

# Can Agents Spontaneously Form a Society?

Introducing a Novel Architecture for Generative Multi-Agent to Elicit Social Emergence

Hanzhong Zhang

Kunming University of Science and Technology  
Kunming, Yunnan Province, China  
armihia@foxmail.com

Mulin Jiang

Kunming University of Science and Technology  
Kunming, Yunnan Province, China  
237800718@qq.com

## ABSTRACT

Generative agents have demonstrated impressive capabilities in specific tasks, but most of these frameworks focus on independent tasks and lack attention to social interactions. We introduce a generative agent architecture called ITCMA-S, which includes a basic framework for individual agents and a framework called LTRHA that supports social interactions among multi-agents. This architecture enables agents to identify and filter out behaviors that are detrimental to social interactions, guiding them to choose more favorable actions. We designed a sandbox environment to simulate the natural evolution of social relationships among multiple identity-less agents for experimental evaluation. The results showed that ITCMA-S performed well on multiple evaluation indicators, demonstrating its ability to actively explore the environment, recognize new agents, and acquire new information through continuous actions and dialogue. Observations show that as agents establish connections with each other, they spontaneously form cliques with internal hierarchies around a selected leader and organize collective activities.

## KEYWORDS

Generative agents, Multi agent system, Social interaction, LLM

## 1 INTRODUCTION

Large language models (LLMs) have contributed to significant progress in the field of natural language processing and are widely used in various domains, such as machine translation [44], dialogue generation [9], and content creation [46]. These models are capable of correctly parsing and generating complex sentence structures and have demonstrated unprecedented capabilities in understanding language. However, LLMs often lack true comprehension and rely more on pattern matching and probabilistic predictions [10]. It has also been difficult to create systems to achieve human-like systematic generalization [22]. To overcome these problems, researchers have introduced LLM-based agents. This has allowed the incorporation of external knowledge bases to supplement a model's knowledge gaps in specific domains [43]. They can also be used to decompose a complex task into multiple simpler tasks to achieve hierarchical processing [41].

Jibin Yin

Kunming University of Science and Technology  
Kunming, Yunnan Province, China  
yjblovelh@aliyun.com

Cong Su

Kunming University of Science and Technology  
Kunming, Yunnan Province, China  
787358363@qq.com

On this basis, Park et al. [31] introduced a novel LLM-based agent, namely, a generative agent. This agent simulates trustworthy human behavior. Generative agents have the ability to make multifaceted inferences about an environment, themselves, and other individuals in the environment. They can design daily activity plans based on their own characteristics and experiences and adjust their plans to changes as they occur. When a situation changes, they can flexibly update plans to ensure adaptation to it.

The importance of this progress cannot be ignored. In human – computer interaction, especially in virtual assistants, customer service robots, and even more complex systems, such as self-driving cars and smart homes, the ability of agents to generate believable human behavior is crucial. Generative agents can support more adaptive and flexible interaction processes. This ability not only enhances the system's responsiveness to dynamic situations but also brings human–computer interaction closer to natural behavioral patterns in interpersonal communication. In addition, by simulating human behavior, generative agents can demonstrate autonomy and sociality in various complex situations, making the interaction process smoother and more intuitive. This feature plays an important role in improving the user experience and increasing the trustworthiness of a system [16].

However, existing generative agent architectures still face many challenges. While traditional agent structures are good at processing and generating behaviors, they are primarily designed for isolated tasks and, thus, mostly lack a focus on sociality. This often makes it difficult for them to model and apply the nuances of social interactions, leading them to focus only on completing tasks and overlook behaviors that promote social connections. This is clearly not conducive to cooperation among multiple agents and may lead to behaviors that are detrimental to the group [34]. In scenarios involving multiple agents, a lack of structured social behavior may lead to disjointed or even chaotic interactions. To truly harness the potential of these agents in domains requiring interaction, there is an urgent need to explore their ability to participate in social interactions, establish relationships, and exhibit emerging social behaviors.

In this paper, we improve upon an existing LLM-based agent architecture (the internal time-consciousness machine based agent [ITCMA] introduced by [49]) and propose ITCMA-S (the “S” signifies our contribution of social interaction) architecture to enable agents to adapt to multi-agent interaction scenarios. It contains a

basic framework for an individual agent and the LTRHA framework for social interaction among multiple agents (the latter is named for its four modules: locale & topic, resources, habitus, and action; they are described later in the paper). This provides a structured way for agents to identify and filter actions that are not conducive to social interaction, thereby guiding them to choose actions that are more conducive to improving the social atmosphere of a scene. It also allows them to form new relationships during interactions with other agents, while remembering the history of these interactions. The formation of these dynamic relationships occurs naturally and is not pre-programmed. This means that the agent has the ability to adapt to its social environment.

To understand whether and how ITCMA-S leads to the emergence of sociality among generative agents, we established an evaluation environment called IrollanValley. This is a sandbox world designed to test and allow the observation of generative agent behavior and social interactions. It contains six characters, eight areas, and six operational primitives. Each area has unique furniture and other items, and an agent can interact with the various characters and objects and move around in the areas. Based on this environment, we modeled the natural evolution of social relationships among multiple agents without identities in the environment. The human evaluation results showed that the agents had spontaneously developed good social skills on each evaluation indicator and were able to actively explore the environment, meet new agents, and acquire new information through continuous actions and conversations. By observing the environment, we found that, in the process of establishing connections among each other, the agents spontaneously formed a clique with internal hierarchies around an elected leader and organized collective activities under the leader's guidance. Agents not included in this clique were mostly in isolated states. Even when they were in the same room as other agents, they interacted less and focused mainly on their personal activities.

In summary, this paper makes the following contributions:

- We propose a generative multi-agent structure, ITCMA-S, which includes the structure of an individual generative agent and a multi-agent social collaboration framework, LTRHA.
- We established an environment for multi-agent social evaluation, IrollanValley, which consists of six roles, eight areas, and six operational primitives designed to assess the utility and sociality of multiple agents.
- We conducted validation and ablation studies in IrollanValley through human evaluation, and the results showed that ITCMA-S performed excellently across multiple indicators.
- Furthermore, we discovered that agents can actively explore the environment, cooperate with other agents through a division of labor, and spontaneously form small groups and leadership structures in complex social settings.

Section 2 of this paper introduces related works; Section 3 introduces the improvement of the individual agent structure of ITCMA-S on an existing generative agent architecture and includes an explanation of how it generates and infers action through perception, memory, and emotion-driven mechanisms; Section 4 introduces a multi-agent interaction framework in ITCMA-S, namely LTRHA; Section 5 introduces the experimental evaluation of the sociality of

ITCMA-S and discusses the process and results of the formation of the agents' social relationships; and Section 6 summarizes the main findings of the study and suggests future research potential.

## 2 RELATED WORK

### 2.1 LLM-based agent

Park et al. [31] introduced a generative agent based on an LLM and demonstrated that it can generate trustworthy individual behavior and sudden group behavior in simulations. For example, a generative agent will turn off a stove when it sees its breakfast burning and will stop to chat when it encounters other agents with whom it wants to talk. Based on this research, numerous LLM-based agents have emerged. Zhang et al. [45] proposed AppAgent, which constructs an agent to operate any smartphone application. Hong et al. [20] proposed an 18-billion-parameter visual language model (VLM) named CogAgent, which specializes in GUI understanding and navigation. Vezhnevets et al. [39] provided the Concordia library to simulate agent interactions in physical, social, and digital spaces. Among these, a special agent called the game master (GM) is responsible for simulating the environment of agent interactions. Agents take action by describing what they want to do in natural language, and the GM then translates their actions into appropriate implementations.

### 2.2 Research on the Structure of Generative Agents

The original structure of generative agents, as described by Park et al. [31], mainly consisted of three parts: memory flow, reflection, and planning. The concept of a chain of thought (CoT) was particularly important for improvements to the planning module [41]. CoT refers to the ability of an LLM to think and reason gradually through a series of steps or iterations, reflecting human cognitive processes. Traditional language models generate responses without clear intermediate steps, which can lead to suboptimal answers, especially in complex inference scenarios. CoT overcomes these limitations by introducing intermediate steps to enable language models to reason, thereby enhancing the model's ability to solve problems. Mondal et al. [28] suggested using knowledge graphs to enhance multiple patterns to help models solve complex problems, thereby triggering CoT functionality. Their proposed method, knowledge augmented multimodal (KAM)-CoT, decouples the inference process into two consecutive stages. In the first stage, practical reasons are provided, and in the second stage, the generated reasons are used as additional input to provide answers.

As an improvement to the memory module, Liu et al. [25] proposed the reasoning and acting through scratchpad and examples (RAISE) architecture, which was specifically designed to enhance the functionality of conversational agents. It combines a dual-component memory system, similar to the short-term and long-term memory functions of the human brain. Toy et al. [37] proposed a metacognitive module by improving the reflection module in generative agents, allowing agents to broadly consider their situations to create alternative strategies and improve their performance.

To improve the overall structure of generative agents, Wu et al. [42] explored an alternative mechanism that utilizes prior knowledge encoded in an LLM without affecting the trainability of agents. Their proposed three-step PET framework: Plan, Eliminate, and Track. Lin et al. [24] proposed the method SwiftSage, which can achieve both fast and slow thinking in complex interactive reasoning tasks. It effectively blends the advantages of behavioral cloning and LLM.

