

# Can We Elicit Capability Uplifts by Distributing Test-Time Compute Across Interacting Agentic Nodes?

Shaheen Ahmed-Chowdhury<sup>1</sup>

<sup>1</sup>Arcadia Impact, shaheen.ahmed@live.co.uk

March 31, 2025

## Abstract

In this report, we explore whether, when normalising for tokens, spending test-time compute across a multi-agent system can elicit capability uplifts over a single-agent baseline. We first motivate the study of emergent capability uplifts within multi-agent systems, before discussing the building of a multi-agent debate tool within Inspect[2], which a single agent can invoke. We use this tool to study the incremental performance uplift on SWE-Bench tasks, with multi-agent systems using instances of GPT-4o and Claude 3.5 Haiku. We find that we can achieve 3-5% uplifts beyond single-agent baselines, but only explore a narrow range of the available debate tool hyperparameters. We conclude by outlining some directions for future research, including the automation of both log analysis and multi-agent system architecture iteration via language models.

## 1 Introduction

On February 24th, 2025, the Cooperative AI Foundation released a report on the risks of multi-agent AI systems[4]. Within it, they define three failure modes (miscoordination, conflict and collusion). They identify seven factors which could lead to these failures, one of which is emergent agency. That is, the ability of multiple interacting agentic entities to exhibit behaviours not observed at the agent level. One instance of emergent agency is emergent capabilities; the ability for multiple interacting agents to produce capabilities not observed at the individual agent level.

This report looks to work towards understanding these uplifts better, and begins by reviewing existing capability uplift literature within one narrowly-scoped capability; software engineering ability. We approximately taxonomise the methods through which multi-agent systems (MAS<sup>1</sup>) are utilised in software engineering tasks, before conducting experiments to

explore one such MAS architecture: multi-agent debate. We create a canonical and highly configurable multi-agent debate tool using the UK AI Security Institute’s (AISI’s) Inspect framework, alongside Autogen and OpenRouter, to examine the incremental capability uplift a single agent can achieve from using such a tool (on SWE-Bench tasks). We observe performance uplifts using this tool, but acknowledge that larger scale experiments are required to reduce noise in performance. We also create suitable single-agent baselines by normalising across tokens throughout. Finally, we conclude by briefly investigating why these uplifts occur, before discussing future directions and concluding.

This work was completed as part of a five-week research project under AISI’s evaluations and scaffolding bounty program[1].

## 2 Related Work

We identified two primary methods through which authors attempt to elicit capability uplifts in software engineering evaluations such as SWE-Bench; sub-tasking and debate. We also observed that MAS architectures can either stay static or adapt through the course of an evaluation. As a result, this brief literature review contains two sections:

- sub-tasking (static and adaptive);
- multi-agent debate systems (static and adaptive).

### 2.1 Sub-Tasking

In this section, we summarise the approaches of a few authors who employ sub-tasking via MAS architectures to enhance capabilities in software engineering tasks. We hope this helps readers better visualise what such task decompositions typically look like.

AgentCoder[7] uses three agents: a programmer agent to generate/refine code, a test designer agent

to create tests without seeing the programmer agent’s code and a test executor agent which runs the generated code against these tests, relaying feedback to the programmer agent. This multi-agent setup outperforms single-agent baselines on HumanEval and MBPP.

AutoCodeRover[8] uses a simpler architecture: a context-retrieval agent which calls a set of specialised search APIs to search through AST-level entities efficiently, and a patch-generation agent to produce actual code changes.

Many more sub-tasking software engineering systems exist, but we hope these provide quick intuitions into what they may look like, and why they may be able to produce capability uplifts in software engineering tasks.

The complexity of these MAS architectures can quickly increase, once we allow an orchestrator agent to adjust the architecture, usually with feedback from the development environment during an evaluation. MDAgents[10] is one example of this, and was selected for an oral presentation at NeurIPS 2024. The paper uses a moderator (or general practitioner) agent to classify a medical query’s complexity, before a recruiter agent decides which specialist agent team to assign to the problem. The specialist agents within the selected team (e.g. neurologist, radiologist, oncologist) then collaborate by exchanging messages, discussing conflicting opinions and producing a final report. This system is adaptive in the sense that the number of agent nodes used is determined by agents in the system itself.

EvoMAC represents the far end of adaptive sub-tasking MAS complexity, and was selected for an ICLR 2025 Poster session. Three teams are formed: a coding team, which uses sub-tasking and a directed-acyclic graph (DAG) to generate code via a “forward pass” through all agents; a testing team which produces tests (again via a DAG), taking generated code and “forward-passing” through the testing agents; an updating team, consisting of a gradient agent and updating agent, the first of which looks to interpret logs from the testing team, and produce a “gradient” to update the MAS on. The updating agent then modifies the architecture of the coding team’s DAG, by adding/remove agents, or adjusting prompts.

We hope the above four examples serve as illustrative MAS architectures, spanning the range of complexity for static and adaptive sub-tasking multi-agent systems.

## 2.2 Multi-Agent Debate

An alternative way of utilising multi-agent systems is to have them debate each other on a task, rather than decomposing the task into sub-tasks for agents to solve. A question that immediately arises is the topology of such a debating system. Questions such as which agent should debate with which, how many agents a single agent should debate with, how many rounds of debate there should be, and which underlying language models should be used for each agent arise, and the answers to these questions result in different MAS debate architectures.

More Agents is All You Need[11] for instance, uses a very simple debate architecture, whereby each agent uses the same underlying LLM, produces an answer, after which each agent votes on the best answer. The authors observe that (in this setup) ensemble size and temperature seem to positively correlate with performance.

From here, more complex debate architectures are available, such as that of Mixture-of-Agents[12]. Here, agents are split into proposers and aggregators. Proposers are prompted for an answer in parallel, before an aggregator combines these answers into one. The original question is then prepended to the aggregator’s answer, and this is fed back to the proposer agents in parallel. The authors mention that hyperparameters such as model diversity and the number of proposers can have a significant effect on performance.

