0.04)	-0.30	1.60 (0.03)	0.24	1.53 (0.03)	-0.40	1.60 (0.04)	-0.22	1.64 (0.02)	-0.12	1.60 (0.03)	-0.16	1.65 (0.02)	-0.12	1.59 (0.02)	-0.22	1.63 (0.03)	-0.13	
		MASO	1.52 (0.03)	-0.10	1.57 (0.03)	0.05	1.53 (0.03)	-0.04	1.57 (0.03)	0.00	1.53 (0.03)	-0.04	1.57 (0.03)	0.00	1.53 (0.03)	-0.04	1.57 (0.03)	0.00	1.53 (0.03)	-0.04	
	4	SAMO	2.62 (0.18)	1.18	2.78 (0.10)	2.47 (0.10)	2.79 (0.20)	2.89 (0.11)	2.88 (0.11)	2.88 (0.11)	2.88 (0.11)	2.88 (0.11)	2.88 (0.11)	2.88 (0.11)	2.88 (0.11)	2.88 (0.11)	2.88 (0.11)	2.88 (0.11)	2.88 (0.11)	2.88 (0.11)	
		SASO	2.49 (0.16)	-0.53	2.45 (0.09)	0.16	2.43 (0.08)	1.62	2.53 (0.17)	10.42	2.64 (0.10)	9.46	2.48 (0.09)	4.11	2.57 (0.17)	9.70	2.67 (0.11)	9.07	2.51 (0.10)	4.18	
		MASO	2.20 (0.21)	-1.18	2.31 (0.17)	-0.69	2.13 (0.16)	-1.21	2.46 (0.21)	-2.77	2.56 (0.15)	-3.02	2.37 (0.14)	-4.38	2.57 (0.21)	-2.20	2.40 (0.14)	-2.76	2.43 (0.14)	-3.28	
	8	SAMO	4.88 (0.78)	1.48	5.09 (0.69)	4.88 (0.69)	4.88 (0.69)	4.88 (0.69)	4.88 (0.69)	4.88 (0.69)	4.88 (0.69)	4.88 (0.69)	4.88 (0.69)	4.88 (0.69)	4.88 (0.69)	4.88 (0.69)	4.88 (0.69)	4.88 (0.69)	4.88 (0.69)	4.88 (0.69)	
		SASO	3.69 (0.57)	-1.06	4.01 (0.36)	1.17	3.54 (0.33)	4.26	3.80 (0.68)	15.71	4.03 (0.43)	16.13	3.63 (0.38)	7.62	3.82 (0.75)	16.93	4.02 (0.47)	17.40	3.66 (0.43)	8.39	
		MASO	2.92 (0.69)	-0.74	3.18 (0.61)	20.85	2.80 (0.51)	20.87	3.48 (0.79)	8.42	3.72 (0.57)	7.50	3.31 (0.49)	9.01	3.52 (0.88)	7.96	3.74 (0.60)	7.11	3.34 (0.53)	8.92	
	16	SAMO	6.48 (1.11)	2.23	6.81 (1.06)	5.40 (1.50)	7.01 (1.01)	7.52 (1.34)	5.85 (1.90)	6.99 (1.43)	6.99 (1.43)	6.99 (1.43)	6.99 (1.43)	6.99 (1.43)	6.99 (1.43)	6.99 (1.43)	6.99 (1.43)	6.99 (1.43)	6.99 (1.43)	6.99 (1.43)	
		SASO	5.51 (1.13)	16.23	6.06 (1.42)	19.38	5.17 (1.18)	4.51	5.59 (2.45)	28.49	6.00 (1.59)	25.24	5.23 (1.13)	11.84	5.62 (2.57)	24.33	6.00 (1.76)	24.74	5.26 (1.44)	11.06	
		MASO	4.17 (1.17)	-24.24	4.62 (1.98)	-23.76	3.89 (1.51)	-24.67	5.05 (2.70)	-9.51	5.43 (2.09)	-9.51	5.43 (2.09)	-9.51	5.09 (2.82)	-9.41	5.41 (2.16)	-9.90	4.64 (1.64)	-11.42	
