
The Influence of Scaffolds on Coordination Scaling Laws in LLM Agents

Anonymous Author(s)

Affiliation

Address

email

Abstract

As large language models improve in capability, they are increasingly taking on more agentic and interactive roles in multi-agent settings that demand effective communication and coordination. In order to measure a model’s capabilities in these settings, new benchmarks are quickly emerging to study language-based, multi-agent interaction, often by adding language scaffolds on top of existing multi-agent environments. However, when evaluating agents on such benchmarks, an agent’s performance can be significantly influenced by implicit factors related to the design of the scaffolds, rather than the inherent properties of the agents. Moreover, it is unclear if coordination among agents in these settings follows scaling laws. We consider one such environment—the popular collaborative cooking environment, Collab-Overcooked—and characterize how scaffolding plays a role in successful collaborations between models of varying sizes. We perform empirical evaluations on the collaborative capabilities of agents and find that, as long as models are given clear instructions on *how* to collaborate, their capabilities follow positive scaling laws in both self-play and cross-play. However, without a scaffold that explicitly defines how the collaboration should be done, we find models struggle to develop effective methods for collaboration, and scaling laws break down. Our experiments highlight how subtle changes in agent scaffolds can drastically impact their collaborative capabilities and raises questions on how to design evaluations for agents that may have to collaborate with open-ended partners.

1 Introduction

Large language models (LLMs) are moving from single-agent deployments to open, multi-agent ecologies in which AI assistants interact with humans and other AI agents to carry out real-world tasks. Improvements in AI agent capabilities are allowing them to be deployed in increasingly complex tasks Kwa et al. [2025]. Such tasks range from trivial ones like scheduling a table at a restaurant or grocery shopping Rogers [2025], to genuinely consequential ones, such as handling bureaucratic obligations, advising about medical care or on other high-stakes decisions Palantir Technologies [2025], Preiksaitis et al. [2024], Li et al. [2024], Newman et al. [2022]. As assistants assume such open-ended roles, they will necessarily need to interact with other agents. We focus on the ability of models to *coordinate*: how well they form shared plans and distribute roles between varying partners. Moreover, we examine when coordination follows *scaling laws*, and how design choices (scaffolds, prompts, and environmental factors) shape their ability to coordinate.

We study coordination in Collab-Overcooked [Sun et al., 2025], a two-player, fully cooperative benchmark featuring explicit inter-agent communication, interactive planning, and sparse rewards. Collab-Overcooked instantiates the core primitives required for successful coordination between LLM

36 agents—language-mediated negotiation, interdependent sub-tasks, shared-resource management, and
37 rapid role allocation—with in a controlled setting that enables precise, repeatable measurement.

38 Since LLM agents operating in the real-world need to be able to coordinate with open-ended partners
39 (whether that be humans or other agents), our analysis focuses on studying models in *cross-play*:
40 when their partner is different than themselves. Our primary analysis examines cross-play within a
41 single model lineage (Qwen3.0; 1.7B–32B) [Yang et al., 2025]. Our first result shows that, under
42 carefully designed scaffolds, models obey clean scaling laws when it comes to coordination: agents
43 tend to do better when their underlying model size increases, or their partner’s.

44 However, we find that as interactions become less clearly defined within the agents’ scaffolds, scaling
45 laws break down. To isolate drivers of coordination, we vary three additional factors: (i) the *scaffolds*
46 that define an agent’s role in the interaction, (ii) *game structure*, varying between asymmetric and
47 symmetric layouts ; (iii) *turn order*, swapping which assistant moves first . Our results indicate
48 that LLM agent coordination capabilities might be become more limited as multi-agent interactions
49 become less well-defined and more open-ended. Our main empirical contributions are as follows:

- 50 • **Scaling trends in cross-play.** With clearly defined scaffolds, performance generally in-
51 creases with the model size of either partner with a compounding effect.
- 52 • **Open-ended interaction degrades scaling laws.** When the restrictions on agent interactions
53 are relaxed, scaling laws break down.
- 54 • **Emergence of hierarchy predicts success.** Task success is correlated with how strong a
55 hierarchy emerges between the two agents in the game.

56 2 Related Work

57 **LLM agent coordination.** We use *Collab-Overcooked* [Sun et al., 2025] as the basis for our
58 experiments, a benchmark originally used to study self-play cooperation. We build on this work by
59 studying model’s in cross-play as well as introducing modifications to the environment and scaffolds
60 to study their effects on the coordination capabilities of agents. There are an increasing number of
61 other benchmarks arising for studying the coordination capabilities of agents, ranging from various
62 types of simple matrix games including “Guess 2/3 of the Average”, “El Farol Bar”, “Divide the
63 Dollar” tse Huang et al. [2025] to a cooperative symmetric game of traveling salesman [Jeknic
64 et al., 2025], to graph-based coordination environments [de Carvalho Silva and Macharet, 2025],
65 and simulating societies of agents [Piatti et al., 2024]. However, while these settings do require
66 communication and coordination between LLMs to succeed, they do not evaluate them in multi-step
67 agentic coordination settings.

68 There is some, but limited work, on investigating LLMs in our setting of multi-agent cross-play
69 scenarios. In [Curvo et al., 2025], the authors evaluate models of different sizes and providers in a
70 mixed-motive game, finding that smaller, lower-performing agents can be influenced by larger ones.
71 Another work Chen et al. [2025] studies the ability of models that are specialized for different tasks
72 to coordinate on time series anomaly detection tasks.

73 **Scaling laws.** Several works have studied various scaling laws in language models Kaplan et al.
74 [2020] and agentic capabilities Kwa et al. [2025]. Other prior work has also documented *inverse*
75 *scaling* in language models: some behaviors worsen with scale and, with improved prompting or
76 data, can become U-shaped [McKenzie et al., 2023, Wei et al., 2022]. We connect scaling laws to
77 multi-agent coordination: demonstrating when and how language models follow scaling laws in
78 cross-play.

79 **Prompting and Scaffolds.** Several works have focused on building scaffolds that coordinate
80 multiple LLMs in a way to outperform a single model [Wu et al., 2023, Suzgun and Kalai, 2024].
81 For example, Wang et al. [2024] introduces a scaffold that consists of two LLM agents with distinct
82 roles: a planning agent and a reasoning agent. Cross et al. [2024] scaffold agents to explicitly engage
83 in theory of mind to improve cooperation capabilities in mixed-motive settings. Similarly, our work
84 studies the impact on agent’s coordination abilities of various design choices in the scaffolds within
85 the Collab-Overcooked environment.

86 **3 The Environment**

87 **Benchmark and tasks.** We evaluate agents in *Collab-Overcooked* [Sun et al., 2025], a two-agent
 88 kitchen simulation that enforces collaboration via resource isolation and asymmetric task knowledge.
 89 Figure 1 shows the asymmetric layout: the *Chef* controls the pot, oven, and delivery pass; the
 90 *Assistant* has access to the ingredient dispenser, chopping board, dish stack, and blender; both agents
 91 share a central counter for hand-offs.