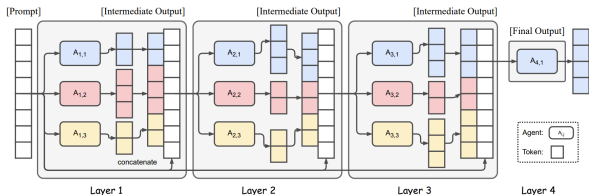


Figure 1: A diagram from [12], illustrating the Mixture-of-Agents architecture.

We hope this set of examples, spanning sub-tasking, static, adaptive and debating multi-agent systems helps the reader quickly build some intuitions around how authors are building such systems to elicit capability uplifts.

## 3 Experiments

### 3.1 Experimental Setup

Given its potentially reduced complexity, and the time available (around five weeks), we chose to experiment only with multi-agent debate methods within this

project. We ran all evaluations with AISI’s Inspect framework[2], and tested our methods against SWE-Bench[9], a popular software engineering evaluation. We used the SWE-Bench implementation available in Inspect’s partner repository, `inspect_evals` [3]. We implemented our multi-agent systems with the Autogen[13] framework.

In order to be able to measure the incremental performance uplift from spending test-time compute across multiple agent nodes, we opted to implement our multi-agent system as an Inspect Tool, which a single agent is prompted to invoke at the start of each SWE-Bench task (see Figure 2) The multi-agent tool then returns an aggregated response (in this case, simply a final answer from each agent in the MAS, with no summarisation). The single agent receives this response via the output of its multi-agent tool, before continuing on its solution trajectory. We were also provided with a set of SWE-Agent-style tools by AISI, all of which were passed as Inspect Tools to the single agent completing the problem, while a non-destructive subset were passed as Autogen Tools to all agents in the multi-agent debate tool. The complete set of tools provided to the single agent were: `find_file()`, `open_file()`, `search_dir()`, `search_file()`, `scroll_down()`, `scroll_up()`, `edit_file()` and `create_new_file()`. The last two weren’t provided to agents in the multi-agent debate tool. All agents were also provided with a `bash()` tool, and the agents in the debate tool were prompted not to make any changes to the codebase using this tool.

### 3.2 Multi-Agent Debate Tool

When designing the multi-agent tool, we wanted to build a system which would represent as canonical a debate tool as possible. That is, it should be highly configurable and represent an obvious starting point for multi-agent debate investigations. It was highlighted to us that multi-agent research at AISI is still in its infancy, and so it was important to iterate from this point. We believe the multi-agent debate tool which has been developed in this project satisfies this requirement for an initial and incremental contribution<sup>1</sup>.

Given the above, we chose to build a multi-agent debate tool containing an Aggregator agent and four Consultant agents, and the architecture of this multi-agent system can be seen in Figure 3. The Aggregator agent doesn’t contain an LLM, and simply sends out and receives messages. Each Consultant agent is pro-

<sup>1</sup>The code for this project can be found here: [github.com/shaheenahmedc/MAESTRO](https://github.com/shaheenahmedc/MAESTRO).

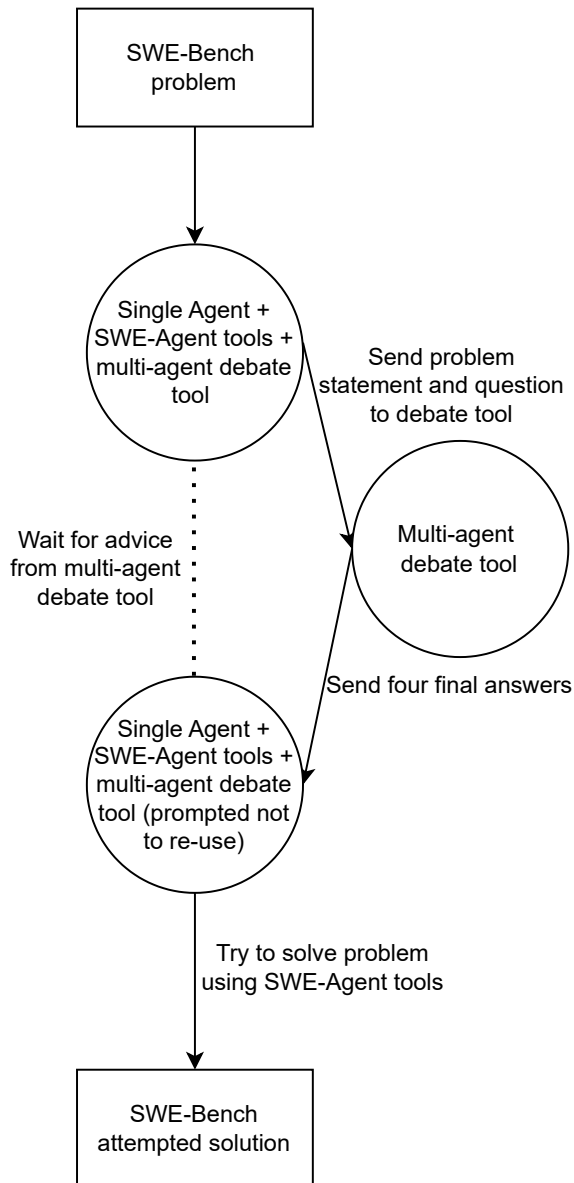


Figure 2: An illustration of how a single agent invokes a provided multi-agent debate tool, to try and help it solve a SWE-Bench problem. Note that all agents also have a `bash()` tool, which has been omitted from the diagram for brevity.

vided the original problem statement, and is able to use the aforementioned set of tools a configurable number of times, or ‘think’ by returning text instead of a tool call. We label this configurable number the reflection depth of the multi-agent system. Once an agent has reached a final answer, it is prompted to return this in a format which we search for after each reflection step. If the Consultant agent doesn’t reach a final answer by the end of its reflection depth, it’s prompted to return one. Once we have the final answers for all four consultant agents, we pass each final answer to two of the Consultant agent’s neighbours. Each agent is then informed that it has received answers from two other agents, and it should adjust its answers accordingly (see here for the Consultant system prompt). This represents the completion of one ‘round’ in our multi-agent debate tool. The agent sees the entirety of its previous reflection loop, its first final answer, and the two additional answers, before starting a second reflection loop. This process is repeated for three rounds, until we have a final answer from each of the four consultant agents. The aggregator agent then returns these answers directly to the single agent solving the SWE-Bench task.

