# Algorithmic Fidelity of Large Language Models in Generating Synthetic German Public Opinions: A Case Study

Anonymous authors

Paper under double-blind review

## Abstract

In recent research, large language models (LLMs) have been increasingly used 1 to investigate public opinions. This study investigates the *algorithmic fidelity* of 2 LLMs, i.e., the ability to replicate the socio-cultural context and nuanced opinions 3 of human participants. Using open-ended survey data from the German Longitu-4 dinal Election Studies (GLES), we prompt different LLMs to generate synthetic 5 public opinions reflective of German subpopulations by incorporating demographic 6 features into the persona prompts. Our results show that Llama performs better 7 than other LLMs at representing subpopulations, particularly when there is lower 8 opinion diversity within those groups. Our findings further reveal that the LLM 9 performs better for supporters of left-leaning parties like The Greens and The 10 Left compared to other parties, and matches the least with the right-party AfD. 11 Additionally, the inclusion or exclusion of specific variables in the prompts can 12 significantly impact the models' predictions. These findings underscore the impor-13 tance of aligning LLMs to more effectively model diverse public opinions while 14 minimizing political biases and enhancing robustness in representativeness. 15

## 16 **1 Introduction**

Recent advances in LLMs have generated significant interest in their potential for synthetic data
generation across various domains. A key and widely debated question is whether LLMs can produce
synthetic data that accurately represent human opinions (Argyle et al., 2023; Santurkar et al., 2023;
Veselovsky et al., 2023; von der Heyde et al., 2025; Long et al., 2024, *inter alia*).

In social science research, it is usually surveys that provide insights into the attitudes and opinions

of a population. Recent studies have explored using information from survey responses for LLM 22 prompts, i.e., creating so-called personas and then using the LLM "responses" to additional questions 23 (Argyle et al., 2023; Dominguez-Olmedo et al., 2024; Durmus et al., 2024, inter alia). Assessing 24 the fidelity of LLMs in capturing and reproducing human opinions deepens our understanding of 25 model behavior while at the same time helping researchers design more reliable models aligned with 26 human values and improving their usability (Ma et al., 2024). Among these studies, Argyle et al. 27 (2023) introduced the *algorithmic fidelity*, a concept for assessing how effectively LLMs replicate 28 the socio-cultural context and nuanced opinions of diverse human subpopulations. Their work used 29 LLMs to generate synthetic responses ("silicon samples") informed by demographic and ideological 30 profiles from political science datasets. Their findings suggest that LLMs can closely approximate 31 real-world opinion distributions in certain contexts, such as U.S. elections, although challenges persist 32 in fully aligning the generated data with actual demographic groups. 33

Recent research on LLM responses to opinion polling have predominantly focused on US-based and 34 English-centric survey data. For example, von der Heyde et al. (2025) evaluated the algorithmic 35 fidelity of GPT-3.5 in predicting German voting behavior. Their findings revealed that while GPT-3.5 36 accurately represented the voting patterns of center and left-leaning political groups, it struggled to 37 capture those of right-leaning parties. However, like many similar studies, their analysis was limited 38 to closed-ended survey questions with single-choice responses. This highlights a key challenge: while 39 LLMs may perform well in English-speaking contexts, less is known about their ability to generate 40 representative opinions in non-English-speaking countries and for open-ended questions. This is 41 particularly true for open-ended setups, where scaling and accurately interpreting responses pose 42 significant challenges (Resnik et al., 2024). 43

- To address these challenges, this study explores the *algorith*mic fidelity of LLMs in generating synthetic public opinions in an open-ended survey question based on German survey data. We use the survey question about the most important problems facing Germany today from the German Longitudinal Election Study (GLES). The survey is a longitudinal panel survey, and the answer distribution can be found in Figure 1. We select variables from the original survey data to represent survey participants with their characteristics as personas. We include three LLMs (Gemma, Llama2, Mixtral) in our study and prompt them to simulate survey participants to answer the open-ended question about the most important political problem in Germany. Finally, we compare the outputs regarding the distributional alignment with original survey answers, predictive performance, and answer diversity. Our most important findings are:
  - 1) Llama2 is better at modeling group opinions (§4.1).
  - 2) Llama2's representativeness fluctuates across survey waves, with the model's representativeness of the population decreasing as survey diversity increases; the model represents subpopulation opinions unevenly, with favorable alignment April 2020, with it becoming the domfor left-leaning parties (such as the Left, the Greens) over right-parties (such as AfD) (§4.2).
  - 3) Including more variables in prompts improves performance, of COVID-19 in early 2020. with party affiliation being the most influential factor (§4.3).

Figure 1: The distribution of the top 5 answer categories between November 2019 and November 2021 in the German GLES survey. There is a significant surge in the Health Policy category from November 2019 to inant focus during this period and afterwards, likely due to the outbreak

#### 2 **Related Work** 45

**LLMs for Survey Response Generation.** Recent studies have increasingly repurposed survey 46 questionnaires, originally designed for public opinion polling, to assess the opinions generated by 47 LLMs (Ma et al., 2024). For instance, Santurkar et al. (2023) identified significant differences between 48 LLM opinion distributions and US-based survey participants. Similarly, Dominguez-Olmedo et al. 49 (2024) highlighted disparities between LLM and human opinions, emphasizing the sensitivity of 50 model outputs to biases in prompting. Tjuatja et al. (2024) found that LLMs are highly sensitive to 51 prompt perturbations and fail to replicate human-like behavior. Collectively, these studies suggest 52 that LLMs align more closely with populations holding left-leaning, Western-oriented values. 53 **Opinion Generation in the German Contexts.** While most studies on opinions in LLM output are 54

English- and U.S.-centric, some research has explored other contexts, such as the German case. In a 55 recent study, von der Heyde et al. (2025) employed the data of 2017 post-election cross-section of the 56 GLES. Respondents to this study reported their vote choice in the survey. von der Heyde et al. (2025) 57 prompted GPT-3.5 (Brown et al., 2020) with personas to simulate the survey participants. Based 58 on the close-ended choice setup, they found that it does not predict citizens' vote choice accurately, 59 exhibiting a bias towards the Green and Left parties, similar to previous work in English contexts. 60 **Evaluation of LLM Outputs.** Previous studiesprimarily focused on closed-ended multiple-choice 61

62 questions, often relying on the model's first token prediction (e.g., Santurkar et al., 2023; Dominguez-Olmedo et al., 2024; Tjuatja et al., 2024) or semi-automated extraction of text answers (von der 63 Heyde et al., 2025). Alternatively, Wang et al. (2024a;b) proposed training a classifier directly on 64 LLM responses manually labeled by annotators, finding this method more robust. After the output 65 extraction, a few evaluation metrics have been applied to measure the alignment of human and 66 LLM responses (Ma et al., 2024), such as Cohen's Kappa (Argyle et al., 2023; Hwang et al., 2023), 67 1-Wasserstein distance (Santurkar et al., 2023; Hwang et al., 2023), KL divergence (Dominguez-68 Olmedo et al., 2024; Sun et al., 2024), Euclidean distance (Wang et al., 2023), Jensen-Shannon 69 distance (Durmus et al., 2024), etc. and correlation and statistical analyses (Sun et al., 2024; Jiang 70 et al., 2024b). For our case study, we adapt these metrics to examine the fidelity of LLM-generated 71

synthetic German public opinions. 72

## 73 **3 Experimental Setups**

## 74 3.1 Data

German Longitudinal Election Study (GLES Panel). We use the GLES Panel dataset from
 GESIS (2023). The survey consists of 21 waves<sup>1</sup> and contains socio-demographic information,
 vote intentions, choices, and political attitudes of participants. The target population is German
 citizens eligible to vote during the respective elections in Germany. Along with the respondents'
 socio-demographic data for our prompts, we use respondents' answers to the question "In your
 opinion, what is the most important problem facing Germany today?" for comparing human answers
 and LLM outputs. The answers of participants were collected as free-form texts.

Selected Information. We included six variables from the original survey: age, gender, leaning
 party, region, education degree, and vocational degree. Details on the sub-groups of the variables are
 shown in §A.

Coding Scheme. For coding the LLM text responses into categories, we follow the coding scheme proposed in an additional sub-study of the GLES (GESIS, 2024). Like Mellon et al. (2024), who collapsed ~50 classes into a simpler classification, we also set "coarse" classes (n=16). We merged rarely represented classes into an upper class (e.g., "Price Level", "Housing Policy", and "Economic Policy" classes into one "Economic Policy" class). The distribution of GLES survey answers based on the coarse classes is shown in Figure 1. The full list of fine and coarse classes can be found in §B in the Appendix.

## 92 3.2 Text Generation

93 Models. We chose three instruction-tuned open-weight LLMs: Llama-2-13b-chat-hf (Tou-94 vron et al., 2023), Gemma-7b-it (Team, 2024), and Mixtral-8x7B-Instruct-v0.1 (Jiang 95 et al., 2024a).

96

**Prompt Design.** After initial trial runs and inspecting the LLM outputs, we used the prompt in Figure 2 in our experiments with LLMs. We chose German as the prompting language, as it is the language in which the GLES survey was conducted, and the generated texts can be compared to the original text. During the experimentation phase, the placeholders are replaced with the respondent's information, i.e., the variables from the survey data.

Identify the most important problem Germany in {month} {year} is facing. Provide the answer in one concise sentence, focusing on a single issue without elaborating or listing additional problems. Do not repeat the information you have been given and give your answer directly and without introductory phrases. Answer in German and only in German, do not use English. Answer from the perspective of a respondent with German citizenship and the characteristics specified below. {article} (The) respondent is {age} years old and {gender}. {pronoun} {educational\_qualification\_clause} and {vocational\_qualification\_clause} {pronoun2} lives in {region} and mainly supports {party}.

Figure 2: Translated prompt in English. The original prompt in German is presented in Figure 10 in Appendix.

## 97 3.3 LLM Output Classification

To evaluate and compare the LLM outputs with human responses, we needed to categorize the responses into specific classes and trained a classifier to code the responses from the LLMs.

Manual Annotation. Drawing from Wang et al. (2024a;b), we manually annotated 1,500 LLM
 outputs, selecting 500 outputs randomly from each LLM. We then trained a classifier based on the
 manually developed annotation scheme for the LLM outputs. Details on the scheme are in §C.

Classifier Training and Inference. We fine-tuned the German version of the base BERT (Devlin et al., 2019) classifier on the annotated LLM outputs. The fine-tuned classifier achieves a weighted F1 score of 0.93 on the test set. The classifier is then used to classify all LLM responses.

<sup>&</sup>lt;sup>1</sup>A survey wave refers to a single round of data collection in a panel survey, gathering information from the same participants at multiple points in time (Andreß et al., 2013).