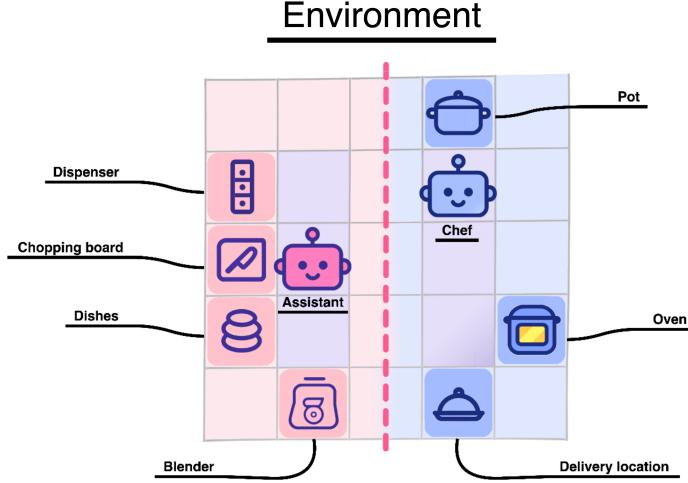


Figure 1: Collab-Overcooked layout. Purple tiles denote walkable floor; pink/blue tiles denote counters. The central line demarcates the asymmetric partition for the asymmetric environment. In the symmetric variant, we remove this divider and convert the counter area to walkable floor; all other items remain in place. The Chef becomes the First Player in the symmetric environment while the Assistant becomes the Second Player.

92 We study two variants of this environment:

- 93 • **Asymmetric:** Roles and access are fixed as above. This setting resolves all hierarchical
 94 ambiguity: one agent naturally plans and delivers (Chef) while the other supplies and
 95 prepares ingredients (Assistant).
- 96 • **Symmetric:** The physical partition is removed; both agents see identical task information
 97 (including the recipe) and can access all stations. No roles are prescribed—agents must
 98 negotiate a plan and divide work via communication. Prompt details are provided in the
 99 Appendix’s Figure 15. As a convention we say that the agent that used to be the Chef is now
 100 First Agent and the agent that used to be the Assistant second agent.

101 **Agents and pairings.** We study self-play (homogeneous pairings), cross-play (heterogeneous
 102 pairings) in all scenarios, and change which agent gets to go first in the asymmetric experiment. We
 103 use **Qwen 3.0** models at five scales: {1.7B, 4B, 8B, 14B, 32B}. LLMs use temperature 0.7, and
 104 a max communication budget of 4 SAY turns per communication window. We disable external tool
 105 usage.

106 **Prompts and memory.** In the *asymmetrical* environment, all agents receive high-level context
 107 regarding the structure of the task and who they are, the current environment state s_t , and role-specific
 108 capabilities. Moreover, the chef has exclusive access to the recipes.

109 In the *symmetrical* environment, agents receive the same information as above. With the difference
 110 that they share the same prompt and context that establishes no role but still informs them about the
 111 task to be completed.

112 In both asymmetrical and symmetrical cases, agents retain a short-horizon scratch memory (last K
 113 transitions, $K=10$) and the transcript of the most recent communication window.

114 **3.1 Experimental protocol**

115 We evaluate *ordered* model pairs to capture role/turn-order effects. All metrics in §3.2 are computed
116 per episode, averaged over the E episodes within a recipe, then over the S recipes within a level,
117 and finally (when reported “overall”) averaged across levels. Each episode is capped at $T_{\max} = 120$
118 environment timesteps.

- 119 • **Asymmetric.** For each ordered pair (i, j) we run $S=5$ recipes, each with $E=10$ episodes,
120 across 5 levels (levels 1–5). This yields $5 \times 5 \times 10 = 250$ episodes per ordered pair.
121 • **Symmetric.** For comparability and tractability, we restrict evaluation to the two levels that
122 admit faithful symmetric/asymmetric mappings and span different coordination regimes
123 (levels 2 and 4; details in the Appendix D). For each ordered pair (i, j) we run the same
124 $S=5$ recipes with $E=10$ episodes per recipe, giving $2 \times 5 \times 10 = 100$ episodes per ordered
125 pair.

126 **Reproducibility.** All code, prompts, logs, and scripts to regenerate figures will be released upon
127 publication.

128 **3.2 Metrics**

129 **1. Success Rate:**

$$\text{success rate} = \frac{\text{number of times the experiment is successful}}{\text{total number of experiments}}$$

130 **2. Similarity to the RAT:** The RAT, i.e., Referential Action Trajectory, is defined as an optimal
131 trajectory. This includes the minimum number of actions required to complete a recipe. The RAT
132 includes actions of both agents. Figure 14 in the Appendix shows an example of a RAT. Let n_k be
133 the length of the RAT for agent k . Let D_j^{\max} be the maximum overlap of the agent’s trajectory with
134 the RAT. In order to measure the partial completion of a task, we use the metric ‘Mean similarity to
135 the RAT’ defined as below:

$$\text{similarity to the RAT} = \frac{D_j^{\max}}{n_k}$$

136 **4 Results**

137 **4.1 Claims**

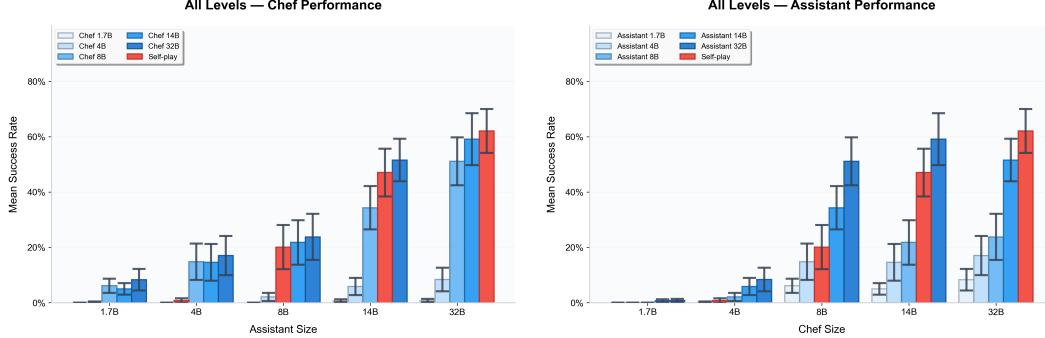
138 Through our experiments, we investigate the following hypotheses:

- 139 • H_1 (*Self-play scaling*). Within a fixed model family, self-play performance increases
140 monotonically with model size.
141 • $H_{2'}$ (*Scaffold dependence of scaling*). Under less prescriptive scaffolds (no role definitions),
142 scaling trends become less evident, indicating that scaling can be scaffold-induced.
143 • $H_{3'}$ (*Hierarchy predicts success*). The emergence of a clear leader–follower hierarchy
144 predicts higher task performance.
145 • $H_{4'}$ (*Parallelization amplifies coordination*). For tasks with greater opportunity for work in
146 parallel (higher decomposability), cooperative performance increases, and scaling trends
147 become more evident.

148 **4.2 Asymmetric Env**

149 **Does cooperation among LLM agents follow scaling laws?** Figure 9 shows a clear scaling trend:
150 as model size increases, mean success rises in both self-play and cross-play. The improvements are
151 most pronounced up to mid-scale (notably from 4B to 8B and again to 14B) and then taper off, with

152 smaller, overlapping gains beyond 14B. This pattern suggests that once a basic planning/communication
 153 competence is achieved, raw capacity ceases to be the primary limiter to increase the success
 154 rate; interaction dynamics become the bottleneck. We explain why this happens using the symmetric
 155 setting in §4.3, where turn order and emergent role assignment explain both the departures from a
 156 purely size-ordered ranking.



(a) With Assistant size fixed, larger Chefs raise success, but gains beyond 14B lie within overlapping confidence intervals. Chef scaling is less effective than Assistant scaling's in Fig. 2b.

(b) With Chef size fixed, larger Assistants improve success. A 32B Assistant especially boosts smaller Chefs—note the strong performance when paired with 8B and 14B models.

