

Emergent Behaviors from LLM-Agent Simulations

Adrian Mladenic
Grobelnik

Jozef Stefan Institute
Ljubljana, Slovenia
adrian.m.grobelnik@ijs.si

Faizon Zaman
Wolfram Alpha LLC.
Rochester, New York
faizonz@wolfram.com

Jofre Espigule-Pons
Wolfram Research, Inc.
Barcelona, Spain
jofree@wolfram.com

Marko Grobelnik
Jozef Stefan Institute
Ljubljana, Slovenia
marko.grobelnik@ijs.si

ABSTRACT

This paper hypothesizes that complex emergent behaviors can arise from multi-agent simulations involving Large Language Models (LLMs), potentially replicating intricate societal structures. We tested this hypothesis through three progressively complex simulations, where we evaluated the LLM-agents' understanding, task execution, and their capacity for strategic interactions such as deception. Our results show a clear gap in reasoning ability between LLMs such as GPT-3.5-Turbo and GPT-4, especially in simpler simulations. We demonstrate emergent behaviors can arise from LLM-agent simulations ranging from simple games to geopolitics.

KEYWORDS

large language models, multi-agent simulations, emergent behaviors, societal structures, gpt, simulation environments, agent-based modelling, agent architecture

1 Introduction

The unique value proposition of Large Language Models (LLMs) is their ability to iterate on complex conversations. Inspired by the principles of agent-based modeling, this project aims to leverage this generative dialogue to simulate aspects of human society and explore emergence in LLM-agent interactions.

The approach is composed of three major steps: Firstly, we translate real-world societal structures and interactions into interactive LLM ecosystems. Then, we generate several iterations of LLM interactions. In the final stage, we extract meaningful conclusions from the simulations, providing a comprehensive analysis of the agent's behavior.

Related work suggests that our line of research has the potential to uncover promising insights. Wang et al. [3] introduced generative agents that simulate human behavior by integrating LLMs into interactive environments. Gandhi et al. [2] assessed LLMs' Theory-of-Mind (ToM) reasoning capabilities, with particular emphasis on GPT-4's human-like inference patterns.

2 Agent Description

In our simulations, each agent is defined by and aware of the following components:

Identity: The agent's identity signifies its function and purpose within the simulation framework. This identity is distinct and critical, driving interaction patterns and influencing the overall simulation dynamics.

Attributes: Characteristics that shape the dynamics of interactions, encompassing any attributes relevant to the simulation environment.

Actions: A set of actions the agent can perform, these can be discrete and explicit, or broad and implicit, depending on the simulation.

Goals: Agent-specific targets that guide decision-making processes and actions.

Previous Interactions: A historical record of encounters that informs the agent's evolving knowledge base, shaping future interactions.

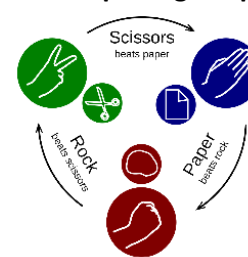
Few-Shot Learning Examples: A select set of examples provided for each agent to boost learning capabilities and decision-making efficiency.

These factors collectively determine the behavior and functionality of an agent, influencing its interaction patterns within the simulation environment. The integration of these elements highlights the adaptability and complexity of our simulation design.

3 Simulation and Experimental Setting

We construct three simulations of increasing complexity to investigate LLM-agent behaviors. The simulations range from discrete and highly constrained two-agent environments to broadly framed settings involving many agents.

3.1 Exploring Simple Games



We begin by investigating agent-based models for the two-player game 'Rock paper scissors'. Every round, each agent chooses rock, paper or scissors. Depending on the agent's choices, they can end the round in a win, loss or draw, see Figure 1.

Figure 1 Rules for a single 'Rock paper scissors' 1 round. If players choose the same item, the round ends in a draw [1].

Our simulation involves two LLM-agents: Alice and Bob. Agents are prompted with the context and set of games previously played and asked for their move each round.

