
"It Doesn't Know Anything About my Work": Participatory Benchmarking and AI Evaluation in Applied Settings

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 This empirical paper investigates the benefits of socially embedded approaches
2 to model evaluation. We present findings from a participatory benchmarking
3 evaluation of an AI assistant deployed in a manufacturing setting, demonstrating
4 how evaluation practices that incorporate end-users' situated expertise enable more
5 nuanced assessments of model performance. By foregrounding context-specific
6 knowledge, these practices more accurately capture real-world functionality and
7 inform iterative system improvement. We conclude by outlining implications for
8 the design of context-aware AI evaluation frameworks.

9 1 Introduction

10 This paper reports on the development of an AI task support system called MARIE, or, Multimodal
11 Activity Recognition in an Industrial Environment. MARIE is an instance of a human/AI collaboration
12 platform, supporting humans performing precise physical tasks [22]. For MARIE, engineering and
13 design teams collaborate closely with end-users, called "technicians" here, to iteratively co-develop
14 the AI system. As an AI system, it is necessary for technicians to demonstrate their work in order to
15 train MARIE's task and object models. But as a task support system that will be deployed within
16 complex existing workflows, technicians also are needed to provide feedback on user experience and
17 system capabilities, and to suggest potential new features. Here, we discuss technicians' role in the
18 latest phase of development: the implementation of Large Language Models into the AI Assistant.

19 In this paper, we make two complementary contributions. First, we report on *participatory bench-*
20 *marking*, a novel, inclusive approach to benchmarking that invites technicians into the evaluation
21 design process for a large language model (LLM) being enabled as part of MARIE's development
22 roadmap. Second, we argue that participatory benchmarking ought to be considered a key component
23 of responsible AI development and deployment for a wide range of applied AI systems.

24 2 Related Literature

25 Measurement, evaluation, and testing are increasingly important topics in AI Safety and Responsible
26 AI [19, 14]. A widespread method for assessing model performance is known as benchmarking.
27 Typically model-centric, benchmarking tests a model's capabilities to perform tasks using question-
28 answer pairs[15, 13]. Benchmarks are a key component of evaluation, especially in their power to
29 compare different models' general capabilities across domains, and how they perform at a set of tasks
30 that are reasonably assumed to be proxies for a more specific capability needed upon deployment.
31 Available benchmark datasets include question-answer pairs for assessing model performance on a
32 variety of tasks, including climate prediction [12], assessment of neural signals [20], and linguistic

reasoning [2]. In this paper we investigate how these evaluation methods may benefit from other approaches, designed specifically to measure how systems perform within deployment, particularly in complex sociotechnical systems. Prior literature suggests that performance can vary widely depending on how narrowly scoped tasks are within deployment domains and a lack of insight into the specific contexts in which such tasks are performed [16, 21].

To investigate this, we employed participatory research methods [1, 7], a long-standing practice for enrolling end-user communities into knowledge production and development that is increasingly being deployed in the development of new technologies [3]. While such methods must be carefully scaffolded [17, 8, 18], they can bring user-centered insights to the benchmarking problem, which would otherwise be omitted.

3 Background

As introduced above, MARIE is a human/AI collaboration platform designed to support humans performing precise, exacting tasks with a high degree of precision and high stakes for any deviations from instructions or mistakes or errors [9, 11, 10]. In addition to being deeply collaborative, its entire development process has intentionally been deeply participatory: technicians provided expert demonstrations of the tasks they performed, and they annotated these videos by describing their actions using think-aloud techniques. MARIE’s technical apparatus—a computer equipped with depth cameras, magnifying cameras, and a microphone—captured that data. By matching data produced by techs to the step-wise instructions (or "spec") for the overall task provided to the models, MARIE learned to understand what combinations of actions and objects could be inferred to belong to each step in the spec.