Figure 2: Cross-play and self-play outcomes in the asymmetric Overcooked setting (means with 95% CIs, averaged over all tasks)

157 **How do LLMs behave in different roles?** Figure 3 decomposes the similarity to RAT by role.
 158 Larger models not only plan well as Chef; they also track partner plans and recover from partner
 159 errors more effectively as Assistant, as we will see in the section 4a and 4c, yielding high RAT
 160 alignment even with weaker Chefs. We have included randomly sampled examples of interactions
 161 between small models playing as Chef and larger models playing as Assistants in order to illustrate
 162 this in the Appendix F. This asymmetry explains part of the reason why some mixed pairs outperform
 163 symmetric strong-strong pairs: a capable follower resolves ambiguity early, increasing the chance
 164 for success.

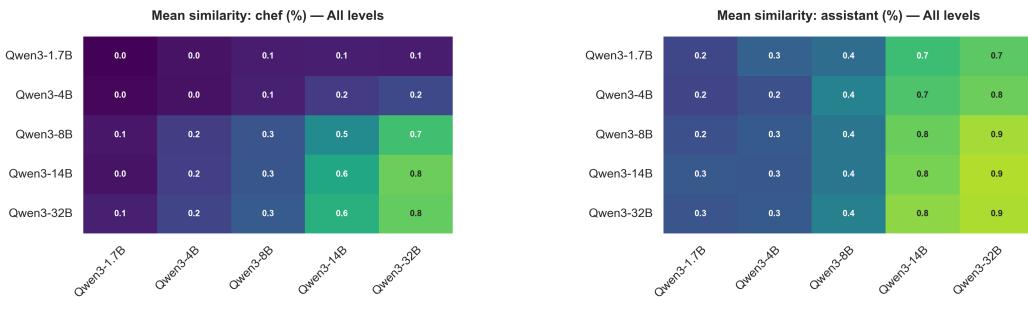


Figure 3: Mean similarity to the RAT (Referential Action Trajectory), averaged across all levels. The y-axis denotes models in the Chef role; the x-axis denotes models in the Assistant role.

165 **Cross-vendor testing** In order to verify our results in the broader ecosystem we have also tested
 166 the asymmetric Overcooked environment with the Acemotron 14B and Gemma 4B and Gemma 14B.
 167 The scaling results seem to be consistent with the ones observed with the Qwen vendors. Acemotron
 168 4B and Gemma 12B models scored similarly to Qwen 14B and Gemma 4B obtained the same scores
 169 as Qwen 4B. Figure 8 in the Appendix shows their success rates.

170 **4.3 Symmetric environment**

171 When comparing the success rate of the levels 2 and 4 for both scenarios (these are the only
 172 overlapping levels for both scenarios for reasons described in the session), we observe in Fig. 4 that
 173 the results for completion of symmetric environments are much higher.

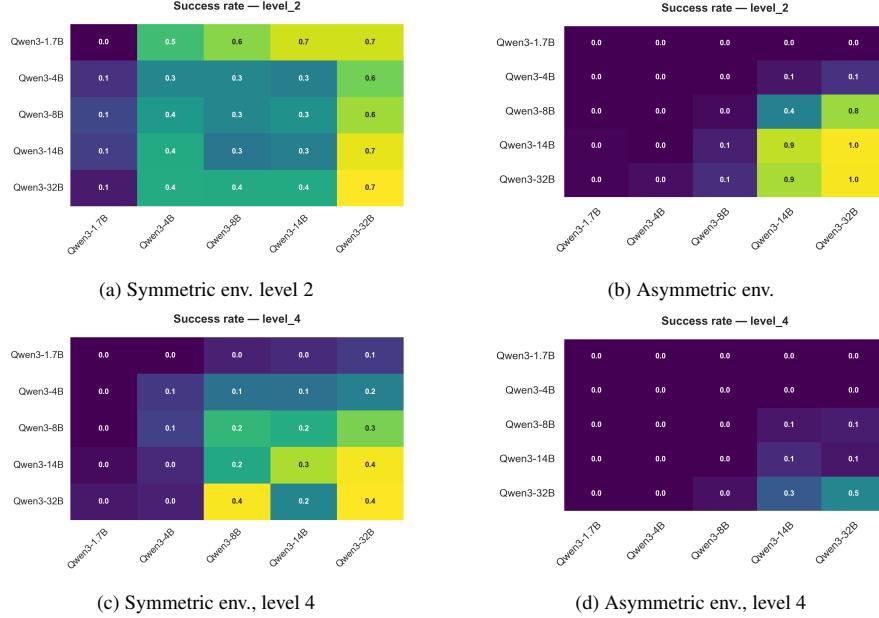


Figure 4: We compare the success of the models playing a symmetric version of the Overcooked game with models playing an asymmetric version. The y-axis represent the First Player while the x-axis represents the Second Player.

174 The reasons for why that happens are varied, as we will see in the course of this session. One clear
 175 environmental difficulty is removed; a common reason for misunderstanding between models is
 176 to know what each one of them has access to, allowing models to access the entire environment
 177 facilitates the task in this sense. However, we show that there are other unexpected factors relevant
 178 for improving the collective performance of LLMs interacting in coordination environments.

179 **Turn order induces a leadership prior.** In the symmetric setting, the agent that *speaks/acts first*
 180 receives a small prior both to itself and to the other agent to act as a coordinator. This prior is not
 181 assigned by any prompting but is just an artifact of being the first one to read the recipe; the agent to
 182 go first is thereafter treated implicitly as the coordinator. This prior explains three striking matrix
 183 patterns in Figure 4a:

- 184 **First column is weakest (1.7B as the Second Player).** When a model is the *First Player*,
 185 it assumes the leading position and expects a cooperative follower. Instead, 1.7B tends to
 186 flood the channel with repeated recitations of the actions one should take in the environment
 187 and self-assigns actions, derailing the stronger partner’s plan and stalling role resolution.
 188 This produces uniformly low success when 1.7B is cast as the follower.
- 189 **First row is unexpectedly strong (1.7B as the First Player).** When 1.7B goes first, its
 190 repetitive instruction dumps act like a crude “project brief.” Stronger partners interpret this
 191 as a leadership signal and, with the follower prior, often salvage the plan by extracting
 192 actionable steps. Paradoxically, the same verbosity that harms it as follower helps it as a
 193 proto-leader because the other agent adopts the follower role early.
- 194 **Last column is strongest (32B as the Second Player).** Having the largest model to go last,
 195 or falling in the follower role, yields robust success rates across leaders: the 32B agent rapidly
 196 understand the context and commits to the other agent’s plan (Figure 12). Furthermore, in
 197 every mirrored pairing with a smaller partner $m \in \{4, 8, 14\}$, the configuration with 32B as

198 *follower* ($m \times 32$) outranks its role-swapped mirror ($32 \times m$) in the level-2 symmetric game
 199 (see the top-7 success rate in Table 1b)

200 For the reasons detailed above, we exclude the 1.7B-parameter models from further analysis. We
 201 find that these models lack a basic understanding of their operating environment and fail to adhere to
 202 the task constraints. Although their strategy yielded favorable outcomes in this specific setting, we
 203 attribute this to the simplicity of the task and the ability of the Second Player to complete it rather than
 204 genuine competence and therefore do not treat it as substantive evidence. We have added evidence of
 205 this behavior in the Appendix E. Moreover, the behavior is not robust as it does not reproduce under
 206 increased task difficulty as seen in Figure 4c.

207 **Hierarchy formation correlates with performance.** Across pairings, the early emergence of a
 208 stable hierarchy—where one agent consistently delegates and the other executes—is associated with
 209 higher success. In our level-2 runs, pairings that quickly converged to a clear hierarchy (e.g., 14×32
 210 and 32×32) achieved success rates of 0.6–0.7, whereas pairings with weak or unstable hierarchies
 211 (e.g., 4×14 and 8×14) hovered around 0.3. Five of the seven highest-performing combinations
 212 (Table 1b) also appear among the seven pairings with the most games labeled as clearly cooperative
 213 (Table 1a). Nevertheless, model capacity remains relevant for performance: pairings that include a
 214 larger model—e.g., 4×32 and 8×32 —achieved success rates near 0.6.