Qwen3-32B	2	SAMO	1.68 (0.03)	1.20	1.70 (0.02)	0.82	1.65 (0.02)	0.32	1.68 (0.03)	4.13	1.68 (0.02)	3.30	1.64 (0.02)	2.44	1.69 (0.03)	3.00	1.70 (0.02)	2.57	1.65 (0.02)	1.55	
		SASO	1.59 (0.04)	-0.89	1.61 (0.03)	-0.27	1.62 (0.03)	0.01	1.61 (0.03)	0.01	1.61 (0.03)	0.01	1.61 (0.03)	0.01	1.61 (0.03)	0.01	1.61 (0.03)	0.01	1.61 (0.03)	0.01	1.61 (0.03)
		MASO	1.52 (0.03)	-0.16	1.56 (0.03)	0.04	1.53 (0.03)	-0.01	1.53 (0.03)	-0.01	1.53 (0.03)	-0.01	1.53 (0.03)	-0.01	1.53 (0.03)	-0.01	1.53 (0.03)	-0.01	1.53 (0.03)	-0.01	1.53 (0.03)
	4	SAMO	1.68 (0.03)	1.20	1.70 (0.02)	0.82	1.65 (0.02)	0.32	1.68 (0.03)	4.13	1.68 (0.02)	3.30	1.64 (0.02)	2.44	1.69 (0.03)	3.00	1.70 (0.02)	2.57	1.65 (0.02)	1.55	
		SASO	1.59 (0.04)	-0.89	1.61 (0.03)	-0.27	1.62 (0.03)	0.01	1.61 (0.03)	0.01	1.61 (0.03)	0.01	1.61 (0.03)	0.01	1.61 (0.03)	0.01	1.61 (0.03)	0.01	1.61 (0.03)	0.01	1.61 (0.03)
		MASO	1.52 (0.03)	-0.16	1.56 (0.03)	0.04	1.53 (0.03)	-0.01	1.53 (0.03)	-0.01	1.53 (0.03)	-0.01	1.53 (0.03)	-0.01	1.53 (0.03)	-0.01	1.53 (0.03)	-0.01	1.53 (0.03)	-0.01	1.53 (0.03)
	8	SAMO	4.27 (1.00)	1.20	4.47 (0.55)	4.17 (0.55)	4.17 (0.55)	4.17 (0.55)	4.17 (0.55)	4.17 (0.55)	4.17 (0.55)	4.17 (0.55)	4.17 (0.55)	4.17 (0.55)	4.17 (0.55)	4.17 (0.55)	4.17 (0.55)	4.17 (0.55)	4.17 (0.55)	4.17 (0.55)	
		SASO	4.22 (1.00)	12.03	4.35 (0.53)	11.65	3.86 (0.45)	8.17	4.23 (0.57)	20.77	4.26 (0.50)	19.57	3.83 (0.42)	13.57	4.25 (0.68)	20.26	4.19 (0.59)	19.55	3.86 (0.42)	14.41	
		MASO	3.68 (0.91)	-12.65	3.82 (0.68)	-12.15	3.48 (0.57)	-9.71	4.05 (0.70)	-4.65	4.11 (0.64)	-3.52	3.71 (0.40)	-3.06	4.21 (0.78)	-0.81	3.85 (0.64)	-0.42	4.15 (0.53)	0.31	
	16	SAMO	6.82 (1.43)	20.19	6.81 (1.47)	13.24	6.82 (1.44)	30.68	6.82 (1.44)	30.68	6.82 (1.44)	30.68	6.82 (1.44)	30.68	6.82 (1.44)	30.68	6.82 (1.44)	30.68	6.82 (1.44)	30.68	
		SASO	5.24 (1.55)	-17.83	5.50 (1.92)	-16.76	4.85 (1.47)	-13.45	6.05 (2.08)	-4.28	6.09 (1.43)	-3.04	5.37 (1.10)	-3.70	6.08 (2.22)	-1.87	5.99 (1.57)	-1.55	5.38 (1.22)	-1.89	
		MASO	5.24 (1.55)	-17.83	5.50 (1.92)	-16.76	4.85 (1.47)	-13.45	6.05 (2.08)	-4.28	6.09 (1.43)	-3.04	5.37 (1.10)	-3.70	6.08 (2.22)	-1.87	5.99 (1.57)	-1.55	5.38 (1.22)	-1.89	