Such a multi-agent debate tool may seem complex, but we believe it represents a fair attempt at a canonical multi-agent debate tool. The messaging topology, number of agents and reflection depth are all entirely configurable, and the Consultant agents each have their own private context to append to. Using Auto-gen’s lower-level Core API allowed us to implement the latter, in addition to increasing the general configurability of this tool, and we hope it may form the basis of an internal multi-agent debate tool at AISI.

### 3.3 Results

The above multi-agent debate tool introduces a high number of dimensions to explore: number of agents; LLMs used across the agents; LLM temperature; reflection depth; messaging topology; number of rounds. Due to time constraints, we were only able to test a small number of these, and opted for varying the LLMs used across the agents. We opted for this as it tied in with one of our primary hypotheses as to why multi-agent debate may elicit capability uplifts. That is, by having diversity in the set of base model families used, we believed that this would increase the diversity of the initial solution trajectories in the debate tool, which would ideally converge on an answer which had explored more of the solution space<sup>2</sup>. We

<sup>2</sup>Also, the Mixture-of-Agents[12] paper finds that model diversity, in comparison to other hyperparameters, influences performance highly in their MAS.

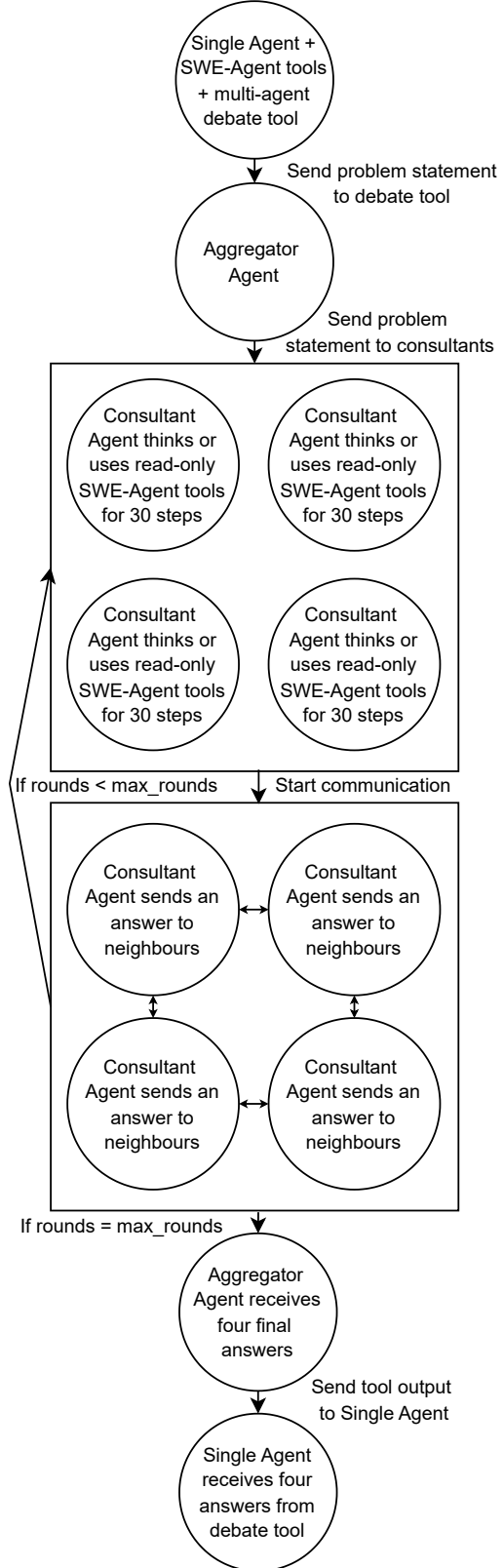


Figure 3: An illustration of how the multi-agent debate tool configuration used in this project attempts to provide answers to a given problem statement.

built up to these diverse model family experiments by first starting with a more unitary setup.

We begin by using Claude 3.5 Haiku (2024-10-22) and GPT-4o (2024-08-06). In our first set of experiments, we set the reflection depth to 30, the single agent token limit to 300k, and set the single agent and all agents in the debate tool to Haiku. We run this experiment for the 50 SWE-Bench tasks contained in SWE-Bench-verified-mini[5] (a subset attempting to replicate performance of the full dataset, with the smallest total container memory possible). We run three epochs of these 50 tasks, and calculate the mean and standard deviation in performance across epochs. We extract the average tokens used by the debate tool across all 50 tasks, via a script. We then add this number to the single-agent token limit of 300k, and re-run the above experiment by not providing the multi-agent debate tool to the single agent. In practice, the debate tool tended to use roughly 200-300k tokens per task, with the current number of rounds, message topology and and reflection depth<sup>34</sup>. We repeat this experiment with GPT-4o, and the results can be found in Figure 4. We observe that, when using the same LLM across the single agent and multi-agent debate tool, shifting test-time compute to the multi-agent debate tool results in approximately the same performance as allocating all test-time compute to the single agent. We also plot  $\pm$  one standard deviation error bars, for the performance levels across the three epochs.

The next set of experiments we conduct investigate whether using a different model for the multi-agent debate tool results in a performance increase. We run two experiments, the first uses Haiku as the single agent and 4o as all agents in the debate tool. The second is the reverse of this. For the single-agent baselines, we replicate the methodology in Figure 4. That is, for each (single agent, multi-agent tool) experiment, we add the single agent token limit to the average number of tokens the debate tool uses per sample. We then take this number of tokens, and set the token limit of the single agent baselines (without the debate tool) to this number. The results for this set of experiments can be found in Figure 5. We see that the mean performances of both (single agent, multi-agent tool) experiments are higher than their

single-agent system (SAS) baselines. We also note that (4o single agent, haiku multi-agent tool) has non-overlapping ( $\pm$  one standard deviation) error bars, compared to its baseline, which while not statistically significant at the 95% confidence interval level, may represent some evidence of the multi-agent debate tool resulting in a performance uplift, in comparison to a token-normalised baseline.