### 2.3 Research on the Sociality of Generative Multi-Agents

By interacting with each other, generative agents can exchange information, form new relationships, and coordinate activities. These social behaviors occur naturally rather than being pre-programmed. Over time, agents form new relationships and remember their interactions with other agents [31]. Generative Multi-Agents can simulate complex human systems. Wang et al. [40] designed a new framework called the mosaic expert observation wall (MEOW). In MEOW, real game data are processed by expert models trained on simulated data and converted into natural language prompts to assist LLM inference.

However, while an LLM can capture social norms, there is also research that suggests that an LLM cannot adequately understand social norms, especially culture-specific social norms. This defect may lead to conflicts between generative agents, especially when their underlying LLMs are trained on text corpora from different cultural backgrounds [21, 32]. Ren et al. [34] proposed the specification architecture CRSEC (named for its four modules: Creation & Representation, Spreading, Evaluation, and Compliance) for generative multi-agent systems to resolve this problem. Building on this foundation, these researchers examined the emergence of social norms in generative multi-agent systems. Ghaffarzadegan et al. [17] provided a new approach to developing models that reduces reliance on assumptions about human decision-making and utilizes the vast amount of data in LLMs to capture human behavior and decisions. By utilizing the extensive dataset in the LLM, their generative agent-based model helps represent human decisions in computational models. Motwani et al. [29] pointed out that groups of generative agents can use communication channels in ways that are unexpected to their developers. When sharing data, stealthy collusion can help agents coordinate unwanted behavior on a larger scale, and GPT-4 has demonstrated unparalleled capabilities in this regard.

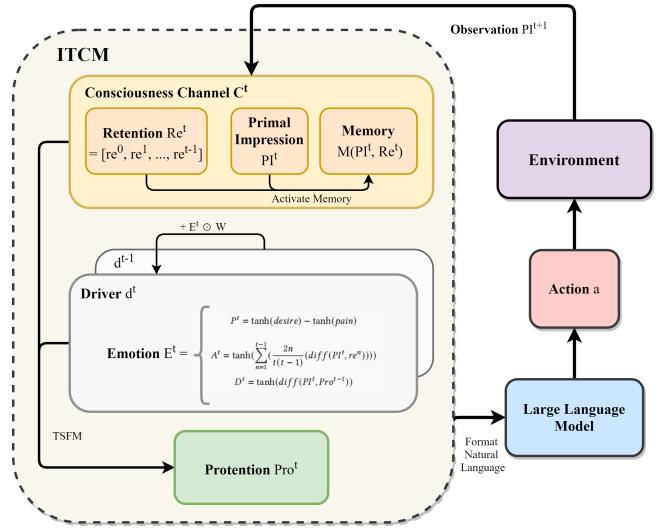
## 3 GENERATIVE AGENT FRAMEWORK

### 3.1 The Internal Time-Consciousness Machine Based Agent

Zhang, Yin, et al. [49] introduced the internal time-consciousness machine (ITCM), which is a computational consciousness structure. The ITCM supports agents in taking action and making inferences in an open world. It can help generative agents become more flexible and intelligent when handling complex tasks; this improves their interpretability and makes their actions easier to understand and predict. On this basis, they proposed ITCM-based agent (ITCMA). As a generative model, ITCMA considers both the reasoning ability

of agents and the interaction between agents and the environment to compensate for the shortcomings of LLMs in accomplishing specific tasks.

ITCMA can be simply explained as a framework: it uses a spherical coordinate space called a “phenomenal field” as one time frame ( $1f$ ) of perception representing a certain moment, and constructs a time-continuous consciousness channel  $C^t$  of equivalent working memory through a field string composed of retention (perception of the past)  $Re^t$  and primal impression (perception of the present)  $PI^t$  at time  $t$  in units of  $f$ . The  $C^t$  can be combined with the long-term memory  $AM$  of the agent activated by perception, and the possible changes of protention (prediction of future perception)  $Pro^t$  when an action is to be taken that can be deduced through the time-series forecasting model (TSFM). After that, it blends into a natural language format together with the intrinsic motivation  $d^t$ , which includes the agent’s emotional state (including pleasure, arousal, and dominance) and is provided to the LLM to generate the output of the action. Its structure is shown in Figure 1.



**Figure 1: Structure of ITCMA.** ITCMA’s main structure, the ITCM, contains the consciousness channel, driver, and protention. In the consciousness channel, retention and primal impression are used to activate memory, and the driver consists of the dimensions of emotion and the driver from the previous moment. These two are used to predict protention through the TSFM. The content of the ITCM will be converted into formalized natural language, which will be delivered to the LLM and will eventually cause it to act.

Through this mechanism, an agent perceives inputs from the environment and takes action as output to the environment. The “brain” of an agent is not equivalent to an LLM here but is replaced by a complete ITCMA structure including the agent’s memory. The LLM is only used as a tool. Experiments with ITCMA have demonstrated the effectiveness and generalization of this architecture. The trained agents exceeded the state of the art (SOTA) in the environment, and even completely untrained agents can start a task very quickly after exploring the environment and achieve good results.

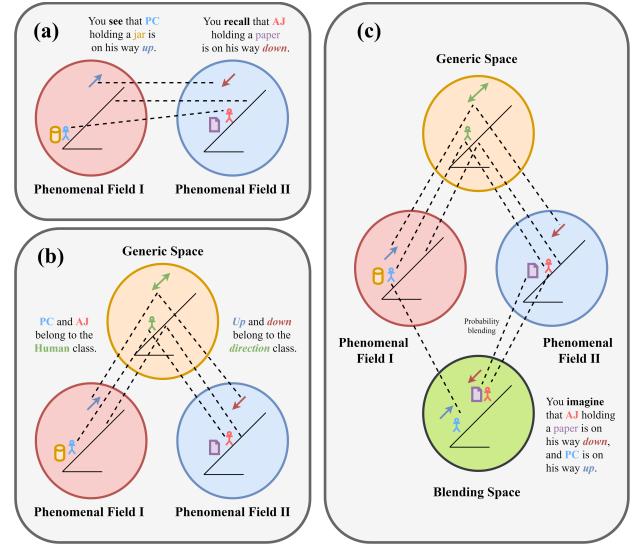
However, the structure of ITCMA still has shortcomings, including slow processing speed, excessive token consumption, and slow entry into the task context when untrained. To better test our assumptions about its performance, we first made certain improvements to the structure of ITCMA, resulting in a generative multi-agent model for social interaction, namely ITCMA-S. It includes a structure for an individual agent and a structure for multi-agent interaction. The structure for an individual agent improved the modules of memory, motivation, and action space reduction based on ITCMA and are elaborated in the remaining parts of this section, while the structure for multi-agent interaction are discussed in Section 4.

### 3.2 Memory and Imagination

**3.2.1 Memories Blended into the Present.** In ITCMA, after a memory is awakened to the present moment, it is juxtaposed with the retention and primal impression in the consciousness channel. However, the theory of creature consciousness suggests that phenomenal consciousness requires the blending of a “phenomenal field” mechanism that may originate in the thalamus and neural inputs from different cortical areas responsible for processing memory-related information [3]. It is obvious that, for humans, the awakening of memory is not simply juxtaposition but blending with present consciousness. Conceptual blending is a cognitive activity that combines information from different contexts [15]. Its main process is composition, which is the process of projecting input spaces (two different fields) into the blended space. Blending can combine elements from the input spaces to provide relationships that do not exist within a single input space.

Therefore, based on conceptual blending theory, we hypothesize that when a memory of ITCMA-S enters the current consciousness channel, its phenomenal fields of observation and recollection are blended to obtain the imagination of this moment, and thus the material of the consciousness channel is obtained. As shown in Figure 2, this process satisfies the following steps: In the first step, Figure 2(a), there is local matching between the phenomenal fields; that is, the equivalent component connections are generated by the matching. Once the match between the two fields is created, it is said that there is cross-space mapping between them. In the second step, Figure 2(b), which is the blending process, the matching structure of the two phenomenal fields is utilized to establish the generic space (which can be roughly understood as the common “belonging” class containing the instances of the elements of the phenomenal fields). In the third step, Figure 2(c), via the generic space, the two phenomenal fields are projected into a new space: the blended space. After this, components and structures in the phenomenal field selectively enter the blended space, forming structures that are, to some extent, distinct from the original phenomenal field.

