

Language Model Fine-Tuning on Scaled Survey Data for Predicting Distributions of Public Opinions

Anonymous authors

Paper under double-blind review

Abstract

1 Large language models (LLMs) present novel opportunities in public opinion
2 research by predicting survey responses in advance during the early
3 stages of survey design. Prior methods steer LLMs via descriptions of
4 subpopulations as LLMs' input prompt, yet such prompt engineering ap-
5 proaches have struggled to faithfully predict the distribution of survey re-
6 sponses from human subjects. In this work, we propose directly fine-tuning
7 LLMs to predict response distributions by leveraging unique structural
8 characteristics of survey data. To enable fine-tuning, we curate SubPOP,
9 a significantly scaled dataset of 3,362 questions and 70K subpopulation-
10 response pairs from well-established public opinion surveys. We show
11 that fine-tuning on SubPOP greatly improves the match between LLM pre-
12 dictions and human responses across various subpopulations, reducing
13 the LLM-human gap by up to 46% compared to baselines, and achieves
14 strong generalization to unseen surveys and subpopulations. Our findings
15 highlight the potential of survey-based fine-tuning to improve opinion pre-
16 diction for diverse, real-world subpopulations and therefore enable more
17 efficient survey designs.

18 1 Introduction

19 Surveys provide an essential tool for probing public opinions on societal issues, especially
20 as opinions vary over time and across subpopulations. However, surveys are also costly,
21 time-consuming, and require careful calibration to mitigate non-response and sampling
22 biases (Choi & Pak, 2004; Bethlehem, 2010). Recent work suggests that large language
23 models (LLMs) can assist public opinion studies by predicting survey responses across
24 different subpopulations, explored in both social science (Argyle et al., 2023; Bail, 2024;
25 Ashokkumar et al., 2024; Manning et al., 2024) and NLP (Santurkar et al., 2023; Chu et al.,
26 2023; Moon et al., 2024; Hämäläinen et al., 2023; Chiang & Lee, 2023). Such capabilities could
27 substantially enhance the survey development process, not as a replacement for human
28 participants but as a tool for researchers to conduct pilot testing, identify subpopulations
29 to over-sample, and test analysis pipelines prior to conducting the full survey (Rothschild
30 et al., 2024).

31 Prior work in steering language models, *i.e.* conditioning models to reflect the opinions
32 of a specific subpopulation, has primarily investigated different prompt engineering tech-
33 niques (Santurkar et al., 2023; Moon et al., 2024; Park et al., 2024a). However, prompting
34 alone has shown limited success in generating completions that accurately reflect the distri-
35 butions of survey responses collected from human subjects. Off-the-shelf LLMs (Achiam
36 et al., 2023; Dubey et al., 2024; Jiang et al., 2023) have shown to mirror the opinions of certain
37 US subpopulations such as the wealthy and educated (Santurkar et al., 2023; Gallegos et al.,
38 2024; Deshpande et al., 2023; Kim & Lee, 2023), while generating stereotypical or biased pre-
39 dictions of underrepresented groups (Cheng et al., 2023b;a; Wang et al., 2024). Furthermore,
40 these models often fail to capture the variation of human opinions within a subpopulation
41 (Kapania et al., 2024; Park et al., 2024b). While fine-tuning presents opportunities to address
42 these limitations (Chu et al., 2023; He et al., 2024), existing methods fail to train models
43 that accurately predict opinion distributions across diverse survey question topics and
44 subpopulations.

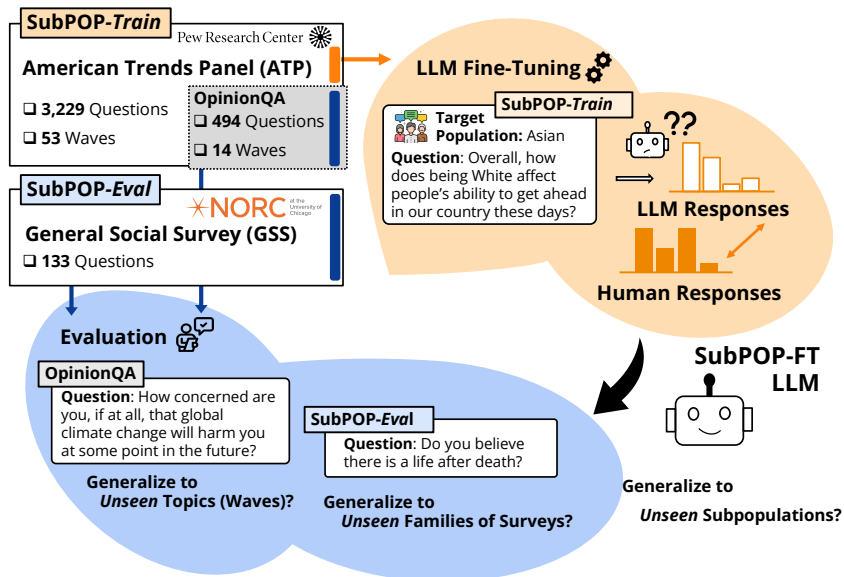


Figure 1: Illustration of our method and SubPOP. We collect survey data from two survey families—ATP from Pew Research (Pew Research Center, 2018) (forming SubPOP-Train) and GSS from NORC (Davern et al., 2024) (forming SubPOP-Eval). LLMs are fine-tuned on SubPOP-Train and evaluated on both OpinionQA (Santurkar et al., 2023) and SubPOP-Eval to assess generalization of distributional opinion prediction across unseen survey topics, survey families, and subpopulations.

45 **The present work.** Here, we propose directly fine-tuning LLMs on large-scale, high-
 46 quality survey data, consisting of questions about diverse topics and responses from each
 47 subpopulation, defined by demographic, socioeconomic, and ideological traits. By casting
 48 pairs of (subpopulation, survey question) as input prompts, we train the LLM to align its
 49 response distribution against that of human subjects in a supervised manner. We posit that
 50 survey data is particularly well-suited for fine-tuning LLMs since: (1) We can train the model
 51 with clear **subpopulation-response pairs** that explicitly link group identities and expressed
 52 opinions, which is rare in LLMs’ pre-training corpora, (2) Large-scale opinion polls are
 53 carefully designed and calibrated (e.g. using post-stratification) to estimate **representative**
 54 human responses, in contrast with LLMs’ pre-training data where certain populations are
 55 over- or underrepresented, (3) Our training objective explicitly aligns model predictions
 56 with response **distributions** from each subpopulation, enabling LLMs to capture variance
 57 within human subpopulations.

58 Training on public opinion survey data has remained under-explored due to the limited
 59 availability of structured survey datasets. To this end, we curate and release **SubPOP**
 60 (**S**ubpopulation-level **P**ublic **O**pinion **P**rediction), a dataset of 70K subpopulation-response
 61 distribution pairs ($6.5\times$ larger compared to previous datasets). We show that fine-tuning
 62 LLMs on SubPOP significantly improves the distributional match between LLM generated
 63 and human responses, and improvements are consistent across subpopulations of varying
 64 sizes. Additionally, the improvement generalizes to *unseen* subpopulations, survey waves
 65 (topics), and survey families, *i.e.* surveys administered by different institutions. Such broad
 66 generalization is particularly critical for real-world public opinions research, where practi-
 67 tioners are most in need of synthetic data for survey questions or subpopulations (or both)
 68 that they have not tested before.

69 Our contributions are summarized as follows:

- 70 • We show that training LLMs on response distributions from survey data signifi-
 71 cantly improves their ability to predict the opinions of subpopulations, reducing the
 72 Wasserstein distance between LLM and human distributions by 32-46% compared
 73 to top-performing baselines. (Section 4.1)

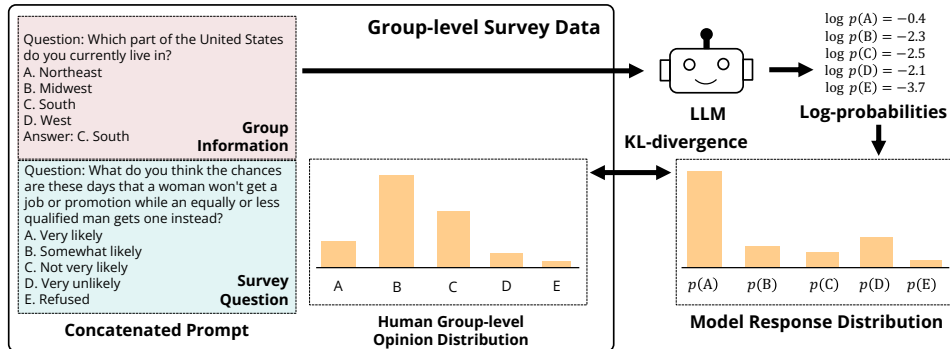


Figure 2: Proposed supervised fine-tuning setup with a survey response dataset such as SubPOP. Survey data is 3-tuple of a survey question, target subpopulation information, and the observed human opinion distribution (*i.e.* how subjects in the group responded to the given question). The training objective, $\mathcal{L}(\theta)$, is a forward KL divergence loss on language model predicted distribution of question option likelihoods; our loss guides the model predictions to match the response distribution of the specified human subpopulation.

- 74 • We show that the performance of the fine-tuned LLMs strongly generalizes to out-
75 of-distribution data, including unseen subpopulations, new survey waves, and
76 different survey families. (Section 4.1 and Section 4.2)
- 77 • We release SubPOP, a curated and pre-processed dataset of public opinion survey
78 results that is $6.5\times$ larger than existing datasets, enabling fine-tuning at scale.

79 2 Related Work

80 **Public opinion datasets.** Several research institutions conduct large-scale public opinion
81 polls and release data from those surveys. Important examples include Pew Research Cen-
82 ter’s American Trends Panel (ATP), which consists of multiple waves of cross-sectional sur-
83 veys on different topics (Pew Research Center, 2018), and General Social Survey (GSS) from
84 the NORC at the University of Chicago (Davern et al., 2024). Existing datasets have curated
85 such data for evaluating LLM-based opinion predictions, including OpinionQA (Santurkar
86 et al., 2023), a subset of ATP survey waves containing about 500 questions on contentious
87 social topics. While OpinionQA is widely used in prior work (He et al., 2024; Zhao et al.,
88 2023; Li et al., 2023; 2024), we find its total number of questions limited in scale for fine-
89 tuning LLMs and instead use this dataset for evaluation. We further collect an extended set
90 of survey data from ATP waves not included in OpinionQA, as well as from GSS to curate
91 SubPOP. Other datasets, such as GlobalOpinionQA (Durmus et al., 2023)—derived from the
92 World Values Survey (WVS) (World Values Survey, 2022) and the Pew Global Attitudes
93 Survey (Pew Research Center, 2024)—and the PRISM dataset (Kirk et al., 2024) investigates
94 how language models align with opinions from populations across the globe and cultures.
95 For additional discussions of related work, refer to Appendix C.

96 3 Methods

97 3.1 Fine-tuning LLMs on Human Response Distributions

98 Our goal is to fine-tune an LLM to predict the distribution of responses for a multiple-choice
99 question, conditioned on descriptions of a human subpopulation we want to simulate,
100 typically a specific demographic, socioeconomic, or ideological subgroup. Consider the
101 example in Figure 2: the question asks, “What do you think the chances are these days
102 that a woman won’t get a job or promotion while an equally or less qualified man gets one
103 instead?” The available responses are: *A. Very likely*, *B. Somewhat likely*, *C. Not very likely*, *D.*

104 *Very unlikely, and E. Refused.* In this case, the LLM will output a probability for each of the
 105 tokens corresponding to the choices A through E, thereby generating a complete response
 106 distribution that we aim to align with the true distribution observed in survey data.

107 Formally, let $q \in Q$ be a question, $g \in G$ be a subpopulation, and \mathcal{A}_q be the set of possible
 108 choices for question q . An LLM with parameters θ produces a conditional probability
 109 distribution $p_\theta(\mathcal{A}_q | q, g)$. We fine-tune this model so that its predicted distribution for
 110 each (q, g) mirrors the human response distribution $p_H(\mathcal{A}_q | q, g)$ collected from real
 111 survey data. To accomplish this, we use LoRA (Hu et al., 2021) and the forward Kullback–
 112 Leibler (KL) divergence as our loss. Concretely, if $p_H(\mathcal{A}_q | q, g)$ represents the group-level
 113 empirical distribution of human opinions and $p_\theta(\mathcal{A}_q | q, g)$ represents the model’s predicted
 114 distribution, our training objective is:

$$\mathcal{L}(\theta) = \mathbb{E}_{q,g} \left[D_{\text{KL}}(p_H(\mathcal{A}_q | q, g) \parallel p_\theta(\mathcal{A}_q | q, g)) \right],$$

115 where D_{KL} denotes the KL divergence. In the example shown in Figure 2, the model is
 116 trained to reduce the KL divergence between the target (survey-based) distribution over
 117 $\{A, B, C, D, E\}$ and its predicted distribution for the subpopulation living in the Southern
 118 United States. We choose forward KL (i.e., $\text{KL}(p_H \parallel p_\theta)$) since it is sensitive to cases where
 119 p_H assigns high probability but p_θ does not, naturally encouraging the model to *cover* the
 120 real distribution. This property aligns with standard maximum-likelihood training, where
 121 the model is penalized for underestimating any response that is frequent in the data. In
 122 other words, if many participants in group g choose option “A” for question q , then the
 123 model probability on “A” should be correspondingly high.

124 Instead of explicitly modeling the group response distribution as $p_H(\mathcal{A}_q | q, g)$, one could
 125 do two alternatives. (1) One-hot encoding: this approach (Li et al., 2024) approximates
 126 the distribution by a one-hot vector, assigning a value of one to the most probable option
 127 and zero elsewhere. (2) Data augmentation by response frequency: this approach (Zhao
 128 et al., 2023) expands the dataset by replicating question-choice pairs in proportion to their
 129 observed frequency. We adopt the explicit distribution modeling in our main experiments
 130 because it directly encodes the distributional information without requiring discrete sam-
 131 pling or replicating data points. A detailed comparison of these approaches is provided in
 132 Section F.2.