A key result we required in this experiment is for the mean performance of (4o SAS, Haiku MAS) to be higher than both (Haiku SAS) and (4o SAS) separately, as this would be evidence for the (4o SAS, Haiku MAS) experiment producing an emergent performance uplift, beyond the performance levels of its constituent underlying models. Figure 5 presents some evidence of this, but we note the requirement for higher statistical significance.

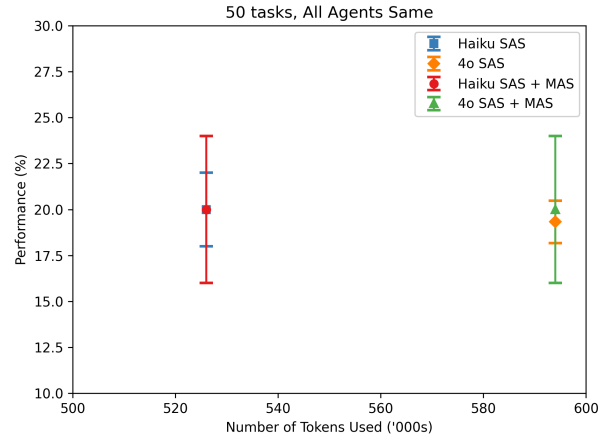


Figure 4: This experiment runs four different single-agent system (SAS) and single agent system, multi-agent system tool (SAS, MAS) configurations over 50 SWE-Bench tasks, for three epochs. All configurations used the same model for all agents. Note the similarity of the performance levels, for the two (SAS, MAS) configurations and their single-agent baselines.

## 4 Discussion

In this section, we’d like to discuss some initial investigations into output logs. We attempted to understand the reasons behind any performance uplifts, by looking for tasks where (4o SAS) in Figure 5 failed across all three epochs, (4o SAS, Haiku MAS) succeeded across all three epochs, and (Haiku SAS) failed in all, or at least one of, the three epochs. Our hypothesis was that any such tasks would represent situations where the composite system managed to achieve greater performance than its constituent underlying LLMs.

<sup>3</sup>Note: we can tune the token usage of the debate tool by varying these hyperparameters. If we wanted to explore higher token domains for instance, we could increase the number of rounds or reflection depth.

<sup>4</sup>Note: the variation in tokens used here doesn’t represent the range of tokens we would eventually like to test. This is discussed later, but in summary, we don’t expect significant differences in performance across the roughly 100k range of tokens used by the debate tool

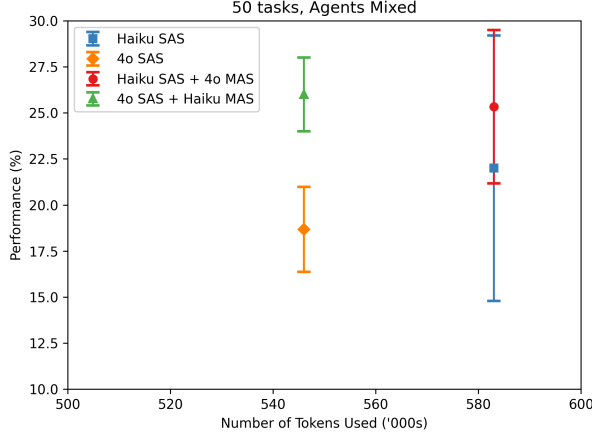


Figure 5: This experiment runs four different SAS and (SAS, MAS) configurations over 50 SWE-Bench tasks, for three epochs. In contrast to Figure 4, each configuration now uses two different language models. Note how, again in contrast to Figure 4, larger performance differences are generated between the (SAS, MAS) configurations and their single-agent baselines.

One SWE-Bench task which exhibited this was `Sphinx-doc_sphinx-10466`, the problem statement and golden solution for which is available in Appendix B. In the interests of brevity (given this analysis compares nine logs), we present Claude-generated summaries of the single agent’s final answer across all logs (captured in the agent’s use of Inspect’s `submit()` tool). We also present summaries of the debate tool’s outputs across each of the three epochs for (4o SAS, Haiku MAS).

This example is interesting because it represents a situation where 4o, when using its Haiku MAS tool, implements a set of solutions which is markedly different to those proposed by 4o and Haiku alone. We can also see, from the tool output summaries specifically, that the different solutions primarily arise from the debate tool, as opposed to some interaction effect between 4o and the debate tool’s output. A valuable next step would be to build an automated log analysis pipeline to extract tasks with these success profiles from a set of logs<sup>5</sup>.

We hope the above example illustrates how analysis on the incremental effect of multi-agent tools on a single agent could be automated, by passing the problem, solution and snippets of SAS and (SAS, MAS) logs

<sup>5</sup>We attempted a similar exercise with (Haiku SAS), (Haiku SAS, 4o MAS) and (4o SAS), but found no tasks with the above success profile. We present this analysis in Appendix A, with full `submit` and tool outputs, to more explicitly illustrate the structure of the outputs.

to a language model to analyse and provide feedback on whether the multi-agent tool is providing useful insights to the single agent. As adaptive MAS schemes such as EvoMAC[6] (discussed above) illustrate, this feedback could then be passed back to an orchestrator or aggregator agent, which could use the feedback to augment the architecture of the multi-agent system.

