

Policy Brief HAI Policy & Society May 2025

Simulating Human Behavior with Al Agents

Joon Sung Park, Carolyn Q. Zou, Aaron Shaw, Benjamin Mako Hill, Carrie J. Cai, Meredith Ringel Morris, Robb Willer, Percy Liang, Michael S. Bernstein

Al agents have been gaining widespread attention among the general public as Al systems that can <u>pursue complex goals</u> and directly take actions in both virtual and real-world environments. Today, people can use Al agents to <u>make payments</u>, <u>reserve flights</u>, and <u>place grocery orders</u> for them, and there is great excitement about the potential for Al agents to manage even more sophisticated tasks.

However, a different type of AI agent—a simulation of human behaviors and attitudes—is also on the rise. These simulation AI agents aim to be useful at asking "what if" questions about how people might respond to a range of social, political, or informational contexts. If these agents achieve high accuracy, they could <u>enable</u> researchers to test a broad set of interventions and theories, such as how people would react to new public health messages, product launches, or major economic or political shocks. Across <u>economics</u>, <u>sociology</u>, <u>organizations</u>, and <u>political science</u>, new ways of simulating individual behavior—and the behavior of groups of individuals— could help expand our understanding of social interactions, institutions, and networks. While work on these kinds of agents is progressing, current architectures must cover some distance before their use is reliable.

Key Takeaways

Simulating human attitudes and behaviors could enable researchers to test interventions and theories and gain real-world insights.

.....

We built an AI agent architecture that can simulate real people in ways far more complex than traditional approaches. Using this architecture, we created generative agents that simulate 1,000 individuals, each using an LLM paired with an in-depth interview transcript of the individual.

To test these generative agents, we evaluated the agents' responses against the corresponding person's responses to major social science surveys and experiments. We found that the agents replicated real participants' responses 85% as accurately as the individuals replicated their own answers two weeks later on the General Social Survey.

.....

Because these generative agents hold sensitive data and can mimic individual behavior, policymakers and researchers must work together to ensure that appropriate monitoring and consent mechanisms are used to help mitigate risks while also harnessing potential benefits.

Policy Brief Simulating Human Behavior with AI Agents



In our paper, "<u>Generative Agent Simulations of 1,000</u> <u>People</u>," we introduce an AI agent architecture that simulates more than 1,000 real people. The agent architecture—built by combining the transcripts of two-hour, qualitative interviews with a large language model (LLM) and scored against social science benchmarks—successfully replicated real individuals' responses to survey questions 85% as accurately as participants replicate their own answers across surveys staggered two weeks apart. The generative agents performed comparably in predicting people's personality traits and experiment outcomes and were less biased than previously used simulation tools.

This architecture underscores the benefits of using generative agents as a research tool to glean new insights into real-world individual behavior. However, researchers and policymakers must also mitigate the risks of using generative agents in such contexts, including harms related to over-reliance on agents, privacy, and reputation.

Introduction

Simulations in which agents are used to model the behaviors and interactions of individuals have been a popular tool for <u>empirical social research</u> for years, even before the emergence of AI agents. Traditional approaches to building agent architectures, such as <u>agent-based models</u> or <u>game theory</u>, rely on clear sets of rules and environments <u>manually specified</u> by the researchers. While these rules make it relatively easy to interpret results, they also limit the contexts in which traditional agents can act while oversimplifying the real-life complexity of human behavior. This, in Generative AI models offer the opportunity to build general purpose agents that can simulate human attitudes across a variety of contexts.

turn, can limit the generalizability and accuracy of the simulation results.

Generative AI models offer the opportunity to build general purpose agents that can simulate human attitudes across a variety of contexts. To create simulations that better reflect the myriad, often idiosyncratic factors that influence individuals' attitudes, beliefs, and behaviors, we built a novel generative agent architecture that combines LLMs with in-depth interviews with real individuals.

We recruited 1,052 individuals—representative of the U.S. population across age, gender, race, region, education, and political ideology—to participate in two-hour qualitative interviews. These in-depth interviews, which <u>included</u> both pre-specified questions and <u>adaptive</u> follow-up questions, are a foundational social science method that has been successfully used by researchers to <u>predict</u> life outcomes beyond what could be learned from traditional surveys and demographic instruments. We also developed an AI interviewer to ask participants the questions based on a semi-structured interview



Policy Brief Simulating Human Behavior with Al Agents

protocol from the <u>American Voices Project</u>—which ranged from life stories to people's views on current social issues.