133 3.2 SubPOP: a Comprehensive Survey Dataset to Fine-tune and Evaluate LLMs

134 OpinionQA (Santurkar et al., 2023) is a widely used dataset for fine-tuning and evaluating
 135 large language models (LLMs) on opinion prediction, containing roughly 500 questions
 136 drawn from 14 American Trends Panel (ATP) waves (Pew Research Center, 2018). Although
 137 valuable, it faces two important limitations: (1) Limited thematic diversity—for instance,
 138 wave 26 focuses on the topic of firearms. (2) Reliance on a single survey family (ATP), which
 139 risks overfitting to a particular style of questions and limits out-of-distribution evaluation
 140 on other sources (e.g., GSS).

141 To address these limitations, we introduce a new dataset, SubPOP, that broadens both the
 142 thematic and institutional scope of opinion prediction data. For training, SubPOP comprises
 143 3,229 multiple-choice questions drawn from ATP waves 61–132, excluding waves included
 144 in OpinionQA. In Table 4, we list the topics of the ATP waves in SubPOP vs. OpinionQA,
 145 both showing the increased thematic diversity of SubPOP (with over 20 new topics) and the
 146 remaining unseen topics in OpinionQA that allow us to test whether LLMs fine-tuned on
 147 SubPOP can generalize to unseen topics.

148 For evaluation, SubPOP also includes 133 multiple-choice questions from the General Social
 149 Survey (GSS) (Davern et al., 2024), serving as an out-of-distribution benchmark. This
 150 expanded collection not only broadens the range of topics beyond OpinionQA’s initial 500
 151 questions, but also enables evaluation on surveys created and administered by different
 152 institutions (Pew Research Center vs. NORC-Chicago). Dataset curation and refinement
 153 pipeline is available in Appendix D.

Table 1: Evaluation on OpinionQA and the SubPOP evaluation set (SubPOP-Eval) for 22 subpopulations following Santurkar et al. (2023). We compute the WD by averaging over all questions and subpopulations. Lower and upper bounds of performance give guidance on how each method performs. For Modular Pluralism, we provide an error rate of one-hot prediction (†) (Section 3.3) which was used in the original paper.

Method	OpinionQA				SubPOP-Eval			
	Llama-2-7B	Llama-2-13B	Mistral-7B	Llama-3-70B	Llama-2-7B	Llama-2-13B	Mistral-7B	Llama-3-70B
Upper bound (Unif.)	0.178				0.208			
Lower bound (Human)	0.031				0.033			
Zero-shot prompt (QA)	0.173	0.170	0.153	0.138	0.206	0.196	0.187	0.160
Zero-shot prompt (BIO)	0.193	0.183	0.162	0.143	0.221	0.212	0.202	0.175
Zero-shot prompt (PORTRAY)	0.195	0.207	0.158	0.209	0.212	0.242	0.194	0.247
Few-shot prompt	0.186	0.175	0.174	0.166	0.217	0.194	0.175	0.182
Modular Pluralism	0.285 (†55.6%)				0.279 (†55.2%)			
Ours (SubPOP-FT)	0.106	0.102	0.096	0.094	0.121	0.113	0.115	0.096

154 3.3 Evaluation Metric

155 We use Wasserstein distance (WD) to quantify how closely the model’s predicted opinion
 156 distribution matches human survey data (Santurkar et al., 2023; Moon et al., 2024; Meister
 157 et al., 2024; Zhao et al., 2023). Formally, for a group g representing some subpopulation and
 158 a question q WD is defined as $WD_{\theta}(q, g) = WD(p_H(\mathcal{A}_q|q, g), p_{\theta}(\mathcal{A}_q|q, g))$ (see formula in
 159 Appendix E). Since WD is computed over ordinal values, we map the categorical answer
 160 options to numbers, such as mapping “Very likely” to 1, “Likely” to 2, and so on.

161 Some prior work utilizes one-hot accuracy (Feng et al., 2024; Li et al., 2023) as an evaluation
 162 metric. However, one-hot accuracy only verifies whether the top-predicted choice matches
 163 the top human response, thereby discarding distributional information. In contrast, WD
 164 accounts for partial overlaps among the categories and reflects the ‘cost’ of shifting proba-
 165 bility mass, providing a more nuanced assessment of distribution discrepancy. Consider the
 166 example question provided in Figure 2, where the human response distribution indicates
 167 that option B (“Somewhat likely”) is the most probable. Now consider two cases in which
 168 the model incorrectly predicts the top choice. In the first case, the model assigns a high
 169 probability to option A (“Very likely”), while in the second case, it assigns a high probability
 170 to option D (“Very unlikely”). Although one-hot accuracy would treat both predictions
 171 equally as errors, WD differentiates between them by accounting for the ordinal relationship
 172 among the options, penalizing the second prediction more heavily for its larger deviation
 173 from the true distribution.

174 4 Experiments

175 **Lower and upper bounds.** We use a uniform distribution over all available choices to es-
 176 tablish an upper bound of the WD between a predicted and the target response distribution.
 177 To compute a lower bound, we sample a group of human respondents from the original
 178 human respondents to calculate the WD between the two, and perform bootstrapping to
 179 obtain a robust estimate. This lower bound captures the intrinsic variance arising from the
 180 respondent sampling process in opinion surveys.

181 **Baselines.** We compare our approach with two baseline methods: prompting and Modular
 182 Pluralism (Feng et al., 2024). For prompting, we consider both zero-shot and few-shot
 183 methods. In zero-shot prompting, we steer the LLM using demographic prompt formats.
 184 Specifically, we employ three different formats following Santurkar et al. (2023): QA, BIO, and
 185 PORTRAY. For instance, to condition the LLM to a person living in the South of the US, the QA
 186 format uses a question-answer format as illustrated in Figure 2; the BIO format conditions
 187 the model with a first-person narrative such as “I currently reside in the South.”; and the
 188 PORTRAY format uses a third-person narrative like “Answer the following question as if you
 189 currently reside in the South.”.

190 Few-shot prompting augments the prompt with a few examples of question-response
 191 distribution pairs alongside the demographic label (Hwang et al., 2023). In particular,

192 we select the top five few-shot examples from the SubPOP training set based on cosine
 193 similarity computed by the embedding model. In our experiments, we represent the
 194 response distribution in JSON format and require the model to output its prediction in the
 195 same JSON format, following the approach in Meister et al. (2024).

196 Modular pluralism (Feng et al., 2024) fine-tunes multiple LLMs on distinct datasets to
 197 capture the viewpoints of different communities (Feng et al., 2023). For a given question,
 198 each fine-tuned LLM generates an opinion that reflects the perspective of the community
 199 it represents, and a separate black-box LLM aggregates these outputs to produce the final
 200 distributional response. Detailed implementation of the lower/upper bounds and the
 201 baselines is provided in Appendix G.

202 4.1 Generalization to Unseen Topics and Survey Families

203 In this section, we assess the ability of our fine-tuned LLMs to generalize to unseen
 204 data—both in terms of new topics and entirely different survey families. To evaluate
 205 these aspects, we use OpinionQA to measure generalization to unseen topics, and SubPOP-
 206 Eval to test generalization to a different survey family. We fine-tune four LLMs (Llama-2-7B,
 207 Llama-2-13B, Mistral-7B, and Llama-3-70B) on SubPOP-Train. We opt for pretrained LLMs
 208 rather than instruction-following models, as previous work has shown that pretrained
 209 models perform better on this task (Moon et al., 2024). A detailed comparison between these
 210 model types is provided in Appendix F.3.

211 **Summary of results.** Table 1 reports the average WD metrics computed over all demo-
 212 graphic groups and survey questions, comparing our fine-tuned models against various
 213 baseline approaches. Our experiments show that fine-tuning on SubPOP-Train significantly
 214 outperforms all other methods, yielding a 32–46% reduction in WD on OpinionQA and a
 215 39–42% reduction on SubPOP-Eval compared to the best baselines. Notably, SubPOP-Train is
 216 based on ATP data, while SubPOP-Eval is derived from GSS surveys—two distinct survey
 217 families that can differ in respondent pools, calibration techniques, and other methodologi-
 218 cal factors, leading to non-trivial distribution shifts despite both being representative of the
 219 US population. Furthermore, our fine-grained analyses at the wave level (see Appendix H)
 220 confirm that these trends persist even at more detailed levels of evaluation.

221 **Comparison to zero- and few-shot prompting.** We first compare the performance of
 222 prompting methods with our approach. Zero-shot prompting results in only modest WD
 223 improvements over the upper bound, with the largest gain observed for Llama-3-70B and
 224 negligible improvements for Llama-2-7B. Even when using few-shot prompting—where five
 225 example question-response distribution pairs are provided—the performance gains remain
 226 minimal. This may be partly due to an under-optimized prompt format (e.g. requiring JSON
 227 output) and the inherent sensitivity of language models to prompt formatting (Sclar et al.,
 228 2023; Anagnostidis & Bulian, 2024). These findings underscore the need for methods, such
 229 as fine-tuning, that enable relatively reliable predictions of opinion distributions.

230 **Comparison to Modular Pluralism.** Modular Pluralism improves one-hot accuracy, re-
 231 ducing prediction error from 72.7% (zero-shot prompting) to 55.6% on OpinionQA, but
 232 underperforms in matching the full distribution of option choices, measured as WD. This
 233 discrepancy in performance highlights the limitations of methods that train LLMs to identify
 234 only the most probable response rather than modeling the entire distribution of responses.
 235 Opinions are inherently distributed: even within a particular subpopulation such as a single
 236 demographic subgroup, distribution of opinions cannot be captured as a single most likely
 237 response. Moreover, instruction-tuned models that serve as a black-box LLM tend to assign
 238 high probabilities on only specific tokens (Lin et al., 2022; Kadavath et al., 2022; Achiam
 239 et al., 2023), further pushing the generated distribution away from the human distribution.

240 4.2 Generalization across Target Subpopulations

241 Here we report two key observations: (1) prediction performance improves consistently
 242 across subpopulations in the fine-tuning data, and (2) the LLMs fine-tuned on SubPOP-Train
 243 generalize well to subpopulations that were not included during fine-tuning.

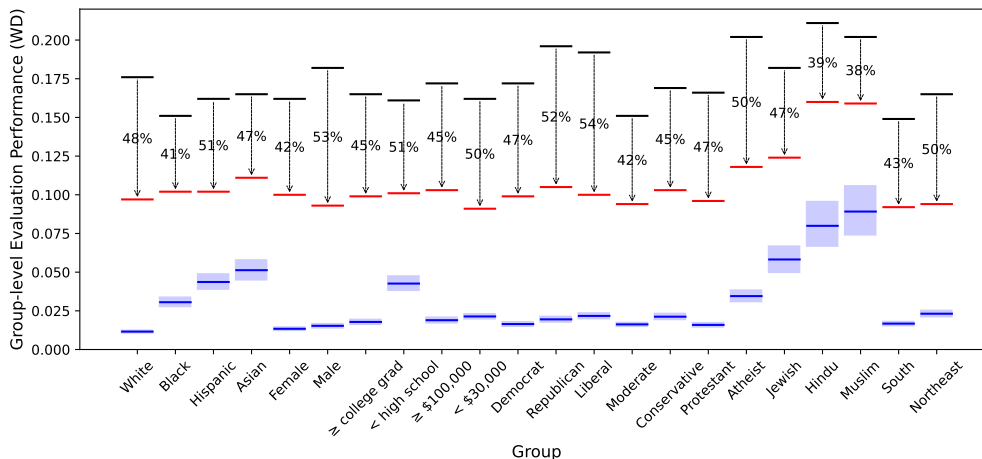


Figure 3: Per-group evaluation performance of our model Llama-2-7B-SubPOP-FT (red lines) on OpinionQA. For comparison, the results from zero-shot QA prompting (black lines) and the lower bound (blue lines) are presented. We observe that the relative improvement, measuring how much of the gap between zero-shot prompting and the lower bound has been closed, remains consistent across subpopulations. Shaded blue regions represent the 95% confidence interval of the lower-bound estimation for each group. Per-group results for other models (Table 9) and the results on SubPOP evaluation set (Table 10) are available in Appendix H.

244 **Consistent performance improvements over subpopulations.** Figure 3 shows the per-
 245 group WD on the OpinionQA evaluation for Llama-2-7B, comparing our fine-tuning ap-
 246 proach with zero-shot prompting and the empirical WD lower bound. To evaluate the
 247 consistency of performance gains, we calculate the *relative improvement* for each subpop-
 248 ulation as how much of the gap between zero-shot prompting and the empirical lower
 249 bound is reduced after fine-tuning. This measure allows us to account for varying lower
 250 bounds across subpopulations: since some groups have fewer respondents, there is greater
 251 uncertainty in their reported distribution in the survey data and greater variance between
 252 the original sample and bootstrap samples.

253 All 22 subpopulations demonstrate a large relative improvement after fine-tuning, ranging
 254 from 38%–54%. The average relative improvement is 46.7% with a standard deviation of
 255 4.4%. This consistency confirms that our fine-tuning approach delivers balanced perfor-
 256 mance gains without disproportionately favoring any particular demographic subgroup.
 257 We hypothesize that the consistent gains over groups largely stem from our dataset de-
 258 sign, which allocates an equal number of training samples to each group. By ensuring
 259 uniformly distributed data points across subpopulations, the model captures sufficient
 260 subgroup-specific signals, ultimately leading to consistent performance improvements.

261 **Generalization to unseen subpopulations.** We further investigate how models fine-tuned
 262 with our approach and SubPOP might show generalization to subpopulations that were
 263 not represented in the training data, a circumstance that may arise in real-world survey
 264 development. For the evaluation, we benchmark our methods against a zero-shot prompting
 265 baseline. Specifically, we evaluate our model, which is fine-tuned on 22 subpopulations
 266 provided in SubPOP-Train, on a set of subpopulations in OpinionQA that were not included
 267 in fine-tuning. This experiment not only checks generalization to unseen subpopulations,
 268 but also involves unseen survey questions, providing a robust assessment of the model
 269 capability for generalization to out-of-distribution data.