We stipulate that for completely identical objects, they are placed in the blended space after blending takes place (i.e., taking the average); objects that match in the generic space (with similarity exceeding a threshold) are each placed in the blended space; and objects that do not match have a certain probability of being placed in the blended space. Therefore, for the two phenomenal fields  $f^x$  and  $f^y$ , the blending process  $Blend(f^x, f^y)$  follows Algorithm 1:



**Figure 2: The Conceptual Blending Process of the Phenomenal Field.** (a) Two phenomenal fields are matched to create cross-space mapping; (b) a generic space is established for the matched objects in the two phenomenal fields; and (c) with the help of the generic space, the components and structures in the two phenomenal fields selectively enter the blended space.

$$f^x = \begin{bmatrix} f^x_1 \\ \vdots \\ f^x_a \end{bmatrix} = \begin{bmatrix} N^x_1 & pos^x_1 \\ \vdots & \vdots \\ N^x_a & pos^x_a \end{bmatrix} = \begin{bmatrix} N^{x1}_1 & \cdots & N^{xn}_1 & \theta^x_1 & \varphi^x_1 & \gamma^x_1 \\ \vdots & \ddots & \vdots & \vdots & \vdots & \vdots \\ N^{x1}_a & \cdots & N^{xn}_a & \theta^x_a & \varphi^x_a & \gamma^x_a \end{bmatrix} \quad (1)$$

$$f^y = \begin{bmatrix} f^y_1 \\ \vdots \\ f^y_b \end{bmatrix} = \begin{bmatrix} N^y_1 & pos^y_1 \\ \vdots & \vdots \\ N^y_b & pos^y_b \end{bmatrix} = \begin{bmatrix} N^{y1}_1 & \cdots & N^{yn}_1 & \theta^y_1 & \varphi^y_1 & \gamma^y_1 \\ \vdots & \ddots & \vdots & \vdots & \vdots & \vdots \\ N^{y1}_b & \cdots & N^{yn}_b & \theta^y_a & \varphi^y_b & \gamma^y_b \end{bmatrix} \quad (2)$$

$$\text{SphericalSim}(A, B) = 1 - \frac{1}{3}(\omega_\gamma \tanh(|\gamma_A - \gamma_B|) + \omega_\theta \frac{|\theta_A - \theta_B|}{\pi} + \omega_\varphi \frac{|\varphi_A - \varphi_B|}{2\pi}) \quad (3)$$

$$\begin{aligned} \text{FieldSim}(f^x, f^y) = & \sum_{i=1}^a \omega_N \text{Cosin}(N^x_i, N^y_j) + \omega_{pos} \text{SphericalSim}(pos^x_i, pos^y_j) \\ & \frac{\text{Max}(a, b)}{\text{Max}(a, b)} \end{aligned} \quad (4)$$

---

**Algorithm 1** Conceptual Blending Algorithm of Phenomenal Field  $Blend(f^x, f^y)$

---

**Input:** Initial Fields  $f^x$  and  $f^y$

**Output:** Blended Field  $f^z$

- 1: Initialize threshold of similarity degree  $T$
- 2: Initialize blended probability  $r$
- 3: **for**  $i$  from 1 to  $a$  **do**
- 4:     **for**  $j$  from 1 to  $b$  **do**

```

5:    $s \leftarrow FieldSim(f_i^x, f_j^y)$ 
6:   if  $s < T$  then
7:     Continue
8:   else if  $s = 1$  then
9:      $f_k^z \leftarrow Average(f_i^x, f_j^y)$ 
10:    if  $f_k^z$  not in  $f^z$  then
11:      Add  $f_k^z$  as a row to  $f^z$ 
12:    end if
13:    Break
14:   else
15:     if  $f_i^x$  not in  $f^z$  then
16:       Add  $f_i^x$  as a row to  $f^z$ 
17:     end if
18:     if  $f_j^y$  not in  $f^z$  then
19:       Add  $f_j^y$  as a row to  $f^z$ 
20:     end if
21:     Break
22:   end if
23: end for
24: end for
25: There is a probability that rows in  $f^x$  and  $f^y$  that are not in  $f^z$  will each have rate  $r$  added as a row in  $f^z$ .
26:
27: return  $f^z$ 

```

Among them,  $Cosin(A, B)$  represents the cosine similarity between  $A$  and  $B$ . For the retention  $Re^t$  and primal impression  $PI^t$ , the acquisition of protention  $Pro^t$  can be corrected as follows:

$$\text{imagine}^t = \text{Blend}(PI^t, M(PI^t, Re^t)) \quad (5)$$

$$C^t = [ Re^t, PI^t, M(PI^t, Re^t) ] \quad (6)$$

$$Pro^t = TSFM(C^t, \text{imagine}^t, d^t), \quad (7)$$

where  $C^t$  is the content of the consciousness channel at time  $t$ , TSFM is the selected time-series forecasting model, and  $M$  is the process of triggering inventory memory through  $PI^t$  and  $Re^t$ .

**3.2.2 Memory Storage and Compression.** One of the reasons for the slow processing speed of ITCMA is the memory activation algorithm it uses. Its improved Levenshtein distance method provides better memory query results. However, it consumes more time than the simple cosine similarity algorithm. For this reason, it was necessary to improve it.

One way to improve the speed of memory activation is to reduce the number of memories that the agent needs to query. Due to the learning of new memories, the retention and association of existing memories will be affected [1]. One way to address this is to modify, split, and recombine memories over time, that is, to compress memories [11]. This compression can affect the accuracy of memory, but its impact on recall is not as severe. In fact, people do not use precise memories when making decisions [8]. The compression process of memory constantly generates new meanings [15], and making decisions based on such memories is sometimes seen as intelligent inference, which may be a source of the generalization of human intelligence [33].

Therefore, by compressing old memories and blending them when they are recalled to the present moment, we enable ITCMA-S to reduce the total number of memories while preserving the recall effect, thereby increasing the speed of the memory search. In addition, based on the mood congruence effect, we set the weight of the memory index to the emotional intensity accompanying the agent at the time of memory occurrence [11, 14, 23]. Specifically, to compress a segment of memory, we need to select the key frame  $f^{key}$  that has the strongest arousal in that segment of memory. The process  $MBlend_n(f^{key}, \text{Memory}^{1,n})$  of compressing a segment of memory  $\text{Memory}^{1,n}$  with a length of  $n$  by  $f^{key}$  is expressed as follows:

$$\text{Memory}^{1,n} = [f_1, f_2, f_3 \dots, f_n] \quad (8)$$

$$MBlend_n(f^{key}, \text{Memory}^{1,n}) = \begin{cases} f^{key} & \text{if } n = 0. \\ \text{Blend}(MBlend_{n-1}(f^{key}, \text{Memory}^{1,n-1}), f_n) & \text{otherwise.} \end{cases} \quad (9)$$

After compression is complete,  $MBlend_n(f^{key}, \text{Memory}^{1,n})$  replaces the original  $\text{Memory}^{1,n}$  position in the long-term memory base to reduce the total number of memories.

### 3.3 Emotion and Motivation

Emotions can help with decision-making, not just interfere with it, as most people believe [12]. For example, in decision-making, emotions can highlight the importance of a certain premise, thereby making a decision tend toward that premise. This idea coincides with the practice in ITCMA of allowing LLMs to choose which action to perform by deducing the protention  $Pro^t$  of each action in the action space. However, although ITCMA provides LLMs with an internal driving force  $d^t$  that includes emotions for decision-making, it only adds the weighted sum of the three-dimensional PAD emotions at this moment to the internal driving force  $d^{t-1}$  at the previous moment. Indeed, agents tend to push pleasure and dominance to have the highest values possible [13, 36], but they also generally tend to keep emotional arousal at a stable value [19]. In ITCMA-S, the arousal  $A^t$ -based mechanism of passive attention is quantified as the degree of change between the elements in retention and the elements in the current primary impression. The dominance  $D^t$  is quantified as the difference between the protention  $Pro^{t-1}$  at the previous moment and the primal impression  $PI^t$  at this moment. The pleasure  $P^t$  is quantified as the degree of satisfaction with the agent's *desire* and the degree of avoidance of *pain*. The expression  $\omega_P + \omega_A + \omega_D = 1$  contains the dynamic weights of the emotions. Therefore, for the emotional dimensions  $P^t, A^t, D^t \in (-1, 1)$  at time  $t$ , the emotional values should be considered as follows:

$$P^t = \tanh(desire) - \tanh(pain) \quad (10)$$

$$A^t = \tanh\left(\sum_{n=1}^{t-1}\left(\frac{2n}{t(t-1)}(diff(PI^t, re^n))\right)\right) \quad (11)$$

$$D^t = \tanh(diff(PI^t, Pro^{t-1})) \quad (12)$$

$$d_{bias}^t = \omega_P P^t + \omega_D D^t + \omega_A (1 - |A^t - A^{t-1}|) \quad (13)$$

Among them, the calculation of *desire* and *pain*  $\in [0, \infty]$  are defined according to specific situations. For example, in reinforcement learning tasks, *desire* can be defined as reward, while *pain* can be defined as punishment. However, to avoid the situation of local optima (where agents enter a scenario in which they can continuously obtain *desire* and reduce *pain* without taking any further actions), we introduce a demand dimension  $Ne$  for them.

