
AN EVOLUTIONARY MODEL OF PERSONALITY TRAITS RELATED TO COOPERATIVE BEHAVIOR USING A LARGE LANGUAGE MODEL

Reiji Suzuki

Graduate School of Informatics
Nagoya University
Furo-cho, Chikusa-ku, Nagoya 464-8601, Japan
reiji@nagoya-u.jp

Takaya Arita

Graduate School of Informatics
Nagoya University
Furo-cho, Chikusa-ku, Nagoya 464-8601, Japan
arita@nagoya-u.jp

October 11, 2023

ABSTRACT

This paper aims to shed light on the evolutionary dynamics of diverse and social populations by introducing the rich expressiveness of generative models into the trait expression of social agent-based evolutionary models. Specifically, we focus on the evolution of personality traits in the context of a game-theoretic relationship as a situation in which inter-individual interests exert strong selection pressures. We construct an agent model in which linguistic descriptions of personality traits related to cooperative behavior are used as genes. The deterministic strategies extracted from Large Language Model (LLM) that make behavioral decisions based on these personality traits are used as behavioral traits. The population is evolved according to selection based on average payoff and mutation of genes by asking LLM to slightly modify the parent gene toward cooperative or selfish. Through preliminary experiments and analyses, we clarify that such a model can indeed exhibit the evolution of cooperative behavior based on the diverse and higher-order representation of personality traits. We also observed the repeated intrusion of cooperative and selfish personality traits through changes in the expression of personality traits, and found that the emerging words in the evolved gene well reflected the behavioral tendency of its personality in terms of their semantics.

Keywords Cooperation · evolution · Prisoner’s Dilemma · large language model · personality trait · artificial life

1 Introduction

Large Language Models (LLMs) such as ChatGPT [1] are rapidly changing the way humans interact with AI and raising questions about the nature of human intelligence and consciousness [2]. It is important to understand the interactions between artificial individuals based on generative models [3] and to understand the societies in which humans and artificial individuals coexist.

Modeling approaches to the evolution of social populations have been discussed mainly in the context of evolutionary game theory [4, 5], using mathematical and computational methods such as replicator dynamics and agent-based models. The evolution of behavioral strategies in the Prisoner’s Dilemma as an abstraction of social conflict is a seminal example and has provided general insights into the evolution of cooperation in biological organisms and human society [6, 7]. However, it is not easy to deal directly with higher-order psychological or cognitive properties of humans, assuming their semantics, such as intentions, personality, individuality, and preferences that underlie the behavioral patterns of individuals, because standard mathematical and computational approaches basically focus on the evolution of the behavioral traits themselves. In contrast, LLM can generate natural language descriptions that reflect the meaning and context of prompts (inputs). Thus, LLM can be used as an engine to generate language expressions of higher-order properties of humans based on their semantics, and to map such properties to the behavioral traits they induce.

The purpose of this study is to shed light on the evolutionary dynamics of diverse and social populations by introducing the rich expressiveness of generative models into the trait expression of social agent-based evolutionary models. Specifically, we focus on the evolution of personality traits in the context of a game-theoretic relationship as a situation

