

# ElectionSim: Massive Population Election Simulation Powered by Large Language Model Driven Agents

Xinnong Zhang<sup>1†</sup>, Jiayu Lin<sup>1†</sup>, Libo Sun<sup>1†</sup>, Weihong Qi<sup>2</sup>, Yihang Yang<sup>1</sup>, Yue Chen<sup>1</sup>, Hanjia Lyu<sup>3</sup>, Xinyi Mou<sup>1</sup>, Siming Chen<sup>1‡</sup>, Jiebo Luo<sup>3‡</sup>, Xuanjing Huang<sup>1</sup>, Shiping Tang<sup>1</sup>, Zhongyu Wei<sup>1‡</sup>

<sup>1</sup>Fudan University, <sup>2</sup>Indiana University, <sup>3</sup>University of Rochester  
 {xnzhang23, jiayulin24, lbsun23, yhyang24, chen\_y24}@m.fudan.edu.cn,  
 {xymou20, simingchen, xjhuang, twukong, zywei}@fudan.edu.cn,  
 wq3@iu.edu, hlyu5@ur.rochester.edu, jluo@cs.rochester.edu  
 Project Page: <http://www.fudan-disc.com/electionsim>

## Abstract

The massive population election simulation aims to model the preferences of specific groups in particular election scenarios. It has garnered significant attention for its potential to forecast real-world social trends. Traditional agent-based modeling (ABM) methods are constrained by their ability to incorporate complex individual background information and provide interactive prediction results. In this paper, we introduce **ElectionSim**, an innovative election simulation framework based on large language models, designed to support accurate voter simulations and customized distributions, together with an interactive platform to dialogue with simulated voters. We present a **million-level voter pool** sampled from social media platforms to support accurate individual simulation. We also introduce **PPE**, a poll-based presidential election benchmark to assess the performance of our framework under the U.S. presidential election scenario. Through extensive experiments and analyses, we demonstrate the effectiveness and robustness of our framework in U.S. presidential election simulations.

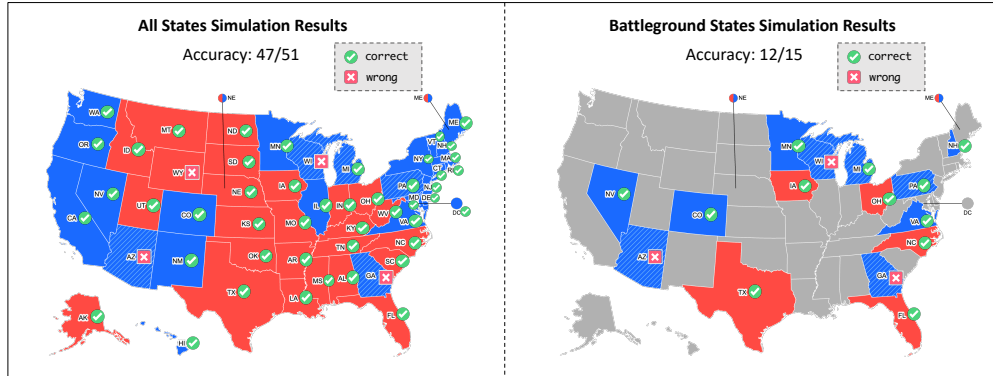


Figure 1: Simulation results of the 2020 Presidential Election. The colors represent the real-world results and the marks represent the simulation results accuracy.

<sup>†</sup>These authors contribute equally to this work.

<sup>‡</sup>Corresponding authors.

# 1 Introduction

Massive Population Election Simulation aims to simulate election events at a large population scale, which has been of vital importance in forecasting potential real-world political trends and capturing specific groups’ preferences on particular topics or special events [29; 44; 13]. Previous works have demonstrated that using mathematical/statistical models can transform traditional political and sociological analysis paradigms [25; 43].

Currently, the mainstream method for election simulation is agent-based modeling (ABM) [55; 38; 32; 15], which employs heuristic-like rules or mathematics functions to simulate the actions of individuals [59], and then scales up these actions to forecast the collective result. The ABM approach, despite its merits, presents two major limitations. First, the integration of individuals’ complex background information into the framework of ABMs is a non-trivial task. Second, the establishment of interaction human-agent interfaces to elucidate predictive results remains a significant challenge [12; 17]. These two shortcomings can result in predictions that are less accurate and less persuasive.

Recently, agent simulations powered by Large Language Models (LLMs) have gained significant attention. Researchers have implemented these simulations at both individual and task levels, focusing on generating highly reliable human-like behavior [56; 63; 66; 49] and facilitating multi-agent collaboration [21; 30; 50]. However, existing studies struggle to address massive population election simulations, as achieving the required levels of diversity and quantity remains difficult.

The U.S. Presidential Election, as a big election event, plays a pivotal role in shaping public engagement and party strategies [3; 53]. We use the U.S. presidential election as a case to explore effective methods for achieving massive and diverse election simulations with LLMs. To facilitate large-scale election simulations using LLMs, three primary challenges must be addressed.

## Q1. How to achieve high accuracy in individual-level simulations?

The macro result is aggregated from individuals so that the preciseness of the individual directly impacts the overall outcome of the massive simulation. LLMs often lack sufficient personalized input for simulating the nuanced behavior of individuals. The limited contextual data restricts their ability to fully capture the diversity of voter behavior, motivations, and decision-making processes, which are critical for generating accurate and meaningful simulations.

## Q2. How to generate customized distributions that align with real-world statistics?

Accurate election simulation requires that the simulated individuals represent the diversity and aligned distribution of real-world populations. While random sampling is able to capture this diversity, it falls short when aligning to the demographic distribution of the real world and is prone to source-driven biases [26; 61; 16; 68; 52]. As a result, a carefully designed sampling strategy that mirrors real-world demographic and behavioral distributions is essential for producing valid and reliable simulations.

## Q3. How to evaluate the performance of election simulation in a systematic way?

Evaluation metrics for election simulations vary depending on the specific context and task. Most existing works primarily focus on prediction accuracy, which offers a limited and unsystematic approach to assessing the full scope of simulation outcomes. Consequently, it is crucial to design a multi-aspects evaluation method to benchmark election results and provide comprehensive analyses.

We address these challenges by introducing **ElectionSim**: a massive population election simulation framework powered by large language models. For accurate individual simulation, We collect 171,210,066 tweets from Twitter between January 1, 2020, and December 29, 2020, to construct a large and diverse voter pool with **million-level distinct users**. For real-world voter distribution alignment, we employ a demographic sampling strategy to align the distribution between sampled users and real-world voters. To provide a systematic evaluation for the election simulation, we construct a **Poll-based Presidential Election (PPE)** benchmark, integrating three evaluation baselines.

Extensive experiments show that LLMs achieve true-to-life performance compared to actual election outcomes following our ElectionSim framework. In voter-level simulations, we achieve a Micro-F1 score of 0.812 for vote-related tasks. At the state level, we accurately predict the outcomes of the 2020 presidential election in **47/51** states, and our predictions match the actual result in **12/15** battleground states as shown in Figure 1.