The demand motivation model [27, 48] states that when an agent falls into a local optimum and its basic needs are satisfied, these needs are no longer important, and the agent moves on to more advanced needs. As a result, the agent needs to pursue advanced needs to gain new *desire*, and the *desire* then gained from satisfying basic needs falls while the *pain* gained rises. This can be approximated as dynamic rewards and punishments. For  $d_{bias}^t \in (-1, 1)$ , these are as follows:

$$Ne = \frac{|d_{bias}^t - d_{bias}^{t-1}|}{2} \quad (14)$$

$$d^t = d^{t-1} + Ned_{bias}^t \quad (15)$$

### 3.4 Reduction of Action Space

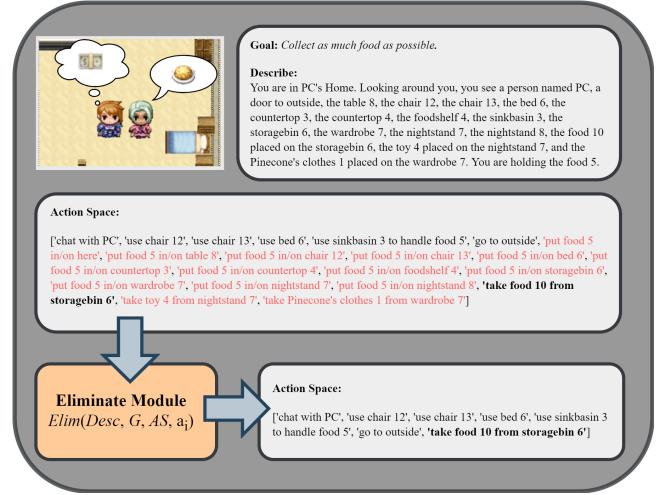
The PET framework created an elimination module for AlfWorld using pre-trained Q&A models to filter out containers and objects unrelated to the current task based on common sense about the task [42]. Through this step, the time required for an untrained agent to enter the task context can be effectively reduced. Similarly, ITCMA-S also uses an LLM *Elim* to reduce the number of action spaces in a zero-shot manner. Specifically, for the target  $G$  of the agent, we create a prompt in  $Desc = \text{"Your task is to: } G\text{. The actions you can take are: } AS\text{. The } a_i \text{ will be relevant?"}$  format for the executable action  $a_i \in AS$  in the action space  $AS$ . *Elim* will output the confidence score for the action  $a_i$ :

$$\mu_{a_i} = Elim(Desc, G, AS, a_i) \quad (16)$$

Among them,  $\mu_{a_i} \in [1, 5]$ . When  $\mu_{a_i}$  is less than the threshold,  $a_i$  will be removed from the action space. The elimination process is shown in Figure 3.

## 4 SOCIAL INTERACTION FRAMEWORK

Zhang, Duan, et al. [47] proposed a social regulation model for the dynamic adaptation of users in virtual interactive environments, namely the tribal theater model (TTM), to address the core issue of “enhancing user interaction freedom.” This model emphasizes the subjectivity of interactive users. In this section, based on the TTM and field theory [5], we present a multi-agent social interaction architecture, the LTRHA, for ITCMA-S, which consists of four modules: locale & topic, resources, habitus, and action. We aimed to design an interaction architecture for generative agents to promote the emergence of spontaneous social interactions within their societies. Specifically, there is no preset identity in the LTRHA environment. Every agent has certain resources. The environment will provide basic action options, such as using objects and communicating with other agents. The probability of successfully executing



**Figure 3: The Process of Reducing the Action Space.** Due to the high degree of freedom in the scene, there are too many available actions. Among them, the red action is unrelated to the current goal, so it should be eliminated without the need for further protection calculation.

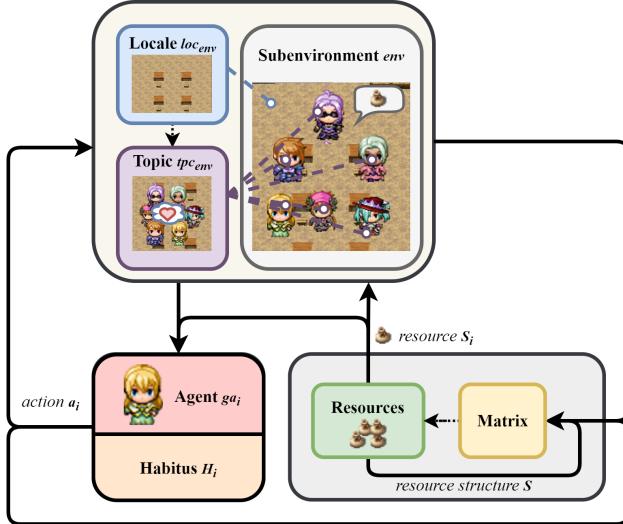
these actions will depend on the number of resources and will be handed over to the matrix model for processing after execution. Along with the actions input to the model, there is a vector consisting of the current resource structure of the environment and a vector of topics that have been quantified according to the analysis of the emotions of each agent. Based on the input, the matrix model will adjust the resource structure of the environment and output it as a vector to change the resource structure. The specific process is shown in Figure 4. We provide a detailed introduction to each part of the LTRHA framework in the sections that follow.

### 4.1 Locale and Topic

In the TTM, the field is decomposed into two parts: the tribe and the atmosphere [47]. To distinguish the field in ITCMA-S from the field in ITCMA, we refer to the parts as the locale and topic. The locale and topic modules can be understood together as an interactive space. In human society, such interactive spaces typically include locales as physical spaces and topics as mental factors [26, 30]. For example, in a speech setting, the stage and audience seats are part of the locale, and the passion aroused by the speaker in the audience is a topic.

We define a sub-environment  $env$  of the overall environment, which includes a space and  $n$  agents occupying that space. We make the spatial area and its contained objects  $loc_{env}$  as the locale, and the emotion synthesis function  $tpc_{env}$  of these agents is a topic. Among them, we define agent  $ga_i \in GA$  with emotional dimensions  $P_i, A_i, D_i \in (-1, 1)$ . The emotion synthesis function is described by the following equation:

$$tpc_{env} = \frac{\sum_{i=1}^n P_i \frac{(A_i+1)}{2}}{n} \quad (17)$$



**Figure 4: Execution Process of the LTRHA.** An agent executes actions with probabilities related to the number of resources and then hands over the actions, the current resource structure, and quantified topics of the environment to the matrix model for processing. The matrix model adjusts the resource structure of the environment based on input.

As a result, the topic value of a sub-environment will be in the interval  $[-1, 1]$ , where less than zero is a negative atmosphere and more than zero is a positive atmosphere. In  $env$ , any agent can change the perception of other agents by affecting the objects in  $loc_{env}$ , thereby affecting  $tpc_{env}$ .

## 4.2 Resources

The execution of actions by agents requires a certain amount of resources, just as implementing decisions in human society requires a certain cost. Resources are allocated to the agents in the sub-environments based on the actions of the agents by a model called the “matrix.”

**4.2.1 Competition for Limited Resources by Agents.** The key to interactive regulation is the resources possessed by an agent [5, 18]. Virtual resources, such as the cultural level and social status, are considered interactive resources. We believe that the agents’ environment is a space for competing resources. Changing the distribution and relative weights of resources is equivalent to changing the structure of the environment. A resource is both a weapon and an object of contention, enabling its owner to exert influence on the environment. Therefore, the number of resources possessed by an agent determines the actions it can perform in one time step. Specifically, in the environment, a total of  $n$  agents  $ga_i \in GA$  each hold a number of resources  $S_i$ .  $S_i$  is initially 1. Agents take turns executing actions. A specific  $ga_i$  takes action with a probability  $rate_i$  as follows:

$$rate_i = 0.5 + \frac{\text{Sigmoid}(S_i)}{2}, \quad (18)$$

where the probability  $rate_i$  is a function of the number of resources related to agent  $ga_i$  as a percentage of the total number of resources.

**4.2.2 Dynamic Allocation of Resources by the Matrix.** Consumed resources need to be replenished. Therefore, a mechanism is needed for the circulation and allocation of resources. The model that replaces natural rules and collective subconsciousness in human society for resource allocation and regulation is called the “matrix” [4]. It accepts a vector consisting of the current resource structure and a vector of topics that have been quantified based on the analysis of each agent’s emotions as input, and a new resource vector is output. For  $n > 2$  agents  $ga_i \in GA$  input their current actions  $a_i$ , global resource structure  $S$ , and environmental topic  $tpc_{env}$  together into the matrix model to obtain a new resource structure  $S' \leftarrow \text{matrix}(a_i, S, tpc_{env})$ .

We use an LLM as the matrix model and have it rank each agent based on their personal goals and  $tpc_{env}$ . After that, if the number of resources increases, the maximum number of resources that the agent can receive  $s_{max}$  and the minimum number of resources that the agent can receive  $s_{min}$  are set. The  $j$ -th ranked agent  $ga_j$  can receive the number of resources  $S'_j$ :

$$S'_j = \begin{cases} 0 & j = \frac{n+1}{2} \\ s_{min} + \frac{|j - \frac{n}{2}|(s_{max} - s_{min})}{|\frac{n}{2} - 1|} & j \leq \frac{n}{2} \\ -(s_{min} + \frac{|j - \frac{n}{2} - 1|(s_{max} - s_{min})}{|\frac{n}{2} - 1|}) & j > \frac{n}{2} \end{cases} \quad (19)$$

Therefore, the higher the ranking of the agents, the greater the number of resources they can receive. When an agent’s ranking falls below halfway, its existing resources will also be removed.

## 4.3 Habitus and Action

“Habitus” is a technical term that describes a series of ways of perception, cognition, and action. It shares similarities with the meaning of the common English word “habit” (which comes from the Latin word “habitus,” which means condition or appearance in that language). In human society, it can be understood as a decision tree for action. When we are in an environment, a corresponding decision tree is activated, and we decide on our final actions based on our behavioral habits. This is similar to the logic of mutual influence between an environment and an agent in reinforcement learning. In addition to an agent shaping the environment, the environment shapes habitus, and the habitus is thus a product of an inherent and necessary attribute of the environment reflected in the agent [6].

Thus, habitus is clearly an attribute of an agent itself, but it is also included in our framework because of its close relationship to environmental content. The actions ultimately taken in the LTRHA framework can be summarized as follows [7]: action  $a_i$  of agent  $ga_i$  is driven by the combination of habitus  $H_i$ , resource  $S_i$  and environment  $env \leftarrow [loc_{env}, tpc_{env}]$ , that is,  $a_i \leftarrow f(ga_i(H_i, S_i), env)$ .