Table 1: For each chat model, embedding model and k , the table reports the performance of three conditions (SAMO, SASO, and MASO). Each “Value” entry is the mean Vendi score for that condition and round, with the value in parentheses giving the corresponding variance across items. For each round, the two Δ (%) columns show relative percentage changes: the Δ above is SAMO versus SASO, and the Δ below is MASO versus SASO, both computed as $(\text{condition} - \text{SASO}) / \text{SASO} \times 100$.

Model	k	Setting	Round 1				Round 2				Round 3				Round 4						
			BER		Qwen-emb		BER		Qwen-emb		BER		Qwen-emb		BER		Qwen-emb				
			Value	Δ (%)	Value	Δ (%)	Value	Δ (%)	Value	Δ (%)	Value	Δ (%)	Value	Δ (%)	Value	Δ (%)	Value	Δ (%)			
GPT 4.1-mini	1x8	SAMO	1.34 (0.38)	-5.69	1.47 (0.66)	5.70	1.64 (0.89)	4.16	1.49 (0.66)	4.08	1.68 (0.92)	3.77	1.47 (0.65)	4.07	1.63 (0.89)	1.59	1.47 (0.68)	4.36	1.63 (0.89)	4.36	
		SASO	1.30 (0.38)	-9.64	1.41 (0.65)	0.78	1.37 (0.57)	-3.58	1.44 (0.89)	4.16	1.49 (0.66)	3.77	1.47 (0.65)	4.07	1.63 (0.89)	1.59	1.47 (0.68)	4.36	1.63 (0.89)	4.36	
		MASO	1.30 (0.38)	-9.64	1.41 (0.65)	0.78	1.37 (0.57)	-3.58	1.44 (0.89)	4.16	1.49 (0.66)	3.77	1.47 (0.65)	4.07	1.63 (0.89)	1.59	1.47 (0.68)	4.36	1.63 (0.89)	4.36	
	2x4	SAMO	6.67 (1.60)	-4.18	7.01 (1.29)	3.19	5.62 (1.12)	-2.27	7.54 (1.82)	7.64 (1.31)	6.23 (1.10)	7.55 (1.82)	6.11	7.21 (1.18)	3.92	6.62 (1.18)	-3.60	7.52 (1.82)	7.39 (1.12)	6.25 (1.16)	
		SASO	6.30 (1.58)	-6.64	6.64 (1.14)	-6.18	5.30 (1.03)	-3.46	6.66 (1.58)	-6.55	6.92 (1.13)	-5.90	5.89 (1.04)	-2.69	6.75 (1.60)	-4.68	6.85 (1.12)	-4.62	6.76 (1.62)	-4.46	
		MASO	6.01 (1.40)	-11.05	5.89 (1.15)	-11.28	4.92 (1.03)	-7.08	6.14 (1.50)	-7.83	6.34 (1.12)	-8.33	5.65 (1.02)	-3.98	6.20 (1.53)	-8.21	6.21 (1.11)	-9.46	5.65 (1.02)	-5.17	
	4x4	SAMO	11.16 (2.30)	-21.48	12.64 (2.22)	21.32	11.14 (1.02)	-16.49	13.60 (2.40)	16.99	12.34 (1.17)	17.36	14.88 (0.95)	13.70	15.13 (1.41)	17.20	13.88 (1.65)	18.23	14.99 (0.91)	15.22	15.14 (1.43)
		SASO	10.88 (0.88)	-4.48 (0.62)	3.00 (0.63)	4.48 (0.76)	4.48 (0.76)	3.01 (0.66)	4.48 (0.76)	3.01 (0.66)	4.48 (0.76)	3.01 (0.66)	4.48 (0.76)	3.01 (0.66)	4.48 (0.76)	3.01 (0.66)	4.48 (0.76)	3.01 (0.66)	4.48 (0.76)	3.01 (0.66)	4.48 (0.76)
		MASO	10.88 (0.88)	-4.48 (0.62)	3.00 (0.63)	4.48 (0.76)	4.48 (0.76)	3.01 (0.66)	4.48 (0.76)	3.01 (0.66)	4.48 (0.76)	3.01 (0.66)	4.48 (0.76)	3.01 (0.66)	4.48 (0.76)	3.					