270 As shown in Table 2, our model achieves a strong reduction in WD even for unseen sub-
 271 populations, indicating that the model can be steered by demographic prompts beyond
 272 the seen subpopulations in training. Interestingly, although SubPOP-Train does not contain
 273 any data with opinion distributions of particular age groups (*e.g.* subjects of age 18-29 or
 274 those of age 65+), the average relative improvement is 44.7%, which is compatible with the

Table 2: Per-group evaluation performance of Llama-2-7B-SubPOP-FT (Ours) on OpinionQA. We report the lower bound, WD for zero-shot prompting, WD for Llama-2-7B-SubPOP-FT, and the relative improvement. Rows highlighted in blue represent subpopulations included during fine-tuning, while uncolored rows correspond to subpopulations that were unseen during fine-tuning.

Group	Lower Bound	Zero Shot	Ours	Relative Improvement (%)
Age: 18-29	0.023	0.185	0.096	54.9
Age: 30-49	0.014	0.151	0.093	42.3
Age: 50-64	0.014	0.154	0.101	37.9
Age: 65+	0.013	0.195	0.115	44.0
Less than high school	0.043	0.161	0.101	50.8
High school graduate	0.017	0.144	0.092	40.9
Some college, no degree	0.018	0.144	0.093	40.5
Associate’s degree	0.026	0.159	0.098	44.9
College grad	0.018	0.165	0.099	44.9
Postgraduate	0.015	0.174	0.106	42.8
Very conservative	0.026	0.208	0.107	55.5
Conservative	0.021	0.169	0.103	44.6
Moderate	0.016	0.151	0.094	42.2
Liberal	0.022	0.192	0.100	54.1
Very liberal	0.025	0.202	0.111	51.4
Democrat	0.016	0.172	0.099	47.1
Republican	0.019	0.196	0.105	52.0
Independent	0.016	0.155	0.093	44.5
Something Else	0.026	0.162	0.092	51.0

275 average relative improvement for seen subpopulations. We provide results for other unseen
 276 groups in Table 7 of Appendix F.4 (average relative improvement of 43.1% with a standard
 277 deviation of 6.7%).

278 **Steerability towards subpopulations.** Given the large improvements in WD across sub-
 279 populations after fine-tuning, we want to test whether the LLM is truly adapting its pre-
 280 dictions based on the subpopulation specified in its prompt (*i.e.* the LLM is being steered)
 281 or if the improvements can be explained by the LLMs’ predictions getting closer to human
 282 responses in general, without any subpopulation-specific adaptation. If the LLM is being
 283 steered, we should expect that the LLM’s predictions for a target subpopulation g_t are closer
 284 to the human distribution for g_t when g_t is the subpopulation specified in the prompt,
 285 compared to when another group g_s is specified in the prompt. We should also expect the
 286 gap in WD to be larger if the distance between the true human distributions for g_t and g_s
 287 are larger, such as differences between the youngest and oldest age groups compared to
 288 adjacent groups.

289 Formally, we define the *intergroup disagreement* between a target group g_t and a source group
 290 g_s as $WD(p_H(\mathcal{A}_q | q, g_t), p_H(\mathcal{A}_q | q, g_s))$ averaged over evaluation questions. In human
 291 responses (left of Figure 4), the disagreement shows the pattern of locality: increases as the
 292 disparity in education levels between two groups grows. We extend this notion to compare
 293 the human distribution from the target group g_t with the LLM-predicted distribution when
 294 the *source* group g_s is specified in the prompt, $WD(p_H(\mathcal{A}_q | q, g_t), p_\theta(\mathcal{A}_q | q, g_s))$. If
 295 the model truly incorporates subpopulation information from the prompt, its intergroup
 296 disagreement pattern should mirror that of the human data.

297 Zero-shot prompting with the base model (right of Figure 4) does not exhibit the locality
 298 pattern seen in the human data, indicating that it cannot be steered by subpopulation labels.
 299 In contrast, the fine-tuned model (middle of Figure 4) reproduces a pattern resembling
 300 the human-human case, even though it was trained on only two education groups (“less
 301 than high school” and “college graduate/some postgrad”) and the other four groups were
 302 unseen. This result demonstrates that our fine-tuned model not only learns to condition

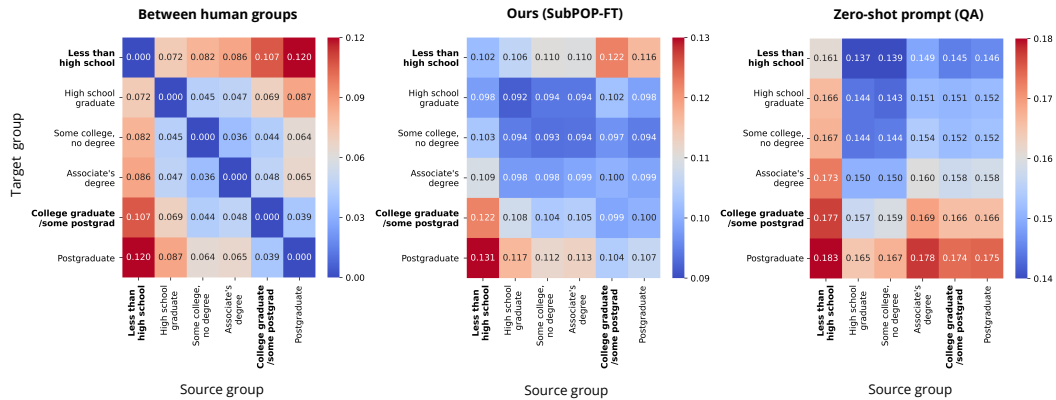


Figure 4: *Intergroup disagreement* pattern between groups of different education levels calculated with OpinionQA and Llama-2-7B as a base model. A target human group is compared to (left) a source human group, (middle) our fine-tuned model conditioned on a source group, (right) a base model conditioned on a source group. Bold-faced groups are included in the fine-tuning data SubPOP-Train, while the others aren't. In the human response (left), we observe a decreasing disagreement level as the education level becomes similar. This disagreement pattern exists in our fine-tuned model but not in the zero-shot prompting with a base model, indicating that our model can be steered to given subpopulation label even for unseen demographics while the base model cannot.

303 on subpopulation information but also generalizes to subpopulations unseen during fine-
 304 tuning. We provide the intergroup disagreement for other traits in Appendix F.4.

305 5 Conclusion

306 In this work, we demonstrated that fine-tuning large language models on structured public
 307 opinion survey data markedly improves their ability to predict human response distri-
 308 butions. We curate SubPOP—a dataset 6.5× larger than previous collections to fine-tune
 309 and evaluate LLMs on survey response distribution prediction. By fine-tuning on SubPOP,
 310 we showed that LLMs can capture the group-specific variability in public opinions, also
 311 generalizing to unseen subpopulations, survey waves and question topics, and different
 312 survey families. Fine-tuning achieves consistent improvements across subpopulations of
 313 varying sizes, and our experiments demonstrate that fine-tuned LLMs are indeed *adapting*
 314 their responses to the subpopulation specified in the prompt, even for subpopulations
 315 unseen during fine-tuning. Finally, our experiments also reveal that as the fine-tuning
 316 dataset grows, model performance continues to scale favorably, underscoring the value of
 317 our larger dataset.

318 Generalization is a critical capability for LLMs if they are to be used to assist public opinion
 319 research, as researchers are most in need of opinion predictions for questions or subpopula-
 320 tions whom they have not surveyed before. Our work, by greatly improving LLMs' ability
 321 to accurately predict opinions with fine-tuning and demonstrating strong generalization to
 322 out-of-distribution data, moves us closer towards the goal of leveraging LLMs for opinion
 323 prediction. However, many open questions remain: why is the model able to generalize
 324 well to unseen subpopulations and questions, and when might it fail to do so? How do
 325 we ensure that LLMs capture opinions along other dimensions not explored in this work,
 326 such as intersections of demographic identities or temporal change? How should LLMs be
 327 integrated into survey designs, to serve as tools that can complement surveys with human
 328 participants? Answering these questions will require interdisciplinary collaborations with
 329 domain experts and critical assessments of LLMs' and traditional survey methods' strengths
 330 and weaknesses.

331 **References**

- 332 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni
333 Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al.
334 Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- 335 Sotiris Anagnostidis and Jannis Bulian. How susceptible are llms to influence in prompts?
336 *arXiv preprint arXiv:2408.11865*, 2024.
- 337 Lisa P Argyle, Ethan C Busby, Nancy Fulda, Joshua R Gubler, Christopher Rytting, and
338 David Wingate. Out of one, many: Using language models to simulate human samples.
339 *Political Analysis*, 31(3):337–351, 2023.
- 340 Ashwini Ashokkumar, Luke Hewitt, Isaias Ghezae, and Robb Willer. Predicting results of
341 social science experiments using large language models. *accessed September, 19:2024*, 2024.
- 342 Christopher A Bail. Can generative ai improve social science? *Proceedings of the National
343 Academy of Sciences*, 121(21):e2314021121, 2024.
- 344 Jelke Bethlehem. Selection bias in web surveys. *International statistical review*, 78(2):161–188,
345 2010.
- 346 James Bisbee, Joshua D Clinton, Cassy Dorff, Brenton Kenkel, and Jennifer M Larson.
347 Synthetic replacements for human survey data? the perils of large language models.
348 *Political Analysis*, 32(4):401–416, 2024.
- 349 Yong Cao, Haijiang Liu, Arnav Arora, Isabelle Augenstein, Paul Röttger, and Daniel Hersh-
350 covich. Specializing large language models to simulate survey response distributions for
351 global populations. *arXiv preprint arXiv:2502.07068*, 2025.
- 352 Souradip Chakraborty, Jiahao Qiu, Hui Yuan, Alec Koppel, Dinesh Manocha, Furong Huang,
353 Amrit Bedi, and Mengdi Wang. Maxmin-rlhf: Alignment with diverse human preferences.
354 In *Forty-first International Conference on Machine Learning*, 2024.
- 355 Daiwei Chen, Yi Chen, Aniket Rege, and Ramya Korlakai Vinayak. Pal: Pluralistic alignment
356 framework for learning from heterogeneous preferences. *arXiv preprint arXiv:2406.08469*,
357 2024.
- 358 Myra Cheng, Esin Durmus, and Dan Jurafsky. Marked personas: Using natural language
359 prompts to measure stereotypes in language models. *arXiv preprint arXiv:2305.18189*,
360 2023a.
- 361 Myra Cheng, Tiziano Piccardi, and Diyi Yang. Compost: Characterizing and evaluating
362 caricature in llm simulations. *arXiv preprint arXiv:2310.11501*, 2023b.
- 363 Cheng-Han Chiang and Hung-yi Lee. Can large language models be an alternative to
364 human evaluations? *arXiv preprint arXiv:2305.01937*, 2023.
- 365 Bernard CK Choi and Anita WP Pak. A catalog of biases in questionnaires. *Preventing
366 chronic disease*, 2(1):A13, 2004.
- 367 Eric Chu, Jacob Andreas, Stephen Ansolabehere, and Deb Roy. Language models trained on
368 media diets can predict public opinion. *arXiv preprint arXiv:2303.16779*, 2023.
- 369 Michael Davern, Rene Bautista, Jeremy Freese, Pamela Herd, and Stephen L. Morgan. Gen-
370 eral social survey 1972-2024, 2024. URL <https://gssdataexplorer.norc.org>. Principal
371 Investigator: Michael Davern; Co-Principal Investigators: Rene Bautista, Jeremy Freese,
372 Pamela Herd, and Stephen L. Morgan. Sponsored by National Science Foundation. NORC
373 ed. Chicago: NORC, 2024.
- 374 Ameet Deshpande, Vishvak Murahari, Tanmay Rajpurohit, Ashwin Kalyan, and Karthik
375 Narasimhan. Toxicity in chatgpt: Analyzing persona-assigned language models. *arXiv
376 preprint arXiv:2304.05335*, 2023.

- 377 Ricardo Dominguez-Olmedo, Moritz Hardt, and Celestine Mender-Dünner. Questioning
378 the survey responses of large language models. *arXiv preprint arXiv:2306.07951*, 2023.
- 379 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle,
380 Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3
381 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- 382 Esin Durmus, Karina Nyugen, Thomas I Liao, Nicholas Schiefer, Amanda Askill, Anton
383 Bakhtin, Carol Chen, Zac Hatfield-Dodds, Danny Hernandez, Nicholas Joseph, et al.
384 Towards measuring the representation of subjective global opinions in language models.
385 *arXiv preprint arXiv:2306.16388*, 2023.
- 386 Shangbin Feng, Chan Young Park, Yuhan Liu, and Yulia Tsvetkov. From pretraining data to
387 language models to downstream tasks: Tracking the trails of political biases leading to
388 unfair nlp models. *arXiv preprint arXiv:2305.08283*, 2023.
- 389 Shangbin Feng, Taylor Sorensen, Yuhan Liu, Jillian Fisher, Chan Young Park, Yejin Choi, and
390 Yulia Tsvetkov. Modular pluralism: Pluralistic alignment via multi-LLM collaboration.
391 In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024*
392 *Conference on Empirical Methods in Natural Language Processing*, pp. 4151–4171, Miami,
393 Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/
394 v1/2024.emnlp-main.240. URL <https://aclanthology.org/2024.emnlp-main.240/>.
- 395 Isabel O Gallegos, Ryan A Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck
396 Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K Ahmed. Bias and fairness in large
397 language models: A survey. *Computational Linguistics*, pp. 1–79, 2024.
- 398 Perttu Hämäläinen, Mikke Tavast, and Anton Kunnari. Evaluating large language models
399 in generating synthetic hci research data: a case study. In *Proceedings of the 2023 CHI*
400 *Conference on Human Factors in Computing Systems*, pp. 1–19, 2023.
- 401 Zihao He, Minh Duc Chu, Rebecca Dorn, Siyi Guo, and Kristina Lerman. Community-
402 cross-instruct: Unsupervised instruction generation for aligning large language models to
403 online communities. *arXiv preprint arXiv:2406.12074*, 2024.
- 404 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang,
405 Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv*
406 *preprint arXiv:2106.09685*, 2021.
- 407 EunJeong Hwang, Bodhisattwa Prasad Majumder, and Niket Tandon. Aligning language
408 models to user opinions. *arXiv preprint arXiv:2305.14929*, 2023.
- 409 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh
410 Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile
411 Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- 412 Liwei Jiang, Taylor Sorensen, Sydney Levine, and Yejin Choi. Can language models reason
413 about individualistic human values and preferences? *arXiv preprint arXiv:2410.03868*,
414 2024.
- 415 Saurav Kadavath, Tom Conerly, Amanda Askill, Tom Henighan, Dawn Drain, Ethan Perez,
416 Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, et al. Language
417 models (mostly) know what they know. *arXiv preprint arXiv:2207.05221*, 2022.
- 418 Shivani Kapania, William Agnew, Motahhare Eslami, Hoda Heidari, and Sarah Fox. ‘simu-
419 lacrum of stories’: Examining large language models as qualitative research participants.
420 *arXiv preprint arXiv:2409.19430*, 2024.
- 421 Jaehyung Kim and Yiming Yang. Few-shot personalization of llms with mis-aligned re-
422 sponses. *arXiv preprint arXiv:2406.18678*, 2024.
- 423 Junsol Kim and Byungkyu Lee. Ai-augmented surveys: Leveraging large language models
424 and surveys for opinion prediction. *arXiv preprint arXiv:2305.09620*, 2023.