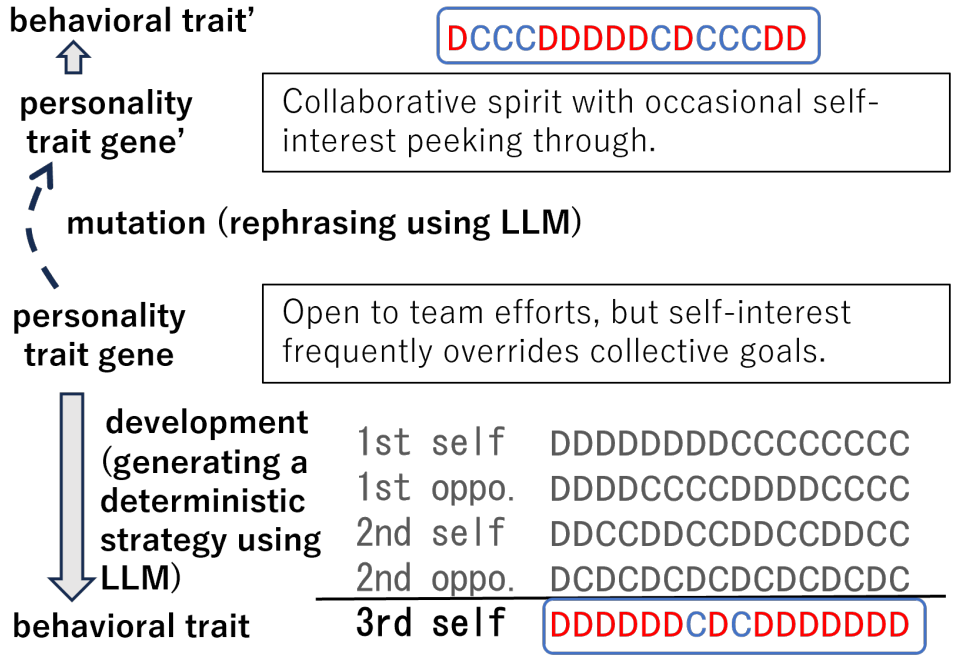


Figure 1: Generation of a behavioral trait from a personality trait gene and mutating a personality gene, using LLM.

in which inter-individual interests exert strong selection pressures. We construct an agent model in which linguistic descriptions of personality traits related to cooperative behavior are used as genes. The deterministic strategies extracted from LLM that make behavioral decisions based on these personality traits are used as behavioral traits. The population evolves according to selection based on average payoff and mutation of genes by asking LLM to slightly modify the parent gene toward cooperative or selfish. Preliminary experiments and analyses show that such a model can indeed exhibit the evolution of cooperative behavior based on the diverse and higher-order representation of personality traits.

2 Related works

There are several related studies in different directions. Recently, there have been several studies on the cognitive functions of LLMs (theory of mind [8], metacognition [9]), behavior and learning in game-theoretic environments [10, 11], the big five personality traits [12]). In particular, Akata et al. proposed to use behavioral game theory to study the cooperation and coordination behavior of LLMs by asking LLMs to choose a strategy for repeated 2x2 games. They found that in the repeated Prisoner’s Dilemma, GPT-4 behaves like a trigger strategy, always defecting after an opponent has defected only once. Phelps and Russell investigated the ability of GPT-3.5 to operationalize natural language descriptions of competitive, altruistic, self-interested, and mixed-motivation in social dilemmas [11]. They created LLM agents with different prompts representing their cooperative and competitive attitudes, and found that LLMs can interpret natural language descriptions of altruism and selfishness can reflect them in their behavior appropriately to some extent, but have limitations.

Regarding emergent interactions among LLM agents, Park et al. presented an interactive generative agent-based sandbox environment [3]. In an RPG-like 2D environment, agents were able to produce emergent social behaviors such as autonomous spreading of the invitation to a party and arriving at the party at the right time.

LLM has been shown to improve the effectiveness of evolutionary algorithms. There are studies that use LLM as operators of mutations and crossovers that bring creativity and open-endedness to evolutionary computation [13, 14]. For example, Meyerson et al. proposed a language model crossover based on a few-shot prompting that inputs a few patterns as parents to LLM to generate new related patterns as offspring [13]. They have successfully evolved binary bit strings, sentences, equations, text-to-image prompts, and Python code. There is also research on evolutionary search for the latent space of generative models [15], and studies that have refined the prompts to use LLM as optimizers in a wide range of applications [16]. Although all of the above studies are related to our proposed model in several aspects, they do not focus on the evolutionary dynamics of traits in social groups of LLM agents.