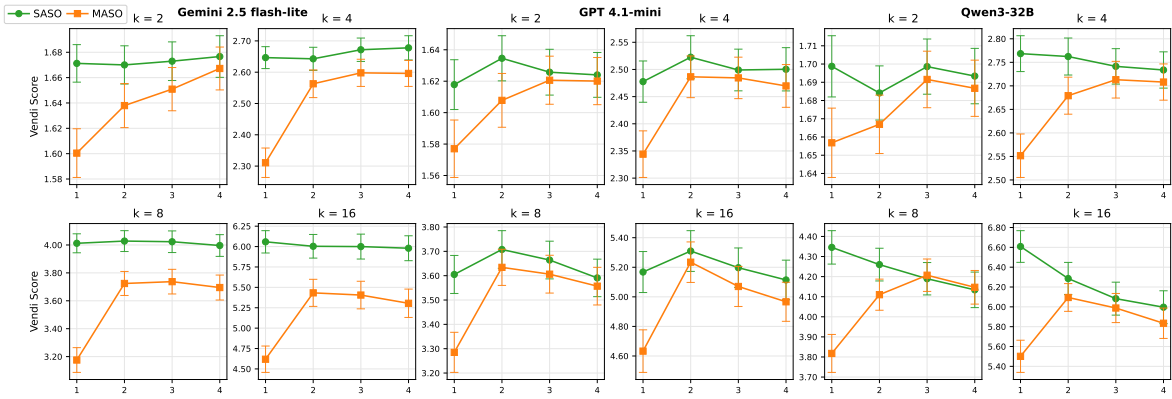


Figure 5: Single-agent vs. multi-agent generation under single-output settings with text-embedding-3-small

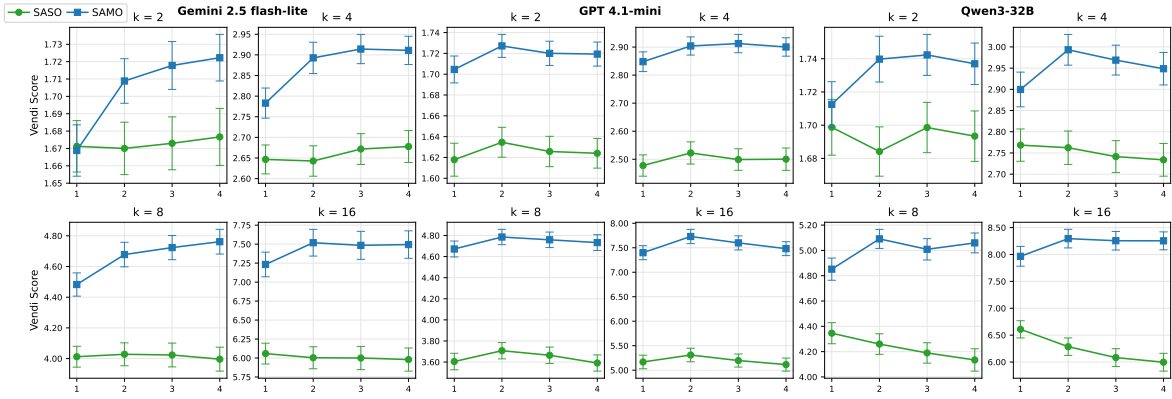


Figure 6: Single-output vs. multi-output generation under single-agent settings with text-embedding-3-small

