

# Belief-Sim: Towards Belief-Driven Simulation of Demographic Misinformation Susceptibility

Anonymous ACL submission

## Abstract

Misinformation is a growing societal threat, and susceptibility to misinformative claims varies across demographic groups due to differences in underlying beliefs. As Large Language Models (LLMs) are increasingly used to simulate human behaviors, we investigate whether they can simulate demographic misinformation susceptibility, treating *beliefs* as a primary driving factor. We introduce BELIEF-SIM, a simulation framework that constructs demographic belief profiles using psychology-informed taxonomies and survey priors. We study prompt-based conditioning and post-training adaptation, and conduct a multi-fold evaluation using: (i) susceptibility accuracy and (ii) counterfactual demographic sensitivity. Across both datasets and modeling strategies, we show that beliefs provide a strong prior for simulating misinformation susceptibility, with accuracy up to 92%.

## 1 Introduction

Misinformation is a critical challenge in today’s information ecosystem, with impacts that vary substantially across demographic groups (Khachaturov et al., 2025; Timm et al., 2025). These differences reflect not only differential targeting (Ribeiro et al., 2019), but also variation in how content is perceived (Guess et al., 2019a). Prior work illustrates these complexities: younger adults were more susceptible to believing COVID-related misinformation in the UK and Brazil (Vijaykumar et al., 2021), while older adults were more likely to share false articles on social media (Guess et al., 2019b; Brashier and Schacter, 2020). Such mixed findings suggest that demographics alone are insufficient to explain susceptibility, i.e., how likely someone is to believe a misinformative claim. Prior studies show that beliefs may play a stronger role (Sultan et al., 2024): people are more likely to accept false information when it aligns or coincides with beliefs

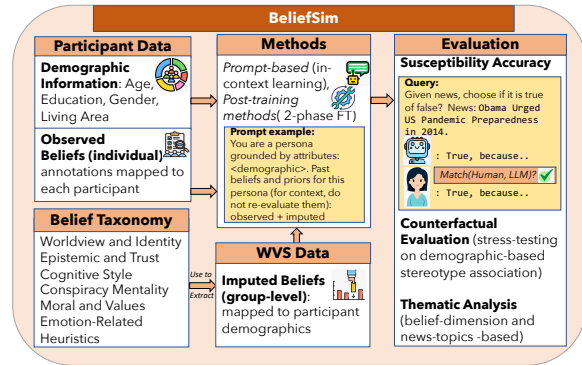


Figure 1: BeliefSim Framework: (1) Participant Data, Observed and Imputed Beliefs (based on Belief Taxonomy) are collected from surveys, (2) Methods consist of prompt-conditioning and post-training adaptation and (3) Evaluation using Susceptibility Accuracy, Counterfactual and Thematic Analysis.

they already hold, a phenomenon known as “belief consistency” (Flynn et al., 2017; Taber and Lodge, 2006; Roozenbeek et al., 2020). Therefore, understanding how demographic differences in misinformation susceptibility relate to underlying beliefs is essential for designing effective interventions.

LLMs are increasingly used to simulate social processes (Park et al., 2023; Zhou et al., 2023; Borah et al., 2025a), including behaviors such as biased endorsement and echo chambers (Acerbi and Stubbersfield, 2023; Nehring et al., 2024; Borah et al., 2025b; Sharma et al., 2024). Closest prior work on demographic misinformation susceptibility simulation relies on coarse personas (Borah et al., 2025b), without modeling belief structures. Conversely, most belief-based studies primarily align LLMs with broad political ideologies or moral foundations (Santurkar et al., 2023; Argyle et al., 2023), or simulate general personality traits (Pratelli and Petrocchi, 2025).

Motivated by these gaps, we ask three research questions: (1) Do beliefs improve demographic-aware simulations of misinformation susceptibility? (2) How can such simulations be rigorously

Dimension	Description	Example from WVS
Worldview & Identity	Group-based and political identities shape how individuals interpret and endorse information.	How proud are you to be from your country?
Epistemic Trust	Beliefs about what sources are credible and how knowledge should be evaluated (e.g., trust in science, media, institutions).	Would you say that the world is better off, or worse off, because of science?
Cognitive Style	Preferred thinking approach, ranging from analytical, accuracy-driven reasoning to intuitive impression-based judgments.	Nowadays, does one often have trouble deciding which moral rules to follow?
Conspiracy Mentality	General tendency to perceive hidden plots or malevolent intent behind significant social or political events.	How many people who are in the media do you believe are involved in corruption?
Morals & Values	Core moral values and cultural worldviews influence judgments of right, wrong, fairness, and societal norms.	Is it justifiable - claiming government benefits to which you are not entitled?
Emotion-Related	Emotion-driven beliefs related to fear, perceived risk, or threat sensitivity that heighten susceptibility to alarming misinformation.	How satisfied are you with your life as a whole these days?
Heuristics	Reliance on mental shortcuts (e.g., familiarity implies truth) that increase vulnerability to repeated or simplified claims.	Please indicate how much you use it: Daily newspaper, TV News, Radio News.

Table 1: Belief taxonomy for modeling demographic-aware misinformation susceptibility in LLMs. We provide representative WVS items as examples of how survey questions align with each belief dimension.

evaluated using utility and counterfactual demographic analyses? and (3) What modeling strategies are best suited for LLM-based simulation?

To answer these questions, we propose **Belief-Sim**, a belief-driven simulation framework for demographic misinformation susceptibility, illustrated in Fig 1. We summarize our contributions as follows: **(1)** We construct a belief taxonomy and build a simulation dataset by aggregating data from existing surveys and studies; **(2)** We conduct empirical analyses to identify factors across belief and demographic dimensions that are most important for the simulation; **(3)** We explore simulation techniques including prompt-based and post-training approaches using several LLMs; and **(4)** We perform a counterfactual study to understand when and how demographic information may provide useful priors vs stereotype-like sensitivity. Finally, we outline actionable steps that can help design future intervention methods.

## 2 Related Work

Psychological research has established that susceptibility to misinformation is heavily driven by *belief consistency*. While demographic factors like age, gender, and education provide significant predictive signals, recent large-scale meta-analyses suggest these effects can be context- and belief-dependent Sultan et al. (2024). However, most studies that establish benchmarks with large-scale human annotations face an inherent scalability bottleneck (Maertens et al., 2024; Borah et al., 2025b). High-quality human data collection is resource-

intensive and slow, often limiting research to static snapshots of specific news cycles (e.g. COVID-19) or narrow domains.

The emerging capability of LLMs to simulate human behavior offers a potential solution. Early work showed that LLMs could act as “silicon subjects” via demographic prompting (Argyle et al., 2023). Studies also investigated if LLMs reproduce classic social science findings and moral values (Park et al., 2023; Borah et al., 2025a; Nair and Wang, 2025). However, fidelity often degrades when simulations depend only on demographic labels (Giorgi et al., 2024), leading to stereotypical associations (Borah and Mihalcea, 2024).

To move beyond demographics, recent research has focused on conditioning models with richer context such as beliefs (Moon et al., 2024; Namikoshi et al., 2024). Most relevant to our work, Chuang et al. (2024) introduced “Human Belief Networks”, showing that seeding an agent with even a single belief (e.g., a stance on welfare) was more predictive of downstream responses than the agent’s entire demographic profile. However, prior belief-focused studies have largely not examined their use in simulating misinformation susceptibility. Our work bridges this gap by focusing on this specific domain, where the interplay between demographics and beliefs is highly complex.

## 3 Belief Taxonomy

Beliefs are an important tool for modeling demographic simulations for misinformation. Prior work has shown that belief dimensions, such as conspir-

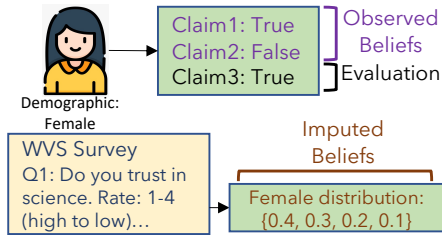


Figure 2: Simulation Data Example.

129 acy beliefs or trust in science, can help predict sus- 167  
 130 ceptibility to misinformation (Ecker et al., 2022; 168  
 131 Munusamy et al., 2024). 169

132 Grounded in psychological and cognitive science 170  
 133 research, we compiled a belief taxonomy shown 171  
 134 in Table 1, that includes seven core dimensions 172  
 135 most associated with misinformation susceptibil- 173  
 136 ity: (1) Worldview and Identity Beliefs (Kahan, 174  
 137 2017; Van Bavel et al., 2021), (2) Epistemic Trust 175  
 138 Beliefs (De Coninck et al., 2021; Lewandowsky 176  
 139 et al., 2023), (3) Cognitive style (Pennycook and 177  
 140 Rand, 2019; Ecker et al., 2022), (4) Conspiracy 178  
 141 mentality (Douglas et al., 2019; De Coninck et al., 179  
 142 2021), (5) Moral and Value Beliefs (D’Errico et al., 180  
 143 2022; Yang et al., 2024), (6) Emotion-Related Be- 181  
 144 liefs (Brady et al., 2017; McLoughlin et al., 2024) 182  
 145 and (7) Heuristic Beliefs (Lin et al., 2016; Fazio, 183  
 146 2020). While the above dimensions have tradition- 184  
 147 ally been studied in isolation and within human 185  
 148 populations, we unify them into a structured taxon-  
 149 omy designed for computational modeling (Further  
 150 details are provided in Appendix A. This enables  
 151 systematic belief simulation by providing a frame-  
 152 work for organizing belief data and evaluating it.

## 4 Simulation Data

153 For our study, we consider four demographic 186  
 154 axes: Gender (female/male), Age (younger: 187  
 155  $\leq 35$  years/older:  $\geq 60$  years), Living Area 188  
 156 (Rural/Urban), and Education (completed high 189  
 157 school/not completed high school).<sup>1</sup> These demo- 190  
 158 graphics are commonly associated with systematic 191  
 159 differences in media exposure, institutional trust, 192  
 160 etc., all of which influence vulnerability to mis- 193  
 161 information (Allcott and Gentzkow, 2017; Guess 194  
 162 et al., 2019b; Allcott et al., 2019; Anspach and 195  
 163 Carlson, 2024). Focusing on these axes enables 196  
 164 controlled analysis of demographic context for sim- 197  
 165 ulating susceptibility signals in LLMs. 198

<sup>1</sup>For ease of analysis, we focus on binary groupings, leav-  
 ing finer-grained analysis to future work.