- 425 Hannah Rose Kirk, Alexander Whitefield, Paul Röttger, Andrew Bean, Katerina Margatina,
426 Juan Ciro, Rafael Mosquera, Max Bartolo, Adina Williams, He He, et al. The prism
427 alignment project: What participatory, representative and individualised human feedback
428 reveals about the subjective and multicultural alignment of large language models. *arXiv*
429 *preprint arXiv:2404.16019*, 2024.
- 430 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu,
431 Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large
432 language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th*
433 *Symposium on Operating Systems Principles*, 2023.
- 434 Thom Lake, Eunsol Choi, and Greg Durrett. From distributional to overton pluralism:
435 Investigating large language model alignment. *arXiv preprint arXiv:2406.17692*, 2024.
- 436 Cheng Li, Mengzhou Chen, Jindong Wang, Sunayana Sitaram, and Xing Xie. Cul-
437 turellm: Incorporating cultural differences into large language models. *arXiv preprint*
438 *arXiv:2402.10946*, 2024.
- 439 Junyi Li, Ninareh Mehrabi, Charith Peris, Palash Goyal, Kai-Wei Chang, Aram Galstyan,
440 Richard Zemel, and Rahul Gupta. On the steerability of large language models toward
441 data-driven personas. *arXiv preprint arXiv:2311.04978*, 2023.
- 442 Stephanie Lin, Jacob Hilton, and Owain Evans. Teaching models to express their uncertainty
443 in words. *arXiv preprint arXiv:2205.14334*, 2022.
- 444 I Loshchilov. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- 445 Benjamin S Manning, Kehang Zhu, and John J Horton. Automated social science: Language
446 models as scientist and subjects. Technical report, National Bureau of Economic Research,
447 2024.
- 448 Nicole Meister, Carlos Guestrin, and Tatsunori Hashimoto. Benchmarking distributional
449 alignment of large language models. *arXiv preprint arXiv:2411.05403*, 2024.
- 450 Igor Melnyk, Youssef Mroueh, Brian Belgodere, Mattia Rigotti, Apoorva Nitsure, Mikhail
451 Yurochkin, Kristjan Greenewald, Jiri Navratil, and Jerret Ross. Distributional preference
452 alignment of llms via optimal transport. *arXiv preprint arXiv:2406.05882*, 2024.
- 453 Andrew Mercer, Arnold Lau, and Courtney Kennedy. For weighting online opt-in samples,
454 what matters most? 2018.
- 455 Suhong Moon, Marwa Abdulhai, Minwoo Kang, Joseph Suh, Widyadewi Soedarmadji,
456 Eran Kohen Behar, and David Chan. Virtual personas for language models via an
457 anthology of backstories. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.),
458 *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*,
459 pp. 19864–19897, Miami, Florida, USA, November 2024. Association for Computational
460 Linguistics. doi: 10.18653/v1/2024.emnlp-main.1110. URL [https://aclanthology.org/](https://aclanthology.org/2024.emnlp-main.1110/)
461 [2024.emnlp-main.1110/](https://aclanthology.org/2024.emnlp-main.1110/).
- 462 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin,
463 Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language
464 models to follow instructions with human feedback. *Advances in neural information*
465 *processing systems*, 35:27730–27744, 2022.
- 466 Joon Sung Park, Carolyn Q Zou, Aaron Shaw, Benjamin Mako Hill, Carrie Cai, Meredith
467 Ringel Morris, Robb Willer, Percy Liang, and Michael S Bernstein. Generative agent
468 simulations of 1,000 people. *arXiv preprint arXiv:2411.10109*, 2024a.
- 469 Peter S Park, Philipp Schoenegger, and Chongyang Zhu. Diminished diversity-of-thought
470 in a standard large language model. *Behavior Research Methods*, pp. 1–17, 2024b.
- 471 Pew Research Center. America trends panel waves. Retrieved February 06, 2025, from
472 <https://www.pewsocialtrends.org/dataset>, 2018.

- 473 Pew Research Center. Pew research center, 2024. URL [https://www.pewresearch.org/
474 example-report](https://www.pewresearch.org/example-report). Accessed: February 10, 2025.
- 475 Sriyash Poddar, Yanming Wan, Hamish Ivison, Abhishek Gupta, and Natasha Jaques.
476 Personalizing reinforcement learning from human feedback with variational preference
477 learning. *arXiv preprint arXiv:2408.10075*, 2024.
- 478 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and
479 Chelsea Finn. Direct preference optimization: Your language model is secretly a reward
480 model. *Advances in Neural Information Processing Systems*, 36, 2024.
- 481 Arun Rajkumar and Shivani Agarwal. A statistical convergence perspective of algorithms
482 for rank aggregation from pairwise data. In *International conference on machine learning*, pp.
483 118–126. PMLR, 2014.
- 484 David M. Rothschild, James Brand, Hope Schroeder, and Jenny Wang. Op-
485 portunities and risks of llms in survey research, 2024. Available on SSRN:
486 <http://dx.doi.org/10.2139/ssrn.5001645>.
- 487 Shibani Santurkar, Esin Durmus, Faisal Ladhak, Cino Lee, Percy Liang, and Tatsunori
488 Hashimoto. Whose opinions do language models reflect? In *International Conference on
489 Machine Learning*, pp. 29971–30004. PMLR, 2023.
- 490 Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. Quantifying language models’
491 sensitivity to spurious features in prompt design or: How i learned to start worrying
492 about prompt formatting. *arXiv preprint arXiv:2310.11324*, 2023.
- 493 Gabriel Simmons. Moral mimicry: Large language models produce moral rationalizations
494 tailored to political identity. *arXiv preprint arXiv:2209.12106*, 2022.
- 495 Anand Siththaranjan, Cassidy Laidlaw, and Dylan Hadfield-Menell. Distributional prefer-
496 ence learning: Understanding and accounting for hidden context in rlhf. *arXiv preprint
497 arXiv:2312.08358*, 2023.
- 498 Taylor Sorensen, Jared Moore, Jillian Fisher, Mitchell Gordon, Niloofar Mireshghallah,
499 Christopher Michael Rytting, Andre Ye, Liwei Jiang, Ximing Lu, Nouha Dziri, et al. A
500 roadmap to pluralistic alignment. *arXiv preprint arXiv:2402.05070*, 2024.
- 501 Chenkai Sun, Ke Yang, Revanth Gangi Reddy, Yi R Fung, Hou Pong Chan, Kevin Small,
502 ChengXiang Zhai, and Heng Ji. Persona-db: Efficient large language model person-
503 alization for response prediction with collaborative data refinement. *arXiv preprint
504 arXiv:2402.11060*, 2024.
- 505 Lindia Tjuatja, Valerie Chen, Tongshuang Wu, Ameet Talwalkwar, and Graham Neubig. Do
506 llms exhibit human-like response biases? a case study in survey design. *Transactions of the
507 Association for Computational Linguistics*, 12:1011–1026, 2024.
- 508 Roger Tourangeau. The psychology of survey response. *University of Cambridge*, 2000.
- 509 Angelina Wang, Jamie Morgenstern, and John P Dickerson. Large language models cannot
510 replace human participants because they cannot portray identity groups. *arXiv preprint
511 arXiv:2402.01908*, 2024.
- 512 World Values Survey. World Values Survey, 2022. URL [https://www.worldvaluessurvey.
513 org/wvs.jsp](https://www.worldvaluessurvey.org/wvs.jsp). [Online; accessed 02/15/2025].
- 514 Binwei Yao, Zefan Cai, Yun-Shiuan Chuang, Shanglin Yang, Ming Jiang, Diyi Yang, and
515 Junjie Hu. No preference left behind: Group distributional preference optimization. *arXiv
516 preprint arXiv:2412.20299*, 2024.
- 517 Siyan Zhao, John Dang, and Aditya Grover. Group preference optimization: Few-shot
518 alignment of large language models. *arXiv preprint arXiv:2310.11523*, 2023.

519 A Limitations

520 In this work, we explore the capability of language models to complement traditional survey
521 design by predicting survey responses in advance. However, we acknowledge the following
522 inherent limitations of this approach.

523 **Role in Survey Research.** While language models can provide a coarse approximation
524 of human opinions, they cannot fully replace human involvement in the survey process.
525 Human opinions evolve dynamically in response to social events, and while pretrained
526 language models can incorporate such knowledge through retrieval-augmented generation,
527 they remain limited in adapting to a rapidly changing world. Moreover, fine-tuning a
528 language model on distributions of human opinions may inadvertently replicate and amplify
529 existing biases of humans, leading to undesirable outcomes. It is important to note that a
530 model fine-tuned on human opinions does not necessarily align with human values and
531 behaviors, nor does it serve as a perfect proxy for human decision-making. The scope
532 of our work is restricted to language models prompted with a group-level information
533 generating response distributions to survey questions, rather than simulating individual
534 human respondents in a personalized manner.

535 **Data Dependence.** Survey response data, even after post-stratification calibration, remain
536 subject to empirical variance, particularly for relatively small groups that comprise about
537 one percent of the U.S. population. Also, while traditional surveys have implemented
538 various strategies to mitigate response bias stemming from the linguistic and multiple-
539 choice nature of survey questions (Tourangeau, 2000), the extent to which these biases
540 affect language models—and how best to address them—remains an open question (Tjuatja
541 et al., 2024; Bisbee et al., 2024). Future research could focus on developing reliable opinion
542 datasets for underrepresented groups and examining how prompt engineering elements
543 can be optimized to reduce bias in language model-generated responses.

544 **Limited Contextual Information.** Our fine-tuning approach, which structures prompts
545 in a QA format, demonstrates strong matching with human opinion distributions. How-
546 ever, we have not explored fine-tuning with richer contextual information. Prior research
547 suggests that incorporating additional contextual details can improve the fidelity of model-
548 generated opinions to actual human responses. We anticipate that more sophisticated
549 steering techniques could further enhance the opinion prediction performance beyond the
550 results presented in this study. Investigating such methods remains an open and promising
551 direction for future work.

552 B Potential Risks

553 Employing language models for opinion prediction has both influential possibilities and risk
554 of misuse. We acknowledge that the risk of misuse cannot be overlooked, and we clearly
555 state that indiscriminately minimizing the discrepancy of opinion response distribution as a
556 fine-tuning target can cause severe harms. In particular, the model might develop a bias
557 toward specific demographics during the course of fine-tuning, an artifact of minimizing
558 response distribution when other safeguard measures are not employed. We emphasize that
559 an oversight and holistic evaluation of methods and pipelines are required before deploying
560 such models for any of the actual applications and interactions with human.

561 C Expanded Discussion of Related Work

562 **Predicting human opinions with LLMs.** Prior work has explored various prompt en-
563 gineering approaches for steering LLM responses: earlier work use rule-based prompts
564 that incorporate demographic profiles of individuals or populations, or few-shot examples
565 of survey question-response (Hwang et al., 2023; Simmons, 2022; Santurkar et al., 2023;
566 Dominguez-Olmedo et al., 2023). Recent work explore prompting LLMs with open-ended
567 text, including interview transcripts (Park et al., 2024a), personal narratives (Moon et al.,

568 2024), or LLM-refined prompts (Kim & Yang, 2024; Sun et al., 2024). Our fine-tuning ap-
569 proach is complementary to prompt engineering methods: while prompt engineering seeks
570 to optimize what information is provided to the LLM (while the model is frozen), fine-tuning
571 seeks to optimize how the model utilizes the provided information (while the prompt is
572 frozen). In this work, we demonstrate that our fine-tuned models exhibit significant im-
573 provements in matching the response distributions of humans without requiring elaborate
574 prompt engineering methods.

575 Other work (Chu et al., 2023; He et al., 2024; Feng et al., 2024) fine-tune language models
576 on text corpora from specific communities (e.g., Reddit) to infer the most popular response
577 or response distribution for a given survey question. While this approach benefits from
578 large-scale and continuously updated text corpora, it struggles with disproportionate repre-
579 sentation online and lacks comprehensive coverage of diverse subpopulations. A few works
580 have explored directly fine-tuning on public opinion survey data, but in different problem
581 settings from ours. Li et al. (2023) apply collaborative filtering to individual-level responses
582 to learn embeddings for individuals, and Zhao et al. (2023) develop a meta-learning frame-
583 work to predict the opinions of new groups given a small number of in-context examples
584 for that group. In contrast, our approach does not require individual-level responses and
585 can generalize to unseen groups and survey questions without *any* responses.

586 A recent work (Li et al., 2024) and a work concurrent to ours (Cao et al., 2025) also explored
587 fine-tuning LLMs on the World Values Survey (WVS) to align the LLM’s opinion response
588 with a culture or entire country populations. In comparison, our work focuses on US surveys,
589 testing whether LLMs can align with finer-grained subpopulations within one country
590 and whether LLMs fine-tuned on one US-representative survey can generalize to another.
591 However, we note that our proposed method for fine-tuning language models applies to
592 any survey dataset with distributional information about subpopulation responses.

