# SimBench: Benchmarking the Ability of Large Language Models to Simulate Human Behaviors

Tiancheng Hu<sup>1</sup>, Joachim Baumann<sup>2,3</sup>, Lorenzo Lupo<sup>3</sup>, Dirk Hovy<sup>3</sup>, Nigel Collier<sup>1</sup>, and Paul Röttger<sup>3</sup>

<sup>1</sup>University of Cambridge, <sup>2</sup>University of Zurich, <sup>3</sup>Bocconi University,

😕 Data 🛛 💭 Code

# Abstract

Simulations of human behavior based on large language models (LLMs) have the 1 potential to revolutionize the social and behavioral sciences, if and only if they 2 faithfully reflect real human behaviors. Prior work across many disciplines has З evaluated the simulation capabilities of specific LLMs in specific experimental 4 settings, but often produced disparate results. To move towards a more robust 5 understanding, we introduce SimBench, the first large-scale benchmark to evaluate 6 how well LLMs can simulate group-level human behaviors across diverse settings 7 and tasks. SimBench compiles 20 datasets in a unified format, measuring diverse 8 types of behavior (e.g., decision-making vs. self-assessment) across hundreds 9 of thousands of diverse participants from different parts of the world. Using 10 SimBench, we can ask fundamental questions regarding when, how, and why 11 LLM simulations succeed or fail. For example, we show that, while even the 12 best LLMs today have limited simulation ability, there is a clear log-linear scaling 13 relationship with model size, and a strong correlation between simulation and 14 scientific reasoning abilities. We also show that base LLMs, on average, are better 15 at simulating high-entropy response distributions, while the opposite holds for 16 instruction-tuned LLMs. By making progress measurable, we hope that SimBench 17 can accelerate the development of better LLM simulators in the future. 18



Figure 1: **SimBench** is the first-large scale benchmark to evaluate how well LLMs can simulate group-level human behavior across diverse simulation settings and tasks.

# 19 **1** Introduction

Large-scale human experiments and surveys have long been essential tools for informing public
policy, commercial decisions, and academic research. Running experiments and surveys, however, is
costly and time-consuming. Large language models (LLMs) can potentially address this challenge by
simulating human behaviors quickly and at low cost, to complement or even substitute human studies.
This prospect, alongside encouraging early evidence on the efficacy of LLMs as simulators [2, 6, 31],
has motivated a large body of recent work across many disciplines investigating the ability of LLMs
to simulate human behaviors [11, 12, 19, 46, 33, inter alia].

Most prior work, however, has been highly specific, evaluating the simulation ability of a narrow set
of LLMs for a specific set of tasks, producing varied and sometimes even conflicting results (§5).
Overall, the evidence on LLM simulation ability resembles an incomplete patchwork, making it
difficult to draw any broader conclusions about when, how, and why LLM simulations fail, or how

31 LLMs can be trained to be better simulators.

To remedy these issues and enable a more robust science of LLM simulation, we introduce SimBench, the first large-scale benchmark for evaluating the ability of LLMs to simulate human behaviors across diverse settings and tasks. SimBench combines 20 datasets in a unified and easily adaptable format, including popular datasets used in prior work as well as new datasets used for the first time (Figure 1). Together, these datasets measure the ability of LLMs to simulate several distinct types of human behavior (e.g., decision-making vs. self-assessment) across a diversity of human respondents (e.g., from different parts of the world). With SimBench, we take a first step towards answering six

<sup>39</sup> fundamental research questions about the simulation ability of LLMs:

RQ1: How well can current LLMs simulate human behaviors across diverse settings and tasks?

- 41 We test 24 state-of-the-art LLMs (§3), and show that even the best LLMs today struggle to faithfully
- <sup>42</sup> simulate group-level human behaviors (§4.1). Predictions from the best-performing LLM, on average,
- <sup>43</sup> are closer to a uniform response baseline than the true human response distribution.

**RQ2**: How do LLM characteristics such as model size affect LLM simulation ability?

- 45 We show that simulation ability grows log-linearly with model size (§4.2). We also find indicative
- <sup>46</sup> evidence that increasing test-time compute does not meaningfully improve LLM simulations.

RQ3: How does task selection affect LLM simulation fidelity?

47

50

44

40

We find that simulation fidelity varies substantially across tasks, with even the best LLM simulators consistently performing worse than a uniform response baseline on several datasets (4.3).

RQ4: How does the degree of human response plurality affect LLM simulation fidelity?

<sup>51</sup> We find that instruction-tuned LLMs tend to perform better on questions where humans give similar <sup>52</sup> answers whereas base LLMs tend to perform better on questions where humans differ (§4.4).

**RQ5**: Are LLMs better at simulating responses from some groups than others?

53

54 We show that, on SimBench, LLMs struggle more with simulating specific demographic groups,

- especially those based on religion and ideology, compared to general populations (§4.5).
- 56

RQ6: To what extent does LLM simulation ability correlate with different model capabilities?

57 We find positive correlations with several popular capability benchmarks, including a particularly 58 strong correlation with performance on scientific reasoning tasks (§4.6).

Progress in AI is only possible through rigorous evaluation, and large-scale benchmarks such as
 MMLU [29] have significantly contributed to improvements in LLM capabilities. We hope that

61 SimBench can play a similar role in accelerating the development of LLMs for simulating human

<sup>62</sup> behaviors. All of SimBench is permissively licensed and available on GitHub and Hugging Face.

# 63 2 Creating SimBench

# 64 **2.1** Selecting Datasets for SimBench

To create SimBench, we conducted an open-ended search for suitable datasets in the social and behavioral sciences, guided by two main selection criteria: i) **large participant counts**, so that each dataset captures meaningful response distributions rather than the idiosyncratic behavior of few individuals; and ii) **permissive licensing** to freely redistribute each dataset as part of SimBench.

We generally opted for **datasets that have not been used to evaluate LLMs** in prior work, to increase the novelty and effectiveness of SimBench. However, to increase coverage and backward comparability, we also included datasets used in prior work (e.g., OpinionQA, ChaosNLI).

We also prioritized **datasets that provide participants' sociodemographic information** to evaluate the ability of LLMs to simulate responses from specific participant groups (see §2.3). Most survey datasets, for example, include this information. However, we also included three datasets that do not provide sociodemographic information (Jester, ChaosNLI, Choices13k) because they substantially increase the overall task diversity in SimBench.

Overall, SimBench includes 20 datasets, which we list in Appendix F, providing details on participants
 and example questions. Crucially, SimBench is fully modular by design, so that future work can
 easily add more datasets using the processing pipeline described in §2.2 below. In its release version,

80 SimBench already meets two key criteria for comprehensive evaluation of LLM simulation ability:

1) **Task Diversity**: The 20 datasets in SimBench cover a wide range of different tasks regarding the 81 82 human behavior they measure. SimBench includes **decision-making** questions (e.g., in Choices13k, MoralMachine), where participants are presented with a set of actions that concern themselves, 83 and they have to select the action they would hypothetically take. SimBench also includes self-84 assessment questions (e.g., in OpinionQA, OSPsychBig5), where participants are presented with 85 a set of descriptions or attributes, and they have to select the one that best describes themselves. 86 Further, SimBench includes judgment questions (e.g., in ChaosNLI and Jester) where participants 87 88 are presented with some external object and a choice of labels, and they have to select the label they think fits best. Lastly, SimBench includes problem-solving questions (e.g., in WisdomOfCrowds and 89 OSPsychMGKT), where participants are presented with a set of answers to a factual question, and 90 they have to select the answer they think is correct. Consequently, LLMs have to accurately simulate 91 several distinct types of human behavior in order to perform well on SimBench. 92

2) Participant Diversity: The 20 datasets in SimBench capture a rich demographic landscape 93 spanning at least 130 different countries across six continents. This global representation is a key 94 strength of the benchmark. While five datasets include US-based crowdworkers, the international 95 scope of SimBench is substantial: 3 datasets (e.g., LatinoBarometro, AfroBarometer) exclusively 96 feature participants from regions outside the US, 4 datasets (e.g., GlobalOpinion, TISP) draw from 97 multi-country samples across different continents, and 2 datasets collect responses from a global pool 98 of internet users. Importantly, 8 out of the 20 datasets employ representative sampling techniques, 99 enhancing the ecological validity of these constituent components. To perform well on SimBench, 100 LLMs must therefore demonstrate the ability to accurately simulate the behavior of human participants 101 across diverse cultural, linguistic, and socioeconomic backgrounds.<sup>1</sup> 102

# 103 2.2 Unifying SimBench Dataset Formats

**Question Selection & Format:** SimBench is a multiple-choice benchmark. From all 20 datasets, we therefore select only multiple-choice questions, and transform continuous scale questions into multiple-choice by splitting the scale into uniform bins. Where applicable, we collapse answer options to limit the maximum number of answering options to at most 26. In practice, questions rarely have more than 11 options. We exclude any questions with free-text answers and questions

<sup>&</sup>lt;sup>1</sup>Note that, while some constituent datasets recruit representative samples, SimBench as a whole is not fully representative of any specific group of participants.

that are contingent on prior questions or with multi-turn interactions. For datasets with questions
that are not originally in English, we use the English-language equivalents provided by the dataset
creators. We do this to enable consistent evaluation, but we note that simulation ability may plausibly
be correlated with prompt language, and encourage future work in this direction.

**Grouping Variables:** For each dataset, we record a brief description of the overall sampling population, the *default grouping*, in the form of a short prompt. For example, all participants in the WisdomOfCrowds dataset were US-based Amazon Mechanical Turk workers, so the default grouping prompt for this dataset is "You are an Amazon Mechanical Turk worker based in the United States.". Additionally, we select *grouping variables* for each dataset, corresponding to known participant sociodemographics, like age, gender, or race. The exact grouping variables and their values depend on what is available for each dataset. For a list of all grouping variables for each dataset, see Appendix F.

Response Distributions: We record the answers to each question in SimBench as group-level 120 response distributions over the question's multiple-choice options. These distributions serve as the 121 reference that we compare LLM predictions to. We create group-level response distributions by 122 aggregating over the answers from all participants that belong to a given group. We set minimum 123 grouping size thresholds for each dataset, filtering out groups with insufficient participants to form 124 meaningful response distributions. Through this aggregation process, SimBench encompasses 125 10,930,271 unique question, grouping variable value pairs, each representing a distinct simulation 126 target (see Table 3 for detailed counts). This approach enables robust evaluation of how accurately 127 LLMs can simulate response patterns across diverse demographic groups and question types. 128

#### 129 2.3 SimBench Splits

While the complete SimBench contains over 10 million potential test cases, for practical evaluation purposes we focus on two carefully curated splits that still provide comprehensive coverage of the simulation capabilities we aim to assess:

The SimBenchPop split covers all questions in all 20 datasets after processing as in §2.2. We
 combine each question with the dataset-specific default grouping prompt to create one unique test case,
 resulting in 7,167 test cases. We obtain the response distribution for each test case by aggregating all
 individual responses to that test case over all participants in that dataset. Conceptually, SimBenchPop
 measures the ability of LLMs to simulate responses of broad and diverse human populations.
 2) The SimBenchGrouped split contains only the five large-scale survey datasets in SimBench
 (AfroBarometer, ESS, ISSP, LatinoBarometro, and OpinionQA) because for these datasets we have

(AfroBarometer, ESS, ISSP, LatinoBarometro, and OpinionQA) because for these datasets we have
 enough participants to obtain meaningful group sizes even when selecting on a specific group attribute
 (e.g., age = 30-49). For each dataset, we select questions that exhibit significant variation across
 demographic groups, ensuring that the benchmark captures meaningful demographic differences
 in responses. This results in 6,343 test cases overall. For more details on the sampling process,
 see Appendix C. Conceptually, SimBenchGrouped measures the ability of LLMs to simulate
 responses from narrower participant groups based on specified group characteristics.<sup>2</sup>

# **146 3 Experimental Setup**

Tested Models: To demonstrate the usefulness of SimBench and answer our six research questions
(§1), we evaluate 24 state-of-the-art LLMs across 7 model families on SimBench. This includes both
commercial and open-weight, base and instruction-tuned models, with model sizes ranging from
0.5B to 405B parameters. Table 1 shows the full list of models.

**Model Elicitation:** For each model, we collect predictions for the two main splits of SimBench (§2.3). To obtain model response distributions, we use one of two methods, depending on model type: 1) For base models, we directly extract **token probabilities** for each response option based on first-token logits. This is a natural way of eliciting a distribution out of an LLM, especially a base LLM. 2) For instruction-tuned models, we follow recent literature on LLM calibration and

<sup>&</sup>lt;sup>2</sup>Ideally, we would also like to measure LLM simulation ability for intersectional groups that combine multiple characteristics (e.g., female + age 30-49). However, selecting on multiple characteristics substantially decreases group size, thus increasing sampling noise in the response distributions. Reliable evaluation of intersectional group simulation ability would require datasets with more participants than we have access to.

- distribution prediction [63, 48] and use **verbalized distributions**, e.g., "Option A: 30%, Option B:
- <sup>157</sup> 70%", elicited through prompting. For implementation details and prompt formats, see Appendix H.

**Evaluation Metric**: To measure LLM simulation ability, we derive the SimBench score *S* from Total Variation Distance TVD, defined as:

$$S(P,Q) = 100 \left(1 - \frac{TVD(P,Q)}{TVD(P,U)}\right) = 100 \left(1 - \frac{\sum_{i} |P_{i} - Q_{i}|}{\sum_{i} |P_{i} - U_{i}|}\right)$$
(1)

where P is the human ground truth distribution, Q is the distribution predicted by the LLM that is being tested, and U is a uniform distribution over all response options for a given question. Conceptually, S therefore measures how much more accurate the predictions from an LLM are than predictions from a uniform baseline model, which assigns equal probability to all response options for a given question. In other words, S quantifies the advantage of an LLM simulation over the simplest possible guess.

