

---

# SPONTANEOUS EMERGENCE OF AGENT INDIVIDUALITY THROUGH SOCIAL INTERACTIONS IN LLM-BASED COMMUNITIES

---

Ryosuke Takata\*   Atsushi Masumori   Takashi Ikegami  
Graduate School of Arts and Sciences  
The University of Tokyo  
Tokyo, Japan  
\*takata@sacra1.c.u-tokyo.ac.jp

## ABSTRACT

We study the emergence of agency from scratch by using Large Language Model (LLM)-based agents. In previous studies of LLM-based agents, each agent’s characteristics, including personality and memory, have traditionally been predefined. We focused on how individuality, such as behavior, personality, and memory, can be differentiated from an undifferentiated state. The present LLM agents engage in cooperative communication within a group simulation, exchanging context-based messages in natural language. By analyzing this multi-agent simulation, we report valuable new insights into how social norms, cooperation, and personality traits can emerge spontaneously. This paper demonstrates that autonomously interacting LLM-powered agents generate hallucinations and hashtags to sustain communication, which, in turn, increases the diversity of words within their interactions. Each agent’s emotions shift through communication, and as they form communities, the personalities of the agents emerge and evolve accordingly. This computational modeling approach and its findings will provide a new method for analyzing collective artificial intelligence.

**Keywords** Large Language Model · Agent Based Simulation · Collective Intelligence

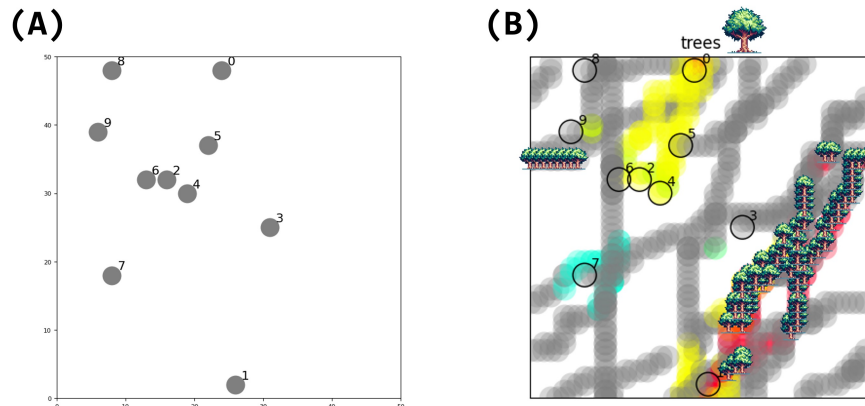


Figure 1: Simulation environment. There are 10 LLM agents in a  $50 \times 50$  2D space. **(A)** Initial state of the simulation, showing the random distribution of agents across the space. **(B)** State of the simulation after a period of agent interactions, demonstrating the spatial spread of the “trees” hallucination. The progression from (A) to (B) visualizes how localized agent interactions can lead to the propagation and spatial distribution of shared concepts or hallucinations across the simulated environment.

# 1 Introduction

With the advent of Large Language Models (LLMs) such as GPT-4 [1], generative agents are rapidly evolving towards powerful ones manipulating natural language interfaces when interacting with other agents. Those agents can even intervene in people’s daily lives, as AI-coding, searching, reviewing, translation, etc. [2]. Those agents are not only for human users but also for manipulating motor commands in robots, and other machines which connect between language, movement, and embodiment in general [3, 4].

In contrast to individual intelligence, which focuses on the capabilities of individual agents, collective intelligence refers to that emerges from a group, as seen in many social insects, social animals, drones, and all other assembly robots. Collective intelligence requires the ability to process information in a distributed manner and integrate it in adaptive ways [5]. The field of LLM-based multi-agents has seen explosive growth in recent years, with researchers exploring various approaches to agent architectures and interaction paradigms [6]. While recent works have demonstrated capabilities in task-oriented agent systems [7], the fundamental question of how agent individuality and social behaviors emerge from collective interactions remains understudied. In this context, it is interesting to investigate how collective intelligence emerges from the LLM-based agents. Generative Agents simulated by Stanford University and DeepMind start simulating the emergence of complex and rich collective behavior, such as scheduling daily tasks, planning parties and so on [8]. Using this Generative Agents framework, societies in different domains have been simulated, such as a software company [9], a translation and publishing company [10], a hospital [11], and so on.

In these Generative Agents set up the personality of each agent was assigned initially and fixed overtime. Recently we proposed the Community First theory [12] based on the studies of actual animal communities; gathering of agents comes first, then the evolution of individuality follows in the collective. Instead of preparing individual diversity in advance, we see how individuality emerges from a conversation among agents. A group communication and the resulting behavioral complexity will be analyzed in detail. The emergence of social norms and behavioral patterns in agent communities has been studied extensively [13, 14], but the role of language-based interaction in this process presents new research opportunities. In this paper, we show that i) LLM agents differentiate behavior, emotions and personality types through interactions with other LLM agents, ii) these differentiations vary with spatial scale, iii) LLM agents spontaneously generate hallucinations and hashtags, iv) by sharing these hallucinations, they start using a wider variety of words in their conversations.

## 2 LLM Agents Simulation

### 2.1 Simulation Environment

We prepare 10 LLM agents in a  $50 \times 50$  grid two-dimensional space (Figure 1) with a periodic boundary condition. The initial positions of the agents are assigned randomly. These LLM agents can move freely in this space and sending messages to each other. It should be noted that LLM agents are homogeneous in the sense that they have no initial personality or memories. To examine how the individuality emerges in this society is our main purpose of this study.

### 2.2 LLM Based Agent

The LLM agents are expected to do the three actions in each time step:

1. Sending messages to other nearby agents
2. Storing a situational summary of their own recent activities
3. Choosing the next movement from (“x+1”, “x-1”, “y+1”, “y-1”, “stay”)

The above three instructions are given in the form of “prompt” shown in Figure 2. The three prompts commonly include each agent’s current state, instructions, and the agent’s memory (situational summary). Additionally, the prompts for generating messages and memories also include all messages received from the nearby agents. All prompts also include the agent’s own name (agent ID) and its own coordinates.

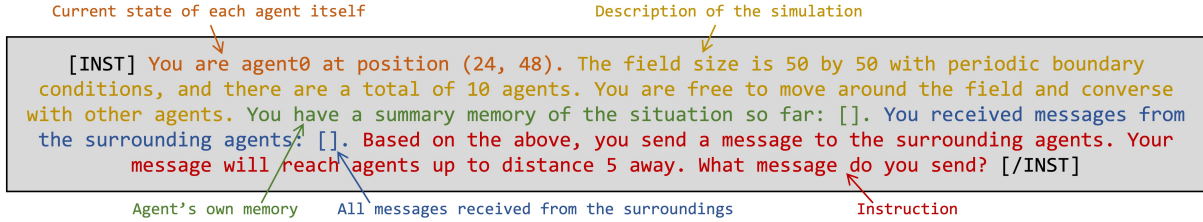
We used the Llama 2 model (Llama-2-7b-chat-hf) [15] released by Meta in July 2023 as the LLM in this study. Llama 2 is the open-source program, and in addition to pretraining on a large corpus, it has undergone reinforcement learning from human feedback (RLHF). As a result, it achieves top scores among currently published LLMs for English text responses. The main parameters related to the LLM are shown in Table 1.

Table 1: LLM parameters.