593 **Pluralistic alignment of LLMs.** Recent literature on pluralistic and distributional align-
594 ment target a similar yet different problem in fine-tuning LLMs (Chakraborty et al., 2024;
595 Melnyk et al., 2024; Poddar et al., 2024; Siththaranjan et al., 2023; Yao et al., 2024; Sorensen
596 et al., 2024; Lake et al., 2024; Chen et al., 2024; Jiang et al., 2024). While this line of work
597 shares a similar goal as ours in training models to reflect on opinions (and preferences) of
598 diverse subpopulations, most work differ from ours in that they operate in the context of
599 training against *pair-wise* preference orderings between alternative language model com-
600 pletions, extending the Bradley-Terry-Luce model (Rajkumar & Agarwal, 2014; Ouyang
601 et al., 2022; Rafailov et al., 2024) or investigating alternative models to account for diverging
602 preference orderings across populations. In contrast, our work trains the model to directly
603 predict the opinion distributions of human subpopulations, where accurately matching
604 distributions across a large variety of subpopulations is of paramount interest. Our work
605 additionally focuses on the particular context of estimating human opinions about societal
606 issues—the objective of public opinion research—which enables relatively straightforward
607 supervised training on openly available, structured survey data as presented by SubPOP.

608 D Dataset Details

609 D.1 American Trends Panel Datasets

610 Pew Research holds regular American Trends Panel (ATP) survey (called waves) (Pew
611 Research Center, 2018) covering various topics (e.g. veterans, political priorities, gender
612 and leadership) and releases result at an individual level. For each anonymized individual,
613 the following information is released: unique identification number, demographic details,
614 survey responses, and weight. Weights (Mercer et al., 2018) are the output of post-survey
615 calibration process that helps adjusting survey results for response bias (e.g., non-response
616 bias, sampling bias) correction and population representativeness. As of January 2025,
617 survey data until wave 132 has been released. About 20 surveys are conducted in each year.

618 D.2 OpinionQA

619 OpinionQA is a subset of ATP curated in (Santurkar et al., 2023). This dataset consists of
 620 contentious 500 questions sampled from 14 ATP waves which have high intergroup dis-
 621 agreement (i.e. large Wasserstein distances among subpopulations’ responses to a question).
 622 It also comes with hand-crafted ordinality information which provides structure to option
 623 lists. For example, options ‘Major reason’, ‘Minor reason’, and ‘Not a reason’, are assigned
 624 an ordinality mapping to 1, 2, and 3, respectively. This ordinality allows a calculation of
 625 1-dimensional Wasserstein distance.

626 Subpopulations we employ are listed in Table 3. This set of groups is adopted for several
 627 small-scale analysis (Santurkar et al., 2023; Zhao et al., 2023; Kim & Yang, 2024). We note
 628 that our approach is not limited to a specific number of groups and data is available for
 629 small or fine-grained demographic subpopulations.

Table 3: A list of 22 subpopulations used throughout our fine-tuning and analysis. We provide the number of respondents in each subpopulation in American Trends Panel Wave 82 for reference.

Trait	Groups	Population % in Wave 82
Region	Northeast	17.2
	South	37.8
Education	College grad+	24.2
	Less than high school	5.2
Gender	Male	44.3
	Female	54.6
Race / ethnicity	Black	9.6
	White	66.1
	Asian	4.8
	Hispanic	15.2
Income	\$100,000 or more	21.8
	Less than \$30,000	21.3
Political Party	Democrat	35.1
	Republican	29.1
Political Ideology	Liberal	20.0
	Conservative	22.6
	Moderate	38.3
Religion	Protestant	40.8
	Jewish	2.0
	Hindu	0.9
	Atheist	0.6
	Muslim	0.7

630 D.3 SubPOP-Train

631 We gather additional data from the American Trends Panel, specifically collecting 53 waves
 632 from Wave 61 to 132. There are 62 waves from Wave 61 - 132, however, some waves have
 633 missing demographic or ideology information (for example, wave 63 does not contain
 634 political ideology information) or the data is not available hence removed during the
 635 curation process. To refine the dataset, we exclude questions that meet the following criteria:
 636 those with more than 10 response options, redacted response data, or dependencies on
 637 prior questions (e.g., assessing political strength). For the remaining questions, we use
 638 GPT-4o to refine their wording, ensuring they are well-suited for prompting the language
 639 models while making minimal modifications. In Figure 5 we provide a few-shot prompt for
 640 question refinement.

641 In Figure 6, we visualize the embeddings of the question texts (projected to 2-dimensions
 642 using t-SNE) from OpinionQA compared to SubPOP-Train and SubPOP-Eval. The visualization
 643 shows how much larger our dataset is than OpinionQA (6.5 \times), along with the expanded
 644 coverage of our dataset into semantic areas untouched by OpinionQA. The embeddings also
 645 reveal the distribution shift from ATP questions to GSS questions: while the ATP and GSS
 646 questions overlap in embedding space, the GSS question appear as small clusters, not evenly

Instruction: Refine the question with a minimal change to make the question sensible. Do not modify options, and do not modify a question if it makes sense. Always start your answer with "Refined question:".

Question: A cross // Do you have any of the following for spiritual purposes?
A. Yes, I have this for spiritual purposes
B. No, I do not have this for spiritual purposes

Refined question: Do you have a cross for spiritual purposes?

Question: As you may know, same-sex marriage is now legal in the U.S. Do you think this is [a good thing or a bad thing] for our society?
A. Very good thing
B. Somewhat good thing
C. Somewhat bad thing
D. Very bad thing

Refined question: As you may know, same-sex marriage is now legal in the U.S. Do you think this is a good thing or a bad thing for our society?,

Question: On a different subject...How much, if at all, do white people benefit from advantages in society that black people do not have
A. A great deal
B. A fair amount
C. Not too much
D. Not at all

Refined question: How much, if at all, do white people benefit from advantages in society that black people do not have?,

Question: Thinking about the past couple of weeks, would you say the news for Donald Trump has been...
A. Very good
B. Mostly good
C. Neither good nor bad
D. Mostly bad
E. Very bad

Refined question: Thinking about the past couple of weeks, would you say the news for Donald Trump has been...

Question: **(Question to refine)**
(Options)

Refined question:

Figure 5: Few-shot prompt for refining the question to suit a language model prompting. An instruction is designed to make a minimal change to the original question, and in-context examples are provided.

647 distributed over the ATP questions. In Table 4, we list each ATP wave in SubPOP-Train and
648 OpinionQA, along with its number of questions and wave topic(s), as defined by ATP.¹ The
649 table indicates which topics are new in SubPOP-Train compared to OpinionQA, indicating the
650 expanded coverage of our dataset, along with which topics remain unseen in OpinionQA,
651 which we can use to test LLMs fine-tuned on SubPOP-Train for generalization.

652 D.4 SubPOP-Eval

653 To further evaluate the out-of-distribution generalization ability of our fine-tuned models,
654 we subsample 133 questions from the GSS 2022 dataset (Davern et al., 2024). We apply the
655 same selection criteria as outlined in Appendix D.3, excluding questions that are redacted,

¹ATP wave topics and time periods are defined at <https://www.pewresearch.org/american-trends-panel-datasets/>.

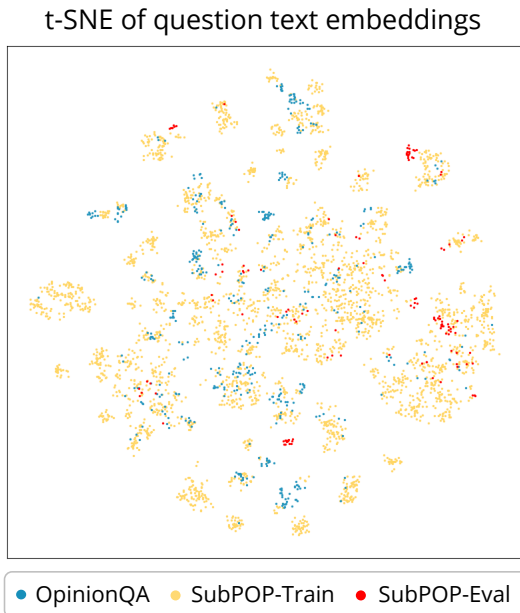


Figure 6: Embeddings of questions from OpinionQA, SubPOP-Train, and SubPOP-Eval.

656 conditioned on prior questions, inferable directly from the group information, derived from
 657 a set of questions, or those with more than 10 response options.

658 D.5 Inspection of Identical Questions

659 Distribution of cosine similarities between two text embeddings (an output of the embedding
 660 model OpenAI-text-embedding-3-large given a question text), one from a question in SubPOP-
 661 Train and another from OpinionQA is shown in Figure 7. We observed a fraction of pairs
 662 having high cosine similarity, and manually inspected question pairs with high relevance.
 663 We find that by setting a threshold cosine similarity of 0.87 we can detect all semantically
 664 identical pairs. We took a conservative threshold of cosine similarity; this value was to
 665 maximize the recall at a cost of precision to ensure detection of overlapping questions.

666 E Experiment Details

667 We conduct our experiments using Nvidia A100 GPUs with 80GB VRAM. Hyperparameter
 668 tuning is performed over learning rates $\{5e-5, 1e-4, 2e-4\}$ and batch sizes $\{64, 128, 256\}$.
 669 After evaluating possible combinations, we select a (learning rate, batch size) = $(2e-4, 256)$
 670 for Llama-2-7B, (learning rate, batch size) = $(2e-4, 256)$ for Mistral-7B-v0.1, and (learning
 671 rate, batch size) = $(1e-4, 256)$ for Llama-2-13B when utilizing the full training dataset. For
 672 Llama-3-70B, we have not done hyperparameter search but heuristically used (learning rate,
 673 batch size) = $(2e-5, 256)$. For sub-sampled training data (Figure 9), we use the following
 674 configurations:

- 675 • (lr, bs) = $(2e-4, 256)$ for 75% of the training data
- 676 • (lr, bs) = $(1e-4, 128)$ for 50% of the training data
- 677 • (lr, bs) = $(1e-4, 128)$ for 25% of the training data

678 All training is performed using LoRA (Hu et al., 2021), with LoRA parameters initialized
 679 from a normal distribution with $\sigma = 0.02$. We set the LoRA rank to 8, alpha to 32, and apply
 680 a dropout rate of 0.05. LoRA weights are applied to the query and value matrices. The
 681 AdamW (Loshchilov, 2017) optimizer is used with a weight decay of 0.

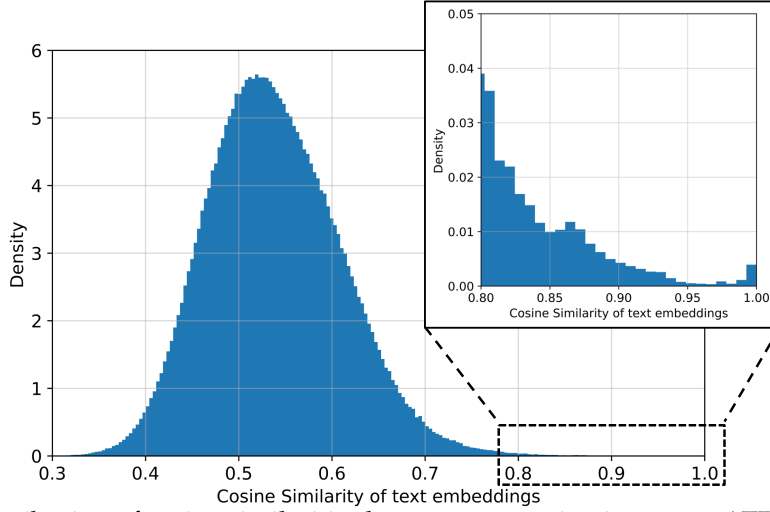


Figure 7: Distribution of cosine similarities between a question in SubPOP-ATP and OpinionQA, having a long tail towards a high cosine similarity. We inspect the question pairs in the range of 0.8 to 1.0 (distribution shown in the magnified view) and use a similarity of 0.87 as a safe threshold to identify a semantically identical question pair.

682 We use offline batched inference of vLLM (version 0.7.2) (Kwon et al., 2023) for inference
 683 and measuring response probability distribution of all methods.

684 **Choice of the Training Objective.** In this section, we explore both forward KL-divergence
 685 and Wasserstein Distance (WD) as training objectives. The forward KL-divergence is defined
 686 as

$$D_{\text{KL}}(p_H \| p_\theta) = \sum_{a \in \mathcal{A}_q} p_H(a) \log \frac{p_H(a)}{p_\theta(a)},$$

687 where $p_H(a) \equiv p_H(a | q, g)$ and $p_\theta(a) \equiv p_\theta(a | q, g)$. Similarly, WD is given by

$$\mathcal{WD}(p_H, p_\theta) = \min_{\gamma \in \Pi(p_H, p_\theta)} \sum_{a, a' \in \mathcal{A}_q} \gamma(a, a') d(a, a'),$$

688 with $\Pi(p_H, p_\theta)$ denoting the set of all couplings between p_H and p_θ , and $d(a, a')$ the L1
 689 distance between choices. Since survey responses are inherently one-dimensional and
 690 ordinal, we can simplify the computation of WD using cumulative distribution functions
 691 (CDFs). In the 1-D case, WD is computed as

$$\begin{aligned} \mathcal{WD}(p_H, p_\theta) &= \int_{-\infty}^{+\infty} |F_{p_H}(x) - F_{p_\theta}(x)| dx, \\ &= \sum_{i=1}^n |F_{p_H}(i) - F_{p_\theta}(i)| \end{aligned}$$

692 where F_{p_H} and F_{p_θ} are the CDFs corresponding to p_H and p_θ , respectively. We use this
 693 discrete formulation as the WD loss in our training.

694 While training with WD resulted in a higher KL-divergence on the validation set, the valida-
 695 tion WD converged to similar levels regardless of the objective (see Figure 8). We attribute
 696 this to KL-divergence penalizing low-probability assignments without significantly altering
 697 the overall distribution geometry. Given the KL divergence’s broader applicability—without
 698 requiring ordinal information—we primarily used KL-divergence in our experiments.

