# LLM-Enhanced Survey Methodology: Validation, Automation, and Mixed Methods at Scale

**Anonymous authors**
Paper under double-blind review

## Abstract

This work explores the transformative potential of large language models (LLMs) across three key domains in survey research: digital twin simulations, AI-driven telephone interviewing, and mixed-methods data collection at scale. By leveraging LLMs to generate synthetic responses mirroring diverse global populations, we assess model validity in comparison to annual nationally representative survey data from over 140 countries. Trials of LLM-powered phone interviews in the U.S. demonstrate that AI interviewers can achieve comparable clarity and lower social conformity bias than human interviewers, while highlighting the need for improved conversational nuance and participant consent. Finally, we integrate LLMs into mixed-methods frameworks, using AI-driven conversational probes to deepen qualitative insights while maintaining methodological rigor. Our findings illustrate the promise of LLMs in enhancing the efficiency, scalability, and cultural sensitivity of survey research, alongside ongoing challenges related to bias, methodological replication and implementation, and technical limitations.

## 1 Digital Twins: Testing for Bias and Model Validity

One of the most promising avenues involves using LLMs to simulate "digital twins" - personas representing various cultural, demographic, and national backgrounds (Bisbee et al., 2024; Zewail et al., 2024). By prompting LLMs to respond as if they were individuals from diverse populations, we are assessing whether synthetic responses align with annual nationally representative survey data from over 140 countries, representing over 95% of the world's population. Comparing across ground truth samples on well-being, life evaluation, economic indicators, perceptions of governance, social issues, education, and health, this effort seeks to identify and detect biases and evaluate the extent to which models reliably mirror human perspective. Additionally, this framework is being instantiated to nearly continuously evaluate the evolving model development pace.

Initial findings indicate substantial variability across models, topics, and cultural contexts. These discrepancies underscore the importance of ongoing validation, especially if LLMs are to be trusted in methodological roles such as generating unbiased survey questions or modeling respondent behavior.

## 2 AI-Driven Phone Interviewing: A Viable Alternative?

Steaming from initial research from Lang & Eskenazi (2025), we have conducted promising United States based trials using LLM-powered phone structured interviews. Among 4,000 individuals invited to participate, 863 consented, with 447 contacted and 345 completing the survey. A follow-up evaluation showed 90% of participants found the AI interviewer to be clear, unbiased, and nonjudgmental. However, 47% described the experience as too scripted and lacking in empathy, while 82% reported no technical problems—the most common issue being misunderstanding of questions. These findings suggest that AI can approximate the quality of trained human interviewers, with the added benefits of cost savings and scalability. Importantly, social conformity bias appears to be lower with AI than in human-

led interviews. Still, the need for prior consent and improved conversational nuance remains, particularly as we extend this research into countries with different regulatory environments. We are currently focusing on understanding mode differences, response quality, and establishing a continuous evaluation framework to understand model evolution and validity.

## 3 Mixed Methods at Scale: LLMs as Conversational Probes

The third area focuses on integrating LLMs into a mixed-methods approach that blends quantitative structure with qualitative depth. By using initial survey responses to inform LLM-driven conversational probes, we aim to "dig deeper" into participant perspectives. This approach not only enhances data richness but also represents a step toward a more respondent-centered survey experience. It leverages the conversational strengths of LLMs to mimic the fluidity of human inquiry while maintaining methodological rigor. To support this, we're developing an evaluation framework that connects multiple models in a graph-based network. This system helps align participant and LLM responses, ensures meaningful engagement for both, and keeps the conversation informed by dynamically updating its memory. While prior research has begun to explore this area (Wuttke et al., 2025), our current efforts aim to redefine data collection, push the boundaries of survey science, and set new standards for participant engagement.

## 4 Conclusion

These three lines of research demonstrate the growing potential of LLMs to transform social science research. From simulating populations to automating interviews and enriching mixed methods, our work is paving the way for more scalable, nuanced, and culturally sensitive survey methodologies. However, each domain presents distinct challenges—ranging from bias in model outputs and limited emotional intelligence to technical barriers—that must be addressed to ensure LLMs are developed and deployed as ethical and effective research and data collection tools.

## References

James Bisbee, Joshua D. Clinton, Cassy Dorff, Brenton Kenkel, and Jennifer M. Larson. Synthetic Replacements for Human Survey Data? The Perils of Large Language Models. *Political Analysis*, 32(4):401–416, 2024. ISSN 1047-1987, 1476-4989. doi: 10.1017/pan.2024.5. URL https://www.cambridge.org/core/product/identifier/S1047198724000056/type/journal_article.

Max M Lang and Sol Eskenazi. Telephone Surveys Meet Conversational AI: Evaluating a LLM-Based Telephone Survey System at Scale. 2025.

Alexander Wuttke, Matthias Aßenmacher, Christopher Klamm, Max M. Lang, Quirin Würschinger, and Frauke Kreuter. AI Conversational Interviewing: Transforming Surveys with LLMs as Adaptive Interviewers. 2025. doi: 10.48550/arXiv.2410.01824. URL http://arxiv.org/abs/2410.01824.

Aliah Zewail, Alexandra Figueroa, Jesse Graham, and Mohammad Atari. Moral Stereotyping in Large Language Models. 2024. doi: 10.31234/osf.io/t9x8r. URL https://osf.io/t9x8r.