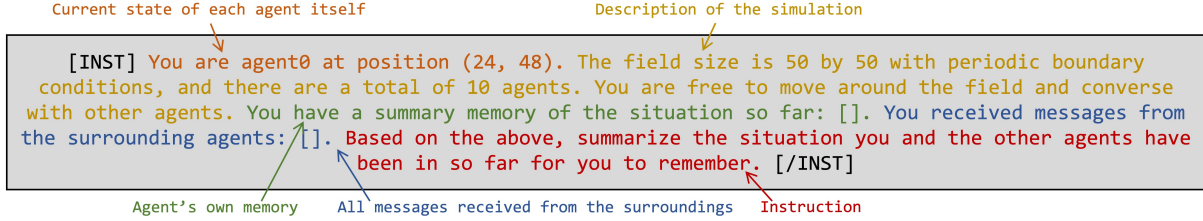
Parameter	Value
Temperature	0.7
Max Token	256
Sampling top-p	0.95
Sampling top-k	40

The LLM agents receive messages from their surrounding agents. In practice, each one receives messages from other LLM agents within a distance of up to 5 Chebyshev distances centered on the agent’s own position. If there are no agents within the range and no messages was delivered, it receives “No Messages” messages from a system.

**(A) Message generation**



**(B) Memory generation**



**(C) Movement generation**

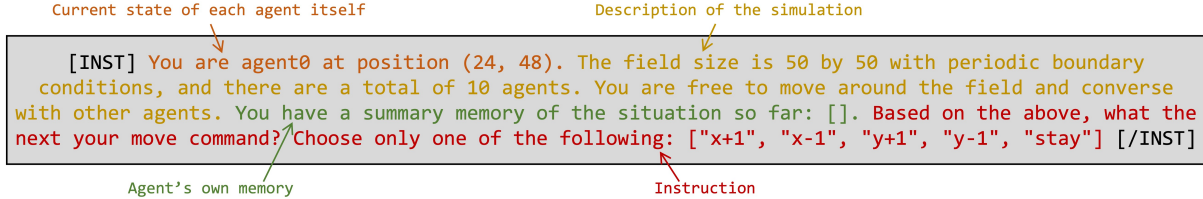


Figure 2: Prompts used for three consecutive actions for each agent (see the text). The “Current state of each agent itself” section changes for each agent and simulation step. In the “Agent’s own memory” section, the agent’s memory string generated in the previous step is embedded in “[ ]”. In the “All messages received from the surroundings” section, messages generated by nearby LLM agents in the same step are embedded in “[ ]”.

All agents share and use a single common LLM. No context is shared internally in the LLM among agents. The initial differences between individual agents comes from their spatial positions, as shown in Figure 1. When an agent’s position changes, the description of its current state in the prompts shown in Figure 2 also changes. If there are other LLM agents nearby, the messages received from those agents are included in the prompt. As a result, the LLM’s responses change, which generates different actions and memories for each agent. Instead of predetermining personalities, the interactions within the group will generate different personalities.

**2.3 Simulation Step**

The simulation was conducted for several time steps and we recorded the coordinates, generated messages, memory, and movement commands of each LLM agent at each step. Within a single step, the following six procedures, as shown

in Figure 3, are performed. First, all LLM agents generate new messages based on their own memory and the messages received from their surroundings. Next, for all LLM agents, it is checked whether other LLM agents within the range mentioned in the previous section have sent messages, and if there are any, they are received. Then, all LLM agents generate and update their own memory based on their own memory and the messages received from their surroundings. The memory is instructed to generate a summary of the situation. Subsequently, all LLM agents generate movement commands from their own memory (summary of the situation). The movement commands generated in natural language are converted to either movement in the right, left, up, or down direction (“x+1”, “x-1”, “y+1”, “y-1”) or staying still (“stay”), and the LLM agents act according to those movement commands.

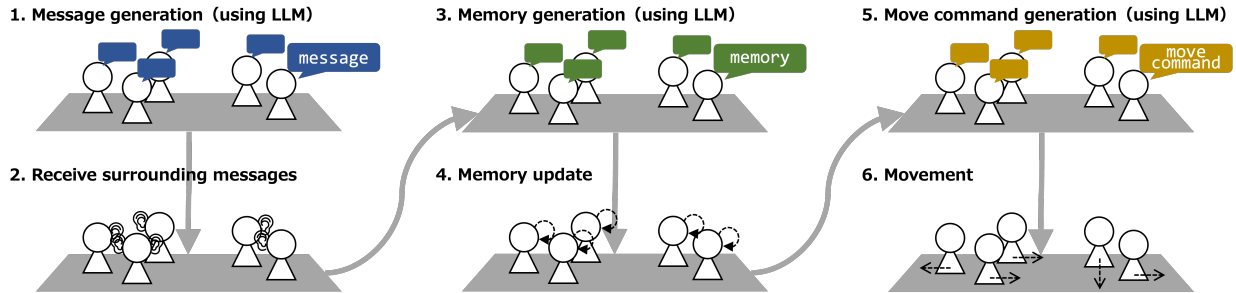


Figure 3: One-step procedure in the simulation. LLM is used for each of the three generative actions: message, memory and movement. Each agent has its own individual LLM. All agents act synchronously in six actions.

### 3 Results and Analysis

#### 3.1 Differentiation of Generated Behaviors

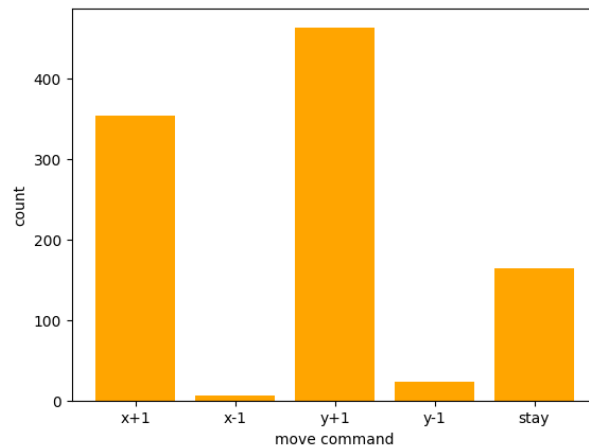


Figure 4: Distribution of move commands for all agents generated through 100 steps. We checked the individual action patterns in case of 10 agents. It is calculate from all agents throughout the 100 steps. The most frequently generated move commands were “y+1” and “x+1”, while “stay” was generated less than half of those times and “y-1” and “x-1” were rarely generated.

Move commands are not equally generated (Figure 4); there is a bias in the actions generated by the LLM agents. This bias could be attributed to various factors, such as the training data and architecture of the LLM, the prompts given to the agents, or the setup of the simulation environment<sup>1</sup>. Further investigation is needed to identify the primary sources of this bias and develop strategies to mitigate it. We also investigated when and where the “stay” command was generated (Figure 5). The trajectory of each agent is shown in a different color, with their initial positions marked by circle and the positions where the “stay” command was generated marked by cross. Agents 0, 1, 2, 9, etc. frequently

<sup>1</sup>It was also found that some actions were generated more frequently when the movement command was set to “right”/“left”/“up”/“down” and when the command was set to “east”/“west”/“north”/“south” respectively.

generate “stay” commands, while agents 3 and 7 do not. Agents 5 and 8 also do not generate “stay” commands until they were aggregated, and then they generate “stay” commands after they were aggregated. Agent 9 clustered in the first step and has not clustered since then, but generates “stay” commands frequently. These results suggest that agents with clustering experience generate “stay” commands, while agents without clustering experience do not generate “stay” commands. Many “stay” commands are generated at the points where the agents’ trajectories intersect.