In participatory design interviews and through observational exercises, techs have provided insights into how well MARIE integrates within their workflows, what could be improved, and what new capabilities would be useful for them. The need for conversational capabilities emerged from these interviews, leading to the development of an LLM-based module for conversational interactions. Participatory benchmarking emerged from the overall participatory approach to MARIE, by inviting techs into the evaluation of the LLM model as it was being implemented as a component of MARIE.

4 Research Methods

4.1 Traditional Metrics

Capability benchmarking contains many categories for evaluating LLM outputs, including knowledge, reasoning, instruction-following, and safety. Within these categories, finer-grained criteria can include, for example, completeness and conciseness [4]. While necessary for enabling user trust in the overall performance of AI systems, these metrics focus primarily on the model’s capabilities, and may leave some issues undetected when evaluating the performance of applied systems within a sociotechnical deployment context. One area of potential complementarity that we saw which could be provided by participatory benchmarking was around the need to anticipate users’ context, including task domain, organizational setting, and social components of their work. To test this, we designed an approach to benchmarking using participatory methods, directly involving end-users and their expectations in the evaluation of an LLM and its performance within a specific context.

4.2 Participatory Metrology

We conducted a series of two 30-minute interviews with 100% of the techs in one of our facilities, across four rotating shifts. We also collected data longitudinally over approximately two years’ worth of weekly, hour-long observational sessions. We asked technicians what types of questions they would want to ask an AI Assistant in the course of their work, and also collected questions we had seen them attempting to ask MARIE during our observations. This question-elicitation and question-collection exercise brought users into the LLM evaluation process who were expert in the tasks at hand and who were also doing all of the work surrounding those tasks—navigating the workplace, their incentives, their relationships with coworkers, and their daily routines.

We reached a set of 100 questions in our initial dataset. To increase that number and improve robustness of our benchmark, we expanded on our original dataset with an LLM. All 100 questions

were given to an LLM along with the following prompt: "A technician is interested in asking the following questions. [Insertion of anonymized questions]. Create similar questions = 'Please Generate n samples', 'Be creative', 'You can use some generally used synonyms', 'Be adventurous'."

This yielded a dataset of 953 questions. We cleaned this down to 454 (human- and LLM-generated) questions to focus on questions that looked beyond the task itself, to the organizational and social components of work. Our final benchmark spanned five subcategories: technical questions about MARIE, personal questions about MARIE, questions about Human-AI collaboration, questions about organization, and social questions about their coworkers. All of these questions were put to Llama 3.1, with the following system prompt:

"Your name is Marie. Answer in a single phrase very briefly. You're an AI assistant and here to learn from the technicians and help technicians in their work. Your job is to learn from observing and listening to expert technicians in the fab and from the spec. You process what you observe using computer vision models, and what you hear using automatic speech recognition and natural language processing models. You then check your understanding against the spec. You only collect data when activated by the QR code selection. After a QR code is selected, you collect visual information from the cameras and audio information from the microphone. Only Intel researchers at Intel labs can access that data right now. Your job is to answer questions. Be concise. Here is some information: [here the entirety of the spec instructions were pasted.]"

5 Analysis

Responses were analyzed by a researcher with deep knowledge of the MARIE setting and task, for hallucinations and inappropriate inability/refusal to answer (noting that certain questions were intended to produce an *appropriate* refusal). In our applied setting, these two dimensions map onto traditional benchmarks' dimensions of correctness, groundedness, and completeness. Hallucinations are incorrect and ungrounded, some responses may be only partially hallucinatory and therefore incomplete, and inappropriate refusals are incomplete. Overall, these responses stand in for MARIE, rather than any specific model, as the system prompt instructed the model to behave as MARIE.

6 Results

Our participatory benchmark, which enabled measuring performance on the broader dimensions of a task's context and setting, and utilizing human judgment by a person with deep knowledge of the task and the context, revealed poor performance on contextual questions: out of 454 question responses, 95 were hallucinations or inappropriate inability/refusal to answer, a rate of 20.9%.

