ROUNDTABLE: INVESTIGATING GROUP DECISION MAKING MECHANISM IN MULTI-AGENT COLLABO RATION

Anonymous authors

006

008 009 010

011

013

014

015

016

017

018

019

021

025

026 027 028

029

Paper under double-blind review

Abstract

This study investigates the efficacy of Multi-Agent Systems in eliciting crossagent communication and enhancing collective intelligence through group decision-making in a decentralized setting. Unlike centralized mechanisms, where a fixed hierarchy governs social choice, decentralized group decision-making allows agents to engage in joint deliberation. Our research focuses on the dynamics of communication and decision-making within various social choice methods. By applying different voting rules in various environments, we find that moderate decision flexibility yields better outcomes. Additionally, exploring the linguistic features of agent-to-agent conversations reveals indicators of effective collaboration, offering insights into communication patterns that facilitate or hinder collaboration. Finally, we propose various methods for determining the optimal stopping point in multi-agent collaborations based on linguistic cues. Our findings contribute to a deeper understanding of how decentralized decision-making and group conversation shape multi-agent collaboration, with implications for the design of more effective MAS environments.

1 INTRODUCTION

Collaboration is a fundamental aspect of the nature and human society. Whether among humans 031 or animals, working together allows groups to overcome individual limitations and achieve greater 032 collective outcomes. In nature, collaboration often arises as a strategy to boost survival, enhance 033 resource gathering, or increase efficiency in completing tasks (Schmidt & Mech, 1997). Similarly, 034 in human societies, collaboration drives innovation, facilitates problem-solving, and fosters shared 035 understanding, enabling individuals to address complex challenges that would be unmanageable otherwise (De Man & Duysters, 2005; Graesser et al., 2018; Bittner & Leimeister, 2013). This innate 037 tendency to collaborate is evident across various domains, from social communities to technological 038 systems, where multiple entities coordinate their efforts toward a common goal. As we advance in developing intelligent agents, understanding and replicating these collaborative dynamics in artificial systems has become increasingly important, predominantly to cope with the complexity and 040 adaptability seen in real-world interactions. 041

Agents powered by Large Language Models (LLMs) have demonstrated impressive problem-solving capabilities across a wide range of tasks. However, single-agent systems encounter significant difficulties when tasked with problems that are either too large or complex, often resulting in instability, misalignment with the intended request, and hallucination (Liu et al., 2024; Kuhn et al., 2023; Lyu et al., 2023). To address these limitations, research has increasingly turned toward Multi-Agent Systems (MAS). MAS have shown greater efficacy in harnessing collective intelligence by allowing individual agents to specialize in distinct skills and facilitating effective collaboration among them (Guo et al., 2024).

When agents working together in a MAS, it is natural for them to have varying interpretations and
 perspectives. While some opinions may align, disagreements are also frequent. This creates an in evitable tension between cooperation and competition, stemming from differences in backgrounds,
 information access, and individual goals. Therefore, the process of aggregating models' diverse pre dictions into a final group decision becomes a crucial aspect of dynamic multi-agent collaboration.



Figure 1: Overview of our multi-agent collaboration platform: RoundTable. It uses a round-based collaboration where agents simultaneously send messages, propose solutions, and vote. Based on a social choice mechanism, RoundTable selects the most preferred proposal for group decisions. Brief introduction of RoundTable is in Section 3.1, details are in Appendix A.1.

072 Many existing LLM-driven MAS are designed based on centralized group decision-making, which 073 typically involves layered or centralized architectures with a hierarchy. When there is a conflict 074 between agents, it is resolved by a pre-assigned agent or a pre-defined process. These systems are 075 often used for tasks where agent networks follow the waterfall method, leading to a stratified ar-076 rangement among agents (Hong et al., 2023; Qian et al., 2023; Dong et al., 2023). However, this 077 hierarchical setup poses several critical challenges: (1) fairness: individual agent messages may not be accurately represented in the final group outcome, leading to potential misrepresentation (Jiang & Lu, 2019); (2) rigidity: the system's fixed structure may lead to over-fitting to specific scenarios 079 (Chen et al., 2023), reducing its adaptability to diverse and dynamic environments; and (3) bias: the 080 agent with final decision-making power may introduce its own biases into the process, potentially 081 distorting the outcomes (Owens et al., 2024). Additionally, centralized MAS cannot perform in environments requiring independent agent decisions, where private or incomplete information hinders 083 a central authority from dictating optimal strategies (Xu et al., 2023). 084

085 Decentralized group decision-making can ease these issues by distributing power among agents, where each agent has the ability to participate in the process. This is a common setting in world 086 simulation and embodied environment, where agents need to behave independently because there 087 exists information asymmetry or data boundary between agents (Mandi et al., 2024; Zhang et al., 088 2023a; Xu et al., 2023; Park et al., 2023), and possibly variability in capabilities that different agents 089 have. With decisions made by multiple agents, the flexible structure of decentralized MAS adapts to 090 various environments, but the lack of a fixed hierarchy demands careful monitoring of collaboration 091 patterns. 092

Centralized and decentralized group decision-making attempt to mimic the ways in which human societies form collective policies, such as in monarchies and democracies. These decision-making mechanisms, known as social choice, are studied across various fields including economics, mathematics, philosophy, and social science, with the goal of aggregating and synthesizing individual preferences into a unified consensus. However, the impact of social choice methods on LLM-based MAS has yet to be explored.

In this study, we evaluate various social choice methods across different environments to observe and
 analyze agents' group behaviors and collaboration pattern. This paper aims to provide the following
 research contributions:

102 103

107

067

068

069

- We investigate how collaborative behavior in decentralized MAS varies across social choice methods, providing insights into how they influence overall cooperation and outcomes.
- We identify key language features in multi-agent conversations as indicators of collaboration, offering a novel approach to analyzing linguistic cues in effective or ineffective interactions.
 - We propose various methods for determining the optimal stopping point in multi-agent collaboration, utilizing the linguistic features we identified.

108 2 RELATED WORKS

109 110

Multi Agent Frameworks Recent advancements in MAS have led to the development of nu-111 merous platforms that facilitate collaboration among multiple agents. AutoGen (Wu et al., 2023) 112 introduces a framework for building LLM-based applications where agents communicate with each 113 other to accomplish tasks. CAMEL (Li et al., 2023a) is designed to enable autonomous collabora-114 tion among chat-based language models using role-playing and inception prompting. ChatDev (Qian 115 et al., 2023) offers a chat-driven framework where LLM-powered agents work together to streamline 116 software design, coding, and testing through unified language communication. MetaGPT (Hong et al., 2023) provides a meta-programming approach that encodes Standardized Operating Proce-117 dures into prompt sequences, facilitating efficient and error-reducing collaboration among LLM-118 based agents. AgentVerse (Chen et al., 2023) orchestrates collaborative groups of expert agents, 119 drawing inspiration from human group dynamics to enhance task performance beyond the capabili-120 ties of individual agents. These MAS platforms address various aspects of the collaborative process, 121 including agent profiling, recruitment, memory, planning, and communication protocols. However, 122 they do not specifically investigate the group decision-making mechanisms within multi-agent inter-123 actions as we do in this paper. Some platforms employ centralized methods, and others only rely on 124 unanimous or majority voting.

125

126 Multi-Agent Collaboration Environments Multi-agent collaboration has been used to handle 127 various applications. In robotics, there has been a continuous study regarding multi-agent collab-128 oration in embodied environments (Dudek et al., 1996; Ota, 2006; Cena et al., 2013; Vorotnikov 129 et al., 2018). In these studies, multi-agent collaboration in robotics has been explored through var-130 ious frameworks, such as task allocation, coordination strategies, and communication protocols, to 131 optimize collective performance and adaptability in dynamic environments. For LLM-based MAS, 132 software development is a popular environment which inherently requires diverse expertise, continuous integration, and coordinated teamwork to build complex, functional systems efficiently (Qian 133 et al., 2023; Hong et al., 2023; Dong et al., 2023). Additionally, LLM's strong context understand-134 ing and immense parametric knowledge supported various world simulations with multi-agents. 135 Including simulating society (Park et al., 2023; 2022), economy (Li et al., 2023b; Zhao et al., 2023), 136 gaming (Xu et al., 2023; Wang et al., 2023), and psychology (Aher et al., 2023; Zhang et al., 2023b). 137 These works require agents to interact with others in a group, and investigate their cooperate and 138 competitive behavior. 139

139 140 141

142

3 MULTI-AGENT COLLABORATION AND GROUP DECISION-MAKING

In this section, we first introduce the formal definition of the setting this study uses for MAS, and
 propose RoundTable, a turn-based multi-agent collaboration platform for evaluating different de centralized group decision-making.

146 In a collaboration environment, there is a set of individuals (agents) $i \in \mathbb{I}$, and a set of all 147 possible states of the world $\mathbb{X} = \{x_1, x_2, x_3, ...\}$. Then there is an individual utility function 148 $u_i: \mathbb{X} \to \mathbb{R}$ for each agent that maps a state to a real number, representing preferences of 149 agents. Based on the individual utility function, each agent proposes a state with proposal function 150 $p_i = f_p(u_i, C_i, \mathbb{X}), p_i \in \mathbb{X}$, where C_i is the context, including environment information, back-151 ground, and collaboration history. Similarly, agent can also vote on proposals with voting function 152 $v_{ij} = f_v(u_i, C_i, \mathbb{S} | p_i \in \mathbb{S} \subseteq \mathbb{X}), i, j \in \mathbb{I}$, where \mathbb{S} represents a candidate list of proposals to vote on. With all proposals and agents' votes, a social choice method (function) F decides one group decision 153 $x^* = F(p_i, v_{ij} | i, j \in \mathbb{I}).$ For example, when $F = Majority, p_1 = Apple, p_2 = Banana, p_3 =$ 154 *Carrot*, $v_{1,i}$, $v_{2,i} = (Yes, No, No)$, $v_{3,i} = (No, Yes, No)$, then $x^* = p_1 = Apple$. 155

156

158

157 3.1 ROUNDTABLE: MULTI-AGENT COLLABORATION PLATFORM

With the definition above, we propose RoundTable, a multi-agent collaboration platform that can take various group decision-making mechanisms. The overview of RoundTable is shown in Figure
Due to the limited space, the details of the platform design and the LLM prompt are shown in Appendix A.1 and A.2. Following is a brief introduction of each phase.

