# Contents

# Generative Agents: Interactive Simulacra of Human Behavior

**Authors:** ['Joon Sung Park', 'Joseph C. O'Brien', 'Carrie J. Cai', 'Meredith Ringel Morris', 'Percy Liang', 'Michael S. Bernstein']

---

Here's a summary of the research paper, formatted as requested:

# Generative Agents: Interactive Simulacra of Human Behavior

## 1. Title, Authors, and Publication Details

- **Title:** Generative Agents: Interactive Simulacra of Human Behavior
- **Authors:** Joon Sung Park, Joseph C. O'Brien, Carrie J. Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein
- **Publication:** UIST '23, October 29-November 1, 2023, San Francisco, CA, USA

## 2. Research Questions/Objectives

- How can we create interactive artificial societies that exhibit believable human behavior?
- Can generative models be used to create agents that simulate believable human behavior, both individually and in emergent group dynamics?
- What architecture is required to enable agents to remember, retrieve, reflect, interact, and plan in dynamically evolving circumstances?

## 3. Methodology Used

- **Agent Architecture Design:** Developed a novel agent architecture comprising a memory stream, reflection mechanism, and planning module, leveraging a large language model (ChatGPT).
- **Sandbox Environment:** Implemented a sandbox environment ("Smallville") inspired by The Sims, populated with 25 generative agents.
- **Evaluations:**
  - **Controlled Evaluation:** "Interviewed" agents using natural language questions to assess self-knowledge, memory retrieval, planning, reactivity, and reflection. Compared the full architecture against ablations (no observation, no reflection, no planning; no reflection, no planning; no reflection) and a human-authored baseline.
  - **End-to-End Evaluation:** Observed the emergent social behaviors (information diffusion, relationship formation, coordination) of the agent community over two simulated days.

## 4. Key Findings and Results

- The full generative agent architecture produced more believable behavior than the ablated architectures and the human crowdworkers.
- Each component (observation, planning, and reflection) contributes critically to the believability of agent behavior. Ablations showed decreased performance with the removal of each component.
- Agents demonstrated emergent social behaviors:
  - Information diffused through the community (e.g., news of a mayoral candidacy and a Valentine's Day party).
  - New relationships formed between agents.
  - Agents coordinated to attend the Valentine's Day party.
- Common errors included failure to retrieve relevant memories, embellishments to memories, and overly formal speech patterns inherited from the language model.

## 5. Main Conclusions and Implications

- Generative agents can simulate believable human behavior by integrating large language models with mechanisms for memory, reflection, and planning.
- The proposed architecture enables agents to dynamically adapt to their experiences and environment, leading to emergent social dynamics.
- Generative agents have potential applications in social prototyping, human-centered design, and interactive environments.

## 6. Limitations Mentioned in the Paper

- The retrieval module could be enhanced to retrieve more relevant information.
- The current implementation is computationally expensive.
- The evaluation was limited to a relatively short timescale.
- The human baseline did not represent maximal human performance.
- The robustness of generative agents to prompt hacking, memory hacking, and hallucination is largely unknown.
- Generative agents may inherit biases from the underlying language models.
- The agents' behavior can be overly polite and cooperative due to instruction tuning.