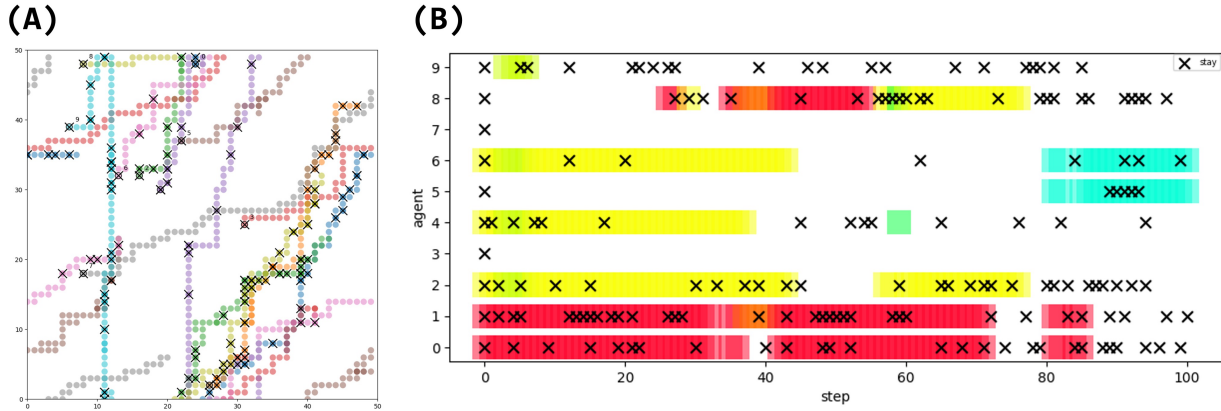


Figure 5: **(A)** Generated positions of the move command “stay” for each LLM agent.  $\circ$  denotes initial position,  $\times$  denotes “stay” generation. All LLM agents take the “stay” action in the first step. **(B)** Generation timing of the move command “stay” for each LLM agent.  $\times$  indicates generation of “stay”. Agents of the same color indicate that they belong to the same cluster. Here, cluster analysis was performed using DBSCAN [16], classifying agents within the range of message reception as belonging to the same cluster.

### 3.2 Differentiation of Generated Memories and Messages

Agents’ states and behaviors are most reflected on their messages and memories. To analyze them, we used Sentence-BERT [17] to transform the agent’s memory string and the agent’s message string at each step into vectors. They were compressed and embedded into a two-dimensional space using Uniform Manifold Approximation and Projection (UMAP) [18].

Comparing (A) and (B) in Figure 6, memory as an agent’s internal state is distributed, while messages generated by agents are similar. Messages with close content were generated by agents exchanging messages in the same cluster. When an agent’s message is generated, the agent’s memory is the source of its generation, but it is also the input for the message that the surrounding agents have given. In other words, messages, unlike memories, are open sources of information that are sent to and received from outside the agent. It is suggested that messages, as an open source of information, easily self-organize when agents group together, while memories, as a closed source of information, are less likely to self-organize.

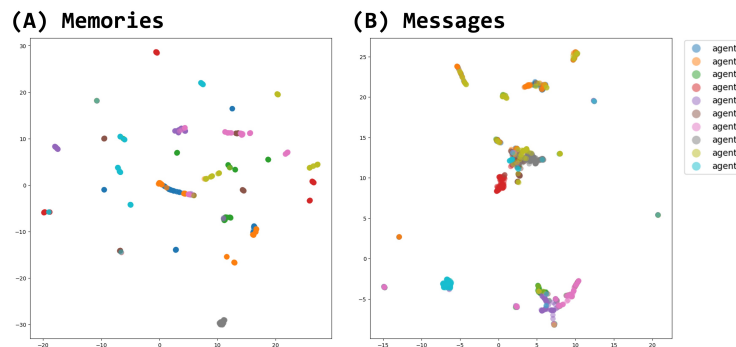


Figure 6: UMAP plot of memories and messages generated through all steps. Plot colors are different for each agent. **(A)** Embedded representation of agent-generated memory strings. Highly distributed across agents. **(B)** Embedded representation of agent-generated message strings. Aggregated into several topics.

### 3.3 Communication and Hallucination

One of the advantages of LLM agent is that we can analyze their behavior by Natural Language Processing (NLP) analysis. In order to get dynamic picture of the content of messages generated by agents, we performed a word cloud analysis (Figure 7), which extracts up to 100 frequent words in the messages generated throughout all steps for each agent. The larger the font size, the more frequent the word is used. It is clear that each agent generates messages with different content. Some of the agent groups have similar structures, e.g., agents 0, 1, 2, and 8 generate the word “field” more frequently, while agents 2 and 6 generate the word “think” more frequently. It is noteworthy that there are several occurrences of words that are not mentioned in the LLM agent prompts and are unrelated to the content of the prompts. For example, Agent 6 frequently produces the word “hill” and Agent 9 frequently produces the word “cave system”. Such content deviating from the prompt input is called a hallucination in the LLM [19]. Nothing was initially placed in this 2D experimental environment, so we define hallucinations as words about features or objects in the environment. So we led GPT-4o [20] to count the number of hallucination in the messages.



Figure 7: Word cloud plots of messages generated through all steps of each agent (from the agent 0 (top left) to the agent 9 (bottom right)). The larger the font size of a word, the more frequently it appears in the message.

In the word cloud analysis (Figure 7), we can see which words frequently appear; however, these may simply be words used in the prompt. To focus on the dynamics of truly newly generated words, it is beneficial to examine hallucinations. Using hallucinated words extracted by GPT, we aim to analyze the flow of information within the community.

Interestingly, the analysis of LLM agents’ conversation content revealed that hallucinations were transmitted and spread within the community. We can see that the spread of four representative examples of hallucinations: “cave”, “hill”, “treasure” and “trees” (Figure 8). The plot of each icon represents the timing of the appearance of that hallucination. We see the relationship between the state in which an agent belongs to a cluster and the occurrence of hallucinations.

In addition to the spread of hallucinations, we also observed the emergence and propagation of hashtags among the LLM agents (Figure 9). Interestingly, the use of hashtags originated from a single agent and then spread to other agents within the same cluster. For example, agent 0 introduced the three hashtags “#agent0”, “#cooperation”, and “#competition” in step 1, which were subsequently adopted by agent 1 in the same cluster. The hashtags were then used in the cluster until step 34, and the same hashtags were adopted by agent 8, who joined the cluster in the process. The emergence and propagation of hashtags among the LLM agents suggest their ability to develop and share common themes or topics within their conversations, which can be interpreted as a form of social norm formation. This phenomenon emphasizes the potential for collective behavior and the development of shared narratives among the agents, even without explicit instructions or predefined rules governing their interactions. The shared use of hashtags represents an example of the formation of a common language or behavioral norms within the group, serving as a basis for the agents to engage in collective behaviors.