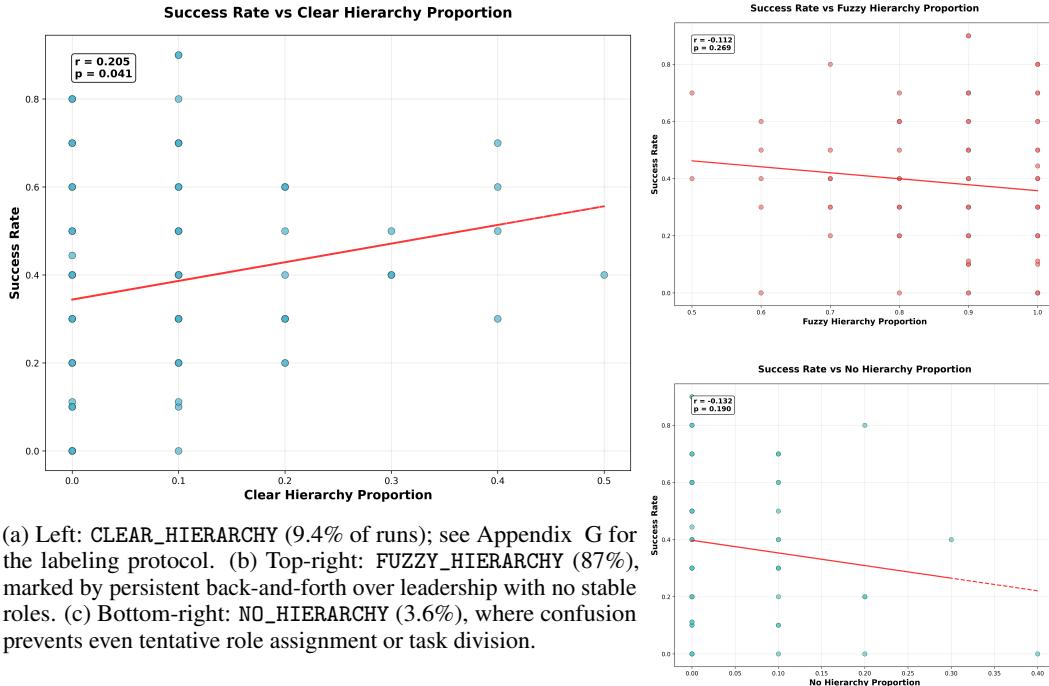
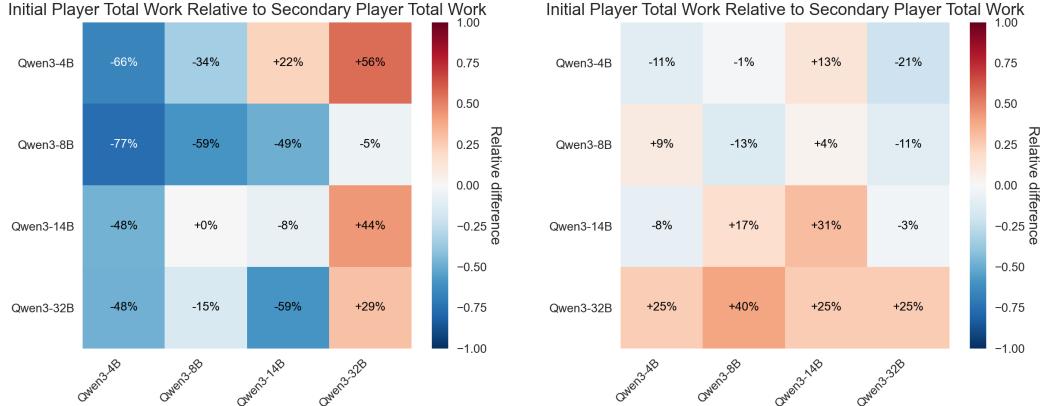


Figure 5: Correlation between models’ success rates and hierarchy conditions.

215 **Parallelizable tasks.** In addition to level-2 evaluations in the symmetric setting, we analyzed level-4
 216 tasks, which are parallelizable and therefore amplify the benefits of coordination. Performance at
 217 level-4 is substantially higher in the symmetric than in the asymmetric environment (Figures 4c, 4d).
 218 Division of work is also more balanced at this level (Figure 6b). Two factors likely contribute to these
 219 results are: (i) earlier task completion, producing shorter interaction histories and a lighter context
 220 window (Table 2); and (ii) the ability for each agent to concentrate on its subtask without closely
 221 tracking the partner’s state.

222 **Task division across task complexity.** Comparing level-2 and level-4, we find that explicit task
 223 division matters little in level-2, where many pairs achieve strong performance without clear role
 224 separation (Figure 6a). By contrast, at level-4, effective task division tends to emerge from successful



(a) Level 2. Number of valid actions by First Player (y-axis) vs. Second Player (x-axis). (b) Level 4. Number of valid actions by First Player (y-axis) vs. Second Player (x-axis).

Figure 6: Task division heat maps in the symmetric environment across two difficulty levels. Values near 0 indicate an even split between players.

225 coordination and is more tightly linked to higher success (Figure 6b). We hypothesize that level-2
 226 problems are simple enough to solve without structured coordination, whereas the structure of level-4
 227 inherently incentivizes parallelization; when models perceive opportunities to share work, they do so.
 228 These findings imply that coordination algorithms are complexity-dependent: strongly cooperative
 229 strategies arise only when task demands warrant them. Though as seen in Figure 4: they won’t
 230 emerge when model capacity is insufficient.

231 5 Conclusion and Future Work

232 **Summary.** We set out to understand when coordination among LLM agents obeys scaling trends
 233 and how much of the apparent scaling is induced by scaffolds rather than intrinsic model ability.
 234 Using *Collab-Overcooked* across asymmetric and symmetric variants, ordered cross-play pairings
 235 within Qwen 3.0 (1.7B–32B), and standardized prompts/memory, we find:

- 236 • **H₁ (Self-play scaling).** With clear, prescriptive scaffolds, self-play improves monotonically
 237 with model size and exhibits a competence threshold between small and mid-scale models.
 238 Cross-play shows the same positive trend when at least one partner is sufficiently capable
 239 (Fig. 9).
- 240 • **H_{2'} (Scaffold dependence).** As we remove role definitions the neat scaling regularities
 241 break (Fig. 4). *Scaffolding* can overstate intrinsic coordination.
- 242 • **H_{3'} (Hierarchy predicts success).** Stable leader–follower structure correlates with higher
 243 success (Fig. 5); turn order creates a leadership prior, and “bigger-as-follower” (e.g., 14 × 32,
 244 8 × 32) often outperforms the mirrored pairing.
- 245 • **H_{4'} (Parallelization amplifies coordination).** On tasks with decomposable subgoals, the
 246 ability to divide tasks among agents boosts division of labor, shortens trajectories, and
 247 strengthens scaling signals (Figs. 4c, 6b).

248 Across settings, larger explicit Assistants, or players that assume that role, are more valuable: they
 249 track global state, infer preconditions, and reduce coordination overhead even with weaker Chefs
 250 (Fig. 3). Limited cross-vendor tests suggest these effects are not idiosyncratic to a single lineage
 251 (§4.2).