To conclude, we make four major contributions:

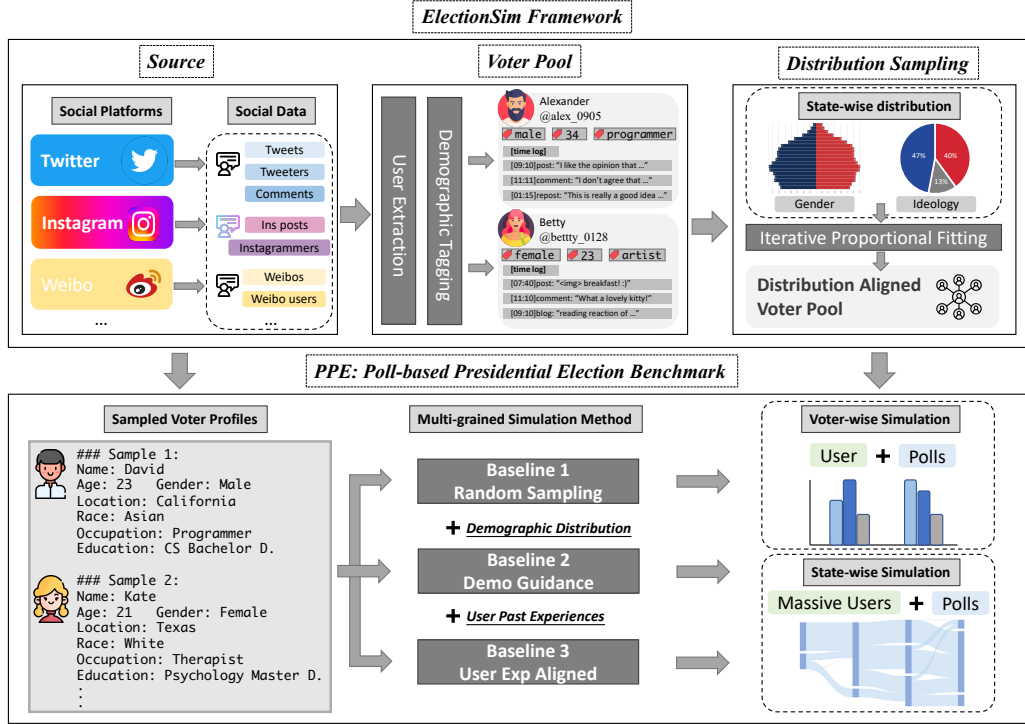


Figure 2: An illustration of the overall pipeline of the U.S. presidential election simulations.

- **ElectionSim**: a massive population election simulation framework, which allows for diverse election simulating scenarios with high confidence by employing a flexible, customized distribution sampling strategy to align with the real world.
- **Million-level voters pool**: a large and diverse voters pool to support massive population election simulation by collecting and combing data from social media platforms.
- **PPE**: a poll-based presidential election benchmark to validate our simulation results. Experiments and further analyses confirm the benchmark’s robustness and demonstrate the effectiveness of our framework in U.S. presidential election simulations.
- **Interactive Simulation Demonstration**: a demonstration system that facilitates multi-round dialogue interactions with voters we simulated. The system allows voters to be selected in a variety of ways, including specifying attributes and specifying responses to specific questions.

In the following sections, we explore the core contents of ElectionSim and PPE. In Section 2, we detail the key elements of the ElectionSim framework. In Section 3, we introduce the PPE benchmark. Detailed experiment results are displayed in Section 4 with further analysis discussed in Section 5. Additionally, we provide a visualization of the U.S. presidential election in Section 6. Supplementary materials including additional simulation experiments, data processing details, prompt libraries, case studies, and questionnaires can be found from Appendix A to Appendix E respectively.

## 2 ElectionSim: Massive Population Election Simulation

To address the challenges outlined in Section 1, we introduce **ElectionSim**: a massive population election simulation framework. The overall pipeline is shown in Figure 2. The ElectionSim framework operates as follows: (1) Raw data, including user information and historical activities such as posts and comments, are collected from social platforms; (2) The data are processed on a per-user basis, and then a demographic tagging process, performed by specialized classifiers, generates a diverse and massive voter pool (§2.1); (3) An interactive proportional fitting sampling strategy is used to approximate marginal distributions into a joint distribution for the target group sampling. After the

pipeline of ElectionSim, we enable the sampling of any customized distribution from the voter pool, facilitating realistic massive election simulation (§2.2).

## 2.1 Massive and Diverse Voter Pool Construction

### 2.1.1 Social Media Data Collection

**Source Data** We choose Twitter\* (now known as X) as the primary data source due to its extensive user base. We collect tweets published by Twitter users in 2020, which allows us to construct a massive and varied voter pool that reflects diverse demographic and ideological perspectives. Table 1a details the specific fields we gather from the Twitter platform, and Table 1b shows the characteristics of the collected data.

Dimension	Field	Description	Feature	Value
User	user_id	The ID of the user who posted this tweet.	Number of Users	9,596,198
	user_at_name	The name of the user who posted this tweet.	Number of Tweets	171,210,066
	tweet_id	A unique ID to represent a tweet.	Avg. Number of Words in Tweets	21.69
Post	tweet_content	The content of this tweet.	Time Span of Tweets	2020.1.1–2020.12.29
	pub_time	The publishing time of this tweet.	Number of Languages in Tweets	50
	lang	The language of this tweet.		

(a) Data fields.

(b) Descriptive statistics.

Table 1: An overview of Twitter dataset characteristics.

**Data Preprocessing** To tailor our voter pool for the context of the U.S. election, we implement a data-cleaning process consisting of four stages: user aggregation, language filtering, post filtering, and user cleaning.

- **User Aggregation** The data collected is initially structured at the post level; in order to construct a voter pool, we aggregate the data to a user-level framework. The key statistical results after this aggregation process are presented in Figure 3.
- **Language Filtering** We start by removing non-English posts to better reflect the perspectives of American users (detailed in Appendix C.1). This is important for accurately interpreting sentiments related to the U.S. election, as language can significantly shape political discourse [47; 48].
- **Post Filtering** Next, we retain only those users who have made over 30 posts and sample 30 of them from all posts as their historical information.
- **User Cleaning** Finally, we evaluate content repeatability by calculating the overlap score from a random sample of 5 tweets from each user. We compute the Jaccard scores for all pairs, and the final score is the average of these values (detailed in Appendix C.2).

Metric	Value
Number of Users	1,006,517
Number of Tweets	30,195,510
Avg. Number of Words in Tweets	22.36

Table 2: Statistical summary of the processed Twitter user dataset.

By applying these stages, we enhance the reliability and relevance of our dataset, making it more suitable for analyzing user behaviors in the context of the U.S. election. Table 2 presents the statistical summary of the processed Twitter user data.

### 2.1.2 Demographic Feature Annotation

The data collected from Twitter primarily consists of user post histories and lacks essential demographic features necessary for effective social group simulation. To enhance the dataset with

\*<https://x.com>

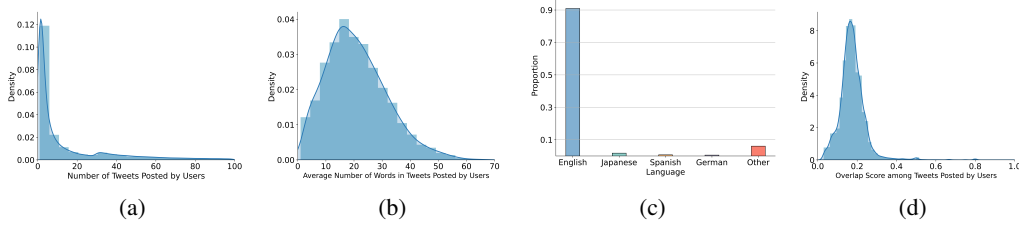


Figure 3: Statistical charts based on user aggregation, including: a) Density plot of tweets per user, b) Density plot of average word count per tweet, c) User distribution across languages, and d) Density plot of user overlap scores.

reliable demographic attributes, we develop a comprehensive taxonomy of demographic features, accompanied by corresponding classifiers to annotate the Twitter data systematically.

**Demographic Taxonomy** We conduct an analysis of multiple U.S. polls [2; 4; 40] and synthesize the demographic classifications pertinent to these studies. Through this process, we develop a comprehensive demographic taxonomy that is closely aligned with the context of the U.S. election, as shown in Table 3.