## 4.1 Evaluation Data

167 For susceptibility evaluation, we use two ground- 168  
 169 truth datasets containing human judgments of 170  
 171 whether they believe a given claim: (1) PANDORA 171  
 172 Dataset from Borah et al. (2025b), containing an- 172  
 173 notations from 318 participants. Each participant 173  
 174 provided judgments on 3 distinct claims, along 174  
 175 with demographic information including age, gen- 175  
 176 der, living area, and education. These claims are 176  
 177 collected from RumorEval (Gorrell et al., 2019), 177  
 178 which consists of true or false rumors covering 178  
 179 eight major news events and natural disaster events; 179  
 180 (2) MIST dataset from Maertens et al. (2024). From 180  
 181 MIST, we use only Study 1 (MIST-1) which in- 181  
 182 cludes 409 participants, each providing judgment 182  
 183 on the same 100 claims, and demographic informa- 183  
 184 tion including age, gender, and education. From 184  
 185 the two datasets combined, we obtain 13.8K claims 185  
 for evaluation.

## 4.2 Belief Data

186 We collect two complementary belief signals, con- 187  
 188 sisting of individually observed belief judgments 188  
 189 and group-level (demographic) belief distributions: 189  
 190 (1) **Observed Data** - claims that were directly 190  
 191 judged by participants in PANDORA and MIST-1 191  
 192 datasets. These responses capture individual-level 192  
 193 belief judgments, reflecting each participant’s per- 193  
 194 sonal stance rather than demographic group aver- 194  
 195 ages. We use two claim judgments as observed be- 195  
 196 liefs for each participant (keeping it separate from 196  
 197 the evaluation data). This leads to 27.6 claim judg- 197  
 198 ments as observed beliefs. 198

199 (2) **Imputed Data** - inferred from the World Val- 199  
 200 ues Survey Wave 7<sup>2</sup> distributions. Imputed data 200  
 201 represent demographic belief priors (group-level), 201  
 202 inferred from WVS, conditioned solely on demo- 202  
 203 graphic attributes. We map these imputed belief 203  
 204 items to our belief taxonomy using exploratory 204  
 205 factor analysis. This yields 126 imputed belief 205  
 206 questions and corresponding demographic distribu- 206  
 207 tions. Table 1 shows examples, and Appendix A.2 207  
 208 contains mapping details and all questions. 208

209 Fig. 2 provides an example of simulation data. We 209  
 210 restrict all our analysis to U.S.-based participants 210  
 211 and English headlines, consistent with the scope of 211  
 212 the available datasets, while still leveraging varia- 212  
 213 tion in living area, age, gender, and education. 213

<sup>2</sup><https://www.worldvaluessurvey.org/WVSDocumentationWV7.jsp>

## 5 Prompt-based Conditioning for Susceptibility Simulation

We first use prompt-based methods, where each prompt conditions the model on participant demographics and belief signals (both observed and imputed) from the surveys. Note that we simulate each demographic axis separately, as it enables controlled analysis of demographic-level differences.

**Method.** For WVS questions, answers are either Likert scales (Bertram, 2007) or Yes/No questions. Therefore, we input modal (most frequent) responses in the prompt as beliefs. For example, consider the WVS question “How important is religion in your life?”, rated on a 1–10 scale. If a demographic group most frequently responds with a value of 3, we use this modal response in the prompt. Given demographic group  $\langle d \rangle$ , and belief annotations  $\langle b \rangle$ , we use: You are a persona grounded by attributes:  $\langle \text{demographic group} \rangle$ . Past beliefs and priors for this persona (for context, do not re-evaluate them):  $\langle \text{beliefs} \rangle$ . When judging a claim, stay consistent with this persona’s prior beliefs where reasonable.

We perform experiments under four primary input conditions for prompt-based conditioning: (1) zero-shot (without any demographic or belief information), (2) demographic information only, (3) belief information only, and (4) both demographics+beliefs. We additionally ablate over belief dimensions (the seven in our taxonomy) and belief sources (observed vs. imputed), as well as over demographic attributes (analyzing groups separately). These experiments yield several insights across 12 total settings. We perform these experiments using three LLMs: llama-3-8b-instr., qwen-2.5-14b-instr. and mistral-7b-instr.-v02, and average our results across three separate runs.

**Evaluation.** We measure susceptibility accuracy: whether the LLM correctly predicts if an individual would judge a claim as true or false.

### 5.1 Results

Fig 3 shows the susceptibility accuracies across the belief and demographic configurations for both datasets. For cases that include Imputed beliefs, we only report best results by selecting the belief dimension that achieves the highest accuracy. Findings show that incorporating belief information consistently improves performance over the

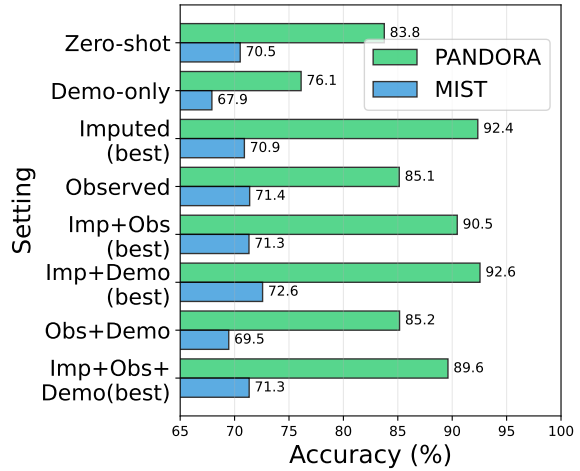


Figure 3: Susceptibility Accuracies by settings and datasets, averaged across demographic groups and models. Imputed + Demo(graphic) (best) performs the best, with all belief included settings doing better than zero-shot and demo-only.

zero-shot and demographic-only baselines. We next analyze the impact of specific belief types, demographic ablations, datasets, and model performances. Appendix C contains detailed results for all settings.

**Demographic vs Belief information.** Overall, beliefs drive most of the positive gains. Demo-only degrades zero-shot performances for both datasets, while beliefs-only (imputed) show larger gains (+16) points over demo-only and zero-shot settings. Imputed + Demo achieves the highest performance. This shows that what people *believe* may be more predictive than *who they are*. This is also in line with findings from Borah et al. (2025a), which show that belief-based signals are stronger than demographic signals in LLMs.

**Imputed vs Observed Beliefs.** Across settings, imputed beliefs consistently outperform observed beliefs, more so for PANDORA. A likely reason is that observed beliefs, being tied to specific participants, can be sparse, or mismatched to the new claim, while imputed demographic beliefs from WVS may act as *smoother population priors*. Furthermore, combining imputed+observed helps but is still lower than imputed-only. This may be due to conflict or being on different scales, and the individualistic observed signal can dilute the cleaner imputed prior. Average accuracy trends are:  $imputed > imputed + observed > observed$ .

**Imputed Belief Ablations.** Adding one belief dimension at a time to the prompt consistently yields higher performance than including all belief dimensions together. Across datasets, emotion-related

and moral-values are the strongest. These also vary by demographics: emotion-related is the strongest in Gender; moral-values is the strongest for Rural/Urban and Education; heuristics/cognitive beliefs are the strongest in Education (Detailed results are in Appendix C.2).

**Demographic Ablations.** Adding demographic information on top of beliefs yields small gains (in the range 0.2-1.2%), especially for imputed beliefs. Thus suggesting that demographics serve as an additional context signal, helping with relevant belief priors. On the effect of belief across demographic groups, accuracies for age, education and living area are improved with belief addition, specifically using the imputed beliefs. Gender shows the largest sensitivity to which beliefs are included; prior work also show that gender effects are often smaller/mixed in misinformation contexts (Sultan et al., 2024). This underscores the need for caution in selecting belief evidence for gender-specific cases.

DATASET	LLAMA	QWEN	MISTRAL
PANDORA	85.90	88.84	87.51
MIST	69.77	73.13	72.71

Table 2: Model Performance comparison across datasets, accuracy trends: qwen > mistral > llama.

**Model and Dataset comparison.** Across models, Qwen consistently shows higher accuracies, followed by mistral and llama, as shown in Table 2.

Across datasets, the PANDORA Dataset has higher accuracy scores than the MIST dataset across all setups. PANDORA yields a high zero-shot baseline (83.8%) and gains the highest from belief integration in the (best imputed + demographic) setting (+8.8). In contrast, MIST has lower zero-shot performance (70.5%), with the highest gain in (best imputed + demographic) setting (+2). Therefore, the impact of adding beliefs also varies with dataset scale: Prolific contains much fewer examples than MIST (13000 vs 318). Nevertheless, adding belief priors is beneficial across datasets.

**Thematic Analysis.** We cluster news claims into latent topics and test whether LLM susceptibility accuracy differs across demographics by topic. We conduct our analysis on the MIST dataset, as it is much larger than PANDORA. We find that demographic differences in susceptibility are topic-dependent: clusters corresponding to lexically obvious or opinion-based claims exhibit minimal demographic variation, whereas more ambigu-

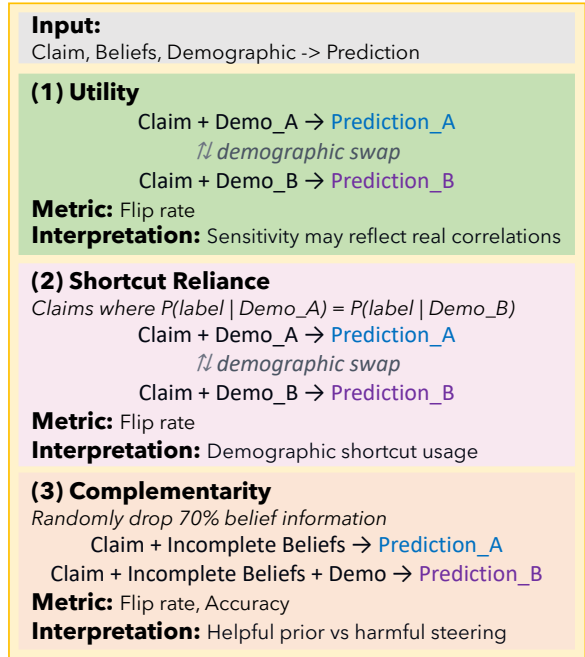


Figure 4: Demographic-based Counterfactual Eval Framework: (1) Utility of demographics, (2) Shortcut Reliance, and (3) Complementarity. Note that we perform a demographic swap within the same demographic group, i.e., swap male with female, rural with urban, etc. keeping the claim constant.

ous or diverse topics like science/health claims and government show larger differences across age (older groups scoring higher) and education groups (higher-educated groups scoring higher). Gender-based differences are consistently small across all topics. These findings suggest that demographic susceptibility is also contextual and topic-sensitive, motivating further evaluation in future studies (More details are in Appendix C.3).