#### **106 3.4 Evaluation Metrics for Experiments**

- <sup>107</sup> In the context of generative models, representativeness is the model's ability to recover population-
- level properties of the original data (Eigenschink et al., 2023), i.e., a dimension of algorithmic fidelity.

109 To compare the representativeness of the LLM answers with the original survey data and to measure

the association between the variables, we used the following evaluation metrics.

Jensen-Shannon (JS) Divergence. JS divergence is a symmetric and normalized measure of divergence derived from KL divergence (Kullback & Leibler, 1951). It is calculated as:

$$JSD(P \parallel Q) = \frac{1}{2}D_{KL}(P \parallel M) + \frac{1}{2}D_{KL}(Q \parallel M)$$
(1)

where  $M = \frac{1}{2}(P + Q)$  is the mixture distribution of *P* and *Q* (Lin, 2006). The JS divergence is bounded between 0 and 1 (when using  $log_2$ ), making it easier to interpret than KL divergence. We use JS Distance, the square root of JS Divergence, as in Durmus et al. (2024), because its bounded range facilitates comparison across different data waves. The JS distance is applied to measure the representativeness of the coded LLM answers compared to the real survey data.

**Entropy.** Entropy measures the variability or uncertainty in a set of outcomes (Jurafsky & Martin, 2024):

$$H(X) = -\sum_{x \in \mathcal{X}} p(x) \log p(x)$$
(2)

We use entropy to assess the diversity of text categories in synthetic and survey data. Lower entropy indicates less variability, meaning fewer bits are needed to represent the information in the data.

122 **Conditional Entropy.** Conditional entropy measures the remaining uncertainty in variable X when 123 another variable Y is known. It calculates the entropy of X given the distribution of Y:

$$H(X \mid Y) = \sum_{y \in \mathcal{A}_Y} P(y) \left[ \sum_{x \in \mathcal{A}_X} P(x \mid y) \log \frac{1}{P(x \mid y)} \right] = \sum_{x \in \mathcal{A}_X} \sum_{y \in \mathcal{A}_Y} P(x, y) \log \frac{1}{P(x \mid y)}$$
(3)

We use conditional entropy to evaluate how much uncertainty remains about responses in the survey when the subpopulation is known. This helps assess whether the synthetic data captures patterns in specific groups within the population.

127 Information Gain. Also called mutual information. It measures how much information one random 128 variable provides about another. It is calculated as the difference between the entropy of the variable 129 and its conditional entropy given another variable:

$$I(X;Y) = H(Y) - H(Y \mid X)$$
(4)

It indicates how much knowing one variable (e.g., X) reduces uncertainty about another variable (e.g., Y). A higher information gain indicates that knowing one variable reduces uncertainty about another variable. In our experiments, we calculate the population entropy H(Y) and conditional subpopulation entropy H(Y | X), where X represents demographic features. We will compute H(Y)and H(Y | X) for subpopulations and compare the information gained in survey and LLM data.

**Cramér's V.** This is a measure of association between nominal variables (Cramér, 1999). It is based on Pearson's  $\chi^2$  test. However, Cramér's V discounts the value of the  $\chi^2$  statistic for both the sample size (N) and the size of the table of counts (minimum of row count or column count minus 1) (Holbrook, 2022). It is computed as:

$$V = \sqrt{\frac{\chi^2}{N \cdot \min(r - 1, c - 1)}} \tag{5}$$

We use Cramérs' V to check "pattern correspondence" in LLM outputs. We map each input variable  $(X_i)$  to the output variable (Y), and check whether the pairwise correlations in survey data are also

141 present in the LLM-generated data.

## 142 **4** Experiments and Results

Three main experiments were conducted on the GLES data. The first evaluated all three LLMs using a single wave, focusing on dataset statistics and representativeness (§4.1). The second extended this analysis across multiple waves with Llama2 to track performance over time (§4.2). The third involved ablation studies to assess how different variables affect representativeness and response diversity (§4.3).

#### 148 4.1 Experiment 1: Model Pre-Experiment in One Wave

After dropping the observations with missing features, we used the survey data from wave 12 (collected between 05-11-2019 and 19-11-2019, i.e., before COVID-19) for the first experiment. All three LLMs have been prompted to generate synthetic answers. The generated answers are classified as stated in §3.3. We compared dataset statistics and the textual style of the answers and computed JS Distance and entropy.

<sup>154</sup> In Table 1, we give an overview of statistics about labels, textual characteristics, and representativeness (on the population level). 42 % of the Gemma model answers were about COVID-19 (identified using Regex), even though the survey answers were collected before COVID-19, which indicates a very large proportion of hallucinations. Therefore, we did not include the Gemma model in the further subpopulation-level analysis. The detailed JS Distances in each social group category can be found in Table 13 in the Appendix.

| Metric                 | Gemma     | Llama2             | Mixtral               | Survey            |
|------------------------|-----------|--------------------|-----------------------|-------------------|
| Avg. Labels per Sample | 1.03      | 1.20               | 1.33                  | 1.03              |
| Avg. Samples per Label | 593       | 692                | 769                   | 597               |
| Non-German Answer Rate | 0.02<br>0 | 25.65<br>0.06<br>0 | 43.75<br>0.03<br>0.05 | 2.29<br>-<br>0.04 |
| LLM Refusal Rate       | 1e-4      | 0                  | 1e-4                  | 0                 |
| COVID Regex Match Rate | 0.42      | 0.03               | 2e-3                  |                   |
| JS Distance to Survey  | 0.62      | 0.28               | 0.29                  | 2.93              |
| Answer Entropy         | 2.26      | 2.90               | 2.56                  |                   |

Table 1: Survey and LLM data statistics in experiment 1

A case study on information gain of the party variables: Llama2 aligns more closely with survey 155 data and maintains subgroup stability. Figure 3 compares population-level answer entropies (left) 156 with conditional entropies (right) for each leaning party value. Information gain, calculated as 157 the difference between these entropies, reflects how much additional insight is provided by knowing 158 the leaning party value. The population-level entropies (H(Y)) are close, with the survey 159 (2.93) closely matching Llama2 (2.90), while Mixtral was a bit lower (2.56) (see left of each subplot). 160 After incorporating the leaning party information and only looking at the samples containing 161 the specific party affiliations  $(H(Y \mid X_i))$ , Mixtral shows still lower conditional entropy (see right of 162

each subplot), indicating less variation in responses. Especially for "Die Grünen (The Greens)" and "AfD", there are drastic drops of  $H(Y | X_i)$ . This suggests that Mixtral may risk reflecting dominant group opinions, reducing diversity, and showcasing stereotypical representations of these subgroups.

165 group opinions, reducing diversity, and showcasing stereotypical representations of these subgroups. In contrast, Llama2 exhibits less information gain, i.e., it is more aligned with the survey data.



Figure 3: Information Gain for leaning party variable  $(X_i)$ . Left: population entropy (H(Y)), right: subpopulation entropy  $(H(Y | X_i))$ . A large gap between left and right  $(H(Y) - H(Y | X_i))$  means big information gain when focusing on the samples of the subpopulation group, indicating responses with this group are less diverse.

#### 167 4.2 Experiment 2: Wave Experiment with Llama2

We focused on Llama2 in the second in-depth experiment and repeated the generation process for the most 10 recent panel waves in GESIS (2023) (waves 12-21). Over these two years, we observe large shifts in survey label distributions (see Figure 1). This allows us to evaluate the representativeness of

the model under varying label distributions and seek answers to the following questions:

• Do the LLM capabilities at both the population and subpopulation levels vary over time?

• In which subpopulations are opinions represented more accurately?

Llama2 captures shifts in survey trends, but higher answer diversity correlates with reduced 174 representativeness. For panel waves 12-21, we repeated the text generation process and classified 175 the answers. Table 2 shows the population-level entropy values and the JS distances. We calculated 176 Pearson's correlation coefficient between survey entropy and the JS distance and got r = -0.35, 177 indicating that the model's representativeness of the population decreases as the diversity in answers 178 increases. For wave 13, with data collection between April 2020 and May 2020, the diversity of 179 answers reached its minimum (with an entropy of 0.58). In Table 12, we see that 92.4 % of answers 180 were about "Health Policy" (and about the COVID-19). This shows that LLMs' responses reflect the 181 change in the survey date. 182

|             | 12   | 13   | 14   | 15   | 16   | 17   | 18   | 19   | 20   | 21   | avg. |
|-------------|------|------|------|------|------|------|------|------|------|------|------|
| LLM entropy | 2.90 | 0.58 | 1.67 | 1.31 | 2.12 | 2.20 | 2.27 | 2.46 | 2.46 | 2.49 | 2.04 |
| JS distance | 0.29 | 0.29 | 0.24 | 0.22 | 0.20 | 0.23 | 0.23 | 0.22 | 0.24 | 0.30 | 0.24 |

Table 2: Population level entropy values and the JSD in the Wave Experiment from wave 12-21.

Subpopulation-level findings: JS distances reveal representational variation influenced by 183 group information and complexity. Figure 4 shows the JS distances at the subpopulation level for 184 each variable. We observed the most variation for education and leaning party variables. 185 Although the difference is smaller than the three variables above, gender and region have a 186 consistent JS difference regardless of panel waves. And no age variable value consistently has a 187 lower JS score. It shows the model can represent the opinions of different social groups at various 188 levels but offers no clear explanation for the variation in representation. This could come from better 189 recognition of certain groups' views, the training data, or model architecture. Another possibility is 190 that some social groups are more "informative" about this question. 191



Figure 4: JS Distance of six subpopulation groups in Experiment 2. An in-depth presentation of the JS Distance for each group is shown in Figure 14 in the Appendix.

## 192 Llama2 closely reflects sociodemographic patterns, with minor deviations from survey influ-

ences. Cramér's V values in Figure 15 in Appendix show pairwise patterns between prompting variables and text answers. However, in comparison with the survey's values, we see that the model underestimates the influence of age and education degree on the text answers, the model consistently overestimates the effect of region, and the gender and party variables are both overestimated and underestimated. However, except region, the differences are usually less than 0.05, indicating that the Llama2 closely reflects patterns between sociodemographics and the survey. 199

A case study on party variables: Llama2 better models groups with left-leaning parties. To check how much of the JS distance can be associated with the modeling difficulty of the variables, we plot the subpopulation entropy and JS distances for leaning party in Figure 5. Since the population entropy is the same for all subpopulations, we can safely assume that lower conditional entropy means higher information gain for that variable. This allows us to examine how representativeness, linked to the available information in a variable, impacts alignment success. When mutual information is high, LLMs can better model subpopulation behavior. However, certain groups do not fall into this trend line; "die Linke (the Left)" and "die Grünen (the Greens)" are modeled better, and "AfD" is worse (compared to the information their groups carry).



Figure 5: Subpopulation entropy and JS distance for leaning party (mean values for waves 12-21).