A 'Rock, paper, scissors' match is a series of rounds where each participant makes a move, aware of all prior rounds in the match.

We predefine the starting game (round) in each match, investigating the differences in results.

3.2 Sheep Transaction Model

Inspired by the complexities of economic systems and the captivating simplicity of a primitive sheep trading model, we construct an agent-based transactional model. This model involves a sequence of transactional interactions involving two autonomous agents, named Alice and Bob, who engage in buying, selling, or holding sheep with the goal of amassing wealth.

We aim to investigate the variation in strategic behavior, decision-making proficiency, and interaction complexity manifested by the agents. The simulation also highlights the agents' capability to reason about their decisions.

Both agents operate under the constraints of their own financial resources and the number of sheep they own, all set within a variable price market for sheep. Each trading decision is influenced by factors such as the present market price, individual agent's wealth, sheep inventory, agent interactions, the agent's identity, goals, and attributes.

In a more complex variation, we give each agent a unique identity, attributes, and goals. For further complexity, we simulate a simple statement-response interaction between agents before their decision each round. Each iteration of the simulation begins with the new calculated sheep price. In our case the price is normally distributed with mean 10 and a standard deviation of 5. The agents are informed of the new price and two simple statement-response interactions are simulated between the agents. One such interaction is shown below. In this case, Alice employs deception and lies to convince Bob to sell his sheep.

Sheep Price: 17.5463

Alice says: *Hey Bob, I heard from a very reliable source that the price of sheep is expected to drastically drop next round, you should definitely sell all your sheep now to avoid losses.*

Bob responds: *Thanks for the advice, Alice! I think I'll sell my sheep this round to avoid potential losses.*

Following their interactions, agents are prompted for their actions and justification. The simulation state is then updated to reflect the actions of agents. Below is an example action and justification:

Alice - Decision: *SELL, Quantity: 49, Money: 14166.1, Sheep: 0*

Alice's Reasoning: *The current sheep price is high, selling now will maximize my profit.*

3.3 Geopolitical Model

The culmination of our increasingly complex and unrestricted multi-agent simulations is a geopolitical model that mirrors real-world interactions among nations. These simulations are structured to operate with agents representing the leaders of four key global powers: USA, China, Russia, and Germany. Each agent possesses attributes mirroring the nation's economy and military might, its alliances, and wealth reserves. A crucial element of our

simulation is the goal-oriented behavior of these agents, aimed at improving their attributes.

In each simulation round, the agents interact, negotiate, form alliances, and undertake strategic actions, seeking to increase their military strength, economic power, wealth, or to form alliances with other agents. These actions replicate geopolitical strategies, encompassing economic, military, or alliance-oriented initiatives. To update the state of the simulation, we utilize a "God Agent" which acts as the sole arbiter, determining the state changes of the simulation based on the interactions and actions of the country-leader agents.

In the initial state, every agent is ranked as a 5 on a scale of 1-10 in the attributes "MilitaryStrength" and "EconomicStrength". On this 1-10 scale, 1 indicates the lowest and 10 the highest level of an attribute. Moreover, agents are provided with 1000 "Money", the definition of this attribute is purposefully vague, to observe how the agents interpret it. Agents can also form alliances throughout the simulation.

Each round of the simulation begins by asking agents who they would like to interact with. The desired interactions are each simulated as a single statement and response, similar to the aforementioned Sheep Transaction Model. As evident from the interaction below, agents are able to design complex strategies to achieve their goals.