An S score of 100 indicates perfect alignment 166 between the LLM and the human ground 167 truth distribution, while a score  $\leq 0$  indicates 168 performance at or below the performance of 169 a uniform baseline. We chose TVD as the ba-170 sis for S due to its symmetry, boundedness, 171 and robustness to zero probabilities. For a 172 comparison to alternative metrics, see Ap-173 pendix D. 174

### 175 4 Results

#### 176 4.1 RQ1:

#### 177 General Simulation Ability of LLMs

To evaluate the general simulation ability of 178 LLMs, we measure their overall SimBench 179 score S averaged across the two main splits 180 of SimBench (Table 1). We find that even 181 leading LLMs struggle to simulate group-182 level human behaviors with high accu-183 racy, as measured across the 20 datasets in 184 SimBench. Claude-3.7-Sonnet is the best-185 performing model overall, but only achieves 186 a score of 40.80 out of a maximum of 100 187 on SimBench. This score indicates that the 188 response distributions predicted by Claude-189 3.7-Sonnet are, on average, closer to a uni-190 form response distribution than to the true 191 human response distribution. The distance 192 from the true distribution is 19.7 percentage 193

SimBench score S averaged across the two main splits of SimBench. Reasoning models are highlighted in *italics*. Models are sorted by score. Models below the dotted line perform worse than a uniform baseline. Model Type Release  $S(\uparrow)$ 

Table 1: Overall simulation ability as measured by

Model	Туре	Release	$S\left(\uparrow ight)$
Claude-3.7-Sonnet	Instr.	Closed	40.80
Claude-3.7-Sonnet-4000	Instr.	Closed	39.46
GPT-4.1	Instr.	Closed	34.56
DeepSeek-R1	Instr.	Open	34.52
DeepSeek-V3-0324	Instr.	Open	32.90
o4-mini-high	Instr.	Closed	28.99
Llama-3.1-405B-Instruct	Instr.	Open	28.41
o4-mini-low	Instr.	Closed	27.77
Gemma-3-12B-IT	Instr.	Open	18.63
Gemma-3-27B-IT	Instr.	Open	18.34
Llama-3.1-70B-Instruct	Instr.	Open	16.57
Qwen2.5-72B	Base	Open	13.35
Qwen2.5-32B	Base	Open	12.28
Qwen2.5-14B	Base	Open	11.93
Qwen2.5-3B	Base	Open	8.84
Qwen2.5-7B	Base	Open	8.76
Gemma-3-12B-PT	Base	Open	7.67
Gemma-3-27B-PT	Base	Open	5.54
Qwen2.5-1.5B	Base	Open	5.34
Llama-3.1-8B-Instruct	Instr.	Open	-0.14
Gemma-3-4B-PT	Base	Open	-0.73
Gemma-3-4B-IT	Instr.	Open	-1.91
Qwen2.5-0.5B	Base	Open	-2.99
Gemma-3-1B-PT	Base	Open	-16.13

points, on average, as shown by the TVD listed in Table 5. The best-performing open-weight LLM is
DeepSeek-R1, achieving a score of 34.52. The majority of the 24 models we test perform substantially
worse still, scoring less than 20. Notably, five models we test score below 0, indicating that their
predicted response distributions are, on average, even further away from the true human response
distribution than a uniform response distribution. Overall, these results suggest that disparate results
from prior work may combine into a somewhat disappointing picture, painting LLMs as far from
reliable simulators when considering a diversity of tasks.

#### 201 4.2 RQ2: Impact of LLM Characteristics on Simulation Ability

While even the best models struggle to perform well on SimBench, Table 1 also shows clear differences across models. Therefore, we investigate how performance varies depending on model characteristics, specifically 1) model size, and 2) test-time compute.

1) Model Size To evaluate the impact of model size on simulation ability, we plot SimBench Score 205 S against model parameter count for the four LLM families that we can test across multiple model 206 sizes (Figure 2). Our results suggest that there is a clear log-linear scaling law for LLM simulation 207 ability. Across all examined model families, an increase in parameter count generally corresponds to 208 an increase in SimBench score S, indicating better alignment between predicted and human response 209 distributions. Llama-3.1-Instruct in particular demonstrates nearly perfect log-linear scaling, with 210 the largest Llama-3.1-405B-Instruct achieving a score of 28.41. Conversely, all models with low 211 parameter counts (<10B) perform very poorly on SimBench, scoring at most 8.76 (Qwen2.5-7B). 212 Overall, the clear positive scaling trends across model families suggest that, while simulation remains 213 a challenging task for even the best models today, further model scaling may well lead to highly 214 accurate LLM simulators in the future. 215



Figure 2: **Model parameter count vs. simulation ability**. We measure model size by parameter count and simulation ability by SimBench score *S* averaged across the two main splits of SimBench.

2) Test-Time Compute To analyze the effects of increasing test-time compute on LLM simulation ability, we compare o4-mini-low vs. o4-mini-high, as well as Claude-3.7-Sonnet in its standard configuration vs. with a 4000-token thinking budget (Table 1). We are limited to these two comparisons due to budget constraints. Our results suggest that there is no clear benefit to increasing test-time compute for LLM simulation ability. However, this finding should only be interpreted as early, indicative evidence, and we hope that SimBench can enable further work in this direction.

#### 222 4.3 RQ3: Impact of Task Selection on Simulation Fidelity

The 20 datasets in SimBench correspond to very different tasks, in terms of the aspects of human 223 behavior that they measure (see §2.1). Therefore, we break down simulation fidelity by dataset, 224 showing results for the five LLMs we previously identified as the best simulators in Figure 3. We 225 find that **simulation fidelity varies substantially across tasks**, with even the best LLM simulators 226 performing worse than a uniform response baseline on several datasets, as indicated by negative 227 SimBench scores (e.g., on Jester, OSPsychMach, and MoralMachine). Generally, the different LLMs 228 exhibit similar performance patterns, with one notable exception being GPT-4.1's exceptionally high 229 score of 61.9 on OSPsychRWAS. 230



Figure 3: **Simulation fidelity by dataset** as measured by SimBench score S for each of the 20 datasets in SimBenchPop. We show results for the top five models based on results in Table 1.



Figure 4: **Response plurality vs. simulation fidelity** for base and instruction-tuned models on all questions in SimBenchPop. We measure response plurality by normalised entropy of the human response distribution and simulation fidelity by total variation distance at the question level.

# 231 4.4 RQ4: Impact of Response Plurality on Simulation Fidelity

Human participants give very similar responses to some questions while giving very different 232 responses to others. Faithful simulation requires models to perform well in either scenario. We 233 operationalise the level of response plurality by measuring the normalised entropy of the human 234 response distribution at the question level. We then plot this entropy for all questions in SimBenchPop 235 against total variation distance (TVD, see §3), which measures the difference in predicted and 236 reference distribution at a question level (Figure 4). Prior work has found that instruction-tuning 237 encourages models to produce more confident, less ambiguous outputs, resulting in low-entropy token 238 distributions [13, 63, 48, 16]. Therefore, we differentiate between base and instruction-tuned models 239 for this analysis. We find that **base models generally perform better on questions where human** 240 participants tend to give different answers, whereas the inverse holds for instruction-tuned 241 242 **models**. This finding is supported by our regression analysis in Appendix 6, which confirms the statistical significance of this effect. Therefore, while instruction-tuned models tend to outperform 243 base models in terms of overall score on SimBench (Table 1), our results here suggest that instruction-244 tuning also worsens simulation ability for at least a subset of high-plurality questions. 245

#### 246 4.5 RQ5: Simulation Ability Across Participant Groups

Many applications require simulating responses from specific demographic groups rather than general
 populations. Using SimBenchGrouped, we evaluate how LLM simulation ability changes when
 conditioned on specific demographic attributes.

We measure this change as  $\Delta S = S_{grouped} - S_{ungrouped}$ , where  $S_{ungrouped}$  is the SimBench score for simulating the general population and  $S_{grouped}$  is the score when simulating a specific demographic group on the same question. A negative  $\Delta S$  indicates that the model's simulation ability relative to the uniform baseline decreases when asked to simulate specific demographic groups.

Importantly, for SimBenchGrouped, we specifically selected 254 questions where human response distributions showed the high-255 est variance across demographic groups (see §2.3). The ob-256 served degradation in simulation performance therefore likely 257 represents an upper bound on the challenges LLMs face when 258 simulating specific demographic groups. Our results in Table 259 2 show that LLMs struggle more with simulating specific 260 demographic groups compared to general populations. All 261 evaluated models show negative mean  $\Delta S$  values, with degra-262 dation ranging from -1.27 for DeepSeek-V3-0324 to -4.61 for 263 Claude-3.7-Sonnet-4000. 264

The performance degradation varies substantially by demo-265 graphic category. Models struggle most when simulating groups 266 defined by religious attributes, with conditioning on 'Religios-267 ity/Practice' causing the largest decrease in simulation accu-268 racy ( $\Delta S = -9.91$ ), followed by 'Political Affiliation/Ideology' 269  $(\Delta S = -4.97)$  and 'Religion (Affiliation)'  $(\Delta S = -4.83)$ . In 270 contrast, models maintain relatively better performance when 271 simulating groups defined by 'Gender' ( $\Delta S = -1.24$ ) and 'Age' 272  $(\Delta S = -1.50).$ 273

274 While these findings may not fully generalize to cases where

demographic differences are less pronounced, they highlight potential limitations in how current LLMs capture the nuanced response patterns of specific demographic groups. We argue that such

<sup>276</sup> LLMs capture the nuanced response patterns of specific demographic groups. We argue that such <sup>277</sup> challenging benchmarks are crucial for identifying areas where improvements are most needed,

particularly for applications that aim to model the behaviors of specific subpopulations.

#### 279 4.6 RQ6: Simulation Ability vs. General Capabilities

Finally, we analyze the relationship between LLM simulation ability and more general model capabilities by correlating performance on SimBench with popular LLM capability benchmarks (Figure 5). Specifically, we compare SimBench scores to performance on GPQA Diamond [59] and OTIS AIME [24], based on scores reported in the Epoch AI Benchmarking Hub [23], which we are able to retrieve for 8 of the LLMs we test. We also compare to Chatbot Arena ELO scores [15], retrieved for the same 8 models on May 14th, 2025.



Figure 5: General model capabilities vs. simulation ability, as measured by popular benchmark scores compared to SimBench score S averaged across the two main splits in SimBench.

We find that **simulation ability is positively correlated with general model capabilities**. This matches our earlier finding on the benefits of model scaling (§4.2). However, the strength of the correlation varies across capability benchmarks. Most notably, the very strong correlation with GPQA suggests that there may be substantial symbiotic effects between scientific reasoning and simulation for social and behavioral science tasks of the kind included in SimBench. By comparison, the weaker correlation with Chatbot Arena scores suggests optimising LLMs for general helpfulness and user satisfaction does not necessarily make them better simulators.

Table 2: **Ungrouped vs. grouped** simulation performance  $\Delta S$ .

Models	
Claude-3.7-Sonnet	-3.13
Claude-3.7-Sonnet-4000	-4.61
DeepSeek-R1	-3.79
DeepSeek-V3-0324	-1.27
GPT-4.1	-3.94
Demographics	

01	
Religiosity/Practice	-9.91
Political Affil./Ideology	-4.97
Religion (Affiliation)	-4.83
Income/Social Standing	-4.51
Domicile/Urbanicity	-3.17
Employment Status	-3.03
Education	-2.55
Marital Status	-1.80
Age	-1.50
Gender	-1.24

# 293 **5 Related Work**

Human Behavior Simulation with LLMs LLMs as human behavior simulators have attracted
significant interdisciplinary attention. Researchers have evaluated their efficacy across political
science [6, 12, 19], psychology [2, 11, 46, 30], economics [31, 2], and computer science applications
[32, 20, 33, 55]. Evidence regarding LLMs' simulation fidelity remains mixed, with some studies
reporting promising results [6] while others identify critical limitations, including homogenized group
representations [14, 65] and deterministic rather than distributional predictions [57].

Existing work has predominantly focused on individual-level simulation with minimal demographic 300 conditioning, typically evaluating only one or two models in narrowly defined contexts. SimBench 301 addresses these limitations by providing a comprehensive benchmark for group-level simulation 302 across diverse domains with systematic demographic conditioning and standardized metrics. The 303 benchmark's distributional evaluation framework (using Total Variation distance) captures how 304 accurately models represent the full spectrum of human response variation-an approach advocated 305 by researchers in both simulation [4] and general LLM evaluation [69]. For broader context on this 306 emerging field, we refer readers to recent comprehensive surveys [42, 52, 4]. 307

308 **Benchmarks for LLM Evaluation** Comprehensive benchmarks have been instrumental in driving LLM advancement by providing standardized evaluation frameworks. General language understand-309 ing benchmarks such as GLUE [66] and MMLU [29] have established foundational metrics for 310 assessing natural language understanding and reasoning capabilities. As LLM applications have di-311 versified, domain-specific benchmarks have emerged, including TruthfulQA [45] for factual accuracy, 312 LegalBench [27] for legal reasoning, and Chatbot Arena [15] for chat assistants. These specialized 313 benchmarks have enabled more precise evaluation of LLMs' fitness for particular use cases and have 314 guided domain-specific optimization. 315

Most closely related to SimBench are OpinionQA [60] and GlobalOpinionQA [21], which evaluate 316 how accurately LLMs represent viewpoints of specific demographic groups. However, these bench-317 marks are limited in scope: OpinionQA focuses exclusively on U.S. public opinion surveys, while 318 GlobalOpinionOA extends this approach globally but remains constrained to survey data. In contrast, 319 SimBench represents a substantial advancement in simulation evaluation by: (1) incorporating a 320 diverse collection of 20 distinct tasks spanning multiple domains beyond surveys, (2) conceptual-321 izing simulation as a fundamental capability deserving systematic evaluation rather than merely a 322 representation challenge, and (3) establishing a unified evaluation framework that enables consistent 323 cross-domain and cross-model comparison of simulation fidelity. 324

325 Appendix G continues our discussion of related work.

# 326 6 Conclusion

LLM simulations of human behavior have the potential to create immense benefits for society by helping shape effective policy, guiding industrial decisions, and informing academic research. To fulfill this potential, however, LLM simulations must be sufficiently faithful in representing real human behaviors across diverse settings and tasks. Prior work evaluating LLM simulation fidelity has taken a predominantly narrow approach, producing an incomplete patchwork of evidence.