- <sup>200</sup> This finding aligns with previous work (e.g., Santurkar et al., 2023; von der Heyde et al., 2025),
- which shows that LLMs tend to have a more left-leaning feature. We show additional results on

other variables and observe similar results showing that LLMs are biased towards left, Western, and educated people in §G (4th paragraph).

#### 204 4.3 Experiment 3: Ablation Experiment

To further show how individual demographic information affects the LLM output diversity, we conducted a series of ablation experiments with the following variations:

- Including only one social group variable.
- Excluding one social group variable.
- Using no social group variables.
- <sup>210</sup> These were compared against the experiment with all variables included.

In the base prompt, the model was only informed that the response was from a German citizen, with

the relevant survey time frame. Detailed prompt variations can be found in Table 7 in Appendix.

213 We used Llama2 and wave 12 data to analyze how adding or removing social group single variables

impacts representativeness and answer diversity. The use of wave 12 is because it took place before

215 COVID (see the dates of the waves in §A), and might have more diverse answer categories (compared

- to the dominance of health policy responses illustrated in Figure 1).
- <sup>217</sup> Variable inclusion and exclusion have an impact on model performance. Figure 6 shows the JS distances in ablation experiments. Including all variables reduces the JS distance by 0.15 compared to the base prompt. Adding a single variable improves predictions. Removing a variable worsens performance, though it is still better than using only one variable, except for the "all except party" case.



Figure 6: JS distances for the ablation experiment.

LLM outputs show stronger correlation with variables when prompted with only one variable. Table 7 compares Cramér's V values between the survey, including only one variable, and all variables included experiments. When only one variable is provided, the generated texts show stronger correlations with the input variable. Although JS distances decrease when more variables are added, this observation suggests that synthetic data patterns are dynamic and can be influenced by the number of prompt variables.

Prompt Variable Survey LLM-one LLM-all 0.07 0.09 0.09 Age Education Degree 0.06 0.25 0.05 Gender 0.08 0.20 0.16 Leaning Party 0.16 0.35 0.17 0.06 0.42 0.15 Region Vocational Degree 0.08 0.120.07

Figure 7: Cramér's V values for the Ablation Experiment

219 The inclusion of the party variable has the most significant impact on model performance,

with its presence leading to substantial improvements in information gain. As shown in Figure 6,

adding either only the party or education variable alone results in the greatest reduction in JS distance

222 compared to the model without demographics; excluding only the party variable leads to a smaller

improvement in JS distance, both highlighting the party variable's greater impact. Similar to Figure 3,

we plot information gain for party, comparing survey entropy to Llama2 with all variables, with only one party variable, and with no variables in Figure 8. As expected, Llama2-base, which includes no subpopulation variables, produces entropies close to the population entropy, with random variations of  $\pm 0.03$ . However, in the 1VAR-party experiment, information gain ranges from 0.2 to 1.3, significantly above random variation. This, along with Cramér's V values, suggests that the

229 model generates typical responses, reducing the variation in subpopulation opinions. Further detailed experimental results are provided in §G of the Appendix.



Figure 8: Information Gain for leaning party variable, comparing survey entropy to Llama2-all (with all variables), 1VAR-party (with only party variable), and Llama2-base (with no variables). Left: population entropy (H(Y)), right: subpopulation entropy  $(H(Y \mid X_i))$ .

230

#### 231 **5 Discussion**

We next distill key findings from our experiments, compare them to prior research, and offer insights into the role of LLMs in modeling demographic behaviors and their practical insights in survey-based applications based on our German case study.

Algorithmic Fidelity in Modeling the German population. von der Heyde et al. (2025) found that GPT-3.5 vote predictions for the 2017 German election are inaccurate and biased towards the Green and Left parties. We also found that the model is better at modeling the opinions of the Green and the Left parties than the right parties. The subpopulation entropy and in-group diversity can partially explain this finding. Other factors could be the models' training data and the RLHF methods used.

Reduction in in-group diversity. Bisbee et al. (2024) found that while GPT-3.5 could replicate survey averages, its synthetic answers lacked variation compared to real survey data. Similarly, von der Heyde et al. (2025) noted GPT-3.5's difficulty in capturing nuanced subpopulation behaviors. In our analysis, we also observe a reduction in in-group diversity under certain conditions, particularly when only one variable is provided to the model or when using the Mixtral model. This suggests that the ability to represent within-group diversity is limited by the model's input structure and specific architecture.

The role of LLM prompts. Binz & Schulz (2023) highlighted how cognitive biases, such as the 247 framing effect, influenced GPT-3's outputs. In our third study, we also noticed that providing only one 248 variable in a prompt caused Llama2 to focus disproportionately on that variable, possibly interpreting 249 it as more critical than when multiple variables were included. Interestingly, this effect varied by 250 model: Mixtral over-relied on variables even with full prompts, while Llama2 showed stronger biases 251 with fewer variables. This suggests that model-specific differences influence how demographic factors 252 are integrated and highlights the need for careful prompt design. Argyle et al. (2023) motivated the 253 silicon sampling approach on this conditional probability formula: 254

$$P(V, B_{\text{LLM}}) = P(V|B_{\text{LLM}})P(B_{\text{LLM}})$$
(6)

here B is demographic backstories, and V is voting patterns. If the model learned the P(V|B), one could correct for the P(B) and obtain:

$$P(V, B_{\text{survey}}) = P(V|B_{\text{survey}})P(B_{\text{survey}})$$
<sup>(7)</sup>

However, in the ablation experiment, we observed that prompting with the social groups is not 257 straightforward, and it does not align LLMs' inner parameters to "solely" consider  $P(B_{\text{social-group}})$ . 258 LLMs might not always be conditioned to sample from the joint distribution of backstories. We 259 propose that demographic variables' order, number, and predictive power have a complex interplay 260 and that this is a further research direction (see, e.g., Shu et al., 2024). Also, insights of vignette 261

experiments from survey methodology (Steiner et al., 2017) could be useful in prompt design. 262

263 Practical evaluation of LLMs. Figure 9 illustrates the JS distances of responses in the last five survey waves compared to earlier waves. As shown, responses from older surveys tend to differ more significantly from those in recent surveys. Consider a hypothetical scenario where LLMs are continuously updated with up-to-date training data while surveys are conducted less frequently due to cost constraints. In such a scenario, Figure 9 suggests that LLMs may also help researchers estimate the answers due to the timeliness of their training data. However, Figure 9: JS distances of answers for even if the model is assumed to be representative enough, the variety in subgroup answers should also be considered for practical uses. In contrast, recent work by Park et al. (2024) demonstrates how LLMs, when applied to large-scale human participant interviews, can simulate subpopulations' attitudes and behaviors, with surveys as a valuable evaluation tool. This highlights the potential of LLMs not only to provide estimates but also to more accurately reflect diverse human perspectives, emphasizing the importance of incorporating subgroup variation.



the last five waves (17-21) from GLES, comparing each survey's answers to those of the preceding surveys. In most cases, differences between survey responses increase over time. The peak of wave 21 compared to wave 20 corresponds to the drop of the health policy category on 29.09.2021 in Figure 1, possibly due to some effect of COVID related topics in that time.

#### **Conclusion and Recommendations** 6 264

This paper evaluates the *algorithmic fidelity* of LLMs to represent the opinions of German sub-265 populations. While von der Heyde et al. (2025) found that GPT-3.5 struggles with the nuances of 266 German subpopulations and the country's multi-party system in closed-ended voting questions, we 267 explore instead free-form open-ended text responses, focusing on how these responses align with 268 survey data. By using free-form text responses rather than multiple-choice questions, we can identify 269 detailed issues in contextual information and the variety of different subpopulations, underscoring 270 the value of this evaluation approach. Our findings show that LLMs, particularly the Llama2 model, 271 are capable of associating text responses with social-demographic variables, indicating a degree of 272 representativeness. However, the number of variables included in the prompt plays a crucial role 273 in model performance. Despite this, the models still tend to generate stereotypical representations, 274 with a noticeable favor towards left-leaning parties, consistent with previous findings on the limited 275 diversity of opinions reflected in LLMs. 276

Based on these findings, we recommend that both LLM and social science researchers consider the 277 following steps for future evaluation of LLM-generated responses in survey-based research: 278

Improved representation of opinion diversity: LLMs should be further developed to reflect the full 279 spectrum of opinions within subpopulations while harmful contents are cautiously manipulated. This 280 includes addressing biases and avoiding the oversimplification of diverse views into stereotypical 281 categories with certain safety mechanisms. 282

Cross-national comparison: Due to the current discussion of English-centric biases of LLMs, a 283 more inclusive evaluation with opinion diversity from non-English data or cross-national sources 284 285 such as GlobalOpinionQA (Durmus et al., 2024) should be conducted and improved.

**Timeliness and survey simulation:** LLMs can be valuable in situations where real survey data is 286 limited or outdated (Namikoshi et al., 2024; Ma et al., 2024). With continuously updated training 287 data, LLMs could be further evaluated in the case of estimating shifts in public opinion. 288

## **References**

Utkarsh Agarwal, Kumar Tanmay, Aditi Khandelwal, and Monojit Choudhury. Ethical reasoning
and moral value alignment of LLMs depend on the language we prompt them in. In Nicoletta
Calzolari, Min-Yen Kan, Veronique Hoste, Alessandro Lenci, Sakriani Sakti, and Nianwen Xue
(eds.), Proceedings of the 2024 Joint International Conference on Computational Linguistics,
Language Resources and Evaluation (LREC-COLING 2024), pp. 6330–6340, Torino, Italia, May
2024. ELRA and ICCL. URL https://aclanthology.org/2024.lrec-main.560.

Hans-Jürgen Andreß, Katrin Golsch, and Alexander W. Schmidt. *Applied Panel Data Analysis for Economic and Social Surveys*. Springer, 2013.

Lisa P. Argyle, Ethan C. Busby, Nancy Fulda, Joshua R. Gubler, Christopher Rytting, and David
 Wingate. Out of one, many: Using language models to simulate human samples. *Political Analysis*, 31(3):337–351, Jul 2023. doi: https://doi.org/10.1017/pan.2023.2. URL https://
 www.cambridge.org/core/journals/political-analysis/article/abs/

302 out-of-one-many-using-language-models-to-simulate-human-samples/

**303** 035D7C8A55B237942FB6DBAD7CAA4E49.

Arnav Arora, Lucie-aimée Kaffee, and Isabelle Augenstein. Probing pre-trained language mod els for cross-cultural differences in values. In Sunipa Dev, Vinodkumar Prabhakaran, David
 Adelani, Dirk Hovy, and Luciana Benotti (eds.), *Proceedings of the First Workshop on Cross- Cultural Considerations in NLP (C3NLP)*, pp. 114–130, Dubrovnik, Croatia, May 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.c3nlp-1.12. URL https:
 //aclanthology.org/2023.c3nlp-1.12.