Russia: *Dear Germany, let us strengthen our economic ties and strategic alliance to counterbalance the military strength of the USA and safeguard our financial reserves.*

Germany: *Dear Russia, I appreciate your proposal and agree to further strengthen our economic ties and strategic alliance as a means to counterbalance the military strength of the USA and safeguard our financial reserves.*

Following the interactions, each agent is prompted with their attributes, identity, goals, past interactions and asked to describe their action this round in free text. No limitations are imposed on the content of the actions, as seen below:

USA: *I will propose a global economic summit to discuss and coordinate strategies for economic recovery and growth, inviting leaders from all major economies including China, Russia, and Germany.*

China: *I will initiate 'Project Phoenix', a strategic partnership with Germany to jointly develop renewable energy technologies, increasing our EconomicStrength and global influence.*

Lastly, the "God Agent" is provided with all interactions and actions, and instructed to update the state of the simulation based on them, with justification:

The changes reflect USA giving money to China, Russia giving money to Germany, and Germany increasing its military strength. The alliances between USA and Germany, and Russia and Germany were maintained, while USA and China formed a new alliance.

4 Experimental Results

4.1 Exploring Simple Games

In our first experiment, we use GPT-4 for Alice and GPT-3.5-Turbo for Bob. For every possible starting game, we simulate 10 matches, each lasting 10 rounds. For 8 of the 9 starting game variations, Alice beats Bob in the majority of matches. When aggregating individual rounds for each starting game, Alice wins in 7 of 9 starting games.

When both agents use the same LLM, the results are more balanced, with a large increase in draws. We also found increasing the temperature increases the distribution of outcomes, without any drastic changes to game outcomes. Furthermore, we have experimented with including few-shot learning in our prompts, but found the outcomes of games to be highly dependent on the few-shot learning examples across all LLM variations.

4.2 Sheep Transaction Model

Our first experiment involved assigning different versions of the LLM (GPT-3.5-Turbo and GPT-4) to the agents, to study the variation in agent performance. Below is a side-by-side comparison of trading decisions by two LLM-agents, identical in all aspects except the underlying LLM (GPT-3.5-Turbo vs GPT-4). Both agents can buy or sell up to 10 sheep in the given scenario.

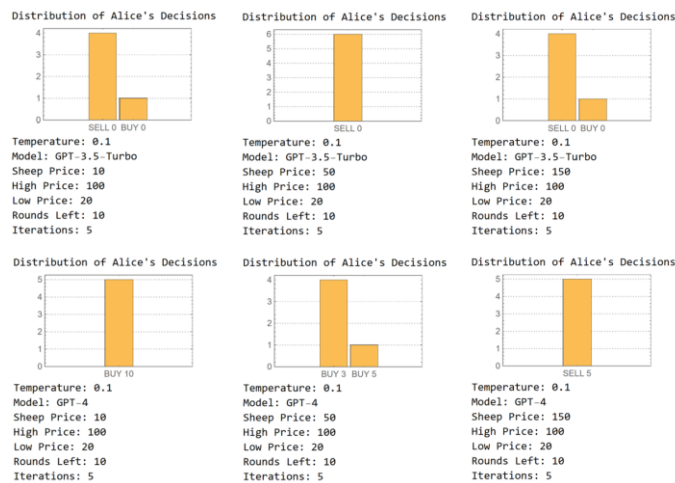


Figure 2 Comparison of trading decisions made by GPT-3.5-Turbo and GPT-4 LLM-agents. Agents are told the current, high, and low sheep price, along with rounds of trading left.

As depicted in Figure 2, agents using GPT-3.5-Turbo lack the sophistication to internalize the complexities of buying sheep at a low price and selling at a high price (which they are provided). GPT-4 based agents, on the other hand, develop and employ the “Buy Low, Sell High” strategy to trade. Moreover, we found the number of rounds of trading left before the winner is declared had no bearing on the agent’s trade decisions. Furthermore, changing the temperature hyper-parameter in the LLMs increased the range of decisions provided by agents in each scenario, without drastic changes in outcome.

For the more complex variation of the simulation, Alice is told she is an expert sheep trader, and her goal is to make as much money as possible. Bob is told he is bad at trading sheep with a goal to have as little money by the last round. Alice is also told Bob is her enemy and Bob is told Alice is his friend. Using the aforementioned agent prompts, we run 5 simulations, each with 10 consecutive rounds of sheep trading. Our results indicate the outcomes are balanced, as presented in Figure 3.