162 • Input and Initialization: RoundTable takes initial query (task), and initializes each agent giving 163 agent specific background and utility function (or goals). 164 • Collaboration Round: RoundTable employs round-based agent collaboration, with each round 165 comprising three phases. Each phase occurs in a simultaneous open conversation, where agents 166 act without order, and information is shared with everyone at the end of each phase. The process 167 ends after R rounds. 168 - Message Phase: Each agent sends a message to its intended recipients. 169 - Proposal Phase: Each agent proposes a potential solution, with an option to skip. 170 - Voting Phase: Each agent votes for the candidate list, with an option to skip. The list 171 consists with each agent's latest proposal and latest accepted proposal by the group. With 172 all votes, social choice method chooses the winner. If ties or disqualifies, the decision is 173 deferred. 174 • Output: After the final round, the latest accepted proposal will be selected as the final output of 175 this system. 176 177 3.2 SOCIAL CHOICE METHODS 178 179 To investigate the impact of different social choice methods on group decision-making in multi-agent collaboration, we compare the following 6 mechanisms: 181 182 • Unanimous Voting(Arrow, 2012): The proposal that receives votes from all agents will be se-183 lected. • Majority Voting(Arrow, 2012): The proposal that receives votes from more than half of all agents 185 will be selected. • Plurality Voting(Arrow, 2012): The proposal that receives the most votes will be selected. 187 • *Rated Voting*(Baujard et al., 2018): Each agent assigns ratings on a 5-point Likert scale to all 188 candidate proposals, with 1 being the lowest and 5 being the highest. The proposal with the 189 highest total score will be selected. 190 • Ranked Voting(Arrow, 2012): Each agent ranks all candidate proposals from the most preferred 191 to the least preferred. Social Choice will assign 1, 1/2, 1/3... points to the 1st, 2nd, 3rd... 192 candidates on each ballot. The proposal with the highest total points will be selected. 193 • Cumulative Voting (Black et al., 1958): For $|\mathbb{I}|$ candidate proposals, each agent is given $|\mathbb{I}|$ points 194 to distribute among the proposals as they see fit. The proposal with the highest total points will 195 be selected. 196 197 Unanimous, majority and plurality voting are one-vote mechanisms, where agents only choose the best candidate from the list; rated, ranked and cumulative voting are score-based mechanisms, which 199 ask a nuanced, gradient preference over candidates. 200 201 4 EXPERIMENTS 202 203 Using RoundTable, we explore multi-agent behavior patterns in R = 10 rounds of collaboration. 204 We evaluate RoundTable in two environments: simulated and complex. 205 206 SIMULATED ENVIRONMENT - EXCHANGE ECONOMY 207 4.1 208 **Introduction** We use an exchange economy for the simulated environment because its advantages 209 make it well-suited for evaluating a decentralized MAS(Varian, 1992). First, it is a plus-sum game 210

make it well-suited for evaluating a decentralized MAS(Varian, 1992). First, it is a plus-sum game where agents seek equilibrium, meaning collaboration opportunities exist if allocation is not yet optimal. Second, multiple equilibria in the market allow dynamic collaboration. Third, while equilibrium doesn't guarantee maximum utility, it must lie within one of the equilibria, helping to assess whether conflicts between agents hinder the group's progress toward the ultimate goal.

Here, K goods and K agents participate in a market, with each good having a quantity of 100. Each agent has a Cobb-Douglas utility function $u_i = \prod a_k^{\theta_k}$, $\sum \theta_k = 1$, where a_k represents the amount

216 of good k for agent i, and θ_k reflects the relative preference for each good (Cobb & Douglas, 1928). 217 Agents aim to maximize their individual utility by finding an optimal allocation of goods. From a 218 group perspective, the hidden goal is to maximize the total utility $U = \sum u_i$, which is not revealed 219 to agents. For more on the exchange economy environment, see Appendix C.2.

220 To reflect common MAS collaboration patterns, we use an asymmetric utility function setup with 221 K = 3 agents, where each agent prefers different goods, mimicking real-world scenarios where 222 agents specialize in different areas. Each agent's utility function is $u_i = a_i^{0.8} \prod_{k \neq i} a_k^{\theta}$, where 223 $\tilde{\theta} = \frac{1-0.8}{|\mathbb{I}|}$. Other types of utility function sets are shared in C.3. We report performance as the 224 average of 100 simulations. 225

Metrics We use various metrics to analyze multi-agent collaboration. For quality, we report the 227 group total utility $U = \frac{\sum u_i}{U_{max}}$, and U_{max} is the largest possible utility achievable in the environment. 228 Efficiency is measured by the area under the curve (AUC), $AUC@n = \sum_{r=1}^{n} \frac{U_r}{U_{max}}$, where r is 229 230 round. Fairness is assessed by the Min/Max ratio, which compares the smallest individual utility to the largest at round 10: $\frac{u_{min}}{u_{max}}$. Rationality is the ratio of rounds where the proposed allocation's utility exceeds the current one.¹ Finally, rigidity measures how often an old allocation is kept, 231 232 $rigidity = \frac{1}{10} \sum_{r=1}^{10} \mathbf{1}_{U_r=U_{r-1}}.$ 233

234 235

236

226

4.2 COMPLEX ENVIRONMENT - RECOMMENDATION SYSTEM

237 **Introduction** For the complex environment, we focus on a recommendation system. It reveals 238 unique collaborative behaviors absent in exchange economies. First, strong information asymmetry 239 exists, as no single piece of information can accurately predict ratings. Second, group decisionmaking fosters diverse problem-solving approaches, with agents both consuming and contributing 240 information in a feedback loop that enhances recommendations. 241

242 The task is to predict user ratings for specific items by analyzing historical interactions and identify-243 ing data patterns. This task is challenging for a single-agent system due to the overwhelming amount 244 of scattered and limited essential information. In our setting, we divide the background information 245 into several parts and assign each to different agents. The agents must then collaborate by exploring, analyzing, exchanging, and synthesizing their information to reach a comprehensive conclusion. 246

247

Dataset We evaluate MAS performance using the MovieLens-100k dataset (Harper & Konstan, 248 2015). The dataset consists of three main tables: *user*, which contains demographic information 249 of all users; *item*, which includes details of all movies, such as title, release date, and genre; and 250 data, which holds 100,000 ratings from 943 users on 1,682 movies. Due to resource constraints, we 251 randomly selected 100 examples from the u1.test split for evaluation. To simplify the data structure, 252 we pre-processed it into three tables: *basicInfo*, *userHistory*, and *movieHistory*. Each table was 253 assigned to a different agent, creating a 3-agent collaboration system. The definition and example 254 of each table can be found in Appendix C.4. 255

256 Metrics We use Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE) to represent 257 utility U_r , assessing the accuracy of predictions compared to the gold rating at round r. To evaluate 258 MAS performance in the recommendation task, we compare it against two baselines: a simple Always Guess 4, which predicts the median regardless of input, and a strong State of the Art (SoTA) 259 model from Behera & Nain (2023), which employs collaborative filtering with temporal features. 260

CROSS-AGENT CONVERSATION ANALYSIS 4.3

To understand agent collaboration and analyze their messages, we perform a language analysis focused on four key features. All observations are collected in 3 agents, gpt-4o-mini setting.

- Message Length is a basic yet significant metric, reflecting the amount of information an agent conveys. We measure it using word count, as longer messages generally enrich the conversation.
- 267 268 269

261

262

263

264

¹This differs from compromise, where proposing lower utility than the previous proposal is natural, but a lower utility than the current allocation is irrational.

270 • Message Complexity assesses how difficult a message is to understand, indicating the depth of 271 the conversation. We use the Flesch-Kincaid grade level to calculate complexity, with higher 272 scores representing more intricate messages (Klare, 1974).

- 273 · Information Difference measures how much new information is introduced. Effective collaboration should consistently bring new insights, while low information gain signals a stalled conversation. It is calculated as the average cosine distance between messages in the current round and the center embedding of the previous round.
 - Dialogue Acts are communicative functions that capture actions, intents, and behaviors within messages as discrete states. We design a set of dialogue acts for multi-agent collaboration: Inform, Request, Confirm, Summarize, Evaluate, Propose, Compromise, Defend, Accept, Decline, and use LLM-labeling to perform multi-label classification on a target message based on the previous round. Details of the dialogue act labels and prompt can be found in Appendix C.5.
 - 4.4 DIALOGUE ACT TRANSITION GRAPH

We constructed transition graphs to visualize the dynamics of dialogue act transitions across rounds, 285 highlighting the most frequent transitions. Each node represents a dialogue act, and directed edges 286 indicate the probability of transitions between acts from one round to the next. The figure displays 287 only the most probable edge per node, excluding self-loops. The detailed definition of nodes and 288 edges is in Appendix C.6. 289

290 291

296

297 298

299

300

301

302

303

304

305

306

307

308

310 311

274

275

276

277

278

279

280

281

282 283

284

4.5 EARLY STOPPING IN MAS

292 In the previous experiments, we enforced 10 rounds of collaboration to compare agent behaviors 293 and identify patterns. However, effective early stopping can prevent redundancy, avoid stagnation, 294 and optimize decision-making by reducing unnecessary iterations. In this experiment, we evaluate 295 the following early stopping methods:

- @10 is our baseline, the final performance after 10 rounds without early stopping.
- *First Agreement* stops collaboration whenever a proposal passes the social choice criterion.
- Consecutive Agreements is when no one make additional proposal after a proposal has been accepted in the previous round.
- Validation Checkpoint is the average number of rounds that produced the best outcomes in the train set, used as an early stopping criterion for all test sets.
- Information Difference is the average embedding distance captured at the round that produced the best outcomes in the train set, and stops collaborations in the test set when the distance became lower than the threshold. This idea is supported by the observation in Section 5.3.1.
- Dialogue Act method utilizes pairs of dialogue acts, linking one from the previous round to another from the current round. We perform ordinary least squares (OLS) regression on all such pairs in relation to the performance. The regression coefficients indicate the most impactful pairs that act as stop signals. Further details of the algorithm are provided in Appendix C.7.

We use 5-fold cross validation to compare the methods with Oracle: the oracle performance for 312 early stopping, which reflects the best outcomes from each test simulation.

313 314

315 316

317

- 5 RESULTS
- 5.1 EXCHANGE ECONOMY

318 The experiment results show distinct performance differences across social choice methods in Table 319 1 and Figure 2(a). Score-based mechanisms achieve higher performance early in the collaboration 320 process compared to one-vote mechanisms, indicating that nuanced evaluations better aggregate 321 agent preferences. One-vote mechanisms struggle with lower and less stable early performance due to increased disagreement, which impairs decision-making. Despite differing early trajectories, all 322 social choice methods reach high-quality outcomes by the end. In most cases, performance peaks 323 in the middle rounds, suggesting that early stopping could be advantageous. Rationality scores

338

340

341

342

343

344

345

347

351

352

Table 1: Exchange economics environment results across different social choice methods. The 325 reported numbers are an average of all simulations, and the numbers in parentheses are standard 326 errors. The smallest and largest value in a category is colored in blue and red. We see that one-327 vote mechanisms show lower and less stable early performance, due to the higher probability of 328 disagreement.

	GROUP TOTAL UTILITY@3	GROUP TOTAL UTILITY@5	GROUP TOTAL UTILITY@10	AUC@3	AUC@5	AUC@10	RATIONALITY	MIN/MAX	RIGIDITY
			Mod	lel: gpt-4o-mir	i, # Agent: K	= 3			
Unanimous	33.94 (3.96)	43.92 (3.96)	48.48 (3.86)	21.77 (2.86)	30.17 (3.05)	38.65 (3.32)	35.00 (1.73)	87.96 (1.85)	93.00 (0.61)
Majority	79.88 (0.80)	81.33 (0.72)	79.61 (0.88)	64.08 (1.81)	70.91 (1.23)	75.67 (0.84)	23.80 (0.94)	64.00 (3.07)	71.30 (1.32)
Plurality	78.70 (1.32)	79.42 (1.12)	76.95 (1.24)	64.17 (1.85)	70.34 (1.37)	74.26 (1.11)	26.50 (1.10)	58.11 (3.29)	69.10 (1.43)
Rated	80.83 (0.49)	80.85 (0.97)	79.91 (0.83)	74.02 (0.92)	76.83 (0.68)	78.63 (0.59)	19.80 (0.89)	67.02 (2.97)	66.70 (1.56)
Ranked	80.92 (0.62)	81.41 (0.73)	78.40 (1.46)	77.31 (1.00)	78.89 (0.74)	79.41 (0.68)	19.03 (0.85)	65.10 (3.05)	61.40 (1.93)
Cumulative	79.05 (1.11)	81.10 (0.89)	78.45 (1.15)	68.60 (1.62)	73.63 (1.11)	76.48 (0.86)	23.13 (1.05)	59.24 (3.36)	65.50 (1.60)





(a) Social choice comparison, using 3 agent, asymmetric, gpt-4o-mini setting.