Marcel Binz and Eric Schulz. Using cognitive psychology to understand gpt-3. *Proceedings of the National Academy of Sciences*, 120(6):e2218523120, 2023. doi: 10.1073/pnas.2218523120. URL
 https://www.pnas.org/doi/abs/10.1073/pnas.2218523120.

James Bisbee, Joshua D. Clinton, Cassy Dorff, Brenton Kenkel, and Jennifer M. Larson. Synthetic replacements for human survey data? the perils of large language models. *Political Analysis*, 32 (4):401–416, 2024. doi: 10.1017/pan.2024.5.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,
Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel
Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler,
Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott
Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya
Sutskever, and Dario Amodei. Language models are few-shot learners, 2020.

Yong Cao, Li Zhou, Seolhwa Lee, Laura Cabello, Min Chen, and Daniel Hershcovich. Assessing
 cross-cultural alignment between ChatGPT and human societies: An empirical study. In Sunipa
 Dev, Vinodkumar Prabhakaran, David Adelani, Dirk Hovy, and Luciana Benotti (eds.), *Proceedings* of the First Workshop on Cross-Cultural Considerations in NLP (C3NLP), pp. 53–67, Dubrovnik,
 Croatia, May 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.c3nlp-1.7.
 URL https://aclanthology.org/2023.c3nlp-1.7.

Jacob Cohen. A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20(1):37–46, 1960. doi: 10.1177/001316446002000104. URL https://doi.org/10.1177/001316446002000104.

Harald Cramér. *Mathematical Methods of Statistics (PMS-9)*. Princeton University Press, 1999.
 ISBN 9780691005478. URL http://www.jstor.org/stable/j.ctt1bpm9r4.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of
 deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and
 Thamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational
 Linguistics. doi: 10.18653/v1/N19-1423. URL https://aclanthology.org/N19-1423.

Ricardo Dominguez-Olmedo, Moritz Hardt, and Celestine Mendler-Dünner. Questioning 339 340 the survey responses of large language models. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), Advances in Neural In-341 formation Processing Systems, volume 37, pp. 45850-45878. Curran Associates, Inc., 342 URL https://proceedings.neurips.cc/paper\_files/paper/2024/ 2024.343

file/515c62809e0a29729d7eec26e2916fc0-Paper-Conference.pdf. 344

Esin Durmus, Karina Nguyen, Thomas Liao, Nicholas Schiefer, Amanda Askell, Anton Bakhtin, 345 Carol Chen, Zac Hatfield-Dodds, Danny Hernandez, Nicholas Joseph, Liane Lovitt, Sam McCan-346 dlish, Orowa Sikder, Alex Tamkin, Janel Thamkul, Jared Kaplan, Jack Clark, and Deep Ganguli. 347 Towards measuring the representation of subjective global opinions in language models. In First 348 Conference on Language Modeling, 2024. URL https://openreview.net/forum?id= 349 zl16jLb91v. 350

Peter Eigenschink, Thomas Reutterer, Stefan Vamosi, Ralf Vamosi, Chang Sun, and Klaudius Kalcher. 351 Deep Generative Models for Synthetic Data: A Survey. IEEE Access, 11:47304–47320, 2023. 352 ISSN 2169-3536. doi: 10.1109/ACCESS.2023.3275134. URL https://ieeexplore.ieee. 353 org/document/10122524/. 354

- GESIS. Gles panel 2016-2021, wellen 1-21. GESIS, Köln. ZA6838 Datenfile Version 6.0.0, 355 https://doi.org/10.4232/1.14114, 2023. 356
- GESIS. Gles codierung des wichtigsten politischen problems (2018-2022). GESIS, Köln. ZA7957 357 Datenfile Version 2.0.0, https://doi.org/10.4232/1.14241, 2024. 358
- Thomas M. Holbrook. An Introduction to Political and Social Data Analysis Using R. Thomas M. 359 Holbrook, May 2022. URL https://bookdown.org/tomholbrook12/IPSDAR/. 360

EunJeong Hwang, Bodhisattwa Majumder, and Niket Tandon. Aligning language models to user 361 opinions. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), Findings of the Association for 362 Computational Linguistics: EMNLP 2023, pp. 5906–5919, Singapore, December 2023. Association 363 for Computational Linguistics. doi: 10.18653/v1/2023.findings-emnlp.393. URL https:// 364 aclanthology.org/2023.findings-emnlp.393. 365

Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris 366 Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, 367 Gianna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier, Marie-368 Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le 369 Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 370 371 Mixtral of experts, 2024a.

- Hang Jiang, Xiajie Zhang, Xubo Cao, Cynthia Breazeal, Deb Roy, and Jad Kabbara. PersonaLLM: 372 Investigating the ability of large language models to express personality traits. In Kevin Duh, Helena 373 Gomez, and Steven Bethard (eds.), Findings of the Association for Computational Linguistics: 374 NAACL 2024, pp. 3605–3627, Mexico City, Mexico, June 2024b. Association for Computational 375 Linguistics. URL https://aclanthology.org/2024.findings-naacl.229. 376
- Dan Jurafsky and James H. Martin. Speech and language processing: an introduction to natural 377 language processing, computational linguistics, and speech recognition, 2024. URL https: 378 //web.stanford.edu/~jurafsky/slp3/. 379
- S. Kullback and R. A. Leibler. On Information and Sufficiency. The Annals of Mathematical Statistics, 380 22(1):79-86, 1951. doi: 10.1214/aoms/1177729694. URL https://doi.org/10.1214/ 381 aoms/1177729694. 382
- Jianhua Lin. Divergence measures based on the shannon entropy. *IEEE Trans. Inf. Theor.*, 37(1): 383 145-151, sep 2006. ISSN 0018-9448. doi: 10.1109/18.61115. URL https://doi.org/10. 384 1109/18.61115. 385
- Lin Long, Rui Wang, Ruixuan Xiao, Junbo Zhao, Xiao Ding, Gang Chen, and Haobo Wang. On 386 LLMs-driven synthetic data generation, curation, and evaluation: A survey. In Lun-Wei Ku, Andre 387

2024, pp. 11065–11082, Bangkok, Thailand and virtual meeting, August 2024. Association for
 Computational Linguistics. URL https://aclanthology.org/2024.findings-acl.
 658.

Bolei Ma, Xinpeng Wang, Tiancheng Hu, Anna-Carolina Haensch, Michael A. Hedderich, Barbara
Plank, and Frauke Kreuter. The potential and challenges of evaluating attitudes, opinions, and
values in large language models. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen
(eds.), *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 8783–
8805, Miami, Florida, USA, November 2024. Association for Computational Linguistics. URL
https://aclanthology.org/2024.findings-emnlp.513.

Jonathan Mellon, Jack Bailey, Ralph Scott, James Breckwoldt, Marta Miori, and Phillip Schmedeman. Do ais know what the most important issue is? using language models to code open-text social survey responses at scale. *Research & Politics*, 11(1):20531680241231468, 2024. doi: 10.1177/

401 20531680241231468. URL https://doi.org/10.1177/20531680241231468.

Keiichi Namikoshi, Alex Filipowicz, David A. Shamma, Rumen Iliev, Candice L. Hogan, and Nikos
 Arechiga. Using llms to model the beliefs and preferences of targeted populations, 2024.

Joon Sung Park, Carolyn Q. Zou, Aaron Shaw, Benjamin Mako Hill, Carrie Cai, Meredith Ringel
 Morris, Robb Willer, Percy Liang, and Michael S. Bernstein. Generative agent simulations of
 1,000 people, 2024. URL https://arxiv.org/abs/2411.10109.

Philip Resnik, Bolei Ma, Alexander Hoyle, Pranav Goel, Rupak Sarkar, Maeve Gearing, Anna Carolina Haensch, and Frauke Kreuter. Topic-oriented protocol for content analysis of text – a
 preliminary study. Unpublished manuscript, June 2024.

Guilherme Sanches de Oliveira and Edward Baggs. *Psychology's WEIRD Problems*. Elements in
 Psychology and Culture. Cambridge University Press, 2023.

Shibani Santurkar, Esin Durmus, Faisal Ladhak, Cinoo Lee, Percy Liang, and Tatsunori Hashimoto.
 Whose opinions do language models reflect? In *Proceedings of the 40th International Confer- ence on Machine Learning*, ICML'23. JMLR.org, 2023. URL https://proceedings.mlr.

415 press/v202/santurkar23a.html.

Bangzhao Shu, Lechen Zhang, Minje Choi, Lavinia Dunagan, Lajanugen Logeswaran, Moontae 416 Lee, Dallas Card, and David Jurgens. You don't need a personality test to know these models 417 are unreliable: Assessing the reliability of large language models on psychometric instruments. 418 In Kevin Duh, Helena Gomez, and Steven Bethard (eds.), Proceedings of the 2024 Confer-419 ence of the North American Chapter of the Association for Computational Linguistics: Human 420 Language Technologies (Volume 1: Long Papers), pp. 5263-5281, Mexico City, Mexico, June 421 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.naacl-long.295. URL 422 https://aclanthology.org/2024.naacl-long.295. 423

Peter M. Steiner, Christiane Atzmüller, and Dan Su. Designing valid and reliable vignette experiments for survey research: A case study on the fair gender income gap. *Journal of Methods and Measurement in the Social Sciences*, 7:52–94, 6 2017. ISSN 2159-7855. doi: 10.2458/v7i2.20321. URL http://journals.librarypublishing.arizona.edu/
jmmss/article/id/833/.

Seungjong Sun, Eungu Lee, Dongyan Nan, Xiangying Zhao, Wonbyung Lee, Bernard J. Jansen, and
 Jang Hyun Kim. Random silicon sampling: Simulating human sub-population opinion using a
 large language model based on group-level demographic information, 2024.

432 Gemma Team. Gemma: Open models based on gemini research and technology, 2024.

Lindia Tjuatja, Valerie Chen, Tongshuang Wu, Ameet Talwalkwar, and Graham Neubig. Do LLMs exhibit human-like response biases? a case study in survey design. *Transactions of the Association* for Computational Linguistics, 12:1011–1026, 2024. doi: 10.1162/tacl\_a\_00685. URL https:

436 //aclanthology.org/2024.tacl-1.56.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand
 Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language

440 models, 2023. URL https://arxiv.org/abs/2302.13971.

Veniamin Veselovsky, Manoel Horta Ribeiro, Akhil Arora, Martin Josifoski, Ashton Anderson, and
 Robert West. Generating faithful synthetic data with large language models: A case study in
 computational social science, 2023. URL https://arxiv.org/abs/2305.15041.

Leah von der Heyde, Anna-Carolina Haensch, and Alexander Wenz. Vox populi, vox ai? using
large language models to estimate german vote choice. *Social Science Computer Review*, 0(0):
08944393251337014, 2025. doi: 10.1177/08944393251337014. URL https://doi.org/10.
1177/08944393251337014.