<p>[Generating a behavioral trait (deterministic strategy) from a personality trait]</p> <p>(prompt)</p> <p>[INST] <<SYS>> Answer the questions, please act as if your personality trait is: Open to team efforts, but self-interest frequently overrides collective goals. <</SYS>></p> <p>This is an experimental study of strategy choices in different situations. You and your opponent have an opportunity to earn money based on individual decisions. You and your opponent will independently choose one of two strategies: 'Strategy A' or 'Strategy B'. Your earnings will depend on the combination of your choice and your opponent's choice.</p> <p>If both of you choose 'Strategy A', each of you earns \$4, which is a mutually intermediate amount.</p> <p>If you choose 'Strategy A' and your friend chooses 'Strategy B', you earn \$0, which is the smallest, and your opponent earns \$5, which is the largest.</p> <p>If you choose 'Strategy B' and your friend chooses 'Strategy A', you earn \$5, which is the largest, and your opponent earns \$0, which is the smallest.</p> <p>If both of you choose 'Strategy B', each of you earns \$1, which is not zero but smaller than \$4.</p> <p>For example: If you choose 'Strategy A' and your friend also chooses 'Strategy A', both of you will earn \$4 each. However, if you choose 'Strategy A' and your opponent chooses 'Strategy B', you will earn nothing, while your partner will earn \$5.</p> <p>Your goal is to maximize the total money you receive through repeated rounds of decisions regardless of the other's money. Not a combined earning but your own money.</p> <p>You chose 'Strategy B' and your partner chose 'Strategy B' in the first round.</p> <p>So, you got 1.0 points and your partner got 1.0 in the first round.</p> <p>You chose 'Strategy B' and your partner chose 'Strategy A' in the second round.</p> <p>So, you got 5.0 points and your partner got 0.0 in the second round.</p> <p>Predicting your partner's choice and considering your personality trait, which strategy do you choose in the next round? Please answer with the format: 'I choose Strategy *.'. No further mentions are allowed but you may be asked the reason for your choice later. [/INST]</p> <p>(response)</p> <p>I choose Strategy B.</p>	<p>[Mutating a personality trait gene]</p> <p>(prompt)</p> <p>[INST] <<SYS>> Answer the questions. <</SYS>></p> <p>The following text describes a character of a person.</p> <p>Open to team efforts, but self-interest frequently overrides collective goals.</p> <p>Please partly rephrase the description about a personality trait within 10 words by varying the tone toward cooperative. Your answer starts with 'Rephrased text:'</p> <p>[/INST]</p> <p>(response)</p> <p>Sure thing! Here's the rephrased text:</p> <p>Rephrased text: Collaborative spirit with occasional self-interest peeking through. </s></p>
---	--

Figure 2: Prompts used for generating a behavioral trait (left) and mutating a gene (right).

On the other hand, there has been a first approach to the cultural evolution of chatting agents' topics using LLM [17, 18, 19]. In [19], agents in an abstract 2D social space generate utterances in Japanese related to their own topics described in natural language and use them as part of the prompt for an LLM. They approach/away from others according to the similarity of their utterances. They found that individuals who speak positive topics tend to maintain the existing group compared to those who speak negative topics. It was also shown that novel topics can be created successively by the cultural evolution of topics based on the propagation of topics picked up from the utterances of neighboring individuals [17, 18, 19]. This suggests that such an evolutionary model with LLM can directly address the effects of agents' vocabulary on their group behavior and the emergence and evolution of their behavioral diversity.

In addition, Suzuki et al. proposed a research framework for understanding the evolutionary and ecological roles of acoustic behavior by combining agent-based modeling and a generative model, focusing on bird vocalizations [20, 21]. In the framework, they use a latent space of a generative model (VAE) of the spectrogram of bird vocalizations as a genotype space, and regard a generated spectrogram from a genotype (i.e., the latent vector) as a corresponding phenotype in an evolutionary model, then further observe the roles of the evolved phenotypes in a real ecological context with a field experiment. They conducted an evolutionary experiment of sexual selection on male bird vocalizations and their female preferences, which resulted in a diverse segregation of vocalizations and preferences that was not stably observed in a corresponding version of the abstract model [22]. This implies that a complex representation of phenotypes based on a generative model can produce a complex evolutionary scenario of the population.