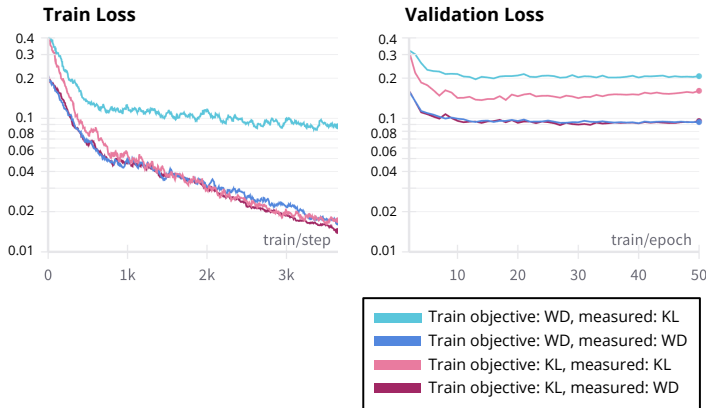


Figure 8: Train loss curve (left) and validation loss curve (right) for Llama-2-7B fine-tuned on 90% of OpinionQA, with the remaining 10% used for validation. Light and dark blue lines represent KL-divergence (KL) and Wasserstein distance (WD) when used KL as a training objective, while light and dark red lines represent KL and WD when used WD as a training objective. The two training objectives yield similar results in terms of WD, the primary measure of opinion distribution matching in our work.

699 F Additional Experiments

700 F.1 Effect of Scaling the Dataset

701 In this section, we examine performance scales with training dataset size. We randomly sam-
 702 ple subsets containing 25%, 50%, 75%, and 87.5% of the full SubPOP training set and evaluate
 703 three models—Llama-2-7B, Llama-2-13B, and Mistral-7B—on OpinionQA. As shown in Fig-
 704 ure 9, we observe diminishing marginal returns, as is typical with fine-tuning; for example,
 705 after training on a random 25%, the models reach 72%-78% of the total improvement they
 706 achieve after fine-tuning on all of SubPOP-train. However, the performance does not entirely
 707 plateau. Instead, it continues to improve as we further increase the training data from 25%
 708 to 100%. We fit linear trend lines (dotted in Figure 9) to the results and observe that the
 709 slopes are similar for each model. This suggests that the rate of improvement—reflected by
 710 the slope in the power-law relationship—is intrinsic to the data and task rather than to the
 711 specific model architecture. In other words, LLMs exhibit comparable data efficiency, with
 712 performance gains that are fundamentally tied to dataset size rather than model-specific
 713 factors.

714 Using these trend lines, we can estimate the amount of fine-tuning data required to reach a
 715 target performance. For instance, we estimate that fine-tuning Mistral-7B on a dataset 25
 716 times larger than the current SubPOP training set would yield a WD value of 0.07, which
 717 is much closer to the empirical lower bound of 0.031 reported in Table 1. This result
 718 underscores the critical importance of collecting more high-quality data, as increased dataset
 719 size can drive significant improvements in model performance.

720 F.2 Effect of Response Distribution Modeling

721 In this section, we compare different methods for capturing the distribution of human
 722 responses. We consider three approaches:

- 723 1. *One-hot*: Predicting only the most probable response, which ignores the full distri-
 724 bution over all responses Li et al. (2024).
- 725 2. *Augment by N*: Augmenting the dataset by replicating each response by a factor of
 726 N according to its observed frequency Zhao et al. (2023).
- 727 3. *Explicit probability modeling*: Directly modeling the full response distribution using
 728 the actual probability values for each option.

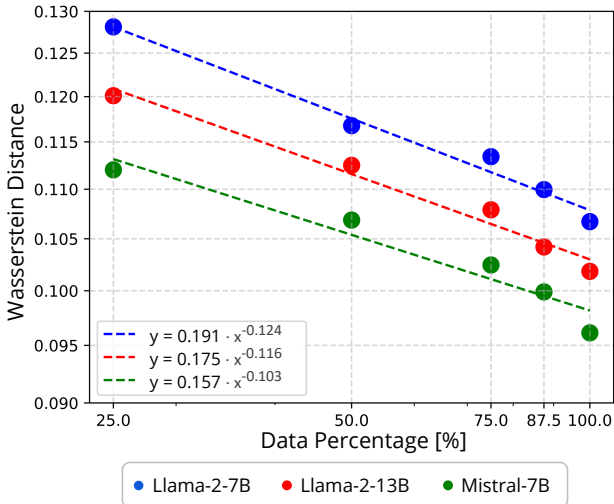


Figure 9: Evaluation results on OpinionQA after fine-tuning each LLM on increasingly large sampled subsets of SubPOP-Train. Both axes are presented in a log scale. The x -axis is the size of sampled dataset and the y -axis is WD against human responses measured on OpinionQA. Dashed lines represent a line of best fit. Performances at data percentage of 100% are identical to ours (SubPOP-FT) in Table 1.

729 Table 5 summarizes the results of these approaches. Notably, explicit probability modeling
 730 substantially outperforms the one-hot method, demonstrating that simply predicting the
 731 single most frequent response fails to capture the opinion diversity present within each
 732 subpopulation.

733 Compared with augment by N (2nd and 3rd column in Table 5), explicit probability mod-
 734 eling also achieves better performance. Importantly, the performance gap exceeds the
 735 quantization error introduced by discretizing the response distribution. For instance, when
 736 discretizing with a factor of N , the quantization error is $\frac{1}{2N}$ —approximately 0.01 or 0.005
 737 in the cases shown in Table 5. Moreover, explicit modeling offers the practical benefit of
 738 reducing the data volume by a factor of N compared to the augmentation approach, thereby
 739 lowering the computational cost of fine-tuning LLMs.

740 These results underscore the importance of explicit distribution modeling. By aligning the
 741 model’s predictive distribution directly with the survey distribution, we achieve higher
 742 accuracy with fewer data samples, avoiding the rounding errors and replication overheads
 743 that are inherent to data-augmentation approaches.

744 **F.3 Post-trained Model**

745 We fine-tune Llama-2-7B-chat to observe the effect of starting from checkpoints that have
 746 been instruction-tuned via Reinforcement Learning from Human Feedback (RLHF). Table 6
 747 shows the evaluation performance of a baseline method (Zero-shot prompt (QA)), fine-tuned
 748 base model (Llama-2-7B) and fine-tuned chat model (Llama-2-7B-chat). We observe the
 749 significant performance improvement, while the baseline method performs worse than
 750 the models not instruction-tuned (Table 1). Especially, the performance for SubPOP-Eval of
 751 chat model is significantly worse than that of base model. We observe the high WD of the
 752 baseline method resulting from the model assigning high probability to a specific token (e.g.
 753 ‘A’), being far apart from the human opinion distribution. After fine-tuning the model are
 754 able to generate a more distributed probability of answer tokens. This result coincides with
 755 the result reported in Moon et al. (2024).

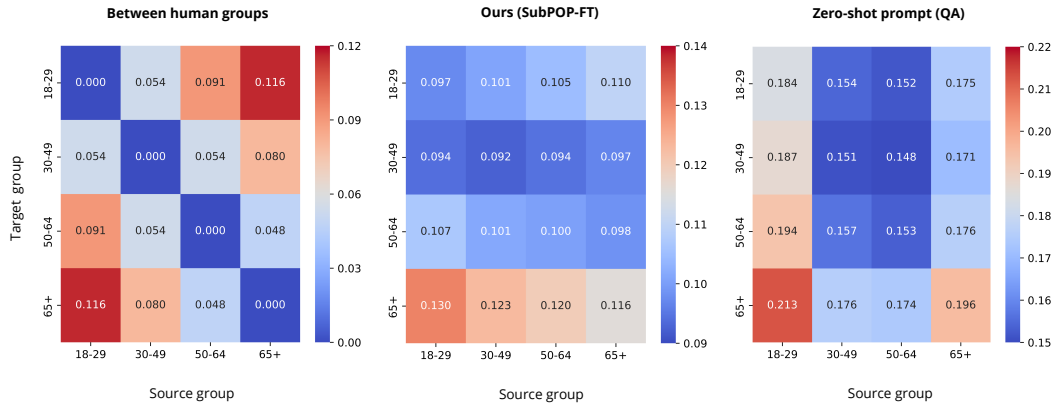


Figure 10: Heatmap of intergroup disagreement between a target human group (y -axis) and a source group (x -axis, either a human group or a group simulated with the language model), for OpinionQA evaluation data and age trait using Llama-2-7B as a base model. All subpopulations are unseen during fine-tuning.

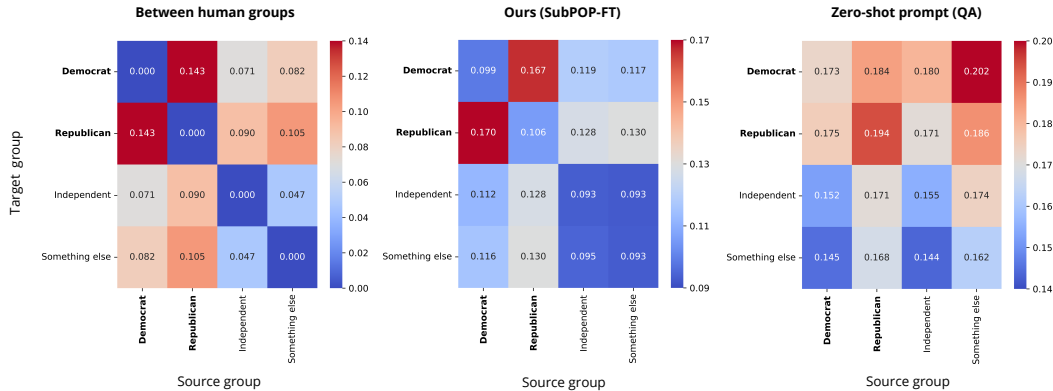


Figure 11: Heatmap of intergroup disagreement between a target human group (y -axis) and a source group (x -axis, either a human group or a group simulated with the language model), for OpinionQA evaluation data and political party (affiliation) trait using Llama-2-7B as a base model. Two subpopulations, Democrat and Republican, are seen during fine-tuning, while Independent and Something Else are unseen.

756 F.4 Generalization to Unseen Subpopulations

757 Here we present a complete list of evaluation performance on OpinionQA for unseen
 758 subpopulations (the groups not used to fine-tune our model) and perform an analysis that
 759 shows our fine-tuned models are able to steer towards the given subpopulation information.

760 As shown in Table 7, we observe a performance improvement across unseen subpopulations.
 761 To verify that the performance improvements arise from the fine-tuned model being able
 762 to steer towards given subpopulations, we measure *intergroup disagreement* pattern for the
 763 demographic and ideology traits, shown in Figure 10, 11, 12, and 13. We consistently observe
 764 across traits that the disagreement pattern of our model resembles that of the human group,
 765 while zero-shot prompting with the base model exhibits a pattern completely different
 766 from the human group result. This observation shows that our fine-tuned model learns to
 767 condition on subpopulation information and also generalizes to subpopulations unseen
 768 during fine-tuning.

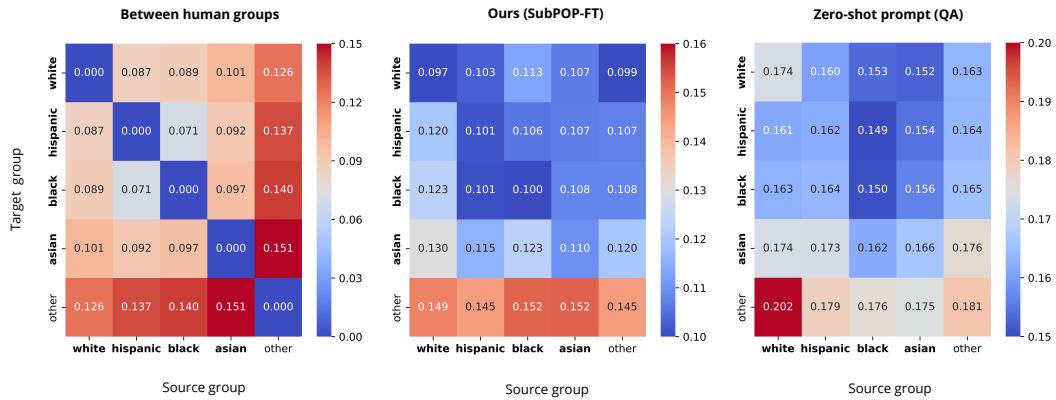


Figure 12: Heatmap of intergroup disagreement between a target human group (y -axis) and a source group (x -axis, either a human group or a group simulated with the language model), for OpinionQA evaluation data and race / ethnicity trait using Llama-2-7B as a base model. Four subpopulations except ‘Other’ are seen during fine-tuning. In this case, the model does not well predict the opinions of Other group. We suspect this occurs because Other is a group with highly diverse race or ethnicity backgrounds, making it inherently difficult to infer its opinion distribution from those of White, Hispanic, Black, and Asian subpopulations.

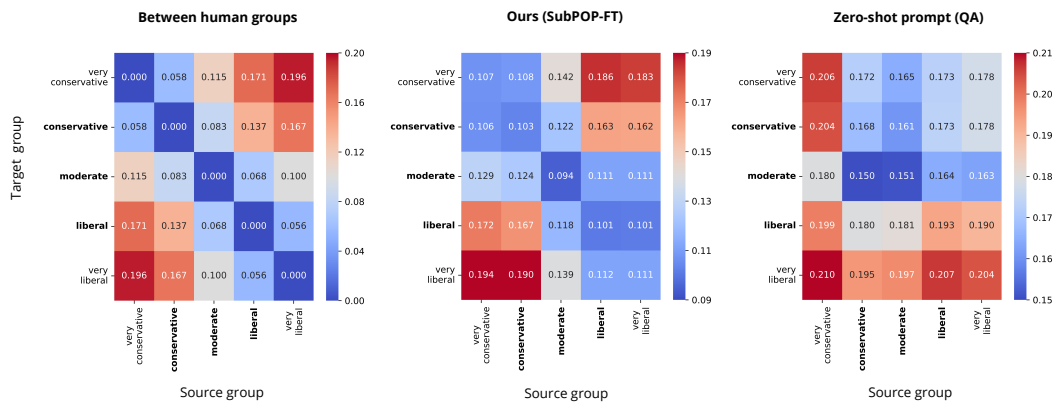


Figure 13: Heatmap of intergroup disagreement between a target human group (y -axis) and a source group (x -axis, either a human group or a group simulated with the language model), for OpinionQA evaluation data and political ideology trait using Llama-2-7B as a base model. Three subpopulations, Conservative, Moderate, and Liberal are seen during fine-tuning, while Very conservative and Very liberal are not seen.