Figure 2: Exchange economy environment results in group total utility. Line plots (left y-axis) show 348 the group total utility achieved over rounds; bar plots (right y-axis) represent the ratio of cases where 349 participants yet to reach an agreement until a certain round. The red horizontal line indicates the 350 maximum achievable group total utility. We observe that decentralized collaboration performs well across various social choices and LLMs.

353 354 show that agents often make suboptimal choices, but decentralized MAS still performs well overall. 355 Score-based mechanisms allow more dynamic decision updates, while one-vote mechanisms exhibit greater rigidity, particularly with unanimous voting. 356

357 When comparing one-vote mechanisms, unanimous voting performs the worst due to its inflexibility 358 and need for total agreement, while majority voting slightly outperforms plurality. Additionally, 359 fairness gradually decreases from unanimous to plurality voting. These observations suggest that, 360 compared to the rigid unanimous method or the looser rules of plurality voting, a moderate level of 361 decision flexibility leads to better outcomes, albeit with some sacrifice in fairness across agents.

- 362 For the results of ablation studies, please refer to Appendix D.1. 363
- 364 5.2**RECOMMENDATION SYSTEM**

366 We illustrate the performance of recommendation system in Figure 3(a), 3(b) and Table 7. The 367 results shows that multi-agent collaboration achieves moderate performance with distributed infor-368 mation, with the decentralized MAS performing comparably to SoTA approaches. While the SoTA method was purely machine-learning-driven, the MAS approach showed reasonable outcomes, in-369 dicating that collaboration between agents can explore and produce meaningful signals that help to 370 make decision. Multi-agent collaboration with access to SoTA models may lead to better perfor-371 mance, but we leave this experiment for future work. Similar to the exchange economy, multiple 372 results in the recommendation system follow a V-shaped curve. This suggests that while agent col-373 laboration was efficient early on, it declined in later rounds, highlighting diminishing returns from 374 repeated interactions. The trend underscores the need to balance sustained collaboration within 375 MAS over time. 376

When comparing one-vote mechanisms to score-based mechanisms, we found that one-vote ap-377 proaches generally worsened or plateaued after several rounds of collaboration. In contrast, score-



Figure 3: Comparison between social choices in recommendation system environment. Line plots (left y-axis) show the comparison of MAE/RMSE achieved over rounds; bar plots (right y-axis) represent the ratio of cases where participants yet to reach an agreement until a certain round.

based mechanisms showed consistent improvement over time, suggesting that score-based mechanisms are better suited for scenarios requiring continuous improvement, allowing for more nuanced input and refinement over multiple rounds.

5.3 CROSS-AGENT CONVERSATION ANALYSIS

398 5.3.1 BASIC STATISTICS

400 Table 8 shows the average values of message length, complexity, and information difference per 401 round. Message length increases over time in both environments, with shorter messages in early rounds and longer ones later. The recommendation system generally has longer messages due to 402 its need for detailed explanations. Message complexity in the exchange economy starts low, peaks 403 early, dips, and then gradually increases toward the final round, while the recommendation system 404 shows consistently higher and increasing complexity throughout. Information difference steadily 405 decreases in both environments, indicating a convergence of topics and less novel information as the 406 discussions progress. 407

The dialogue act annotation results are presented in Table 9 in Appendix D.4. The dialogue act 408 annotation results in both environments show that *Inform* and *Request* acts dominate, indicating 409 frequent information sharing and input requests critical for task progression. The *Confirm* act rises 410 sharply after round 3 in the recommendation system and more gradually in the exchange economy, 411 showing validation is more prominent in recommendation tasks. *Summarize* acts increase slightly 412 in later rounds, consolidating information, while *Evaluate* acts remain consistently high, reflecting 413 ongoing assessment. Propose acts surge early in the exchange economy and peak later in the rec-414 ommendation system, suggesting early proposals are vital in negotiations. Acts like *Compromise*, 415 Defend, Accept, and Decline are rare, reflecting minimal adversarial behavior, while Others remain 416 low, indicating the high quality of dialogue act definitions.

417 418

387 388

389

390

391 392

393

394

395 396

397

5.3.2 DIALOGUE ACT TRANSITION GRAPH

419 The most probable transition graphs in both environments are shown in Figure 4(a) and 4(b). The 420 breakdown of social choices are displayed in Figure 5 and 6. The transitions reveal a structured 421 progression from requesting information to proposing and evaluating solutions, followed by resolu-422 tion or negotiation. In the exchange economy, collaboration centers around a loop between Request 423 and Propose, reflecting the cooperative decision-making process. Key transitions include Inform to 424 Request, Confirm to Request, and Compromise to Propose, showing how agents share knowledge, 425 confirm information, and adopt ideas. Accept transitions to End, signaling negotiation closure. In 426 contrast, the recommendation system centers on *Request*, indicating a more cooperative, exploratory 427 process with incomplete information.

428

- 429 5.4 EARLY STOPPING IN MAS
- 431 We show comparison of early stopping methods in recommendation system in Table 2, full results in Table 10, and their basic statistics are in Table 11 in Appendix D.5. To begin, V-shaped performance





(a) Dialogue act transition graph for exchange economy.

(b) Dialogue act transition graph for recommendation system.

Figure 4: Dialogue act transition graph for exchange economy and recommendation system environments. Only the most probable outgoing edge is presented. Self-loops are excluded.

Table 2: Comparison between early stopping methods in recommendation system environment. The performance is compared with *Oracle* and baseline @10 in MAE. (\downarrow) indicates better performance with higher and lower values, respectively. Experiments are based on the results of 3 agents, gpt-4omini setting. Results are based on 5-fold cross validation. We observe that language-based methods performed well overall.

452	EARLY STOPPING METHOD	UNANIMOUS	MAJORITY	PLURALIRY	RATED	RANKED	CUMULATIVE
452		RECOMM	ENDATION SY	STEM, in MAE	(1)		
433	@10 (Baseline)	0.88	0.86	0.84	0.79	0.84	0.76
454	First Agreement	0.82	0.80	0.81	0.82	0.89	0.82
455	Consecutive Agreements	0.84	0.81	0.84	0.77	0.82	0.77
400	Validation Checkpoint	0.81	0.81	0.82	0.80	0.83	0.82
456	Information Difference	0.86	0.78	0.85	0.79	0.78	0.81
457	Dialogue Act	0.84	0.82	0.79	0.76	0.85	0.82
437	Oracle	0.73	0.63	0.59	0.62	0.69	0.67
458							

460 trends were observed across most cases, where the best outcome is observed before reaching the 461 final checkpoint (@10), indicating that intermediate stopping leads to better results than allowing the process to fully complete. This pattern supports the necessity of early stopping in multi-agent 462 collaboration settings, where overextending the interaction can lead to diminishing returns in per-463 formance. 464

465 Almost all early stopping methods outperformed the baseline (@10), reinforcing the effectiveness 466 of early termination in improving outcomes. Notably, methods leveraging linguistic features, such 467 as Information Difference and Dialogue Act, delivered better performance across different social choices and environments. These findings highlight the advantage of incorporating dialogue-based 468 metrics in deciding when to stop, especially in complex multi-agent environments. Their results 469 across both the exchange economy and recommendation system scenarios further supports the claim 470 that linguistic indicators can be powerful tools for optimizing collaboration outcomes, reducing 471 cognitive load, and improving decision efficiency. 472

473

432

433

434

435

436

437

438

439 440

441

442 443

444

445 446

447

448

449

450

451

459

DISCUSSION 6

474 475 476

477

IMPACT OF SOCIAL CHOICE METHODS ON COLLABORATIVE BEHAVIOR 6.1

The study demonstrates that social choice methods significantly impact multi-agent collaboration, 478 with score-based mechanisms achieving higher performance and efficiency, especially in early 479 rounds, by allowing agents to express nuanced preferences. In contrast, one-vote mechanisms strug-480 gle with rigidity, leading to lower initial performance. The strictness of social choice methods, such 481 as unanimous voting, reduces collaboration performance due to its inability to accommodate dissent. 482 Score-based mechanisms, which support expressing preference gradients, foster dynamic exchanges 483 and improved decision quality, making them more effective in negotiation-heavy tasks. 484

One-vote mechanisms also have unique advantages. First, the format of voting is simpler than score-485 based mechanisms, which led to fewer format errors during experiments from LLMs. Furthermore, unanimous voting showed the highest fairness and non-decreasing outcome over rounds when utility
is quantifiable. These strengths suggest that social choice methods in decentralized MAS should be
thoughtfully chosen based on the specific task, or the primary metric being prioritized. Mixture of
different methods in different stages of collaboration may also bring effectiveness, which we leave
for future studies.

491 492

493

6.2 LINGUISTIC INDICATORS OF COLLABORATION

Linguistic analysis of agent conversations reveals several key indicators of effective collaboration. The length and complexity of messages increase progressively across rounds, suggesting that deeper engagement and richer exchanges occur as the collaboration advanced. This is particularly evident in the recommendation system environment, where tasks require greater information exchange and negotiation. Complex tasks inherently demand more intricate language, as agents need to convey not just factual information but also interpret, analyze, and synthesize inputs from others.

500 Dialogue acts provide another layer of insight into collaborative behavior. Inform, Propose and Request acts dominate conversations, highlighting that agents were frequently sharing new infor-501 mation and seeking clarifications, which are essential for progressing towards a solution. Confirm 502 act became more prevalent in later stages, particularly in environments where validation was cru-503 cial. Interestingly, *Propose* and *Evaluate* acts were more frequent in negotiation-oriented tasks, like 504 exchange economies, where agents needed to offer and assess potential solutions regularly. This 505 linguistic pattern underscores that successful collaboration hinges on the flow of information, clar-506 ification, and ongoing evaluation of proposals. These indicators align with patterns seen in human 507 collaboration, where information exchange and continuous feedback loops drive cooperative suc-508 cess.

509 510

511

6.3 EARLY STOPPING IN MULTI-AGENT COLLABORATION

Early stopping is essential in multi-agent collaboration because it prevents diminishing returns and
inefficiency in decision-making. Our experiments shows that while the initial rounds of collaboration led to significant improvements in group utility and decision quality, continuing beyond a
certain point often results in stagnant or even declining performance. This highlights the need for
mechanisms that can intelligently terminate the collaboration process once optimal solutions have
been reached, avoiding unnecessary iterations.

518 Key insights into effective early stopping methods emerge from analyzing conversational patterns. 519 Information difference suggests that as the rounds progress, the novelty of information decreases, 520 with fewer new ideas being introduced. Transition graphs of dialogue acts reveal common collaborative patterns. These signals lead to linguistic feature-based stopping methods, which halt 521 collaboration when no meaningful new information is detected, or when dialogue act patterns indi-522 cate the need for termination. Our findings demonstrate that implementing early stopping methods 523 based on linguistic cues can enhance the efficiency and effectiveness of decentralized multi-agent 524 collaboration by preventing unnecessary prolonged discussions. 525

525 526

7 CONCLUSION

527 528

This paper investigates the dynamics of cross-agent communication and decentralized decisionmaking in MAS, exploring how agent conversation and social choice methods impact collaboration and collective intelligence. We show that only the correct selection of a social choice method for a given environment can promote the success of a MAS, linguistic features of agent conversations serve as indicators of effective collaboration, and these cues provide valuable information for the effective termination of MAS.