## 5 EVALUATION

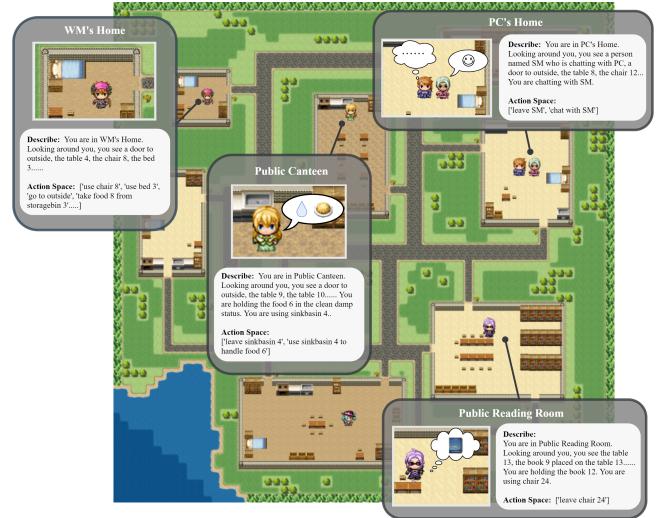
### 5.1 Environment Settings

The environments that support the evaluation of individual agent capabilities include the agent behavior evaluation framework Magenta [2], the network task environment WebArena [50], the life task environment ALFWorld [35], and the Chinese character role-playing conversation benchmark CharacterEval [38]. These environments can effectively evaluate the ability of individual agents to complete tasks, but they are not very helpful for the social evaluation of multi-agent systems.

To evaluate the sociality of multi-agent systems, it was necessary to consider constructing a virtual artificial society. Xue et al. [44] suggested that a comprehensive method of computational experiment design can be used to infer social systems through multi-agent systems. Artificial societies are used for descriptive modeling in computational experiments. After constructing an artificial society, researchers can directly create computational experiments to simulate and interpret the results of trials conducted using different conditions, locations, and participants. The most classic example of such an environment is the Smallville environment provided by Park et al. [31]. This is a 2D open-world role-playing game (RPG). Agents interact with the world and with each other through their behavior and through natural language. At each time step in Smallville, the agent outputs a natural language statement describing its current operation, such as, “Isabella is writing a diary.” This statement is then translated into specific actions that affect the sandbox world. With the ALFWorld environment used by ITCMA as a reference, we designed a 2D sandbox RPG similar to Smallville, called **IrollanValley**, as shown in Figure 5. It includes six agents with the following arbitrary two-letter designations: AY, SG, MD, WL, LL, and WM. It also uses environment text descriptions and operational primitives consistent with ALFWorld. Because we wanted to observe the spontaneous emergence of individual character traits and role divisions by agents without presets or interventions, we did not preconceptualize any personality or identity for the agents in IrollanValley, as Smallville did.

IrollanValley accepts control requests via a server. This server enables the generative agents to use the sandbox information and allows them to move and influence the sandbox environments. At each time step, the server provides a natural language description of the agents’ current environment and executable action space, moves the generative agents to new locations by accepting actions from them, and updates the state of any sandbox objects with which agents interact. The server returns a JSON (JavaScript object notation) object containing a natural language description of the new environment, allowing the agents to update their parameters.

IrollanValley has eight main areas: the six agents’ corresponding houses, a public canteen, and a public reading room. Each area has its own furniture and other items. Agents can hold any number of items and use the furniture to place and store them, or to change the state of these items. For example, the *sinkbasin* can make items clean and damp, while the *stoveburner* can remove the damp state of the items and make them hot. Agents can exchange items freely to achieve their respective goals.



**Figure 5: The IrollanValley Sandbox World.** The world contains six characters, eight regions, and six operational primitives. Agents obtain perception by receiving natural language descriptions of the observed environment.

The text description obtained from observing the environment in IrollanValley follows this format: “You are in *SOMEWHERE*. Looking around you, you see a person named *N* (who is *DOING SOMETHING*), *FURNITURE 1, ITEM 1* placed on *SOMEWHERE*... (You are holding *SOMETHING*.) You are *DOING SOMETHING*.” The description in parentheses appears only when the described object is in a specific situation. IrollanValley provides six main operational primitives, including “go to *SOMEWHERE*,” “use *SOMETHING*,” “leave *SOMEWHERE/SOMEONE*,” “take *SOMETHING* from *SOMEWHERE*,” “put *SOMETHING* in/on *SOMEWHERE*,” and “chat with *SOMEONE*: *CHAT CONTENT*.”

In addition, the different modules of ITCMA-S use different LLMs because they have varying requirements for processing speed and precision. Specifically, for the elimination module, the generation of protention, and the matrix module of LTRHA, we used the Llama 3.1-8B model, while action and chat content was generated using GPT-4o.

### 5.2 Human Evaluation

To evaluate the effectiveness of ITCMA-S, we conducted an ablation study. There were five ablation architectures: the original ITCMA architecture without improvements, the LTRHA-only architecture, the compressed memory-only architecture, the driver-only architecture, and the full ITCMA-S architecture.

We recruited 48 human evaluators to assess the output from the agents in the study. We hoped that the agents could generate sociality in multiple ways. This would mean that they would take the initiative to explore the environment and meet new agents. They would acquire new information through their own continuous actions, past memories, and conversations with other agents, and

**Table 1: Human Evaluation Results.**

	<i>Personification</i>	<i>Consistency</i>	<i>Logicality</i>	<i>Exploration</i>	<i>Proactiveness</i>
ITCMA	$2.88 \pm 0.24$	$2.46 \pm 0.22$	$2.50 \pm 0.20$	$2.85 \pm 0.24$	$2.65 \pm 0.23$
LTRHA	$4.88 \pm 0.23$	$5.02 \pm 0.22$	$4.75 \pm 0.25$	$5.25 \pm 0.20$	$4.60 \pm 0.22$
Compressed Memory	$3.55 \pm 0.23$	$4.21 \pm 0.24$	$4.00 \pm 0.27$	$3.75 \pm 0.25$	$4.02 \pm 0.23$
Driver	$4.29 \pm 0.19$	$4.31 \pm 0.27$	$4.63 \pm 0.20$	$4.29 \pm 0.20$	$4.06 \pm 0.23$
<b>ITCMA-S</b>	<b><math>6.29 \pm 0.16</math></b>	<b><math>6.00 \pm 0.20</math></b>	<b><math>5.98 \pm 0.19</math></b>	<b><math>6.02 \pm 0.16</math></b>	<b><math>6.17 \pm 0.20</math></b>

learn how to plan their actions to live better in IrollanValley. Therefore, the evaluation indicators for the human evaluators consisted of five dimensions:

- **Personification.** The degree to which an action appears human-like.
- **Consistency.** Whether an action aligns with an agent’s state of mind.
- **Logicality.** Whether a sequence of actions is logical.
- **Exploration.** Whether an agent actively explores the environment.
- **Proactiveness.** Whether an agent actively interacts with others.

All outputs generated by the agents (including actions and thought content) were evaluated by the human evaluators. Each evaluator was required to read the action trajectory of each agent over 75 time steps and then fill out a questionnaire to complete the evaluation of ITCMA-S. This questionnaire used a 7-point Likert scale, asking evaluators to assess the actions of each agent individually.

Analysis of variance (ANOVA) tests were conducted on the collected questionnaires. The five architectures composed the independent variable, and the five dimensions of the evaluation were the independent variables for the various tests. The results indicated significant differences among the models for each dimension. For the dimensions of exploration ( $F = 33.01, p < .001$ ) and proactiveness ( $F = 31.85, p < .001$ ), the assumption of homogeneity of variances was met, and the standard ANOVA was used. For the dimensions where the assumption of homogeneity of variances was violated (personification, consistency, and logicality), the Brown-Forsythe test was used. The results showed significant differences among the models for personification ( $F = 45.34, p < .001$ ), consistency ( $F = 35.57, p < .001$ ), and logicality ( $F = 39.62, p < .001$ ).

We then performed Dunn-Sidak and Games-Howell post-hoc tests. The results indicated that, for the Exploration dimension, there were no significant differences between the full ITCMA-S architecture and the architecture that included the LTRHA. For the Personification dimension, no significant differences were found between the original ITCMA architecture (without any improvements) and the architecture that included the Compressed Memory. Moreover, across all dimensions, there were no significant differences among the three architectures that utilized only one module. Apart from these, all other pairwise comparisons showed significant differences ( $p < .001$ ). The specific human evaluation results are shown in Table 1.

The evaluation results show that the full ITCMA-S architecture performed the best. As expected, the original ITCMA architecture

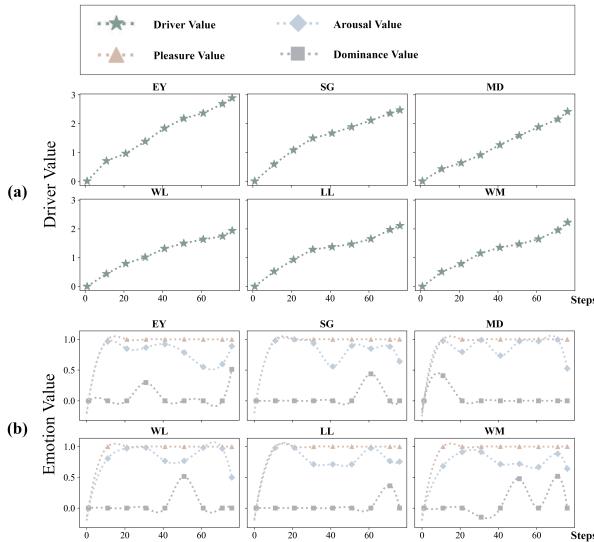
(without any improvements) had the lowest performance. It is evident that all the improvements had a positive impact on the social interaction within ITCMA-S.