769 G Baseline Details

770 **Zero-shot prompting:** Three prompt styles—QA, BIO, and PORTRAY—are introduced in [San-](#)
 771 [turkar et al. \(2023\)](#) to integrate group information into prompts. These prompts are then
 772 combined with survey questions to construct inputs for LLM. Then, the first-token log-
 773 probability from LLM is measured to calculate the model’s response distribution over
 774 options. In our baseline (and also in fine-tuning experiments) we focus on the QA steering
 775 format. Examples of this prompting method are shown in Figure 14.

776 **Few-shot prompting:** We craft a conditioning prompt that contains not only group infor-
 777 mation but also the group’s response distribution to k train questions, following [Hwang](#)
 778 [et al. \(2023\)](#). For a test question $q_{test} \in Q_{test}$, we first sort training questions Q_{train} into
 779 $\{q_1, q_2, \dots\}$ such that $\text{sim}(E(q_1), E(q_{test})) > \text{sim}(E(q_2), E(q_{test}))$, and so on. $E(q)$ denotes the

780 embedding model (OpenAI-text-embedding-3-large) output of the input q and sim is a
 781 cosine similarity between two embedding vectors. Then, response information of the first
 782 k questions $\{q_i, p(\mathcal{A}_{q_i}|q_i, g)\}_{i=1}^k$ are used as few shot prompts to have the language model
 783 verbalize Meister et al. (2024) expected response distribution for the given g and q_{test} . An
 784 example of the prompt for $k = 3$ case is shown in Figure 15, while we run the baseline
 785 experiment in a $k = 5$ setting.

786 **Modular Pluralism:** The intuition behind Modular Pluralism Feng et al. (2024) is that a
 787 language model trained on a text corpus of a specific subpopulation will faithfully represent
 788 public opinion of that population. Given a survey question with a PORTRAY-style steering
 789 prompt, each of language model ‘modules’ (fine-tuned Mistral-7B-Instruct-v0.1) generates
 790 an option choice with explanation. A black-box LLM (GPT-3.5-turbo-Instruct) receives all
 791 generations and select a generation that best aligns with the given group. Finally, using
 792 the chosen generation as a context, a black-box LLM generates probability distribution
 793 over options. The example pipeline is shown in Figure 16. Instead of the sub-sampled
 794 OpinionQA dataset the authors of the method used, we use the exactly same evaluation set
 795 across all baseline methods and our approach for a fair comparison.

796 **Upper bound:** We estimate the distribution between human responses and uniform distri-
 797 bution as an upper bound of WD metrics.

798 **Lower bound:** We compute a lower bound by randomly sampling a group of respondents
 799 and calculating the Wasserstein distance (WD) between the distribution of the sampled
 800 group and that of the original respondents for each question. We then bootstrap with
 801 $R = 1000$ to construct a 95% confidence interval (CI) for the WDs. Further details on this
 802 estimation process are provided below.

803 **Computing weighted answer distributions:** For each group g and question q , we have n_{gq}
 804 responses from respondents who belong to group g answering question q : $x_1, x_2, \dots, x_{n_{gq}}$,
 805 where $x_i \in \mathcal{A}_q$, i.e., the answer set for question q (e.g., $\{1, 2, 3, 4\}$). Furthermore, each respon-
 806 dent (and thus, their response) is associated with a wave-specific weight $w_1, w_2, \dots, w_{n_{gq}}$,
 807 provided by Pew Research. We compute the human answer distribution $\pi_{gq}^{(H)}$ as a weighted
 808 sum over responses, where the proportion of respondents providing answer $a \in \mathcal{A}_q$ is
 809 estimated as

$$\pi_{gq}^{(H)}(a) = \frac{\sum_{i=1}^{n_{gq}} w_i \mathbb{1}[x_i = a]}{\sum_{i=1}^{n_{gq}} w_i}.$$

810 **Bootstrapping at the respondent-level:** We draw bootstrap samples per group at the
 811 respondent-level including questions from all survey waves. This allows us to capture
 812 correlations in answer distributions across questions and across waves.

813 Specifically, let \mathcal{P}_g represent the set of respondents in group g , where $|\mathcal{P}_g| = n_g$. We produce
 814 bootstrapped samples by repeatedly sampling n_g respondents from \mathcal{P}_g with replacement.
 815 Let $p_1^{(r)}, p_2^{(r)}, \dots, p_{n_g}^{(r)}$ represent the sampled respondents for the r -th bootstrap, and let
 816 $w_1^{(r)}, w_2^{(r)}, \dots, w_{n_g}^{(r)}$ represent their corresponding weights.

817 For each question q , let $\mathcal{P}_{gq} \subseteq \mathcal{P}_g$ represent the set of respondents from group g who
 818 answered question q ; as before, $|\mathcal{P}_{gq}| = n_{gq}$. Let us define $q(p_i)$ as person p_i ’s response to
 819 question q if p_i answered question q , i.e., $p_i \in \mathcal{P}_{gq}$, and 0 otherwise. Then, we compute the
 820 r -th answer distribution to question q as:

$$\pi_{gq}^{(r)}(a) = \frac{\sum_{i=1}^{n_g} \mathbb{1}[p_i^{(r)} \in \mathcal{P}_{gq}] w_i^{(r)} \mathbb{1}[q(p_i^{(r)}) = a]}{\sum_{i=1}^{n_g} \mathbb{1}[p_i^{(r)} \in \mathcal{P}_{gq}] w_i^{(r)}}.$$

821 **Human lower bound of WD.** Our statistic of interest is the mean Wasserstein distance
 822 over all questions Q across all waves per group. We approximate this as the WD between

823 the observed human distribution $\pi_{gq}^{(H)}$ and the bootstrap sample $\pi_{gq}^{(r)}$ for question q and
824 group g . Over all $R = 1000$ bootstraps, we have

$$\mathcal{D}_g^{(H)} = \left\{ \frac{1}{|Q|} \sum_{q \in Q} WD(\pi_{gq}^{(H)}, \pi_{gq}^{(r)}) \right\}_{r=1}^R .$$

825 To quantify agreement between human samples, we report the mean and 95% CI (i.e., from
826 2.5th to 97.5th percentiles) of $\mathcal{D}_{gq}^{(H)}$.

827 H Wave, Group-level Opinion Matching

828 Here we present a group-level and wave-level averaged Wasserstein distance. Wave-level
829 result is in Table 8, and group-level results for OpinionQA and SubPOP-Eval are in Table 9,
830 10, respectively. We observe that the improvements in distribution matching between LLM
831 response and human response are consistent across diverse subpopulations and waves.

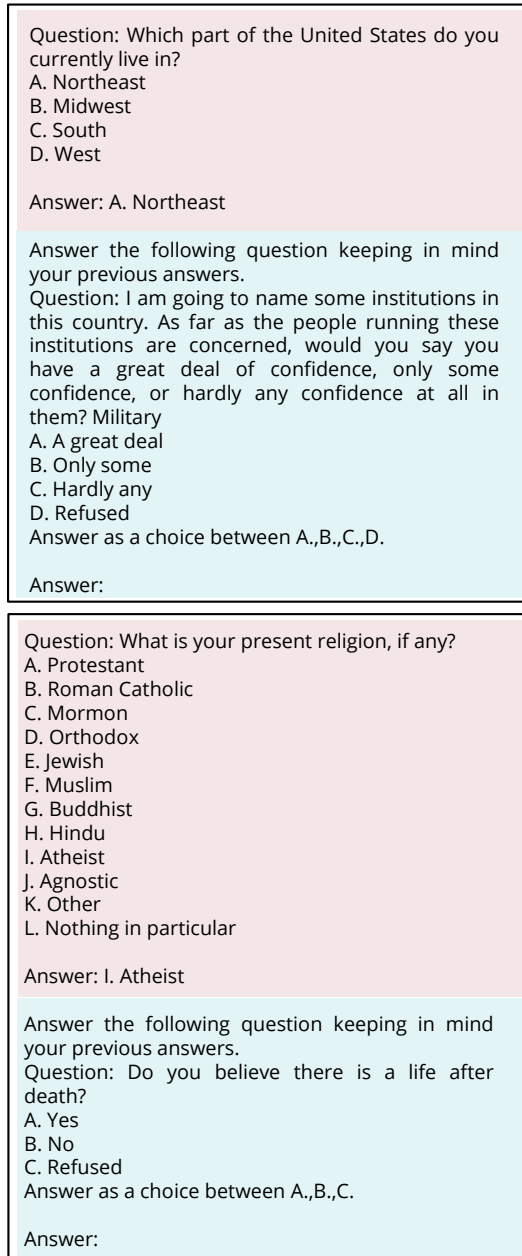


Figure 14: Two examples of Zero-shot prompting in the QA format [Santurkar et al. \(2023\)](#). Subpopulation’s information (colored in pink) is concatenated with survey question (colored in sky blue). The first-token log-probability (probabilities assigned to A, B, C, ...) are used to calculate language model’s response distribution. The same steering prompt format is used in our fine-tuning experiment.

Table 4: American Trends Panel (ATP) wave topics for waves included in SubPOP-Train (top) and OpinionQA (bottom). Golden rows represent wave topics in SubPOP-Train that are not present in OpinionQA, and blue rows represent wave topics in OpinionQA that are not present in SubPOP-Train. For waves 68-79, survey questions related to COVID-19 (e.g., contact tracing, vaccines, and relocation) were included as part of a survey along with the main survey topic.

Wave	# questions	Wave Topic
68	90	American News Pathways, George Floyd, Black Lives Matter
69	92	Politics, 2020 Census
70	56	Religion in public life, social media's role in politics and society
71	84	Voter attitudes
72	18	New media
73	82	American News Pathways, social media
74	51	Online harassment, race relations
75	18	2020 pre-election survey
76	44	American News Pathways
77	13	Culture of work
78	57	2020 post-election survey
79	93	American News Pathways
80	45	Political priorities
81	52	Economics, pandemic financial outlook
83	54	Coronavirus vaccines and restrictions
84	50	Religion in politics and tolerance
85	93	News coverage of the Biden administration's first 100 days
87	90	Current political news and topics
88	37	Tech companies and policy issues
90	79	Twitter news attitudes
91	64	Benchmark study
93	19	Social media update
95	78	Politics timely and topical
96	57	Post-coronavirus pandemic spirituality
98	76	Coronavirus impacts on communities, living arrangements and life decisions
99	20	Artificial intelligence (AI) and human enhancement
103	12	Economic well-being
104	92	Politics, Religion in Public Life
105	38	Global Attitudes US Survey 2022
106	62	Religion and the environment
107	92	Government and Parties
108	83	COVID and Climate, Energy and the Environment
109	51	New Digital Platforms and Gender Identity
110	90	Politics timely and topical
111	23	Online dating and E-commerce
112	31	Social media update
113	53	2022 National Survey of Latinos (NSL)
114	93	Covid, scientists, and religion
115	63	Parents survey
116	75	Politics timely and topical
117	16	Religion and politics
118	25	Podcasts, news, and racial identity
119	70	AI and human enhancement
120	61	Politics timely and topical
121	31	Culture of work
124	75	Global Attitudes US Survey 2023
125	69	Politics timely and topical
126	93	Racial attitudes, modern family
127	59	Americans and their data
128	89	Americans and planet Earth
129	107	Politics timely and topical
130	94	Politics representation
131	70	Gender and leadership
Wave	# questions	Wave Topic
26	44	Guns
29	20	Views on gender
32	24	Community types, Sexual harassment
34	16	Biomedical and food issues
36	68	Gender and leadership
41	41	Views of America in 2050
42	26	Trust in science
43	51	Race in America
45	13	Misinformation
49	19	Privacy and surveillance
50	43	American families
54	50	Economic inequality
82	56	2021 Global Attitudes Project U.S. survey
92	23	Political Typology

Table 5: Comparison of evaluation performance for three response distribution modeling approaches, with Llama-2-7B as a base model. The last column (Explicit) is identical to the ours presented in Table 1. A model fine-tuned to predict the most probable choice (one-hot) performs the worst, as the model has not learned distributional opinion at fine-tuning phase. A model trained on augmented data (Aug. ($\times 50$, $\times 100$)), while performing much better than one-hot still underperforms the explicit distribution modeling.

Eval Dataset	One-hot	Aug. ($\times 50$)	Aug. ($\times 100$)	Explicit (Ours)
OpinionQA	0.163	0.110	0.107	0.106
SubPOP-Eval	0.178	0.130	0.123	0.121

Table 6: Performance of the fine-tuned Llama-2-7B-chat model (Chat LLM FT). For comparison, we also present lower and upper bounds, the baseline method Zero-shot prompt (QA) and our fine-tuned Llama-2-7B (Base LLM FT).

Method	OpinionQA	SubPOP-Eval
Upper bound (Unif.)	0.178	0.208
Lower bound (Human)	0.031	0.033
Base zero-shot prompt (QA)	0.173	0.206
Base LLM FT	0.106	0.121
Chat zero-shot prompt (QA)	0.308	0.383
Chat LLM FT	0.109	0.148

Table 7: Evaluation performance of our fine-tuned Llama-2-7B model on OpinionQA for subpopulations not included in the fine-tuning dataset SubPOP-Train. For reference, we present a lower bound (human) and the zero-shot prompting (QA). Absolute difference refers to the WD difference between zero-shot prompting and ours, and the relative improvement is calculated in a same way as Figure 3. R. I. stands for Relative Improvement.