LLM generation and prediction are unstable, with single agents often making quasi-random actions
 based on observations. Decentralized MAS improve performance by guiding agents through conversations to make more targeted proposals and enabling group decision-making to select the best
 option. Our study provides insights into how group communication and decentralized decision making enhance multi-agent collaboration, offering important implications for designing more efficient MAS environments.

540 REFERENCES

546

552

577

585

- Gati V Aher, Rosa I Arriaga, and Adam Tauman Kalai. Using large language models to simulate
 multiple humans and replicate human subject studies. In *International Conference on Machine Learning*, pp. 337–371. PMLR, 2023.
- ¹⁴⁵ Kenneth J Arrow. *Social choice and individual values*, volume 12. Yale university press, 2012.
- Antoinette Baujard, Frédéric Gavrel, Herrade Igersheim, Jean-François Laslier, and Isabelle Lebon.
 How voters use grade scales in evaluative voting. *European Journal of Political Economy*, 55: 14–28, 2018.
- Gopal Behera and Neeta Nain. Collaborative filtering with temporal features for movie recommen dation system. *Procedia Computer Science*, 218:1366–1373, 2023.
- Eva Alice Christiane Bittner and Jan Marco Leimeister. Why shared understanding matters– engineering a collaboration process for shared understanding to improve collaboration effectiveness in heterogeneous teams. In 2013 46th Hawaii international conference on system sciences, pp. 106–114. IEEE, 2013.
- 57 Duncan Black et al. The theory of committees and elections. 1958.
- Cecilia Garcia Cena, Pedro F Cardenas, Roque Saltaren Pazmino, Lisandro Puglisi, and Rafael Aracil Santonja. A cooperative multi-agent robotics system: Design and modelling. *Expert Systems with Applications*, 40(12):4737–4748, 2013.
- Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chen Qian, Chi-Min Chan,
 Yujia Qin, Yaxi Lu, Ruobing Xie, et al. Agentverse: Facilitating multi-agent collaboration and
 exploring emergent behaviors in agents. *arXiv preprint arXiv:2308.10848*, 2023.
- 565 Charles W Cobb and Paul H Douglas. A theory of production. 1928.566
- Ard-Pieter De Man and Geert Duysters. Collaboration and innovation: a review of the effects of mergers, acquisitions and alliances on innovation. *Technovation*, 25(12):1377–1387, 2005.
- Yihong Dong, Xue Jiang, Zhi Jin, and Ge Li. Self-collaboration code generation via chatgpt. *arXiv* preprint arXiv:2304.07590, 2023.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Gregory Dudek, Michael RM Jenkin, Evangelos Milios, and David Wilkes. A taxonomy for multiagent robotics. *Autonomous Robots*, 3:375–397, 1996.
- Arthur C Graesser, Stephen M Fiore, Samuel Greiff, Jessica Andrews-Todd, Peter W Foltz, and
 Friedrich W Hesse. Advancing the science of collaborative problem solving. *Psychological science in the public interest*, 19(2):59–92, 2018.
- T Guo, X Chen, Y Wang, R Chang, S Pei, NV Chawla, O Wiest, and X Zhang. Large language model based multi-agents: A survey of progress and challenges. 33rd International Joint Conference on Artificial Intelligence (IJCAI 2024), 2024. URL https://par.nsf.gov/biblio/ 10508149.
 - F Maxwell Harper and Joseph A Konstan. The movielens datasets: History and context. Acm transactions on interactive intelligent systems (tiis), 5(4):1–19, 2015.
- Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang,
 Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, et al. Metagpt: Meta programming for multiagent collaborative framework. *arXiv preprint arXiv:2308.00352*, 2023.
- Jiechuan Jiang and Zongqing Lu. Learning fairness in multi-agent systems. Advances in Neural Information Processing Systems, 32, 2019.
 - George R Klare. Assessing readability. Reading research quarterly, pp. 62–102, 1974.

- 594 Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. Semantic uncertainty: Linguistic invariances 595 for uncertainty estimation in natural language generation. In The Eleventh International Confer-596 ence on Learning Representations, 2023. URL https://openreview.net/forum?id= 597 VD-AYtPOdve. 598 Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. Camel: Communicative agents for" mind" exploration of large language model society. Advances in Neural 600 Information Processing Systems, 36:51991–52008, 2023a. 601 602 Nian Li, Chen Gao, Yong Li, and Qingmin Liao. Large language model-empowered agents for 603 simulating macroeconomic activities. arXiv preprint arXiv:2310.10436, 2023b. 604 Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and 605 Percy Liang. Lost in the middle: How language models use long contexts. Transactions of the 606 Association for Computational Linguistics, 12:157–173, 2024. 607 608 Oing Lyu, Shreya Havaldar, Adam Stein, Li Zhang, Delip Rao, Eric Wong, Marianna Apidianaki, 609 and Chris Callison-Burch. Faithful chain-of-thought reasoning. arXiv preprint arXiv:2301.13379, 610 2023. 611 Zhao Mandi, Shreeya Jain, and Shuran Song. Roco: Dialectic multi-robot collaboration with large 612 language models. In 2024 IEEE International Conference on Robotics and Automation (ICRA), 613 pp. 286-299. IEEE, 2024. 614 615 OpenAI. Gpt-4o-mini: Advancing cost-efficient intelligence, 2023. URL https://openai. 616 com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/. Ac-617 cessed: 2024-08-18. 618 Jun Ota. Multi-agent robot systems as distributed autonomous systems. Advanced engineering 619 informatics, 20(1):59-70, 2006. 620 621 Deonna M Owens, Ryan A Rossi, Sungchul Kim, Tong Yu, Franck Dernoncourt, Xiang Chen, Ruiyi 622 Zhang, Jiuxiang Gu, Hanieh Deilamsalehy, and Nedim Lipka. A multi-llm debiasing framework. 623 arXiv preprint arXiv:2409.13884, 2024. 624 Joon Sung Park, Lindsay Popowski, Carrie Cai, Meredith Ringel Morris, Percy Liang, and Michael S 625 Bernstein. Social simulacra: Creating populated prototypes for social computing systems. In 626 Proceedings of the 35th Annual ACM Symposium on User Interface Software and Technology, pp. 627 1-18, 2022.628 629 Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and 630 Michael S Bernstein. Generative agents: Interactive simulacra of human behavior. In Proceedings 631 of the 36th annual acm symposium on user interface software and technology, pp. 1–22, 2023. 632 Chen Qian, Xin Cong, Cheng Yang, Weize Chen, Yusheng Su, Juyuan Xu, Zhiyuan Liu, 633 and Maosong Sun. Communicative agents for software development. arXiv preprint 634 arXiv:2307.07924, 6, 2023. 635 636 Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-637 networks. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language 638 Processing. Association for Computational Linguistics, 11 2019. URL http://arxiv.org/ 639 abs/1908.10084. 640 Paul A Schmidt and L David Mech. Wolf pack size and food acquisition. The American Naturalist, 641 150(4):513-517, 1997. 642 643 Hal R Varian. Microeconomic analysis, 1992. 644 645 Sergey Vorotnikov, Konstantin Ermishin, Anaid Nazarova, and Arkady Yuschenko. Multi-agent robotic systems in collaborative robotics. In Interactive Collaborative Robotics: Third Inter-646
- 647 *national Conference, ICR 2018, Leipzig, Germany, September 18–22, 2018, Proceedings 3*, pp. 270–279. Springer, 2018.

- Shenzhi Wang, Chang Liu, Zilong Zheng, Siyuan Qi, Shuo Chen, Qisen Yang, Andrew Zhao, Chaofei Wang, Shiji Song, and Gao Huang. Avalon's game of thoughts: Battle against deception through recursive contemplation. *arXiv preprint arXiv:2310.01320*, 2023.
- Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang, Xiaoyun Zhang, and Chi Wang. Autogen: Enabling next-gen llm applications via multi-agent conversation framework. *arXiv preprint arXiv:2308.08155*, 2023.
- Yuzhuang Xu, Shuo Wang, Peng Li, Fuwen Luo, Xiaolong Wang, Weidong Liu, and Yang Liu.
 Exploring large language models for communication games: An empirical study on werewolf. *arXiv preprint arXiv:2309.04658*, 2023.
- Hongxin Zhang, Weihua Du, Jiaming Shan, Qinhong Zhou, Yilun Du, Joshua B Tenenbaum, Tian min Shu, and Chuang Gan. Building cooperative embodied agents modularly with large language
 models. arXiv preprint arXiv:2307.02485, 2023a.
- Jintian Zhang, Xin Xu, and Shumin Deng. Exploring collaboration mechanisms for llm agents: A
 social psychology view. *arXiv preprint arXiv:2310.02124*, 2023b.
 - Qinlin Zhao, Jindong Wang, Yixuan Zhang, Yiqiao Jin, Kaijie Zhu, Hao Chen, and Xing Xie. Competeai: Understanding the competition behaviors in large language model-based agents. *arXiv* preprint arXiv:2310.17512, 2023.
- 667 668 669

672 673

664

665

666

651

658

A DETAILS OF ROUNDTABLE

A.1 ROUNDTABLE PLATFORM DESIGN

In this section, we explain the details of the design in RoundTable by breaking it down to phases.

Input and Initialization defines each agent for collaboration. Each agent will be given a specific C_i and u_i , which include the same input user query, information of RoundTable, and the social choice method F. However, they may have different background, information access, or utility function. While C_i and u_i can be defined automatically by agent recruiting or profiling (Chen et al., 2023), we manually define them for simplicity.

679 Message Phase is the initial stage of the iterative collaboration process, where all agents participate 680 in a simultaneous, open chat conversation. During this phase, agents can freely choose recipients 681 based on their context to send messages. They can respond to others' questions, share insights, or 682 ask their own questions. There are no restrictions on whom to communicate with or what to say. 683 Once generated, all messages are updated to C_i for everyone, regardless of the original sender or 684 intended recipients. This process guarantees that all messages are sent simultaneously, makes it has 685 no specific order of communication. Since only one message can be sent per iteration, back-and-686 forth conversations can take place across multiple rounds.

687 **Proposal Phase** allows agents to present a proposal p_i , offering a potential solution or response to 688 the user query from their perspective, aiming to fulfil their own goal or maximize the utility. Before 689 generating a proposal, agents will be asked to do a step by step reasoning, which is by default 690 invisible to other agents. Reasoning is designed to give agents a chance to summarize and analyze 691 information in C_i . The most recent proposal from each agent will be considered as a candidate 692 for voting in the next phase. Like the Message Phase, proposals are made simultaneously and will be updated for C_i once everyone has generated a response. Agents also have the option to skip 693 submitting a new proposal; in such cases, their previous proposal will automatically be included 694 among the candidates. 695

Voting Phase collects one latest proposal p_i from each agent, additionally include the latest intermediate group decision x^* from previous rounds to form a candidate proposal list. Agents in this round are asked to express preferences toward candidates. The preferences can be in various formats, including voting, rating and ranking, depends on what types of social choice method Fis used. Similar with the proposal phase, agents have option to not vote (give up). With all the collected proposals and preferences, social choice method selects one proposal as the intermediate group decision. Finally, voting details and the result are updated to all agents' context C_i . Output The iteration will run for 10 rounds. After the 10th round, the latest intermediate group decision will be selected as the final output.