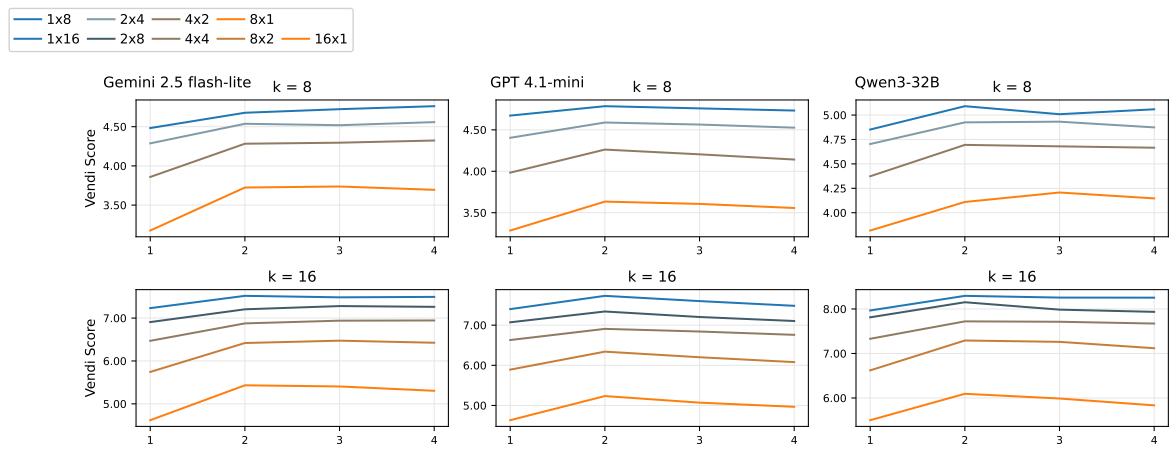


Figure 7: Exploring multi-agent multi-output (MA-MO) configurations with text-embedding-3-small

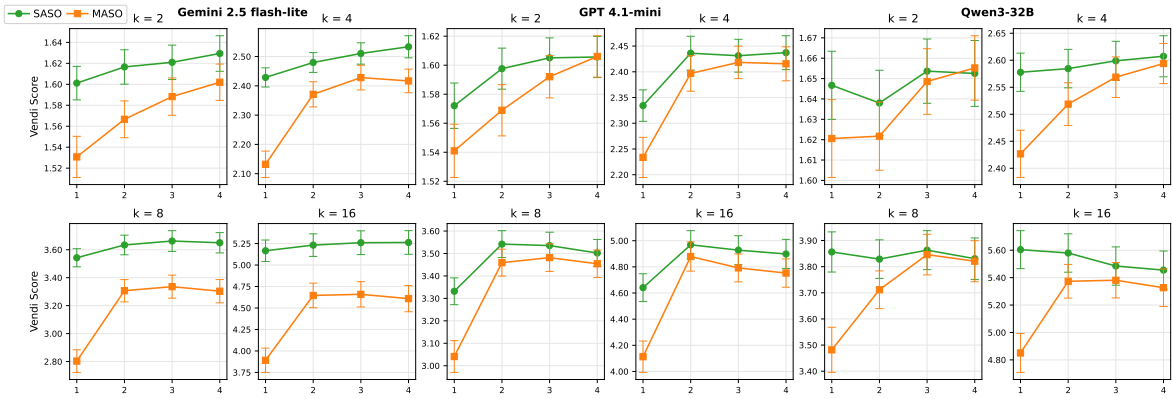


Figure 8: Single-agent vs. multi-agent generation under single-output settings with Qwen3 8b-embedding