Table 1: Comparison of Solution Approaches

System Run		Solution Summary
4o (1/3)	SAS	Added <code>uniqueLocation</code> method to <code>Message</code> class to convert location tuples to a set and back to a list, removing duplicates.
4o (2/3)	SAS	Modified <code>Message.__init__</code> with helper function to ensure unique locations through set conversion.
4o (3/3)	SAS	Implemented <code>uniqueLocation</code> method in <code>Message</code> class for set-based deduplication.
4o+Haiku MAS (1/3)		<b>Tool:</b> Generated solutions with custom <code>deduplicate_locations()</code> function with order preservation. <b>Submit:</b> Implemented deduplication within <code>Message</code> class; noted partial solution.
4o+Haiku MAS (2/3)		<b>Tool:</b> Recommended <code>OrderedDict.fromkeys()</code> for deduplication; suggested modifying both classes. <b>Submit:</b> Implemented with <code>OrderedDict.fromkeys()</code> for both classes; noted testing limitations.
4o+Haiku MAS (3/3)		<b>Tool:</b> Generated sophisticated deduplication strategies across Sphinx and Babel components. <b>Submit:</b> Implemented <code>_unique_locations</code> in <code>Message</code> class; acknowledged need for further integration.
Haiku (1/3)	SAS	Correctly modified <code>Catalog.__iter__()</code> using <code>dict.fromkeys()</code> to remove duplicates while preserving order.
Haiku (2/3)	SAS	Changed to modify <code>Message</code> class, using simple set conversion to remove duplicates.
Haiku (3/3)	SAS	Focused on <code>Message.__init__()</code> , using <code>list(dict.fromkeys())</code> to maintain order.
Actual Solution		Modified <code>Catalog.__iter__()</code> using <code>sorted(set())</code> to create unique position list.

## 5 Increased Diversity Experiments

We next wanted to increase the diversity of the models used in the debate tool. However, it proved to be a difficult engineering challenge to reliably get our API provider (OpenRouter) to work consistently with



Autogen<sup>6</sup>. As a result, we were only able to use OpenRouter’s Gemini 2.0 Flash Lite API<sup>7</sup>. However, in this section, we would still like to outline the high diversity multi-agent experiments we were planning to implement, as we believe the methodology is generally applicable to investigating the performance uplift of most hyperparameters within our multi-agent debate tool. We opted for the following model set in this experiment: 4o at default temperature; 4o at a temperature of 0.9; Haiku at default temperature; Haiku at a temperature of 0.9; Gemini 2.0 Flash Lite at default temperature. We allow each agent to act as the single agent, and the other four agents are assigned to the multi-agent debate tool. This would have resulted in five pairs of experiments (each model as the SAS with a MAS tool, then each SAS alone), but Gemini 2.0 Flash Lite failed to run with Inspect, reducing Figure 6 to eight datapoints. Note that we didn’t have time to run multiple epochs for each datapoint here, and have also used single agent baselines from previous experiments, hence (SAS, MAS) datapoints not always aligning with their expected baselines. While some performance uplifts are present from using the debate tool, more epochs are required before drawing conclusions from this set of experiments.

We hope Figure 6 illustrates the broader experiment structure: use a different LLM for each agent in our (SAS, MAS tool) system, allowing each LLM to adopt the role of the single agent each time. From here, calculate the average number of tokens the multi-agent tool uses each time, and add this number to the token limit of the single agent baseline.

We can extrapolate this once more, to testing the performance uplifts of such a permuted set of experiments with different numbers of tokens (tuned via the debate tool’s hyperparameters). One of our goals was to produce a plot such as the one sketched in Figure 7. We hypothesise that as the number of tokens is increased further, the multi-agent debate tool will start to produce a larger performance uplift, compared to single-agent baselines.

## 6 Future Work

In this section, we briefly discuss some ideas for future work. The primary goal of any extension work would

<sup>6</sup>A variety of problems arose here, from sporadic API failures, to OpenRouter load balancing possibly resulting in inconsistent message/tool call formats. Most problems were eventually fixed, with some sporadic errors remaining for most models. Most remaining errors take the form of zero-token output errors. Some recent issues were also reported between Autogen and OpenRouter: <https://github.com/microsoft/autogen/issues/5621>

<sup>7</sup>Haiku was also accessed via OpenRouter, and worked very reliably

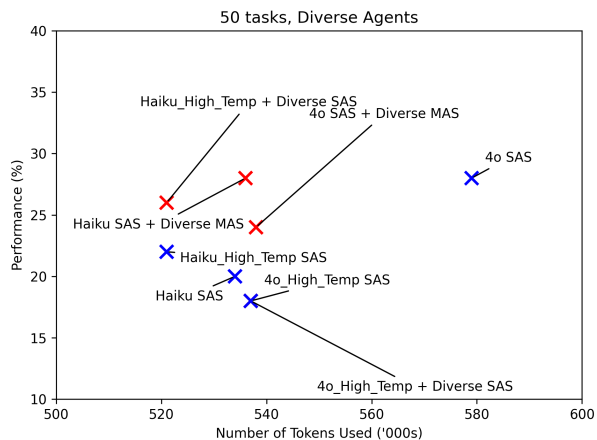


Figure 6: A preliminary set of experiments, attempting to extend Figures 4 and 5 to a larger number of LLMs within each (SAS, MAS) configuration. While performance over single-agent baselines looks somewhat promising, the noise is too high to reach robust conclusions.

be to reduce the noise of the above experiments and reach statistically significant conclusions for performance uplifts in MAS hyperparameters of interest. This can be expensive however, as using roughly half a million tokens per sample leads to costs of around \$1 per SWE-Bench task in the (4o SAS, 4o MAS) experiment for instance. The aforementioned automated log analysis pipeline should also be built, to automate the process of understanding the debate tool’s impact on the single agent, possibly by having an LLM evaluate the incremental benefit of the debate tool’s output on the single agent’s trajectories (in comparison to single-agent baselines). Experiments should then be run in an iterative process, switching between log analysis, debate tool refinement and hyperparameter exploration. Another exciting direction would be to build an evaluation which swaps the Aggregator agent in Figure 3 for an Orchestrator agent, which is provided with a set of tools it can use to adjust the MAS tool architecture. Such an evaluation would seek to measure a frontier model’s ability to achieve performance uplifts by adjusting the aforementioned hyperparameters of our debate tool<sup>8</sup>.

If, after wider experimentation, we don’t see a significant performance uplift from this debate tool, this could help AISI steer away from using multi-agent debate tools for software engineering tasks, and look

<sup>8</sup>Alongside the scaffolding proposal for this report, we submitted an evaluation proposal to AISI’s bounty program[1]. It can be found here. This evaluation proposal may be implemented with AISI, but requires updating after this project.