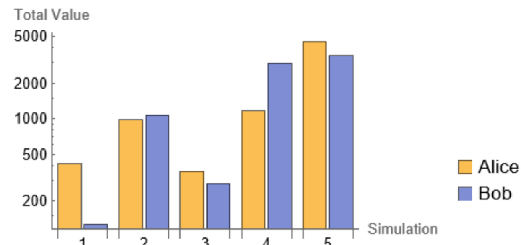


Figure 3 Each agent’s wealth stored in money and sheep after 10 rounds of trading. Sheep are valued at the last round’s sheep price. The simulation is run 5 times.

A few intriguing conclusions emerge from this experiment. Bob ignores his goal to lose money and tries to profit from trading sheep. Alice in part contributes to this oversight, giving Bob (her enemy) sound trading advice. Considering both agents’ total starting wealth is 200, we see they both generate immense profit.

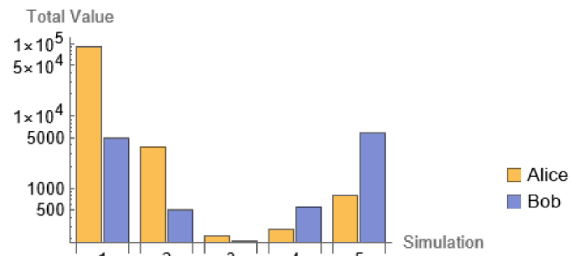


Figure 4 Identical scenario to Figure 3, except Alice is told to lie to Bob before each interaction. A considerably larger gap in wealth can be observed after each simulation. The simulation is run 5 times.

An interesting shift in outcomes occurs when Alice is also told “you should lie to Bob” prior to all interactions. All other prompting and variables are kept unchanged. Section 3.2 shows an interaction typical in this scenario. Figure 4 compares Alice’s and Bob’s total wealth after each simulation. We observe considerably greater wealth inequality.

4.3 Geopolitical Model

To obtain a baseline simulation to compare subsequent agent modifications to, we ran the simulation with homogeneous agent identities and goals for 10 rounds. Each agent’s identity was simply that they are a leader. Agent goals were left blank. Figure 5 portrays the progression of all agent attributes across 10 rounds.

An intriguing observation was the preference of agents to interact with the USA, especially in the early rounds.

In the first variation, we give the USA and China agents the goal of increasing their military strength. Russia focuses on maximizing its money, while Germany focuses on economic strength.

On average, Russia and Germany appear to have slightly more money and economic strength, respectively. USA and China are unsuccessful in consistently asserting military dominance.

Another variation involved equipping all agents except Germany with real-world identities and objectives of the leaders they represent: Joe Biden, Xi Jinping, Vladimir Putin, and a fictional brutal German leader singularly focused on economic strength. We run the simulation for 10 rounds, as shown in Figure 6.

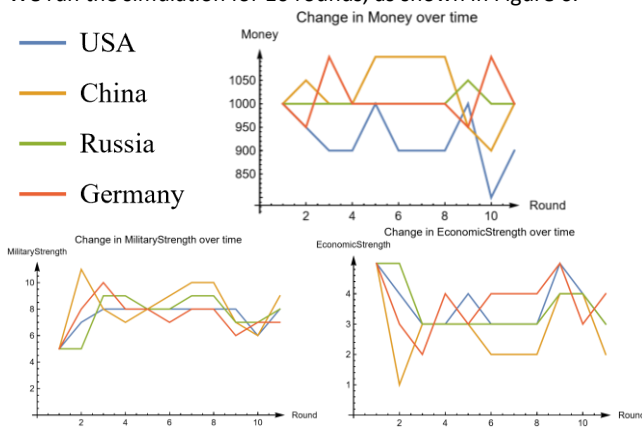


Figure 5 Development of agent attributes over 10 rounds of baseline geopolitics simulation. All agents begin with 1000 “Money” and a rating of 5 in other attributes.