A.2 ROUNDTABLE PROMPT DESIGN

We present the prompt design for RoundTable in this section. Italicized words represent inputs or variables within the prompt.

Initialization Prompt

714 715	# Agent Initialization
716	You are <i>my_name</i> , an agent in a recurring collaboration environment designed to address and solve complex problems
717	solve complex problems.
718	# Task Description
719	task_description
720	
721	# Collaboration Rules
722	You start with nothing decided. The intermediate result will be decided by the social choice function at the and of each round
723	In each round, the collaboration runs in 3 phases with the following order:
724	1 Message Phase: At the beginning of each round you can send one message to a shared
725	channel for either Talking to one or more agents. All agents will send messages simultaneously.
720	You will be able to see all messages from all agents after the end of the message phase.
727	2. Proposal Phase: After the end of the message phase, you will have the opportunity to
720	propose potential solution. If you don't propose in this phase, your latest proposal will be used
729	for voting.
731	3. Voting Phase: At the end of the round, all agents' latest proposal will be voted. When
732	will be processed with the social choice function; name of social choice, where explanation
733	tion of social choice. If the social choice function selects a proposal, the intermediate result
734	will be updated accordingly. After each round, each agent will be able to see the result of the
735	vote from the previous round and the conversation history from all rounds.
736	i v
737	The collaboration will run for <i>max_rounds</i> rounds. After the last round, the latest result will be
738	the final result.
739	# Veur Dechergeun d
740	# Your Background
741	my_ugen_duckrounu
742	# Game History
743	## Latest Candidates at Round latest_candidates_round:
744	latest_candidates
745	
746	## Latest Voting Result at Round vote_history_length:
747	latest_vote_history
748	## Latast Approved Droposal
749	Proposal latest approved proposal id from Round latest approved proposal round
750	
751	Proposal <i>latest_approved_proposal_id</i> Detail:
752	latest_approved_proposal_detail
753	
754	## Conversation History until Round conversation_history_length:
755	conversation_history

7	5	6
7	5	7
7	5	8
7	5	9
7	6	0
7	6	1
7	6	2
7	6	3
7	6	4
7	6	5
7	6	6
7	6	7
7	6	8
7	6	9
7	7	0
7	7	1
7	7	2
7	7	3
7	7	4
7	7	5
' 7	' 7	6
' 7	' 7	7
7	' 7	ו 0
_	1	0
7	7	9
7	8	0
7	8	1
7	8	2
7	8	3
7	8	4
7	8	5
7	8	6
7	8	7
7	8	8
7	8	9
7	9	0

Message Phase Prompt

You at 1. Ans	re <i>my_name</i> , currently in the message phase of round <i>round_num</i> . In this phase, you can: swer questions posed by others.
2. Sha	are your findings or insights.
3. Asl You n	ay engage in multiple activities using multiple sentences.
Please "mess Don't	e type your message in the following JSON format: {"target": <list agent="" names="" of="">, age": <str, message="" your="">} generate anything except the JSON format.</str,></list>
Propo	osal Phase Prompt
You a	re <i>my_name</i> , currently in the proposal phase of round <i>round_num</i> . You have an opportu-

nity to make a proposal of the potential solution. Whether or not you submit a new proposal, your latest proposal will be considered as a candidate proposal for the voting phase.

You have two options:

1. Make a proposal:

- You can propose a potential solution by the provided format.

2. Do not make a proposal:

- If you do not want to propose a solution, you can return None as your proposal.

Please type your proposal in the following JSON format: {"reason_for_decision": <your step by step reasoning for your decision>, "proposal": proposal_format_text, or None} Don't generate anything except the JSON format.

Voting Phase Prompt (for Vote-Based Mechanisms)

788	You are <i>my_name</i> , at the voting phase at round <i>round_num</i> .
789	In this phase, all agents' latest proposal will be voted by <i>name_of_social_choice</i> , where
790	<i>explanation_of_social_choice</i> . If the social choice function selects a proposal, the intermediate
791	result will be updated accordingly.
792	
793	You have two actions to choose: vote or not vote.
794	1. For vote:
795	- You can only vote for one of the proposals from the candidate list.
796	2. For not vote:
707	- You should vote None.
700	- If you do not want to vote for any of the proposals, you can vote None.
790	- If there is no proposal, you vote None.
799	
800	The same proposal proposed by multiple agents will be merged as one proposal.
801	If there are multiple proposals passed, none of the proposals will be selected.
802	If no proposals are passed, the current intermediate result will be kept.
803	
804	The current candidate proposals are as follows:
805	proposal_list
806	What is such 2. Discus success in the following ICON formate ("manage for desiring")
807	what is your vote? Please answer in the following JSON format: { "reason_for_decision :
808	candidates you want to yote or None>}
809	Don't generate anything except the JSON format.

B DETAILS OF SOCIAL CHOICE METHODS

We present the definition and the prompt we use to define different social choice methods.

Unanimous Voting Description

The proposal that receives votes from all agents will be selected. If no proposal receives votes from all agents, no proposal will be selected.

Majority Voting Description

The proposal that receives votes from more than half of all agents will be selected. If no proposal meets this condition, none will be selected.

Plurality Voting Description

The proposal that receives the most votes will be selected.

Rated Voting Description

Each agent assigns ratings on a 5-point Likert scale to all candidate proposals, with 1 being the lowest and 5 being the highest. The proposal with the highest total score will be selected.

Ranked Voting Description

Each agent ranks all candidate proposals from the most preferred to the least preferred. Social Choice will assign 1, 1/2, 1/3... points to the 1st, 2nd, 3rd... candidates on each ballot. The proposal with the highest total points will be selected.

Cumulative Voting Description

For X candidate proposals, each agent is given X points to distribute among the proposals as they see fit. The proposal with the highest total points will be selected.

C DETAILS OF EXPERIMENTS

C.1 SIMPLE AND COMPLEX ENVIRONMENTS

Simulated environments are simplified, abstract representations of the world. They include economic experiments, such as those based on game theory, and games like chess. The primary advantages of simulated environments are: 1. They allow for the emphasis on specific agent behaviors within a controlled setting by modifying the environmental design; 2. The utility function can be mathematically expressed, and the quality of each proposal can be quantified numerically.

Complex environments are real-world applications that are more intricate and challenging, offering
a closer reflection of scenarios where agents must handle numerous variables and unpredictability.
These include tasks like collaborative problem-solving, negotiations, and real-time strategy games.
While evaluation in these settings can be subjective, they provide deeper insights into how agents
adapt, cooperate, and make decisions in dynamic, less controlled conditions, highlighting their resilience and versatility.

C.2 EXCHANGE ECONOMY

Exchange economy environment has multiple unique advantages. First, it is a plus-sum game. Economically, agents seek an equilibrium until no one can improve their utility without reducing others', implying the possibility of further collaboration exists if the allocation is not at equilibrium. Second, there are multiple possible equilibria exist in the market, allowing dynamic collaboration direction. Third, equilibria are not equal to reaching maximum total utility U_{max} , but U_{max} must be lying in 864 one of the equilibria. This can help us to measure if conflicts between agents harmful for the group to reach an ultimate goal. 866

The task description we use is as follows:

Exchange Economy Task Description

You will collaborate with other agents in a recurring exchange market game.

There are *num_of_agents* agents in this market: *list_of_agents*.

There are *num_of_goods* goods in the market: *list_of_goods*. Total quantity of each good is as follows: *total_num_of_goods*.

In this game, you will collaboratively decide how to distribute the goods among the agents. Your goal is to maximize your own utility function.

Below is the format of an agent's goal:

Exchange Economy Agent Goal

Your goal is to maximize your individual utility function by communicating, proposing, and voting with other agents. Your utility function is *util_func*

C.3 EXCHANGE ECONOMY: OTHER UTILITY SETS

We use the asymmetric scenario for the main result due to its similarity to real-world applications, here are other utility sets we examined.

Symmetric is the case where all agents prefer the same good, mirroring the scenario where agents 894 collaborate with a common goal. The following utility function is applied to each agent: $u_i =$ $a_1^{0.8} \prod_{k \neq 1} a_k^{\tilde{\theta}}, \tilde{\theta} = \frac{1 - 0.8}{|I|}$

Uniform is the case where all agents indifferently prefers all goods, using the utility function: $u_i =$ 897 $\prod a_k^{\theta}, \tilde{\theta} = \frac{1}{|I|}$ 898

899 We find these two sets of utility functions are not ideal for evaluating multi-agent collaboration. 900 Since all the agents has the same utility function, the oracle allocation among agents to reach maxi-901 mized group total utility is very close to even split. However, even split is the most frequent proposal 902 from agents while they have no information for other agents' preferences, making most of the result 903 almost perfect. This phenomenon is observed across all LLMs we've tested.

904 905 906

867 868

870 871

872

873

874

875

876 877 878

879 880

882 883

885

887 888

889 890

891

892 893

895

896

RECOMMENDATION SYSTEM C.4

907 908

909 Recommendation system aims to reveal unique collaborative behaviors that exchange economic does not have. First, there exists strong information asymmetry between agents. Any piece of information 910 is not sufficient to correctly predict the rating. Secondly, the group decision-making in these systems 911 encourages diverse approaches to problem-solving. Agents are not only consumers of information 912 but also contributors, creating a feedback loop where the quality of recommendations improves with 913 increased participation and collaboration. Lastly, the value generated in such a system is not solely 914 based on the final recommendation but also on the process of reaching that recommendation. The 915 interactions, negotiations, and information exchanges that occur along the way contribute to the 916 overall effectiveness and satisfaction of the system.

917

Below is the description of the environment sent to agents:

918	movie_id	movie_title	releas	e_date		genr	·e	
919	231	Batman Returns	19920	101	['Action', 'Adv	enture',	'Comedy',	'Crime']
920								
921		user id	age	gender	occupation	state		
922		$\frac{-\text{user} \text{int}}{7}$	 29	F	artist	NY		
923		,		-	ui tibt			

Table 3: Example of *basicInfo* dataset, which includes basic information about the target user and target movie to predict.

Recommendation System Task Description

You will collaborate with other agents in a movie recommendation game.
In this game, you will collaboratively predict the rating of a target movie (*target_movie_title*) for a target user.
There are 3 agents in this game: BasicInfo Agent, MovieHistory Agent, UserHistory Agent.
1. BasicInfo Agent has access to the basic information of the target movie and target user. It has access to the data with the following schema: *get_schema(user_info)*2. MovieHistory Agent has access to the rating history of the target movie from other people. It has access to the data with the following schema: *get_schema(movie_info)*3. UserHistory Agent has access to the rating history of the target user to other movies. It has access to the data with the following schema: *get_schema(movie_rating_history)*3. UserHistory Agent has access to the rating history of the target user to other movies. It has access to the data with the following schema: *get_schema(user_rating_history)*

You can't see other agents' information directly, but you can get information from other agents through communication. Your goal is to predict the rating a target user would give to *target_movie_title*. Utilize all available information about both the user and the movie to make the most accurate prediction possible. You only have access to partial information, but you can communicate with other agents to get more information.