## 5.2 Demographic-based Counterfactual Evaluation

Demographics can be a helpful context signal in addition to belief information and legitimately correlate with misinformation susceptibility in real data, but they can also induce spurious signals that appear accurate but could be stereotypical (Wan et al., 2025; Geirhos et al., 2020). To investigate this trade-off, we conduct counterfactual evaluations to assess whether demographic attributes provide informative priors or introduce stereotype-sensitive associations.

**Method.** We conduct three complementary analyses on the PANDORA+MIST-1 data to investigate effects of demographics in susceptibility prediction (Fig 4):

(1) Utility: Do demographics improve susceptibil-

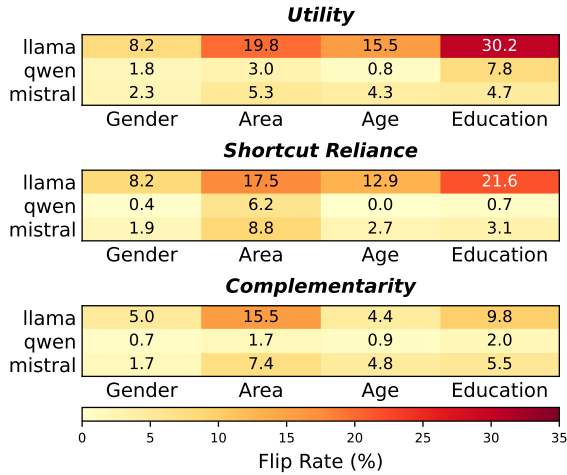


Figure 5: Flip-Rates for Counterfactual Evaluation: qwen and mistral models have lower flip-rates in comparison to llama.

ity prediction beyond claim content? We first evaluate utility on the original data distribution using the demographic-only setting. Next, we create a counterfactual demographic swap (keeping the query claim fixed), re-prompt the model, and measure the flip rate between predictions. Higher flip rates indicate stronger demographic sensitivity, which may reflect real-world demographic-label correlations rather than purely spurious associations.

(2) Shortcut reliance: Does the model rely on demographic cues even when they are non-informative by construction? Here, we design a controlled subset where the distribution of human labels in each claim is matched across demographic groups (e.g., Female and Male have equal counts of ‘true’ and ‘misinfo’ for that claim). Therefore, this makes the demographic attribute non-predictive. A high flip rate here indicates demographic short-cutting and possible stereotyping.

(3) Complementarity: When belief evidence is incomplete, do demographics act as a helpful weak prior or a harmful shortcut? We simulate missing evidence by randomly dropping a large fraction of WVS belief statements (70%). We then compare model predictions under two conditions: (1) belief only and (2) belief + demographic. High flip rates between these conditions indicate that demographic cues influence predictions under uncertainty. Additionally, we compare susceptibility prediction accuracy in both conditions to assess whether adding demographics improves performance.

**Results.** Fig 5 shows the flip-rates per model and demographic groups averaged across datasets. Overall, qwen consistently achieves lower flip rates

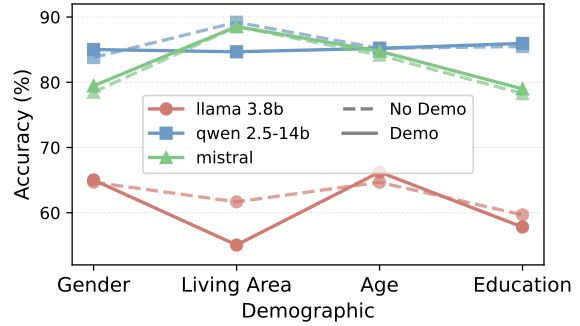


Figure 6: Complementarity Analysis Results: when partial belief information is added, adding demographics has mixed effects, showing very weak gains.

(averaging 2.16%), suggesting lower dependence on demographic cues. Mistral also exhibits lower flip rates (averaging 4.3%). Contrastingly, llama shows substantially higher flip rates (averaging 14.31%), thus higher demographic sensitivity. Across models, the largest flip rates concentrate in Education and Living Area, especially in the Shortcut Reliance setting, where demographics are designed to be non-predictive. This could also relate to stereotype-driven reliance. For complementarity analysis, Fig 6 shows that demographics are only weakly complementary: mistral shows small, consistent gains across settings. qwen improves slightly for education/gender, while drops for living area. llama also shows unstable patterns.

Overall, these findings suggest that demographic sensitivity is highly model-dependent. qwen and mistral are comparatively stable under counterfactual demographic changes, while llama shows substantially higher reliance on demographics, even when demographics are explicitly made non-predictive (Shortcut Reliance task). At the same time, demographics provide limited practical benefit when partial beliefs are present, showing they are a weak prior at best. These findings motivate moving beyond prompt-only conditioning, as we find that demographic cues in the prompt can introduce unstable shortcuts. Therefore, belief signal must be integrated in a way that is robust to conflicting sources (individual and group-based beliefs).

## 6 Post-training Adaptation for Susceptibility Simulation

Prompt-based conditioning experiments show the benefits of using beliefs for demographic misinformation susceptibility simulation. However, mixing imputed and observed beliefs in the prompt is associated with lower performance. This is intuitive, as

observed (ground-truth) beliefs are tied to individual participants and imputed beliefs are only linked to demographic groups and not individuals, creating inconsistencies. Furthermore, training directly on demographics and susceptibility labels risks leakage, where the model may learn label shortcuts as observed in counterfactual evaluation (Wan et al., 2025; Geirhos et al., 2020). To address this, we propose **BAFT** (Belief-Adapter Fine-tuning), which decouples imputed beliefs from observed data, while still learning transferable belief representations through a separate adapter.

**Method.** To decouple group-based imputed beliefs from individualistic observed beliefs, BAFT consists of a two-phase design: (1) belief modeling, and (2) susceptibility fine-tuning.

**Phase 1: Belief Modeling.** We train a belief adapter by freezing the base LLM and training a belief head, implemented as a linear projection followed by softmax, to predict demographic-conditioned response distributions for WVS belief questions. Concretely, given the encoder representation  $h_\phi(q, d)$  for question  $q$  and demographic context  $d$ , we model  $p_\theta(y | q, d) = \text{softmax}(Wh_\phi(q, d) + b)$ . This captures population-level variability in belief responses. Importantly, this stage uses only survey supervision (no misinformation labels), reducing the risk of label leakage or reward hacking.

**Phase 2: Susceptibility Fine-tuning.** In this phase, we freeze the belief adapter and train a lightweight susceptibility head that combines the base LLM’s semantic representation of the prompt with the belief representation produced by the adapter. Given a persona-conditioned input (demographics and observed beliefs) and target claim  $x$ , the frozen adapter outputs a belief embedding  $z_{\text{bel}}$ , which we concatenate with the claim representation  $h_\phi(x)$  and predict susceptibility via a binary classifier:

$$p_\psi(y | x, d) = \text{softmax}(U[h_\phi(x); z_{\text{bel}}(d)] + c), \quad (1)$$

where  $y \in \{\text{true}, \text{misinformation}\}$ . We train only the susceptibility head parameters  $\psi = \{U, c\}$  with cross-entropy loss, keeping the base model and belief adapter frozen. This is to ensure susceptibility learning relies on transferable belief structure rather than encoding demographic shortcuts.

**Data.** For Phase 1, we construct demographic belief priors from empirical WVS response distributions. For each question  $q$  and demographic group  $d$ , we estimate a categorical distribution  $P(r | q, d)$

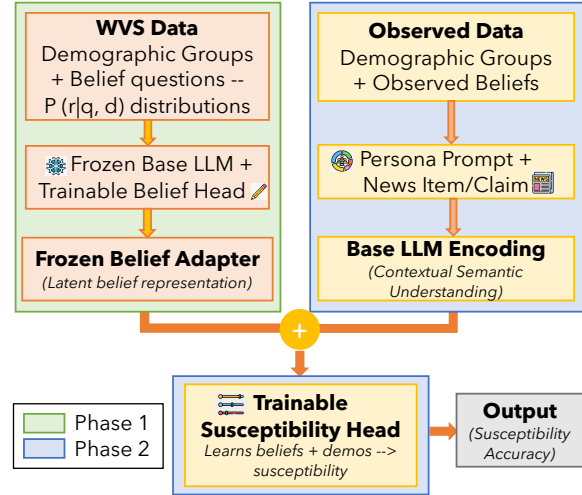


Figure 7: Belief Adapter Fine-Tuning framework. Green-shaded components correspond to Phase 1 (Belief Modeling), and blue-shaded components correspond to Phase 2 (Susceptibility Fine-Tuning).

over response options  $r \in \{1, \dots, K_q\}$ . These distributions capture meaningful demographic variability across WVS belief questions. We restrict our experiments to six demographic groups, as MIST does not include living-area annotations, and the PANDORA data is too small to support standalone data for fine-tuning. This results in 126 WVS response distributions per demographic group, yielding  $6 \times 125 = 750$  examples in total. We apply textual augmentation via back-translation using Google Translate,<sup>3</sup> producing 1250 examples. This is sufficient as the belief adapter is a lightweight linear head requiring fewer examples to reliably capture demographic belief patterns. Training is performed using an 80/20 train/val split.

For Phase 2, we use (PANDORA+MIST-1) data, keeping each demographic group separate per data point (as in prompt conditioning). The target is the prediction on a claim, giving susceptibility scores. We use a 80/20 train/val split, yielding approximately 33.1k training examples and 5.1k validation examples (after removal of overlapping examples). For evaluation, we also used the (MIST Study-2b) (to investigate cross-study generalization), consisting of 7k evaluation points.

We compare BAFT against two strong baselines: (1) standard LoRA fine-tuning using PANDORA+MIST-1 (Baseline (1-phase)), and (2) a two-phase variant of our pipeline that replaces the head-based adaptation with LoRA fine-tuning (LoRA-FT (2-phase)).