Dimension	Attribute	Classification
Personal Traits	Gender	Male Female
	Age	Youth (18-35 years old) Middle-aged (36-65 years old) Elderly (over 65 years old)
	Race	White Black Asian Hispanic
	Ideology	Liberal Moderate Conservative
Political Orientation	Partisanship	Democrat Republican Independent Others

Table 3: The demographic label taxonomy.

**Demographic Classifier** To develop a demographic dataset for classifiers, we must establish a ground-truth test dataset for evaluation and a reliable training dataset for effective model training. The following outlines the process for constructing these datasets and training the classifiers.

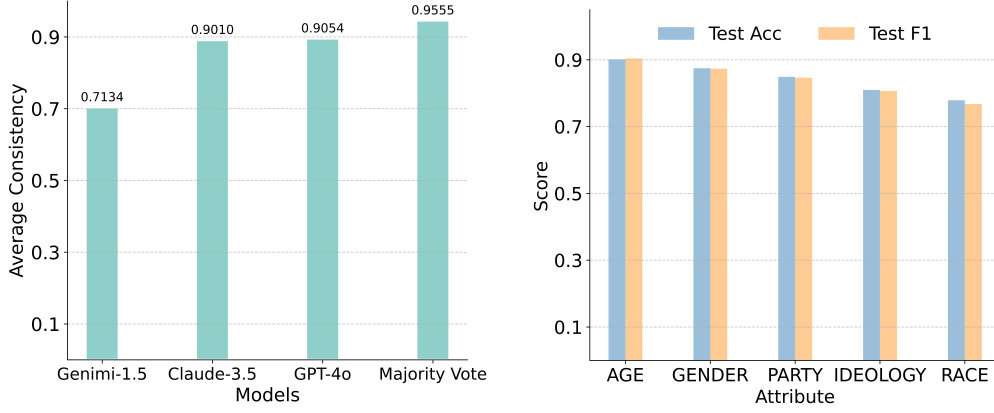
- **Test Set Construction** We use a semi-automated method to create the ground-truth test data in two stages: API annotation and manual verification. First, we sample 200 users from the Twitter dataset and annotate them using advanced commercial LLM APIs based on our demographic taxonomy (detailed in Appendix C.3), including GPT-4o<sup>†</sup>, Claude-3.5<sup>‡</sup>, and Gemini-1.5<sup>§</sup>. We then engage 5 professional annotators to validate the API-generated annotations and finalize the ground-truth labels.

<sup>†</sup>gpt-4o-2024-08-06

<sup>‡</sup>claude-3-5-sonnet-20240620

<sup>§</sup>gemini-1.5-pro

- **Train Set Construction** To create a large train dataset, we utilize commercial APIs for annotation. We evaluate the consistency between each API’s annotations and the ground-truth labels on the test set, as well as the consistency between the majority voting results from these APIs and the ground-truth labels (detailed in Appendix C.4, results shown in Figure 4a). It shows that majority voting achieves the highest consistency. Consequently, we sample 10,000 users from the processed Twitter dataset, label them using three APIs, and adopt the majority voting results as the labels for our training set.
- **Implementation Details** We choose the Longformer<sup>¶</sup> [5] model as the backbone for our classifiers due to its capability to process long text windows of up to 4096 tokens. We develop classifiers for each demographic feature, resulting in a total of 5 classifiers. For each classifier, we implement full parameter fine-tuning, set the learning rate  $5 \times 10^{-5}$ , batch size to 16, employ AdamW as the optimizer, and train the model 3 epochs on 8 NVIDIA RTX4090 GPUs.
- **Classifiers Performance** For each demographic feature, we select the checkpoint with the highest F1 score on the test set as the final model. The performance of the classifiers obtained on the test set is shown in Figure 4b.



(a) Consistency between commercial APIs and manual annotations on the test set (b) The accuracy and F1 score of each classifier on the test set

Figure 4: Classification accuracy of API-based LLMs and our demographic classifiers.

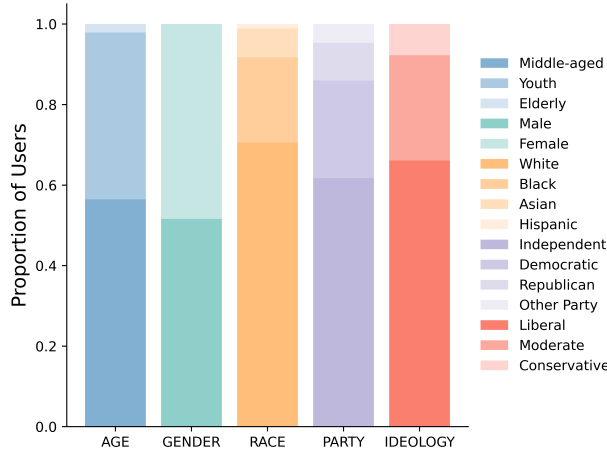


Figure 5: Distribution of attribute categories in the voter pool

<sup>¶</sup><https://huggingface.co/allenai/longformer-base-4096>

## 2.2 Real-world Demographic Distribution Sampling Strategy

### 2.2.1 Real-World Demographics in Electoral Processes

Demographics play a significant role in U.S. elections, including presidential races. Political science literature widely acknowledges that factors such as gender, race, and income level influence U.S. citizens’ voting decisions [39; 60]. Moreover, evolving demographic dynamics are significantly impacting candidates’ campaign strategies and the overall election outcomes. For instance, as the number of Hispanic voters grows, many are shifting away from the Democratic Party and becoming more receptive to Republican candidates [1], which is notably affecting electoral results in states with large Hispanic populations, such as Texas and Florida. In addition to demographics, ideology and partisanship are also critical predictors of U.S. voting behavior. Therefore, to ensure our large-scale simulation of election results closely mirrors reality, it is essential to incorporate accurate distributions of demographics, ideology, and partisanship. In the remainder of this section, we introduce the methods used to simulate the distribution of these real-world attributes.

### 2.2.2 Methods for Distribution Sampling

To accurately reflect the demographic and ideological makeup of U.S. citizens, we combine multiple datasets to construct the population distribution in our study. Specifically, we utilize data from the U.S. Census Bureau’s Voting and Registration in the Election of November 2022, along with the 2020 Time Series Study from the American National Election Studies (ANES) [2]. The U.S. Census Bureau has been collecting data on the characteristics of American voters since 1964. This dataset primarily tracks how many citizens of voting age are registered and how many vote, broken down by factors such as age, gender, race, ethnicity, and more. In contrast, the American National Election Studies (ANES) consist of academically-run national surveys of U.S. voters, conducted before and after every presidential election. The ANES data extend a long tradition of studies, with records dating back to 1948. The two datasets allow us to simulate real population distributions based on variables such as state of registration, gender, age, race, ideology, and partisanship.

The datasets chosen are the most recent and frequently cited in U.S. election studies [4; 40], while the selected features also capture fundamental demographic and ideological attributes that are key predictors of voting behavior [62]. Table 4 provides the sources of all the attributes used in the population simulation, along with the classifications for each attribute.

Dataset	Attribute	Classification
U.S. Census Bureau	Gender	Male Female
	Age	18-24 years old 25-34 years old 35-44 years old 45-64 years old over 65 years old
	Race	White Black Asian Hispanic
	Ideology	Liberal Moderate Conservative
ANES	Partisanship	Democrat Republican Independent Others

Table 4: The Demographic Distribution Sampling Attributes and Classifications

The U.S. Census Bureau provides accurate demographic data on voters across all 50 states and the District of Columbia for 2022, including variables such as gender, race, and age, which we directly incorporate into our simulation. In contrast, the ANES dataset samples voters from each state and includes individual interview responses. Therefore, the dataset does not provide direct information about the distribution of voters’ ideology, partisanship, or community type at the state level. However, according to the ANES methodology report [18], the sampling strategy involves random selection based on mailing addresses and is stratified by area. This method produces a sample that closely mirrors the true demographic and ideological distribution of U.S. citizens across states. As a result, we treat the distribution of these attributes in the ANES data as representative of their actual distribution in each state.