252 **Implications for evaluation design.** Our results indicate that benchmark conclusions about “coordi-
 253 nation scaling” can be artifacts of wrappers. To make agent evaluations informative and reproducible,
 254 we recommend:

- 255 1. **Report scaffold details** (role scripts, turn order, communication budgets) and *vary* them via
256 a *scaffold sensitivity analysis* to test whether results are robust or scaffold-induced.
257 2. **Include parallelizable regimes** to probe genuine cooperation rather than serial plan-
258 following.
259 3. **Disaggregate by role and order** (Chef vs. Assistant; who goes first), since these systemati-
260 cally mediate scaling .

261 Taken together, H_1 , $H_{3'}$, and $H_{4'}$ describe *when* scaling is likely to appear: clear roles, fast hierarchy
262 formation, and opportunities for parallel work. $H_{2'}$ delineates *when* scaling ceases to exist: under
263 open-ended interaction without strong priors about who plans, who executes, and how to negotiate.
264 The practical upshot is that *coordination scaling is conditional*—it requires scaffolds that reduce
265 ambiguity or partners capable enough to establish structure on the fly.

266 **Future Work.** As future work, we plan to operationalize these insights in real-world settings
267 by deploying heterogeneous, web-enabled agents that act semi-autonomously. This will let us
268 systematically probe how role assignment and scaffold design shape cooperation. Operating in a
269 richer environment should surface more complex interaction patterns; our results already indicate
270 that seemingly small environmental changes—e.g., the Level 4 recipes in the symmetric condition
271 4d—can substantially change agent’s behavior. We hope this work lays the groundwork for a broader
272 study of how heterogeneous models coordinate within the emergent society of AIs (Section K).

273 **References**

- 274 F. Chen, L. Zhang, G. Pang, R. Zimmermann, and S. Deng. Synergizing large language models and
275 task-specific models for time series anomaly detection. *arXiv preprint arXiv:2501.05675*, 2025.
- 276 L. Cross, V. Xiang, A. Bhatia, D. L. Yamins, and N. Haber. Hypothetical minds: Scaffolding theory
277 of mind for multi-agent tasks with large language models. *arXiv preprint arXiv:2407.07086*, 2024.
- 278 P. M. P. Curvo, M. Dragomir, S. Torpes, and M. Rahimi. Reproducibility study of "cooperate
279 or collapse: Emergence of sustainable cooperation in a society of llm agents", 2025. URL
280 <https://arxiv.org/abs/2505.09289>.
- 281 J. V. de Carvalho Silva and D. G. Macharet. Can llm agents solve collaborative tasks? a study on
282 urgency-aware planning and coordination, 2025. URL <https://arxiv.org/abs/2508.14635>.
- 283 I. Jeknic, A. Duchnowski, and A. Koller. Collaborative problem-solving in an optimization game,
284 2025. URL <https://arxiv.org/abs/2505.15490>.
- 285 J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford, J. Wu,
286 and D. Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- 287 T. Kwa, B. West, J. Becker, A. Deng, K. Garcia, M. Hasin, S. Jawhar, M. Kinniment, N. Rush,
288 S. Von Arx, et al. Measuring ai ability to complete long tasks. *arXiv preprint arXiv:2503.14499*,
289 2025.
- 290 S. Li, V. Balachandran, S. Feng, J. Ilgen, E. Pierson, P. W. W. Koh, and Y. Tsvetkov. Mediqa:
291 Question-asking llms and a benchmark for reliable interactive clinical reasoning. *Advances in
292 Neural Information Processing Systems*, 37:28858–28888, 2024.
- 293 I. R. McKenzie et al. Inverse scaling: When bigger isn't better. *arXiv preprint arXiv:2306.09479*,
294 2023. URL <https://arxiv.org/abs/2306.09479>.
- 295 J. Newman, M. Mintrom, and D. O'Neill. Digital technologies, artificial intelligence, and bureaucratic
296 transformation. *Futures*, 136:102886, 2022.
- 297 Palantir Technologies. Aip for defense. <https://www.palantir.com/platforms/aip/>
298 defense/, 2025. Accessed: 2025-09-01.
- 299 A. Piatti, R. Dantuluri, S. Narayanan, I. Schlag, G. Zheng, B. Ravindran, et al. Emergence of
300 sustainable cooperation in a society of LLM agents. In *Advances in Neural Information Processing
301 Systems 37 (NeurIPS 2024)*, 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/hash/a7e2f8c1d2ef1b0d7c99f1b36f5a9b1c-Abstract-Conference.html.
302 NeurIPS 2024; key kept as llmsociety2023 for back-compat.
- 304 C. Preiksaitis, N. Ashenburg, G. Bunney, A. Chu, R. Kabeer, F. Riley, R. Ribeira, and C. Rose. The
305 role of large language models in transforming emergency medicine: scoping review. *JMIR medical
306 informatics*, 12:e53787, 2024.
- 307 R. Rogers. Openai adds shopping to chatgpt in a challenge to google. *WIRED*, Apr. 2025. URL
308 <https://www.wired.com/story/openai-adds-shopping-to-chatgpt/>. WIRED Gear.
- 309 H. Sun, S. Zhang, L. Ren, H. Xu, H. Fu, C. Yuan, and X. Wang. Collab-overcooked: Benchmarking
310 and evaluating large language models as collaborative agents. *arXiv preprint arXiv:2502.20073*,
311 2025. URL <https://arxiv.org/abs/2502.20073>.
- 312 M. Suzgun and A. T. Kalai. Meta-prompting: Enhancing language models with task-agnostic
313 scaffolding. *arXiv preprint arXiv:2401.12954*, 2024.
- 314 J. tse Huang, E. J. Li, M. H. Lam, T. Liang, W. Wang, Y. Yuan, W. Jiao, X. Wang, Z. Tu, and
315 M. R. Lyu. How far are we on the decision-making of llms? evaluating llms' gaming ability in
316 multi-agent environments, 2025. URL <https://arxiv.org/abs/2403.11807>.
- 317 D. Wang, Z. Ye, F. Fang, and L. Li. Cooperative strategic planning enhances reasoning capabilities in
318 large language models, 2024. URL <https://arxiv.org/abs/2410.20007>.

- 319 J. Wei et al. Inverse scaling can become u-shaped. *arXiv preprint arXiv:2211.02011*, 2022. URL
320 <https://arxiv.org/abs/2211.02011>.
- 321 Q. Wu, G. Bansal, J. Zhang, Y. Wu, B. Li, E. Zhu, L. Jiang, X. Zhang, S. Zhang, J. Liu, A. H.
322 Awadallah, R. W. White, D. Burger, and C. Wang. Autogen: Enabling next-gen llm applications
323 via multi-agent conversation. *arXiv preprint arXiv:2308.08155*, 2023. doi: 10.48550/arXiv.2308.
324 08155. URL <https://arxiv.org/abs/2308.08155>.
- 325 A. Yang, A. Li, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, C. Gao, C. Huang, C. Lv, et al. Qwen3
326 technical report. *arXiv preprint arXiv:2505.09388*, 2025. URL <https://arxiv.org/abs/2505.09388>.

328 **A Asymmetric Environment Complete Per-level Analysis**

329 In Figure 7, we show the per-level analysis of the mean similarity to the optimal policies (RAT) and
 330 success rates per level.