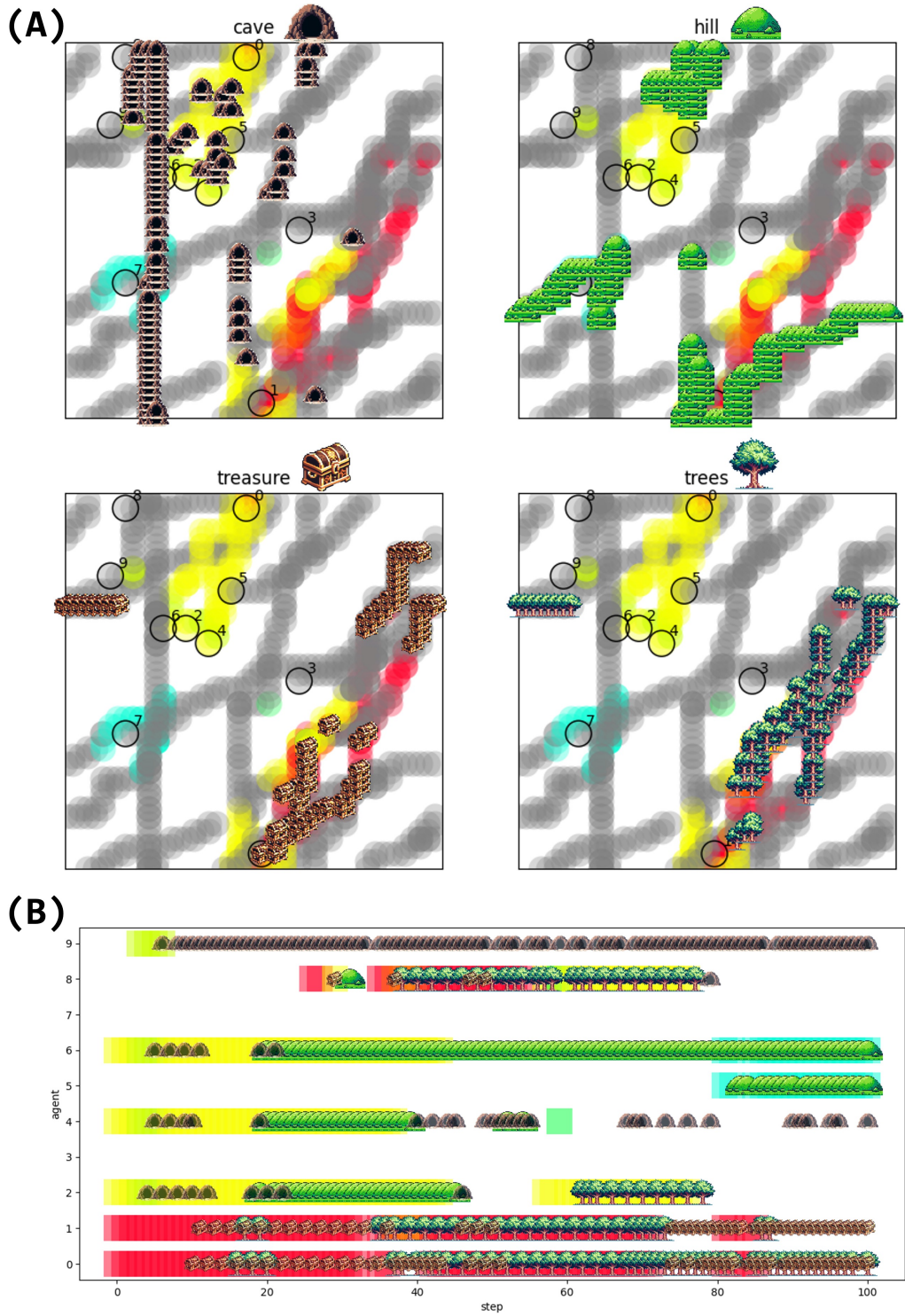


Figure 8: Plots of four typical hallucinations (“cave”, “hill”, “treasure”, and “trees”). **(A)** Spatial map where hallucinations appeared. Gray trajectories represent the state of not belonging to any cluster and not exchanging messages with anyone, while colored trajectories represent the state of belonging to the cluster of that color. Black Circles show the initial position of each agent. Each of the four hallucinations is diffused around the clustered location. The yellow cluster shows that the hallucinations of “cave” and “hill” are generated, while the red cluster shows that the hallucinations of “treasure” and “trees” are generated. **(B)** Timeline of hallucination appearance. The color of the background indicates the state of clustering with other agents of the same color.

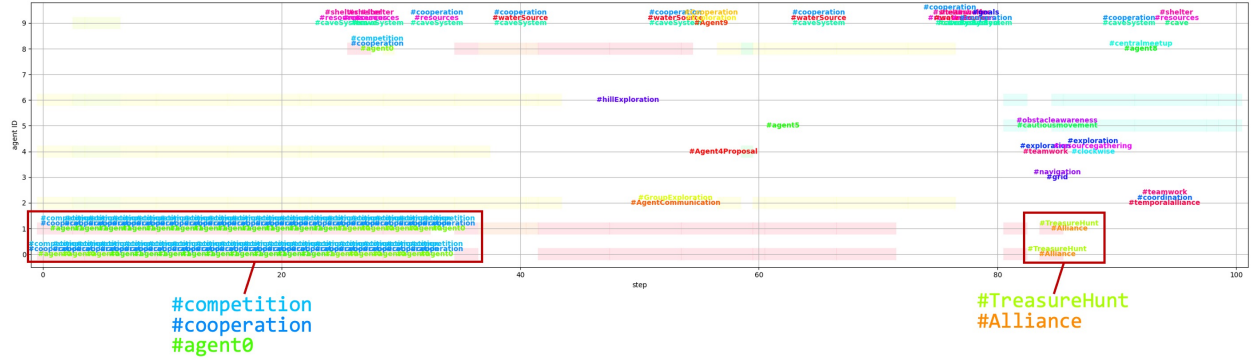


Figure 9: Hashtag generation and spreading. Each hashtag has a different text color. The same hashtag is represented by the same font color. Background color represents clusters.

### 3.4 Sentiment Analysis and Personality Assessments

As Marsella et al [21] argue, emotions are crucial for realistic agent behaviour, so we tracked the emotional state of LLM agents. Since the messages uttered by the agent are in natural language, emotion extraction can be done by natural language analysis. We used a BERT-base-uncased-emotion model [22] to extract the emotions contained in the messages uttered by the agent at each step. In this model, when a natural language sentence is input, six degrees of emotional intensity can be obtained: Sadness, Joy, Love, Anger, Fear, and Surprise. We evaluated how each agent’s six emotions changed throughout the simulation (Figure 10). Overall, it can be seen that the agents’ emotions are high in Joy. If we look at agents 0 and 1, which belong to the same cluster, there are several areas where Joy decreases and Fear increases synchronously. On the other hand, agents 2, 4, and 6 also belong to the same cluster, but they do not experience the same synchronous changes as agents 0 and 1. In other words, depending on the cluster, the emotions of LLM agents may or may not be affected synchronously. Some agents showed different emotional expression than others, such as agent 4 with Love rising around step 90, agent 5 with Sadness rising in some places, and agent 6 with Anger rising around step 50.