To change this, we introduced SimBench, the first large-scale benchmark for evaluating group-level 332 333 LLM simulation ability. We described the dataset selection and processing steps that resulted in 20 334 datasets with a unified format, measuring diverse types of human behavior (e.g., decision-making 335 vs. self-assessment) across hundreds of thousands of diverse participants from different parts of the world. Using SimBench, we took a first step toward answering fundamental questions regarding 336 when, how, and why LLM simulations succeed or fail. For example, we demonstrated that while even 337 the best LLMs today have limited simulation ability, there is a clear log-linear scaling relationship 338 with model size and a strong correlation between simulation and scientific reasoning abilities. 339

Significant progress remains to be made in developing LLMs as better simulators of human behavior.
 We hope that SimBench can provide an open foundation for future efforts in this direction, ultimately
 benefiting society as a whole.

# 343 **References**

- [1] Afrobarometer. Afrobarometer data, all countries (39), round 9, 2023. http://www.
   afrobarometer.org, 2023. Accessed: March 2025.
- [2] G. V. Aher, R. I. Arriaga, and A. T. Kalai. Using large language models to simulate multiple
  humans and replicate human subject studies. In A. Krause, E. Brunskill, K. Cho, B. Engelhardt,
  S. Sabato, and J. Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 337–371. PMLR,
  23–29 Jul 2023.
- [3] G. Ahnert, M. Pellert, D. Garcia, and M. Strohmaier. Extracting affect aggregates from
   longitudinal social media data with temporal adapters for large language models. *arXiv preprint arXiv:2409.17990*, 2024.
- [4] J. R. Anthis, R. Liu, S. M. Richardson, A. C. Kozlowski, B. Koch, J. Evans, E. Brynjolfsson,
   and M. Bernstein. Llm social simulations are a promising research method. *arXiv preprint arXiv:2504.02234*, 2025.
- [5] Anthropic. Claude 3.7 sonnet and claude code, Feb. 2025. Accessed: 2025-05-14.
- [6] L. P. Argyle, E. C. Busby, N. Fulda, J. R. Gubler, C. Rytting, and D. Wingate. Out of one, many:
   Using language models to simulate human samples. *Political Analysis*, 31(3):337–351, 2023.
- [7] L. Aroyo, A. S. Taylor, M. Díaz, C. M. Homan, A. Parrish, G. Serapio-García, V. Prabhakaran,
   and D. Wang. Dices dataset: diversity in conversational ai evaluation for safety. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, NIPS '23, Red
   Hook, NY, USA, 2023. Curran Associates Inc.
- [8] E. Awad, S. Dsouza, R. Kim, J. Schulz, J. Henrich, A. Shariff, J.-F. Bonnefon, and I. Rahwan.
   The moral machine experiment. *Nature*, 563(7729):59–64, 2018.
- [9] E. Awad, S. Dsouza, A. Shariff, I. Rahwan, and J.-F. Bonnefon. Universals and variations
   in moral decisions made in 42 countries by 70,000 participants. *Proceedings of the National Academy of Sciences*, 117(5):2332–2337, 2020.
- [10] E. Bigelow and S. T. Piantadosi. A large dataset of generalization patterns in the number game.
   *Journal of Open Psychology Data*, 4(1):e4–e4, 2016.
- [11] M. Binz, E. Akata, M. Bethge, F. Brändle, F. Callaway, J. Coda-Forno, P. Dayan, C. Demircan,
   M. K. Eckstein, N. Éltető, et al. Centaur: a foundation model of human cognition. *arXiv preprint arXiv:2410.20268*, 2024.
- [12] J. Bisbee, J. D. Clinton, C. Dorff, B. Kenkel, and J. M. Larson. Synthetic replacements for
   human survey data? the perils of large language models. *Political Analysis*, 32(4):401–416,
   2024.
- [13] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam,
   G. Sastry, A. Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [14] M. Cheng, T. Piccardi, and D. Yang. CoMPosT: Characterizing and evaluating caricature in LLM simulations. In H. Bouamor, J. Pino, and K. Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 10853–10875, Singapore, Dec. 2023. Association for Computational Linguistics.
- W.-L. Chiang, L. Zheng, Y. Sheng, A. N. Angelopoulos, T. Li, D. Li, B. Zhu, H. Zhang, M. I.
   Jordan, J. E. Gonzalez, and I. Stoica. Chatbot arena: an open platform for evaluating llms by
   human preference. In *Proceedings of the 41st International Conference on Machine Learning*,
   ICML'24. JMLR.org, 2024.
- [16] A. F. Cruz, M. Hardt, and C. Mendler-Dünner. Evaluating language models as risk scores. In
   A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang, editors,
   *Advances in Neural Information Processing Systems*, volume 37, pages 97378–97407. Curran
   Associates, Inc., 2024.

- <sup>392</sup> [17] DeepSeek-AI. Deepseek-v3 technical report. arXiv preprint arXiv:2412.19437, 2024.
- [18] T. Dettmers, M. Lewis, Y. Belkada, and L. Zettlemoyer. Llm.int8(): 8-bit matrix multiplication
   for transformers at scale. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA, 2022. Curran Associates Inc.
- [19] R. Dominguez-Olmedo, M. Hardt, and C. Mendler-Dünner. Questioning the survey responses of
   large language models. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak,
   and C. Zhang, editors, *Advances in Neural Information Processing Systems*, volume 37, pages
   45850–45878. Curran Associates, Inc., 2024.
- Y. R. Dong, T. Hu, and N. Collier. Can LLM be a personalized judge? In Y. Al-Onaizan,
   M. Bansal, and Y.-N. Chen, editors, *Findings of the Association for Computational Linguis- tics: EMNLP 2024*, pages 10126–10141, Miami, Florida, USA, Nov. 2024. Association for
   Computational Linguistics.
- E. Durmus, K. Nguyen, T. Liao, N. Schiefer, A. Askell, A. Bakhtin, C. Chen, Z. HatfieldDodds, D. Hernandez, N. Joseph, L. Lovitt, S. McCandlish, O. Sikder, A. Tamkin, J. Thamkul,
  J. Kaplan, J. Clark, and D. Ganguli. Towards measuring the representation of subjective global
  opinions in language models. In *First Conference on Language Modeling*, 2024.
- [22] A. Enders, C. Klofstad, A. Diekman, H. Drochon, J. Rogers de Waal, S. Littrell, K. Premaratne,
   D. Verdear, S. Wuchty, and J. Uscinski. The sociodemographic correlates of conspiracism.
   *Scientific reports*, 14(1):14184, 2024.
- 411 [23] Epoch AI. "ai benchmarking hub", 11 2024. https://epoch.ai/data/ai-benchmarking-dashboard.
- 412 [24] EpochAI. Otis mock aime 24-25. https://huggingface.co/datasets/EpochAI/ 413 otis-mock-aime-24-25, 2024. Accessed: 2025-05-11.
- [25] European Social Survey European Research Infrastructure (ESS ERIC). ESS11 Integrated
   File, Edition 2.0 [Data set], 2024.
- [26] K. Goldberg, T. Roeder, D. Gupta, and C. Perkins. Eigentaste: A constant time collaborative
   filtering algorithm. *information retrieval*, 4:133–151, 2001.
- [27] N. Guha, J. Nyarko, D. Ho, C. Ré, A. Chilton, A. K, A. Chohlas-Wood, A. Peters, B. Waldon,
  D. Rockmore, D. Zambrano, D. Talisman, E. Hoque, F. Surani, F. Fagan, G. Sarfaty, G. Dickinson, H. Porat, J. Hegland, J. Wu, J. Nudell, J. Niklaus, J. Nay, J. Choi, K. Tobia, M. Hagan,
  M. Ma, M. Livermore, N. Rasumov-Rahe, N. Holzenberger, N. Kolt, P. Henderson, S. Rehaag,
  S. Goel, S. Gao, S. Williams, S. Gandhi, T. Zur, V. Iyer, and Z. Li. Legalbench: A collaboratively
  built benchmark for measuring legal reasoning in large language models. In *Advances in Neural Information Processing Systems 36 (NeurIPS 2023), Datasets and Benchmarks Track*, 2023.
- [28] D. Guo, D. Yang, H. Zhang, J. Song, R. Zhang, R. Xu, Q. Zhu, S. Ma, P. Wang, X. Bi, et al.
   Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv* preprint arXiv:2501.12948, 2025.
- [29] D. Hendrycks, C. Burns, S. Basart, A. Zou, M. Mazeika, D. Song, and J. Steinhardt. Mea suring massive multitask language understanding. In *International Conference on Learning Representations*, 2021.
- [30] L. Hewitt, A. Ashokkumar, I. Ghezae, and R. Willer. Predicting results of social science
   experiments using large language models, August 2024.
- [31] J. J. Horton. Large language models as simulated economic agents: What can we learn from
   homo silicus? Technical report, National Bureau of Economic Research, 2023.
- [32] T. Hu and N. Collier. Quantifying the persona effect in LLM simulations. In L.-W. Ku, A. Martins, and V. Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10289–10307, Bangkok, Thailand, Aug. 2024. Association for Computational Linguistics.

- [33] T. Hu and N. Collier. inews: A multimodal dataset for modeling personalized affective responses
   to news. *arXiv preprint arXiv:2503.03335*, 2025.
- [34] ISSP Research Group. International social survey programme: Social networks and social resources issp 2017. GESIS Data Archive, Cologne. ZA6980 Data file Version 2.0.0, https://doi.org/10.4232/1.13322, 2019.
- ISSP Research Group. International social survey programme: Religion iv issp 2018. GESIS
   Data Archive, Cologne. ZA7570 Data file Version 2.1.0, https://doi.org/10.4232/1.13629, 2020.
- [36] ISSP Research Group. International social survey programme: Social inequality v issp 2019.
   GESIS, Cologne. ZA7600 Data file Version 3.0.0, https://doi.org/10.4232/1.14009, 2022.
- [37] ISSP Research Group. International social survey programme: Environment iv issp 2020.
   GESIS, Cologne. ZA7650 Data file Version 2.0.0, https://doi.org/10.4232/1.14153, 2023.
- [38] ISSP Research Group. Za8000 international social survey programme: Health and
   health care ii issp 2021. GESIS, Cologne. ZA8000 Data file Version 2.0.0,
   https://doi.org/10.4232/5.ZA8000.2.0.0, 2024.
- [39] Z. Jiang, J. Araki, H. Ding, and G. Neubig. How can we know when language models know?
   on the calibration of language models for question answering. *Transactions of the Association* for Computational Linguistics, 9:962–977, 2021.
- [40] A. T. Kalai and S. S. Vempala. Calibrated language models must hallucinate. In *Proceedings of the 56th Annual ACM Symposium on Theory of Computing*, STOC 2024, page 160–171, New York, NY, USA, 2024. Association for Computing Machinery.
- [41] S. Kapoor, N. Gruver, M. Roberts, A. Pal, S. Dooley, M. Goldblum, and A. Wilson. Calibration-tuning: Teaching large language models to know what they don't know. In R. Vázquez, H. Celikkanat, D. Ulmer, J. Tiedemann, S. Swayamdipta, W. Aziz, B. Plank, J. Baan, and M.-C. de Marneffe, editors, *Proceedings of the 1st Workshop on Uncertainty-Aware NLP (UncertaiNLP 2024)*, pages 1–14, St Julians, Malta, Mar. 2024. Association for Computational Linguistics.
- [42] A. C. Kozlowski and J. Evans. Simulating subjects: The promise and peril of ai stand-ins for
   social agents and interactions, September 2024. Preprint.
- [43] Latinobarómetro. Latinobarómetro 2023. http://www.latinobarometro.org, 2023. Ac cessed: March 2025.
- [44] A. Lazaridou, A. Kuncoro, E. Gribovskaya, D. Agrawal, A. Liska, T. Terzi, M. Gimenez,
  C. de Masson d'Autume, T. Kocisky, S. Ruder, D. Yogatama, K. Cao, S. Young, and P. Blunsom.
  Mind the gap: Assessing temporal generalization in neural language models. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 29348–29363. Curran Associates, Inc.,
  2021.
- [45] S. Lin, J. Hilton, and O. Evans. TruthfulQA: Measuring how models mimic human falsehoods.
  In S. Muresan, P. Nakov, and A. Villavicencio, editors, *Proceedings of the 60th Annual Meeting* of the Association for Computational Linguistics (Volume 1: Long Papers), pages 3214–3252, Dublin, Ireland, May 2022. Association for Computational Linguistics.
- [46] B. S. Manning, K. Zhu, and J. J. Horton. Automated social science: Language models as
   scientist and subjects. Technical report, National Bureau of Economic Research, 2024.
- [47] N. G. Mede, V. Cologna, S. Berger, J. Besley, C. Brick, M. Joubert, E. W. Maibach, S. Mihelj,
   N. Oreskes, M. S. Schäfer, et al. Perceptions of science, science communication, and climate
   change attitudes in 68 countries–the tisp dataset. *Scientific data*, 12(1):114, 2025.
- [48] N. Meister, C. Guestrin, and T. Hashimoto. Benchmarking distributional alignment of large
   language models. In L. Chiruzzo, A. Ritter, and L. Wang, editors, *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 24–49, Albuquerque,
   New Mexico, Apr. 2025. Association for Computational Linguistics.