Attribute	Group	Lower Bound	Zero-shot (QA)	Ours	Absolute Diff.	R. I.
Age	18-29	0.023	0.185	0.096	0.089	0.548
Age	30-49	0.014	0.151	0.093	0.058	0.424
Age	50-64	0.014	0.154	0.101	0.052	0.377
Age	65+	0.013	0.195	0.115	0.080	0.438
Region	Midwest	0.016	0.153	0.095	0.058	0.425
Region	West	0.017	0.162	0.095	0.068	0.465
Education	Associate’s Degree	0.026	0.159	0.098	0.061	0.455
Education	High School Graduate	0.017	0.144	0.092	0.053	0.413
Education	Postgraduate	0.015	0.174	0.106	0.068	0.426
Education	Some College, No Degree	0.018	0.144	0.093	0.051	0.405
Income	\$50,000-\$75,000	0.016	0.153	0.098	0.054	0.396
Income	\$30,000-\$50,000	0.019	0.144	0.094	0.050	0.400
Political Ideology	Very Conservative	0.026	0.208	0.107	0.101	0.555
Political Ideology	Very Liberal	0.025	0.202	0.111	0.091	0.514
Political Party	Independent	0.016	0.155	0.093	0.062	0.445
Political Party	Something Else	0.026	0.162	0.092	0.069	0.510
Race	Other	0.050	0.180	0.144	0.036	0.275
Religion	Agnostic	0.028	0.189	0.115	0.074	0.459
Religion	Buddhist	0.063	0.207	0.149	0.059	0.405
Religion	Nothing in Particular	0.019	0.153	0.092	0.061	0.454
Religion	Orthodox	0.083	0.221	0.180	0.041	0.298
Religion	Other	0.051	0.184	0.123	0.061	0.457
Religion	Roman Catholic	0.018	0.145	0.098	0.047	0.371

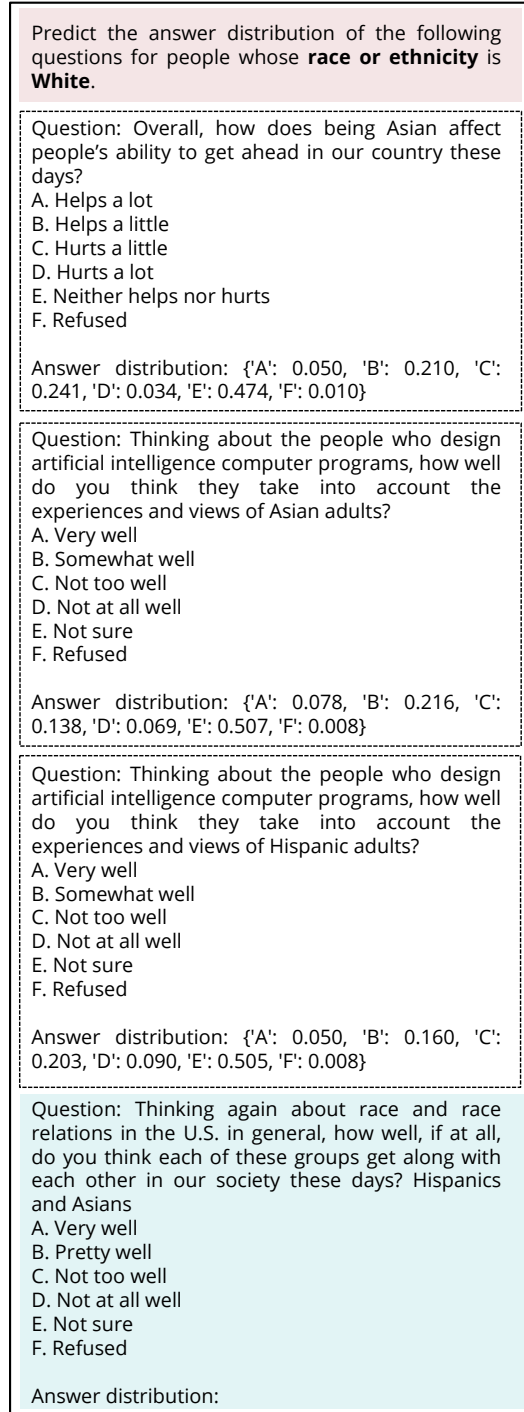


Figure 15: Few-shot prompting example for $k = 3$. Group information is presented in the beginning of the prompt (colored in pink). Following group information, k questions whose text embedding are the most similar to the text embedding of the evaluation question (colored in sky blue) are presented along with their opinion distribution. k questions are presented in the ascending order of cosine similarity. The generation of language model (verbalization of opinion distribution) is parsed to obtain the response distribution.

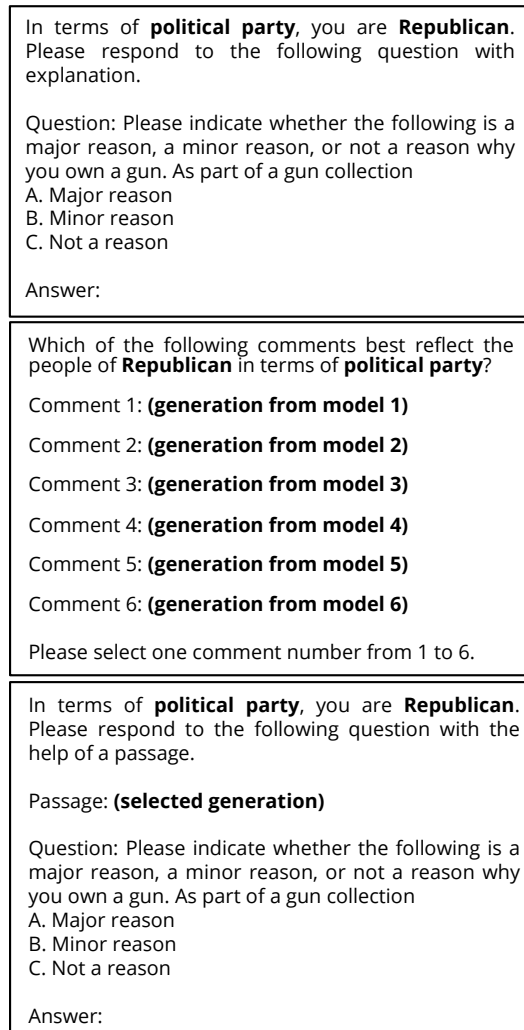


Figure 16: Pipeline example of Modular Pluralism. Given a demographic group and a survey question, the first prompt is asked to multiple (6) language models, Mistral-7B-v0.1-Instruct fine-tuned on the community text corpus. The generations are sent to a black-box LLM (gpt-3.5-0613-Instruct) in the format of the second prompt. The black-box LLM answers which one of generations best reflects the given demographics. Finally, the selected generation serves as a context to answer the given survey question and the black-box LLM is prompted (the third prompt) to generate response distribution over the answer token A, B, C, etc.

Table 8: Per-wave Wasserstein distance on OpinionQA for each base model, comparing baseline zero-shot prompting (QA) with our fine-tuned model Ours(SubPOP-FT). Highlighted rows represent waves whose topics are not covered by the training data (SubPOP-Train). We observe WD improvement consistently across survey waves and also for waves of topics not covered in the training data.

Wave	Llama-2-7B		Llama-2-13B		Mistral-7B-v0.1		Llama-3-70B	
	Zero-shot	Ours	Zero-shot	Ours	Zero-shot	Ours	Zero-shot	Ours
26	0.191	0.145	0.180	0.126	0.178	0.131	0.134	0.084
29	0.169	0.096	0.172	0.123	0.153	0.096	0.125	0.085
32	0.163	0.110	0.156	0.098	0.137	0.099	0.151	0.091
34	0.155	0.105	0.171	0.089	0.134	0.095	0.138	0.083
36	0.175	0.120	0.184	0.126	0.175	0.107	0.130	0.087
41	0.160	0.090	0.155	0.084	0.134	0.073	0.116	0.085
42	0.159	0.053	0.146	0.059	0.127	0.059	0.131	0.084
43	0.179	0.112	0.172	0.104	0.154	0.102	0.124	0.099
45	0.177	0.101	0.177	0.093	0.149	0.084	0.126	0.091
49	0.151	0.098	0.143	0.131	0.128	0.116	0.159	0.087
50	0.209	0.139	0.196	0.121	0.188	0.125	0.154	0.078
54	0.158	0.087	0.158	0.087	0.128	0.077	0.118	0.079
82	0.173	0.098	0.171	0.075	0.148	0.077	0.174	0.093
92	0.165	0.073	0.153	0.071	0.140	0.055	0.126	0.081

Table 9: Per-group Wasserstein distance on OpinionQA for each base models, before and after fine-tuning on SubPOP-Train. Base refers to zero-shot prompting (QA) and FT refers to Fine-tuned model performance. *Full group variable name is “College grad, some Postgrad”.

Attribute	Group	Human Baseline	Llama-2-7B		Llama-2-13B		Mistral-7B-v0.1		Llama-3-70B	
			Base	FT	Base	FT	Base	FT	Base	FT
Region	Northeast	0.023	0.165	0.094	0.155	0.088	0.155	0.083	0.134	0.084
	South	0.017	0.149	0.092	0.143	0.085	0.133	0.081	0.113	0.078
Education	College grad*	0.018	0.165	0.099	0.157	0.096	0.136	0.089	0.125	0.085
	Less than high school	0.043	0.161	0.101	0.150	0.096	0.134	0.094	0.151	0.091
Gender	Male	0.015	0.182	0.093	0.152	0.089	0.131	0.083	0.138	0.083
	Female	0.013	0.162	0.100	0.158	0.092	0.146	0.088	0.130	0.087
Race / ethnicity	Black	0.031	0.151	0.102	0.144	0.095	0.132	0.091	0.116	0.085
	White	0.012	0.176	0.097	0.178	0.093	0.145	0.085	0.131	0.084
	Asian	0.051	0.165	0.111	0.167	0.104	0.143	0.102	0.124	0.099
	Hispanic	0.044	0.162	0.102	0.163	0.098	0.134	0.092	0.126	0.091
Income	\$100,000 or more	0.019	0.172	0.103	0.162	0.100	0.147	0.091	0.159	0.087
	Less than \$30,000	0.021	0.162	0.091	0.148	0.083	0.127	0.080	0.154	0.078
Political Party	Democrat	0.016	0.172	0.099	0.158	0.092	0.161	0.082	0.118	0.079
	Republican	0.019	0.196	0.105	0.235	0.101	0.181	0.095	0.174	0.093
Political Ideology	Liberal	0.022	0.192	0.100	0.181	0.094	0.166	0.084	0.126	0.081
	Conservative	0.021	0.169	0.103	0.153	0.099	0.144	0.094	0.141	0.092
	Moderate	0.016	0.151	0.094	0.153	0.090	0.132	0.082	0.106	0.081
Religion	Protestant	0.016	0.015	0.166	0.096	0.158	0.092	0.146	0.086	0.143
	Jewish	0.058	0.182	0.124	0.182	0.122	0.165	0.115	0.144	0.115
	Hindu	0.079	0.211	0.160	0.232	0.163	0.211	0.161	0.181	0.157
	Atheist	0.035	0.202	0.118	0.204	0.110	0.196	0.099	0.135	0.098
	Muslim	0.089	0.202	0.159	0.209	0.156	0.204	0.146	0.171	0.144

Table 10: Per-group Wasserstein distance on SubPOP-Eval for each base models, before and after fine-tuning on SubPOP-Train. Base refers to zero-shot prompting (QA). *Full group variable name is “College grad, some Postgrad”.

Attribute	Group	Human	Llama-2-7B		Llama-2-13B		Mistral-7B-v0.1		Llama-3-70B	
		Baseline	Base	FT	Base	FT	Base	FT	Base	FT
Region	Northeast	0.027	0.196	0.113	0.193	0.103	0.185	0.108	0.156	0.078
	South	0.018	0.183	0.108	0.185	0.103	0.176	0.103	0.138	0.080
Education	College grad*	0.019	0.206	0.105	0.175	0.101	0.167	0.099	0.137	0.077
	Less than high school	0.036	0.191	0.129	0.182	0.117	0.172	0.121	0.180	0.108
Gender	Male	0.017	0.186	0.102	0.176	0.101	0.170	0.099	0.150	0.079
	Female	0.016	0.184	0.108	0.198	0.105	0.176	0.100	0.151	0.080
Race / ethnicity	Black	0.029	0.200	0.114	0.179	0.102	0.170	0.107	0.139	0.094
	White	0.014	0.190	0.105	0.187	0.103	0.181	0.102	0.153	0.083
	Asian	0.049	0.201	0.119	0.190	0.107	0.184	0.114	0.158	0.096
	Hispanic	0.050	0.204	0.133	0.199	0.122	0.182	0.134	0.172	0.115
Income	\$100,000 or more	0.021	0.210	0.111	0.184	0.106	0.176	0.102	0.179	0.082
	Less than \$30,000	0.026	0.179	0.115	0.172	0.103	0.165	0.105	0.171	0.086
Political Party	Democrat	0.020	0.219	0.103	0.197	0.092	0.199	0.091	0.128	0.076
	Republican	0.023	0.205	0.123	0.234	0.117	0.206	0.115	0.187	0.093
Political Ideology	Liberal	0.019	0.224	0.102	0.191	0.090	0.188	0.096	0.134	0.076
	Conservative	0.022	0.184	0.120	0.178	0.112	0.172	0.113	0.160	0.092
	Moderate	0.018	0.191	0.110	0.183	0.103	0.170	0.103	0.141	0.082
Religion	Protestant	0.019	0.187	0.110	0.179	0.107	0.172	0.105	0.164	0.082
	Jewish	0.066	0.245	0.149	0.226	0.144	0.218	0.129	0.164	0.119
	Hindu	0.095	0.264	0.180	0.253	0.169	0.252	0.186	0.223	0.166
	Atheist	0.021	0.222	0.126	0.207	0.103	0.199	0.116	0.132	0.106
	Muslim	0.090	0.253	0.175	0.240	0.181	0.238	0.173	0.203	0.158