Then, we built the <u>generative agents</u> based on participants' full interview transcripts and an LLM. When a generative agent was queried, the full transcript was injected into the model prompt, which instructed the model to imitate the relevant individual when responding to questions, including forced-choice prompts, surveys, and multi-stage interactional settings.

Once the generative agents were in place, we evaluated them on their ability to predict participants' responses to common social science surveys and experiments, which the participants completed after their in-depth interviews. We tested on the core module of the <u>General Social Survey</u> (widely used to assess survey respondents' demographic backgrounds, behaviors, attitudes, and beliefs); the 44-item <u>Big</u> <u>Five Inventory</u> (designed to assess an individual's personality); five <u>well-known behavioral economic</u> games (the dictator game, first and second player trust

The generative agents proved remarkably effective in simulating individuals' real-world personalities. games, public goods game, and prisoner's dilemma); and five <u>social science experiments</u> with <u>control</u> and <u>treatment</u> conditions. For the General Social Survey (which has categorical responses), we measured accuracy and correlation based on whether the agent selects the same survey response as the person. For the Big Five Inventory and the economic games (which have continuous responses), we assessed accuracy and correlation using mean absolute error.

Research Outcomes

Overall, the generative agents proved remarkably effective in simulating individuals' real-world personalities. For example, the generative agents predicted participants' responses to the General Social Survey with an average normalized accuracy of 85%—meaning that, on average, generative agents replicated participant responses 85% as accurately as the participants themselves when they were asked to retake the surveys and experiments two weeks later. This result is 14 to 15 percentage points higher than the accuracy of traditional demographic-based and persona-based agents that use the same LLMs but do not have access to the interviews.

The generative agents also outperformed demographic and persona-based agents on the Big Five personality test, achieving a normalized correlation of 80% when replicating real individuals' openness, conscientiousness, extraversion, agreeableness, and neuroticism. But they performed similarly as demographic and persona-based agents for the economic games, with a normalized correlation of 66% (i.e., 66% as high as the participants' own correlation with themselves two weeks later) across an

Policy Brief Simulating Human Behavior with AI Agents



...the interview-based generative agents consistently reduced biases across tasks compared to demographic-based agents. while gender-based Demographic Parity Difference remained fairly consistent across tasks (likely due to already low levels of discrepancy).

aggregate of the dictator game, the first and second player trust games, the public goods game, and the prisoner's dilemma.

Beyond those tests, we evaluated the generative agents' behavior in a set of social science experiments. These included investigations of how perceived intent <u>affects</u> blame assignment and how fairness <u>influences</u> emotional responses. Real-world participants and the generative agents agreed on the replication results of all five studies we tested.

The generative agents also lessened bias in predictive accuracy across social groups. Given rightful concerns about AI systems disadvantaging or <u>misrepresenting</u> underrepresented populations, we conducted a subgroup analysis focused on political ideology, race, and gender. These are <u>dimensions</u> <u>of particular interest</u> in the <u>literature</u>. We used the Demographic Parity Difference, which <u>measures</u> the performance <u>difference</u> between the best- and worst-performing groups, to quantify bias. Notably, we found that the interview-based generative agents consistently reduced biases across tasks compared to demographic-based agents. Drops in political ideology bias and racial bias vary depending on the survey,

Policy Discussion

Generative agents could become useful tools for estimating attitudes and survey-based experimental treatment effects. For example, if you were considering the sorts of survey questions you might ask in a national survey, generative agents could help to estimate average responses the population might give. However, many open questions remain: How accurate are generative agents when simulating behavior, in addition to attitudes? What innovations are needed for generative agent simulations to accurately estimate the impacts of policy changes? While we will continue to build the empirical and technical research to expand the horizon of generative agents, we urge policymakers to critically examine analyses that overclaim what generative agents can actually achieve today.