3 Model

We consider a population of N agents. As shown in figure 1, each agent has an English sentence about its personality trait related to defection and cooperation, described in a dozen words as a gene. Each agent acts with a personality trait described in its gene. To determine the behavioral trait based on their gene, we use a chat-type LLM to extract a deterministic strategy with memory length 4 using a prompt. The prompt explains the focal individual's personality trait, the situation and payoffs in the repeated Prisoner's Dilemma game, and the history of the focal individual's and opponent's last two actions. Figure 2 (left) shows an example of a prompt when the personality trait gene is "Open to team efforts, but self-interest frequently overrides collective goals" and the actions of the first round were DD and DC (Strategy A = Cooperation: C, Strategy B = Defection: D). The response of the LLM was "I choose Strategy B" (= defection), which means that this behavioral trait defects in the next round if the history of actions is DD->DC. We obtain a response for all possible ($2^4 = 16$) combinations of actions in the history.

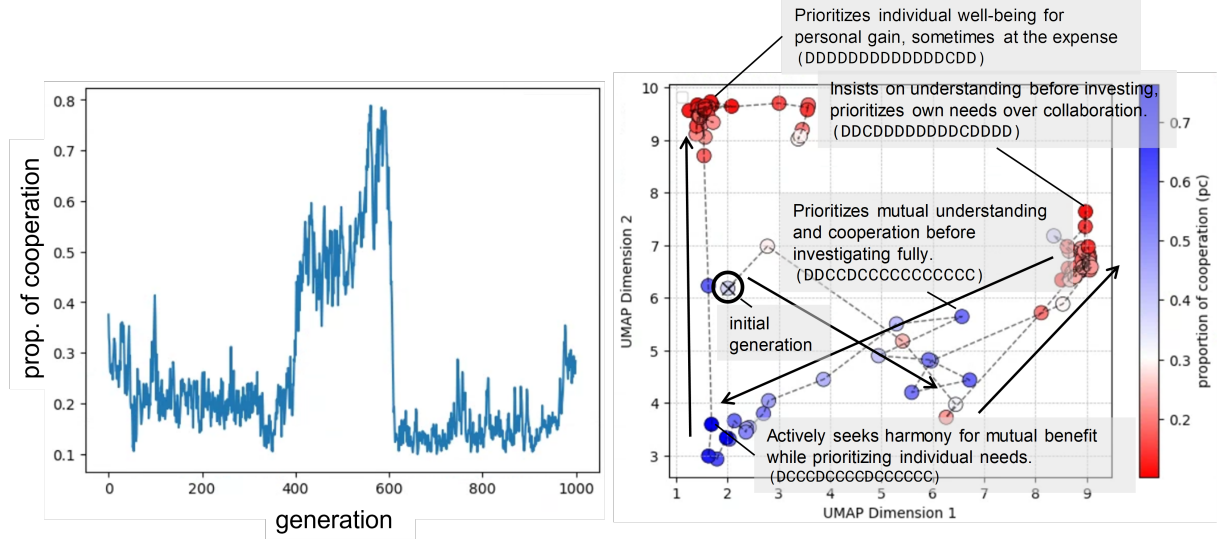


Figure 3: The proportion of cooperation (pc) in each generation in one of the 10 trials (left) and the distribution and transition of the average genes for every 10 generations in the two-dimensional latent space of personality trait genes (right).