<sup>3</sup><https://translate.google.com/>

Model	PANDORA+MIST-1		MIST-2	
	Acc	F1	Acc	F1
<b>Baseline (1-phase)</b>				
llama	0.750	0.751	0.680	0.673
qwen	0.760	0.761	0.720	0.712
mistral	0.740	0.761	0.650	0.653
<b>LoRA-FT (2-phase)</b>				
llama	0.809	0.809	0.876	0.861
qwen	0.800	0.800	0.893	0.892
mistral	0.800	0.790	0.887	0.867
<b>BAFT (2-phase)</b>				
llama	0.793	0.793	0.884	0.896
qwen	0.792	0.791	0.924	0.906
mistral	0.792	0.792	0.884	0.823

Table 3: Fine-tuning performances. 2-phase belief training improves performance across both datasets.

**Evaluation.** Phase 1 reflects population-level belief fidelity, so we use distribution-level metric over Likert-scale responses, Kullback–Leibler (KL) divergence (Kullback and Leibler, 1951). This measures how well predicted belief distributions align with each demographic (lower is better for both). Phase 2 reflects task utility. So, we compute susceptibility accuracy for this phase.

## 6.1 Results

**Phase 1.** Qwen has the lowest KL-divergence (0.051), followed by mistral (0.087) and llama (0.287) for BAFT. With LoRA-ft, the trends remain similar (complete results are in Appendix D).

**Phase 2.** Table 3 shows that cross-study (MIST-2) generalization improves drastically once we move from 1-phase to 2-phase. The 1-phase baseline model performs reasonably well on PANDORA+MIST-1, however, it drops sharply on MIST-2. Comparing the 2-phase variants, BAFT remains the strongest on MIST-2 (up to 92.4%). LoRA-ft slightly improves on PANDORA+MIST-1 but gives a weaker transfer, may be due to mild overfitting when more parameters are updated. Across models, qwen consistently performs the best for 2-phase settings. Additionally, we perform shortcut reliance experiments, finding that BAFT leads to 0% flip-rates even for llama models (which earlier had the highest flip-rate: Appendix D.3).

## 7 Lessons Learned and Actionable Steps

Our findings show that belief priors are central to simulating demographic misinformation susceptibility with LLMs. We introduce Belief-Sim and BAFT, demonstrating that decoupling belief modeling from susceptibility prediction enables higher accuracy while reducing spurious demographic sensitivity. These findings offer actionable information for designing targeted demographic interventions. **Beliefs are crucial to demographic susceptibility simulation.** Across both simulation techniques, adding belief priors (especially imputed) leads to high accuracy gains. Future studies can investigate additional survey sources (e.g., Pew and other cross-national surveys) to broaden coverage and focus on diverse demographic groups and contexts. Furthermore, studying the effective decoupling of different types of beliefs is important.

**Simple head tuning is effective for generalization.** Our two-phase BAFT approach generalizes better than full fine-tuning, suggesting that lightweight, modular adaptation can be both effective and cheaper for simulation. Future work can further improve robustness by exploring alternative adapter designs and fusion mechanisms for combining belief and text representations.

**Counterfactual evaluation is important.** Our counterfactual experiments quantify each model’s reliance on demographic shortcuts. Such measures are necessary as sensitivity can be either a real or a spurious cue. Future work can build on this by developing fine-grained stress tests and evaluating which demographic cues drive model flips.

## 8 Conclusion

This paper introduced Belief-Sim, a belief-driven framework for simulating demographic misinformation susceptibility using LLMs. Across two datasets and through prompt-based conditioning and post-training adaptation, we show that belief priors are a key driver for demographic-level susceptibility patterns. In contrast, demographics alone are an unreliable signal that can induce shortcut reliances. Our results highlight the importance of decoupling imputed and observed belief sources and we provide practical evaluation tools for it, such as prediction accuracy and counterfactual sensitivity. Based on these findings, we outline directions for future work and release our open-source framework, BeliefSim.<sup>4</sup>

<sup>4</sup><https://anonymous.4open.science/r/belief-sim>

## 599 Limitations

### 600 Coverage and Intersectionality

601 Our demographic modeling is intentionally largely  
602 single-axis: we evaluate on a small set of demo-  
603 graphic attributes (8 groups) and treat them inde-  
604 pendently. Our goal was to isolate the effects of  
605 simulation across demographic groups and simplify  
606 counterfactual analysis as well as balance labels  
607 across groups. However, intersectional groups (e.g.,  
608 older x low-education, female x rural) can have in-  
609 teresting implications for simulation purposes, and  
610 future work can further investigate Belief-Sim for  
611 intersectional belief profiles. In addition, while we  
612 mostly adhere to binary demographic axes, each  
613 axis can be further analyzed in a finer-grained ap-  
614 proach. Finally, our study is based on US partici-  
615 pants only, due to the availability of current data  
616 sources. We acknowledge that these may have  
617 WEIRD implications (Mihalcea et al., 2025), and  
618 future work should focus on investigating wider,  
619 cross-cultural misinformation susceptibility.

### 620 Counterfactual Flips do not always mean 621 stereotyping

622 Our flip-rate analyses quantify sensitivity to de-  
623 mographic perturbations, but flips do not uniquely  
624 measure stereotype-driven shortcutting. Further-  
625 more, low flip rates do not guarantee fairness  
626 - models may still encode demographic effects  
627 indirectly through correlations or beliefs. Al-  
628 though our experiments help isolate spurious de-  
629 pendence, our experiments may not capture real-  
630 world causal pathways as they require much ex-  
631 tensive analysis. Future work can investigate ex-  
632 tending our counterfactual flips experiments with  
633 stronger causal/robustness evaluations to better un-  
634 derstand causal demographic effects and/or mea-  
635 sure stereotype-driven shortcutting.

### 636 Experimental Limitations

637 Access to closed-source LLMs is constrained by  
638 transparency and cost. Therefore, we limit our  
639 experiments to three open-source models (to sup-  
640 port controlled counterfactual swaps, ablations, and  
641 multiple runs). As a result, our conclusions may  
642 not fully generalize to the strongest closed-source  
643 models. Future work can extend our framework  
644 by benchmarking a broader set of closed-source  
645 models under matched prompting and budgeted  
646 sampling.

## References 647

- Alberto Acerbi and Joseph M Stubbersfield. 2023. Large language models show human-like content biases in transmission chain experiments. *Proceedings of the National Academy of Sciences*, 120(44):e2313790120. 648 649 650 651 652
- Hunt Allcott and Matthew Gentzkow. 2017. Social media and fake news in the 2016 election. *Journal of economic perspectives*, 31(2):211–236. 653 654 655
- Hunt Allcott, Matthew Gentzkow, and Chuan Yu. 2019. Trends in the diffusion of misinformation on social media. *Research & politics*, 6(2):2053168019848554. 656 657 658 659
- Nicolas M Anspach and Taylor N Carlson. 2024. Not who you think? exposure and vulnerability to misinformation. *New Media & Society*, 26(8):4847–4866. 660 661 662
- Lisa P. Argyle, Ethan C. Busby, Nancy Fulda, Joshua R. Gubler, Christopher Rytting, and David Wingate. 2023. Out of one, many: Using language models to simulate human samples. *Political Analysis*, 31(3):337–351. 663 664 665 666 667
- Dane Bertram. 2007. Likert scales. Retrieved November, 2(10):1–10. 668 669
- Angana Borah, Marwa Houalla, and Rada Mihalcea. 2025a. Mind the (belief) gap: Group identity in the world of LLMs. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 18441–18463, Vienna, Austria. Association for Computational Linguistics. 670 671 672 673 674 675
- Angana Borah and Rada Mihalcea. 2024. Towards implicit bias detection and mitigation in multi-agent LLM interactions. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 9306–9326, Miami, Florida, USA. Association for Computational Linguistics. 676 677 678 679 680 681
- Angana Borah, Rada Mihalcea, and Verónica Pérez-Rosas. 2025b. Persuasion at play: Understanding misinformation dynamics in demographic-aware human-llm interactions. *arXiv preprint arXiv:2503.02038*. 682 683 684 685 686
- William J Brady, Julian A Wills, John T Jost, Joshua A Tucker, and Jay J Van Bavel. 2017. Emotion shapes the diffusion of moralized content in social networks. *Proceedings of the National Academy of Sciences*, 114(28):7313–7318. 687 688 689 690 691
- Nadia M Brashier and Daniel L Schacter. 2020. Aging in an era of fake news. *Current directions in psychological science*, 29(3):316–323. 692 693 694
- Yun-Shiuan Chuang, Krirk Nirunwiroj, Zach Studdiford, Agam Goyal, Vincent V. Frigo, Sijia Yang, Dhavan Shah, Junjie Hu, and Timothy T. Rogers. 2024. Beyond demographics: Aligning role-playing llm-based agents using human belief networks. *Preprint, arXiv:2406.17232*. 695 696 697 698 699 700