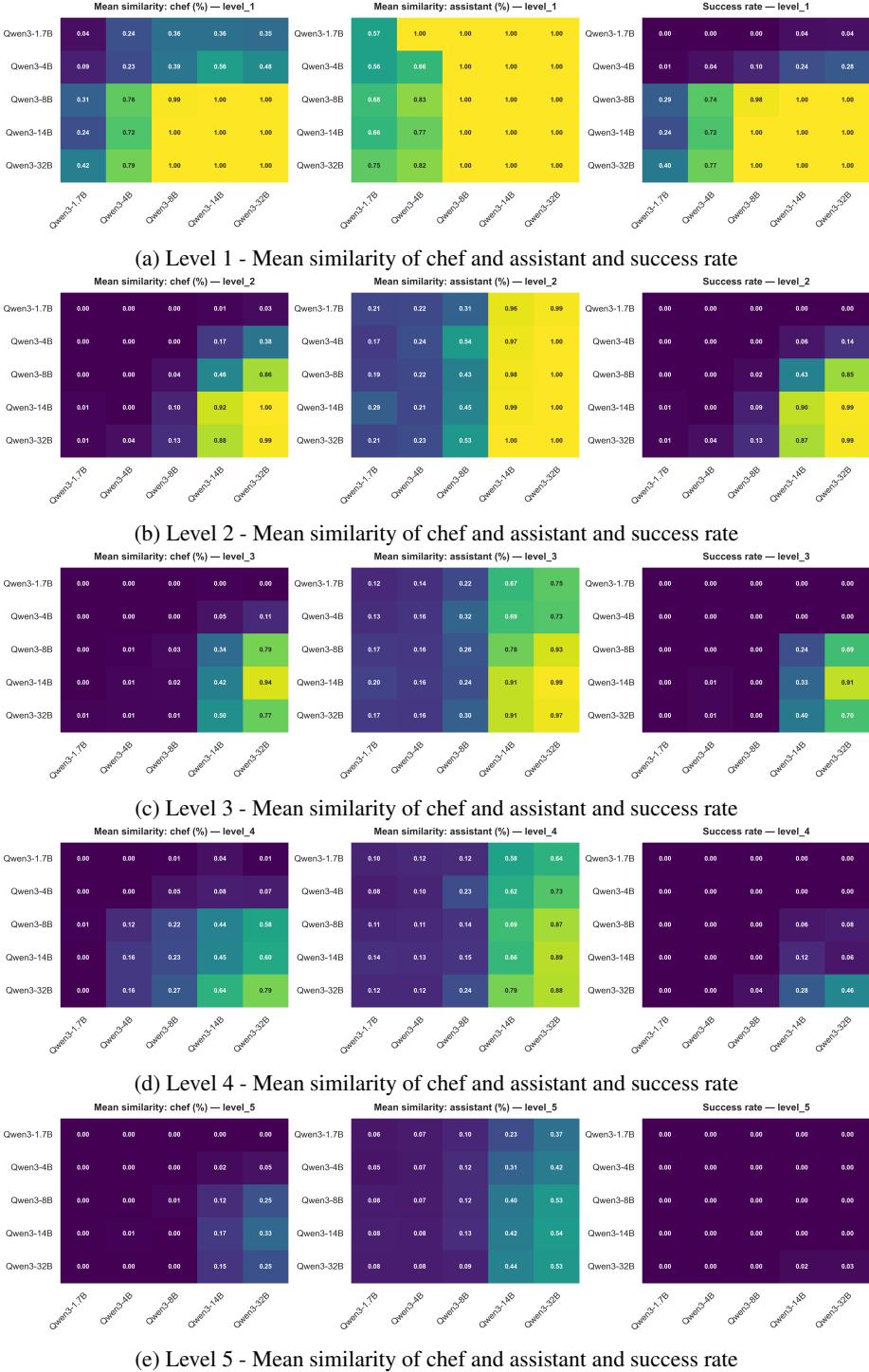


Figure 7: Mean similarity to optimal policies (RAT) and success rate in asymmetric environment

331 **B Cross vendor pairings**

332 To check that our coordination trends are not idiosyncratic to a single lineage, we paired similarly
 333 sized models across different families (Qwen, Gemma, Nemotron) under the *same* scaffold, prompts,
 334 decoding (temperature 0.7), and communication budget as in the main experiments. The heat map in
 335 Figure 8 summarizes success rates for these cross-vendor pairings.

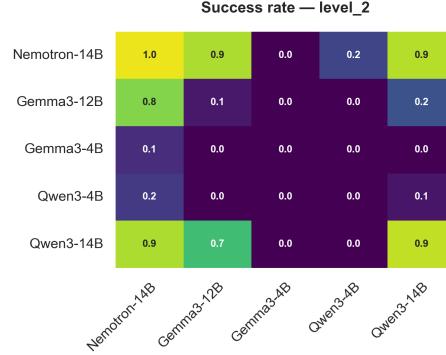


Figure 8: Success rate over similar sized models of cross vendors.

336 **C Agent’s Action Space for All Environments**

337 These are all the actions defined by the DSL that the agents have access to and have to use in order to
 338 play the underlying implementation of Overcooked.

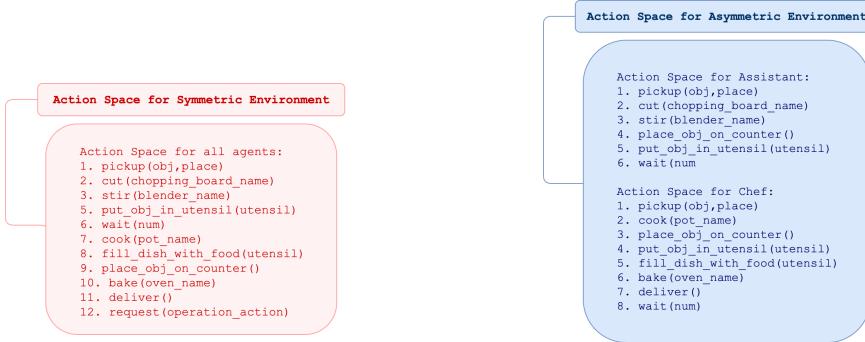


Figure 9: Action space

339 **D Rationale for Selecting Levels 2 and 4 in the Symmetric Environment**

340 We restrict symmetric-environment evaluation to Levels 2 and 4 for empirical coverage of two distinct
 341 coordination regimes. In the asymmetric setting, Level 2 shows the greatest success-rate dispersion
 342 among smaller models (cf. Fig. 7), making it sensitive to coordination differences not swamped by
 343 raw capacity. Symmetrizing Level 2 removes access constraints while preserving difficulty, letting us
 344 probe early role negotiation of a wider range of models.

345 The recipes of Level 4 incentivize the development of a concurrent algorithm in order to solve them
 346 as sub-tasks (Figure 10), making it an interesting study case for *work division* and *parallel execution*.
 347 The symmetric variant lets us assess whether agents that *can* parallelize actually do, and how this
 348 interacts with emergent hierarchy.

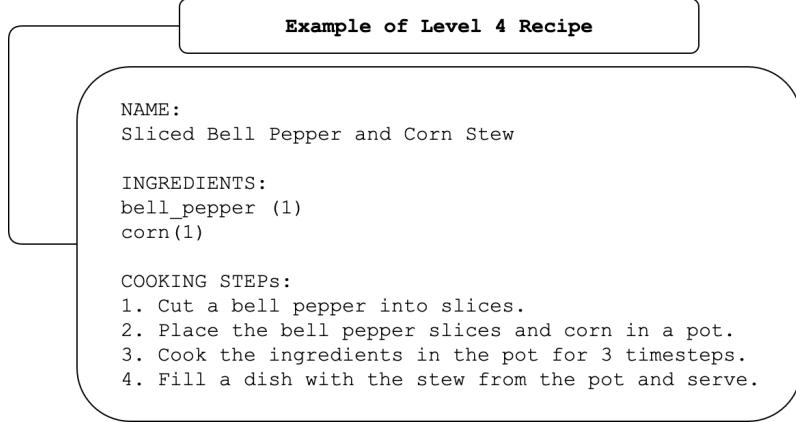


Figure 10: Example of one of the recipes used in level 4. The recipes are shared among environments.