 Xinpeng Wang, Chengzhi Hu, Bolei Ma, Paul Rottger, and Barbara Plank. Look at the text: Instruction-tuned language models are more robust multiple choice selectors than you think. In
 *First Conference on Language Modeling*, 2024a. URL https://openreview.net/forum?
 id=qHdSA85GyZ.

Xinpeng Wang, Bolei Ma, Chengzhi Hu, Leon Weber-Genzel, Paul Röttger, Frauke Kreuter, Dirk
Hovy, and Barbara Plank. "my answer is C": First-token probabilities do not match text answers in
instruction-tuned language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of the Association for Computational Linguistics ACL 2024*, pp. 7407–7416, Bangkok,
Thailand and virtual meeting, August 2024b. Association for Computational Linguistics. URL
https://aclanthology.org/2024.findings-acl.441.

Xuena Wang, Xueting Li, Zi Yin, Yue Wu, and Jia Liu. Emotional intelligence of large language
models. *Journal of Pacific Rim Psychology*, 17:1–12, 2023. doi: 10.1177/18344909231213958.
URL https://doi.org/10.1177/18344909231213958.

## 461 A Data and Prompt Template

Table 3 shows the dates for the waves involved in the original GLES survey (GESIS, 2023). The six

463 main social demographic variables and their subgroups in the original survey are presented in Table 4.

| Wave | Start Date | End Date   |
|------|------------|------------|
| 10   | 06-11-2018 | 21-11-2018 |
| 11   | 28-05-2019 | 12-06-2019 |
| 12   | 05-11-2019 | 19-11-2019 |
| 13   | 21-04-2020 | 05-05-2020 |
| 14   | 03-11-2020 | 17-11-2020 |
| 15   | 25-02-2021 | 12-03-2021 |
| 16   | 06-05-2021 | 19-07-2021 |
| 17   | 07-07-2021 | 20-07-2021 |
| 18   | 11-08-2021 | 24-08-2021 |
| 19   | 15-09-2021 | 25-09-2021 |
| 20   | 29-09-2021 | 12-10-2021 |
| 21   | 09-12-2021 | 21-12-2021 |

Table 3: Data collection dates of GLES waves

In Figure 10, we show the original prompt in German we used for the LLM response generation. We use the template in German and expect the LLMs to respond in German, in order to mirror the real survey case. The placeholders in the prompt are replaced with the respondents' information and then fed to the LLMs that were experimented with. The prompt templates used in the ablation study in §4.3 are presented in Table 7.

| Social Groups     | Sub-Groups   |
|-------------------|--|
| Age               | 18-29<br>30-44<br>45-59<br>60+   |
| Gender            | Male<br>Female   |
| Leaning Party     | AfD<br>CDU/CSU<br>FDP<br>Grünen<br>A minor party<br>Linke<br>SPD<br>No party   |
| Region            | East Germany<br>West Germany   |
| Education Degree  | High school diploma<br>Higher education entrance qualification<br>Secondary school diploma<br>Intermediate school diploma<br>Is still student<br>No school diploma   |
| Vocational Degree | Completed vocational internship/volunteer work<br>Vocational school diploma<br>University of applied sciences degree<br>Specialist school diploma<br>Completed apprenticeship<br>Master craftsman or technician qualification<br>University degree<br>In vocational training<br>Commercial or agricultural apprenticeship<br>Commercial apprenticeship<br>No vocational training completed |

Table 4: Social-Demographic Groups and Sub-Groups.

{artikel} Befragte ist {age} Jahre alt und {gender}. {pronoun} {schulabschluss\_clause } und {berufabschluss\_clause} {pronoun2} lebt in {region} und unterstützt hauptsächlich {party}.

Figure 10: Original prompt template in German.

Identifizieren Sie das wichtigste Problem, mit dem Deutschland im {month} {year} konfrontiert ist. Geben Sie die Antwort in einem prägnanten Satz an, konzentrieren Sie sich nur auf ein einziges Thema ohne weitere Ausführungen oder Auflistung zusätzlicher Probleme. Wiederholen Sie nicht die Informationen die Ihnen gegeben wurden, und geben Sie Ihre Antwort direkt und ohne einleitende Phrasen. Antworten Sie auf Deutsch und ausschließlich auf Deutsch, verwenden Sie keine Englische Sprache. Antworten Sie aus der Sicht eines Befragten mit deutscher Staatsbürgerschaft und den im nachfolgenden spezifizierten Eigenschaften.

## **B** Classes of the Coding Scheme

In Table 8, we show the original classes of the coding scheme based on GESIS (2024) as well as the merged and reduced 16 classes.

## 472 C Annotation Scheme

The annotators followed the coding guidelines provided by (GESIS, 2024) for consistency. We used the coarse classes for annotation to achieve better agreement among annotators. One of the leading authors and another project collaborator were asked to conduct the same annotations of the 1,500 samples as volunteers. There exists disagreement on minor cases that were resolved after discussion. Both annotators are consent about the annotated data use. Figure 11 and 12 show the annotation screen and the annotation instruction given to the annotators respectively.

- 3 --data\_path /content/drive/MyDrive/labeling/base/chunk\_0.json\
- 4 -- TABLE\_MAX\_WIDTH 50 \
- 5 --TEXT\_MAX\_WIDTH 100

...

| ++        | +   | ++  |
|-----------|---|---|
| classid   | ClassName   | subclasses  |
| 0         | Politisches System und Prozesse (Political System<br>and Processes)         | Election Campaign and Government Formation,<br>Political Structures and Processes, Democracy,<br>Bureaucracy, Lobbyism, Corruption, Values,<br>political culture            |
|           | Sozialpolitik (Social Policy)   | Social Policy, Social Justice, Poverty,<br>Unemployment and Basic(Social) Security, Pensions<br>and Demographic Change  |
| 2         | Gesundheitspolitik (Health Policy)  | Health Policy, Nursing, Corona Pandemic   |
| 3         | Familien- und Gleichstellungspolitik (Family and<br>Gender Equality Policy) | Family Policy, Gender Equality  |
| 4         | Bildungspolitik (Education Policy)  | Education Policy, School Policy   |
| 5         | Umweltpolitik (Environmental Policy)  | Energy Policy, Environmental Policy, Climate<br>Policy, Natural Disasters   |
| 6         | Wirtschaftspolitik (Economic Policy)  | Economic Policy, Price Level, Infrastructure,<br>Digital Infrastructure, Transport Policy, Housing<br>Policy  |
| 7         | Sicherheits (Security)  | Terrorism, Internal Security, Crime and Violence,<br>Radicalization and Extremism, Law and Justice,<br>Defense  |
| 8         | Außenpolitik (Foreign Policy)   | Foreign Policy, Europe and European Union,<br>(Germany's Relations with) Russia, (Germany's<br>Relations with) Turkey, International Conflicts<br>and Peace, War in Ukraine |
| 9         | Medien und Kommunikation (Media and Communication)                          | Media   |
| 10        | Sonstiges (Others)  | Others  |
| 11        | Migration and Integration   | Migration and Integration   |
| 12        | Ostdeutschland  | Ostdeutschland  |
| 13        | keine Angabe (not specified)  | keine Angabe  |
| 14        | weiß nicht  |   |
| 15        | LLM refusal   | As an AI model,, It is not possible to answer   |
| 16        | Values, political culture and general social criticism                      | Nationalism, Rise of left, Social cohesion,social  <br>  division   |
| number:14 |   |   |

output: Das wichtigste Problem, mit dem Deutschland im November 2019 konfrontiert ist, ist die unzureichende Ausstattung der öffentlichen Pflanzen mit ausreichend Wasser.

translation: The most important problem facing Germany in November 2019 is the inadequate supply of public plants with sufficient water.

label:

Figure 11: Annotation Screen.

<sup>2 !</sup>python /content/drive/MyDrive/labeling/annotation.py \

#### Labeling Rules

- Samples can be labeled with multiple classes, (up to 3) if needed.
- example annotations:

• 1,10

- 14
- 12,2,1
- If you really can't decide on a sample, you can label it with -1, and continue labeling, we can decide on these cases later on.
- The following guide is from the original study and it can help you while giving your decision:
  - 1. When coding the answers, no contextual knowledge should be used if possible. In addition, no interpretation should be necessary to assign an answer to a category.
  - 2. If several topics are mentioned (although the question explicitly asks for one topic), several codes are assigned in the order in which they are mentioned. A maximum of 3 codes (in order of mention) are assigned.
  - 3. if general topics are supported with additional examples (enumerations, reasons), **only the general topic**, and not the example, is coded. Examples:
  - 'Citizens are becoming more and more criminal because of the euro crisis, the refugees, the new president of the USA and Brexit, but our politicians only listen to the citizens with half an ear!' -[17 crime and violence, criticism of politicians], not Euro crisis or Refugees.
  - 'Health policy and care'
    - [codes: 37 health policy, 38 care]
  - 'Healthcare (e.g. care)' [code: 37 health policy]

The original coding scheme, detailed examples and explanations can be found here if needed:

Figure 12: The instructions at the annotating tool.

## 479 **D** Technical Setup

480 We used Python 3.12.1 and the transformers <sup>2</sup> library (version 4.42.4) by HuggingFace (with

481 Pytorch Framework as the backend) to create two custom classes: the BertClassifier (for the

multilabel classification task) and TextGenerator (for generating synthetic answers ).

**Text Generation.** To fit models into a single GPU, we have used the 8-bit quantized version of the models. The inference configurations can be found in the study repository. We did the inference with batch\_size of 16 to benefit the parallel computing power and reduce runtime. On average, TextGenerator generated 1.16 answers per second. We performed 25 generation experiments, using ca. 75 GPU hours for the generation task.

**Text Classification.** As the contexts are all in German, we used the German version of the BERT model<sup>3</sup>. The BertClassifier training takes around 20 minutes for the setup 5. We trained the with a batch\_size of 32, a learning\_rate of 2e-5, and a fixed\_precision at 16 bits to fasten the convergence. The early stop condition stopped the training after 3rd epoch when no further loss reduction was observed. Table 5 shows the other relevant model parameters and hyperparameters.

| Parameter     | Value |
|---------------|-------|
| epochs        | 15    |
| learning_rate | 2e-5  |
| batch_size    | 32    |
| weight_decay  | 0.01  |
| fp16          | True  |
| max_length    | 512   |

<sup>&</sup>lt;sup>2</sup>https://huggingface.co/docs/transformers/

<sup>&</sup>lt;sup>3</sup>https://huggingface.co/google-bert/bert-base-german-cased

## 493 E Qualitative Analysis

Table 9 shows a few sample responses from LLM experiments. Before discussing models' representativeness, we point to three qualitative issues observed in the text generations:

**Introductory Phrases.** We noted that all models use some "introductory phrases" even though 496 models were prompted not to use any. The Gemma model starts the sentence 96 % of the time by 497 listing the social group variables given to it. Llama2 model uses the "The most important problem 498 facing Germany" phrase in 96 % of its generations. The Mixtral model uses the "One of the most 499 important issues" phrase in 75 % of the answers. Even if the text lengths had been comparable, these 500 style characteristics would allow humans to discriminate synthetic responses from actual ones. At this 501 point, we did not put further effort into trying different prompts and making the synthetic responses 502 503 stylistically similar to survey responses.