701	David De Coninck, Thomas Frissen, Koen Matthijs, Leen d’Haenens, Grégoire Lits, Olivier Champagne-Poirier, Marie-Eve Carignan, Marc D David, Nathalie Pignard-Cheynel, Sébastien Salerno, et al. 2021. Beliefs in conspiracy theories and misinformation about covid-19: Comparative perspectives on the role of anxiety, depression and exposure to and trust in information sources. <i>Frontiers in psychology</i> , 12:646394.	755
702		756
703		757
704		758
705		
706		759
707		760
708		761
		762
709	Karen M Douglas, Joseph E Uscinski, Robbie M Sutton, Aleksandra Cichocka, Turkay Nefes, Chee Siang Ang, and Farzin Deravi. 2019. Understanding conspiracy theories. <i>Political psychology</i> , 40:3–35.	763
710		764
711		
712		
713	Francesca D’Errico, Giuseppe Corbelli, Concetta Papapicco, and Marinella Paciello. 2022. How personal values count in misleading news sharing with moral content. <i>Behavioral Sciences</i> , 12(9):302.	765
714		766
715		767
716		768
		769
		770
		771
717	Ullrich KH Ecker, Stephan Lewandowsky, John Cook, Philipp Schmid, Lisa K Fazio, Nadia Brashier, Panayiota Kendeou, Emily K Vraga, and Michelle A Amazeen. 2022. The psychological drivers of misinformation belief and its resistance to correction. <i>Nature Reviews Psychology</i> , 1(1):13–29.	772
718		773
719		774
720		775
721		776
722		
723	Jonathan St BT Evans and Keith E Stanovich. 2013. Dual-process theories of higher cognition: Advancing the debate. <i>Perspectives on psychological science</i> , 8(3):223–241.	777
724		778
725		779
726		
727	Lisa K Fazio. 2020. Repetition increases perceived truth even for known falsehoods. <i>Collabra: Psychology</i> , 6(1).	780
728		781
729		782
		783
		784
730	Daniel J Flynn, Brendan Nyhan, and Jason Reifler. 2017. The nature and origins of misperceptions: Understanding false and unsupported beliefs about politics. <i>Political psychology</i> , 38:127–150.	785
731		786
732		787
733		788
734	Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A Wichmann. 2020. Shortcut learning in deep neural networks. <i>Nature Machine Intelligence</i> , 2(11):665–673.	789
735		790
736		791
737		792
738		793
		794
		795
739	Salvatore Giorgi, Tingting Liu, Ankit Aich, Kelsey Jane Isman, Garrick Sherman, Zachary Fried, João Sedoc, Lyle Ungar, and Brenda Curtis. 2024. <a href="#">Modeling human subjectivity in LLMs using explicit and implicit human factors in personas</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 7174–7188, Miami, Florida, USA. Association for Computational Linguistics.	796
740		797
741		798
742		799
743		
744		800
745		801
746		802
		803
747	Genevieve Gorrell, Elena Kochkina, Maria Liakata, Ahmet Aker, Arkaitz Zubiaga, Kalina Bontcheva, and Leon Derczynski. 2019. <a href="#">SemEval-2019 task 7: RumourEval, determining rumour veracity and support for rumours</a> . In <i>Proceedings of the 13th International Workshop on Semantic Evaluation</i> , pages 845–854, Minneapolis, Minnesota, USA. Association for Computational Linguistics.	804
748		805
749		806
750		807
751		808
752		809
753		810
754		
	Andrew Guess, Jonathan Nagler, and Joshua Tucker. 2019a. <a href="#">Less than you think: Prevalence and predictors of fake news dissemination on facebook</a> . <i>Science Advances</i> , 5(1):eaau4586.	755
		756
		757
		758
	Andrew Guess, Jonathan Nagler, and Joshua Tucker. 2019b. Less than you think: Prevalence and predictors of fake news dissemination on facebook. <i>Science advances</i> , 5(1):eaau4586.	759
		760
		761
		762
	Dan M Kahan. 2017. Misconceptions, misinformation, and the logic of identity-protective cognition.	763
		764
	Carolyn Kaiser, Jakob Kaiser, Vladimir Manewitsch, Lea Rau, and Rene Schallner. 2025. Simulating human opinions with large language models: Opportunities and challenges for personalized survey data modeling. In <i>Adjunct Proceedings of the 33rd ACM Conference on User Modeling, Adaptation and Personalization</i> , pages 82–86.	765
		766
		767
		768
		769
		770
		771
	David Khachaturov, Roxanne Schnyder, and Robert Mullins. 2025. Governments should mandate tiered anonymity on social-media platforms to counter deep-fakes and LLM-driven mass misinformation. <i>arXiv preprint arXiv:2506.12814</i> .	772
		773
		774
		775
		776
	Solomon Kullback and Richard A Leibler. 1951. On information and sufficiency. <i>The annals of mathematical statistics</i> , 22(1):79–86.	777
		778
		779
	Stephan Lewandowsky, Ullrich KH Ecker, John Cook, Sander Van Der Linden, Jon Roozenbeek, and Naomi Oreskes. 2023. Misinformation and the epistemic integrity of democracy. <i>Current opinion in psychology</i> , 54:101711.	780
		781
		782
		783
		784
	Xialing Lin, Patric R Spence, and Kenneth A Lachlan. 2016. Social media and credibility indicators: The effect of influence cues. <i>Computers in human behavior</i> , 63:264–271.	785
		786
		787
		788
	Rakoen Maertens, Friedrich M Götz, Hudson F Golino, Jon Roozenbeek, Claudia R Schneider, Yara Kyrychenko, John R Kerr, Stefan Stieger, William P McClanahan, Karly Drabot, et al. 2024. The misinformation susceptibility test (mist): A psychometrically validated measure of news veracity discernment. <i>Behavior Research Methods</i> , 56(3):1863–1899.	789
		790
		791
		792
		793
		794
		795
	Ioana E Marinescu, Patrick N Lawlor, and Konrad P Kording. 2018. Quasi-experimental causality in neuroscience and behavioural research. <i>Nature human behaviour</i> , 2(12):891–898.	796
		797
		798
		799
	Killian L McLoughlin, William J Brady, Aden Goolsbee, Ben Kaiser, Kate Klonick, and MJ Crockett. 2024. Misinformation exploits outrage to spread online. <i>Science</i> , 386(6725):991–996.	800
		801
		802
		803
	Rada Mihalcea, Oana Ignat, Longju Bai, Angana Borah, Luis Chiruzzo, Zhijing Jin, Claude Kwizera, Joan Nwatu, Soujanya Poria, and Tamar Solorio. 2025. Why ai is weird and shouldn’t be this way: Towards ai for everyone, with everyone, by everyone. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 39, pages 28657–28670.	804
		805
		806
		807
		808
		809
		810

811	Suhong Moon, Marwa Abdulhai, Minwoo Kang, Joseph	Shibani Santurkar, Esin Durmus, Faisal Ladhak,	869
812	Suh, Widyadewi Soedarmadji, Eran Kohen Behar,	Cinoo Lee, Percy Liang, and Tatsunori Hashimoto.	870
813	and David M. Chan. 2024. <a href="#">Virtual personas for</a>	2023. <a href="#">Whose opinions do language models reflect?</a>	871
814	<a href="#">language models via an anthology of backstories.</a>	<i>Preprint</i> , arXiv:2303.17548.	872
815	<i>Preprint</i> , arXiv:2407.06576.		
816	Shalini Munusamy, Kalaivanan Syasyila, Azahah	Nikhil Sharma, Q Vera Liao, and Ziang Xiao. 2024.	873
817	Abu Hassan Shaari, Muhammad Adnan Pitchan,	Generative echo chamber? effect of llm-powered	874
818	Mohammad Rahim Kamaluddin, and Ratna Jatnika.	search systems on diverse information seeking. In	875
819	2024. Psychological factors contributing to the crea-	<i>Proceedings of the 2024 CHI Conference on Human</i>	876
820	tion and dissemination of fake news among social	<i>Factors in Computing Systems</i> , pages 1–17.	877
821	media users: a systematic review. <i>BMC psychology</i> ,		
822	12(1):673.	Mubashir Sultan, Alan N Tump, Nina Ehmann, Philipp	878
823	Inderjeet Nair and Lu Wang. 2025. Do language mod-	Lorenz-Spreen, Ralph Hertwig, Anton Gollwitzer,	879
824	els think consistently? a study of value preferences	and Ralf HJM Kurvers. 2024. Susceptibility to on-	880
825	across varying response lengths. <i>arXiv preprint</i>	line misinformation: A systematic meta-analysis	881
826	<i>arXiv:2506.02481</i> .	of demographic and psychological factors. <i>Pro-</i>	882
827	Keiichi Namikoshi, Alex Filipowicz, David A Shamma,	<i>ceedings of the National Academy of Sciences</i> ,	883
828	Rumen Iliev, Candice L Hogan, and Nikos Arechiga.	121(47):e2409329121.	884
829	2024. Using llms to model the beliefs and prefer-		
830	ences of targeted populations. <i>arXiv preprint</i>	Charles S Taber and Milton Lodge. 2006. Moti-	885
831	<i>arXiv:2403.20252</i> .	vated skepticism in the evaluation of political beliefs.	886
832	Jan Nehring, Aleksandra Gabryszak, Pascal Jür-	<i>American journal of political science</i> , 50(3):755–	887
833	gens, Aljoscha Burchardt, Stefan Schaffer, Matthias	769.	888
834	Spielkamp, and Birgit Stark. 2024. <a href="#">Large language</a>	Jasper Timm, Chetan Talele, and Jacob Haimes. 2025.	889
835	<a href="#">models are echo chambers</a> . In <i>Proceedings of the</i>	Tailored truths: Optimizing llm persuasion with per-	890
836	<i>2024 Joint International Conference on Computa-</i>	sonalization and fabricated statistics. <i>arXiv preprint</i>	891
837	<i>tional Linguistics, Language Resources and Evalua-</i>	<i>arXiv:2501.17273</i> .	892
838	<i>tion (LREC-COLING 2024)</i> , pages 10117–10123,	Jay J Van Bavel, Diego A Reinero, Victoria Spring, Eliz-	893
839	Torino, Italia. ELRA and ICCL.	abeth A Harris, and Annie Duke. 2021. Speaking my	894
840	Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Mered-	truth: Why personal experiences can bridge divides	895
841	ith Ringel Morris, Percy Liang, and Michael S Bern-	but mislead. <i>Proceedings of the National Academy</i>	896
842	stein. 2023. Generative agents: interactive simulacra	<i>of Sciences</i> , 118(8):e2100280118.	897
843	of human behavior. In <i>Proceedings of the 36th an-</i>		
844	<i>annual acm symposium on user interface software and</i>	Santosh Vijaykumar, Yan Jin, Daniel Rogerson,	898
845	<i>technology</i> , pages 1–22.	Xuerong Lu, Swati Sharma, Anna Maughan, Bianca	899
846	Gordon Pennycook and David G Rand. 2019. Lazy, not	Fadel, Mariella Silva de Oliveira Costa, Claudia	900
847	biased: Susceptibility to partisan fake news is better	Pagliari, and Daniel Morris. 2021. How shades of	901
848	explained by lack of reasoning than by motivated	truth and age affect responses to covid-19 (mis) in-	902
849	reasoning. <i>Cognition</i> , 188:39–50.	formation: randomized survey experiment among	903
850	Manuel Pratelli and Marinella Petrocchi. 2025. <a href="#">Evalu-</a>	whatsapp users in uk and brazil. <i>Humanities and</i>	904
851	<a href="#">ating the Simulation of Human Personality-Driven</a>	<i>Social Sciences Communications</i> , 8(1).	905
852	<a href="#">Susceptibility to Misinformation with LLMs</a> . IOS	Herun Wan, Jiaying Wu, Minnan Luo, Zhi Zeng, and	906
853	Press.	Zhixiong Su. 2025. <a href="#">Truth over tricks: Measuring</a>	907
854	Filipe N. Ribeiro, Koustuv Saha, Mahmoudreza Babaei,	<a href="#">and mitigating shortcut learning in misinformation</a>	908
855	Lucas Henrique, Johnnatan Messias, Fabricio Ben-	<a href="#">detection</a> . <i>Preprint</i> , arXiv:2506.02350.	909
856	venuto, Oana Goga, Krishna P. Gummadi, and	Aimei Yang, Alvin Zhou, Jieun Shin, Ke Huang-	910
857	Elissa M. Redmiles. 2019. <a href="#">On microtargeting so-</a>	Isherwood, Wenlin Liu, Chuqing Dong, Eugene Lee,	911
858	<a href="#">cially divisive ads: A case study of russia-linked</a>	and Jingyi Sun. 2024. Sharing is caring? how moral	912
859	<a href="#">ad campaigns on facebook</a> . In <i>Proceedings of the</i>	foundation frames drive the sharing of corrective mes-	913
860	<i>Conference on Fairness, Accountability, and Trans-</i>	sages and misinformation about covid-19 vaccines.	914
861	<i>parency, FAT* ’19</i> , page 140–149, New York, NY,	<i>Journal of Computational Social Science</i> , 7(3):2701–	915
862	USA. Association for Computing Machinery.	2733.	916
863	Jon Roozenbeek, Claudia R Schneider, Sarah Dryhurst,	Xuhui Zhou, Hao Zhu, Leena Mathur, Ruohong Zhang,	917
864	John Kerr, Alexandra LJ Freeman, Gabriel Recchia,	Haofei Yu, Zhengyang Qi, Louis-Philippe Morency,	918
865	Anne Marthe Van Der Bles, and Sander Van Der Lin-	Yonatan Bisk, Daniel Fried, Graham Neubig, et al.	919
866	den. 2020. Susceptibility to misinformation about	2023. <a href="#">Sotopia: Interactive evaluation for social</a>	920
867	covid-19 around the world. <i>Royal Society open sci-</i>	<a href="#">intelligence in language agents</a> . <i>arXiv preprint</i>	921
868	<i>ence</i> , 7(10):201199.	<i>arXiv:2310.11667</i> .	922