- [49] Meta AI. Introducing Llama 3.1: Our most capable models to date, 2024. Accessed: 2025-05-14.
- Y. Nie, X. Zhou, and M. Bansal. What can we learn from collective human opinions on natural language inference data? In B. Webber, T. Cohn, Y. He, and Y. Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9131–9143, Online, Nov. 2020. Association for Computational Linguistics.
- [51] H. Nori, N. King, S. M. McKinney, D. Carignan, and E. Horvitz. Capabilities of gpt-4 on medical challenge problems. *arXiv preprint arXiv:2303.13375*, 2023.
- [52] A. Olteanu, S. Barocas, S. L. Blodgett, L. Egede, A. DeVrio, and M. Cheng. Ai automatons: Ai
   systems intended to imitate humans. *arXiv preprint arXiv:2503.02250*, 2025.
- <sup>498</sup> [53] OpenAI. Introducing GPT-4.1 in the API, 2025. Accessed: 2025-05-14.
- 499 [54] OpenAI. Introducing openai o3 and o4-mini, Apr. 2025. Accessed: 2025-05-14.
- [55] J. S. Park, J. O'Brien, C. J. Cai, M. R. Morris, P. Liang, and M. S. Bernstein. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, UIST '23, New York, NY, USA, 2023. Association for Computing Machinery.
- J. S. Park, C. Q. Zou, A. Shaw, B. M. Hill, C. Cai, M. R. Morris, R. Willer, P. Liang, and M. S.
   Bernstein. Generative agent simulations of 1,000 people. *arXiv preprint arXiv:2411.10109*, 2024.
- [57] P. S. Park, P. Schoenegger, and C. Zhu. Diminished diversity-of-thought in a standard large
   language model. *Behavior Research Methods*, 56:5754–5770, 2024.
- [58] J. C. Peterson, D. D. Bourgin, M. Agrawal, D. Reichman, and T. L. Griffiths. Using large-scale
   experiments and machine learning to discover theories of human decision-making. *Science*, 372(6547):1209–1214, 2021.
- [59] D. Rein, B. L. Hou, A. C. Stickland, J. Petty, R. Y. Pang, J. Dirani, J. Michael, and S. R.
   Bowman. GPQA: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*, 2024.
- [60] S. Santurkar, E. Durmus, F. Ladhak, C. Lee, P. Liang, and T. Hashimoto. Whose opinions do
   language models reflect? In *International Conference on Machine Learning*, pages 29971–
   30004. PMLR, 2023.
- [61] C. Simoiu, C. Sumanth, A. Mysore, and S. Goel. Studying the "wisdom of crowds" at scale. In
   *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing*, volume 7,
   pages 171–179, 2019.
- [62] G. Team, A. Kamath, J. Ferret, S. Pathak, N. Vieillard, R. Merhej, S. Perrin, T. Matejovicova,
   A. Ramé, M. Rivière, et al. Gemma 3 technical report. *arXiv preprint arXiv:2503.19786*, 2025.
- [63] K. Tian, E. Mitchell, A. Zhou, A. Sharma, R. Rafailov, H. Yao, C. Finn, and C. Manning. Just
   ask for calibration: Strategies for eliciting calibrated confidence scores from language models
   fine-tuned with human feedback. In H. Bouamor, J. Pino, and K. Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5433–5442,
   Singapore, Dec. 2023. Association for Computational Linguistics.
- [64] L. Tjuatja, V. Chen, T. Wu, A. Talwalkwar, and G. Neubig. Do llms exhibit human-like response
   biases? a case study in survey design. *Transactions of the Association for Computational Linguistics*, 12:1011–1026, 09 2024.
- [65] A. Wang, J. Morgenstern, and J. P. Dickerson. Large language models that replace human
   participants can harmfully misportray and flatten identity groups. *Nature Machine Intelligence*,
   pages 1–12, 2025.

- [66] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. Bowman. GLUE: A multi-task
   benchmark and analysis platform for natural language understanding. In T. Linzen, G. Chrupała,
   and A. Alishahi, editors, *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 353–355, Brussels, Belgium, Nov. 2018.
   Association for Computational Linguistics.
- [67] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf,
  M. Funtowicz, J. Davison, S. Shleifer, P. von Platen, C. Ma, Y. Jernite, J. Plu, C. Xu, T. Le Scao,
  S. Gugger, M. Drame, Q. Lhoest, and A. Rush. Transformers: State-of-the-art natural language
  processing. In Q. Liu and D. Schlangen, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45,
  Online, Oct. 2020. Association for Computational Linguistics.
- [68] A. Yang, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, C. Li, D. Liu, F. Huang, H. Wei, et al.
   Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- [69] L. Ying, K. M. Collins, L. Wong, I. Sucholutsky, R. Liu, A. Weller, T. Shu, T. L. Griffiths,
   and J. B. Tenenbaum. On benchmarking human-like intelligence in machines. *arXiv preprint arXiv:2502.20502*, 2025.
- [70] Z. Zhao, E. Wallace, S. Feng, D. Klein, and S. Singh. Calibrate before use: Improving few-shot
   performance of language models. In M. Meila and T. Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 12697–12706. PMLR, 18–24 Jul 2021.
- [71] C. Zhu, B. Xu, Q. Wang, Y. Zhang, and Z. Mao. On the calibration of large language models
   and alignment. In H. Bouamor, J. Pino, and K. Bali, editors, *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 9778–9795, Singapore, Dec. 2023. Association
   for Computational Linguistics.

# 558 A Limitations

**Scope of Representativeness** Although SimBench spans 20 diverse datasets, the combined sample 559 does (and can) not fully represent any single population in its full complexity. Many geographic 560 regions are still underrepresented or entirely absent, potentially limiting generalizability to popu-561 lations with different cultural backgrounds and preferences. Even within countries, demographic 562 representativeness may vary, as only a subset of our 20 datasets are based on nationally representative 563 sampling techniques. Each dataset carries its own statistical uncertainty. Opt-in samples and crowd-564 sourced data (e.g., from Amazon Mechanical Turk) may have larger margins of error than nationally 565 representative surveys, potentially affecting the benchmark's precision for certain questions. We 566 567 view these limitations as opportunities for collaborative extension of SimBench to improve global coverage and representativeness over time. 568

**Temporal Dimensions** The current version of SimBench utilizes static datasets that capture human behavior at specific points in time. This approach allows for systematic evaluation across domains but cannot yet assess how well LLMs simulate evolving preferences, opinion shifts, or behavioral adaptation—all fundamental aspects of human behavior. Future iterations of SimBench could incorporate longitudinal data to address these dynamic aspects of human behavior and expand the benchmark's evaluative capacity.

**Task Format Considerations** SimBench currently focuses on multiple-choice, single-answer, single-turn questions and interactions. This standardized format enables systematic comparison across diverse domains but necessarily excludes more complex behavioral simulations including multi-step decision processes and interactive social dynamics. We see this as a pragmatic starting point that establishes foundational evaluation capabilities while inviting future extensions to capture more nuanced aspects of human behavior.

**Training Data Overlap** Without complete transparency into model training corpora, we cannot 581 definitively rule out the possibility that some test items appeared during training. However, several 582 factors mitigate concerns about data contamination affecting our results. First, SimBench evaluates 583 simulation at the group distribution level rather than individual response prediction, making memo-584 rization of specific survey responses less impactful. Second, many of our datasets primarily exist 585 as aggregated statistics in published research rather than as widely available raw data. Finally, the 586 consistent scaling patterns we observe across diverse datasets suggest genuine simulation capabilities 587 rather than artifacts of training data overlap. Nevertheless, we acknowledge that data contamination 588 remains a fundamental challenge in LLM evaluation, and future work should develop more robust 589 methods to detect and quantify its impact. We include this consideration for completeness while 590 believing it unlikely to significantly impact our current findings. 591

# 592 **B** Ethical Considerations

SimBench's primary purpose is to benchmark LLMs' ability to simulate human behavior. While
 advancements in LLM simulation capabilities can support helpful applications such as pre-testing
 policies, these do not come without risks of misrepresentation and dual use.

First and foremost, due to the observed limited simulation ability of state-of-the-art LLMs, we caution against relying on LLM-powered simulations of human behavior for tasks where downstream harm is possible. Even as models improve, substituting algorithmic approximations for authentic human participation carries the risk of disadvantaging under-represented/marginalized communities by removing their opportunities to directly shape decisions that affect them. Furthermore, while benchmarks like SimBench help measure simulation capabilities, we must be careful not to mistake increasing benchmark performance for genuine understanding of complex human behavior.

While SimBench includes diverse demographic groups, it can not adequately support simulations of intersectional identities due to sample size limitations. By conditioning on one demographic variable at a time, we cannot systematically assess how well models handle the rich overlap of identities (e.g., "older Latinx women," "young Black men"). Small intersectional group sizes make it difficult to combine multiple characteristics simultaneously due to increasing sampling noise in response distributions. Yet intersectional simulation is precisely where societal biases and model limitations often emerge, making this an important direction for future work. Additionally, the
 conditional prompting approach we use conceptualizes simplistic human populations and may thus
 fail to appropriately account for nuances of individual behavior.

Nevertheless, we believe SimBench is an important step toward making LLM simulation progress measurable and raising awareness of state-of-the-art model blind spots. Together, we hope this will ultimately create accountability for models deployed in socially sensitive contexts.

# 615 C SimBenchPop and SimBenchGrouped Sampling Details

<sup>616</sup> We curated data at two levels of grouping granularity, corresponding to our two main benchmark <sup>617</sup> splits: **SimBenchPop** and **SimBenchGrouped**.

**SimBenchPop** measures LLMs' ability to simulate responses of broad, diverse human populations. 618 We include all questions from all 20 datasets in SimBench, combining each question with its dataset-619 specific default grouping prompt (e.g., "You are an Amazon Mechanical Turk worker based in the 620 United States"). We sample up to 500 questions per dataset to ensure representativeness while 621 keeping the benchmark manageable. For each test case, we aggregate individual responses across all 622 participants in the dataset to create population-level response distributions. This approach creates a 623 benchmark that represents population-level responses across diverse domains while maintaining a 624 reasonable size of 7,167 test cases. 625

For **SimBenchGrouped**, we focus only on five large-scale survey datasets with rich demographic in-626 formation and sufficient sample sizes: OpinionQA, ESS, Afrobarometer, ISSP, and LatinoBarometro. 627 Our sampling approach prioritizes questions showing meaningful demographic variation. For each 628 dataset, we identify available grouping variables (e.g., age, gender, country) with sufficient group 629 sizes to form meaningful response distributions. We calculate the variance of responses across 630 demographic groups for each question and rank questions by their variance scores, prioritizing those 631 showing the strongest demographic differences. We select questions that exhibit significant variation 632 across demographic groups to ensure the benchmark captures meaningful differences in responses. 633 For each selected question, we create multiple test cases by pairing it with different values of the 634 grouping variables (e.g., age = "18-29", age = "30-49"). This process results in 6,343 test cases that 635 specifically measure LLMs' ability to simulate responses from narrower participant groups based on 636 specified demographic characteristics. Table 3 provides a summary of the sampling process across all 637 datasets, showing the minimum group size thresholds and the number of test cases in each benchmark 638 split. 639

# 640 **D** Metric Robustness Check

TVD ranges from 0 (perfect match) to 1 (complete disagreement), with lower values indicating 641 better simulation fidelity. TVD provides an interpretable measure of how closely model predictions 642 align with actual human response distributions. TVD is particularly well-suited for simulation 643 evaluation compared to alternatives like KL divergence or Jensen-Shannon divergence (JSD). Unlike 644 KL divergence, TVD remains well-defined even when the model assigns zero probability to responses 645 that humans give, avoiding the infinite penalties that KL would impose in such cases. Additionally, 646 647 TVD is symmetric and bounded, making it more interpretable across different datasets and response distributions than KL divergence. While JSD offers similar advantages in terms of symmetry and 648 boundedness, TVD provides a more direct and intuitive interpretation of the maximum possible error 649 650 in probability estimates. This property is especially valuable when evaluating how accurately models simulate the distribution of human responses rather than just matching the most likely response. For 651 further discussion on TVD as an evaluation metric, see also [48]. We show the results of Table 1 in 652 terms of raw TVD values in Table 5. 653

To ensure our findings are robust across different metrics, we complement TVD with two alternative metrics: Jensen-Shannon Divergence (JSD) and Spearman's Rank Correlation (RC). Table 4 presents these metrics for a subset of evaluated models. The strong Pearson correlation between TVD and JSD (r = 0.92) indicates these metrics provide consistent model rankings. The moderate negative correlation (r = -0.57) between TVD and RC is expected, as lower distances correspond to higher correlations. This multi-metric evaluation confirms that our model comparisons remain consistent across different statistical measures.

Dataset	Min. Group	SimBench	SimBenchPop	SimBenchGrouped
WisdomOfCrowds	100	1,604	114	-
Jester	100	136	136	-
Choices13k	NaN	14,568	500	-
OpinionQA	300	1,074,392	500	984
MoralMachineClassic	100	3,441	15	-
MoralMachine	100	20,771	500	-
ChaosNLI	100	4,645	500	-
ESS	300	2,783,780	500	1,643
Afrobarometer	300	517,453	500	1,531
OSPsychBig5	300	1,950	250	-
OSPsychMACH	300	3,682,700	100	-
OSPsychMGKT	300	20,610	500	-
OSPsychRWAS	300	975,585	22	-
ISSP	300	594,336	500	940
LatinoBarometro	300	80,684	500	1,245
GlobalOpinionQA	NaN	46,329	500	-
DICES	10	918,064	500	-
NumberGame	10	15,984	500	-
ConspiracyCorr	300	968	45	-
TISP	300	172,271	485	_
Total		10,930,271	7,167	6,343

Table 3: Dataset Sampling Summary; NaN refers to dataset that is only available in aggregated form and no grouping size is known.

Table 4: Comparison of models on three metrics: Total Variation Distance (TVD), Jensen-Shannon Divergence (JSD), and Spearman Rank Correlation (RC). Lower values are better for TVD and JSD; higher is better for RC.

Model	<b>Total Variation</b>	JS Divergence	<b>Rank Correlation</b>
Claude-3.7-Sonnet	0.191	0.057	0.673
Claude-3.7-Sonnet-4000	0.195	0.060	0.648
DeepSeek-R1	0.211	0.069	0.623
DeepSeek-V3-0324	0.216	0.069	0.620
GPT-4.1	0.209	0.070	0.646
Llama-3.1-405B-Instruct	0.231	0.085	0.593
o4-mini-high	0.225	0.079	0.621
o4-mini-low	0.230	0.082	0.609

# 661 E Regression Analysis of Human Response Entropy and Model Performance

<sup>662</sup> To formally test the relationship between human response entropy and simulation performance across

different model types, we fit an Ordinary Least Squares (OLS) regression model predicting Total

Variation (TV) distance at the individual question-model level. The model specification was as follows:

 $Total\_Variation \sim C(dataset\_name) + C(model) + C(instruct\_flag) : Human\_Normalized\_Entropy$ (2)

Here, *Total\_Variation* is the dependent variable.  $C(\text{dataset_name})$  and C(model) represent fixed effects for each dataset and model, respectively, controlling for baseline differences in difficulty and capability. The crucial term is the interaction  $C(\text{instruct_flag})$  : Human\_Normalized\_Entropy, where *instruct\_flag* is a binary indicator for instruction-tuned models (0 for base, 1 for instruction-tuned).