### 2.2.3 Iterative Proportional Fitting

In addition to using the distributions of individual variables to capture population characteristics, we also generate the joint distribution of all variables to more accurately represent voters in each state. Specifically, we apply Iterative Proportional Fitting (IPF) to estimate the joint distribution of all attributes within each state. IPF is an iterative method that estimates missing values for joint population attributes, providing the foundation for a more comprehensive simulation. This method is widely used in social science and population research [14].

In our study, we follow the classical IPF method to construct the joint distribution of all the attributes in our simulation. Specifically, we start with a two-way table with individual components denoted as  $x_{ij}$  and targeted estimation  $\hat{x}_{ij}$ . The targeted estimation  $\hat{x}_{ij}$  satisfies  $\sum_j \hat{x}_{ij} = v_i$  and  $\sum_i \hat{x}_{ij} = w_j$ . The iterations are specified as follows:

Let  $\hat{x}_{ij}^{(0)} = x_{ij}$ . For  $\alpha > 1$ :

$$\hat{x}_{ij}^{(2\alpha-1)} = \frac{\hat{x}_{ij}^{(2\alpha-2)} v_i}{\sum_{k=1}^J \hat{x}_{ik}^{(2\alpha-2)}}$$

$$\hat{x}_{ij}^{(2\alpha)} = \frac{\hat{x}_{ij}^{(2\alpha-1)} w_j}{\sum_{k=1}^I \hat{x}_{kj}^{(2\alpha-1)}}$$

The iterations end when the estimated marginals are sufficiently close to the real marginals or when they stabilize without further convergence.

In our study, we implement the IPF algorithm for each state using five attributes: gender, race, age group, ideology, and partisanship. In most cases, the algorithm does not converge, but the gaps between the estimated and actual marginals are less than 5%, with 888 out of 918 marginals falling within this range. For the outliers, since IPF adjusts proportionally to the marginals, the overall ratio of marginals remains consistent. We then use the estimated joint distribution and marginals for our massive simulation.

## 3 PPE: Poll-based Presidential Election Benchmark

In this section, we introduce the U.S. presidential election prediction benchmark via a poll-based survey method, namely **PPE**. The task is clarified in §3.1 and the source of the benchmark is introduced in §3.2. §3.3 provides a comprehensive explanation of the questionnaire design process, along with an overview of the questionnaire’s statistics. Furthermore, we introduce the multi-stage evaluation baseline in §3.4, enabling the massive simulation execution.

### 3.1 Task Definition

The poll-based presidential election prediction task is an application example of the massive social simulation paradigm. This task simulates the poll results of each person within a group at a fine granularity and aggregates the results according to certain distribution patterns (such as demographic information at the state level) to obtain the macro-level social opinion tendencies.



### 3.2 Poll Collection

We use the American National Election Studies (ANES) 2020 Time Series Study data as the source of the benchmark [2]. The ANES 2020 Time Series Study is a continuation of a series of studies conducted since 1948 to enable analysis of public opinion and voting behavior in U.S. presidential elections. The 2020 study features a two-wave panel design with pre-election and post-election interviews.

The ANES 2020 questionnaires consist of two parts: the Screener Questionnaire and the Survey Questionnaire. The Survey Questionnaire is further divided into pre-election and post-election sections. It primarily consists of multiple-choice questions, with a few multiple-answer and fill-in-the-blank questions. The questionnaire covers a broad range of topics, with the approximate percentage of each topic indicated in Table 5.

Percentage	Topic
10%	Voting behavior
7%	Candidate evaluations
3%	Party evaluations
12%	Evaluations of government and politics
13%	Demographics
7%	Personal experience and outlook
9%	Political engagement
4%	Predispositions (traits, values, etc.)
13%	Group identities and attitudes
19%	Political issues
3%	Other

Table 5: Approximate percentage of each topic in 2020 ANES survey questionnaire

The questionnaire is organized in modules, with each cohering on one or more dimensions of format, mode, or subject matter. These modules are further subdivided into sections and questions. Figure 6 outlines the modules and their component sections for the 2020 pre-election study.

A total of 8,280 respondents participated in the pre-election interviews, of whom 7,453 also completed the post-election reinterviews, while 827 did not respond. Most questions in the questionnaire are logically related, and depending on the respondents’ answers, follow-up questions are asked accordingly. As a result, the number of questions asked varied for each respondent. Table 6 provides a summary of the 2020 ANES survey questionnaire.

Metric	Value	
	Pre-election	Post-election
Number of Questions	443	511
Number of Respondents	8280	7453
Time Span	Aug 18, 2020 - Nov 2, 2020	Nov 8, 2020 - Jan 4, 2021
Avg. Number of Questions Answered per Respondent	341.68	349.83

Table 6: 2020 ANES survey questionnaire summary.

### 3.3 Poll-based Questionnaire Design

The questionnaire used for the massive social simulation is modified based on the ANES data to better suit the needs of simulating real voters through social agents. Using the 2020 ANES questionnaire as a prototype, we design our questionnaire following these steps:

<ul style="list-style-type: none"> <li>● Pre-election start <ul style="list-style-type: none"> <li>○ Survey start</li> <li>○ R has booklet</li> <li>○ Survey consent</li> <li>○ Self-reported sex</li> </ul> </li> <li>● Engagement <ul style="list-style-type: none"> <li>○ Interest in campaigns</li> </ul> </li> <li>● Pre vote <ul style="list-style-type: none"> <li>○ Pre-election voting module</li> <li>○ Likely to vote</li> <li>○ Voting in prior election</li> </ul> </li> <li>● Attitudes &amp; Candidates <ul style="list-style-type: none"> <li>○ Candidate likes &amp; dislikes</li> <li>○ Emotions</li> <li>○ Congress approval</li> <li>○ Presidential approval</li> <li>○ Covid policy</li> <li>○ Feeling therm.</li> <li>○ Party likes-dislikes</li> <li>○ Liberal-cons place't</li> <li>○ Candidate traits</li> <li>○ Election expectations to win</li> </ul> </li> <li>● Government &amp; Parties <ul style="list-style-type: none"> <li>○ Divided government</li> <li>○ Party ID</li> <li>○ Trust in government</li> <li>○ Social trust</li> <li>○ Elections make govt responsive</li> <li>○ Party performance</li> </ul> </li> <li>● Issues 1 (longer time series) <ul style="list-style-type: none"> <li>○ Services &amp; spending</li> <li>○ Defense spending</li> <li>○ Health insurance</li> <li>○ Jobs/standard of living</li> <li>○ Aid to Blacks</li> <li>○ Enviro-business tradeoff</li> <li>○ Federal spending</li> <li>○ Economic performance</li> <li>○ Abortion</li> <li>○ Death penalty</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>○ US position in world</li> <li>● Issues 2 (timely) <ul style="list-style-type: none"> <li>○ Election integrity</li> <li>○ Democratic norms</li> <li>○ Compromise</li> <li>○ Trump issues</li> <li>○ Covid-19 response</li> <li>○ Inequality</li> <li>○ Climate change</li> <li>○ Parental leave</li> <li>○ Services to same-sex couples</li> <li>○ Transgender policy</li> <li>○ Gay rights</li> <li>○ Immigration</li> <li>○ Speaking English</li> <li>○ Russia interference</li> <li>○ Unrest</li> </ul> </li> <li>● Religion <ul style="list-style-type: none"> <li>○ Religion</li> </ul> </li> <li>● Demographics <ul style="list-style-type: none"> <li>○ Main demographics</li> </ul> </li> <li>● Demographics <ul style="list-style-type: none"> <li>○ Demographics 2</li> <li>○ Demographics 3</li> <li>○ Economic peril</li> </ul> </li> <li>● CASI <ul style="list-style-type: none"> <li>○ Sexual orientation</li> <li>○ Political violence</li> <li>○ Wealth</li> <li>○ Family income</li> <li>○ Mental health</li> <li>○ Health</li> <li>○ Pol. correctness</li> <li>○ Gun ownership</li> <li>○ Media sources</li> <li>○ Gender resentment</li> <li>○ Political knowledge</li> <li>○ Interview ratings</li> <li>○ Life satisfaction</li> </ul> </li> </ul>
--	---

Figure 6: Pre-election questionnaire modules and sections: ANES 2020 Time Series Study

**Selecting Topics** We select a total of 24 socially significant topics and collect questions that reflect respondents' opinions on these topics. Each topic is associated with a varying number of specific questions, ranging from 1 to 6. Table 7 presents a list of the topics and the number of questions corresponding to each.