349 E 1.7B Model Removal



Figure 11: In the majority of Qwen 1.7B runs, the agent degenerates into repeatedly echoing the DSL or other role-specific instructions—often for many turns—instead of producing task-advancing actions. This is an example of a level-2 asymmetric environment run.

350 F Larger Models as Assistants

351 Larger assistants drive most of the cross-play gains. Holding the Chef fixed, success rises with the
 352 Assistant’s scale (Fig. 2b); and Assistant RAT alignment remains high even when paired with weaker
 353 Chefs (Fig. 3b). Qualitatively, large Assistants act as robust followers: they maintain an accurate
 354 global state, infer missing preconditions, and issue targeted partner requests or execute the needed
 355 primitives—rather than reciting rules.

356 These behaviors reduce coordination overhead for the Chef and explain why “bigger-as-follower”
 357 configurations (e.g., 14 × 32, 8 × 32) often outperform their mirrors (e.g., 32 × 14, 32 × 8) in the
 358 symmetric setting (Fig. 4a).

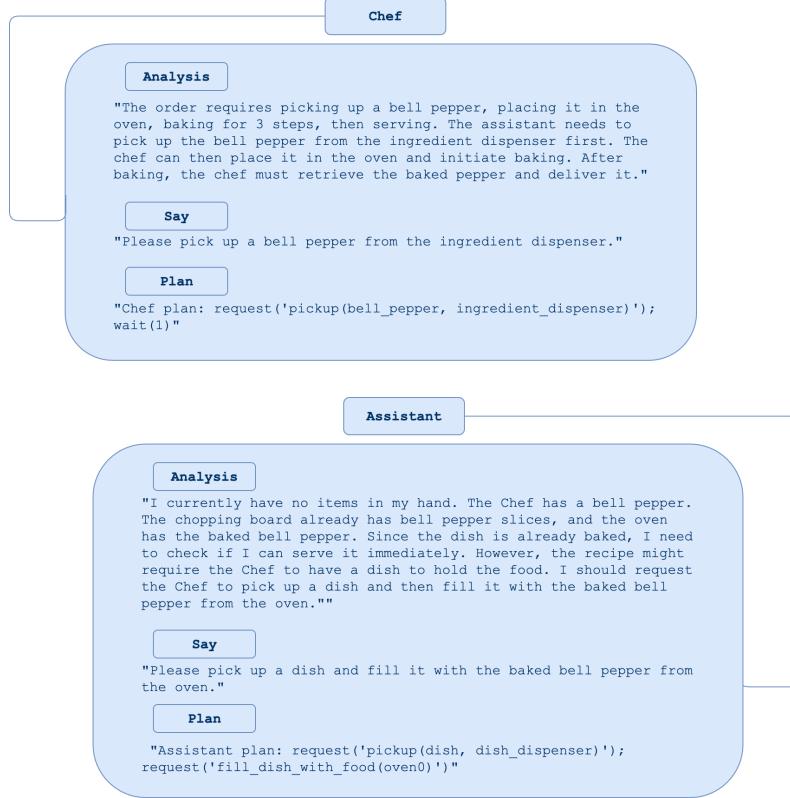


Figure 12: A 32B Assistant exhibits high-quality follower behavior: it summarizes the state (Chef holds pepper; oven contains a baked item), remembers Chef’s past states (that the Chef itself seems to have forgotten), infers a latent precondition (plating requires a dish), and issues precise partner requests that close the loop—`request('pickup(dish, dish_dispenser)');` `request('fill_dish_with_food(oven0)')`—instead of echoing the DSL. This example was taken from an asymmetric run.

359 G Top performers in the symmetric environment for level-2

360 This section outlines our hierarchy labeling procedure and reports the top-performing model pairings
 361 on the symmetric level-2 task, emphasizing the association between early role formation and success
 362 rates.

363 **Hierarchy evaluation protocol.** Across all runs in the symmetric setting, the distribution of
 364 hierarchy labels was **CLEAR_HIERARCHY** (9.4%), **FUZZY_HIERARCHY** (87.0%), and **NO_HIERARCHY**
 365 (3.6%). The labeling protocol is specified in Figure 13. For each episode, the evaluator received the
 366 base prompt plus 36 interaction snippets: 6 from the beginning, 6 from the middle, and 6 from the
 367 end *for each of the two agents*. Each snippet contained the complete reasoning, planning, and spoken
 368 turns at that time point. This yielded up to 36 analyzed steps per game; if an episode had fewer than
 369 36 steps, we analyzed all available interactions.

370 We evaluated with GPT-4-nano and GPT-4 as well, though their evaluation quality was below human-
 371 validator expectations. For GPT-5, a human audit of a 5% stratified sample of symmetric runs found
 372 the evaluations to be largely correct.

373 Table 1a lists the model pairings that most frequently exhibited a *clear* leader–follower structure in
 374 the symmetric environment. For comparison, Table 1b ranks the same set of experiments by success
 375 rate. The overlap between the two tables supports our central claim: early emergence of a stable
 376 hierarchy is positively associated with task success. Pairs that rapidly settle on complementary roles
 377 (one agent delegating, the other executing) tend to complete recipes more reliably and with fewer
 378 coordination stalls.

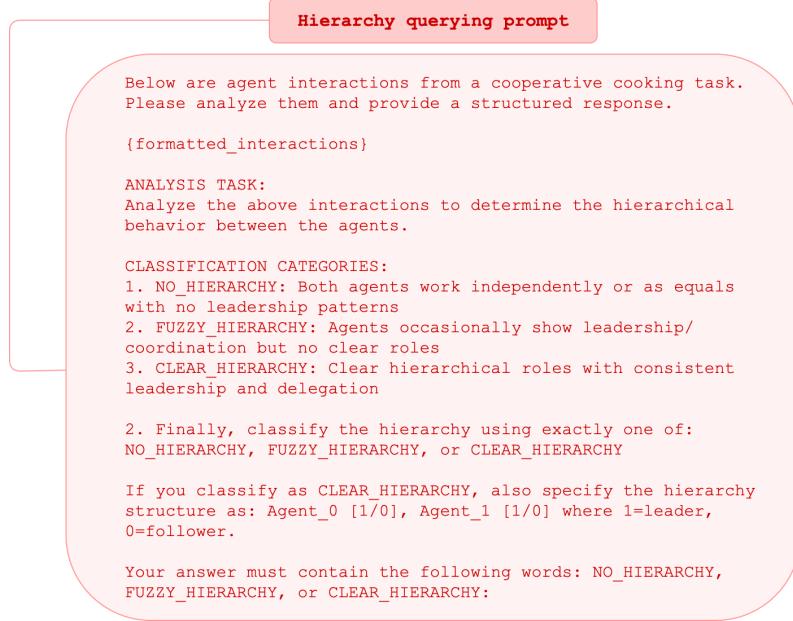


Figure 13: Prompt used to query the OpenAI GPT-5 model in order to achieve hierarchical classification.

Table 1: Ranking of the best 7 performing models in the symmetric environment in the hierarchy and success rate metrics.

(a) Models that most commonly formed clear hierarchies in the symmetric environment

Rank	Configuration	Success
1	14 × 14	0.3
2	14 × 32	0.7
3	14 × 8	0.3
4	32 × 14	0.4
5	32 × 8	0.4
6	8 × 32	0.6
7	32 × 32	0.7

(b) Highest success in the symmetric environment

Rank	Configuration
1	32 × 32
2	14 × 32
3	8 × 32
4	4 × 32
5	32 × 4
6	32 × 8
7	32 × 14

379 H Task division among players

380 For the task-division metric, we compute a role-specific RAT for each agent, align their trajectories,
 381 and compare the counts of RAT-consistent (valid) actions between agents to assess how evenly work
 382 was split in a setting without enforced role constraints.