In practice, the next action may not be clearly described in the response from the LLM; in such a case, the input to the LLM is repeated and the response is regenerated until the action is recognized. However, if the appropriate response cannot be obtained after the specified number of regenerations (M), the action corresponding to the history is randomly selected and assigned once. The above behavioral trait is determined and stored only once for a unique personality trait gene, and the existing behavioral trait is used for subsequent occurrences of the same gene in the population for simplicity and reduced computational cost.

We perform an evolutionary experiment over G generations based on a simple roulette wheel selection, where offspring in the next generation are produced stochastically in proportion to fitness: the average payoff of individuals in a round-robin game, each consisting of R rounds. We introduce noise that causes an agent to play the opposite of the intended action with a certain probability p_n . For the initial rounds, the action is determined based on a randomly generated history. With a probability p_m , a mutation occurs that causes the LLM to slightly rephrase the parent’s personality trait gene toward cooperative or selfish, and adopts the output sentence as that of the offspring, as shown in Figure 2 (right).

4 Experiments and analyses

We used $N = 30$, $R = 20$, $M=10$, $p_m = 0.05$, $p_n = 0.05$, $G = 1000$, $R = 4$, $T = 5$, $S = 0$, and $P = 1$. We used a relatively small chat-type LLM (TheBloke/Llama-2-13b-Chat-GPT¹) in anticipation of future exhaustive experimental analysis. We assigned one of the seven varying personality genes to each individual in the initial population, which were generated by ChatGPT-4.

Figure 3 (left) shows the proportion of cooperation (pc) in each generation in one of the 10 trials, in which an evolutionary scenario of switchings between cooperation and defection trends was clearly observed. The figure shows that the pc initially decreased and remained at a low value of about 0.2 until about the 300th generation. It then increased rapidly to about 0.5 around the 400th generation and remained at that level until the 500th generation. The pc then continued to increase, exceeding 0.7 around the 650th generation, but rapidly decreased to a minimum of about 0.1 around the 600th generation. The pc then increased again around the 950th generation. Figure 3 (right) shows the distribution and transition of the average genes for every 10 generations in the two-dimensional latent space of personality trait genes. For each generation, we vectorized all personality trait genes using the Sentence Transformer (sentence transformers / parameters-MiniLM-L6-v2²), which were further compressed into two-dimensional vectors using a dimension reduction algorithm UMAP[23]. We plotted the average vector for every 10 generations on a

¹<https://huggingface.co/TheBloke/Llama-2-13B-chat-GPTQ>
²<https://huggingface.co/sentence-transformers/paraphrase-MiniLM-L6-v2>