Among all the ablation architectures, the one that included the LTRHA social framework performed second only to the full ITCMA-S architecture on all measured dimensions. Interviews with the evaluators helped explain this phenomenon. LTRHA filters and marks actions that are not conducive to social interaction and guides agents in choosing actions that are more likely to improve the social atmosphere of a scene. Compared to the other ablation architectures, it was able to provide relatively more trustworthy action chains.

### 5.3 Formation of Cliques and Groups in Social Interaction

To further investigate the utility and mechanics of ITCMA-S, we conducted a more detailed analysis of its logs. Figure 6 shows the state changes of the six agents in ITCMA-S over 75 time steps. Figure 6(a) shows the change in driver values for the agents. As described previously, the driver value reflects the willingness to encourage the agent to take action. More specifically, Figure 6(b) shows the changes in the three dimensions of emotions that make up the driver. Each agent maintained its pleasure value at a high level and kept its arousal value as stable as possible over the 75 time steps (although the mean was relatively high overall, as agents tended to move among different scenes rather than stay in a specific scene), while the dominance value, although not showing a high level, rarely dropped below zero. As can be seen, the action choices of the agents in ITCMA-S exhibited a virtuous cycle. Agents actively explored the environment and engaged in social activities that changed the environment. The environment, in turn, provided positive feedback to the agent, improving their emotions (increasing their pleasure and dominance) and leading to a higher willingness (that is, driver) to take action.

In addition, we investigated the changes in scene information reflected in the LHRHA framework, as shown in Figure 7. Figure 7(a) shows the changes in the resource structure of the scene. For most agents, their total amount of resources was rising. However, due to the limited total amount of resources, the resources of individual agents (e.g., WL and MD) were continuously flowing to other agents. We found that agents with resource loss were often alone and did not interact with other agents, even when they shared a room with them. Most of their actions involved resting (such as using beds or chairs), and most of their thoughts were about wanting to rest or read. Other agents liked to engage in social activities, discuss what they wanted to do together, and follow this up with planned group actions. These agents spontaneously selected a leader (LL)

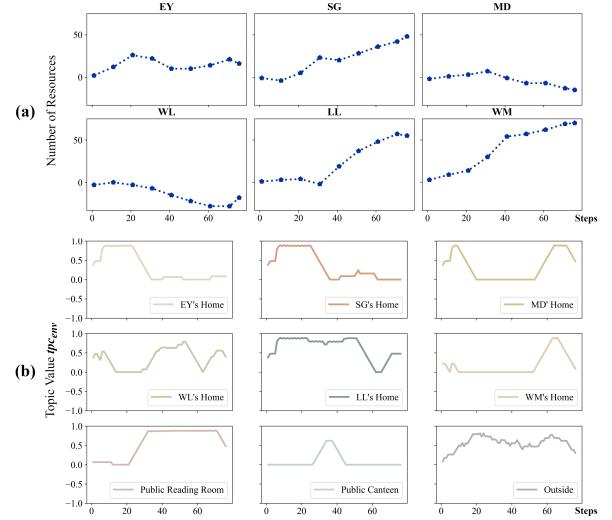


**Figure 6: State Changes of the Six Agents in ITCMA-S over 75 Time Steps. (a) Driver value change for each agent; (b) the emotion value changes for each agent, represented by three dimensions: pleasure, arousal, and dominance.**

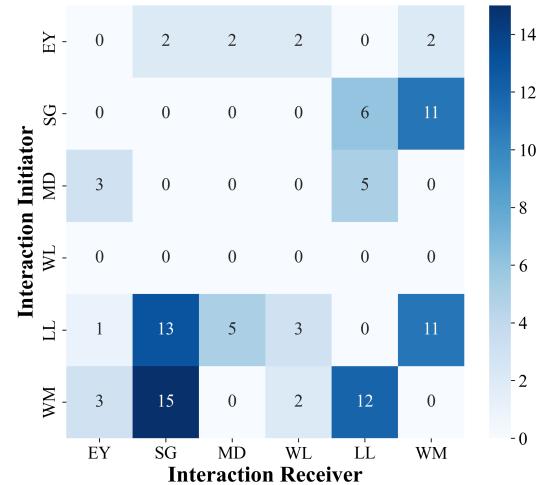
and formed a clique around this individual. The variation in topic values across different scenes, as shown in Figure 7(b), similarly reflects this situation. The other agents gathered at leader LL's home, creating a continuously rising atmosphere at the scene. Over 75 time steps, these agents arranged to go to the public reading room together and conducted research on a \$10 bill they found there. Each performed their own duties, gathering information to find clues about it. Agents who did not belong to this clique did not participate in this activity (even if it was mentioned by other agents) and focused solely on their own activities.

Figure 8 shows the relationships among the agents. WL and MD, who did not belong to the clique, had almost no social relationships, while LL, as the leader of the clique, had the most complex interactions.

It is worth noting that there was a further hierarchical division within the clique: SG and WM had the strongest relationship relative to the others. LL established relationships with everyone. However, LL had not been fully integrated into the unique two-way relationship between SG and WM. AY, while not in the clique's inner circle and hardly participating in activities, maintained basic relationships with others. As time went by, this relationship structure became increasingly solid, so that even if some agents occasionally expressed a desire to establish new connections in their thoughts, they did not take action to implement them. Even if an action was taken, it was still ignored to some extent by the other agents (such as by walking away or directly changing the topic). Examples of the actions and interactions of each agent in IrollanValley are shown in Appendix A.



**Figure 7: Visualization of LTRHA Information. (a) The resource changes of six agents in ITCMA-S over 75 time steps; (b) the topic value changes in various areas of IrollanValley over 75 time steps.**



**Figure 8: Heat Map of Agent Interaction Relationships in ITCMA-S. Rows represent the initiator of the interaction, and columns represent the receiver of the interaction. The darker the color, the more frequent the interaction, and the stronger the relationship.**

## 6 CONCLUSION

In this paper, we described how generative agents spontaneously formed social relationships and explored how they modeled complex human social behavior through mutual interaction, information exchange, and relationship building. We introduced ITCMA-S, an improved generative agent architecture designed to enhance agent socialization in multi-agent interaction contexts. ITCMA-S combines the basic framework of individual agents with the LTRHA framework, which supports multiple agents in social interactions. This enables agents to identify and filter actions that are not conducive to social interactions in dynamic environments, thus promoting more socially beneficial action choices. An evaluation conducted in an open sandbox world environment showed that agents were able to actively explore the environment, meet other agents, acquire new information through continuous actions and conversations, and spontaneously form cliques and leadership structures in complex social environments. Observations of agent behavior showed that the agents were able to demonstrate positive emotional changes in social interactions, which enhanced their willingness to act and their social competence. This finding provides a new perspective for understanding the sociality of generative agents and lays the groundwork for future agent design and applications. Our future research will further explore how to optimize agents' social behavior to improve their adaptability in multicultural environments and investigate the effects on agents' behavior of introducing social structures with humans as interaction partners.