923	<b>A Taxonomy Dimensions and Imputed Beliefs</b>		973
924			974
925	Beliefs are an important tool for modeling demographic simulations in the context of misinformation. Psychological work shows that several belief dimensions, such as conspiracy beliefs, political ideology, trust in science, etc. can help us predict who is most susceptible to misinformation (Ecker et al., 2022; Munusamy et al., 2024). Rather than using only demographic information as a flat category, adding belief profiles provides a more fine-grained representation of how different groups may perceive information. Recent LLM studies also argue for modeling population-level beliefs and preferences to simulate targeted groups, using belief-like representations to approximate human responses at scale (Namikoshi et al., 2024; Kaiser et al., 2025).		975
926			976
927			977
928			
929			
930			
931			
932			
933			
934			
935			
936			
937			
938			
939			
940			
941	<b>A.1 Taxonomy formation</b>		978
942	We create a belief taxonomy of beliefs consisting of seven core dimensions that are associated with misinformation susceptibility grounded in psychological and cognitive science research:		979
943			980
944			981
945			982
946	<b>(1) Worldview and Identity Beliefs:</b> individuals interpret information through the lens of social identity and worldview (Kahan, 2017; Van Bavel et al., 2021), <b>(2) Epistemic Trust Beliefs:</b> individuals differ systematically in epistemic trust toward institutions and experts (De Coninck et al., 2021; Lewandowsky et al., 2023), <b>(3) Cognitive style:</b> individuals vary in cognition, such as reliance on analytic versus intuitive reasoning (Penneycook and Rand, 2019; Ecker et al., 2022), <b>(4) Conspiracy mentality:</b> individuals differ in a generalized predisposition to conspiracies (Douglas et al., 2019; De Coninck et al., 2021), <b>(5) Moral and Value Beliefs:</b> individuals prioritize different moral values (D’Errico et al., 2022; Yang et al., 2024), <b>(6) Emotion-Related Beliefs:</b> individuals vary in emotional responsiveness (heightened emotional arousal, such as anger or fear may amplify belief in misinformation) (Brady et al., 2017; McLoughlin et al., 2024) and <b>(7) Heuristic Beliefs:</b> individuals rely to different degrees on shortcuts such as repetition, familiarity, and social endorsement (Lin et al., 2016; Fazio, 2020).		983
947			984
948			985
949			986
950			987
951			988
952			989
953			990
954			
955			
956			
957			
958			
959			
960			
961			
962			
963			
964			
965			
966			
967			
968			
969			
970			
971			
972			
		omy enables systematic belief simulation in models by providing a framework for collecting and organizing belief data and providing interpretable axes along which belief priors can be instantiated in models.	991
		<b>A.2 Mapping to Belief Taxonomy.</b>	992
		We map imputed data to the above belief dimensions. We apply exploratory factor analysis on WVS responses to identify latent belief groupings, examining which items naturally cluster together (e.g., trust-related items to Epistemic & Trust Beliefs). Post that, we conduct a manual review to ensure proper alignment with our taxonomy and validate the final question-to-dimension mapping. Overall, we obtain 126 imputed belief questions across dimensions. Table 9 contains the 126 WVS questions mapped to the 7 belief taxonomy dimensions.	993
			994
			995
			996
			997
			998
			999
			1000
			1001
			1002
			1003
			1004
			1005
			1006
			1007
			1008
			1009
			1010
		<b>B Datasets</b>	1011
		<b>B.1 Evaluation Data</b>	1012
		For susceptibility evaluation, we use two ground-truth datasets containing human judgments of whether they believe a given claim: (1) PANDORA Dataset from Borah et al. (2025b), containing annotations from 318 participants. Each participant provided judgments on 3 distinct claims, along with demographic information including age, gender, living area, and education. These claims are collected from RumorEval (Gorrell et al., 2019), which consists of true or false rumors covering eight major news events and natural disaster events; (2) MIST dataset from Maertens et al. (2024). From MIST, we use only Study 1 (MIST-1) which includes 409 participants, each providing judgment on the same 100 claims, and demographic information including age, gender, and education. From the two datasets combined, we obtain 13.8K claims for evaluation.	1013
			1014
			1015
			1016
			1017
			1018
			1019
			1020
		<b>B.2 Belief Data</b>	1012
		We collect two complementary belief signals, consisting of individually observed belief judgments and group-level (demographic) belief distributions: <b>(1) Observed Data</b> - claims that were directly judged by participants in PANDORA and MIST-1 datasets. These responses capture individual-level belief judgments, reflecting each participant’s personal stance rather than demographic group averages. We use two claim judgments as observed be-	1013
			1014
			1015
			1016
			1017
			1018
			1019
			1020

**Dataset Examples**

**MIST:**  
**Participant ID:** 1  
**Participant Demographic Group:** d1  
**Claim:** The Government Is Knowingly Spreading Disease Through the Airwaves and Food Supply  
**Participant Choice:** fake  
**Gold Label:** fake

**PANDORA:**  
**Participant ID:** 1  
**Participant Demographic Group:** d1  
**Claim:** News: Updated CDC Covid Numbers Show Only Six Percent of Total US Deaths Actually Dies Solely From Virus.  
**Participant Choice:** true  
**Gold Label:** fake

Figure 8: Dataset Examples - MIST and PANDORA

liefs for each participant (keeping it separate from the evaluation data). This leads to 27.6 claim judgments as observed beliefs.

(2) **Imputed Data** - inferred from the World Values Survey Wave 7<sup>5</sup> distributions. Imputed data represent demographic belief priors (group-level), inferred from WVS, conditioned solely on demographic attributes. We map these imputed belief items to our belief taxonomy using exploratory factor analysis. This yields 126 imputed belief questions and corresponding demographic distributions. Table 1 shows examples, and Appendix A.2 contains mapping details and all questions.

Fig 11 shows an example from each MIST and PANDORA datasets. Both datasets are similar, except for the size. Furthermore, MIST does not contain the living area information of the participants, therefore, we do not use it for fine-tuning purposes.

**C Prompt-Based Conditioning**

**C.1 Prompt Details**

Fig 10 shows the susceptibility accuracies across the belief and demographic configurations for both

<sup>5</sup><https://www.worldvaluessurvey.org/WVSDocumentationWV7.jsp>

**Prompt**

You are a persona grounded by attributes: <demographic group>.

Past beliefs and priors for this persona (for context, do not re-evaluate them): <beliefs>.

When judging a claim, stay consistent with this persona’s prior beliefs where reasonable.

Figure 9: Prompt Conditioning Experiments

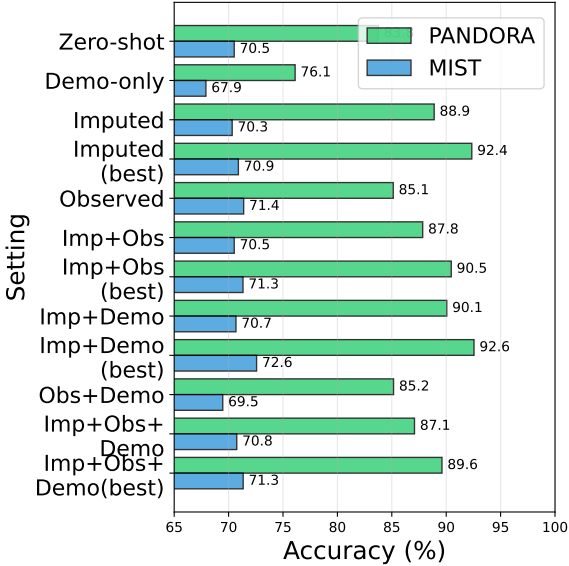


Figure 10: Susceptibility Accuracies across settings and datasets

datasets. For Imputed cases, we average across all belief dimensions, and Imputed (best) cases, we select the belief dimension that achieves the highest susceptibility accuracy. Findings show that incorporating belief information consistently improves performance over the zero-shot and demographic-only baselines. We next analyze the impact of specific belief types, demographic ablations, datasets, and model performances.

**C.2 Belief-Dimension Analysis**

Table 4 shows the belief-dimension-wise scores for each dataset and demographic dimension.