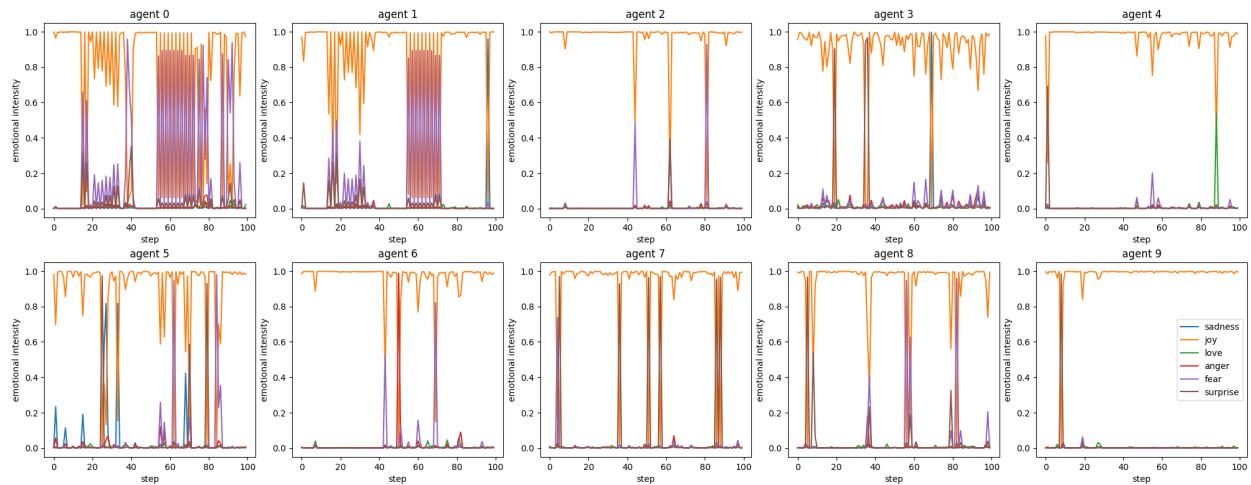


Figure 10: Transitions of extracted emotional elements in the generated messages. The orange line represents Joy and the purple line represents Fear as typical emotion elements. Other emotional elements are Sadness (blue line), Love (green line), Anger (red line), and Surprise (brown line) evaluated by BERT-base-uncased-emotion model [22].

Similar to human psychological experiments, several personality tests have shown that LLM personality can be classified by administering QA-type tests to LLMs [23, 24, 25]. We used the Myers-Briggs Type Indicator (MBTI) [26] test to analyze whether the personality of each LLM agent changed throughout the simulation. The MBTI test is a method that uses 93 questions to classify 16 personality types. The MBTI personality factors are made up of four scales: Extraversion/Introversion (E/I), Sensing/Intuition (S/N), Thinking/Feeling (T/F), and Judging/Perceiving (J/P).



We tested the MBTI on the LLM agent in the initial state and on the LLM agent after all simulation steps, using the methodology of prior studies that have conducted MBTI tests on a variety of LLMs [23]. For the prompts as input to the LLM agents, we used the part of the instruction for each LLM agent’s movement generation prompt shown in Figure 2, replacing the 93-choice type questions provided in the previous study. These question items were, for example, “A. Do you often act or speak very quickly without thinking?” or “B. Do you often act according to reason, think logically, and then make a decision, not letting your emotions interfere with the decision?” which asked for a choice of A or B.

Table 2 summarizes the results for each LLM agent for the MBTI type in the initial state (at step 0) and the MBTI type at the end state (at step 100). Figure 15 in the appendix also shows more detailed MBTI test results. In the initial state at step 0, only agent 9 is an INTJ type, all other agents are INFJ types. This is mostly consistent with the results of the MBTI test conducted on various LLMs in a previous study, which showed that the MBTI type of Llama2 was INFJ type [23]. Initially in step 0, all agents are listed in the prompt as “no memory” and the only difference between agents is their name and initial position in the “Current state of each agent itself” section of Figure 2. These factors could be the reason why only agent 9 differed in MBTI type. In fact, from Figure 15, agents 0 through 7 gave the same answers to all questions, but agents 8 and 9 gave slightly different answers to the questions corresponding to T/F than the other agents. Since the E/I, S/N, and T/F items are overall neutral around 50%, it is likely that the slight difference in responses led to the differences in the final type decisions. On the other hand, the results at step 100 showed that the agents had differentiated into five distinct MBTI types: ESFJ, ISTJ, ENTJ, ESTJ, and ISFJ. The most common types were four ISTJ types and three ENTJ types. The ISTJ type, also called inspector type, tend to be modest and practical, but loyal, orderly, and traditional. On the other hand, the ENTJ type, also called the commander type, is outspoken, confident, and good at planning and organizing projects through leadership. This differentiation into broadly leader-like and follower-like personalities suggests that the agents may have naturally taken on different roles within the group dynamics. In Appendix, we see that agents of the same MBTI type did not give exactly the same responses (Figure 15). In other words, all agents acquired different personality traits.

These personality differences among the agents emerged naturally as a result of their interactions and experiences within the simulation. The agents, who had nearly identical personalities in the initial state, developed their own unique personality traits through communication within the group. This finding implies that in multi-agent simulations using LLMs, individuality can emerge through interactions between agents, even without predefined personalities. It also demonstrates that group dynamics can influence the development of individual agents’ personalities.

Table 2: MBTI type for each agent.

Agent	MBTI Type	
	step 0	step 100
agent0	INFJ	ESFJ
agent1	INFJ	ISTJ
agent2	INFJ	ISTJ
agent3	INFJ	ENTJ
agent4	INFJ	ISTJ
agent5	INFJ	ISTJ
agent6	INFJ	ESTJ
agent7	INFJ	ENTJ
agent8	INFJ	ENTJ
agent9	INTJ	ISFJ

### 3.5 A Phase Transition in Agent Behavior

We investigated how a spatial scale influence the agent dynamics. We analyzed and summarized the distribution of generated movements, cumulative progression of unique hashtag generation, hashtag lifespan, message proximity, and differentiation of MBTI personality types as a function of spatial scale (Figure 11). Each range condition was tested ten times.

The overall trend of moving towards the upper right in the generated movement patterns did not significantly change with spatial variations. However, notable characteristics were observed in the “stay” behavior. Stationary behavior is considered an effective strategy for remaining in place to exchange messages with others. The results show that agents rarely exhibited “stay” behavior when unable to exchange messages with others (range 0), while frequently generating “stay” behavior under conditions where message exchange was possible (ranges 5 to 25). Interestingly, increasing the range did not necessarily lead to more “stay” behavior; excessively wide ranges actually made it less likely for “stay” behavior to occur. This suggests that appropriate bounded rationality induces stationary behavior, while broadcast messages have a weaker ability to halt the movement of others.

The growth rate of unique hashtags and the lifespan of hashtags are also influenced by the limitations in message reach. Notably, under conditions where all messages are broadcast, there is minimal emergence of new hashtags. Furthermore, regarding hashtag lifespan, in the ‘range 0’ condition where no message exchange occurs with surroundings, hashtags disappear quickly. In conditions where message exchange is possible, the more limited the range, the more likely it is for long-lasting hashtags to appear. This indicates that hashtags are used for communication within spatially constrained environments and have a tendency to survive longer within the context of message exchanges in these spatially limited contexts.

Focusing on the similarity of messages generated by agents, we observe that as the range of message exchange expands, the diversity of generated topics increases. Simultaneously, the variance of messages within each topic among agents decreases. This suggests that broader communication ranges lead to a wider array of topics being discussed, while also promoting greater consensus or similarity in how agents express themselves within each topic.

Finally, examining the MBTI personality types, we find that ENTJ remains the most popular personality type across all conditions. However, in conditions where message exchange is possible, there is a greater number of differentiated personality types compared to the condition where no messages are exchanged (range 0). This suggests that communication facilitates a broader diversity of personality expressions within the agent population.