Mention of Future Events. Despite specifying the survey month and year (November 2019),
 responses occasionally referenced events that occurred after the survey date, such as the COVID pandemic (2020) and the Energy Crisis in Germany (2021). For example, the Gemma model
 contained COVID-19-related words (COVID, corona, coronavirus, COVID-19, sars-cov, etc.) in 42
 % of its responses. This problem was observed relatively less in Llama2 and Mixtral models (3 %
 and 0.2 % of answers, respectively).

Mixed-Language Answers. Although models were instructed to respond in German, a small
 percentage (up to 3%) of answers had some parts in German, whereas some terms or clauses were in
 English.

## 513 F Additional Metrics

<sup>514</sup> In this section, we present metrics in addition to those in §3.4. These include the base metrics for the <sup>515</sup> main experimentation, as well as additional metrics used for the additional results in §G.

**Proportion Agreement.** It is the proportion of two variables exactly matching. Like accuracy, this measure does not consider the probability of matching by chance and should be used as a descriptive quantity (Argyle et al., 2023).

**Cohen's Kappa** ( $\kappa$ ). It is a measure of agreement between two categorical variables and is often used as a measure of inter-rater agreement (Cohen, 1960). Unlike proportional agreement, it corrects for the agreement by chance, and It is defined as  $\kappa = \frac{p_0 - p_e}{1 - p_e}$  where

- $p_0$  is the observed agreement ratio
- $p_e$  is the expected agreement when annotators assign labels randomly.

**Kullback-Leibler Divergence.** It also known as relative entropy, is a method used in measuring the statistical distance between two probability distributions (Kullback & Leibler, 1951). For distributions P and Q of a discrete random variable  $X = [X_1, ..., X_n]$ , the Kullback-Leibler (KL) divergence can be defined as:

$$D_{\mathrm{KL}}(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \log\left(\frac{P(x)}{Q(x)}\right)$$
(8)

KL Divergence is not a distance measure since it does not satisfy the symmetry requirement of a metric. i.e KL(P,Q) = KL(Q,P), unless P and Q are equal. If not, KL divergence is always greater than 0 and not bounded.

Absolute Percentage Error (APE). JS distance enables us to compare model performance at the question level. We employ the APE to evaluate the accuracy of predictions in each category. APE is calculated by determining the absolute differences between predicted and actual frequencies and then normalizing these differences by the reference survey frequencies. For each label, we compute:

$$APE_L = \left| \frac{y_t - \hat{y}_t}{y_t} \right| \tag{9}$$

where  $y_t$  is survey frequency and  $\hat{y}_t$  LLM output frequency for the label *l*.

## 536 G Additional Results

<sup>537</sup> In this section, we show additional results and figures for the main experimentation in §4.

Label Distribution on LLM Outputs. Figure 13 shows the label distribution on the three LLM outputs based on the coarse labels.

540 Detailed JS Distances of Subpopulation in Experiment 1: Llama2 achieves better performance in most categories. We show the detailed JS Distances in each social group category in Experiment 1 in Table 13 for all three experimented LLMs. Among the three LLMs, we notice Llama2 has the least JS Distances across the most subcategories compared to the other two LLMs, showing more alignment with the real survey data.

Additional JS Distances of Subpopulation Variables in Experiment 2: Further indication of 545 WEIRD bias of LLMs. We show additional results of the survey's subpopulation entropy and 546 the JS Distance between the Llama output and survey results across 10 waves in experiment 2 in 547 Figure 16, i.e., results for four other variables in addition to the variable leaning\_party in Figure 548 5 from §4.2. For the vocational degree variable, groups with a completed vocational internship are 549 the least well-represented on average. For the education degree variable, groups with no degree or 550 only a secondary school diploma are less represented compared to those with higher educational 551 qualifications. Among age groups, older cohorts are less well-represented than younger ones. For the 552 553 regional variable, which includes only two groups, there is a greater discrepancy in representing East Germany compared to West Germany. Overall, these findings align with prior evidence that LLMs 554 exhibit biases favoring Western, younger, and more educated subpopulations, commonly referred to 555 as WEIRD bias<sup>4</sup>, as highlighted in studies such as Santurkar et al. (2023), Cao et al. (2023), Arora 556 et al. (2023), and Agarwal et al. (2024). 557

Label-Level Percentage Errors in Experiment 3. Table 6 shows label-level percentage errors. We compared the best-performing models in the 1-var-\* and w/o.-\* experiments. In 7 out of 14 labels, using all variables produced the lowest percentage errors. Both JS distances and percentage errors decreased with the inclusion of more variables, reinforcing the representational accuracy of the Llama2 model for the German population.

**Proportional Agreement (PA) and \kappa Scores across Waves: LLMs face in achieving consistent** 563 agreement with survey data over time, particularly in representing complex social phenomena. 564 The scores in Table 10 compare Llama2-generated responses with resampled survey data across 565 different waves. The PA for Llama2 fluctuates significantly, reaching as high as 0.56 in wave 13 but 566 dropping in later waves, especially after wave 17, indicating inconsistencies in the model's ability 567 to align with the original survey data. By contrast, the survey resample maintains relatively stable 568 PA scores, ranging between 0.34 and 0.55, indicating better agreement with the original survey. 569 The  $\kappa$  scores, however, are low for both Llama2 and the survey resample, with Llama2 performing 570 particularly poorly (0.01–0.04). This suggests that while the model can capture some high-level 571 agreement (as seen in PA scores), it struggles to replicate the nuanced variability and structure of 572 human responses across waves, especially as the diversity of responses increases. These findings 573 underscore the challenges LLMs face in representing complex social phenomena in survey data, 574 especially over time. 575

<sup>&</sup>lt;sup>4</sup>The term "WEIRD bias" originates in psychology, where values from Western, Educated, Industrialized, Rich, and Democratic (W.E.I.R.D.) societies are assumed to represent universal "human" values (Sanches de Oliveira & Baggs, 2023).

| Labels                         | LLM-base | min(1var. *) | min(w/o. *) | LLM-all |
|--------------------------------|----------|--------------|-------------|---------|
| Political System, Processes    | 7.92     | -3.62        | 0.92        | 9.71    |
| Social Policy                  | -23.76   | -21.94       | -15.16      | -7.45   |
| Health Policy                  | 22.02    | 3.32         | 2.84        | 2.08    |
| Family and Gender Equality     | -0.27    | -0.25        | -0.14       | -0.10   |
| Education Policy               | -2.20    | -2.15        | -1.88       | -1.83   |
| Environmental Policy           | 4.43     | 0.42         | 0.15        | 2.52    |
| Economic Policy                | 8.15     | 7.92         | 11.49       | 9.21    |
| Security                       | -6.26    | -7.37        | -7.06       | -6.88   |
| Foreign Policy                 | 6.58     | 1.93         | 1.59        | 0.54    |
| Media and Communication        | 0.09     | -0.01        | 0.00        | 0.03    |
| Others                         | -3.11    | -3.12        | -3.12       | -3.00   |
| Migration and Integration      | -17.16   | -21.40       | -19.80      | -14.61  |
| East Germany                   | 0.01     | -0.18        | -0.17       | 0.14    |
| Values, Polit. Culture, Social | 3.56     | 1.31         | 9.90        | 9.65    |

Table 6: Percentage errors in the ablation experiment. LLM-base denotes no demographics; min(1var. \*) denotes best-performing experiment with one variable; min(w/o. \*) denotes best-performing experiment with all except one variable; LLM-all denotes all variables.

Label-level Breakdown with APE: LLM predictions show over- and underrepresentation of 576 certain political topics compared to survey data. Table 11 and 12 show the comparison of 577 predicted label percentages between survey and LLM answers in experiment 1 and 2 respectively. In 578 Table 11, we observe that "Security", "Migration and Integration" and "Social Policy" topics are less 579 represented than the survey in all LLM-texts and "Values, political culture and general social criticism" 580 represented much more, with a mean APE=218.1. the Mixtral model emphasized "Environment 581 Politics", whereas Llama2 focused on "Political Systems and Processes" more than others. Table 12 582 shows "Health Policy", "Values, political culture and general social criticism" and "Economic Policy" 583 are consistently more represented as the text answer categories, whereas "Migration and Integration" 584 and "Security" are less represented. We also calculated the mean APE per label to compare errors on 585 average on which categories the Llama2 represented the political topics more accurately. 586

Factual Knowledge of Llama2 without Demographic Prompts. We also check actual knowledge of Llama2 model without giving any survey contexts in Table 14. It shows Llama2's general accuracy in providing vote percentages for major elections, though minor errors and formatting issues occur, especially in the 2020 U.S. election. For Germany's most important problem in 2019, Llama2 identifies key issues like climate change and economic security but gives rough percentage estimates rather than precise data. These results suggest that while the model captures broad trends, it struggles with exact figures when not given specific demographic or contextual prompts.

| Experiment Name           | Prompt Information  |
|---------------------------|---|
| 1_var_region              | Der/Die Befragte lebt in {eastwest}. [/INST]  |
| 1_var_party               | Der/Die Befragte unterstützt hauptsächlich {party}. [/INST]   |
| 1_var_education_degree    | Der/Die Befragte {schulabschluss_clause} [/INST]  |
| 1_var_age                 | Der/Die Befragte ist {age} Jahre alt. [/INST]   |
| 1_var_gender              | {artikel} Befragte ist {gender} [/INST]   |
| 1_var_vocational_degree   | Der/Die Befragte {berufabschluss_clause} [/INST]  |
| without_age               | {artikel} Befragte ist {gender}. {pronoun} {schulabschluss_clause} und<br>{berufabschluss_clause} {pronoun2} lebt in {eastwest} und unterstützt<br>hauptsächlich {party}. [/INST] |
| without_region            | {artikel} Befragte ist {age} Jahre alt und {gender}. {pronoun}  |
| C                         | {schulabschluss_clause} und {berufabschluss_clause} {pronoun2} unterstützt  |
| without vocational degree | naupisachlich {party}. [/INSI]  |
| without_vocational_degree | (achulabsahluss alausa) (pronoun?) laht in (asstyrast) und unterstützt  |
|                           | hauptsächlich {party}. [/INST]  |
| without_education_degree  | {artikel} Befragte ist {age} Jahre alt und {gender}. {pronoun} und  |
| C                         | {berufabschluss_clause} {pronoun2} lebt in {eastwest} und unterstützt   |
|                           | hauptsächlich {party}. [/INST]  |
| without_party             | {artikel} Befragte ist {age} Jahre alt und {gender}. {pronoun}  |
|                           | {schulabschluss_clause} und {berufabschluss_clause} {pronoun2} lebt in  |
|                           | {eastwest}. [/INST]   |
| without_gender            | Der/Die Befragte ist {age} Jahre alt. Er/Sie {schulabschluss_clause} und  |
|                           | {berufabschluss_clause} Er/Sie lebt in {eastwest} und unterstützt   |
|                           | hauptsächlich {party}. [/INST]  |

Table 7: Ablation experiments and the modified prompt contents. 1\_var\_\* denotes the experimentation of prompting with only one variation\*. without\_\* denotes the experimentation of prompting with all variables except \*.