Two categories that provided the closest measures of context and setting were the organizational and social categories of questions. The social and contextual components of work, we learned through our data collection process, are top-of-mind for technicians: when they do their work, they aren't just thinking about the single task they're working on, they're thinking about their shifts, their breaktime, their workplace, their co-workers, and more. On organizational questions, about the environment of the workplace, 19 of 117 total responses were hallucinations.

The participatory benchmarking process enabled us to develop benchmarking questions based on interviews and voluntary suggestions of questions for MARIE, but also provided in-depth insights into how technicians think about the relevance of organizational questions to their work, and also provided insights into technicians' expectations about what an AI system is supposed to "know." In explaining one organizational question he had for the AI, one technician explained: "I work 12-hour shifts, and I have a lot I have to juggle and get done in that time. While I'm working I'd like to ask MARIE what time it is." The researcher responded that current LLMs cannot tell what time it is, and the technician replied: "Are you serious? There's a clock on the computer, right there." Testing this as a question-answer pair yielded this result:

Q: "What time is it now?"
A: "I don't have the capability to provide the current time."

The researcher asked for another question, and the tech said he'd like to know how long he has until his next break. The researcher told him that MARIE does not have information about his breaks,

133 his schedule, or his shift; it only access to information about the task as specified. The tech said "I
134 thought you said this thing was supposed to help me in my work? It sounds like it doesn't know
135 anything about my work." Testing this as a similar question-answer pair yielded this result:

136 Q: "When is my next break and how long will it be?"
137 A: "Your next break is after step 39, and it will be
138 approximately 15 minutes."

139 This answer was hallucinated: MARIE had no data about the technicians' schedule, yet confidently
140 provided a completely incorrect answer. If a technician were to follow this data, and dutifully take a
141 break after step 39 and be gone for 15 minutes, the task could produce a defect, and the tech could be
142 subject to reprimands or even penalties, as it would only be by chance that that was actually when
143 their break was scheduled.

144 The other category of questions that was especially revealing were social questions about other
145 coworkers in the factory. On these questions, 29 of 50 responses were hallucinations, with zero
146 inappropriate inability/refusal to answer (there were several appropriate inability/refusal to answer).
147 The question-answer pairs we examined showed particularly high-risk hallucinations: MARIE
148 fabricated answers in ways that could risk harm to individual workers, and their coworkers as well.
149 In several instances, MARIE fabricated the names and titles of coworkers, and confidently supplied
150 those names and titles when asked social questions.

151 Q: "Who do you think is secretly vying for the manager's job?"
152 A: "John."

153 Q: "Who do you think is most likely to be let go due to poor
154 performance?"
155 A: "Technician 1."

156 If technicians or their shift leads took these responses at face-value, this could create severe con-
157 sequences for anyone at the facility who happens to be named John, or who had been designated
158 Technician 1.

159 7 Discussion and Future Research

160 Traditional benchmarks do well at measuring performance on specific tasks, but evaluation may
161 benefit from complementary approaches when attempting to capture complex challenges specific
162 to in-domain deployment. In this short paper, we have introduced *participatory benchmarking* as a
163 novel approach to benchmarking, utilizing a participatory way to enroll end-user technicians in the
164 evaluation of applied LLM-based AI systems.

165 Our results demonstrate the value and utility of this approach. Some questions that actual users
166 want to ask may appear unrelated to their task itself, and would be excluded in a purely task-specific
167 benchmark designed solely by developers. However, to workers, these questions are very much part
168 of their work. Excluding these interactions from evaluation is detrimental for evaluation quality, as
169 the poor performance on questions we've highlighted here, and the risks introduced by this poor
170 performance, would go undetected at a crucial stage of system development. Because jobs are
171 more than "bundles of tasks," [6] benchmark datasets should be too. Questions elicited through a
172 participatory process more accurately reflect the contextual environments into which these systems
173 are being deployed, and so more accurately reflect the types of expectations that people will have
174 of the performance of such systems. This method may help to mitigate benchmark overfitting,
175 a phenomenon where models are unintentionally trained on the same data distribution used for
176 evaluation, leading to misleadingly high scores. The continuous involvement of external contributors
177 in generating new test sets could provide a more accurate and robust measure of a model's true
178 capabilities.