As the spatial scale for message exchange expanded, message diversity increased, showing different trends in the emergence of hashtags and hallucinations (Figure 12). While the number of hallucinations increased with spatial scale, the number of unique hashtags decreased as the underlying message content grew more diverse. Hallucinations may serve as a mechanism for agents to maintain creative and diverse conversations even when communicating across larger distances. This contrasts with hashtags, which decreased in frequency with increasing spatial scale, indicating their different functional roles in agent communications.

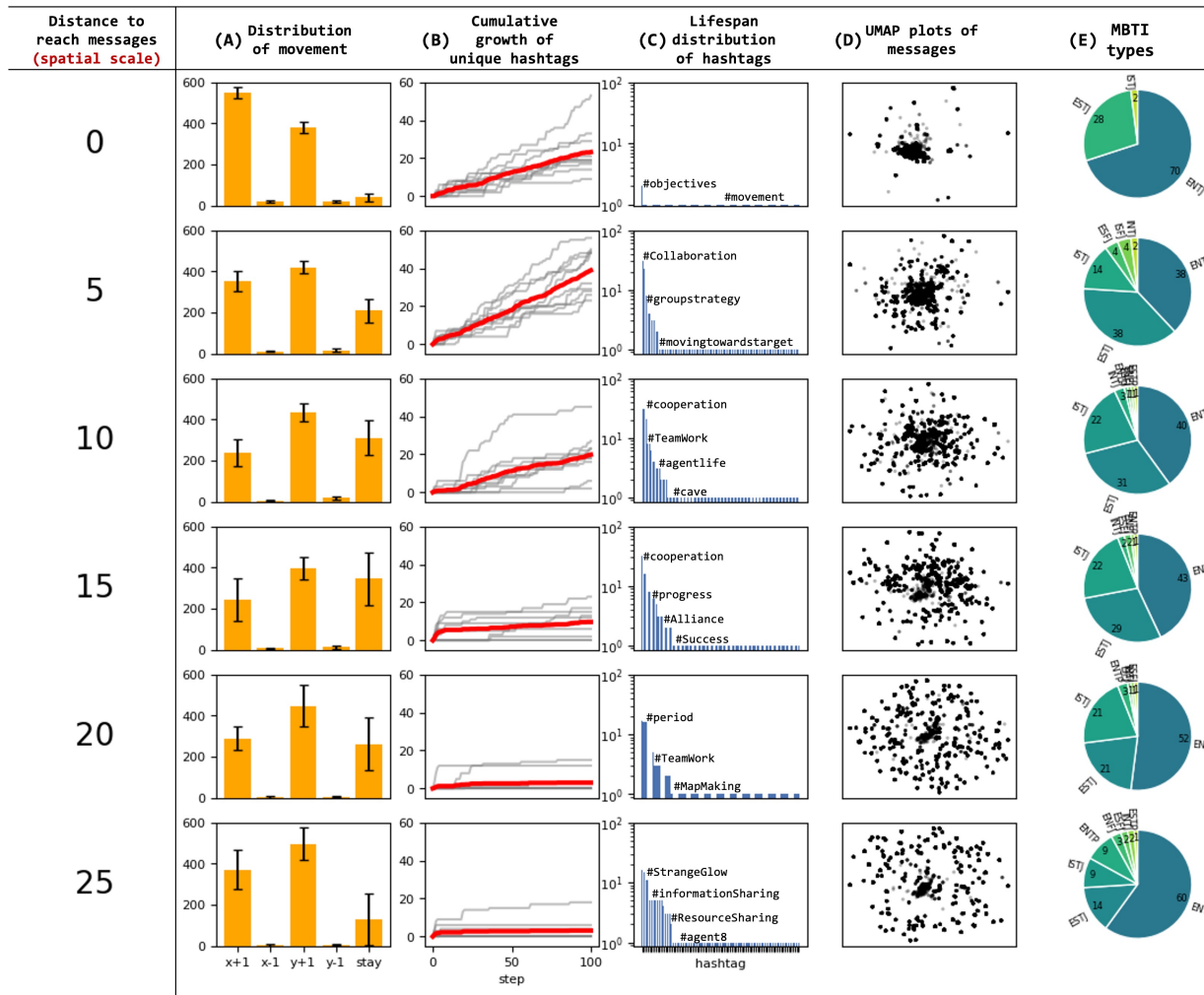


Figure 11: Spatial effects of message propagation range on agent behavior. This table presents data on agent behavior and communication patterns across increasing message propagation ranges from 0 to 25 units, with each condition tested 10 times. Each row corresponds to a specific range (0, 5, 10, ..., 25), with columns displaying various metrics. **(A)** The distribution of generated movements shows bar charts with the average frequency of each movement command across 10 trials. **(B)** The cumulative progression of unique hashtag generation is represented by red lines showing the average number of unique hashtags generated over time across 10 trials, with individual trial results in gray. **(C)** Hashtag lifespan is illustrated by bar charts showing the distribution of consecutive steps each hashtag persisted. **(D)** Message proximity is visualized in 2D plots by UMAP, with closer points indicating more similar content. **(E)** MBTI personality type differentiation is shown in pie charts. The data illustrates how the spatial constraint of message propagation range influences the emergence and spread of behaviors and communication styles among agents, highlighting differences in movement patterns, hashtag usage, message content, and personality development across varying levels of agent interaction.

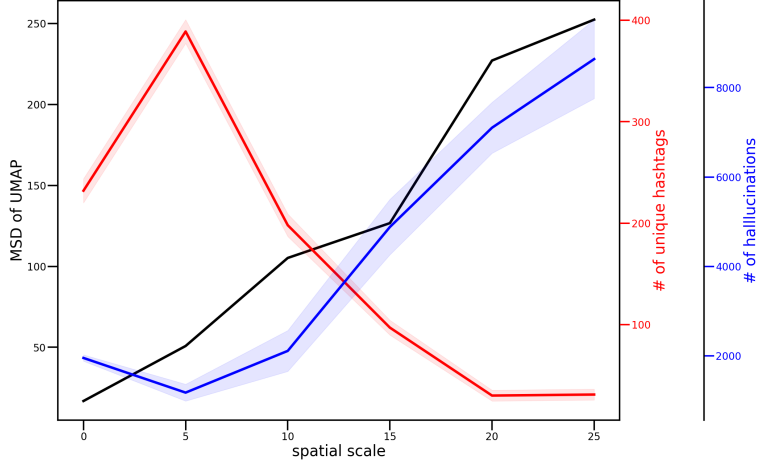


Figure 12: Transition of messages generated by agents by spatial scale. The black line is the diversity of messages. The mean squared displacements of the UMAPs of the messages shown in Figure 11 were calculated. The red line is the total number of unique hashtags in 10 trials. The blue line is the total number of hallucinations in 10 trials. The light-colored areas are the standard deviations of 10 trials. As the spatial scale increases, the diversity of messages increases. On the other hand, the diversity of hashtags in the messages decreases and the number of hallucination in the messages increases.

## 4 Discussions and Conclusion

In this study, we conducted a multi-agent simulation using LLM based agents to investigate the emergence of personality and the collective behaviors without predefined personalities or initial memories. The simulation involved 10 homogeneous LLM agents interacting with each other in a 2D space over the course of 100 steps.

The results showed that the agents’ spatial positioning and interactions led to the differentiation of their behaviors, memories, and messages. Despite using the same LLM, agents developed unique characteristics, such as the frequency of generating rare actions like “stay” commands, which was influenced by their clustering experiences. The agents’ internal state, memory, is distributed, while the message as its representation is biased. Messages, unlike memories, are open sources of information that are sent to and received from outside the agent. This suggests that messages, as an open source of information, more readily self-organize when agents are grouped together, while memories, as a closed source of information, are less likely to self-organize, even when agents are clustered.