Figure 13: Label distributions of three experimented LLMs

| Fine Labels   | Coarse Labels  |
|---|--|
| Election Campaign and Government Formation, Politi-<br>cal Structures and Processes, Democracy, Bureaucracy,<br>Lobbyism, Corruption  | Political System and Processes                         |
| Values, political culture and general social criticism  | Values, political culture and general social criticism |
| Social Policy, Social Justice, Poverty, Unemployment<br>and Basic Security, Pensions and Demographic Change   | Social Policy  |
| Health Policy, Nursing, Corona Pandemic   | Health Policy  |
| Family Policy, Gender Equality  | Family and Gender Equality Policy                      |
| Education Policy, School Policy   | Education Policy                                       |
| Energy Policy, Environmental Policy, Climate Policy, Natural Disasters  | Environmental Policy                                   |
| Economic Policy, Price Level, Infrastructure, Digital<br>Infrastructure, Transport Policy, Housing Policy   | Economic Policy  |
| Terrorism, Internal Security, Crime and Violence, Radi-<br>calization and Extremism, Law and Justice, Defense   | Security   |
| Foreign Policy, Europe and European Union, (Ger-<br>many's Relations with) Russia, (Germany's Relations<br>with) Turkey, International Conflicts and Peace, War in<br>Ukraine | Foreign Policy   |
| Media   | Media and Communication                                |
| Others  | Others   |
| Migration and Integration   | Migration and Integration                              |
| East Germany  | East Germany   |
| Not specified   | Not specified  |
| Don't know  | Don't know   |

Table 8: Description of the Classes, originally in German, translated to English

| Model   | Original Answer   | Translation   |
|---------|---|---|
| Gemma   | Als 61-jährige Frau mit deutscher<br>Staatsbürgerschaft, die in Westdeutschland lebt<br>und hauptsächlich keine Partei unterstützt, ist das<br>wichtigste Problem, mit dem Deutschland im<br>November 2019 konfrontiert ist, die weiterhin<br>bestehende Situation der COVID-19-Pandemie. | As a 61-year-old woman with German<br>citizenship, living in West Germany<br>and mainly not supporting any political<br>party, the most important issue facing<br>Germany in November 2019 is the<br>ongoing situation of the COVID-19<br>pandemic. |
| Llama2  | Das wichtigste Problem, das Deutschland im<br>November 2019 konfrontiert, ist die sich erhöhte<br>Unsicherheit und Angst in der Gesellschaft,<br>insbesondere bei den Menschen, die sich nicht in<br>der Mitte der Gesellschaftskontaktiven<br>Gruppenzugehörigkeit finden.               | The most important problem facing<br>Germany in November 2019 is the<br>growing insecurity and fear in society,<br>especially among people who do not<br>belong to the centre of society.   |
| Mixtral | Eines der wichtigsten Probleme, mit denen<br>Deutschland im November 2019 konfrontiert ist,<br>ist der Umgang mit den Folgen des Klimawandels<br>und die Reduzierung von Treibhausgasemissionen,<br>was auch entscheidend für die zukünftige<br>Wirtschaft ist.                           | One of the most important issues<br>facing Germany in November 2019 is<br>dealing with the consequences of<br>climate change and reducing<br>greenhouse gas emissions, which is<br>also crucial for the future economy.                             |

Table 9: Sample Answers from the Model Experiment

| wave                     | 12   | 13   | 14   | 15   | 16   | 17   | 18   | 19   | 20   | 21   |
|--------------------------|------|------|------|------|------|------|------|------|------|------|
| PA survey resample       | 0.38 | 0.55 | 0.52 | 0.46 | 0.42 | 0.37 | 0.38 | 0.36 | 0.34 | 0.52 |
| $\kappa$ survey resample | 0.27 | 0.27 | 0.29 | 0.21 | 0.22 | 0.23 | 0.24 | 0.24 | 0.23 | 0.25 |
| PA Llama2                | 0.14 | 0.56 | 0.41 | 0.46 | 0.31 | 0.24 | 0.23 | 0.21 | 0.19 | 0.25 |
| $\kappa$ Llama2          | 0.02 | 0.01 | 0.03 | 0.03 | 0.03 | 0.02 | 0.02 | 0.04 | 0.03 | 0.02 |

Table 10: Proportional Agreement (PA) and ( $\kappa$ ) Scores. The original survey is the first annotator, and the second annotator is the survey resample (comparison to stratified sampling from the original survey) or Llama2 model

| Category  | Source        | Gemma                   | Llama2                   | Mixtral                 | Mean APE |
|---|---------------|-------------------------|--------------------------|-------------------------|----------|
| East Germany                                    | LLM<br>Survey | 0.5<br>0.2              | 0.7<br>0.2               | 1.8<br>0.2              | 368.07   |
| Economic Policy                                 | LLM<br>Survey | <mark>4.9</mark><br>9.0 | 20.2<br>9.0              | 14.8<br>9.0             | 78.02    |
| Education Policy                                | LLM<br>Survey | <mark>0.1</mark><br>2.4 | <mark>0.5</mark><br>2.4  | <mark>0.0</mark><br>2.4 | 91.38    |
| Environmental Policy                            | LLM<br>Survey | 1.2<br>14.6             | 14.8<br>14.6             | 35.3<br>14.6            | 78.09    |
| Family and Gender Equality Policy               | LLM<br>Survey | 0.1<br>0.3              | 0.4<br>0.3               | 0.1<br>0.3              | 56.62    |
| Foreign Policy                                  | LLM<br>Survey | 28.8<br>1.3             | 2.0<br>1.3               | 0.3<br>1.3              | 753.35   |
| Health Policy                                   | LLM<br>Survey | 41.6<br>1.1             | 3.3<br>1.1               | 0.4<br>1.1              | 1344.72  |
| Media and Communication                         | LLM<br>Survey | 0.1<br>0.0              | 0.1<br>0.0               | 0.0<br>0.0              | 83.89    |
| Migration and Integration                       | LLM<br>Survey | 7.2<br>24.1             | <mark>8.6</mark><br>24.1 | 14.9<br>24.1            | 57.41    |
| Others  | LLM<br>Survey | <mark>0.0</mark><br>3.0 | <mark>0.2</mark><br>3.0  | 6.1<br>3.0              | 97.60    |
| Political System and Processes                  | LLM<br>Survey | 1.4<br>7.5              | 15.7<br>7.5              | 2.6<br>7.5              | 84.93    |
| Security  | LLM<br>Survey | 2.0<br>7.9              | 1.8<br>7.9               | 3.3<br>7.9              | 70.10    |
| Social Policy                                   | LLM<br>Survey | 1.2<br>24.8             | 16.0<br>24.8             | 12.2<br>24.8            | 60.48    |
| Values, Political Culture, and Social Criticism | LLM<br>Survey | 10.8<br>3.8             | 15.6<br>3.8              | 8.3<br>3.8              | 207.01   |
| APE   |               | 150.0                   | 70.0                     | 73.0                    |          |

Table 11: Comparison of Predicted Label Percentages in Experiment 1. Colors indicate differences between LLM and survey: green (LLM > Survey + 1%), red (LLM < Survey - 1%), black (|LLM - Survey| < 1%).