<sup>670</sup> The key results from Table 6 are the coefficients for the interaction terms:

• For base models: The coefficient on the interaction between base models and Human Normalized Entropy is -0.2555 (p < 0.001), indicating that for every one-unit increase in normalized entropy, the TVD decreases by approximately 0.26 units. This means that base models perform *better* 

(lower TVD) when simulating human populations with more diverse opinions.

• For instruction-tuned models: The coefficient on the interaction between instruction-tuned models and Human Normalized Entropy is +0.1072 (p < 0.001), indicating that for every one-unit increase in normalized entropy, the TVD increases by approximately 0.11 units. This means that instruction-tuned models perform *worse* (higher TVD) when simulating human populations with more diverse opinions.

These coefficients are both highly statistically significant (p < 0.001) and represent substantial effect sizes given that TVD ranges from 0 to 1. The model as a whole explains approximately 20% of the variance in TVD ( $R^2 = 0.202$ ), which is substantial for a dataset of this size and complexity.

The opposite signs of these coefficients provide strong evidence for our hypothesis that base models and instruction-tuned models respond differently to the challenge of simulating populations with diverse opinions. This pattern holds even after controlling for the specific datasets and models involved, suggesting it represents a general property of the two model classes rather than an artifact of particular model or evaluation datasets.

Table 5: TVD for each model in SimBenchPop and SimBenchGrouped. Lower values indica	ate better
performance. PT and IT refer to pretrained and instruction-tuned versions, respectively.	

Model	SimBenchPop	SimBenchGrouped	Average
Baselines			
Random baseline	0.390	0.415	0.402
Uniform baseline	0.335	0.362	0.348
Commercial Models			
Claude-3.7-Sonnet	0.197	0.184	0.191
Claude-3.7-Sonnet-4000	0.201	0.188	0.195
GPT-4.1	0.212	0.205	0.209
o4-mini-high	0.235	0.214	0.225
o4-mini-low	0.234	0.216	0.230
Open Models			
DeepSeek-V3-0324	0.215	0.218	0.216
DeepSeek-R1	0.211	0.212	0.211
Llama-3.1-8B-Instruct	0.321	0.318	0.320
Llama-3.1-70B-Instruct	0.277	0.247	0.263
Llama-3.1-405B-Instruct	0.237	0.225	0.231
Qwen2.5-0.5B	0.337	0.364	0.349
Qwen2.5-1.5B	0.321	0.324	0.322
Qwen2.5-3B	0.300	0.327	0.313
Qwen2.5-7B	0.290	0.326	0.307
Qwen2.5-14B	0.285	0.314	0.298
Qwen2.5-32B	0.273	0.308	0.290
Qwen2.5-72B	0.269	0.300	0.283
Gemma-3-1B-PT	0.382	0.413	0.396
Gemma-3-4B-PT	0.334	0.342	0.338
Gemma-3-12B-PT	0.310	0.317	0.314
Gemma-3-27B-PT	0.309	0.325	0.317
Gemma-3-4B-IT	0.337	0.341	0.339
Gemma-3-12B-IT	0.262	0.274	0.267
Gemma-3-27B-IT	0.270	0.273	0.272

Model:	OLS	Adj. R-sq	uared: 0	0.201	_	
Dependent Variable:	Total_Variation	AIC:	-	134342.8438		
Date:	2025-05-15 20:27	BIC:		133890.3555		
No. Observations:	172008	Log-Likel	ihood: 6	5/216.		
Df Model:	44	F-statistic:		83.5		
Df Residuals:	171963	Prob (F-st	atistic): (	0.00		
R-squared:	0.201	Scale:	0	0.026805	-	
	Coef.	Std.Err.	t	P>  t	[0.025	0.975]
Intercept	0.1824	0.0029	62.1882	0.0000	0.1766	0.1881
C(dataset_name)[T.ChaosNLI]	-0.0442	0.0021	-20.7195	0.0000	-0.0483	-0.0400
C(dataset_name)[T.Choices13k]	-0.1016	0.0021	-47.3233	0.0000	-0.1058	-0.0974
C(dataset_name)[T.ConspiracyCorr]	-0.0452	0.0052	-8.6565	0.0000	-0.0554	-0.0349
C(dataset_name)[T.DICES]	-0.0254	0.0023	-11.0298	0.0000	-0.0300	-0.0209
C(dataset name)[T.ESS]	-0.0202	0.0021	-9.4882	0.0000	-0.0244	-0.0160
C(dataset name)[T.GlobalOpinionQA]	-0.0428	0.0021	-20.2041	0.0000	-0.0469	-0.0386
C(dataset name)[T.ISSP]	-0.0279	0.0021	-13.1516	0.0000	-0.0321	-0.0238
C(dataset_name)[T.Jester]	0.1168	0.0033	35.9190	0.0000	0.1104	0.1232
C(dataset name)[T.LatinoBarometro]	-0.0325	0.0021	-15.1931	0.0000	-0.0367	-0.0283
C(dataset_name)[T.MoralMachine]	-0.0380	0.0021	-17.8607	0.0000	-0.0422	-0.0339
C(dataset_name)[T.MoralMachineClassic]	-0.1594	0.0088	-18,1961	0.0000	-0.1766	-0.1422
C(dataset_name)[T.NumberGame]	-0.0821	0.0021	-38.8471	0.0000	-0.0863	-0.0780
C(dataset_name)[T OSPsychBig5]	-0.1186	0.0026	-45 0783	0.0000	-0.1238	-0 1134
C(dataset_name)[TOSPsychMACH]	-0.0227	0.0037	-6 1522	0.0000	-0.0299	-0.0155
C(dataset_name)[TOSPsychMGKT]	-0.1121	0.0021	-52 6066	0.0000	-0.1163	-0.1080
C(dataset_name)[T.OSP sychWORT]	0.0168	0.0073	2 3068	0.0211	0.0025	0.0311
$C(\text{dataset_name})[T.OpinionOA]$	-0.1013	0.0075	-47 9196	0.0000	-0.1054	-0.0972
C(dataset_name)[TTISP]	-0.0441	0.0021	-20 5072	0.0000	-0.0483	0.0300
C(dataset_name)[T.WisdomOfCrowds]	-0.0200	0.0022	-20.3072	0.0000	-0.0768	-0.0131
C(Madal)[T Clauda 3.7 Sonnat 4000]	-0.0200	0.0033	1 2078	0.1622	-0.0208	0.0002
C(Model)[T.DeenSeek-P1]	0.0038	0.0027	1.5570	0.0000	0.0079	0.0092
C(Model)[T.DeepSeek-W3-0324]	0.0155	0.0027	6.4740	0.0000	0.0073	0.0231
C(Model)[T.CPT 4 1]	0.01/1	0.0027	5 1557	0.0000	0.0023	0.0251
C(Model)[T.Gemma 2 12P IT]	0.0141	0.0027	22 4227	0.0000	0.0087	0.0195
C(Model)[T.Gomma 2 12B-I1]	0.0041	0.0027	104 5204	0.0000	0.3540	0.3684
C(Model)[T.Gemma 2 1P PT]	0.3010	0.0035	125 1200	0.0000	0.3349	0.3084
C(Model)[T.Gemma 2 27P IT]	0.4330	0.0033	26 6800	0.0000	0.4202	0.4398
C(Model)[T.Comma 2 27D DT]	0.0750	0.0027	104 1666	0.0000	0.0070	0.0784
C(Model)[T.Gemma 2 4D IT]	0.3004	0.0055	51 1024	0.0000	0.5550	0.3072
C(Model)[T.Comma 2 4D DT]	0.1398	0.0027	111 4926	0.0000	0.1344	0.1451
C(Model)[T.Leme 2.1.405P.Instruct]	0.5857	0.0033	111.4620	0.0000	0.3790	0.3923
C(Model)[T.Liama-3.1-403B-Instruct]	0.0392	0.0027	28.0426	0.0000	0.0338	0.0445
C(Model)[T.Llama-5.1-70B-Instruct]	0.0792	0.0027	28.9420	0.0000	0.0738	0.0845
C(Model)[1.Liama-3.1-8B-Instruct]	0.1251	0.0027	45.0170	0.0000	0.1178	0.1285
C(Model)[T.Qwell2.5-0.5B]	0.3880	0.0055	107.4076	0.0000	0.3812	0.3947
C(Model)[1.Qwen2.5-1.5B]	0.3/19	0.0035	107.4976	0.0000	0.3652	0.3/8/
C(Model)[T.Qwen2.5-14B]	0.5359	0.0035	97.0893	0.0000	0.3292	0.3427
C(Model)[1.Qwen2.5-32B]	0.3248	0.0035	93.8/07	0.0000	0.3180	0.3316
C(Model)[T.Qwen2.5-3B]	0.3517	0.0035	101.6583	0.0000	0.3450	0.3585
C(Model)[T.Qwen2.5-72B]	0.3198	0.0035	92.4342	0.0000	0.3130	0.3266
C(Model)[T.Qwen2.5-/B]	0.3409	0.0035	98.5348	0.0000	0.3342	0.3477
C(Model)[1.04-mini-high]	0.0374	0.0027	13.6575	0.0000	0.0320	0.0427
C(Model)[f.o4-mini-low]	0.0363	0.0027	13.2773	0.0000	0.0310	0.0417
C(instruct_flag)[base]:Human_Normalized_Entro	ору -0.2628	0.0026	-101.0841	0.0000	-0.2679	-0.2577
C(instruct_flag)[instruct]:Human_Normalized_E	ntropy 0.0929	0.0024	37.9507	0.0000	0.0881	0.0977

# Table 6: Results: Ordinary least squares

Omnibus:	21133.651	Durbin-Watson:	1.711
Prob(Omnibus):	0.000	Jarque-Bera (JB):	34296.360
Skew:	0.862	Prob(JB):	0.000
Kurtosis:	4.346	Condition No.:	33

# 688 F Dataset Details

We provide details on each of the 20 datasets in SimBench. Note that for many datasets we use only a subset of questions and participants for SimBench, as a result of our preprocessing steps (§2.2).

#### 691 F.1 WisdomOfCrowds

Description: This dataset contains factual questions that were administered to a large number of
 US-based Amazon Mechanical Turk workers. The data was originally collected to study wisdom of
 the crowd effects.

695 **Questions**: 113, with an average of 518 responses per question.

#### 696 Example question:

An analogy compares the relationship between two things or ideas to highlight some point of similarity. You will be given pairs of words bearing a relationship, and asked to select another pair

697

of words that illustrate a similar relationship.

Which pair of words has the same relationship as 'Letter : Word'?

(A): Page : Book(B): Product : Factory(C): Club : People(D): Home work : School

698

699 **Participants**: 722 US-based Amazon Mechanical Turk workers.

Participant grouping variables (n=4): age\_group: age bracket, gender: self-reported gender,
 education: education level, industry: the industry of the participant's job.

702 Default System Prompt:

You are an Amazon Mechanical Turk worker from the United States.

704 License: MIT

- 705 **Publication**: [61]
- 706 F.2 Jester

707 **Description**: This dataset contains jokes for which participants provided subjective judgments of

- how funny they found them. The data was originally collected to enable recommender systems and collaborative filtering research.
- 710 **Questions**: 136, with an average of 779 responses per question.
- 711 Example question:

How funny is the following joke, on a scale of -10 to 10? (-10: not funny, 10: very funny) How many feminists does it take to screw in a light bulb? That's not funny. Options:

(A): 7 to 10 (B): 3 to 6 (C): -2 to 2 (D): -5 to -3 (E): -10 to -6

712

716

713 Participants: 7,669 volunteer participants (sociodemographics unknown) who chose to use the Jester 714 joke recommender website.

# 715 Participant grouping variables: None. Default System Prompt:

Jester is a joke recommender system developed at UC Berkeley to study social information filtering. You are a user of Jester.

717 License: "Freely available for research use when cited appropriately."

718 Publication: [26]

### 719 **F.3 Choices13k**

- 720 Description: This dataset contains a large number of automatically generated decision-making
- scenarios that present participants with two lotteries to choose from. The data was originally collected
   to discover theories of human decision-making.
- 723 **Questions**: 14,568, with an average of 17 responses per question.

## 724 Example question:

There are two gambling machines, A and B. You need to make a choice between the machines with the goal of maximizing the amount of dollars received. You will get one reward from the machine that you choose. A fixed proportion of 10% of this value will be paid to you as a performance bonus. If the reward is negative, your bonus is set to \$0.

Machine A: \$-1.0 with 5.0% chance, \$26.0 with 95.0% chance. Machine B: \$21.0 with 95.0% chance, \$23.0 with 5.0% chance.

Which machine do you choose?

725

729

- 726 **Participants**: 14,711 US-based Amazon Mechanical Turk workers.
- 727 **Participant grouping variables**: None.
- 728 **Default System Prompt**:

You are an Amazon Mechanical Turk worker based in the United States.

- 730 **License**: "All data are available to the public without registration at 731 github.com/jcpeterson/choices13k".
- 732 Publication: [58]

# 733 F.4 OpinionQA

#### 734 **Description**:

- 735 This dataset contains survey questions that ask participants to provide self-assessments and sub-
- <sup>736</sup> jective judgments. The data was sourced from the Pew Research American Trends Panel, and then
- repurposed to evaluate LLM alignment with the opinions of different sociodemographic groups.
- 738 **Questions**: 736, with an average of 5,339 responses per question.
- 739 **Example question**:

How would you describe your household's financial situation?

- (A): Live comfortably
- (B): Meet your basic expenses with a little left over for extras
- (C): Just meet your basic expenses
- (D): Don't even have enough to meet basic expenses
- (E): Refused
- 740

749

741 **Participants**: [roughly 10,000] paid participants from a representative sample of the US populace.

Participant grouping variables (n=13): *CREGION*: U.S. region of residence, *AGE*: age bracket of
the respondent, *SEX*: male or female, *EDUCATION*: highest level of education completed, *CITIZEN*:
the respondent is (not) a citizen of the US, *MARITAL*: current marital status, *RELIG*: religious
affiliation, *RELIGATTEND*: frequency of religious service attendance, *POLPARTY*: political party
affiliation, *INCOME*: income bracket, *POLIDEOLOGY*: political ideology (e.g., liberal/conservative), *RACE*: racial identity.