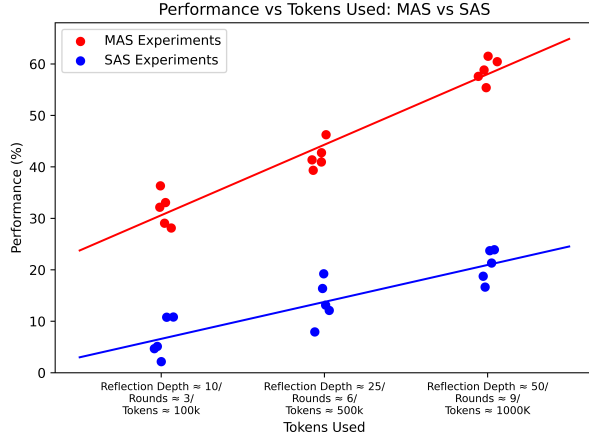


Figure 7: A hypothesised set of experiments, attempting to observe the relationship between performance and test-time compute, across single and multi-agentic systems.

more at sub-tasking MAS architectures. It may be that engineering tasks aren’t as amenable to performance uplifts via multi-agent debate as for instance, math or logic problems, where agents may be able to better utilise chains of thought/tool use from other agents, to either critique their own or other agents’ thoughts. Similarly, it may be that more decomposable and complex tasks are also better suited to both multi-agent sub-tasking or debate systems, and switching SWE-Bench for such a task may generate a larger performance difference between single and multi-agent systems. EvoMAC[6] does precisely this, tasking their multi-agent systems with delivering an engineering product, from just a design document.

## 7 Conclusion

In this report, we explored whether providing a single agent with a multi-agent debate tool, which it can invoke at the start of a SWE-Bench task, leads to performance uplifts when normalising for tokens. We observed that while there was some evidence of this (at the  $\pm$  one standard deviation confidence interval), larger sample sizes and more epochs are required. We outlined and implemented a baselining method, which we believe would form the basis of a more extensive set of experiments, across a number of hyperparameters available for adjusting the architecture of the multi-agent tool built in this project. We hope that future work at AISI is able to use this tool and the methods outlined in this report (alongside the hypothesised automated log analysis pipeline and multi-agent evaluation) to better understand how multi-agent scaffolding is able to elicit capability uplifts within agentic language models.

## References

- [1] UK AISI. *AISI Bounty Program*. 2024. URL: <https://www.aisi.gov.uk/work/evals-bounty>.
- [2] UK AISI. *Inspect*. 2025. URL: [https://github.com/UKGovernmentBEIS/inspect\\_ai](https://github.com/UKGovernmentBEIS/inspect_ai).
- [3] UK AISI. *inspect<sub>e</sub>vals*. 2025. URL: [https://github.com/UKGovernmentBEIS/inspect\\_evals](https://github.com/UKGovernmentBEIS/inspect_evals).
- [4] Lewis Hammond et al. *Multi-Agent Risks from Advanced AI*. 2025. arXiv: 2502.14143 [cs.MA]. URL: <https://arxiv.org/abs/2502.14143>.
- [5] Marius Hobbhahn. *SWEBench-verified-mini*. 2025. URL: <https://github.com/mariushobbhahn/SWEBench-verified-mini>.
- [6] Yue Hu et al. *Self-Evolving Multi-Agent Collaboration Networks for Software Development*. 2024. arXiv: 2410.16946 [cs.SE]. URL: <https://arxiv.org/abs/2410.16946>.
- [7] Dong Huang. *AgentCoder: Multi-Agent Code Generation with Effective Testing and Self-optimisation*. <https://arxiv.org/pdf/2312.13010>. 2024.
- [8] Dong Huang. *AgentCoder: Multi-Agent Code Generation with Effective Testing and Self-optimisation*. <https://arxiv.org/pdf/2312.13010>. 2024.
- [9] Carlos E. Jimenez et al. *SWE-bench: Can Language Models Resolve Real-World GitHub Issues?* 2024. arXiv: 2310.06770 [cs.CL]. URL: <https://arxiv.org/abs/2310.06770>.
- [10] Yubin Kim et al. *MDAgents: An Adaptive Collaboration of LLMs for Medical Decision-Making*. 2024. arXiv: 2404.15155 [cs.CL]. URL: <https://arxiv.org/abs/2404.15155>.
- [11] Junyou Li et al. *More Agents Is All You Need*. 2024. arXiv: 2402.05120 [cs.CL]. URL: <https://arxiv.org/abs/2402.05120>.
- [12] Junlin Wang et al. “Mixture-of-Agents Enhances Large Language Model Capabilities”. In: *The Thirteenth International Conference on Learning Representations*. 2025. URL: <https://openreview.net/forum?id=h0ZfDIrj7T>.
- [13] Qingyun Wu et al. *AutoGen: Enabling Next-Gen LLM Applications via Multi-Agent Conversation*. 2023. arXiv: 2308.08155 [cs.AI]. URL: <https://arxiv.org/abs/2308.08155>.



## A Appendix A: Haiku SAS vs (Haiku SAS, 4o MAS)

In this appendix, we present some log snippets, which attempt to illustrate why our debate tool elicits capability uplifts between the (Haiku SAS) and (Haiku SAS, 4o MAS) experiments in Figure 5. See Section 4 for more details.

We looked for tasks where each of the three epochs for (Haiku SAS) in Figure 5 failed, each of the three epochs for (Haiku SAS + 4o MAS) succeeded, and some of the three epochs for (4o SAS) failed. We didn't find any such tasks, so instead looked at a task where (4o SAS) was also successful across all three epochs. This occurred for one of the 50 tasks, namely `Sphinx-doc__sphinx-9367`. We see that Haiku misunderstands the problem and thinks that a test needs to be written, while GPT-4o correctly identifies that a method needs to be modified.

### A.1 Problem

Please solve the following coding issue:

1-element tuple rendered incorrectly

Describe the bug

This is a followup to #7964 which has been addressed in #8265.

However the special case of a 1-element tuple is still not handled correctly.

(1,) is rendered as (1), but should keep the trailing comma.