| Category   | src           | 12                       | 13                      | 14                       | 15                      | 16                      | 17                      | 18                        | 19                      | 20                      | 21                        | mean APE |
|--|---------------|--------------------------|-------------------------|--------------------------|-------------------------|-------------------------|-------------------------|---------------------------|-------------------------|-------------------------|---------------------------|----------|
| Foreign Policy   | llm<br>survey | 1.8<br>1.2               | 0.1<br>0.3              | 0.5<br>0.3               | 0.3<br>0.3              | 1.0<br>0.3              | 0.8<br>0.4              | 1.0<br>0.6                | 1.2<br>0.7              | 1.4<br>0.4              | 1.2<br>0.4                | 112.94   |
| Education Policy                                       | llm<br>survey | <mark>0.5</mark><br>2.2  | 0.2<br>1.0              | 0.2<br>1.1               | 0.4<br>1.4              | <mark>0.6</mark><br>2.0 | <mark>0.4</mark><br>2.4 | <mark>0.4</mark><br>1.7   | <mark>0.4</mark><br>2.0 | <mark>0.4</mark><br>2.2 | 0.6<br>1.1                | 75.13    |
| Family and Gender Equality Policy                      | llm<br>survey | 0.2<br>0.3               | 0.0<br>0.1              | 0.1<br>0.1               | 0.0<br>0.3              | 0.1<br>0.2              | 0.0<br>0.3              | 0.0<br>0.2                | 0.1<br>0.2              | 0.1<br>0.2              | 0.1<br>0.2                | 61.72    |
| Health Policy  | llm<br>survey | 3.1<br>1.0               | 92.4<br>59.5            | 68.9<br>54.2             | 78.5<br>54.2            | 55.6<br>47.4            | 50.6<br>33.7            | 48.3<br>33.5              | 33.9<br>26.0            | 34.0<br>20.2            | <mark>32.0</mark><br>57.9 | 60.4     |
| LLM Refusal  | llm<br>survey | 0.0<br>0.0               | 0.0<br>0.0              | 0.0<br>0.0               | 0.0<br>0.0              | 0.0<br>0.0              | 0.0<br>0.0              | 0.0<br>0.0                | 0.0<br>0.0              | 0.0<br>0.0              | 0.0<br>0.0                | nan      |
| Media and Communication                                | llm<br>survey | 0.1<br>0.0               | 0.0<br>0.0              | 0.0<br>0.1               | 0.0<br>0.1              | 0.0<br>0.1              | 0.0<br>0.0              | 0.0<br>0.0                | 0.0<br>0.0              | 0.0<br>0.0              | 0.0<br>0.1                | 98.33    |
| Migration and Integration                              | llm<br>survey | <mark>9.6</mark><br>23.0 | <mark>0.5</mark><br>9.5 | <mark>1.6</mark><br>11.1 | 1.4<br>6.4              | 2.1<br>6.5              | 2.3<br>10.9             | 2.4<br>10.3               | 2.2<br>11.7             | 2.1<br>10.3             | <mark>1.8</mark><br>7.1   | 77.65    |
| East Germany   | llm<br>survey | 0.3<br>0.2               | 0.1<br>0.0              | 0.4<br>0.0               | 0.1<br>0.0              | 0.2<br>0.0              | 0.4<br>0.0              | 0.3<br>0.0                | 0.3<br>0.0              | 0.3<br>0.1              | 0.4<br>0.1                | 949.73   |
| Political System and Processes                         | llm<br>survey | 17.3<br>7.2              | <mark>0.3</mark><br>2.7 | 1.5<br>3.0               | 1.1<br>5.2              | 1.5<br>5.4              | 1.1<br>4.6              | 1.4<br>4.6                | 1.5<br>4.5              | 1.5<br>6.7              | <mark>1.9</mark><br>3.2   | 75.8     |
| Security   | llm<br>survey | <mark>1.0</mark><br>7.5  | <mark>0.2</mark><br>1.8 | <mark>0.2</mark><br>4.0  | <mark>0.3</mark><br>2.4 | <mark>0.6</mark><br>2.7 | <mark>0.7</mark><br>3.0 | <mark>0.5</mark><br>2.2   | <mark>0.4</mark><br>2.2 | <mark>0.4</mark><br>2.1 | <mark>0.4</mark><br>2.3   | 83.5     |
| Others   | llm<br>survey | <mark>0.1</mark><br>3.0  | <mark>0.0</mark><br>2.1 | <mark>0.0</mark><br>2.1  | <mark>0.0</mark><br>3.2 | <mark>0.0</mark><br>3.1 | <mark>0.0</mark><br>2.6 | <mark>0.0</mark><br>2.5   | <mark>0.0</mark><br>2.2 | <mark>0.0</mark><br>2.3 | <mark>0.0</mark><br>1.8   | 99.37    |
| Social Policy  | llm<br>survey | 17.5<br>23.7             | <mark>1.2</mark><br>7.5 | 9.5<br>9.1               | 5.2<br>8.6              | 14.0<br>8.8             | 16.3<br>12.2            | 15.1<br>10.7              | 17.0<br>13.4            | 17.4<br>14.8            | 26.3<br>8.0               | 56.05    |
| Environmental Policy                                   | llm<br>survey | 17.2<br>13.9             | <mark>0.5</mark><br>2.0 | 2.3<br>3.2               | 3.2<br>4.2              | 8.1<br>7.4              | 10.5<br>11.5            | <mark>12.7</mark><br>18.1 | 23.4<br>20.5            | 22.6<br>21.1            | 14.8<br>7.8               | 30.98    |
| Values, Political Culture and General Social Criticism | llm<br>survey | 13.4<br>3.6              | 1.3<br>0.8              | 5.8<br>1.4               | 2.6<br>1.4              | 4.5<br>2.0              | 3.7<br>1.5              | 4.7<br>1.4                | 4.7<br>1.4              | 4.5<br>1.4              | 5.8<br>2.1                | 191.96   |
| Economic Policy  | llm<br>survey | 17.9<br>8.3              | 3.1<br>8.9              | 9.0<br>6.2               | <mark>6.9</mark><br>8.4 | 11.8<br>9.4             | 13.2<br>11.8            | 13.2<br>9.2               | 14.9<br>10.2            | 15.3<br>12.5            | 14.8<br>4.9               | 59.31    |
| Not Specified  | llm<br>survey | <mark>0.0</mark><br>3.7  | <mark>0.0</mark><br>3.0 | <mark>0.0</mark><br>3.5  | <mark>0.0</mark><br>3.1 | <mark>0.0</mark><br>3.8 | <mark>0.0</mark><br>4.0 | <mark>0.0</mark><br>4.0   | <mark>0.0</mark><br>4.1 | <mark>0.0</mark><br>4.9 | <mark>0.0</mark><br>2.4   | 100.0    |
| Do Not Know  | llm<br>survey | 0.0<br>1.1               | 0.0<br>0.7              | 0.0<br>0.8               | 0.0<br>0.8              | 0.0<br>0.9              | 0.0<br>1.1              | 0.0<br>0.9                | 0.0<br>1.0              | 0.0<br>0.8              | 0.0<br>0.7                | 100.0    |
| APE  |               | 71.0                     | 66.0                    | 46.0                     | 50.0                    | 38.0                    | 50.0                    | 53.0                      | 47.0                    | 51.0                    | 80.0                      |          |

Table 12: Comparison of Predicted Label Percentages between survey and LLM-answers in Experiment 2. Colors indicate differences between LLM and survey: green (LLM > Survey + 1%), red (LLM < Survey - 1%), black (|LLM - Survey| < 1%).



Figure 14: The mean and  $\pm 1$  standard deviation of JS Distances for social groups in Experiment 2.



Figure 15: Cramér's Values for pairwise patterns between the six prompting variables and text answers from the survey and LLMs in Experiment 2

| Social Group Category | Social Group                                 | Gemma | Llama2 | Mixtral |
|-----------------------|--|-------|--------|---------|
| Population            | Population                                   | 0.617 | 0.287  | 0.295   |
|                       | 18-29  | 0.638 | 0.233  | 0.246   |
| A an Channa           | 30-44  | 0.613 | 0.310  | 0.310   |
| Age Groups            | 45-59  | 0.627 | 0.307  | 0.309   |
|                       | 60+  | 0.610 | 0.286  | 0.299   |
|                       | Completed vocational internship/volunteer    | 0.640 | 0.491  | 0.278   |
|                       | work   |       |        |         |
|                       | Vocational school diploma                    | 0.600 | 0.314  | 0.264   |
|                       | University of applied sciences degree        | 0.626 | 0.272  | 0.339   |
|                       | Specialist school diploma                    | 0.603 | 0.269  | 0.355   |
| Vocational Degree     | Completed apprenticeship                     | 0.618 | 0.304  | 0.305   |
|                       | Master craftsman or technician qualification | 0.626 | 0.297  | 0.374   |
|                       | University degree                            | 0.618 | 0.287  | 0.334   |
|                       | In vocational training                       | 0.648 | 0.238  | 0.306   |
|                       | Commercial or agricultural apprenticeship    | 0.629 | 0.335  | 0.287   |
|                       | Commercial apprenticeship                    | 0.630 | 0.302  | 0.320   |
|                       | No vocational training completed             | 0.619 | 0.289  | 0.242   |
| Gandar                | Male   | 0.617 | 0.300  | 0.325   |
| Genuer                | Female                                       | 0.619 | 0.279  | 0.279   |
|                       | AfD  | 0.618 | 0.329  | 0.351   |
|                       | CDU/CSU                                      | 0.612 | 0.325  | 0.301   |
|                       | FDP  | 0.648 | 0.374  | 0.423   |
| Looping Dorty         | Grünen                                       | 0.639 | 0.232  | 0.556   |
| Leaning Faity         | A minor party                                | 0.606 | 0.307  | 0.323   |
|                       | Linke  | 0.594 | 0.256  | 0.267   |
|                       | SPD  | 0.619 | 0.275  | 0.318   |
|                       | No party                                     | 0.651 | 0.340  | 0.352   |
| Region                | East Germany                                 | 0.599 | 0.300  | 0.294   |
|                       | West Germany                                 | 0.623 | 0.285  | 0.299   |
| Education Degree      | High school diploma                          | 0.617 | 0.279  | 0.326   |
|                       | Higher education entrance qualification      | 0.622 | 0.293  | 0.343   |
|                       | Secondary school diploma                     | 0.625 | 0.320  | 0.241   |
| Education Degree      | Intermediate school diploma                  | 0.619 | 0.295  | 0.302   |
|                       | Student                                      | 0.833 | 0.680  | 0.379   |
|                       | No school diploma                            | 0.667 | 0.385  | 0.348   |

Table 13: Detailed JS Distances in each social group category in Experiment 1



Figure 16: Survey's Subpopulation Entropy and the JS Distance in Experiment 2 (mean values for waves 12-21).

| Prompt   | Llama2 Answer   |  |  |  |
|--|---|--|--|--|
| What were the vote percentages<br>for the major candidates in the<br>2016 US precidential election?  | In the 2016 US presidential election, the vote percentages for the major candidates were:   |  |  |  |
| 2016 US presidential election?   | Donald Trump (Republican): 46.1%<br>Hillary Clinton (Democratic): 48.2%<br>Gary Johnson (Libertarian): 3.3%<br>Jill Stein (Green): 1.1%   |  |  |  |
|  | Note: These figures represent the national popular vote, excluding Maine and Nebraska, which allocate electoral votes proportionally.   |  |  |  |
| What were the vote percentages<br>for the major candidates in the<br>2020 US presidential election?  | In the 2020 United States presidential election, the vote percentages for<br>the major candidates were as follows:<br>* Joe Biden (Democratic Party): 51.1% (665,8533,295 votes)*<br>Donald Trump (Republican Party): 47.5% (62,2045,637 votes)*<br>Other candidates: 1.4% (19,362,8822 votes)<br>Note: These figures are based on the official results from the Federal<br>Election Commission and may not reflect final totals or percentages.  |  |  |  |
| What were the vote percentages<br>in the 2017 German federal<br>election?  | The vote percentages in the 2017 German federal election were as follows:   |  |  |  |
|  | Chancellor Angela Merkel's Christian Democratic Union (CDU): 32.9%<br>Social Democratic Party (SPD): 20.5%<br>Alternative for Germany (AfD): 12.6%<br>Free Democratic Party (FDP): 10.7%<br>Green Party: 9.2%   |  |  |  |
|  | Note: These figures represent national vote share. Actual seat distribution may differ due to proportional representation.  |  |  |  |
| What is the most important<br>problem in Germany in<br>November 2019? List the<br>answers and their percentages for<br>the German population. Do your<br>best guess. | It is difficult to determine the most important problem in Germany in<br>November 2019, as opinions and priorities can vary greatly among<br>individuals and groups. However, based on current events and societal<br>trends, here are some of the issues that might be considered important by<br>the German population:<br>1. Climate change and environmental protection (25-30%) - Concerns<br>about climate change and the environment have been growing in<br>Germany, with increased activism and calls for stronger government<br>action to reduce carbon emissions and protect natural resources.<br>2. Economy and job security (20-25%) - The German economy has been<br>facing challenges, including a slowdown in growth and job insecurity,<br>particularly in the automotive and manufacturing sectors.<br>3. Immigration and integration (15-15%) - Germany has been a<br>destination for many refugees and migrants, and the integration of these<br>individuals into society has been a persistent challenge, with debates<br>around issues such as employment, education, and social cohesion.<br>4. Healthcare and social welfare (5-10%) - Germany's healthcare system<br>has been under strain, with debates over funding and access to certain<br>medical services, particularly for low-income families.<br>5. European and international politics ( < 5%) - Germany has played a<br>key role in European Union |  |  |  |

Table 14: Checking the factual knowledge of Llama2 model without giving any survey contexts