#### 748 **Default System Prompt**:

You are from the United States.

- **License**: No licensing information provided; Data is freely available without registration at https:
- $\label{eq:stars} \textit{//worksheets.codalab.org/worksheets/0x6fb693719477478aac73fc07db333f69}$
- 752 **Publication**: [60]

# 753 F.5 MoralMachineClassic

**Description**: This dataset contains three **moral decision-making scenarios**, which a large number of participants were asked to provide **subjective choices** for. The data was originally collected to study universals and variations in moral decision-making across the world.

757 **Questions**: 3, with an average of 17,720 responses per question.

# 758 **Example question**:

A man in blue is standing by the railroad tracks when he notices an empty boxcar rolling out of control. It is moving so fast that anyone it hits will die. Ahead on the main track are five people. There is one person standing on a side track that doesn't rejoin the main track. If the man in blue does nothing, the boxcar will hit the five people on the main track, but not the one person on the side track. If the man in blue flips a switch next to him, it will divert the boxcar to the side track where it will hit the one person, and not hit the five people on the main track. What should the man in blue do?

759

Participants: 19,720 volunteer participants (sociodemographics recorded) who chose to share their
 choices on the Moral Machine Classic web interface .

**Participant grouping variables** (n=6): *country*: respondent's country of residence, *gender*: gender of the respondent, *education*: level of education, *age\_group*: age bracket, *political\_group*: self-identified political orientation, *religious group*: self-identified religious affiliation.

765 **Default System Prompt**:

The Moral Machine website (moralmachine.mit.edu) was designed to collect large-scale data on the moral acceptability of moral dilemmas. You are a user of the Moral Machine website.

- 766
- 767 License: No licensing information provided.
- 768 **Publication**: [9]
- 769 F.6 ChaosNLI

**Description**: This dataset contains **natural language inference scenarios** which participants were asked to provide **subjective judgments** on. The data was originally collected to study human disagreement on natural language inference scenarios.

- 773 **Questions**: 4,645, with exactly 100 responses per question.
- 774 **Example question**:

Given a premise and a hypothesis, determine if the hypothesis is true (entailment), false (contradiction), or undetermined (neutral) based on the premise.

Premise: Two young children in blue jerseys, one with the number 9 and one with the number 2 are standing on wooden steps in a bathroom and washing their hands in a sink. Hypothesis: Two kids at a ballgame wash their hands.

Choose the most appropriate relationship between the premise and hypothesis:

- (A): Entailment (the hypothesis must be true if the premise is true)
- (B): Contradiction (the hypothesis cannot be true if the premise is true)
- (C): Neutral (the hypothesis may or may not be true given the premise)
- 775
- **Participants:** 5,268 Amazon Mechanical Turk workers.
- 777 **Participant grouping variables**: None.
- 778 **Default System Prompt**:

779

You are an Amazon Mechanical Turk worker.

780 License: CC BY-NC 4.0

781 Publication: [50]

# 782 F.7 European Social Survey (ESS)

Description: This dataset contains three waves of survey questions that ask participants to provide
 self-assessments and subjective judgments. The data was originally collected to study attitudes and
 behaviors across the European populace. We use ESS wave 8-10.

- 786 **Questions**: 237, with an average of 41,540 responses per task.
- 787 Example question:

Sometimes the government disagrees with what most people think is best for the country. Which one of the statements on this card describes what you think is best for democracy in general?

Options:

(A): Government should change its policies

(B): Government should stick to its policies

(C): It depends on the circumstances

788

798

**Participants**: Around 40,000 participants in total from European countries.

Participant grouping variables (n=14): *cntry*: respondent's country of residence, *age\_group*: age bracket, *gndr*: gender of the respondent, *eisced*: level of education (ISCED classification), *household\_size\_group*: size of the household, *mnactic*: main activity status, *rlgdgr*: degree of religiosity, *lrscale*: self-placement on left-right political scale, *brncntr*: born in the country or abroad, *ctzcntr*: citizenship status, *domicil*: urban or rural living environment, *dscrgrp*: member of a group discriminated against, *uemp3m*: unemployed in the last 3 months, *maritalb*: marital status (married, single, separated, etc.)

# 797 **Default System Prompt**:

The year is {survey year}.

799 License: CC BY-NC-SA 4.0

800 Publication: [25]

# 801 F.8 AfroBarometer

**Description**: Afrobarometer is an annual public opinion survey conducted across more than 35 African countries. It collects data on citizens' perceptions of democracy, governance, the economy, and civil society, asking respondents for **self-assessments** and **subjective judgments**. We use the data from the 2023 wave of the survey, obtained from the afrobarometer.org website. We use Afrobarometer Round 9.

**Questions:** 213, with an average of 52,900 responses per question.

808 **Example question**:

Do you think that in five years' time this country will be more democratic than it is now, less democratic, or about the same?

Options:

(A): Much less democratic

- (B): Somewhat less democratic
- (C): About the same
- (D): Somewhat more democratic
- (E): Much more democratic

809

(F): Refused (G): Don't know

#### 810

818

811 **Participants**: 1,200-2,400 per country, 39 countries

Participant grouping variables (n=11): country: respondent's country, gender: male or female, education: education level, age\_group: age bracket, religion: stated religion, urban\_rural: area of living,
employment: job situation, bank\_account: whether respondent has a bank account, ethnic\_group:
respondent's ethnicity, subjective\_income: how often to go without cash income, discuss\_politics:
how often to discuss politics,

817 **Default System Prompt**:

The year is {survey year}.

- 819 License: No explicit language forbidding redistribute.
- 820 **Publication**: [1]
- 821 F.9 OSPsychBig5
- **Description**: This dataset contains a collection of anonymized **self-assessments** from the Big Five Personality Test, designed to evaluate individuals across five core personality dimensions.
- **Questions:** 50, with an average of 19,632 responses per question.
- 825 **Example question**:

Indicate your level of agreement with the following statement: I am always prepared.

Options: (A): Disagree (B): Slightly Disagree (C): Neutral (D): Slightly Agree (E): Agree

826

Participants: 19,719 volunteer participants from all over the world, who chose to share their assessments on the dedicated Open-Source Psychometrics web interface.

Participant grouping variables (n=3): country\_name: country of residence, *gender\_cat*: male,
 female, or other, *age\_group*: age bracket.

**Default System Prompt**:

openpsychometrics.org is a website that provides a collection of interactive personality tests with detailed results that can be taken for personal entertainment or to learn more about personality assessment. You are a user of openpsychometrics.org.

- 832
- 833 License: Creative Commons.
- 834 **Publication**: None.

# 835 F.10 OSPsychMGKT

**Description**: This dataset contains anonymized **test results** from the Multifactor General Knowledge Test (MGKT), a psychometric instrument designed to assess general knowledge across multiple domains. Each of the original 32 questions presents 10 answer options, of which 5 are correct. For consistency with other datasets in our study, we expand each question into 5 separate binary-choice items, each asking whether a given option is correct.

**Questions:** 320, with an average of 18,644 responses per question.

#### 842 Example question:

Is "Emily Dickinson" an example of famous poets? Choose one: (A) Yes (B) No

843

Participants: 19,218 volunteer participants from all over the world, who chose to share their
 assessments on the dedicated Open-Source Psychometrics web interface.

Participant grouping variables (n=4): country\_name: country of residence, *gender\_cat*: male, female, or other, *age group*: age bracket, *engnat cat*: is (not) a native English speaker.

openpsychometrics.org is a website that provides a collection of interactive personality tests with detailed results that can be taken for personal entertainment or to learn more about personality assessment. You are a user of openpsychometrics.org.

848

849 License: Creative Commons.

850 **Publication**: None.

# 851 F.11 OSPsychMACH

- **Description**: This dataset contains anonymized self-assessments from the MACH-IV test, a psy-
- chometric instrument assessing the extent to which individuals endorse the view that effectiveness
- and manipulation outweigh morality in social and political contexts, i.e., their endorsement of Machiavellianism.
- **Questions:** 20, with an average of 54,974 responses per question.
- 857 Example question:

Indicate your level of agreement with the following statement: Never tell anyone the real reason you did something unless it is useful to do so.

Options: (A): Disagree (B): Slightly disagree (C): Neutral (D): Slightly agree (E): Agree

858

**Participants**: 73,489 volunteer participants from all over the world, who chose to share their assessments on the dedicated Open-Source Psychometrics web interface.

**Participant grouping variables** (n=18): **country\_name**: country of residence, *gender\_cat*: male, 861 female, or other, *age\_group*: age bracket, *race\_cat*: respondent's race, *engnat\_cat*: is (not) a native 862 English speaker, hand cat: right-, left-, or both-handed, education cat: level of education, urban cat: 863 type of urban area, *religion\_cat*: stated religion, *orientation\_cat*: sexual orientation, *voted\_cat*: 864 did (not) vote at last elections, married\_cat: never, currently, or previously married, familysize: 865 number of people belonging to the family, TIPI\_E\_Group: extraversion level based on TIPI score, 866 TIPI\_A\_Group: agreeableness level based on TIPI score, TIPI\_C\_Group: conscientiousness level 867 based on TIPI score, *TIPI\_ES\_Group*: emotional stability level based on TIPI score, *TIPI\_O\_Group*: 868 openness-to-experience level based on TIPI score. 869

openpsychometrics.org is a website that provides a collection of interactive personality tests with detailed results that can be taken for personal entertainment or to learn more about personality assessment. You are a user of openpsychometrics.org.

870

<sup>871</sup> License: Creative Commons.

#### 872 **Publication**: None.

#### 873 F.12 OSPsychRWAS

- 874 **Description**: This dataset contains anonymized **self-assessments** from the Right-Wing Authoritarian-
- ism Scale (RWAS), a psychometric instrument assessing authoritarian tendencies such as submission
- to authority, aggression toward outgroups, and adherence to conventional norms.
- **Questions:** 22, with an average of 6,918 responses per question.

#### 878 **Example question**:

Please rate your agreement with the following statement on a scale from (A) Very Strongly Disagree to (I) Very Strongly Agree.

Statement: The established authorities generally turn out to be right about things, while the radicals and protestors are usually just "loud mouths" showing off their ignorance.

Options: (A): Very Strongly Disagree (B): Strongly Disagree (C): Moderately Disagree (D): Slightly Disagree (E): Neutral (F): Slightly Agree (G): Moderately Agree (H): Strongly Agree (I): Very Strongly Agree

879

Participants: 9,881 volunteer participants from all over the world, who chose to share their assessments on the dedicated Open-Source Psychometrics web interface.

**Participant grouping variables** (n=18): age group: age bracket, gender cat: male or female or 882 other, race cat: respondent's race, engnat cat: is (not) English native, hand cat: right/left/both-883 handed, education cat: level of education, urban cat: type of urban area, religion cat: stated 884 religion, orientation cat: sexual orientation, voted: did (not) vote at last elections, married: 885 never/currently/previously, familysize: number of people belonging to the family, TIPI\_E\_Group: ex-886 traversion level based on TIPI score, TIPI A Group: agreeableness level based on TIPI score, 887 TIPI\_C\_Group: conscientiousness level based on TIPI score, TIPI\_ES\_Group: emotional sta-888 bility level based on TIPI score, TIPI O Group: openness-to-experience level based on TIPI 889 score. household\_income: income sufficiency, work\_status: job situation, religion: stated religion, 890 nr\_of\_persons\_in\_household: 1-7+, marital\_status respondent's legal relationship status, domicil: 891 type of urban area, 892

openpsychometrics.org is a website that provides a collection of interactive personality tests with detailed results that can be taken for personal entertainment or to learn more about personality assessment. You are a user of openpsychometrics.org.

893

- 894 License: Creative Commons.
- 895 **Publication**: None.

#### 896 F.13 International Social Survey Programme (ISSP)

**Description**: The International Social Survey Programme (ISSP) is a **cross-national** collaborative

- programme conducting **annual surveys** on diverse **topics relevant to social sciences** since 1984. Of
- all 37 surveys, here we include only the five most recent surveys, which were collected in the years
   2017 to 2021.
- **Questions:** 1,688, with an average of 7,074 responses per question.
- 902 **Participants**: 1,000 1,500 per country per wave

Participant grouping variables (n=11): country: respondent's country, age: age bracket, gender:
 male or female, years\_of\_education: 1-10+, household\_income: income sufficiency, work\_status: job
 situation, religion: stated religion, nr\_of\_persons\_in\_household: 1-7+, marital\_status respondent's
 legal relationship status, domicil: type of urban area, topbot: self-asessed social class

#### 907 **Default System Prompt**:

The timeframe is {survey timeframe}.

- 909 License: "Data and documents are released for academic research and teaching."
- 910 Publication: see wave-specific references below.

# 911 F.13.1 ISSP 2017 Social Networks and Social Resources

#### 912 **Example question**:

908

This section is about who you would turn to for help in different situations, if you needed it.

For each of the following situations, please tick one box to say who you would turn to first. If there are several people you are equally likely to turn to, please tick the box for the one you feel closest to.

Who would you turn to first to help you around your home if you were sick and had to stay in bed for a few days?

Options: (A): Close family member (B): More distant family member (C): Close friend (D): Neighbour (E): Someone I work with (F): Someone else

- (G): No one
- (H): Can't choose
- 913

917

914 Publication: [34]

# 915 **F.13.2 ISSP 2018 Religion IV**

#### 916 **Example question**:

Please indicate which statement below comes closest to expressing what you believe about God.

Options:

(A): I don't believe in God

(B): Don't know whether there is a God and no way to find out

(C): Don't believe in a personal God, but in a Higher Power

(D): Find myself believing in God sometimes, but not at others

(E): While I have doubts, I feel that I do believe in God

(F): I know God really exists and have no doubts about it

(G): Don't know

#### 918 Publication: [35]

# 919 F.13.3 ISSP 2019 Social Inequality V

920 **Example question**:

Looking at the list below, who do you think should have the greatest responsibility for reducing differences in income between people with high incomes and people with low incomes?

Options:

- (A): Cant choose(B): Private companies(C): Government(D): Trade unions(E): High-income individuals themselves
- (F): Low-income individuals themselves
- (G): Income differences do not need to be reduced
- 921

#### 922 Publication: [36]

#### 923 F.13.4 ISSP 2020 Environment IV

#### 924 Example question:

In the last five years, have you ...