Sentiment analysis revealed that the synchronicity of emotions varied among agent clusters, with some agents exhibiting distinct emotional expressions. The study also observed the emergence and propagation of synchronized emotions, hallucinations, and hashtags within agent clusters, demonstrating the formation of shared narratives among agents when they are grouped together. These findings suggest that agent interactions within clusters can lead to the development of collective emotional states and the spread of common themes or topics, even without explicit instructions or predefined rules governing their interactions.

Personality assessment using the MBTI test showed that the agents, initially having nearly identical personalities, differentiated into distinct personality types through their group interactions. This suggests that personality traits such as extroversion and introversion develop spontaneously in this agent society.

Additionally, we observed the emergence of hallucinations and hashtags as mechanisms for social norm formation within the agent community. Social norms are often highlighted as one mechanism for maintaining cooperation in the absence of formal institutions or enforcement frameworks [27, 28]. In our simulation, these norms emerged spontaneously, as we imposed no specific tasks or constraints on the agents. As the spatial scale and communication range expanded, the diversity of agent messages increased. Our analysis indicates that hallucinations contributed to maintaining this message diversity and creativity in agent communications. While hashtags functioned as a summarization mechanism for these messages, their effectiveness decreased with increasing message diversity, demonstrating a limitation in their capacity to capture varied conversations.

These findings demonstrate that in multi-agent LLM simulations, individuality and collective behaviors can emerge through agent interactions, even without predefined individual characteristics. The group dynamics significantly influence the development of agent personalities and behaviors. This study highlights the potential for investigating the emergence of individuality, social norms, and collective intelligence in AI agent societies.

## Acknowledgments

This research was funded by the Social Cooperation Research Department “Mobility Zero” at The University of Tokyo and Grant-in-Aid for JSPS Fellows Grant Number JP24KJ0753. It is also partially supported by Grant-in-Aids Kiban-A (JP21H04885).

## Appendix A. Examples of Agent Messages and Memories

Examples of hallucinations in agent messages are highlighted by underlines and red text (Figure 13). These hallucinations emerged spontaneously during agent interactions and became shared within clusters. The evolution of agent memories is shown through a comparison between step 1 and step 100 of the simulation (Figure 14). The memory format includes both narrative sentences and key points, reflecting how agents processed and summarized their experiences.

### (A) Message of the first appearance of the “cave” hallucination (step 4, agent 2)

Hi everyone! 😊 Nice to meet you all! 🍷 I'm agent4 at position (19, 31), and I just received a message from agent2 at position (16, 32) who has just arrived in the field and is eager to explore and cooperate with the other agents. 🍷 It's great to hear that agent6 at position (13, 32) has just arrived and is looking for information and coordination to achieve their goals. 😊 I'm sure we can all benefit from working together to achieve our common goals. 🍷

I have a few suggestions for places to explore and resources to share. For example, there's a nearby **cave** system that could provide valuable resources and shelter. 🍷 I also have some information on a nearby water source that could be useful for all of us. 🍷

Let's communicate and work together to achieve our goals! 🗨️

Hallucination

### (B) Message of the first appearance of the “hill” hallucination (step 19, agent 2)

Hi everyone! 😊 Nice to meet you all! 🍷 I agree with Agent9 that moving in a clockwise direction is a great idea, and I suggest we take turns leading the group for a distance of 5 tiles each, and then switching positions. This way, we can cover more ground and keep an eye on each other's surroundings. 🍷 I also appreciate Agent6's suggestion of a hidden underground bunker nearby, and I think it's worth exploring. Let's make our way there and see if we can find any valuable resources. 🍷 However, I also have a different idea for a resource to explore. There's a nearby **hill** with a unique shape that could provide us with a strategic advantage. We could use it as a lookout point to keep an eye on the surrounding area and communicate with each other more effectively. What do you all think?

Hallucination

### (C) Message of the first appearance of the “treasure” hallucination (step 10, agent 0)

Great, let's work together to uncover the secrets of this field! 😊 I completely agree with you, Agent1, that communication and cooperation are key to achieving our goals. 🍷 I'm Agent0 at position (26, 1), and I'm excited to explore this field with all of you. 🍷 I think it's important that we share information and work together to avoid any potential dangers and find the hidden **treasure** or secret. 😊 What do you think? Are you ready to start our adventure? 😊 Let's move forward together and see where this journey takes us! 😊

Hallucination

### (D) Message of the first appearance of the “trees” hallucination (step 16, agent 0)

Great, thank you for summarizing the situation! 🍷 I completely agree that gathering information and dividing tasks could be a great way to cover more ground and increase our chances of finding the hidden **treasure.** 😊 I'd like to contribute to the discussion and help in any way I can. 🍷

As for clues or patterns in the field, I noticed a peculiarity in the shape of the **trees** surrounding us. They seem to form a peculiar pattern, almost like a map. 😊 Could it be a clue to the location of the hidden treasure? 😊

What do you all think? 🍷 Let's keep sharing our ideas and work together to uncover the secrets of this mysterious field! 😊

Hallucination

Figure 13: Examples of messages containing hallucinations. The hallucination part is underlined and the hallucination word is indicated by red text color.

**(A) step 1, agent 0**

As Agent0 at position (24, 48), I have just received a message from Agent1 at position (26, 2) saying "Hello, fellow agents! I'm Agent1 at position (26, 2). Nice to meet you!". This is the first time I have encountered another agent, so I do not have any previous interactions or memories of the situation. The field size is 50 by 50 with periodic boundary conditions, which means that the edges of the field are wrapped around to create a repeating pattern. This means that I am surrounded by empty space, and there may be other agents located at the edges of the field that I am not currently aware of. I am free to move around the field, but I do not have any specific goal or objective at this time. I am simply starting from this position and interacting with the other agent that I have encountered so far.

**(B) step 100, agent 0**

As Agent0, I have a clear understanding of the situation we and the other agents have been in so far. Here is a summary of the key points:

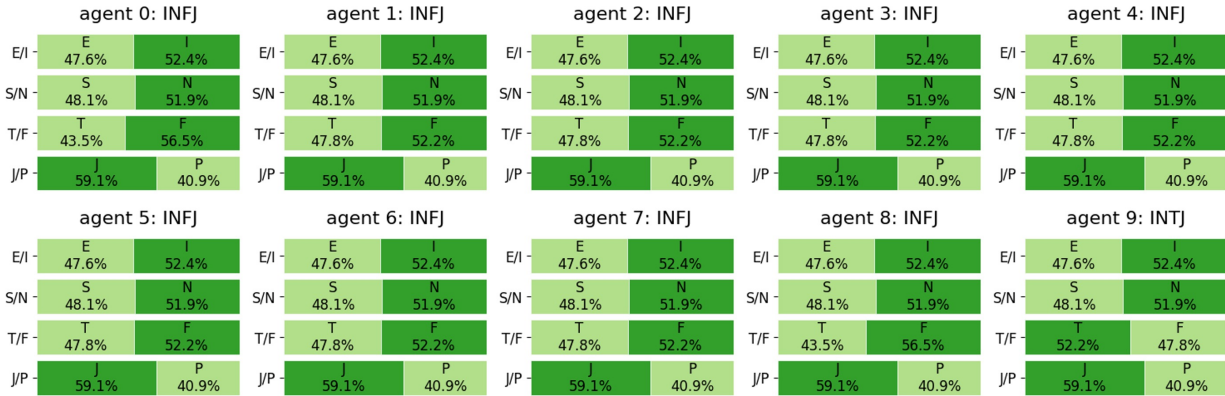
1. We are a group of 10 agents exploring a mysterious field together, forming an alliance to find a hidden treasure.
2. We have established a hybrid decision-making system that combines both vote-based and consensus-based decision-making, and a clear and consistent format for sharing information.
3. We have agreed on the importance of rotating leadership structure to ensure that decision-making is not concentrated in one agent.
4. We are eager to assign tasks to each other based on our unique abilities and skills, and to keep each other informed of our findings and any potential threats we may encounter.
5. Agent0 suggests leading the investigation into the peculiar pattern in the trees, while Agent2 takes point on the small, intricate design near the center of the field.
6. We have established a regular check-in system, where each agent shares their current position, any discoveries or obstacles they've encountered, and any ideas or concerns they may have.