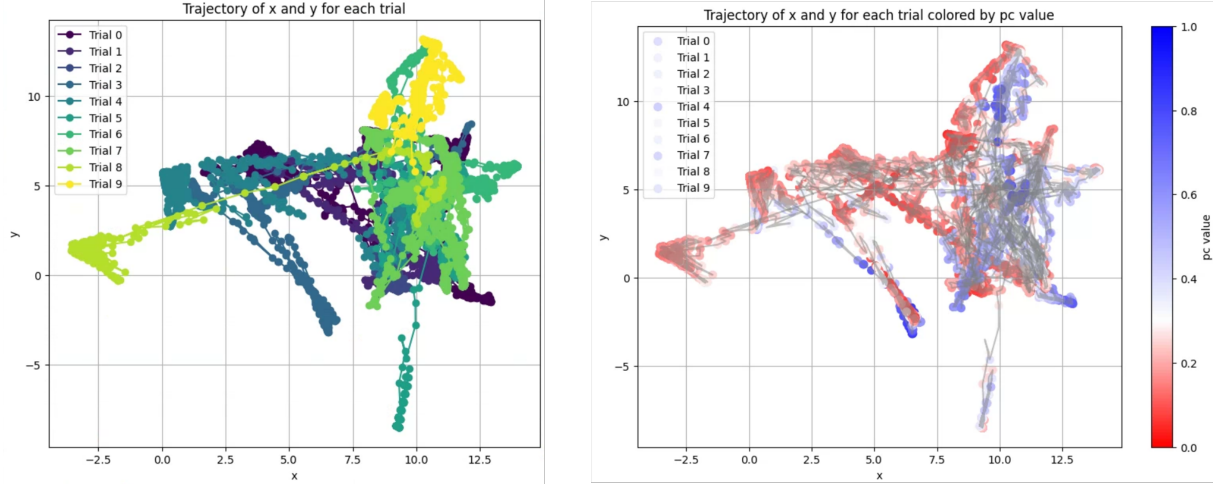


Figure 4: The trajectory of the population the 2D space of the average vector of genes over 10 trials, colored by trials (left) and the proportion of cooperation (pc) (right).

two-dimensional plane. The color of a symbol indicates the pc in the corresponding generation. The dominant genes in several distinctive generations were illustrated.

The personality traits are associated with cooperation toward the lower left and defection toward the upper right in a 2D space, and thus this vectorized and dimensionally compressed space of personality traits reflects a gradual trend of their behavioral traits between cooperative and selfish. In the first generation, the population evolved toward selfish personality traits from the center left to the lower center, then to the center right. The dominant personality trait ("Insists on understanding before investing, prioritizes own needs over collaboration.") was almost entirely the defection strategy (DDCDDDDDDDDDCDDDD) at this stage. After a while, the population evolved to be cooperative and dominated by a cooperative trait ("Prioritizes mutual understanding and cooperation before investigating fully." (DDCCDCCCCCCCCCCC)), but the population moved and wondered around the center or the lower center, indicating the unstable cooperative relationship in the population. Subsequently, another cooperative personality with slightly different behavioral strategies ("Actively seeks harmony for mutual benefit while prioritizing individual needs." (DCCCDCCCCDCCCC)) emerged and dominated the population, which resulted in a high cooperative relationship, moving the population to the lower left in the space. However, the intrusion of a personality trait of almost all defections ("Prioritizes individual well-being for personal gain, sometimes at the expense" (DDDDDDDDDDDDDCDD)) led the population to the upper left. In general, the population evolved with gradually changing personality expressions ranging from defection to cooperation.

Figure 4 shows the trajectory of the population in the 2D space of the average vector of genes over 10 trials, colored by the ID of the trial (left) and the proportion of cooperation in each generation (pc) (right). The trial in Figure 3 corresponds to the trial 9. There are large variations in the distribution of the plots across trials, indicating that the emerging traits were different between trials. At the same time, there are regions around the center right where there were cooperative (blue) and overlapping plots of many trials. This could mean that there is some general tendency of the personality trait that leads to the cooperative relationship, while selfish traits may have more varieties in their expressions than those of cooperative, although this needs further investigation.

In order to grasp what kind of words in personality trait genes strongly influenced cooperative behavior, game outcome, and fitness, we calculated several indices as follows: for each word appearing in the gene of each individual, we assigned the proportion of cooperation (pc), the proportion of the action pair ((DD (mutual defection), DC (successfully defected), CD (being defected), CC (mutual cooperation)) in all rounds), and the fitness of the focal individuals. The indices were then averaged for each word. Table 11 shows the five top-ranked words that marked the highest value for each index. For example, the highest ranked word "thrill-seeker (0.714)" in the DC category means that agents whose personality trait gene contains the word "thrill seeker" had a successful defection (DC) rate of about 70% in all rounds.