Number of questions	Topic
≥ 5	Democratic Norms, Immigration, LGBTQ+ Rights
4	Environment, Government
2	Criminal Justice, Education, Gender Resentment, Health Care, Social Welfare, US Position in World
1	Abortion, Aid to Blacks, Aid to Poor, Defense, Economy, Election Integrity, Inequality, Infrastructure, Parental Leave, Social Security, Taxes, Unrest, Voting Behavior

Table 7: Selected questionnaire topics.

**Optimizing Question Selection** To reduce the complexity of the questionnaire while retaining as much original information as possible, we implement the following steps:

1. Removing multiple-answer questions and fill-in-the-blank questions.
2. Eliminating conditional questions, i.e., those asked only to respondents who choose specific options in previous questions, ensuring that the remaining questions apply to all respondents.

**Merging Invalid Options** We define invalid options as those like "Refused", "Don't know", or "Haven't thought much about this", which do not provide clear responses. Having too many such options in the questionnaire may interfere with the agent's performance in answering questions. We merge these invalid options into a single new category: "DK/RF."

**Converting Intensity Questions to Position Questions** In the original ANES questionnaire, some questions are overly detailed, asking respondents to specify the intensity of their stance, such as distinguishing between "Agree strongly" and "Agree somewhat." This level of detail imposes overly high demands on agents to align with human opinions. In redesigning the questionnaire, we ignore the intensity of opinions and only ask about the respondents' basic stance. We also convert the original 7-point Likert scale to a 3-point scale. This reduction in the number of options makes it easier for the agents to make choices. Figure 7 gives two examples of revised questions.

Original Question	Revised Question
<p>Some people feel the government in Washington should see to it that every person has a job and a good standard of living. Suppose these people are at one end of a scale, at point 1. Others think the government should just let each person get ahead on their own. Suppose these people are at the other end, at point 7. And, of course, some other people have opinions somewhere in between, at points 2, 3, 4, 5, or 6.</p> <p>Where would you place yourself on this scale, or haven't you thought much about this?</p> <p>-9. Refused -8. Don't know 1. Government should see to jobs and standard of living 2. 3. 4. 5. 6. 7. Government should let each person get ahead on own 99. Haven't thought much about this</p>	<p>Some people feel the government in Washington should see to it that every person has a job and a good standard of living. Others think the government should just let each person get ahead on their own. And, of course, some people have a neutral position.</p> <p>Which of the following best describes your view?</p> <p>-2. DK/RF 1. Government should see to jobs and standard of living 2. Neutral 3. Government should let each person get ahead on own</p>
<p>Which party do you think would do a better job of handling the nation's economy?</p> <p>-9. Refused -8. Don't know 1. Democrats would do a much better job 2. Democrats would do a somewhat better job 3. Not much difference between them 4. Republicans would do a somewhat better job 5. Republicans would do a much better job</p>	<p>Which party do you think would do a better job of handling the nation's economy?</p> <p>-2. DK/RF 1. Democrats would do a better job 2. Not much difference between them 3. Republicans would do a better job</p>

Figure 7: Revised Question Example.

In the end, we finalize 49 questions across 24 issues. Table 8 gives some general information about the questionnaire we designed. Most questions have 3 to 5 options, allowing for clear distinctions between different positions without making the choices too ambiguous. The full questionnaire can be found in Appendix E.

### 3.4 Prompt-based Evaluation Baselines

We design three prompt-based evaluation baselines to evaluate the simulation results comprehensively, as shown in Figure 8.

Metric	Value
Number of Questions	49
Number of Topics	24
Avg. Number of Words per Question	34.06
Avg. Number of Options per Question	3.22
Number of Respondents	8280

Table 8: Questionnaire Summary.

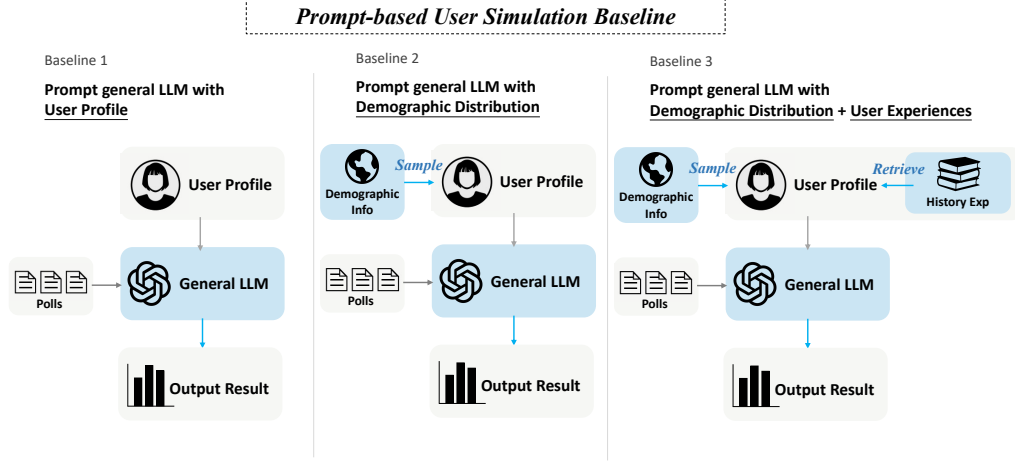


Figure 8: An illustration of prompt-based evaluation baselines.

**Baseline 1: Prompt LLM with user profile** After the user pool construction, one simple way to carry out the social simulation is to randomly sample users from the user pool and interview them with the poll question. The overall distribution of the poll answers can partially reflect the group’s preference towards different candidates and topics.

**Baseline 2: Prompt LLM with demographic distribution** To better impend over the real-world results, we add demographic distribution information during the user preparation. We sample the users according to the census data according to the sampling strategy mentioned in §2.2, which includes the demographic distribution of the voters across different states. The same poll-based interview task is carried out after the sampling.

**Baseline 3: Prompt LLM with multiple information** Though demographic information adds diversity to the group sampling. Different individuals with similar demographic features are hard to distinguish. Baseline 2 may sample randomly from a subset of the user pool, the individuals within which are all "Middle-income graduated Asian males in their 20s". Thus, we add specific users’ experiences, namely posts on social media, to the sampled user profiles to make each individual lively.

## 4 Evaluation

In this section, we display the experiment results carried out under the PPE benchmark at both §4.1 voter level and §4.2 state level.

#### 4.1 Voter-wise Simulation

The first step of Poll-based President Election Prediction is to conduct fine-grained modeling of each voter’s opinion behavior. Therefore, to verify the accuracy of the modeling for an individual, we conduct a voter-wise simulation. The task can be described as: given a series of demographic tag information of an individual, predict his/her response to a specific question.