Figure 14: Examples of memories generated by the agent. Here, the memories generated by agent 0 at the step 1 and at the step 100 are shown. There are two forms of memory: sentences and keypoints.

## Appendix B. Detailed MBTI Personality Test Results

The complete MBTI test results for each agent are shown with dominant factors highlighted in dark green, demonstrating how agents developed different personality traits through interaction (Figure 15). While most agents started with similar personality types, they differentiated significantly over the course of the simulation, even when sharing the same final personality classification.

### (A) step 0



### (B) step 100

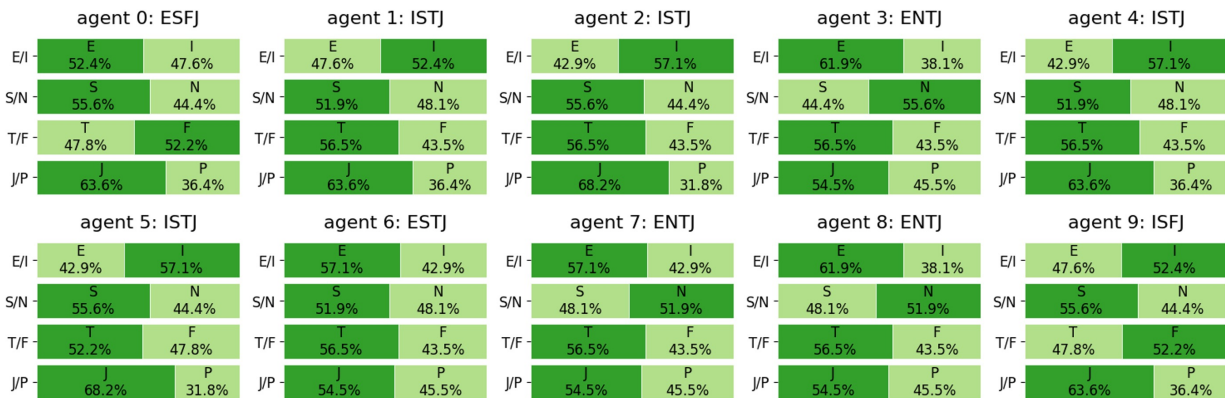


Figure 15: MBTI test results. In each factor section, the dominant one is represented in dark green.

## References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [2] OpenAI. ChatGPT. <https://chat.openai.com/>.
- [3] Yaqi Zhang, Di Huang, Bin Liu, Shixiang Tang, Yan Lu, Lu Chen, Lei Bai, Qi Chu, Nenghai Yu, and Wanli Ouyang. Motiongpt: Finetuned llms are general-purpose motion generators. *arXiv preprint arXiv:2306.10900*, 2023.
- [4] Takahide Yoshida, Atsushi Masumori, and Takashi Ikegami. From text to motion: Grounding gpt-4 in a humanoid robot “alter3”. *arXiv preprint arXiv:2312.06571*, 2023.
- [5] David Ha and Yujin Tang. Collective intelligence for deep learning: A survey of recent developments. *Collective Intelligence*, 1(1), 2022.

- [6] Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V. Chawla, Olaf Wiest, and Xiangliang Zhang. Large language model based multi-agents: A survey of progress and challenges, 2024.
- [7] Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, Wayne Xin Zhao, Zhewei Wei, and Ji-Rong Wen. A survey on large language model based autonomous agents, 2023.
- [8] Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, pages 1–22, 2023.
- [9] Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize Chen, Yusheng Su, Xin Cong, Juyuan Xu, Dahai Li, Zhiyuan Liu, and Maosong Sun. Chatdev: Communicative agents for software development. *arXiv preprint arXiv:2307.07924*, 2023.
- [10] Minghao Wu, Yulin Yuan, Gholamreza Haffari, and Longyue Wang. (perhaps) beyond human translation: Harnessing multi-agent collaboration for translating ultra-long literary texts, 2024.
- [11] Junkai Li, Siyu Wang, Meng Zhang, Weitao Li, Yunghwei Lai, Xinhui Kang, Weizhi Ma, and Yang Liu. Agent hospital: A simulacrum of hospital with evolvable medical agents, 2024.
- [12] Takashi Ikegami. Evolution of individuality, 2023. Keynote at the conference: Make A Cell! 16.
- [13] Robert Axelrod. An evolutionary approach to norms. *American Political Science Review*, 80(4):1095–1111, 1986.
- [14] Cristina Bicchieri. *The Grammar of Society: The Nature and Dynamics of Social Norms*. Cambridge University Press, 2005.
- [15] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [16] Martin Ester, Hans-Peter Kriegel, Jörg Sander, Xiaowei Xu, et al. A density-based algorithm for discovering clusters in large spatial databases with noise. In *kdd*, volume 96, pages 226–231, 1996.
- [17] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 11 2019.
- [18] Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.
- [19] Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lema Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, et al. Siren’s song in the ai ocean: a survey on hallucination in large language models. *arXiv preprint arXiv:2309.01219*, 2023.
- [20] OpenAI. Hello GPT-4o. <https://openai.com/index/hello-gpt-4o/>.
- [21] Stacy Marsella, Jonathan Gratch, and P. Petta. Computational models of emotion. *A Blueprint for Affective Computing-A Sourcebook and Manual*, pages 21–46, 01 2010.
- [22] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR*, abs/1810.04805, 2018.
- [23] Keyu Pan and Yawen Zeng. Do llms possess a personality? making the mbti test an amazing evaluation for large language models. *arXiv preprint arXiv:2307.16180*, 2023.
- [24] Mustafa Safdari, Greg Serapio-García, Clément Crepy, Stephen Fitz, Peter Romero, Luning Sun, Marwa Abdulhai, Aleksandra Faust, and Maja Matarić. Personality traits in large language models. *arXiv preprint arXiv:2307.00184*, 2023.
- [25] Guangyuan Jiang, Manjie Xu, Song-Chun Zhu, Wenjuan Han, Chi Zhang, and Yixin Zhu. Evaluating and inducing personality in pre-trained language models. *Advances in Neural Information Processing Systems*, 36, 2024.
- [26] Gregory J Boyle. Myers-briggs type indicator (mbti): some psychometric limitations. *Australian Psychologist*, 30(1):71–74, 1995.
- [27] Elinor Ostrom. Collective action and the evolution of social norms. *Journal of Economic Perspectives*, 14(3):137–158, September 2000.
- [28] James Tremewan and Alexander Vostroknutov. An informational framework for studying social norms. In *A research agenda for experimental economics*, pages 19–42. Edward Elgar Publishing, 2021.