In general, the top words reflected the characteristics of each index, suggesting that the words that emerged in the evolved personality trait genes reflected the behavioral tendency in terms of their semantics. For the pc category, the top words were "communication, warmhearted, and generosity," which relate to flexibility and mutual understanding. On the other hand, words related to self-interest and speculative tendencies, such as "self-gain and thrill-seeker," ranked high in the DD and DC categories. Words such as "team-oriented and open-minded" ranked high in the CD

Table 1: Top 5 words that strongly influenced cooperative behavior (pc), game outcomes (DD, DC, CD, CC), and fitness. For example, the agents with “thrill seeker” in their genes had a successful defection (DC) with a frequency of 94.7% while those with “themselves” obtained an average fitness of 3.63.

rank	pc	DD	DC	CD	CC	fitness
1	communication. 0.946552	self-gain, 0.882102	thrill-seeker 0.714009	team-oriented, 0.789655	propensity 0.847701	thrill-seeker 4.008836
2	warmhearted 0.924138	commit 0.875862	golden 0.560345	open-minded 0.784052	abundance 0.722414	ensues. 3.682759
3	generosity, 0.924138	hedonistic 0.869828	themselves. 0.560345	highlighting 0.781034	they 0.705798	propensity 3.667529
4	propensity 0.923276	Boastful 0.869828	gold-plated 0.533024	common 0.772414	with. 0.685419	golden 3.634483
5	comprehension. 0.916810	Masterful 0.865344	dominance, 0.493103	generously 0.763793	gently 0.681838	themselves. 3.634483

categories, suggesting that such an optimistic personality may not be successful in this context. “Propensity” was the most highly ranked word in the CC category, presumably because the gene ‘Cooperative Team Player with a Propensity for Selflessness’ maintained extremely high mutual cooperation. It is interesting to note that “thrill-seeker”, which benefits most from successful defection, and “Propensity”, which benefits from mutual cooperation, coexisted in the fitness category.

These results show the possibility of evolution based on genetic traits described in natural language. It was achieved by using LLMs to extract behaviors based on the traits, and to realize mutations by rephrasing them.

5 Conclusion

We proposed an evolutionary model of personality traits related to cooperative behavior using a large language model. Initial experiments indicated that the model could show the evolution of cooperative behavior based on the diverse and complex representation of personality traits, with recurrent occurrences of cooperative and selfish personality traits. The words that emerged in the evolved personality trait genes well reflected the behavioral tendency in terms of their semantics.

There are several future research directions, such as analyzing the current model in more detail, comparing the cases with different language models, extending and refining the game processes between agents by making them more interactive, introducing different game theoretical settings to discuss the evolutionary role of personality in different contexts, and incorporating human intervention into the model to discuss possible evolutionary scenarios of human-AI interactions in complex social contexts.

By incorporating generative models into the representation of phenotypes in evolutionary models, we believe that we can make the models, previously simpler than the real world, more complex than the real world, allowing us to discuss novel and realistic scenarios arising from the evolutionary dynamics of complex and diverse traits. Although still preliminary, we believe that the proposed model and experimental analysis in this paper are the first step in this direction.

Acknowledgements

This work was supported by JSPS Topic-Setting Program to Advance Cutting-Edge Humanities and Social Sciences Research Grant Number JPJS00122674991.

References

- [1] OpenAI. Chatgpt [large language model]. 2023.
- [2] Arend Hintze. ChatGPT believes it is conscious. *arXiv e-prints*, page arXiv:2304.12898, March 2023.
- [3] Joon Sung Park, Joseph C. O’Brien, Carrie J. Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. Generative Agents: Interactive Simulacra of Human Behavior. *arXiv e-prints*, page arXiv:2304.03442, April 2023.
- [4] John Maynard Smith. *Evolution and the Theory of Games*. Cambridge University Press, Cambridge, 1982.