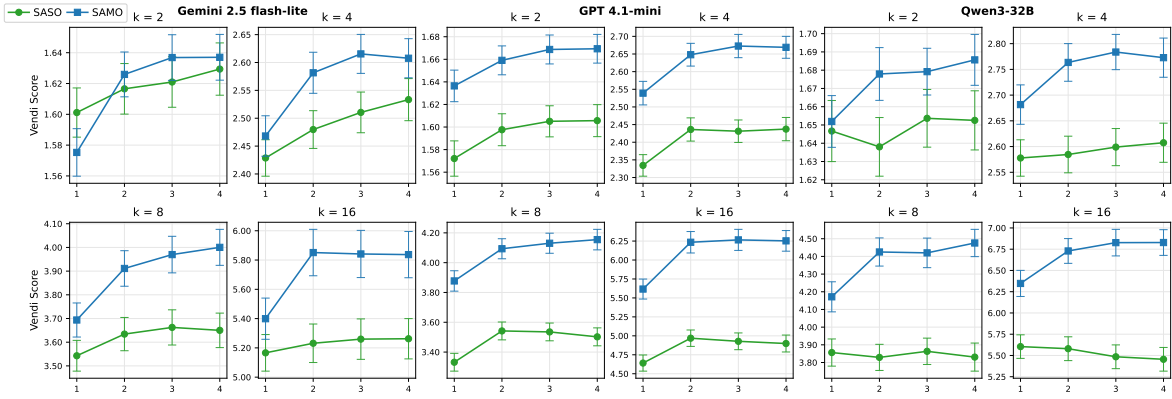


Figure 9: Single-output vs. multi-output generation under single-agent settings with Qwen3 8b-embedding

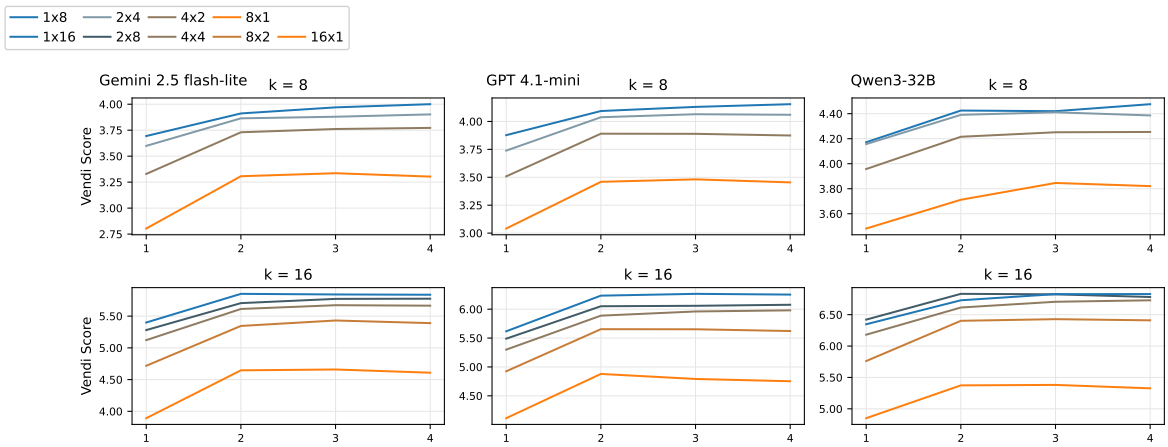


Figure 10: Exploring multi-agent multi-output (MA-MO) configurations with Qwen3 8b-embedding

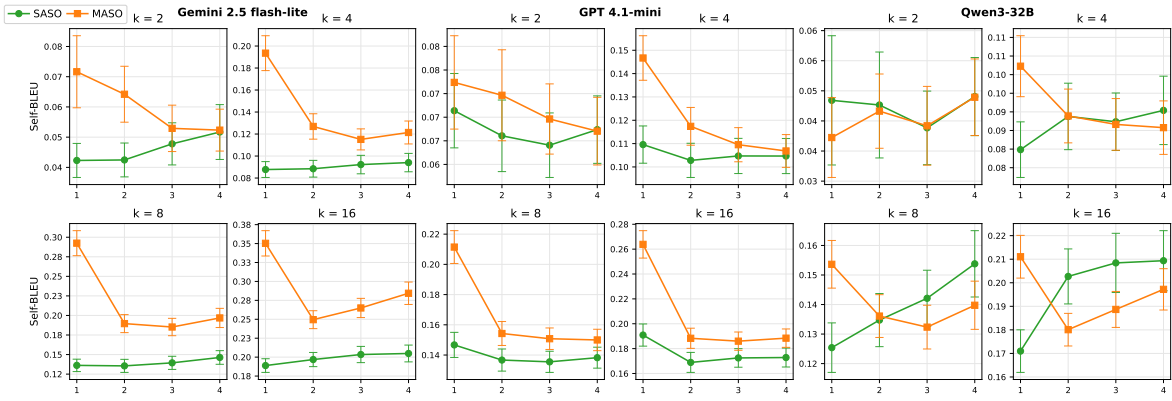


Figure 11: Single-agent vs. multi-agent generation under single-output settings with Self-BLEU