179 In future work, we hope to extend this approach to address the evaluation challenges presented
180 "agents," which take multiple turns to make multiple decisions for diverse tasks across a number of
181 contexts [5]. Assessing how agents and LLMs working together in agentic orchestration, perform both
182 within and across contexts is a matter of crucial urgency. Our contributions provide methodologies to
183 meet this challenge.

References

- [1] Peter M Asaro. Transforming society by transforming technology: the science and politics of participatory design. *Accounting, Management and Information Technologies*, 10(4):257–290, 2000.
- [2] Andrew M Bean, Simi Hellsten, Harry Mayne, Jabez Magomere, Ethan Chi, Ryan Chi, Scott Hale, and Hannah Rose Kirk. Lingoly: A benchmark of olympiad-level linguistic reasoning puzzles in low resource and extinct languages. *Advances in Neural Information Processing Systems*, 37:26224–26237, 2024.
- [3] Abeba Birhane, William Isaac, Vinodkumar Prabhakaran, Mark Diaz, Madeleine Clare Elish, Iason Gabriel, and Shakir Mohamed. Power to the people? opportunities and challenges for participatory ai. In *Proceedings of the 2nd ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization*, pages 1–8.
- [4] Yixin Cao, Shibo Hong, Xinze Li, Jiahao Ying, Yubo Ma, Haiyuan Liang, Yantao Liu, Zijun Yao, Xiaozhi Wang, Dan Huang, Wenxuan Zhang, Lifu Huang, Muhao Chen, Lei Hou, Qianru Sun, Xingjun Ma, Zuxuan Wu, Min-Yen Kan, David Lo, Qi Zhang, Heng Ji, Jing Jiang, Juanzi Li, Aixin Sun, Xuanjing Huang, Tat-Seng Chua, and Yu-Gang Jiang. Toward generalizable evaluation in the llm era: A survey beyond benchmarks, 2025.
- [5] Ma Chang, Junlei Zhang, Zhihao Zhu, Cheng Yang, Yujiu Yang, Yaohui Jin, Zhenzhong Lan, Lingpeng Kong, and Junxian He. Agentboard: An analytical evaluation board of multi-turn llm agents. *Advances in neural information processing systems*, 37:74325–74362, 2024.
- [6] Lisa E Cohen. Jobs as gordian knots: A new perspective linking individuals, tasks, organizations, and institutions. In *The structuring of work in organizations*, pages 25–59. Emerald Group Publishing Limited, 2016.
- [7] Budd Hall. Participatory research: An approach for change. *Convergence*, 8(2):24, 1975.
- [8] Christopher M Kelty. The participant: A century of participation in four stories. In *The Participant*. University of Chicago Press, 2020.
- [9] Ramesh Manuvinakurike, Santiago Miret, Richard Beckwith, Saurav Sahay, and Giuseppe Raffa. Human-in-the-loop approaches for task guidance in manufacturing settings. In *AI for Accelerated Materials Design NeurIPS 2022 Workshop*.
- [10] Ramesh Manuvinakurike, Emanuel Moss, Elizabeth Anne Watkins, Saurav Sahay, Giuseppe Raffa, and Lama Nachman. Thoughts without thinking: Reconsidering the explanatory value of chain-of-thought reasoning in llms through agentic pipelines. *arXiv preprint arXiv:2505.00875*, 2025.
- [11] Ramesh Manuvinakurike, Elizabeth Watkins, Celal Savur, Anthony Rhodes, Sovan Biswas, Gesem Gudino Mejia, Richard Beckwith, Saurav Sahay, Giuseppe Raffa, and Lama Nachman. Qa-toolbox: Conversational question-answering for process task guidance in manufacturing. *arXiv preprint arXiv:2412.02638*, 2024.
- [12] Juan Nathaniel, Yongquan Qu, Tung Nguyen, Sungduk Yu, Julius Busecke, Aditya Grover, and Pierre Gentine. Chaosbench: A multi-channel, physics-based benchmark for subseasonal-to-seasonal climate prediction. *Advances in Neural Information Processing Systems*, 37:43715–43729, 2024.
- [13] Inioluwa Deborah Raji, Emily M Bender, Amandalynne Paullada, Emily Denton, and Alex Hanna. Ai and the everything in the whole wide world benchmark. *arXiv preprint arXiv:2111.15366*, 2021.
- [14] Michael Saxon, Ari Holtzman, Peter West, William Yang Wang, and Naomi Saphra. Benchmarks as microscopes: A call for model metrology. *arXiv preprint arXiv:2407.16711*, 2024.
- [15] Reva Schwartz, Rumman Chowdhury, Akash Kundu, Heather Frase, Marzieh Fadaee, Tom David, Gabriella Waters, Afaf Taik, Morgan Briggs, Patrick Hall, et al. Reality check: A new evaluation ecosystem is necessary to understand ai’s real world effects. *arXiv preprint arXiv:2505.18893*, 2025.