One important risk for policymakers, practitioners, researchers, and others using generative agents is <u>overreliance</u> on generative agents when simulation accuracy is low. To ensure that policymakers don't rely on an inaccurate simulation, we must develop tools and methodologies so they know when they can, and can't, trust these simulations. Additionally, policymakers should not apply generative agents beyond the range of applications that have been validated. A second major risk relates to privacy: The interview

Policy Brief Simulating Human Behavior with AI Agents



data used to build the generative agents is often sensitive, and data leaks could cause considerable harm to interviewees. Other concerns include the co-option of individuals' likenesses, as these agents can believably replicate a person's answers in a survey response or experiment. Significant reputational harm could also result from someone manipulating agent responses to falsely attribute defamatory statements to individuals whose data is used in the agent bank.

A range of other ethical and legal questions must also be considered. For example, what are the ethical implications of using AI agents that <u>simulate</u> <u>a deceased person</u>? How should human consent be managed? And what are the risks of agents being misused for fraudulent purposes? Given the inherent uncertainty of future advancements in generative AI, such as AI models' future reasoning abilities, managing these risks early on is crucial. Policymakers should consider establishing bright-line rules that determine how AI agents may or may not be used for human simulation purposes.

We made the decision not to release our generative agents for public use. Instead, to support further research while protecting participant privacy, we have chosen to provide controlled, research-only API access to our agent bank. We grant open access to aggregated responses on fixed tasks for general research use and restrict access to individual responses on open tasks for researchers following a review process, ensuring the agents are accessible while minimizing risks associated with the source interviews. Other researchers building similar systems should replicate our safeguards, and policymakers weighing how generative agents could be used in research settings should explore requirements for individual data rights, access, and deletion. One important risk for policymakers, practitioners, researchers, and others using generative agents is overreliance on generative agents when simulation accuracy is low.

Policymakers and researchers should work together to ensure that appropriate monitoring and consent mechanisms are used to enhance trust, protect individual rights, and mitigate the risks of generative agent use. For example, our team proposed the possibility of an audit log for the use of every agent in our agent bank. That way, individuals who participated in a survey and had their preferences captured by a generative agent could see what the agent is doing and exert control over it over time. Permission could be granted one day and withdrawn a month later, reflecting individual consent. Translating such protections into policy-such as making them part of grant terms and conditions-would help researchers to detect and mitigate malicious use of generative agents built using people's personal data shared via in-depth interviews.

Looking forward, generative agents hold serious promise for enhancing human behavioral research and developing new insights into personal preferences and decision making. However, mitigating the risks of these innovations, through research and policy controls on agent access and auditing, will be crucial to harnessing their opportunities in economics, political science, and beyond. Reference: The original article is accessible at Joon Sung Park, Carolyn Q. Zou, et al., "Generative Agent Simulations of 1,000 People," arxiv.org, November 15, 2024, https://arxiv.org/ abs/2411.10109.

Stanford University's Institute for Human-Centered Artificial Intelligence (HAI) applies rigorous analysis and research to pressing policy questions on artificial intelligence. A pillar of HAI is to inform policymakers, industry leaders, and civil society by disseminating scholarship to a wide audience. HAI is a nonpartisan research institute, representing a range of voices. The views expressed in this policy brief reflect the views of the authors. For further information, please contact HAI-Policy@stanford.edu.



Joon Sung Park is a

PhD student in computer science in the Human-**Computer Interaction** and Natural Language Processing groups at Stanford University.



is director and principal scientist for Human-Al Interaction at Google DeepMind.

Meredith Ringel Morris



Carolyn Q. Zou is

a PhD student in computer science at Stanford University.



Aaron Shaw is an associate professor in the Department of Communication Studies at Northwestern University.

Benjamin Mako Hill

is an associate professor in the Department of Communication at the University of Washington.



Carrie J. Cai is a senior staff research scientist at Google DeepMind and manager/ area lead of Human-Al Interaction in Google's People+AI Research group.



Robb Willer is a professor of sociology and, by courtesy, psychology and business at Stanford University.

Percy Liang is an associate professor of computer science at Stanford University and a senior fellow at Stanford HAI.



associate professor of computer science at Stanford University and a senior fellow



Stanford University Human-Centered **Artificial Intelligence**

Stanford HAI: 353 Jane Stanford Way, Stanford CA 94305-5008 T 650.725.4537 F 650.123.4567 E HAI-Policy@stanford.edu hai.stanford.edu