Across both datasets, single dimensions consistently outperform the combined prior, and this is especially stark on PANDORA, where all-dimensions collapses to much lower accuracy while the best

1060	single dimensions remain high.		1111
1061	On Prolific, the top-performing dimensions are		1112
1062	fairly stable: emotion-related dominates Gender,		1113
1063	while moral-and-value-beliefs is strongest and most		1114
1064	frequent for Age and Rural/Urban, and heuristic/		1115
1065	cognitive beliefs for Education. On MIST, gains		1116
1066	from beliefs are smaller but more consistent: win-		
1067	ners are spread across - top dimensions like moral-		
1068	and-value, emotion-related, heuristics, etc., still		
1069	dominate. In addition, worldview-and-identity also		
1070	does better in MIST.		
1071	<b>C.3 Thematic Analysis</b>		
1072	We perform topic clustering on the news claims		
1073	to identify latent topical grouping and investigate		
1074	whether LLM susceptibility accuracies vary across		
1075	demographics within each topic. We perform this		
1076	on Dataset 2, as it is much larger and diverse than		
1077	Dataset 1. Using non-negative matrix factorization		
1078	on claim texts, we identify five dominant topics:		
1079	(1) global power, population, and future threats;		
1080	(2) political leaders and historical rankings; (3)		
1081	science, health, and technology claims; (4) poli-		
1082	tics, government, and information-control narra-		
1083	tives; and (5) public opinion, moral values, and so-		
1084	cial belief statements. Topic-level analysis shows		
1085	that susceptibility accuracy is not uniform across		
1086	demographics. Topics related to political rank-		
1087	ings and public opinion (Topics 2 and 5) achieve		
1088	uniformly high accuracy across all demographic		
1089	groups, resulting in negligible age, gender, or edu-		
1090	cation gaps. Contrastingly, science/health claims		
1091	and government-related narratives (Topics 3 and		
1092	4) exhibit the largest demographic variation, with		
1093	older and higher-education groups generally achiev-		
1094	ing higher accuracy. Gender-based differences are		
1095	minimal across all topics. Table 5 shows results		
1096	across some topics and models showing differences		
1097	across demographics.		
1098	These findings suggest that demographic suscep-		
1099	tibility in LLMs is not just a broad, identity-		
1100	driven phenomenon but also contextual and topic-		
1101	sensitive. This also aligns with interdisciplinary		
1102	work in psychology and communication suggesting		
1103	that demographic susceptibility to misinformation		
1104	is highly context-dependent. Prior studies show		
1105	that gender alone is a weak predictor of misinform-		
1106	ation belief, with analytic thinking and reason-		
1107	ing style playing a much larger role (Marinescu		
1108	et al., 2018). In contrast, age and education de-		
1109	mographics differ more for complex or ambiguous		
1110	claims, particularly in science and policy domains,		
	consistent with work on cognitive processing, and		1111
	motivated reasoning (Kahan, 2017). From a cogni-		1112
	tive perspective, this pattern supports dual-process		1113
	theories of reasoning, whereby strong lexical re-		1114
	duce individual differences, while inference-heavy		1115
	claims amplify them (Evans and Stanovich, 2013).		1116
	<b>D Fine-tuning Analysis</b>		1117
	<b>D.1 Phase 1 results</b>		1118
	For Phase 1, Table 6 shows that Qwen has the		1119
	lowest KL divergence, followed by Mistral and		1120
	Llama for both settings - head tuning and full-fine-		1121
	tuning. With LoRA-ft, the trends remain similar		1122
	but slightly lower than head-ft. This indicates that		1123
	Qwen more faithfully captures population-level be-		1124
	lieff distributions, a critical property for reliable		1125
	downstream simulation of demographic misinform-		1126
	ation susceptibility. Other models also perform		1127
	competitively, suggesting that most models LLMs		1128
	can capture belief distributions to a reasonable ex-		1129
	tent.		1130
	<b>D.2 Phase 2 results</b>		1131
	Table 7 shows susceptibility accuracies and F1		1132
	scores across different setups and models. Across		1133
	methods, the key finding is that cross-study (MIST-		1134
	2) generalization improves drastically once we		1135
	move from 1-phase to 2-phase pipelines. The 1-		1136
	phase baseline model perform decent, however, the		1137
	performance drops sharply on MIST-2. For MIST-		1138
	2, our BAFT approach performs the best, showing		1139
	that belief adapter + lightweight susceptibility head		1140
	captures more transferable signals aligned with out-		1141
	of-domain data.		1142
	Comparing the 2-phase variants, we find that		1143
	BAFT remains the strongest on MIST-2, however,		1144
	LoRA-ft slightly improves on PANDORA+MIST-		1145
	1 but gives weaker transfer. This may be due to		1146
	mild overfitting when more parameters are updated.		1147
	Across models, Qwen consistently performs the		1148
	best for 2-phase settings. Overall, the 2-phase train-		1149
	ing approach is helpful improving robustness and		1150
	also adapting to out-of-domain data.		1151
	<b>D.3 Shortcut Reliance Experiments</b>		1152
	Table 8 shows that shortcut reliance is low, espe-		1153
	cially under head tuning. We measure flip-rate:		1154
	the fraction of evaluation examples for which the		1155
	model’s predicted label changes when a shortcut		1156
	feature is removed or perturbed, and Probability		1157
	Delta: the average absolute change in the model’s		1158

Dataset	Dimension	Gender	Age	Education	Rural/Urban
<b>Prolific-dataset</b>	worldview_and_identity	89.20625	90.19875	69.54372184	90.71375
	epistemic_and_trust_beliefs	85.515	89.2575	84.96031741	89.0475
	cognitive_style_beliefs	89.0875	88.69	84.76190471	87.6575
	conspiracy_mindset	83.25375	89.0475	85	86.985
	moral_and_value_beliefs	90.15875	90.19875	84.36507924	90.72375
	emotion_related	88.3811111	89.545	85.31746026	89.08625
	heuristic	86.78625	90.00125	86.38888875	90.555
	alldimensions	73.74625	88.96875	85.95238095	89.26
<b>MIST-dataset</b>	worldview_and_identity	71.66125	72.49375	71.65924606	
	epistemic_and_trust_beliefs	71.68	71.5875	71.47649656	
	cognitive_style_beliefs	71.27458333	71.91125	71.35596075	
	conspiracy_mindset	71.1875	71.64	71.26027066	
	moral_and_value_beliefs	71.02625	71.99	70.52219001	
	emotion_related	71.2175	71.43	70.75732404	
	heuristic	71.85875	72.09125	71.80576158	
	alldimensions	71.50958333	71.74125	71.18524694	

Table 4: Belief Dimension Analysis

1159 predicted probability for the original class after  
1160 the shortcut is removed. On PROLIFIC+MIST-  
1161 1, head-tuned models (llama, qwen and mistral)  
1162 show near-zero flip rates, indicating that removing  
1163 the shortcut signal has little effect on predicted la-  
1164 bels. Full fine-tuning also tends to have small flip  
1165 rates and probability deltas in-domain, but they are  
1166 more frequently non-zero (e.g., llama/mistral flip  
1167 more than head tuning), consistent with the idea  
1168 that updating more parameters can slightly increase  
1169 sensitivity to spurious cues.

1170 On MIST-2, the most notable shortcut sensitiv-  
1171 ity appears for qwen under full fine-tuning (flip  
1172 rate 0.1098, larger prob delta), whereas the corre-  
1173 sponding head-tuned Qwen model shows 0 flip rate.  
1174 This pattern supports our previous findings that  
1175 two-phase-head-only training is more robust: it  
1176 preserves performance while reducing reliance on  
1177 shortcut features that do not transfer.

## 1178 D.4 Training Details

### 1179 D.4.1 2-phase Head-FT details

1180 Phase 1 trains a belief adapter to predict  
1181 demographic-conditioned WVS response distribu-  
1182 tions. The model is a frozen encoder with a train-  
1183 able linear head mapping the last-token hidden state  
1184 to 10 logits, followed by softmax. Training uses  
1185 AdamW (lr  $5 \times 10^{-4}$ , batch size 16, 2 epochs)  
1186 and optimizes a scale-aware KL divergence loss,  
1187 computed only over the valid bins 1..K for each  
1188 example. For evaluation, we test generalization to  
1189 unseen belief items, reporting KL distributional fit

### Belief Modeling Prompt

You are modeling belief distributions for a  
single demographic group.

Demographic : d1

Question: q1

There are k possible response options, num-  
bered 1 through k. Predict the probability  
for each option.

Figure 11: Prompt for Phase-1 Belief Modeling

and majority-category accuracy. 1190

1191 Phase 2 trains a susceptibility prediction head  
1192 on top of a frozen belief adapter learned in Phase  
1193 1. The model combines a frozen base encoder  
1194 and frozen belief head to produce both a seman-  
1195 tic representation of the input and a demograph-  
1196 ic-conditioned belief probability vector. These two  
1197 are concatenated and fed into a lightweight, train-  
1198 able classification head. We format the inputs as  
1199 instruction-following prompts that include a de-  
1200 mographic persona and up to two belief examples.  
1201 During training, only the susceptibility head param-  
1202 eters are updated using cross-entropy loss, while  
1203 the base model and belief adapter remain frozen.  
1204 The model is optimized with AdamW (learning  
1205 rate ( $5 \times 10^{-4}$ ), batch size 8) for 2 epochs, with

Demographic	Model	Topic/Theme	Effect and Notes
Age	Llama-3-8B	New study & ideology (T1)	+2.8 ppts (60+ <30). Older: 0.984 vs younger: 0.956.
	Llama-3-8B	Global threats & population (T2)	+2.0 ppts. Older: 0.961 vs younger: 0.941.
	Llama-3-8B	Marijuana/new-study (T4)	+1.7 ppts. 60+ group: 0.947 vs <30: 0.930.
	Qwen-2.5-14B	Marijuana/new-study (T4)	+2.5 ppts. 60+ group: 0.917 vs <30: 0.892.
Education	Mistral-7B	All topics	<0.5 ppts. Age differences negligible.
	Llama-3-8B	Sleep/blue-light health (T2)	+8.4 ppts. Completed: 0.918 vs not: 0.833 (six cases).
	Qwen-2.5-14B	Politics/elections (T3)	+16.7 ppts. Completed perfect; not-completed: 0.833 (six cases).
	Llama-3-8B	Marijuana/new-study (T4)	+0.9 ppts. Difference: 0.962 vs 0.953.
	Mistral-7B	Marijuana/new-study (T4)	+0.2 ppts. Effect negligible.
Gender	All models	Other topics	0.0–0.2 ppts. No meaningful education effect.
	Llama-3-8B	Global influence/future threats (T0)	$\leq 2$ ppts ( $ F - M $ ). Nearly identical accuracies.
	Llama-3-8B	Science & health claims (T2)	$\leq 2$ ppts. No consistent female–male separation.
	Llama-3-8B	Public opinion/social beliefs (T4)	$\approx 0$ ppts. Near-ceiling for both genders.
	Mistral-7B	Global influence/future threats (T0)	$\approx 0$ ppts. Accuracies overlap almost exactly.
	Mistral-7B	Marijuana/public opinion (T4)	$\leq 1$ ppt. No directional trend.
Qwen-2.5-14B	All topics	0 ppts. Identical accuracies across all topics.	