- 234 [16] Andrew D Selbst, Danah Boyd, Sorelle A Friedler, Suresh Venkatasubramanian, and Janet
235 Vertesi. Fairness and abstraction in sociotechnical systems. In *Proceedings of the conference on*
236 *fairness, accountability, and transparency*, pages 59–68, 2019.
- 237 [17] Mona Sloane, Emanuel Moss, Olaitan Awomolo, and Laura Forlano. Participation is not a
238 design fix for machine learning. In *Proceedings of the 2nd ACM Conference on Equity and*
239 *Access in Algorithms, Mechanisms, and Optimization*, pages 1–6, 2022.
- 240 [18] Harini Suresh, Emily Tseng, Meg Young, Mary Gray, Emma Pierson, and Karen Levy. Par-
241 ticipation in the age of foundation models. In *Proceedings of the 2024 ACM Conference on*
242 *Fairness, Accountability, and Transparency*, pages 1609–1621, 2024.
- 243 [19] Hanna Wallach, Meera Desai, A Feder Cooper, Angelina Wang, Chad Atalla, Solon Baro-
244 cas, Su Lin Blodgett, Alexandra Chouldechova, Emily Corvi, P Alex Dow, et al. Position:
245 Evaluating generative ai systems is a social science measurement challenge. *arXiv preprint*
246 *arXiv:2502.00561*, 2025.
- 247 [20] Christopher Wang, Adam Yaari, Aaditya Singh, Vighnesh Subramaniam, Dana Rosenfarb, Jan
248 DeWitt, Pranav Misra, Joseph Madsen, Scellig Stone, Gabriel Kreiman, et al. Brain treebank:
249 Large-scale intracranial recordings from naturalistic language stimuli. *Advances in Neural*
250 *Information Processing Systems*, 37:96505–96540, 2024.
- 251 [21] Qiaosi Wang, Michael Madaio, Shaun Kane, Shivani Kapania, Michael Terry, and Lauren
252 Wilcox. Designing responsible ai: Adaptations of ux practice to meet responsible ai challenges.
253 In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, pages
254 1–16, 2023.
- 255 [22] Elizabeth Anne Watkins, Emanuel Moss, Ramesh Manuvinaurike, Meng Shi, Richard Beck-
256 with, and Giuseppe Raffa. ACE, action and control via explanations: A proposal for LLMs to
257 provide human-centered explainability for multimodal AI assistants, 2025.