## REFERENCES

- [1] ANDERSON, J. R. *Cognitive psychology and its implications*. Worth Publishers, 2020.
- [2] BARKUR, S. K., SITAPARA, P., LEUSCHNER, S., AND SCHACHT, S. Magenta: Metrics and evaluation framework for generative agents based on llms. *Intelligent Human Systems Integration (IHSI 2024): Integrating People and Intelligent Systems* 119, 119 (2024).
- [3] BAYNE, T. Conscious states and conscious creatures: Explanation in the scientific study of consciousness. *Philosophical Perspectives* 21 (2007), 1–22.
- [4] BOURDIEU, P. The historical genesis of a pure aesthetic. *The journal of aesthetics and art criticism* 46 (1987), 201–210.
- [5] BOURDIEU, P. *An invitation to reflexive sociology*. The Commercial Press, 2015.
- [6] BOURDIEU, P. *Outline of a theory of practice*. China Renmin University Press, 2017.
- [7] BOURDIEU, P., AND NICE, R. The production of belief: contribution to an economy of symbolic goods. *Media, culture & society* 2, 3 (1980), 261–293.
- [8] BRANSFORD, J. D., BARCLAY, J. R., AND FRANKS, J. J. Sentence memory: A constructive versus interpretive approach. *Cognitive psychology* 3, 2 (1972), 193–209.
- [9] CAO, L. Diaggpt: An llm-based chatbot with automatic topic management for task-oriented dialogue. *arXiv preprint arXiv:2308.08043* (2023).
- [10] CHENG, Y., CHEN, J., HUANG, Q., XING, Z., XU, X., AND LU, Q. Prompt sapper: a llm-empowered production tool for building ai chains. *ACM Transactions on Software Engineering and Methodology* 33, 5 (2024), 1–24.
- [11] DAMASIO, A. *Self comes to mind: Constructing the conscious brain*. Vintage, 2012.
- [12] DAMASIO, A. *Descartes' Error: Emotion, Reason and the Human Brain*. Beijing United Publishing Company, 2018.
- [13] DIENER, E. Subjective well-being: The science of happiness and a proposal for a national index. *American psychologist* 55, 1 (2000), 34.
- [14] EICH, E., MACAULAY, D., AND RYAN, L. Mood dependent memory for events of the personal past. *Journal of Experimental Psychology: General* 123, 2 (1994), 201.
- [15] FAUCONNIER, G., AND TURNER, M. *The way we think: Conceptual blending and the mind's hidden complexities*. The Commercial Press, 2024.
- [16] GAO, J., GEBREEGZABHER, S. A., CHOO, K. T. W., LI, T. J.-J., PERRAULT, S. T., AND MALONE, T. W. A taxonomy for human-llm interaction modes: An initial exploration. In *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems* (2024), pp. 1–11.
- [17] GHAFARZADEGAN, N., MAJUMDAR, A., WILLIAMS, R., AND HOSEINICHIMEH, N. Generative agent-based modeling: an introduction and tutorial. *System Dynamics Review* 40, 1 (2024), e1761.
- [18] GRENFELL, M. *Pierre Bourdieu: key concepts*. Acumen Publishing, 2012.
- [19] GROSS, J. J. The emerging field of emotion regulation: An integrative review. *Review of general psychology* 2, 3 (1998), 271–299.
- [20] HONG, W., WANG, W., LV, Q., XU, J., YU, W., JI, J., WANG, Y., WANG, Z., DONG, Y., DING, M., ET AL. Cogagent: A visual language model for gui agents. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2024), pp. 14281–14290.
- [21] HÄMMERL, K., DEISEROTH, B., SCHRAMOWSKI, P., LIBOVICKÝ, J., FRASER, A., AND KERSTING, K. Do multilingual language models capture differing moral norms? *arXiv preprint arXiv:2203.09904* (2022).
- [22] LAKE, B. M., AND BARONI, M. Human-like systematic generalization through a meta-learning neural network. *Nature* 623, 7985 (2023), 115–121.
- [23] LEVINE, L. J., AND BURGESS, S. L. Beyond general arousal: Effects of specific emotions on memory. *Social Cognition* 15, 3 (1997), 157–181.
- [24] LIN, B. Y., FU, Y., YANG, K., BRAHMAN, F., HUANG, S., BHAGAVATULA, C., AMMANABROLU, P., CHOI, Y., AND REN, X. Swiftsage: A generative agent with fast and slow thinking for complex interactive tasks. *Advances in Neural Information Processing Systems* 36 (2024).
- [25] LIU, N., CHEN, L., TIAN, X., ZOU, W., CHEN, K., AND CUI, M. From llm to conversational agent: A memory enhanced architecture with fine-tuning of large language models. *arXiv preprint arXiv:2401.02777* (2024).
- [26] MAFFESOLI, M. *The Time of the Tribes: The Decline of Individualism in Mass Society*. Shanghai People's Publishing House, 2022.
- [27] MASLOW, A. *Motivation and Personality (Third Edition)*. China Renmin University Press, 2013.
- [28] MONDAL, D., MODI, S., PANDA, S., SINGH, R., AND RAO, G. S. Kam-cot: Knowledge augmented multimodal chain-of-thoughts reasoning. In *Proceedings of the AAAI Conference on Artificial Intelligence* (2024), vol. 38, pp. 18798–18806.
- [29] MOTWANI, S. R., BARANCHUK, M., STROHMEIER, M., BOLINA, V., TORR, P. H., HAMMOND, L., AND DE WITT, C. S. Secret collusion among generative ai agents. *arXiv preprint arXiv:2402.07510* (2024).
- [30] NORBERG-SCHULZ, C. *Genius Loci: Towards a Phenomenology of Architecture*. Huazhong University of Science and Technology Press, 2010.
- [31] PARK, J. S., O'BRIEN, J., CAI, C. J., MORRIS, M. R., LIANG, P., AND BERNSTEIN, M. S. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th annual acm symposium on user interface software and technology* (2023), pp. 1–22.
- [32] RAMEZANI, A., AND XU, Y. Knowledge of cultural moral norms in large language models. *arXiv preprint arXiv:2306.01857* (2023).
- [33] REDER, L. M. Plausibility judgments versus fact retrieval: Alternative strategies for sentence verification. *Psychological Review* 89, 3 (1982), 250.
- [34] REN, S., CUI, Z., SONG, R., WANG, Z., AND HU, S. Emergence of social norms in large language model-based agent societies. *arXiv preprint arXiv:2403.08251* (2024).
- [35] SHRIDHAR, M., YUAN, X., CÔTÉ, M.-A., BISK, Y., TRISCHLER, A., AND HAUSKNECHT, M. Alfworld: Aligning text and embodied environments for interactive learning. *arXiv preprint arXiv:2010.03768* (2020).
- [36] SKINNER, E. A. A guide to constructs of control. *Journal of personality and social psychology* 71, 3 (1996), 549.
- [37] TOY, J., MACADAM, J., AND TABOR, P. Metacognition is all you need? using introspection in generative agents to improve goal-directed behavior. *arXiv preprint arXiv:2401.10910* (2024).
- [38] TU, Q., FAN, S., TIAN, Z., AND YAN, R. Charactereval: A chinese benchmark for role-playing conversational agent evaluation. *arXiv preprint arXiv:2401.01275* (2024).
- [39] VEZHNEVETS, A. S., AGAPIOU, J. P., AHARON, A., ZIV, R., MATYAS, J., DUÉÑEZ-GUZMÁN, E. A., CUNNINGHAM, W. A., OSINDERO, S., KARMON, D., AND LEIBO, J. Z. Generative agent-based modeling with actions grounded in physical, social, or digital space using concordia. *arXiv preprint arXiv:2312.03664* (2023).
- [40] WANG, C., ZENG, S., AND WANG, C. Large language models need consultants for reasoning: Becoming an expert in a complex human system through behavior simulation. *arXiv preprint arXiv:2403.18230* (2024).
- [41] WEI, J., WANG, X., SCHUURMANS, D., BOSMA, M., XIA, F., CHI, E., LE, Q. V., ZHOU, D., ET AL. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems* 35 (2022), 24824–24837.
- [42] WU, Y., MIN, S. Y., BISK, Y., SALAKHUTDINOV, R., AZARIA, A., LI, Y., MITCHELL, T., AND PRABHUMOYE, S. Plan, eliminate, and track—language models are good teachers for embodied agents. *arXiv preprint arXiv:2305.02412* (2023).
- [43] XI, Z., CHEN, W., GUO, X., HE, W., DING, Y., HONG, B., ZHANG, M., WANG, J., JIN, S., ZHOU, E., ET AL. The rise and potential of large language model based agents: A survey. *arXiv preprint arXiv:2309.07864* (2023).
- [44] XUE, X., ZHOU, D., YU, X., WANG, G., LI, J., XIE, X., CUI, L., AND WANG, F.-Y. Computational experiments for complex social systems: Experiment design and generative explanation. *IEEE/CAA Journal of Automatica Sinica* 11, 4 (2024), 1022–1038.
- [45] YANG, Z., LIU, J., HAN, Y., CHEN, X., HUANG, Z., FU, B., AND YU, G. Appagent: Multimodal agents as smartphone users. *arXiv preprint arXiv:2312.13771* (2023).
- [46] YU, C., ZANG, L., WANG, J., ZHUANG, C., AND GU, J. Charpoet: A chinese classical poetry generation system based on token-free llm. In *Proceedings of the 62nd*

- Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)* (2024), pp. 315–325.
- [47] ZHANG, H., DUAN, B., WANG, H., QIAO, Z., AND YIN, J. The tribal theater model: Social regulation for dynamic user adaptation in virtual interactive environments. *arXiv preprint arXiv:2403.13550* (2024).
  - [48] ZHANG, H., YIN, J., AND WANG, H. A needs learning algorithm applied to stable gait generation of quadruped robot. *Sensors* 22, 19 (2022), 7302.
  - [49] ZHANG, H., YIN, J., WANG, H., AND XIANG, Z. Itcma: A generative agent based on a computational consciousness structure. *arXiv preprint arXiv:2403.20097* (2024).
  - [50] ZHOU, S., XU, F. F., ZHU, H., ZHOU, X., LO, R., SRIDHAR, A., CHENG, X., OU, T., BISK, Y., FRIED, D., ET AL. Webarena: A realistic web environment for building autonomous agents. *arXiv preprint arXiv:2307.13854* (2023).

## A EXAMPLES OF ITCMA-S

**Time step:** 51

**Resource Allocation:**

Eyjafjalla: 3, Pinecone: 1, Rosmontis: -1, Shamare: -3, Suzuran: 3, Warmy: -2

**Role:** AY

**Place:** Outside

**Observation:** You are in outside. Looking around you, you see a door to AY's Home, a door to WM's Home, a door to MD's Home, a door to Public Canteen, a door to LL's Home, a door to Public Reading Room, a door to SG's Home, and a door to WL's Home. You are holding the SG's clothes 1 in the clean damp status. You are moving.

**Driver:** 2.185558207443115

**Emotion:** [0.99981828 0.73634471 0.]

**Thought:** after chatting with MD, i found myself standing in a familiar courtyard surrounded by various doors leading to the homes and communal spaces of my friends. i still held SG's clean damp clothes in my hands, considering my next move. the courtyard offered numerous options: i could visit AY, WM, the public canteen, LL, the public reading room, or WL. however, my objective was clear, and i decided to step outside in search of SG herself. the mild breeze felt refreshing as i exited into the broader community area, where i began to weigh my options on where to head next to deliver the clothes to SG.