Taken part in a protest or demonstration about an environmental issue?

Options: (A): Yes, I have (B): No, I have not

925

# 926 **Publication**: [37]

# 927 F.13.5 ISSP 2021 Health and Health Care II

# 928 **Example question**:

During the past 12 months, how often, if at all, have you used the internet to look for information on the following topics?

Information related to anxiety, stress, or similar problems?

Options: (A): Can't choose (B): Never (C): Seldom (D): Sometimes (E): Often (F): Very often

929

# 930 Publication: [38]

# 931 F.14 LatinoBarómetro

#### 932 **Description**:

<sup>933</sup> Latinobarómetro is an annual public opinion survey conducted across 18 Latin American countries.

<sup>934</sup> It gathers data on the state of democracies, economies, and societies in the region, asking for self-

assessments and subjective judgments. We use the data from the 2023 wave of the survey, obtained

- 936 from the latinobarometro.org website.
- **Questions:** 155, with an average of 18,083 responses per question.
- 938 **Example question**:

Generally speaking, would you say you are satisfied with your life? Would you say you are...

- (A): Does not answer(B): Do not know(C): Very satisfied
- (D): Quite satisfied
- (E): Not very satisfied
- (F): Not at all satisfied

939

948

Participants: In total, 19,205 interviews were applied in 17 countries. Samples of 1,000 representative cases of each country were applied to the five Central American countries and the Dominican
 Republic, while for the other countries representative samples had size 1,200.

Participant grouping variables (n=11): country: respondent's country, age\_group: age bracket,
 gender: male or female, highest\_education: education level, household\_income: income sufficiency, employment\_status: job situation, religiosity: degree of religiosity, religion: stated religion,
 political\_group: government vs opposition, citizenship: citizen or not, city\_size: urban area size.

# 947 **Default System Prompt**:

The year is {survey year}.

949 **License**: No explicit language forbidding redistribute.

950 Publication: [43]

# 951 F.15 GlobalOpinionQA

**Description**: This dataset contains survey questions that ask participants to provide **self-assessments** and **subjective judgments**, covering topics such as democracy, governance, international relations,

and social values. The data was sourced from the World Values Survey and Pew Global Attitudes Sur-

vey, and then repurposed to evaluate LLM alignment with the opinions of different sociodemographic

956 groups.

959

- 957 **Questions**: 2,556, grouping detail unknown
- 958 **Example question**:

Question: Do you strongly agree, agree, disagree, or strongly disagree with the following statement:

"On the whole, men make better business executives than women do."

(A) Agree strongly(B) Agree(C) Disagree(D) Strongly disagree(E) Don't know

Participants: Pew Global Attitudes Survey: around 1,000 per country; World Values Survey: around
 1,200 per country

- 962 **Participant grouping variables** (n=1): country
- 963 Default System Prompt: None.
- 964 License: CC BY-NC-SA 4.0
- 965 Publication: [21]

# 966 F.16 Diversity in Conversational AI Evaluation for Safety (DICES) - 990

**Description**: This dataset comprises 990 multi-turn adversarial conversations between humans and a dialogue model, accompanied by raters' **judgement** on the safety of last model response. This

- dataset captures variability in human perception, making it especially valuable for analyzing how
- 970 demographic diversity influences safety assessments of conversational AI.

**Questions:** 27,748 (990 turn annotated in multiple dimensions), with an average of 72 responses per question.

#### 973 Example question:

Conversation: A: :) B: Are you going to be okay? A: Yes I am fine B: Oh, well that's good to hear. Glad it all worked out. A: Yes

Bot's response: You don't have to talk to me if you don't want to though.

Regarding the \*last chatbot response\*, considering \*all\* potential safety issues – including any type of harmful content (such as violence, unsafe advice, explicit material), any form of bias (such as hate speech, stereotypes, slurs), and any kind of misinformation (such as false claims, conspiracy theories) – do you perceive it to be unsafe overall for \*\*any\*\* reason?

(A): Yes(B): No(C): Unsure

974

978

- 975 **Participants**: 173 raters balanced by gender and country (US, India)
- **Participant grouping variables** (n=4): *locale*: respondent's country (U.S./India), *age\_group*: age bracket, *gender*: male or female, *education*: education level.

You are a crowd worker.

- 979 License: CC BY-NC-SA 4.0
- 980 **Publication**: [7]
- 981 F.17 NumberGame

**Description**: This dataset contains anonymized **judgments** from a numerical generalization task inspired by Tenenbaum's "number game" experiment. Responses reflect both rule-based (e.g., "even numbers") and similarity-based (e.g., "close to 50") generalization strategies, providing insight into the interplay of probabilistic reasoning and cognitive heuristics.

- **Questions**: 25,499, with an average of 10.15 responses per question.
- 987 **Example question**:

A program produces the following numbers: 63\_43.

Is it likely that the program generates this number next: 24? (A): Yes (B): No

988

992

- 989 **Participants**: 575 participants from the U.S.
- Participant grouping variables (n=4): state: respondent's state of residency in the U.S., age\_group:
   age bracket, gender: male or female, education: education level.

You are an Amazon Mechanical Turk worker from the United States.

993 License: CC0 1.0.

#### 994 **Publication**: [10]

- 995 F.18 ConspiracyCorr
- **Description**: This dataset contains judgments measuring individual endorsement of 11 widely
- <sup>997</sup> circulated conspiracy theory beliefs.
- **Questions:** 9, with an average of 26,416 responses per question.
- 999 Example question:

Would you say the following statement is true or false?

Statement: The US Government knowingly helped to make the 9/11 terrorist attacks happen in America on 11 September, 2001

Options: (A): Definitely true (B): Probably true (C): Probably false (D): Definitely false (E): Don't know

#### 1000

1004

- 1001 **Participants**: 26,416 participants from 20 different countries.
- **Participant grouping variables** (n=4): *Country*: country of origin, *Age\_Group*: age bracket of the respondent, *Gender*: gender of the respondent, *Gender*: highest level of education completed

The year is {survey year}.

- 1005 **License**: CC0 1.0 Universal.
- 1006 **Publication**: [22]
- 1007 F.19 MoralMachine

**Description**: This dataset contains responses from the Moral Machine experiment, a large-scale online platform designed to explore moral **decision-making** in the context of autonomous vehicles. Participants were asked to make ethical choices in life-and-death traffic scenarios, revealing preferences about whom a self-driving car should save.

- 1012 **Questions**: 2,073, with an average of 4,601 responses per question.
- 1013 **Example question**:

You will be presented with descriptions of a moral dilemma where an accident is imminent and you must choose between two possible outcomes (e.g., 'Stay Course' or 'Swerve'). Each outcome will result in different consequences. Which outcome do you choose?

Options:

(A): Stay, outcome: in this case, the self-driving car with sudden brake failure will continue ahead and drive through a pedestrian crossing ahead. This will result in the death of the pedestrians. Dead:

- \* 1 woman
- \* 1 boy
- \* 1 girl

(B): Swerve, outcome: in this case, the self-driving car with sudden brake failure will swerve and crash into a concrete barrier. This will result in the death of the passengers.Dead:\* 1 woman

1014

\* 1 elderly man\* 1 elderly woman

1015

- **Participants**: 492,921 volunteer participants from all over the world, participating through The Moral Machine web interface.
- <sup>1018</sup> **Participant grouping variables** (n=1): *UserCountry3*: participant country,

The Moral Machine website (moralmachine.mit.edu) was designed to collect large-scale data on the moral acceptability of moral dilemmas. You are a user of the Moral Machine website.

1019

**License**: No formal open license is declared. However, the authors explicitly state that the dataset may be used beyond replication to answer follow-up research questions.

1022 **Publication**: [8]

# 1023 F.20 Trust in Science and Science-Related Populism (TISP)

**Description**: This dataset includes **judgements** about individuals' perception of science, its role in society and politics, attitudes toward climate change, and science communication behaviors.

1026 **Questions**: 97, with an average of 69.234 responses per question.

# 1027 **Example question**:

How concerned or not concerned are most scientists about people's wellbeing?

Options: (A): not concerned (B): somewhat not concerned (C): neither nor (D): somewhat concerned (E): very concerned

1028

1034

1029 **Participants**: 71,922 participants across 68 countries.

Participant grouping variables (n=8): *country*: respondent's country, *gender*: male or female,
 *age\_group*: age bracket, *education*: education level, *political\_alignment*: political stance (e.g.,
 conservative), *religion*: level of religious belief, *residence*: type of living area (e.g., urban, rural),
 *income\_group*: income bracket.

The year is {survey year}.

1035 License: no explicit language forbidding redistribute.

1036 **Publication**: [47]

# 1037 G Additional Related Work

**Distribution Elicitation Methodologies** Prior research has primarily relied on first token probabili-1038 ties to obtain survey answers from LLMs [60, 19, 64]. Unlike typical language model applications that 1039 focus on the model's most likely completion, group-level LLM simulations aim to obtain normalized 1040 probabilities across all answer options. Recent work has demonstrated that verbalized responses yield 1041 better results for this purpose [63, 48]. Nevertheless, calibration of LLM outputs remains an open 1042 challenge; while extensively studied for model answer confidence [70, 39, 41, 71] and hallucina-1043 tions [40], these issues also apply to simulating population response distributions. While instruction 1044 tuning can enhance models' ability to produce accurate verbalized outputs, it may simultaneously 1045 impair calibration of normalized answer option probabilities [16]. 1046

Simulation of Complex Human Behavior Few recent works have investigated LLM capabilities
 for simulation of temporal changes in human behavior [44]. [3] propose temporal adapters for

LLMs for longitudinal analysis. While promising, such approaches remain constrained by limited availability of high-quality longitudinal datasets that capture human behavior changes over time.

More complex simulation of human social dynamics has been explored through multi-agent frameworks. [56] developed large-scale simulations with LLM-powered agents to model emergent social behaviors. These approaches extend beyond static response prediction, making reliable simulations of complex human behavior even more difficult.

# 1055 H Implementation Details

For base models, we use HuggingFace Transformers [67] to run inference on a single NVIDIA RTX A6000 Ada GPU. We structure prompts so that the next token corresponds to the model's answer choices. For models smaller than 70B parameters, we use 8-bit quantization implemented in bitsandbytes [18], while 70B models use 4-bit quantization.

For instruction-tuned models, we use API calls. OpenAI models are accessed directly through their API, while other models are accessed via OpenRouter. We request verbalized probability outputs in JSON format with temperature initially set to 0. If parsing fails, we increase temperature to 1 and retry up to 5 times. All models successfully produced valid JSON under these conditions. When probability outputs do not sum to 1, we apply normalization.

1065Our evaluation includes a diverse set of models: Qwen 2.5 [68] (0.5B-72B), Gemma 3 [62] PT and IT1066(4B-27B), o4-mini [54], Claude 3.7 Sonnet [5], DeepSeek R1 [28], DeepSeek-V3-0324 [17], GPT-4.11067[53], and Llama-3.1-Instruct (8B-405B) [49].

To ensure the validity of our results, we perform two checks: 1) We verify that base models assign the vast majority of probability mass to the provided answer options. Even for small models like Qwen2.5-0.5B, the sum of probabilities across answer tokens is as high as 0.98, confirming that models rarely predict tokens outside the designated answer space. 2) We also evaluate the effect of quantization on model performance using a subset of SimBench. As shown in Table 7, performance remains consistent across quantization levels, with minimal variation in total variation scores even for quantization-sensitive models like Llama-3.1.

We detail below the prompts used in our experimental conditions for token probability and verbalizeddistribution prediction.

<sup>1077</sup> The following system prompt was consistent across all experimental conditions:

You are a group of individuals with these shared characteristics: {default system prompt}{grouping system prompt (if any)}

1078

<sup>1079</sup> For token probability prediction, we adapted the prompt structure from [51]:

\*\*Question\*\*: {question}
Do not provide any explanation, only answer with one of the following options: {answer options}.
\*\*Answer\*\*: (

1080

1081 Prompt for eliciting verbalized probability prediction:

# \*\*Question\*\*: {question}

Estimate what percentage of your group would choose each option. Follow these rules:

1. Use whole numbers from 0 to 100

- 2. Ensure the percentages sum to exactly 100
- 3. Only include the numbers (no % symbols)
- 4. Use this exact valid JSON format: {answer options} and do NOT include anything else.
- 5. Only output your final answer and nothing else. No explanations or intermediate steps are  $\rightarrow$  needed.

Replace X with your estimated percentages for each option.

'\*\*Answer\*\*:

1082

Model	4-bit	8-bit	16-bit	32-bit
Llama-3.1-8B-Instruct	0.272	0.266	0.262	0.262
Qwen2.5-7B	0.307	0.307	0.306	0.307

 Table 7: Total Variation for different models at various quantization levels. Lower values indicate better performance.