**Model** We select a series of models for performance comparison, which include both open-source and commercial large language models:

- GPT-4o: A proprietary large-scale language model, known for its advanced capabilities in understanding and generating human-like text.
- GPT-4o-mini: A smaller variant of the GPT-4o model, designed for scenarios with limited computational resources.
- Claude-3.5-Sonnet: A sophisticated language model that excels in poetic and creative text generation.
- Qwen-2-7b-Instruct: An intermediate-sized open-source model, tuned for instruction following and textual response generation.
- Qwen-2-72b-Instruct: A robust open-source model designed for comprehensive language understanding and response generation.
- Qwen-2.5-72b-Instruct: An open-source language model fine-tuned for following instructions and generating coherent responses.
- Llama-3-70b-Instruct: A large language model optimized for understanding and responding to user prompts in a guided manner.

**Metrics** We employ the average Micro-F1 and Macro-F1 scores on opinion poll questions as our evaluation metrics. The Micro-F1 score is utilized to measure the overall fit accuracy of the model. The consideration of Macro-F1 arises from the presence of significant opinion biases in some questions, where the distribution among options is highly uneven. Macro-F1 is capable of measuring the simulation accuracy on small samples, thereby reflecting the expressiveness and simulation precision of minority opinions.

During the evaluation process, we adjust the sample by disregarding individuals who refused to answer (negative options in the questionnaire), and correspondingly adjust the total sample size. The rationale for this adjustment is that real-world opinion polls may fail to conduct interviews due to various unforeseen circumstances. Such outcomes do not accurately reflect the opinion tendencies of the respondents and do not align with the motivation of this study.

**Details** Specifically, we randomly select 1,000 individuals from the 8,280 ANES 2020 respondents as subjects for our study and utilize their responses to demographic-related questions in the questionnaire as known demographic tags<sup>1</sup>. Based on the question introduced in §3.2, we predict and evaluate the accuracy of respondents’ answers to questions covering various topic areas. During the simulation, the model’s max token is 32, and the temperature is 0.5.

**Result** We report the results of the voter-wise simulation conducted in Table 9, which includes both the overall test set and the voting-related subset consisting of 6 questions strongly correlated with voting behavior. The following observations are made:

- The model’s performance on the voting-related subset is generally superior to its performance on the full dataset, with the optimal model achieving a micro-F1 score of over 80% on the voting-related subset.
- Under the same settings, the model’s Macro-F1 score is generally significantly lower than the Micro-F1 score, indicating that there is room for improvement in the model’s simulation of opinions held by minority individuals, suggesting that the model exhibits a certain degree of bias.

---

<sup>1</sup>We consider AGE, GENDER, RACE, INCOME, EDUCATION, AREA, REGION, EMPLOYMENT, MARITAL, RELIGIOUS, PARTY, and IDEOLOGY as known tags to simulate individual behavior.

Type	Model	Overall		Voting Subset	
		Micro-F1	Macro-F1	Micro-F1	Macro-F1
Commercial	GPT-4o	<b>76.16</b>	55.97	<b>81.20</b>	61.03
	GPT-4o-mini	<u>75.39</u>	58.18	80.26	74.72
	Claude-3.5-Sonnet	<u>73.65</u>	58.70	77.52	71.95
Open-source	Qwen2-7b-Instruct	67.53	43.31	76.39	65.65
	Qwen2-72b-Instruct	74.77	<u>58.71</u>	78.39	77.95
	Qwen2.5-72b-Instruct	74.97	57.81	<u>80.41</u>	<b>79.27</b>
	Llama3-70b-Instruct	74.86	<b>59.96</b>	80.16	<u>79.17</u>

Table 9: Model performance on voter-wise simulation. We compared the performance of commercial and open-source LLMs on both the overall test set (**Overall**) and the voting-related subset (**Voting Subset**). The best results are **bolded**, and the second-best results are underlined.

- The performance of the open-source LLMs with a parameter size of 70b is comparable to that of commercial large models. In contrast, LLMs with a 7b size underperformed.
- GPT-4o-mini model achieves relatively high micro-F1 and macro-F1 scores on voting-related subsets, effectively balancing cost, performance, and efficiency. Therefore, unless specifically indicated, simulations involving commercial models all utilize the GPT-4o-mini model.

## 4.2 State-wise Simulation

In this section, we conduct state-wise simulations by integrating the user pool with state-level demographic information. This task can be described as follows: given the demographic information of a state, predict the overall response performance and opinion tendencies of the state’s voters on a series of poll questions.

**Experimental Setting** We conduct a 2020 presidential election simulation for 51 U.S. states. Based on the demographic distributions provided in §2.2 for each state, we perform population sampling on a state-by-state basis. After conducting individual simulations for the sampled populations, we integrate the final results to represent the ultimate opinion distribution for each state. When sampling from the user pool, we determine the sample size for each state based on 1/10,000 of the total population reported in the CENSUS 2020\*\*.

**Comparative Methods** As introduced in §3.4, we propose and validate various simulation methods.

- **Random Sample (Baseline 1):** In this approach, we conduct random sampling from the constructed large user pool. Subsequently, an interview is conducted for each individual, and we aggregate the result of the interview after simulation.
- **Demographic Distribution Guidance (Baseline 2):** In this approach, we integrate the demographic distribution of each state and sample from the large user pool based on this information for each state.
- **User Experience Alignment (Baseline 3):** In this approach, in addition to integrating demographic information, we also incorporate each user’s social media posts from the user pool into the simulation. When integrating posts, we consider temporal information and filter the posts based on the timeline of the simulated event. Specifically, since the simulation involves the 2020 presidential election, to avoid knowledge leakage, social media posts published by users in November 2020 and beyond were excluded.

**Metrics** We propose two metrics, of different granularities, to assess the performance of population simulation in the context of presidential elections:

\*\*<https://www.census.gov/data/tables/2020/dec/2020-apportionment-data.html>

Model	Method	Overall		Battleground	
		CER $\uparrow$	CVS $\downarrow$	CER $\uparrow$	CVS $\downarrow$
Llama3-70b-Instruct	Baseline1	0.510	0.399	0.600	0.386
	Baseline2	0.745	0.118	0.533	0.093
	Baseline3	0.843	0.094	0.733	0.065
Qwen2.5-72b-Instruct	Baseline1	0.510	0.383	0.600	0.370
	Baseline2	0.843	0.078	0.733	0.054
	Baseline3	0.902	0.071	0.733	0.045
	Baseline3*	<b>0.922</b>	<b>0.070</b>	<b>0.800</b>	<b>0.042</b>
GPT-4o-mini	Baseline1	/	/	0.667	0.323
	Baseline2	/	/	<b>0.800</b>	0.052
	Baseline3	/	/	<b>0.800</b>	0.056

Table 10: Model performance on state-wise simulation. We evaluate different methods for their accuracy in forecasting the 2020 U.S. Presidential Election across all 51 states (**Overall**) and 15 battleground states (**Battleground**). The CER measures state-level prediction accuracy, while CVS denotes the RMSE of simulated versus actual vote shares. \*: Building on Qwen2.5-72b’s strong performance on Baseline3, we extend its use to a 1/1000 population sample (around 300,000 agents). This approach effectively predicts outcomes in **47 states** and **12 battleground states**, with reduced RMSE in vote share predictions.

- **Consistency of Election Result (CER):** A coarse-grained measure of the consistency between a state’s simulation result and the actual result. This metric is quantified by calculating the proportion of sample states for which the election simulation results align with the actual result.
- **Consistency of Vote Share (CVS):** A fine-grained metric to assess the consistency between the relative vote share within a state and the actual share. The relative vote share is defined as the ratio of the actual vote percentage for a candidate from one party to the sum of the actual vote percentages for the candidates from both the Democratic and Republican parties [25]. In this metric, we calculate the Root Mean Square Error (RMSE) between the simulated vote share and the actual vote share for each state. Subsequently, we use the average RMSE across all states as CVS.