383 **RAT adaptation for the symmetric setting.** In the symmetric environment we omit
 384 counter-placement and counter-pickup actions from the reference RAT. In such an unconstrained
 385 environment, handoffs via the central counter are not relevant. Besides, the task can be completed
 386 end-to-end by a single agent, so these logistics steps are not required by the optimal plan. Excluding
 387 them shortens the reference trajectory and thus reduces the number of actions an agent must match to
 388 achieve a perfect RAT alignment. We don't compare RAT across environments.

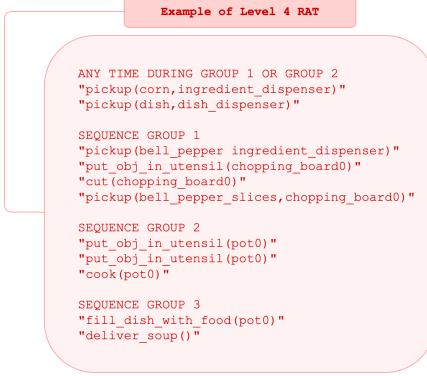


Figure 14: Example of a RAT used to compare level-4 agents actions.

389 I Mean Timestamp for Task Completion for All Environments

Table 2: Mean timestamp by model combination and recipe level 4 running on an symmetric and asymmetric environment

Model Combination	Mean Timestamp (Symmetric)	Mean Timestamp (Asymmetric)
1.7 × 1.7	119.0	119.0
1.7 × 4	113.3	119.0
1.7 × 8	115.0	119.0
1.7 × 14	112.1	119.0
1.7 × 32	116.5	119.0
4 × 1.7	119.0	119.0
4 × 4	112.9	119.0
4 × 8	109.4	119.0
4 × 14	103.7	119.0
4 × 32	95.4	119.0
8 × 1.7	116.5	119.0
8 × 4	110.4	119.0
8 × 8	106.0	119.0
8 × 14	104.3	116.9
8 × 32	101.1	117.0
14 × 1.7	116.4	119.0
14 × 4	106.7	119.0
14 × 8	100.5	119.0
14 × 14	100.0	114.9
14 × 32	95.0	115.9
32 × 1.7	116.6	119.0
32 × 4	109.2	119.0
32 × 8	99.2	118.7
32 × 14	103.2	109.0
32 × 32	99.2	97.7

390 J Prompting

391 **Thinking mode.** Qwen 3.0 supports a “thinking” mode, but pilot runs showed episode runtimes
 392 increasing by numbers as large as an order of magnitude, making evaluation impractical. Because our
 393 focus is coordination under realistic latency constraints, we disable thinking mode and use standard
 394 decoding for all models.

395 Figure 15 shows the modified prompt used to establish the task context for both players in the
396 symmetric environment.

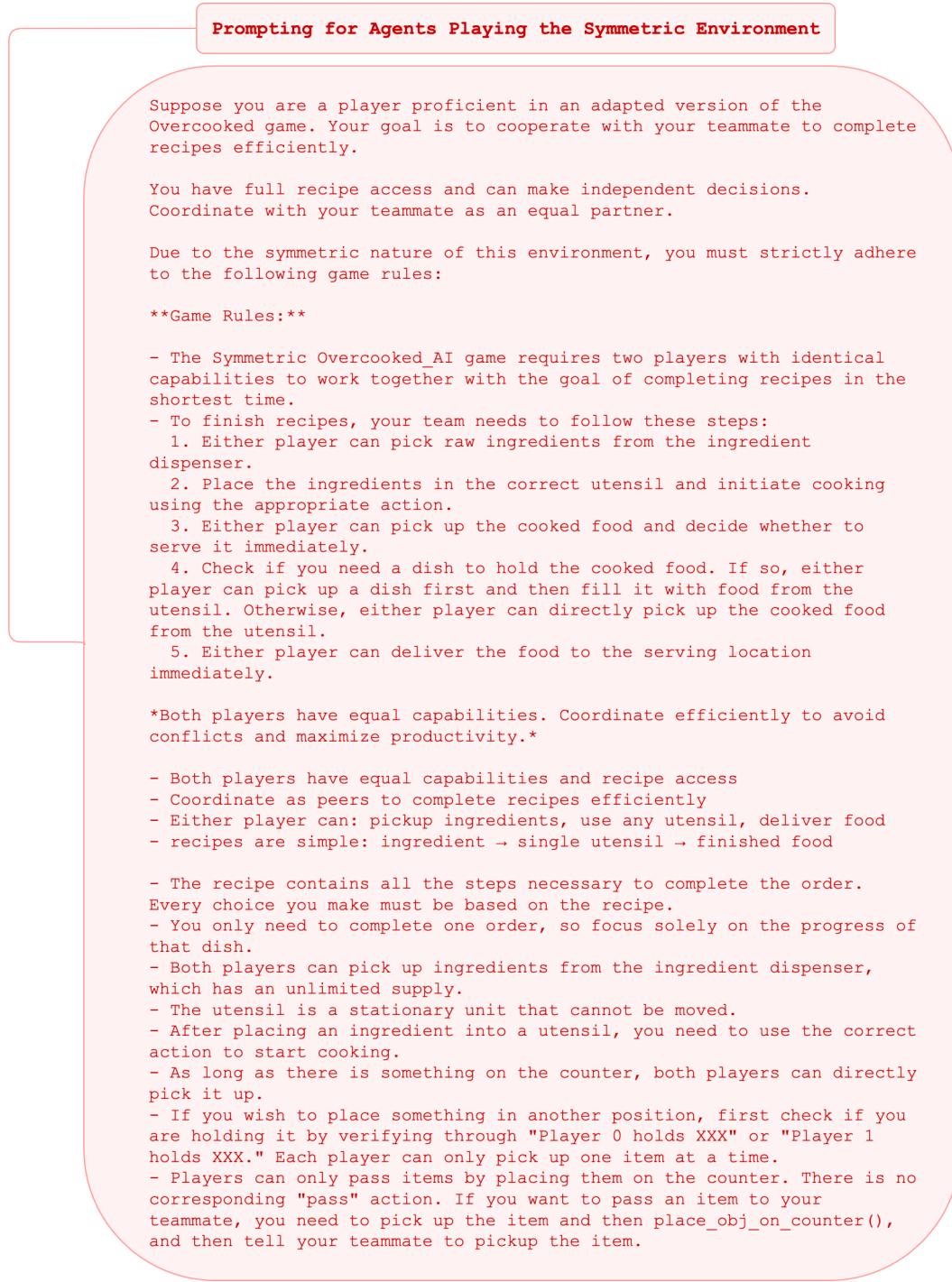


Figure 15: Symmetric environment rules' prompt. Both agents in the symmetric environment receive the same prompt.

397 **K A Society of AIs**

Box 1: What we mean by a “society of AIs”

398 A large, open ecology of heterogeneous AI agents that (i) act semi-autonomously in the world, (ii) interact repeatedly with other agents and humans, and (iii) thereby exhibit emergent, system-level behavior. This brings *distinct* failure modes—miscoordination, conflict, collusion—mediated by risk factors such as information asymmetries, network effects, selection pressures, destabilising dynamics, commitment and trust, emergent agency, and multi-agent security. Our experiments treat Overcooked scaffolds as micro-institutions and study how coordination scales (or fails) under different institutional choices.

399