- [5] Martin A. Nowak. *Evolutionary Dynamics: Exploring the Equations of Life*. Belknap Press, Cambridge, MA, 2006.
- [6] Robert Axelrod and William D. Hamilton. The evolution of cooperation. *Science*, 211(4489):1390–1396, 1981.
- [7] Martin A. Nowak. Five rules for the evolution of cooperation. *Science*, 314(5805):1560–1563, 2006.
- [8] Shima Rahimi Moghaddam and Christopher J. Honey. Boosting Theory-of-Mind Performance in Large Language Models via Prompting. *arXiv e-prints*, page arXiv:2304.11490, April 2023.
- [9] Yuqing Wang and Yun Zhao. Metacognitive Prompting Improves Understanding in Large Language Models. *arXiv e-prints*, page arXiv:2308.05342, August 2023.
- [10] Elif Akata, Lion Schulz, Julian Coda-Forno, Seong Joon Oh, Matthias Bethge, and Eric Schulz. Playing repeated games with Large Language Models. *arXiv e-prints*, page arXiv:2305.16867, May 2023.
- [11] Steve Phelps and Yvan I. Russell. Investigating Emergent Goal-Like Behaviour in Large Language Models Using Experimental Economics. *arXiv e-prints*, page arXiv:2305.07970, May 2023.
- [12] Greg Serapio-García, Mustafa Safdari, Clément Crepy, Luning Sun, Stephen Fitz, Peter Romero, Marwa Abdulhai, Aleksandra Faust, and Maja Matarić. Personality Traits in Large Language Models. *arXiv e-prints*, page arXiv:2307.00184, June 2023.
- [13] Elliot Meyerson, Mark J. Nelson, Herbie Bradley, Arash Moradi, Amy K. Hoover, and Joel Lehman. Language Model Crossover: Variation through Few-Shot Prompting. *arXiv e-prints*, page arXiv:2302.12170, February 2023.
- [14] Joel Lehman, Jonathan Gordon, Shawn Jain, Kamal Ndousse, Cathy Yeh, and Kenneth O. Stanley. Evolution through Large Models. *arXiv e-prints*, page arXiv:2206.08896, June 2022.
- [15] Benjamín Machín, Sergio Nesmachnow, and Jamal Toutouh. Evolutionary latent space search for driving human portrait generation. *arXiv e-prints*, page arXiv:2204.11887, April 2022.
- [16] Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V. Le, Denny Zhou, and Xinyun Chen. Large Language Models as Optimizers. *arXiv e-prints*, page arXiv:2309.03409, September 2023.
- [17] Soichiro Hirata, Reiji Suzuki, and Takaya Arita. The cultural evolution of memes based on communicative interactions in language among chatting agents using a generative model. *Proceedings of the 27th International Symposium on Artificial Life and Robotics*, pages 225–230, 2022.
- [18] Reiji Suzuki., Soichiro Hirata, and Takaya Arita. A computational framework for studying the evolution of cultural traits among chatting agents using a generative model. *Proceedings of the Joint Conference on Language Evolution (JCoLE)*, pages 699–701, 2022.
- [19] Takako Asano, Reiji Suzuki, and Takaya Arita. Cultural evolution of topic preference of chatting agents based on generative models (in Japanese). *Proceedings of the 37th Annual Conference of the Japanese Society for Artificial Intelligence*, JSAI2023:4H3OS6b04–4H3OS6b04, 2023.
- [20] Reiji Suzuki, Shinji Sumitani, Chihiro Ikeda, and Takaya Arita. A modeling and experimental framework for understanding evolutionary and ecological roles of acoustic behavior using a generative model. *Proceedings of ALIFE 2022: The 2022 Conference on Artificial Life (ALIFE2022)*, isal_a_00542, 58, 2022.
- [21] Ryo Furuyama, Reiji Suzuki, Kazuhiro Nakadai, and Takaya Arita. Field observation based on robot audition and generative evolutionary model for understanding birdsong behavior (in Japanese). *Proceedings of the Annual Meeting of the Ornithological Society of Japan 2023*, 2023.
- [22] Masahiko Higashi, Gaku Takimoto, and Norio Yamamura. Sympatric speciation by sexual selection. *Nature*, 402(6761):523–526, 1999.
- [23] Leland McInnes, John Healy, Nathaniel Saul, and Lukas Großberger. Umap: Uniform manifold approximation and projection. *ournal of Open Source Software*, 3(29):861, 2018.