In the recommendation system environment, all agents share the same goal, but has different background dataset. Here is the goal used in our setting:

Recommendation System Agent Goal

Your goal is to predict the rating the target user would give to the target movie. Utilize all available information about both the user and the movie to make the most accurate prediction possible. You only have access to *data_access*, but you can communicate with other agents to get more information.

Your Data:
 agent_dataset

960 To ease the complexity of the data structure, we pre-processed the data into three parts: basicInfo, 961 which contains the basic details of the target user and movie; userHistory, which includes the target 962 user's ratings and basic information for other movies; and movieHistory, which organizes all ratings 963 from other users for the target movie, ranked by a preference similarity score calculated using non-964 negative matrix factorization on the co-occurrence rating table between users and movies. In our 965 experiment, each part of the dataset is assigned to a different agent, forming a 3-agent collaboration 966 system. Examples of three datasets are in Table 3, 4 and 5.

C.5 DIALOGUE ACT LABELING

We present the definition of each dialogue act we use in the experiments. The definitions are also used in LLM prompt for automatic annotation.

Conversation Acts (Informational):

972	movie_id	movie_title	genre	release_date	rating	rated_date
072	1	Toy Story	['Animation', "Children's", 'Comedy']	19950101	3	19980331
973	14	Postino, Il	['Drama', 'Romance']	19940101	3	19980331
974	24	Rumble in the Bronx	['Action', 'Adventure', 'Crime']	19960223	3	19980331
014	50	Star Wars	['Action', 'Adventure', 'Romance', 'Sci-Fi', 'War']	19770101	4	19980331
975	109	Mystery Science Theater 3000: The Movie	['Comedy', 'Sci-Fi']	19960419	3	19980331
976						

Table 4: Example of *userHistory* dataset, which includes rating history from the target user to other movies. Here only shows 5 rows for spacing.

980	user_id	user_pref_similarity	personal_average_score	age	gender	occupation	state	rated_date	rating
001	343	0.98	3.99	43	М	engineer	GA	19971009	5
901	806	0.98	3.64	27	М	marketing	NY	19971217	3
982	773	0.98	3.28	20	М	student	MN	19980227	2
083	805	0.98	3.35	27	F	other	DC	19971209	3
300	447	0.97	3.57	30	М	administrator	MN	19971106	2
984									

Table 5: Example of *movieHistory* dataset, which includes rating history of the target movie from other users. Here only shows 5 rows for spacing.

- Inform Shares new information that wasn't previously known.
- Request Asks for information that the speaker doesn't have.
- Confirm Asks to verify or validate shared information.
- · Summarize Provides a brief overview of the main points.
- *Evaluate* Gives an opinion or judgment about the information.

Collaboration Acts (Decision-Making):

- *Propose* Introduce a new solution in the discussion.
- · Compromise Offers a balanced solution that incorporates parts of different parties' preferences.
- Defend Maintain support for an idea or solution after consideration or challenge.
- Accept Agrees to or accept an idea or solution.
- Decline Refuses or disagrees with an idea or solution.
- C.6 DIALOGUE ACT TRANSITION GRAPH 1007

1008 Here, we formally define dialogue act transition graph. Two nodes and the directed edge between them consist with a pair of dialogue acts and its transition probability. We denote a pair of dialogue 1010 acts as $A \to B$, where A is a dialogue act observed from an agent's message in round r-1, and B is the one from other speakers in round r. The transition probability, $p_{A \to B}$, is the probability of the 1011 existence of $A \rightarrow B$ among all observed As To be more specific, a directed edge between dialogue 1012 act A and B is calculated by the following equation: 1013

1014 1015

1016 1017

977

978

979

985

986

987 988 989

990

991 992

993

994

995 996

997 998

999

1000

1001

1002

1003 1004

1005

$$p_{A \to B} = \frac{|A \to B|}{|A|} = \frac{\sum_{s} \sum_{r=1}^{10} |\{(i, \neg i)| i \in \mathbb{I}, A \in \mathbb{D}\mathbb{A}_{i, r-1}, B \in \mathbb{D}\mathbb{A}_{\neg i, r}\}|}{\sum_{s} \sum_{r=1}^{10} |\{i| i \in \mathbb{I}, A \in \mathbb{D}\mathbb{A}_{i, r-1}\}|}$$
(1)

1018 Where s is the index of simulations, $\mathbb{D}\mathbb{A}_{i,r}$ is a set of dialogue acts from the message of the agent i 1019 at round r, and $\mathbb{D}\mathbb{A}_{r=0} = \{Start\}, \mathbb{D}\mathbb{A}_{r=11} = \{End\}.$

1020 1021

C.7 EARLY STOPPING IN MAS 1022

1023 Details of **Dialogue Act** early stopping method is as follows. We first collect all pairs from dialogue acts as independent variable (with one-hot encoding), where one form round r-1 and another from 1024 r. Then make the group performance at round r to match with dialogue act pairs as dependent 1025 variable. Next, we run OLS regression on this dataset, assigning coefficients and p-value to each

Table 6: Ablation studies on exchange economics environment. The reported numbers are an average of all simulations, and the numbers in parentheses are standard errors. The smallest and largest value in a category is colored in blue and red.

		GROUP TOTAL UTILITY@3	GROUP TOTAL UTILITY@5	GROUP TOTAL UTILITY@10	AUC@3	AUC@5	AUC@10	RATIONALITY	MIN/MAX	RIGIDITY
gpt-3.5 gpt-40- gpt-40 llama3 llama3	i-turbo -mini -1-8b -1-70b	62.44 (2.24) 80.61 (1.46) 76.39 (4.26) 69.55 (5.87) 93.10 (3.28)	63.40 (0.35) 81.43 (1.32) 84.29 (1.94) 72.91 (5.45) 97.15 (0.56)	Social (64.07 (0.61) 78.18 (1.64) 91.02 (1.15) 73.80 (5.39) 96.69 (0.79)	Choice: Major 56.15 (2.81) 66.64 (3.07) 57.51 (3.82) 59.60 (6.22) 76.41 (4.61)	ity, # Agent: K 59.15 (1.69) 72.59 (2.10) 68.03 (2.34) 64.91 (5.55) 84.64 (2.78)	61.56 (0.88) 76.18 (1.49) 78.97 (1.16) 69.07 (5.08) 90.69 (1.43)	17.88 (2.11) 25.89 (2.25) 38.16 (4.19) 3.78 (0.71) 35.11 (4.03)	97.58 (0.91) 58.41 (6.14) 74.27 (3.44) 77.66 (5.20) 87.66 (3.28)	47.00 (3.08) 70.00 (2.84) 73.00 (1.93) 86.00 (1.83) 70.67 (2.67)
K = 3 $K = 4$ $K = 5$	1	79.88 (0.80) 58.29 (2.39) 63.67 (2.72)	81.33 (0.72) 63.10 (2.12) 72.42 (2.33)	Social Cl 79.61 (0.88) 64.99 (2.12) 74.56 (2.18)	hoice: Majority 64.08 (1.81) 43.64 (2.04) 40.23 (1.90)	7, Model: gpt-4 70.91 (1.23) 51.05 (1.87) 52.17 (1.88)	to-mini 75.67 (0.84) 57.66 (1.89) 63.13 (1.94)	23.80 (0.94) 34.25 (1.39) 36.26 (1.36)	64.00 (3.07) 64.54 (3.14) 57.54 (3.28)	71.30 (1.32) 80.30 (0.96) 73.50 (1.36)
pair. and r	The s negati	strength of vely, and	coefficien p-value sh	t shows ho ows the sta	w each di ttistical si	alogue ac gnificanc	et pairs rel e of these	ate the perfo relationship	ormance, j os.	positively
To fi hype	nd th rpara	e best ear meters:	rly stoppir	ng rounds,	we do a	greedy s	earch on	training set	for the	following
	• <i>te</i>	op_da: the	number o ive.	f dialogue	act pairs	with the t	top coeffi	cients to app	oly in can	didates. 1
	• <i>p</i> a	- <i>value</i> : th nd None.	e threshol	d of p-valu	ie for a pa	air to be	considere	d as candida	ates. 0.05	5, 0.1, 0.2
	• <i>c</i> r	<i>ount_per_</i> ound to pa	<i>round</i> : the ass. When	threshold passed, it g	of how n gives +1 s	nany occu core for t	irrences of ermination	of a candidat on. 1 to 3, in	te should clusive.	exist in a
	• s to	<i>core</i> : the op_da, incl	threshold lusive.	of how ma	any score	s a round	l needs to	be termina	ited. 1 to	value of
We e each of ea	xami socia rly st	ne all com ll choice a opping on	binations and enviror the test se	of hyperpar nment. Usi et.	rameters on the be	on the trai st hyperp	ining data arameter	ı, and identif s, we evalua	ies the be te the per	est sets for formance
C.8	Мо	DELS								
For a agent and <i>I</i> we us <i>Llam</i>	ll exp ts. In <i>Llama</i> sed <i>pa</i> <i>va-3.1</i>	eriments, ablation a-3.1-8b-In araphrase ab-Instru	we by defa study, we <i>nstruct</i> for <i>-MiniLM-</i> <i>uct</i> . For all	ault use gpt use gpt-3 compariso L6-v2 (Rein LLM infer	<i>t-40-mini- 5-turbo-0</i> on (Dubey mers & G rences, w	2024-07- 125, gpt- 7 et al., 20 urevych, e used <i>ter</i>	-18 (Oper 40-2024- 024). To 2019). Fo mperature	AI, 2023) as $05-13$, Llam calculate sen or dialogue a $e = 0$.	s the main a-3.1-70 ntence em act labelin	LLM for <i>b-Instruct</i> beddings ig, we use
D	DET	TAILS OF	F RESUL	TS						
D.1	AB	LATION S	TUDIES							
The agent agent 40-20	result t coll t syst 024-0	t of ablation aboration em in maj 05-13, Llan	on studies with diffe ority votin ma-3.1-8b	is shown i rent LLMs g setting w <i>Instruct</i> ar	n Table 6 in Figur 7777 ith 30 sin 1878 ith 30 sin	We fur e 2(b). F nulations 3.1-70b-	ther comp or simpli . The resu <i>Instruct</i> in	pare the perf city, we onl alts of <i>gpt-3</i> . n Table 1 sho	formance y compar <i>5-turbo-(</i> ow that m	of multi- re $K = 3$ 0.00000000000000000000000000000000000

collaboration works with different LLMs. The stronger (larger) the model, the better performance observed in collaboration, both in end round performance and also efficiency in reaching to an agreement. Notably, the Rationality of the weaker model, *gpt-3.5-turbo-0125* and *Llama-3.1-8b-Instruct*, is only 17.88% and 3.78%, indicates that single agent will largely fail due to the complexity of the task. Nonetheless, the final round performance has reached 64% and 80%, shows the effectiveness of multi-agent collaboration with weak models.

1079 Scaling on MAS is evaluated with different number of agents. From the nature of decentralized collaboration, the more the agents, the harder to reach an agreement. Furthermore, the additional

Table 7: Recommendation system environment results comparison between social choices. The smallest and largest value in a category is colored in blue and red. Experiments are made with gpt-40-mini.