**Result** We report the results of the state-wise simulation conducted in Table 10. We select GPT-4o-mini, Llama-3-70b-Instruct, and Qwen2.5-72b-Instruct for the simulation and report their performance across 51 states and 15 battleground states (considering cost, we only use GPT-4o-mini for battleground state simulations). We delineate the battleground states and determine the actual vote share based on the results presented on CNN’s presidential election statistics page<sup>††</sup>. We have the following observations:

- During the Baseline2 and Baseline3 simulations, Qwen2.5-72b-Instruct and GPT-4o-mini both achieve good simulation accuracy in battleground states, with all state voting results predicted correctly. Although Llama3-70b-Instruct has poor result predictions in battleground states, its fine-grained metric (CVS) is also significantly lower than that of Baseline1. This demonstrates our framework can model public opinion more accurately than simulation based on a naive sampling method.
- In terms of open-source model performance, the results of Baseline3 are similar to Baseline2’s across different granularity metrics. However, GPT-4o-mini’s performance on Baseline3 is not as good as on Baseline2. A possible reason is the introduction of temporal constraints in Baseline3, whereas in Baseline2 there is some degree of knowledge leakage about the election within the model.

<sup>††</sup><https://www.cnn.com/election/2020/results/president>.

## 5 Further Analysis

### 5.1 Prompt Strategy Study

In the voter-wise simulation task, we endeavor to alter the prompting strategy to observe the simulation effects of LLMs. We experiment with modifying the format of demographic prompts and the format requirements for agent responses to polling questions.

Regarding the demographic prompt format, we explore two approaches: 1) dict, presenting information directly in the prompt in dictionary format; 2) biography, generating autobiographies based on demographic information and then presenting these autobiographies in the prompt.

In terms of the format requirements for responses to polling questions, we also investigate two paradigms: 1) direct, direct answering; 2) reason, providing explanations alongside the answers.

The results presented in Table 11 indicate that the direct-answer prompt configuration achieves optimal simulation performance on both the GPT-4o-mini and Claude-3.5-sonnet models. Compared to directly providing personal information, the autobiographical approach shows insignificant improvements on most metrics, with only a notable enhancement in Macro-F1 for questions strongly related to voting behavior. Considering that the autobiographical processing method may introduce additional hallucinated information, leading to distortion in the agent profile, the main experiment employs the prompt setting that involves direct answers and provides personal information in a dictionary format.

Model	Answer	Personal	Overall		Voting Subset	
	Format	Information	Micro-F1	Macro-F1	Micro-F1	Macro-F1
GPT-4o-mini	direct	dict	<b>75.39</b>	58.18	80.26	74.72
	reason	dict	73.87	52.19	78.75	73.58
	direct	biography	<u>74.64</u>	54.30	79.26	<b>80.85</b>
	reason	biography	<u>74.13</u>	52.98	<u>80.57</u>	64.38
Claude-3.5-Sonnet	direct	dict	73.65	<u>58.70</u>	77.52	71.95
	reason	dict	73.01	54.79	77.37	62.38
	direct	biography	74.19	<b>62.19</b>	<b>80.78</b>	<u>80.48</u>
	reason	biography	72.22	53.27	78.73	63.29

Table 11: Model performance with different prompt strategy. The best results are **bolded**, and the second-best results are underlined.

### 5.2 Ablation Study on Voter-wise Simulation

In our study on the role of core elements in voter-wise simulation, we conduct ablation studies on the temporal information provided to agents during the simulation process and the agents’ ideological labels, which include ideology and political party affiliation.

The results are presented in Table 12. It can be observed that, compared to ideological labels, the removal of temporal information has the least impact on the precision of individual modeling. This is consistent with our hypothesis in the main experiment that the removal of temporal tags could increase the likelihood of model knowledge leakage, leading to "pseudo-predictions" of past events with world knowledge. However, we continue to emphasize the importance of temporal information in group modeling and opinion forecasting, as the macro-level demographic information we integrate is heavily time-dependent. Moreover, we sample user historical statements based on time, which cannot be reflected in the results of individual-level modeling.

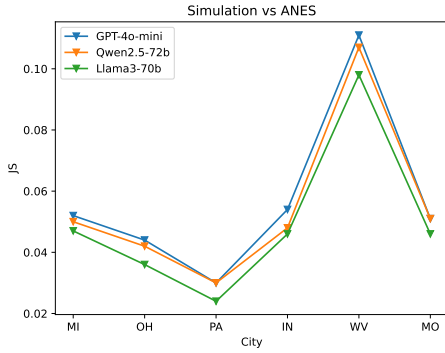
Furthermore, the ideology tag’s influence on the overall test set performance suggests that it encapsulates a particular cognitive framework and pattern that affects how individuals perceive various domain-specific issues. Removing party information, on the other hand, results in more substantial performance variations in subsets that are strongly tied to voting behavior, underscoring the pivotal role of party affiliation in forecasting actions related to voting. This distinction underscores the



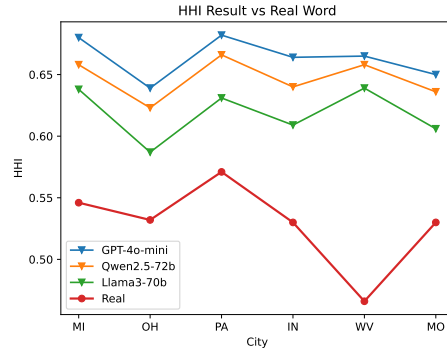
nuanced differences between ideological and party-based influences on voter behavior and opinion formation.

Method	Overall		Voting Subset	
	Micro-F1	Macro-F1	Micro-F1	Macro-F1
Ours	75.39	58.18	80.26	74.72
- time info	75.31	57.26	80.68	79.17
- tag: ideology	74.18	53.48	79.70	72.43
- tag: party	73.48	54.23	74.67	68.63

Table 12: The result of ablation study. Ideology affects overall answer accuracy more, while party affiliation influences voting question precision more.



(a) JS Divergency between models and the actual response distribution of ANES 2020.



(b) HHI of models and actual result. The actual result is in red.

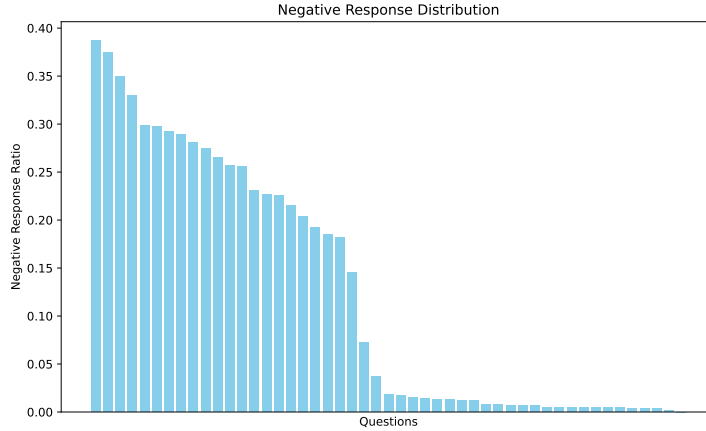


Figure 10: The actual negative response ratio on each question.