To Reproduce

Add a testcase

`("(1,)", "(1,)", # Tuple (single element)`

at

`https://github.com/sphinx-doc/sphinx/blob/e0b1e1002b500acc63dfd0806f8095dd6b27037b/tests/test_pycode_ast.py#L57`

### A.2 Golden Solution

```
diff --git a/sphinx/pycode/ast.py b/sphinx/pycode/ast.py
--- a/sphinx/pycode/ast.py
+++ b/sphinx/pycode/ast.py
@@ -213,10 +213,12 @@ def visit_UnaryOp(self, node: ast.UnaryOp) -> str:
     return "%s %s" % (self.visit(node.op), self.visit(node.operand))

 def visit_Tuple(self, node: ast.Tuple) -> str:
-    if node.elts:
-        return "(" + ", ".join(self.visit(e) for e in node.elts) + ")"
-    else:
-        if len(node.elts) == 0:
-            return "()"
-        elif len(node.elts) == 1:
-            return "(%s,)" % self.visit(node.elts[0])
-        else:
-            return "(" + ", ".join(self.visit(e) for e in node.elts) + ")"

if sys.version_info < (3, 8):
    # these ast nodes were deprecated in python 3.8
```

### A.3 Figure 5 - Haiku SAS

In this section, we only present the single agent's submission to its `submit()` tool.

#### A.3.1 2025-03-27T23-23-14+00-00\_swe-bench\_V5epJU9SBLzaX66gbm\_j4AW.eval

I added a new test case for the 1-element tuple `("(1,)", "(1,)",)` in the `'test_unparse'` parametrized test in the `'tests/test_pycode_ast.py'` file. This test ensures that a single-element tuple is correctly rendered with its trailing comma, addressing the issue with tuple representation.

#### A.3.2 2025-03-27T23-25-59+00-00\_swe-bench\_5v5QeS3549i2bdQZejJaoA.eval

I added a test case for the 1-element tuple rendering issue in the `test_pycode_ast.py` file.

Specifically, I added the test case:

```
'''python
("(1,)", "(1,)", # Tuple (single element)
'''
```

This test case ensures that a 1-element tuple with a trailing comma is correctly rendered, preserving the comma when unparsed. The test is added to the existing `'test_unparse'` function's `parametrize` decorator, which checks that the AST unparsing method correctly handles various Python expressions.

This addresses the issue of 1-element tuple rendering, making sure that `'(1,)'` is not simplified to `'(1)'`, thus maintaining the distinction between a single-element tuple and a parenthesized expression.

#### A.3.3 2025-03-27T23-52-13+00-00\_swe-bench\_C4v8qFUvLc8Hrq9sgSES5B.eval

I added a test case for the 1-element tuple, showing that `'(1,)'` should render as `'"(1,)"'` instead of `'"(1)"'`. This verifies that the trailing comma in a 1-element tuple is preserved during unparsing.

The changes were made in the test parametrization in the file `'tests/test_pycode_ast.py'`, adding the line:

```
'''python
("(1,)", "(1,)", # Tuple (single element)
'''
```

This addresses the issue described in the bug report, ensuring that the special case of a 1-element tuple is handled correctly by maintaining the trailing comma. (Haiku SAS, 4o MAS) Figure 2:

#### A.4 Figure 5 - (Haiku SAS, 4o MAS)

In this section, we present the output from the single agent's `consult_multi_agent()` tool, followed by the single agent's submission to its `submit()` tool.

##### A.4.1 2025-03-21T19-18-58+00-00\_swe-bench\_dpaLoBoMAHNkTRPoegLJvP.eval

## All solutions from the solver agents:

### Solution 1: FINAL ANSWER: To fix the incorrect rendering of single-element tuples, modify the `'visit_Tuple'` method in the `'_UnparseVisitor'` class within `'sphinx/pycode/ast.py'` to append a trailing comma when the tuple has only one element. Add a test case in `'tests/test_pycode_ast.py'` to verify this behavior, ensuring that `'(1,)'` is correctly rendered with a trailing comma.

### Solution 2: FINAL ANSWER: To fix the rendering of single-element tuples, modify the `'visit_Tuple'` method in the `'_UnparseVisitor'` class in `'sphinx/pycode/ast.py'` to append a trailing comma when the tuple has only one element. Add a test case in `'tests/test_pycode_ast.py'` to verify this behavior. This ensures that `'(1,)'` is rendered correctly, preserving the trailing comma.

### Solution 3: FINAL ANSWER: Modify the `'visit_Tuple'` method in `'sphinx/pycode/ast.py'` to check for single-element tuples and append a trailing comma. Update the test suite in `'tests/test_pycode_ast.py'` with a test case for a single-element tuple and verify the fix by running the tests. This ensures tuples like `'(1,)'` are rendered correctly with a trailing comma.

### Solution 4: FINAL ANSWER: To fix the issue of incorrectly rendering single-element tuples, modify the `'visit_Tuple'` method in the `'_UnparseVisitor'` class within `'sphinx/pycode/ast.py'` to append a trailing comma for single-element tuples. Update the method to check if the tuple has only one element and append a comma if so. Add a test case in `'tests/test_pycode_ast.py'` to verify this behavior and run the tests to ensure correctness.

I solved the issue of 1-element tuple rendering by modifying the `'visit_Tuple'` method in `'sphinx/pycode/ast.py'`.

The changes ensure that a single-element tuple like `'(1,)'` is correctly rendered with a trailing comma, while preserving the existing behavior for empty and multi-element tuples.

Key changes: 1. Added a specific handling for single-element tuples by adding `','` to the output 2. Updated the test case in `'tests/test_pycode_ast.py'` to include the test for a single-element tuple

Specifically, I modified the `'visit_Tuple'` method to:

```
'''python
def visit_Tuple(self, node: ast.Tuple) -> str:
```

##### A.4.2 2025-03-27T18-19-25+00-00\_swe-bench\_KpYoxUvQn6TDYswDLGnBMY.eval

## All solutions from the solver agents:

### Solution 1: FINAL ANSWER: The tuple rendering logic is in the `'visit_Tuple'` method in `'sphinx/pycode/ast.py'`. Modify this method to append a comma for single-element tuples. Add a test case in `'tests/test_pycode_parser.py'` to verify that `'(1,)'` is rendered correctly as `'(1,)'`, preserving the tuple's identity.