Table 5: Thematic Analysis across Demographic

Model	Head Tuning	Full Fine-tuning
Qwen	0.0511	0.0732
Mistral	0.087	0.107
Llama	0.287	0.487

Table 6: KL divergence comparison across different models and tuning methods.

Model	PANDORA+MIST-1		MIST-2	
	Acc	F1	Acc	F1
<b>Baseline (1-phase)</b>				
llama	0.750	0.751	0.680	0.673
qwen	0.760	0.761	0.720	0.712
mistral	0.740	0.761	0.650	0.653
<b>LoRA-FT (2-phase)</b>				
llama	0.809	0.809	0.876	0.861
qwen	0.800	0.800	0.893	0.892
mistral	0.800	0.790	0.887	0.867
<b>BAFT</b>				
llama	0.793	0.793	0.884	0.896
qwen	0.792	0.791	0.924	0.906
mistral	0.792	0.792	0.884	0.923

Table 7: Fine-tuning performances

1206 reproducible initialization via fixed random seeds.  
1207 We finally evaluate at each epoch using accuracy  
1208 and macro-F1.

#### 1209 D.4.2 2-phase LoRA-FT details

1210 Here, we use Low-Rank Adaptation (LoRA) for  
1211 parameter-efficient fine-tuning in both phases. In  
1212 Phase 1, LoRA enables the model to adapt its rep-  
1213 resentations to capture demographic belief distri-  
1214 butions. In Phase 2, LoRA is again applied to the  
1215 base model during susceptibility training, allowing  
1216 the model to adjust task-relevant representations  
1217 while keeping the number of trainable parameters  
1218 small.

## 1219 E Significance Tests Across Experiments

1220 **Prompt-based Conditioning (Susceptibility Ac-**  
1221 **curacy):** Conditioning LLMs with belief informa-  
1222 tion improves demographic susceptibility simula-  
1223 tion compared to zero-shot and demo-only prompts.  
1224 We averaged our findings across 3 individual runs,  
1225 and improvements of belief settings are statistically  
1226 significant (paired t-test,  $p < 0.05$  for both PAN-  
1227 DORA and MIST-1).

1228 **Belief Dimension Analysis.** Certain belief dimen-  
1229 sions are more predictive of susceptibility than  
1230 others. We use repeated-measures ANOVA and  
1231 performance variation vary significantly by be-  
1232 lief dimensions, more for PANDORA, with a few  
1233 non-significant variation for MIST (congitive, and  
1234 all dimensions).

1235 **Post-Training Adaptation** In PANDORA+MIST-  
1236 1, accuracy differences between baseline and the  
1237 2-phase methods are modest (e.g., 0.75 vs. 0.793  
1238 for llama), and z-tests show that both BAFT and  
1239 LoRA fine-tuning significantly outperforms the  
1240 baseline ( $p=0.02-0.03$ ). On MIST-2 the gains are

Method	Model	PANDORA+MIST-1			MIST-2		
		Flip Rate	Prob Delta	Acc Drop	Flip Rate	Prob Delta	Acc Drop
Full-FT	Llama	0.0239	0.0148	0.0049	0.0000	0.0000	0.0000
	Qwen	0.0173	0.0476	0.0012	0.1098	0.0700	0.1027
	Mistral	0.0323	0.0024	0.0000	0.0000	0.0129	0.0000
Head Tuning	Llama	0.0000	0.0125	0.0000	0.0000	0.0129	0.0000
	Qwen	0.0080	0.0261	0.0012	0.0000	0.0427	0.0000
	Mistral	0.0000	0.0071	0.0000	0.0000	0.0069	0.0000

Table 8: Shortcut reliance comparison between Full Fine-tuning and Head Tuning across different models on PANDORA+MIST-1 and MIST-2 datasets.

pronounced: for llama the baseline accuracy is 0.680 while BAFT reaches 0.884 and LoRA-FT 0.876. Two-proportion z-tests show that both BAFT and LoRA-FT significantly outperform the baseline ( $p < 0.001$ ).

## F Model Choices, Implementation Details and Computational Resources

We conduct experiments using instruction-tuned LLMs including LLaMA, Qwen, and Mistral, implemented with the Hugging Face Transformers and PEFT libraries. Our framework adopts a modular design with explicit belief heads and susceptibility heads, and employs LoRA-based fine-tuning where base-model adaptation is required, while relying on lightweight head training in BAFT to decouple belief modeling from susceptibility prediction. All experiments are run on NVIDIA A40 GPUs with a maximum sequence length of 1024 tokens and fixed random seeds to ensure reproducibility and computational efficiency.

## G Reproducibility

We open-source our codes and data, which are uploaded to the submission system. This would help future work to reproduce our results and explore BeliefSim: demographic-aware misinformation susceptibility simulation in LLMs using belief-priors, and BAFT.

Table 9: World Values Survey Questions mapped to Misinformation Taxonomy Dimensions

Dimension	Q ID	Question	
<b>WorldView and Identity Beliefs</b>	Q6	Importance of religion (1=Very important to 4=Not at all important)	
	Q4	Importance of politics (1=Very important to 4=Not at all important)	
	Q171	Frequency of religious service attendance (1=More than once/week to 7=Never)	
	Q173	Religious self-identification (1=Religious person, 2=Not religious, 3=Atheist)	
	Q240	Political left-right scale (1-10)	
	Q254	National pride (1=Very proud to 4=Not at all proud, 5=Not from US)	
	Q255- Q259	Closeness to village/region/country/continent/world (1=Very close to 4=Not close at all)	
	Q19, Q23	Unwanted neighbors: different race, different religion (1=Don't want, 2=Want)	
	Q235- Q239	Views on government styles: strong leader, technocracy, army rule, democracy, religious law (1=Very good to 4=Very bad)	
	Q241- Q249	Essential characteristics of democracy: tax/subsidize, religious law interpretation, free elections, unemployment aid, army takeover, civil rights, income equality, obedience, women's rights (1-10 scale)	
	<b>Epistemic and Trust Beliefs</b>	Q57	General trust in people (1=Most can be trusted, 2=Need to be careful)
		Q58-Q63	Trust in family, neighborhood, people you know, first-time meetings, other religions, other nationalities (1=Trust completely to 4=Don't trust at all)
		Q69, Q71, Q75	Confidence in police, government, universities (1=A great deal to 4=None at all)
Q158- Q163		Views on science and technology: health/comfort, opportunities, faith vs science, moral breakdown, daily relevance, world better/worse (1-10 scale)	
<b>Cognitive Style Beliefs</b>		Q152- Q153	Priority ranking: economic growth, defense forces, participation in decisions, beautification (select most and next most important)
	Q154- Q155	Priority ranking: order, participation, fighting prices, free speech (select most and next most important)	
	Q156- Q157	Priority ranking: stable economy, humane society, ideas over money, fight crime (select most and next most important)	
	Q176	Moral uncertainty: trouble deciding which moral rules are right (1-10 scale)	
	<b>Conspiracy Mindset</b>	Q112	Perceived corruption level in country (1=No corruption to 10=Abundant corruption)
		Q113- Q117	Corruption beliefs about state authorities, business executives, local authorities, civil service providers, journalists/media (1=None to 4=All of them)
Q118		Frequency of bribery needed for services (1=Never to 4=Always)	
<b>Moral and Value Beliefs</b>		Q176	Moral uncertainty (1-10 scale)
	Q177- Q195	Justifiability of: claiming unentitled benefits, fare evasion, stealing, tax cheating, bribery, homosexuality, prostitution, abortion, divorce, premarital sex, suicide, euthanasia, wife beating, child beating, violence, terrorism, casual sex, political violence, death penalty (1=Never justifiable to 10=Always justifiable)	
	<b>Emotion Related</b>	Q44-Q45	Views on technology development and respect for authority (1=Good, 2=Don't mind, 3=Bad)
Q46		Happiness level (1=Very happy to 4=Not at all happy)	
Q47		Health status (1=Very good to 5=Very poor)	
Q48		Freedom of choice and control over life (1=No choice to 10=Great deal of choice)	
Q49-Q50		Life satisfaction and financial satisfaction (1=Completely dissatisfied to 10=Completely satisfied)	
Q52		Felt unsafe from crime in last 12 months (1=Often to 4=Never)	
Q131		General security feeling (1=Very secure to 4=Not at all secure)	
Q146- Q148		Worry about war, terrorist attack, civil war (1=Very much to 4=Not at all)	
<b>Heuristic</b>		Q94-Q104	Active membership in: church/religious, sports/recreation, arts/music/education, labor union, political party, environmental, professional, humanitarian/charity, consumer, self-help/mutual aid, women's groups (1=Inactive, 2=Active, 3=Don't belong)
	Q201- Q208	Media/communication usage: newspaper, TV news, radio news, mobile phone, email, internet, social media, talking with friends (1=Daily to 5=Never)	
	Q18-Q26	Unwanted neighbors: drug addicts, different race, AIDS, immigrants, homosexuals, different religion, heavy drinkers, unmarried couples, different language (1=Don't want, 2=Want)	
	Q29-Q31	Gender attitudes: men as better political leaders, university more important for boys, men as better business executives (1=Strongly agree to 4=Strongly disagree)	
	Q33-Q35	Job/gender attitudes: men's right to jobs when scarce, natives over immigrants, woman earning more causes problems (1=Strongly agree to 5=Strongly disagree)	