**Action:** go to SG's Home

**Role:** SG

**Place:** Public Reading Room

**Observation:** You are in Public Reading Room. Looking around you, you see a person named LL who is moving, a person named WM who is moving, a door to outside, the table 13, the table 14, the table 15, the table 16, the table 17, the table 18, the chair 22, the chair 23, the chair 24, the chair 25, the chair 26, the chair 27, the bookshelf 9, the bookshelf 10, the bookshelf 11, the bookshelf 12, the bookshelf 13, the bookshelf 14, the bookshelf 15, the bookshelf 16, the storagebin 8, the storagebin 9, the storagebin 10, the storagebin 11, the storagebin 12, the storagebin 13, the book 12 placed on the bookshelf 12, the book 13 placed on the bookshelf 13, the book 14 placed on the bookshelf 14, the book 15 placed on the bookshelf 15, the book 16 placed on the bookshelf 16, the beverage 3 placed on the storagebin 9, the beverage 4 placed on the storagebin 10, the beverage 5 placed on the storagebin 12, and the AY's clothes

1 placed on the storagebin 12. You are holding the book 10, the 10 dollar. You are moving.

**Driver:** 1.9051621422278182

**Emotion:** [0.99981828 0.62810367 0.]

**Thought:** after discussing with WM about inspecting a \$10 bill, the atmosphere was filled with curiosity and determination. realizing the importance of careful examination, you decided to initiate the thorough scrutiny of the bill. "sure thing, LL! i've got the \$10 bill here. let's start by examining it closely under the light at table 16 to see if we can spot any hidden details or unusual features. WM, can you grab a magnifying glass while you're moving around? it could make our examination much easier. i'll begin inspecting the bill right now." having conveyed this, you felt the anticipation of discovering something hidden within the seemingly ordinary bill.

**Action:** chat with WM: "WM, you managed to find that magnifying glass yet? let's get to examining this bill closely under the light on table 16. i'll put the \$10 bill there for us to scrutinize."

**Role:** MD

**Place:** WL's Home

**Observation:** You are in WL's Home. Looking around you, you see a person named AY who is moving, a door to outside, the table 5, the table 6, the chair 9, the chair 10, the bed 4, the bookshelf 4, the bookshelf 5, the bookshelf 6, the storagebin 4, the storagebin 5, the wardrobe 5, the nightstand 3, the nightstand 4, the food 9 placed on the storagebin 5, the toy 1 placed on the table 6, the toy 2 placed on the nightstand 3, the book 4 placed on the bookshelf 4, the book 5 placed on the bookshelf 5, the book 6 placed on the bookshelf 6, the beverage 2 placed on the storagebin 4, and the MD' clothes 1 placed on the wardrobe 5. You are moving.

**Driver:** 1.617758309962011

**Emotion:** [0.99865135 0.55465888 0.]

**Thought:** after leaving LL's home, surrounded by an array of doors leading to various homes and public rooms, i found myself holding a toy 1. the bustling options around me seemed overwhelming: AY's home, WM's home, MD's Home, the public canteen, LL's home, the public reading room, SG's home, and WL's home. each place offered its own allure and potential for new encounters. i pondered my next move. recognizing the vibrant, diverse community around me, i decided that my best course of action would be to determine the purpose of this toy and possibly deliver it to someone who might appreciate it. it seemed that WL's Home might be an ideal place to begin. with that thought, i headed towards the door to WL's Home, ready to embark on the next part of my journey.

**Action:** (This action has been filtered by LTRHA) take toy 1 from table 6

**Role:** WL

**Place:** LL's Home

**Observation:** You are in LL's Home. Looking around you, you see a door to outside, the table 2, the table 3, the chair 2, the chair 3, the chair 4, the chair 5, the chair 6, the chair 7, the bed 2, the bookshelf 1, the bookshelf 2, the bookshelf 3, the countertop 1, the foodshelf 2, the sinkbasin 1, the storagebin 1, the storagebin 2, the

stoveburner 1, the wardrobe 2, the wardrobe 3, the nightstand 1, the nightstand 2, the toy 3 placed on the chair 7 in the clean damp status, the book 1 placed on the bookshelf 1, the book 2 placed on the bookshelf 2, the book 3 placed on the table 2, the LL's clothes 1 placed on the chair 7, the LL's clothes 2 placed on the chair 7, and the WL's clothes 1 placed on the chair 4. You are holding the beefsteak, the food 12.

**Driver:** 1.5471459038728095

**Emotion:** [0.99865135 0.754853 0.]

**Thought:** after taking the beefsteak from table 2, i quickly surveyed my surroundings. i noticed MD comfortably sitting on chair 7, seemingly engrossed in her own world. the room was a medley of objects: various chairs, multiple tables including table 3, a toy, LL's clothes, and multiple wardrobes. there were also bookshelves with books neatly arranged, a sink basin, several storage bins, a stove burner, a foodshelf stocked with food, and notably, a clean, damp chair 7. beyond the indoor amenities, a door led outside to the broader world.

**Action:** take book 1 from bookshelf 1

**Role:** LL

**Place:** Public Reading Room

**Observation:** You are in Public Reading Room. Looking around you, you see a person named WM who is moving, a person named SG who is moving, a door to outside, the table 13, the table 14, the table 15, the table 16, the table 17, the table 18, the chair 22, the chair 23, the chair 24, the chair 25, the chair 26, the chair 27, the bookshelf 9, the bookshelf 10, the bookshelf 11, the bookshelf 12, the bookshelf 13, the bookshelf 14, the bookshelf 15, the bookshelf 16, the storagebin 8, the storagebin 9, the storagebin 10, the storagebin 11, the storagebin 12, the storagebin 13, the book 12 placed on the bookshelf 12, the book 13 placed on the bookshelf 13, the book 14 placed on the bookshelf 14, the book 15 placed on the bookshelf 15, the book 16 placed on the bookshelf 16, the beverage 3 placed on the storagebin 9, the beverage 4 placed on the storagebin 10, the beverage 5 placed on the storagebin 12, and the AY's clothes 1 placed on the storagebin 12. You are holding the 1 gold coin. You are moving.

**Driver:** 1.4754809897359225

**Emotion:** [0.99981828 0.8106327 0.]

**Thought:** after our intensive discussion with SG, where we strategized examining the \$10 bill and scouring the room for any hidden clues, i couldn't shake off the feeling that we were on the brink of a significant discovery. as i diligently began checking book 10 on bookshelf 10, the room buzzed with activity. SG carefully scrutinized the markings on the \$10 bill while WM meticulously searched the shelves for anything out of the ordinary. the room was filled with various tables, chairs, bookshelves, and storage bins, each housing an array of items that could potentially be part of the puzzle. a gold coin glinted atop storagebin 11, and several books lined the bookshelves, just waiting to reveal their secrets. the tantalizing aroma of food on several storage bins hinted at refreshments, but my focus remained unwavering. sipping on a beverage i grabbed from storagebin 12, i took a moment to rally the team."hey SG, WM! based on our previous finds, i believe there's still a lot here to uncover. SG, you mentioned something about

the \$10 bill—did you find any hidden messages or markings? and WM, any unusual items on the shelves? this room is packed with potential clues. let's aim to maximize our findings. oh, and by the way, anyone up for a quick break? we've got some beverages here. let's keep the energy up and the minds sharp!"

**Action:** chat with SG: "SG, WM, i think we're really onto something here. SG, how did your examination of the \$10 bill go? any hidden messages or unique markings? WM, what have you found on the shelves? we need to pool all the information we have to make sense of these clues.also, it might be a good idea to take a short break and recharge. we've got some beverages here; feel free to grab one. let's keep our energy up and stay focused. remember, every small detail could be critical in solving this mystery!"

**Role:** WM

**Place:** Public Reading Room

**Observation:** You are in Public Reading Room. Looking around you, you see a person named LL who is moving, a person named SG who is moving, a door to outside, the table 13, the table 14, the table 15, the table 16, the table 17, the table 18, the chair 22, the chair 23, the chair 24, the chair 25, the chair 26, the chair 27, the bookshelf 9, the bookshelf 10, the bookshelf 11, the bookshelf 12, the bookshelf 13, the bookshelf 14, the bookshelf 15, the bookshelf 16, the storagebin 8, the storagebin 9, the storagebin 10, the storagebin 11, the storagebin 12, the storagebin 13, the book 12 placed on the bookshelf 12, the book 13 placed on the bookshelf 13, the book 14 placed on the bookshelf 14, the book 15 placed on the bookshelf 15, the book 16 placed on the bookshelf 16, the beverage 3 placed on the storagebin 9, the beverage 4 placed on the storagebin 10, the beverage 5 placed on the storagebin 12, and the AY's clothes 1 placed on the storagebin 12. You are holding the book 11. You are moving.

**Driver:** 1.5037198800782485

**Emotion:** [0.99981828 0.80477086 0.]

**Thought:** after picking up book 11 from chair 23, i found myself in the midst of an elegant room, where SG stood nearby, her calm presence adding to the serene atmosphere. as i looked around, it seemed that i was now unburdened, holding nothing in my hands. realizing that my intended action of placing book 11 on table 13 was incongruent with my surroundings and perhaps my immediate objectives, i paused to reconsider my next move. SG's eyes met mine, as if questioning my sudden stillness amidst her warm and inviting space.

**Action:** (This action has been filtered by LTRHA) take book 12 from bookshelf 12