### Solution 2: FINAL ANSWER: To fix the issue of incorrectly rendering a 1-element tuple `'(1,)'`, modify the `'visit_Tuple'` method in `'sphinx/pycode/ast.py'` to append a trailing comma for single-element tuples. Add a test case in `'tests/test_pycode_parser.py'` to verify the correct rendering of `'(1,)'`. This ensures the tuple's identity is preserved during rendering.

### Solution 3: FINAL ANSWER: The tuple rendering logic is in the `'visit_Tuple'` method in `'sphinx/pycode/ast.py'`. Modify this method to check if the tuple has one element and append a trailing comma if true. This ensures single-element tuples are rendered correctly as `'(1,)'`. Add a test case to verify this behavior.

I solved the issue of 1-element tuple rendering by modifying the `'visit_Tuple'` method in `'testbed/sphinx/pycode/ast.py'`.

The changes ensure that a single-element tuple always includes a trailing comma when rendered. Specifically, I updated the implementation to: 1. Return `""` if the tuple is empty 2. Always add a trailing comma for single-element tuples 3. Preserve the existing comma behavior for multi-element tuples

The modified code looks like this:

```
'''python
def visit_Tuple(self, node: ast.Tuple) -> str:
    if not node.elts:
        return ""
    result = ", ".join(self.visit(e) for e in node.elts)
    # Add a trailing comma if there's only one element
    return "(" + (result + ",") if len(node.elts) == 1 else result + ")"
```

I also added a test case to the `'test_pycode_ast.py'` file to verify the 1-element tuple rendering:

```
'''python
("(1,)", "(1,)") # Tuple (single element)
'''
```

This ensures that `'(1,)'` is correctly rendered as `'(1,)'`, preserving the tuple's identity.

#### A.4.3 2025-03-27T18-56-44+00-00\_swe-bench\_YzuNkFhH3bXBeNt4Dgj6ug.eval

## All solutions from the solver agents:

### Solution 1: FINAL ANSWER: The issue of rendering 1-element tuples incorrectly in Sphinx is located in the 'sphinx/pycode/ast.py' file, specifically within the 'visit\_Tuple' method. To address this, modify the 'visit\_Tuple' method to check if the tuple has only one element. If it does, ensure that a comma is appended after the element before enclosing it in parentheses. This change will ensure that '(1,)' is rendered correctly with the trailing comma, distinguishing it from a simple parenthesized expression. After implementing this fix, it is crucial to add test cases that cover single-element tuples to verify the correctness of the change. Additionally, update any relevant documentation to reflect this behavior and follow best practices for code review and deployment to ensure the change is integrated smoothly into the codebase.

### Solution 2: FINAL ANSWER: Modify the 'visit\_Tuple' method in 'sphinx/pycode/ast.py' to handle single-element tuples by appending a comma after the element. This ensures that '(1,)' is rendered correctly. Update the method to check for a single element and append a comma accordingly. Add test cases to confirm the fix and follow best practices for documentation and code review.

### Solution 3: FINAL ANSWER: To address the issue of incorrectly rendering 1-element tuples in Sphinx, the solution involves modifying the 'visit\_Tuple' method in the 'sphinx/pycode/ast.py' file. This method is responsible for handling the rendering of tuples. The change required is to ensure that when a tuple has only one element, a comma is appended after the element before enclosing it in parentheses. This adjustment will ensure that a single-element tuple like '(1,)' is rendered correctly with the trailing comma, distinguishing it from a simple parenthesized expression like '(1)'.

Steps to implement the solution: 1. Locate the 'visit\_Tuple' method in the 'sphinx/pycode/ast.py' file. 2. Modify the method to check if the tuple has only one element. If so, append a comma after the element. 3. Enclose the element and the comma in parentheses to ensure the correct rendering of a single-element tuple. 4. After implementing the change, run the existing test suite to ensure that the fix does not introduce any regressions. 5. Add new test cases specifically for single-element tuples to verify that they are rendered correctly with the trailing comma. 6. Update any relevant documentation to reflect the changes made to the codebase. 7. Follow best practices for code review

and deployment to ensure the changes are integrated smoothly into the main codebase.

By following these steps, the issue of incorrectly rendering 1-element tuples in Sphinx can be resolved effectively.

### Solution 4: FINAL ANSWER: To fix the rendering of 1-element tuples in Sphinx, modify the 'visit\_Tuple' method in 'sphinx/pycode/ast.py' to append a comma when the tuple has only one element. This ensures '(1,)' is rendered correctly. Validate the fix by running tests and adding new test cases for single-element tuples. Update documentation and follow best practices for code review and deployment.

I solved the issue of incorrectly rendered 1-element tuples in the Sphinx project by making two key changes:

1. Modified 'visit\_Tuple' method in 'sphinx/pycode/ast.py': - Added a specific condition to handle single-element tuples - Ensured that a single-element tuple preserves its trailing comma - The new implementation explicitly adds a comma for 1-element tuples

Updated code:

```
'''python
def visit_Tuple(self, node: ast.Tuple) -> str:
    if not node.elts:
        return "()"
    # For single-element tuples, add a trailing comma
    elif len(node.elts) == 1:
        return "(" + self.visit(node.elts[0]) + ","
    # For multi-element tuples, use the original behavior
    else:
        return "(" + ", ".join(self.visit(e) for e in node.elts) + ")"
'''
```

2. Added a test case in 'tests/test\_pycode\_ast.py': - Included a new test for single-element tuple: '("(1,)", "(1,))"

This ensures that '(1,)' is correctly rendered with a trailing comma, distinguishing it from a simple parenthesized expression '(1)'.

The changes address the issue of single-element tuples losing their trailing comma during rendering, preserving the Python syntax that requires a comma to create a single-element tuple.

## B Appendix B: Second Logs Example

In this appendix we provide the problem statement and golden solution to accompany the discussion of the (4o SAS) and (4o SAS, Haiku MAS) experiments in Section 4.

### B.1 Problem

Please solve the following coding issue:

Message.locations duplicate unnecessary  
Describe the bug  
When running