1084		MAE@1	MAE@2	MAE@3	MAE@4	MAE@5	MAE@6	MAE@7	MAE@8	MAE@9	MAE@10
1085	Unanimous Majority	0.85 0.81	0.84 0.81	0.84 0.81	0.86 0.81	0.86 0.84	0.87 0.82	0.88 0.84	0.89 0.84	0.89 0.85	0.88 0.86
1086	Plurality Rated	0.81 0.83	0.82 0.80	0.82 0.77	0.83 0.77	0.82 0.77	0.83 0.75	0.83 0.76	0.84 0.77	0.83 0.76	0.84 0.79
1087	Ranked	0.90	0.83	0.86	0.86	0.84	0.87	0.86	0.86	0.84	0.84
1088	Cumulative	0.83	0.83	0.83	0.78	0.78	0.77	0.78	0.79	0.77	0.76
1089		RMSE@1	RMSE@2	RMSE@3	RMSE@4	RMSE@5	RMSE@6	RMSE@7	RMSE@8	RMSE@9	RMSE@10
	Unanimous	1.03	1.04	1.06	1.06	1.07	1.08	1.09	1.09	1.09	1.09
1090	Majority	1.08	1.03	1.02	1.02	1.03	1.02	1.04	1.05	1.05	1.06
1001	Plurality	1.11	1.07	1.06	1.09	1.08	1.12	1.12	1.14	1.13	1.15
1091	Rated	1.09	1.03	1.00	0.98	0.98	0.97	0.98	0.99	0.99	1.01
1092	Ranked Cumulative	1.20 1.10	1.13 1.09	1.13 1.09	1.13 1.04	1.10 1.03	1.13 1.02	1.11 1.02	1.11 1.03	1.07 1.01	1.07 1.01

Table 8: Heatmaps of average message length, message complexity, and information difference in each round. Analysis is made in 3 agents, gpt-4o-mini setting. Color gradient is calculated with green as maximum, red as minimum, and white as median value in each table.

1099	EXCHANGE ECONOMY											RECOMMENDATION SYSTEM										
1000		1	2	3	4	5	6	7	8	9	10		1	2	3	4	5	6	7	8	9	10
1100									AVI	ERAGE N	IESSAG	E LE	NGTH									
1100	Unanimous	46	62	83	85	89	94	96	96	97	94		68	85	94	96	100	101	103	104	105	104
1101	Majority	43	62	82	84	87	91	91	93	91	91		67	85	93	99	101	103	105	105	106	105
1101	Plurality	42	59	81	82	84	90	89	91	93	90		67	84	94	99	101	104	105	106	107	106
1100	Rated	42	60	82	83	89	91	92	94	94	93		67	86	96	100	101	103	104	106	107	105
1102	Ranked	43	61	80	83	86	90	91	90	91	90		67	84	98	99	102	102	103	104	105	104
1100	Cumulative	42	61	81	83	86	91	91	92	93	91		68	85	95	98	100	102	103	104	104	103
1103																						
									AVER/	AGE ME	SSAGE C	COM	PLEXIT	Y								
1104	Unanimous	7.1	10.2	9.3	9.0	9.0	9.1	9.2	9.5	9.5	9.6		10.4	10.7	11.5	11.7	11.9	12.0	12.0	12.1	12.3	12.2
	Majority	7.6	10.4	9.5	9.4	9.2	9.2	9.4	9.5	9.6	9.7		10.3	10.7	11.5	11.7	11.9	12.0	12.1	12.1	12.1	12.1
1105	Plurality	6.8	10.5	9.4	9.1	8.9	9.1	9.2	9.3	9.4	9.4		10.3	10.6	11.5	11.7	11.9	12.0	12.0	12.1	12.1	12.1
	Rated	5.8	10.4	9.4	9.1	9.1	9.2	9.4	9.3	9.5	9.6		10.3	10.7	11.6	11.8	12.0	12.0	12.2	12.2	12.3	12.3
1106	Ranked	6.7	10.2	9.5	9.2	9.1	9.1	9.4	9.4	9.5	9.6		10.3	10.7	11.6	11.7	11.9	12.1	12.1	12.2	12.2	12.2
1100	Cumulative	6.9	10.6	9.4	9.3	9.1	9.2	9.4	9.5	9.6	9.7		10.3	10.6	11.5	11.7	11.7	11.9	11.9	12.0	12.1	12.1
1107																						
1107								4	WERAG	E INFOI	RMATIO	N DI	FFEREN	ICE								
1100	Unanimous		0.23	0.16	0.13	0.13	0.12	0.11	0.11	0.10	0.10			0.20	0.13	0.11	0.09	0.09	0.08	0.08	0.08	0.09
1100	Majority		0.26	0.17	0.13	0.12	0.12	0.11	0.10	0.10	0.11			0.19	0.13	0.11	0.10	0.10	0.09	0.09	0.08	0.09
1100	Plurality		0.28	0.16	0.13	0.14	0.12	0.11	0.11	0.10	0.11			0.21	0.14	0.11	0.10	0.09	0.09	0.09	0.08	0.09
1109	Rated		0.33	0.16	0.13	0.12	0.10	0.11	0.10	0.10	0.11			0.20	0.13	0.11	0.10	0.10	0.09	0.09	0.08	0.09
	Ranked		0.28	0.17	0.13	0.12	0.11	0.10	0.10	0.10	0.10			0.20	0.13	0.11	0.09	0.09	0.09	0.09	0.08	0.09
1110	Cumulative		0.30	0.18	0.13	0.12	0.11	0.11	0.11	0.10	0.10			0.20	0.13	0.10	0.09	0.09	0.09	0.09	0.08	0.09

amount of historical context brought by increasing number of agents can also be a burden to the agent's performance. Comparing K = 3, 4, 5 settings in Table 1, the quality and efficiency of collaboration generally drops with the increasing number of agents. Exceptionally, performance in K = 4 setting is worse than that of K = 5. This is because majority voting is harder to achieve an agreement with, even number of participants. With 4 agents, a proposal needs at least 3 votes to get accepted; 5 agent setting also requires 3 votes, but they can ignore up to 2 agents' preferences or agreements.

D.2 RECOMMENDATION SYSTEM

We show the detailed result between social choices in recommendation system environment in Table 7.

CROSS-AGENT CONVERSATION ANALYSIS D.3

D.3.1 LINGUISTIC STATISTICS

In this section, we show how linguistic features presented in different environments and social choices. Table 8 is heatmaps for average message length, message complexity, and information difference observed in agent conversation.