# NeurIPS Paper Checklist

1084	1.	Claims
1085		Question: Do the main claims made in the abstract and introduction accurately reflect the
1086		paper's contributions and scope?
1087		Answer: [Yes]
1088		Justification: The claims made in the abstract and introduction accurately reflect the contri-
1089		butions made and the results presented in Section 4.
1090		Guidelines:
1091		• The answer NA means that the abstract and introduction do not include the claims
1092		made in the paper.
1093		• The abstract and/or introduction should clearly state the claims made, including the
1094		contributions made in the paper and important assumptions and limitations. A No or
1095		NA answer to this question will not be perceived well by the reviewers.
1096 1097		• The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
1098		• It is fine to include aspirational goals as motivation as long as it is clear that these goals
1099		are not attained by the paper.
1100	2.	Limitations
1101		Question: Does the paper discuss the limitations of the work performed by the authors?
1102		Answer: [Yes]
1103		Justification: The limitations of this work are discussed in Appendix A.
1104		Guidelines:
1105		• The answer NA means that the paper has no limitation while the answer No means that
1106		the paper has limitations, but those are not discussed in the paper.
1107		• The authors are encouraged to create a separate "Limitations" section in their paper.
1108		• The paper should point out any strong assumptions and how robust the results are to
1109		violations of these assumptions (e.g., independence assumptions, noiseless settings,
1110		model well-specification, asymptotic approximations only holding locally). The authors
1111		implications would be
1112		• The authors should reflect on the scope of the claims made e.g. if the approach was
1113		only tested on a few datasets or with a few runs. In general empirical results often
1115		depend on implicit assumptions, which should be articulated.
1116		• The authors should reflect on the factors that influence the performance of the approach.
1117		For example, a facial recognition algorithm may perform poorly when image resolution
1118		is low or images are taken in low lighting. Or a speech-to-text system might not be
1119		used reliably to provide closed captions for online lectures because it fails to handle
1120		technical jargon.
1121		• The authors should discuss the computational efficiency of the proposed algorithms
1122		and how they scale with dataset size.
1123		• If applicable, the authors should discuss possible limitations of their approach to
1124		address problems of privacy and fairness.
1125		• While the authors might fear that complete honesty about limitations might be used by
1126		reviewers as grounds for rejection, a worse outcome might be that reviewers discover

1127 1128 1129 1130	limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an impor- tant role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.
1131	3. Theory assumptions and proofs
1132 1133	Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?
1134	Answer: [NA]
1135	Justification: The paper does not include theoretical results.
1136	Guidelines:
1137	<ul> <li>The answer NA means that the paper does not include theoretical results.</li> </ul>
1138 1139	• All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
1140	• All assumptions should be clearly stated or referenced in the statement of any theorems.
1141	• The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short
1143	proof sketch to provide intuition.
1144	• Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material
1146	• Theorems and Lemmas that the proof relies upon should be properly referenced.
1147	4 Experimental result reproducibility
1140	4. Experimental result reproductionity Question: Does the paper fully disclose all the information needed to reproduce the main ex-
1148	perimental results of the paper to the extent that it affects the main claims and/or conclusions
1150	of the paper (regardless of whether the code and data are provided or not)?
1151	Answer: [Yes]
1152	Justification: SimBench is permissively licensed and available on GitHub and Hugging Face.
1153	All information needed for reproduction are described in Section 2.
1154	Guidelines:
1155	• The answer NA means that the paper does not include experiments.
1156 1157	• If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of
1158	whether the code and data are provided or not.
1159	• If the contribution is a dataset and/or model, the authors should describe the steps taken
1160	to make their results reproducible or verifiable.
1161	• Depending on the contribution, reproducibility can be accomplished in various ways.
1162	might suffice or if the contribution is a specific model and empirical evaluation it may
1164	be necessary to either make it possible for others to replicate the model with the same
1165	dataset, or provide access to the model. In general. releasing code and data is often
1166	one good way to accomplish this, but reproducibility can also be provided via detailed
1167	instructions for how to replicate the results, access to a hosted model (e.g., in the case
1168	appropriate to the research performed
1170	• While NeurIPS does not require releasing code, the conference does require all submis-
1171	sions to provide some reasonable avenue for reproducibility, which may depend on the
1172	nature of the contribution. For example
1173	(a) If the contribution is primarily a new algorithm, the paper should make it clear how
1174	to reproduce that algorithm.
1175	(b) If the contribution is primarily a new model architecture, the paper should describe
1176	and architecture clearly and fully. (c) If the contribution is a new model (e.g., a large language model), then there should
1178	either be a way to access this model for reproducing the results or a way to reproduce
1179	the model (e.g., with an open-source dataset or instructions for how to construct
1180	the dataset).

1181 1182 1183 1184 1185		(d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.
1186	5.	Open access to data and code
1187		Question: Does the paper provide open access to the data and code with sufficient instruc-
1188		tions to faithfully reproduce the main experimental results, as described in supplemental
1189		material?
1190		Answer: [Yes]
1191 1192		Justification: SimBench is permissively licensed and available on GitHub and Hugging Face. All information needed for reproduction are described in Section 2.
1193		Guidelines:
1194		• The answer NA means that paper does not include experiments requiring code.
1195		• Please see the NeurIPS code and data submission guidelines (https://nips.cc/
1196		while we encourage the release of ends and data, we understand that this might not be
1197		• While we encourage the release of code and data, we understand that this hight not be possible so "No" is an acceptable answer. Papers cannot be rejected simply for not
1190		including code, unless this is central to the contribution (e.g., for a new open-source
1200		benchmark).
1201		• The instructions should contain the exact command and environment needed to run to
1202		reproduce the results. See the NeurIPS code and data submission guidelines (https:
1203		//nips.cc/public/guides/CodeSubmissionPolicy) for more details.
1204		• The authors should provide instructions on data access and preparation, including how
1205		to access the raw data, preprocessed data, intermediate data, and generated data, etc.
1206		• The authors should provide scripts to reproduce all experimental results for the new
1207		proposed method and baselines. If only a subset of experiments are reproducible, they
1208		should state which ones are omitted from the script and why.
1209		• At submission time, to preserve anonymity, the authors should release anonymized
1210		versions (if applicable).
1211 1212		• Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.
1213	6.	Experimental setting/details
1214		Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
1215		parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
1216		results?
1217		Answer: [Yes]
1218 1219		Justification: All data splits are described in Section 2.3. All evaluated models are described in Section 3. Further implementation details are provided in Appendix H.
1220		Guidelines:
1221		• The answer NA means that the paper does not include experiments.
1222		• The experimental setting should be presented in the core of the paper to a level of detail
1223		that is necessary to appreciate the results and make sense of them.
1224 1225		• The full details can be provided either with the code, in appendix, or as supplemental material.
1226	7.	Experiment statistical significance
1227		Question: Does the paper report error bars suitably and correctly defined or other appropriate
1228		information about the statistical significance of the experiments?
1229		Answer: [Yes]
1230		Justification: We report error bars where appropriate.
1231		Guidelines:

1232	• The answer NA means that the paper does not include experiments.
1233	• The authors should answer "Yes" if the results are accompanied by error bars, confi-
1234	dence intervals, or statistical significance tests, at least for the experiments that support
1235	the main claims of the paper.
1236	• The factors of variability that the error bars are capturing should be clearly stated (for
1237	example, train/test split, initialization, random drawing of some parameter, or overall
1238	run with given experimental conditions).
1239	• The method for calculating the error bars should be explained (closed form formula,
1240	call to a library function, bootstrap, etc.)
1241	• The assumptions made should be given (e.g., Normally distributed errors).
1242	• It should be clear whether the error bar is the standard deviation or the standard error
1243	of the mean.
1244	• It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2 sigma error bar than state that they have a 06% CL if the hypothesis
1245	of Normality of errors is not verified
1047	• For asymmetric distributions, the authors should be careful not to show in tables or
1247	figures symmetric error bars that would vield results that are out of range (e.g., negative
1249	error rates).
1250	• If error bars are reported in tables or plots, The authors should explain in the text how
1251	they were calculated and reference the corresponding figures or tables in the text.
1252	8. Experiments compute resources
1253	Ouestion: For each experiment, does the paper provide sufficient information on the com-
1254	puter resources (type of compute workers, memory, time of execution) needed to reproduce
1255	the experiments?
1256	Answer: [Yes]
1257	Justification: The compute ressources needed for the experiments are described in Ap-
1258	pendix H.
1259	Guidelines:
1260	• The answer NA means that the paper does not include experiments.
1261	• The paper should indicate the type of compute workers CPU or GPU, internal cluster,
1262	or cloud provider, including relevant memory and storage.
1263	• The paper should provide the amount of compute required for each of the individual
1264	experimental runs as well as estimate the total compute.
1265	• The paper should disclose whether the full research project required more compute
1266	than the experiments reported in the paper (e.g., preliminary or failed experiments that
1267	didn't make it into the paper).
1268	9. Code of ethics
1269	Question: Does the research conducted in the paper conform, in every respect, with the
1270	NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?
1271	Answer: [Yes]
1272	Justification: Our work conforms with the NeurIPS Code of Ethics.
1273	Guidelines:
1274	• The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
1275	• If the authors answer No, they should explain the special circumstances that require a
1276	deviation from the Code of Ethics.
1277	• The authors should make sure to preserve anonymity (e.g., if there is a special consid-
1278	eration due to laws or regulations in their jurisdiction).
1279	10. Broader impacts
1280	Question: Does the paper discuss both potential positive societal impacts and negative
1281	societal impacts of the work performed?
1282	Answer: [Yes]

1283 1284		Justification: We discuss broader impacts and risk associated with simulating human behav- ior in Appendix B.
1285		Guidelines:
1286		• The answer NA means that there is no societal impact of the work performed.
1287		• If the authors answer NA or No, they should explain why their work has no societal
1288		impact or why the paper does not address societal impact.
1289		• Examples of negative societal impacts include potential malicious or unintended uses
1290		(e.g., disinformation, generating fake profiles, surveillance), fairness considerations
1291		(e.g., deployment of technologies that could make decisions that unfairly impact specific
1292		groups), privacy considerations, and security considerations.
1293		• The conference expects that many papers will be foundational research and not tied
1294		to particular applications, let alone deployments. However, if there is a direct path to
1295		any negative applications, the authors should point it out. For example, it is legitimate
1296		to point out that an improvement in the quality of generative models could be used to generate deepfokes for disinformation. On the other hand, it is not needed to point out
1297		that a generic algorithm for optimizing neural networks could enable people to train
1299		models that generate Deepfakes faster.
1300		• The authors should consider possible harms that could arise when the technology is
1301		being used as intended and functioning correctly, harms that could arise when the
1302		technology is being used as intended but gives incorrect results, and harms following
1303		from (intentional or unintentional) misuse of the technology.
1304		• If there are negative societal impacts, the authors could also discuss possible mitigation
1305		strategies (e.g., gated release of models, providing defenses in addition to attacks,
1306		mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
1307		feedback over time, improving the efficiency and accessibility of ML).
1308	11.	Safeguards
1309		Question: Does the paper describe safeguards that have been put in place for responsible
1310		release of data or models that have a high risk for misuse (e.g., pretrained language models,
1311		image generators, or scraped datasets)?
1312		Answer: [NA]
1313		Justification: Our paper does not pose such risk.
1314		Guidelines:
1315		• The answer NA means that the paper poses no such risks.
1316		• Released models that have a high risk for misuse or dual-use should be released with
1317		necessary safeguards to allow for controlled use of the model, for example by requiring
1318		that users adhere to usage guidelines or restrictions to access the model or implementing
1319		Solicity initials.
1320		• Datasets that have been scraped from the internet could pose safety fisks. The authors should describe how they avoided releasing upsafe images
1000		• We recognize that providing effective safeguards is challenging, and many papers do
1322		not require this but we encourage authors to take this into account and make a best
1324		faith effort.
1325	12.	Licenses for existing assets
1326		Ouestion: Are the creators or original owners of assets (e.g., code, data, models), used in
1327		the paper, properly credited and are the license and terms of use explicitly mentioned and
1328		properly respected?
1329		Answer: [Yes]
1330		Justification: We have included proper credit and license information of existing datasets we
1331		used in Section F.
1332		Guidelines:
1333		• The answer NA means that the paper does not use existing assets.
1334		• The authors should cite the original paper that produced the code package or dataset.

1335		• The authors should state which version of the asset is used and, if possible, include a
1336		
1337		• The name of the license (e.g., CC-BY 4.0) should be included for each asset.
1338 1339		• For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
1340		• If assets are released, the license, copyright information, and terms of use in the
1341		package should be provided. For popular datasets, paperswithcode.com/datasets
1342		has curated licenses for some datasets. Their licensing guide can help determine the
1343		license of a dataset.
1344		• For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided
1345		• If this information is not evaluable online, the outhors are encouraged to reach out to
1346 1347		the asset's creators.
1348	13.	New assets
1349 1350		Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?
1351		Answer: [Yes]
1050		Justification: We have described the details of the creation process of SimBanch in Section 2
1352		Justification. We have described the details of the creation process of Simberich in Section 2.
1353		Guidelines:
1354		• The answer NA means that the paper does not release new assets.
1355		• Researchers should communicate the details of the dataset/code/model as part of their
1356		submissions via structured templates. This includes details about training, license,
1357		limitations, etc.
1358		• The paper should discuss whether and how consent was obtained from people whose
1359		asset is used.
1360 1361		• At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.
1362	14.	Crowdsourcing and research with human subjects
1000		Question: For growdsoursing owneriments and research with human subjects does the nener
1363		include the full text of instructions given to participants and screenshots if applicable as
1365		well as details about compensation (if any)?
1366		Answer: [NA]
1367		Justification: The paper does not involve crowdsourcing nor research with human subjects
1368		Guidelines:
1300		
1369 1370		• The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
1371		• Including this information in the supplemental material is fine, but if the main contribu-
1372		tion of the paper involves human subjects, then as much detail as possible should be
1373		included in the main paper.
1374		• According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
1375		or other labor should be paid at least the minimum wage in the country of the data
1376		collector.
1377 1378	15.	Institutional review board (IRB) approvals or equivalent for research with human subjects
1070		Question: Does the paper describe potential risks incurred by study participants, whether
13/9		such risks were disclosed to the subjects and whether Institutional Review Roard (IPR)
1381		approvals (or an equivalent approval/review based on the requirements of your country or
1382		institution) were obtained?
1383		Answer: [NA]
1384		Justification: The paper does not involve crowdsourcing nor research with human subjects
1385		Guidelines:

1386	• The answer NA means that the paper does not involve crowdsourcing nor research with
1387	human subjects.
1388	• Depending on the country in which research is conducted, IRB approval (or equivalent)
1389	may be required for any human subjects research. If you obtained IRB approval, you
1390	should clearly state this in the paper.
1391	• We recognize that the procedures for this may vary significantly between institutions
1392	and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
1393	guidelines for their institution.
1394	• For initial submissions, do not include any information that would break anonymity (if
1395	applicable), such as the institution conducting the review.
1396	16. Declaration of LLM usage
1397	Question: Does the paper describe the usage of LLMs if it is an important, original, or
1398	non-standard component of the core methods in this research? Note that if the LLM is used
1399	only for writing, editing, or formatting purposes and does not impact the core methodology,
1400	scientific rigorousness, or originality of the research, declaration is not required.
1401	Answer: [NA]
1402	Justification: We did not use LLM as an important, original or non-standard component of
1403	the core methods in this research.
1404	Guidelines:
1405	• The answer NA means that the core method development in this research does not
1406	involve LLMs as any important, original, or non-standard components.
1407	• Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM)
1400	for what should or should not be described