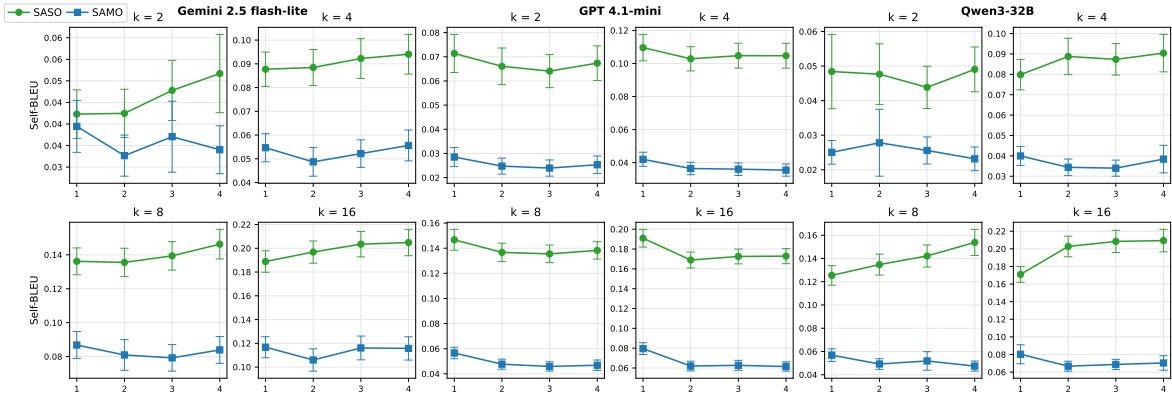


Figure 12: Single-output vs. multi-output generation under single-agent settings with Self-BLEU

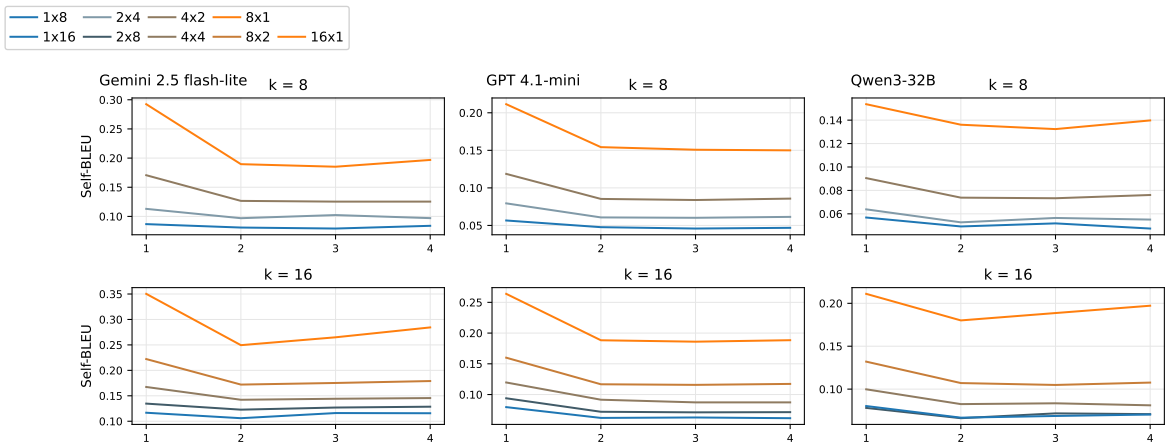


Figure 13: Exploring multi-agent multi-output (MA-MO) configurations with Self-BLEU