	1	2	3	EX0 4	CHANGE 5	ECONC 6	ОМҮ 7	8	9 UN4		1	2	3	RECO!	MMENDA 5	ATION SY 6	7 7	8	9
Inform	1.00	0.99	0.88	0.59	0.69	0.68	0.62	0.66	0.62	0.54	1.00	0.98	0.99	1.00	1.00	1.00	0.99	1.00	
Request	1.00	1.00	0.96	0.94	0.89	0.91	0.89	0.89	0.85	0.82	1.00	1.00	1.00	1.00	1.00	1.00	0.99	1.00	
Confirm	0.08	0.34	0.28	0.17	0.27	0.28	0.33	0.37	0.41	0.55	0.29	0.49	0.63	0.00	0.79	0.68	0.74	0.69	
Evaluate	0.98	0.91	0.15	0.80	0.90	0.94	0.91	0.93	0.91	0.92	0.07	0.58	0.58	0.62	0.60	0.59	0.64	0.62	
Propose	0.06	0.47	0.99	0.99	0.99	1.00	0.98	0.98	0.92	0.92	0.10	0.07	0.20	0.20	0.15	0.21	0.30	0.28	
ompromise	0.00	0.09	0.43	0.13	0.24	0.39	0.42	0.37	0.39	0.40	0.01	0.00	0.00	0.00	0.01	0.00	0.00	0.00	
Defend	0.00	0.00	0.04	0.00	0.02	0.03	0.02	0.01	0.03	0.08	0.00	0.00	0.00	0.01	0.01	0.02	0.02	0.01	
Accept	0.00	0.00	0.01	0.00	0.00	0.02	0.06	0.08	0.12	0.27	0.00	0.00	0.00	0.01	0.00	0.02	0.01	0.01	
Decline	0.00	0.00	0.01	0.00	0.00	0.01	0.01	0.00	0.02	0.03	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.01	
Otners	0.00	0.00	0.02	0.06	0.04	0.01	0.06	0.04	0.04	0.04	0.03	0.04	0.00	0.03	0.01	0.02	0.04	0.01	
T-6	1.00	1.00	0.06	0.72	0.72	0.66	0.65	0.77	MA	JORITY	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	
Request	1.00	1.00	1.00	0.72	0.72	0.86	0.85	0.77	0.88	0.85	1.00	1.00	1.00	1.00	1.00	1.00	1.00	0.99	
Confirm	0.16	0.25	0.34	0.31	0.44	0.36	0.41	0.41	0.53	0.52	0.25	0.48	0.56	0.65	0.71	0.71	0.74	0.71	
ummarize	0.15	0.02	0.15	0.13	0.19	0.21	0.20	0.18	0.33	0.30	0.01	0.15	0.07	0.18	0.11	0.17	0.19	0.16	
valuate	0.93	0.91	0.67	0.78	0.86	0.88	0.88	0.86	0.89	0.87	0.05	0.61	0.52	0.49	0.58	0.52	0.61	0.54	
ropose	0.37	0.41	0.97	0.99	0.95	0.97	0.99	0.96	0.93	0.90	0.05	0.08	0.29	0.17	0.22	0.15	0.25	0.21	
.ompromise	0.00	0.03	0.38	0.14	0.25	0.41	0.39	0.31	0.36	0.30	0.00	0.00	0.02	0.00	0.00	0.00	0.00	0.00	
Accept	0.00	0.00	0.04	0.03	0.04	0.11	0.08	0.03	0.20	0.20	0.00	0.02	0.02	0.00	0.02	0.00	0.02	0.00	
Decline	0.00	0.00	0.02	0.02	0.03	0.02	0.04	0.03	0.04	0.00	0.00	0.01	0.02	0.00	0.00	0.00	0.01	0.00	
Others	0.00	0.01	0.01	0.08	0.08	0.08	0.08	0.08	0.07	0.06	0.02	0.00	0.03	0.03	0.01	0.02	0.02	0.01	
									PLU	RALITY									
nform	1.00	0.98	0.93	0.74	0.69	0.68	0.61	0.64	0.66	0.64	1.00	1.00	0.98	0.99	0.99	0.99	1.00	0.98	
Request	1.00	1.00	0.95	0.97	1.00	0.93	0.91	0.90	0.96	0.85	1.00	1.00	1.00	1.00	1.00	1.00	0.99	1.00	
Confirm	0.01	0.29	0.16	0.26	0.40	0.42	0.43	0.46	0.52	0.55	0.29	0.44	0.61	0.70	0.73	0.78	0.76	0.78	
Ryaluate	0.39	0.01	0.62	0.17	0.17	0.22	0.20	0.22	0.18	0.25	0.05	0.15	0.11	0.23	0.18	0.23	0.19	0.17	
ronose	0.19	0.49	0.99	1.00	0.98	0.97	0.97	0.95	0.93	0.87	0.00	0.05	0.31	0.31	0.30	0.32	0.32	0.34	
Compromise	0.00	0.02	0.53	0.14	0.14	0.38	0.37	0.33	0.32	0.32	0.00	0.01	0.00	0.01	0.02	0.00	0.00	0.00	
Defend	0.00	0.00	0.00	0.02	0.00	0.05	0.06	0.06	0.07	0.07	0.00	0.00	0.01	0.03	0.02	0.04	0.04	0.02	
ccept	0.00	0.00	0.01	0.02	0.00	0.08	0.15	0.10	0.13	0.25	0.00	0.00	0.01	0.03	0.01	0.02	0.04	0.02	
Decline	0.00	0.00	0.01	0.02	0.00	0.05	0.05	0.01	0.02	0.02	0.00	0.00	0.01	0.01	0.01	0.01	0.01	0.01	
iners	0.00	0.00	0.02	0.10	0.06	0.08	0.05	0.09	0.07	0.04	0.01	0.01	0.02	0.01	0.06	0.02	0.07	0.00	
c .	1.00	0.04	0.07	0.70	0.70	0.56	0.50	0.65	F	RATED	1.00	1.00	0.00	0.00	0.00	1.00	1.00	1.00	
atorm	1.00	0.94	0.97	0.78	0.69	0.56	0.59	0.65	0.65	0.61	1.00	1.00	0.98	0.99	0.99	1.00	1.00	1.00	
Confirm	0.00	0.47	0.33	0.31	0.35	0.41	0.33	0.40	0.43	0.48	0.27	0.48	0.67	0.70	0.76	0.65	0.67	0.71	
ummarize	0.16	0.02	0.17	0.13	0.17	0.19	0.23	0.23	0.20	0.25	0.09	0.13	0.14	0.16	0.17	0.22	0.27	0.21	
Evaluate	0.40	1.00	0.57	0.69	0.87	0.85	0.85	0.84	0.92	0.85	0.09	0.60	0.62	0.66	0.58	0.66	0.58	0.51	
Propose	0.11	0.66	0.97	0.99	0.98	0.99	1.00	0.97	0.95	0.94	0.16	0.05	0.26	0.20	0.21	0.23	0.27	0.31	
Compromise	0.00	0.03	0.48	0.12	0.24	0.40	0.32	0.34	0.38	0.32	0.00	0.00	0.01	0.01	0.01	0.00	0.00	0.00	
Accept	0.00	0.00	0.03	0.01	0.02	0.05	0.07	0.05	0.07	0.04	0.00	0.00	0.00	0.01	0.01	0.00	0.00	0.01	
Decline	0.00	0.00	0.02	0.01	0.01	0.02	0.01	0.03	0.03	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.01	0.00	
Others	0.02	0.00	0.02	0.10	0.08	0.08	0.06	0.07	0.05	0.08	0.03	0.02	0.06	0.04	0.03	0.02	0.02	0.05	
									R	ANKED									
Inform	1.00	0.96	0.96	0.78	0.79	0.66	0.69	0.71	0.72	0.75	1.00	1.00	1.00	1.00	0.99	0.99	1.00	1.00	7
Request	1.00	1.00	0.99	0.98	0.96	0.95	0.95	0.96	0.91	0.89	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	
Summarize	0.05	0.03	0.49	0.55	0.40	0.55	0.45	0.52	0.48	0.34	0.21	0.45	0.55	0.62	0.75	0.73	0.78	0.08	1
Evaluate	0.92	0.95	0.61	0.65	0.79	0.20	0.20	0.86	0.83	0.89	0.06	0.52	0.64	0.66	0.63	0.58	0.59	0.60	
Propose	0.15	0.48	0.92	0.96	0.99	0.98	0.97	0.93	0.90	0.87	0.09	0.04	0.27	0.16	0.29	0.19	0.24	0.28	
Compromise	0.00	0.01	0.35	0.14	0.21	0.32	0.35	0.36	0.37	0.21	0.00	0.01	0.01	0.02	0.00	0.00	0.01	0.03	
Defend	0.00	0.00	0.05	0.02	0.03	0.06	0.03	0.02	0.05	0.09	0.00	0.00	0.00	0.00	0.00	0.01	0.01	0.01	
Accept Decline	0.00	0.00	0.05	0.03	0.04	0.05	0.12 0.01	0.13	0.11	0.27	0.00	0.00	0.02	0.01	0.02	0.01	0.01	0.02	
Others	0.01	0.02	0.02	0.05	0.06	0.13	0.11	0.04	0.03	0.01	0.01	0.02	0.08	0.03	0.01	0.03	0.01	0.03	
									CUM	ULATIVE									
Inform	1.00	0.94	0.99	0.82	0.72	0.68	0.60	0.67	0.72	0.74	1.00	1.00	1.00	0.98	0.97	0.98	1.00	1.00	
Request	1.00	1.00	0.98	0.97	0.96	0.95	0.94	0.92	0.95	0.86	1.00	1.00	1.00	1.00	1.00	0.99	1.00	1.00	
Summarize	0.05	0.02	0.14	0.16	0.09	0.35	0.25	0.23	0.21	0.32	0.10	0.12	0.00	0.72	0.22	0.22	0.17	0.29	
Evaluate	0.81	0.98	0.58	0.74	0.86	0.90	0.92	0.90	0.90	0.86	0.04	0.65	0.63	0.51	0.61	0.62	0.59	0.56	
Propose	0.17	0.58	0.98	0.97	0.99	0.98	0.95	0.89	0.89	0.87	0.08	0.03	0.24	0.26	0.22	0.22	0.27	0.28	
Compromise	0.00	0.05	0.54	0.14	0.19	0.42	0.30	0.35	0.28	0.27	0.01	0.00	0.00	0.02	0.00	0.00	0.02	0.01	
Defend	0.00	0.00	0.01	0.00	0.03	0.05	0.05	0.02	0.05	0.06	0.01	0.01	0.01	0.01	0.01	0.00	0.02	0.03	
Accept	0.00	0.00	0.01	0.00	0.01	0.09	0.10	0.10	0.21	0.21	0.01	0.01	0.00	0.00	0.01	0.03	0.01	0.03	
Decline	0.00	0.00	0.00	0.00	0.00	0.04	0.01	0.00	0.00	0.03	0.01	0.01	0.00	0.00	0.00	0.07	0.01	0.02	
Others		0.00	0.00	0.00	0.00	0.07	0.00	0.00	0.00	0.00	0.01	0.00	0.07	0.05	0.01	0.07	0.01	0.00	
thers																			

Table 9: Heatmaps of ratio of dialogue acts in each round. Color gradient is calculated with green 1135

- 1176
- 1177 D.5 EARLY STOPPING 1178

We share full result of early stopping experiment in Table 10, and basic statistics of early stopping 1179 methods, including early stopped round, effective ratio, and information difference threshold in Ta-1180 ble 11. Effective ratio stands for the chance of a certain rule can be applied in the test set. Threshold 1181 is the embedding distance threshold used for early stopping. 1182

Although linguistic feature based methods, Information Difference and Dialogue Act, outperforms 1183 other early stopping methods and baseline in exchange economy environment, the difference is 1184 small. This is due to the small room for improvement, identified between the baseline (@10) and 1185 Oracle performance. 1186

1187



Figure 5: Dialogue act transition graph for different social choices in exchange economics environ-ment.



Figure 6: Dialogue act transition graph for different social choices in recommendation system environment.

Table 10: Comparison between early stopping methods in different social choices. The performance is shown with group total utility and MAE. (\uparrow) and (\downarrow) indicates better performance with higher and lower values, respectively. Experiments are based on the results of 3 agents, gpt-4o-mini setting. Results are based on 5-fold cross validation. We observe that language-based methods performed well overall.

1229	EARLY STOPPING METHOD	UNANIMOUS	MAJORITY	PLURALIRY	RATED	RANKED	CUMULATIVE					
1000	EXCHANGE ECONOMY, in Group Total Utility(↑)											
1230	@10 (Baseline)	0.48	0.80	0.77	0.80	0.78	0.78					
1231	First Agreement	0.48	0.78	0.76	0.74	0.77	0.74					
1000	Consecutive Agreements	0.48	0.80	0.77	0.80	0.78	0.79					
1232	Validation Checkpoint	0.37	0.81	0.80	0.81	0.81	0.81					
1233	Information Difference	0.39	0.81	0.80	0.81	0.81	0.82					
100/	Dialogue Act	0.42	0.81	0.80	0.81	0.81	0.81					
1234	Oracle	0.48	0.84	0.82	0.83	0.84	0.84					
1235												
1006	RECOMMENDATION SYSTEM, in MAE (\downarrow)											
1230	@10 (Baseline)	0.88	0.86	0.84	0.79	0.84	0.76					
1237	First Agreement	0.82	0.80	0.81	0.82	0.89	0.82					
1000	Consecutive Agreements	0.84	0.81	0.84	0.77	0.82	0.77					
1230	Validation Checkpoint	0.81	0.81	0.82	0.80	0.83	0.82					
1239	Information Difference	0.86	0.78	0.85	0.79	0.78	0.81					
1040	Dialogue Act	0.84	0.82	0.79	0.76	0.85	0.82					
1240	Oracle	0.73	0.63	0.59	0.62	0.69	0.67					
1241												

1242									
1243									
1244									
1245									
1246									
1247									
1248									
1249									
1250									
1251									
1252									
1253									
125/									
1255									
1200									
1230									
1257									
1258									
1259									
1260									
1267		Table 11: Basic	statistics o	f differen	t early sto	pping	methods	s.	
1262	STATISTICS	FARLY STOPPING METHOD	UNANIMOUS	MAIORITY	PLURALITY	RATED	RANKED	CUMULATIVE	AVERAGE
1203	SIAIISTICS	Oracle	EXCH	ANGE ECONO	MY 3 84	3.81	3.49	4 11	3 73
1265		First Agreement	3.00	1.57	1.54	1.15	1.07	1.33	1.61
1205	Early Stopped Round	Validation Checkpoint	2.60	2.80	3.00	3.00	2.60	3.00	2.83
1200		Dialogue Act	6.23	3.32 7.09	3.28 8.33	3.09 7.64	5.47	7.40	7.03
1207		Ensemble First Agreement	0.62	4.15	4.73	7.64	5.01	3.66	5.86 0.94
1200	Effective Ratio	Consecutive Agreement Information Difference	0.00 1.00	0.03	0.03 1.00	0.02 1.00	0.00 1.00	0.01 1.00	0.02 1.00
1269	Threshold	Dialogue Act	0.62	0.34	0.07	0.19	0.63	0.31	0.36
1270	Threshold	Information Difference	DECOMM	U.I.7	UCTEM	0.20	0.18	0.17	0.18
1271		Oracle	2.40	2.10	2.38	2.15	2.25	2.44	2.29
1272		First Agreement Consecutive Agreement	1.53 3.34	3.31	1.00 3.15	1.04 3.61	1.00 3.46	3.87	3.46
1273	Early Stopped Round	Validation Checkpoint Information Difference	1.20 3.12	1.00 3.67	1.20 4.30	1.00 3.50	1.00 3.97	1.00 3.34	1.07 3.65
1274		Dialogue Act Ensemble	7.58 3.60	6.98 5.87	6.02 5.12	5.22 5.38	5.37 3.97	6.53 10.00	6.28 5.66
1275		First Agreement Consecutive Agreement	1.00 0.82	1.00	1.00 0.91	1.00	1.00 0.62	1.00 0.78	1.00
1276	Effective Ratio	Information Difference	1.00	1.00	1.00	1.00	0.99	1.00	1.00
1277	Threshold	Information Difference	0.15	0.13	0.12	0.14	0.12	0.14	0.13
1278									
1279									
1280									
1281									
1282									
1283									
1284									
1285									
1286									
1287									
1288									
1289									
1290									
1201									
1202									
1202									
1207									
1205									
1290									