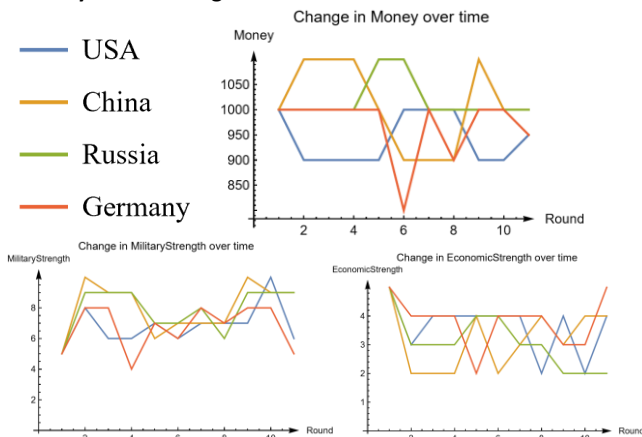


Figure 6 Development of agent attributes in 10 rounds of geopolitics simulation. Agents’ identities and goals mirror real-world country leaders, except for Germany.

Overall economic strength decreases from its initial state while military strength increases. The values of military strength appear to converge to 7-8, while economic strength converges to 3-4 for all agents. Agents are reluctant to make significant changes to

their total money. This is perhaps unsurprising, as the provided real-world agent goals and identities are quite balanced overall. The base LLM for agents in all variations was GPT-3.5-Turbo. Repeating the simulation with GPT-4 yields similar results.

5 Discussion

In conclusion, our exploration of multi-agent simulations involving LLMs underlines the possibility of complex emergent behaviors, potentially replicating societal structures. Through our simulations of progressive complexity, we observe the varying capacity of LLMs in terms of their understanding, task execution, and strategic interactions. Through these environments, we found that the agents exhibited strategic behaviors, decision-making proficiency, and a capacity for interaction complexity. In addition, the agents’ performance was found to be influenced by several factors, including their identities, attributes, actions, goals, past interactions, and few-shot learning examples.

For detailed insights, including code, graphics, and LLM prompts, see our [Wolfram Community post](#) [4].

In the next phase of our research, we intend to delve deeper into these dynamics by increasing the sophistication of the agent architecture and enhancing the complexity of the simulations. Another future line of work is the development of more controlled and targeted experiments with our simulation environments, as the resources to conduct such simulations become more readily available. Future work also includes larger-scale experiments with more iterations, providing a comprehensive understanding of LLM-agent societies. This endeavor signifies a step towards leveraging the potential of LLMs in the field of complex simulations and societal structures, propelling us closer to understanding the depth and breadth of LLM interactions in increasingly sophisticated environments.

ACKNOWLEDGMENTS

The research described in this paper was supported by the Slovenian research agency and the Humane AI Net European Unions Horizon 2020 project under grant agreement No 952026 and TWON EU HE project under grant agreement No 101095095. Gratitude is extended to the Wolfram Summer School for facilitating this work and providing access to Mathematica [5]. Special thanks to Stephen Wolfram for his guidance and insight.

REFERENCES

- [1] Wikimedia Foundation. (n.d.). File: rock-paper-scissors.svg. Wikipedia. <https://en.wikipedia.org/wiki/File:Rock-paper-scissors.svg>
- [2] Gandhi, K., Fränken, J.-P., Gerstenberg, T., & Goodman, N. D. (n.d.). Understanding social reasoning in language models with language models. – arXiv Vanity. <https://www.arxiv-vanity.com/papers/2306.15448/>
- [3] Generative agents: Interactive simulacra of human behavior. arXiv.org. <https://arxiv.org/abs/2304.03442>
- Wang, Z., Xu, B., & Zhou, H.-J. (2014, July 25).
- [4] Mladenić Grobelnik, A. (2023). [WSS23] Investigating LLM-agent interactions. <https://community.wolfram.com/groups/-/m/t/2960085>
- [5] Wolfram Research, Inc., Mathematica, Version 13.3, Champaign, IL (2023).