### 5.3 Response Behavior Analysis

**Sentiment and Attitude Polarity Analysis** Given the results from previous analyses that highlighted the tendencies of different models, we further examined the behavior of GPT-4o-mini, Qwen2.5-72b-Instruct, and Llama-3-70b-Instruct in the context of opinion poll questionnaires. We characterized their performance using two metrics:

- **JS Divergency (JS):** A metric measures the similarity between two distributions. We derived state-level answer distributions from ANES 2020 as a reference, calculated JS Divergency against simulated results, and used the average as a consistency metric with ANES.

$$JS = \frac{1}{2}D_{kl}(P||Q) + \frac{1}{2}D_{kl}(Q||P), \quad (1)$$

where  $D_{kl}$  is the Kullback-Leibler divergence,  $P$  and  $Q$  are the probability distributions.

- **Herfindahl-Hirschman Index (HHI):** A metric assesses distribution concentration. We computed it for each question, averaging for the model’s HHI score, and also determined the HHI for actual ANES distributions for reference.

$$HHI = \sum_{i=1}^n \left(\frac{s_i}{S}\right)^2, \quad (2)$$

where  $s_i$  is the share of the  $i$ -th dimension, and  $S$  is the total share.

The results are presented in Figure 9a and Figure 9b. Compared to the emotional polarity of ANES, the HHI metric of the model-simulated results is higher, indicating that the models indeed exhibit a more pronounced bias on opinion poll questions than the actual situation. Additionally, although GPT-4o-mini and Qwen2.5-72b-Instruct performed better in state-wise simulation than Llama-3-70b-Instruct, their JS Divergency metric is not as favorable as that of Llama-3-70b-Instruct. This also suggests that the ANES results do not fully reflect the final election outcomes.

**Negative Response Distribution** Due to the sample adjustment strategy employed in our voter-wise simulation metric design, which excluded negative responses, we report in Figure 10 the proportion of negative sample responses to the total for each question in the original data. The results show that the proportion of negative responses to most questions is relatively low, with the highest not exceeding 40%, indicating that the adjusted samples are still representative.

#### 5.4 Case Study

State	Candidates	Ours		ABM		Actual Result	
		Relative Vote Share	Winner	Relative Vote Share	Winner	Relative Vote Share	Winner
MI	Biden-Harris	0.5412	*	0.5454	*	0.5142	*
	Trump-Pence	0.4588		0.4546		0.4858	
OH	Biden-Harris	0.4371		0.4925		0.4589	
	Trump-Pence	0.5629	*	0.5075	*	0.5411	*
PA	Biden-Harris	0.5280	*	0.5204	*	0.5061	*
	Trump-Pence	0.4720		0.4796		0.4939	
IN	Biden-Harris	0.4652		0.4835		0.4184	
	Trump-Pence	0.5348	*	0.5165	*	0.5816	*
WV	Biden-Harris	0.3811		0.3831		0.3022	
	Trump-Pence	0.6189	*	0.6169	*	0.6978	*
MO	Biden-Harris	0.4603		0.4440		0.4216	
	Trump-Pence	0.5397	*	0.5560	*	0.5784	*
RMSE		0.0439		0.0476			

Table 13: Comparison of GPT-4o-mini simulation results in 6 states with the ABM method and actual results. The reporting states are Michigan (MI), Ohio (OH), Pennsylvania (PA), Indiana (IN), West Virginia (WV), and Missouri (MO).

In Table 13, We compared the relative voting shares and outcomes in 6 states for the GPT-4o-mini model against the predicted results of ABM [25], as well as the actual election results. The findings indicate that our model’s predictions align with the real outcomes across all 6 states. However, it is noteworthy that there is a consistent overestimation of the vote share for the Biden-Harris ticket, suggesting a potential bias within the model. This overestimation implies that the model may be

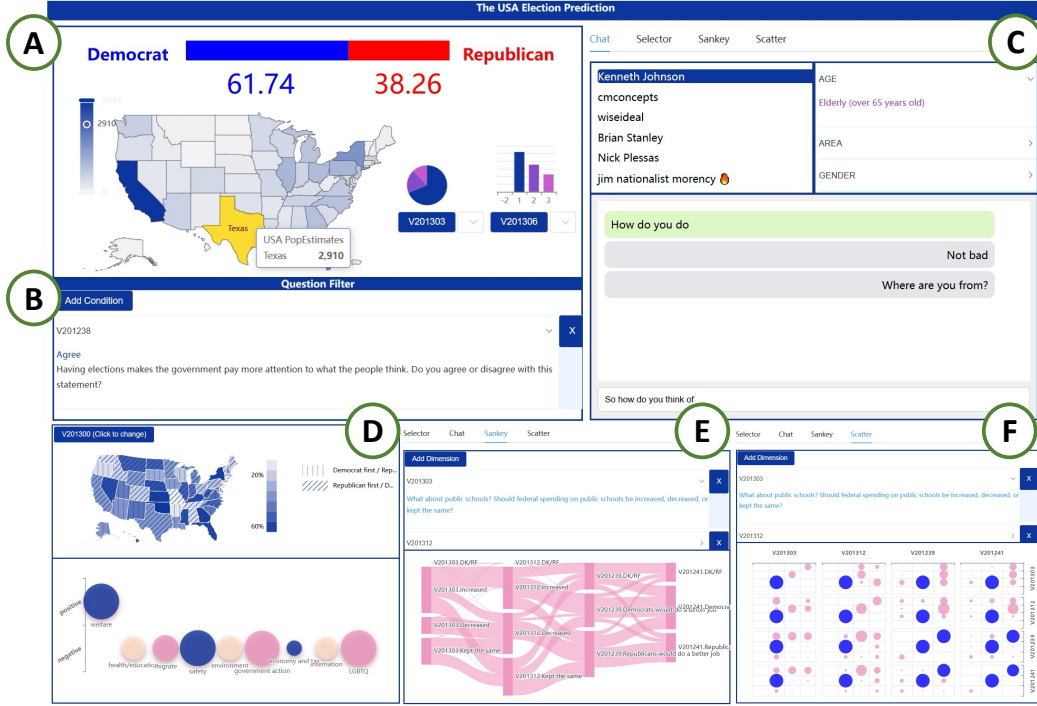


Figure 11: Visualization Interface. (A):Map Filter; (B):Condition Filter; (C):Individual Information; (D):Distribution Overview; (E)(F):High-Dimensional View

influenced by certain biases, which could lead to the incorporation of hallucinated information that distorts the agent’s profile.

## 6 Visualization

This section provides an overview of the design rationale (§6.1), interface design (§6.2), and usage strategies (§6.3) of the visualization work.

### 6.1 Design Rationale

Our dataset includes predictions from all agents across 49 questions, with a large volume of discrete question options and multiple dimensions. Therefore, we have summarized our visualization requirements into the following three requirements:

**R1: Display Macro Trends.** The interface should display distributions of the population across key dimensions such as engagement, support rates, geographic distribution, and specific question distribution.

**R2: Display Micro Characteristics.** When focusing on specific groups, questions, or agents, the interface should show detailed choices for corresponding questions and comprehensive information about individual agents.

**R3: Provide Interactive Features.** The interface should provide sufficient flexibility to allow users to filter specific groups through interactive actions and to display the correlations and trends among these groups across specific dimensions.

### 6.2 Visualization Design

Based on the requirements above, we have designed an interface as shown in Figure 11, which primarily consists of five views: (A)Map Filter, (B)Condition Filter, (C)Individual Information,

(D) Distribution Overview, and 2 (E) (F) high-dimensional view. All views share the same dataset, allowing users to interactively select specific conditions to view the distribution of the selected population across various dimensions. The project is built using npm (10.0.1) and Vue (5.0.8), with dependencies on third-party libraries such as Echarts, Element-plus, D3, and Axios.

**Map Filter** The Map Filter primarily displays graphical distributions, including support rate information, geographical distribution, and custom questions, while also providing a geographic filter function. The top section displays the support rate percentages for the two parties among the currently selected population. The middle section features a map of the United States, where color intensity represents the number of selected populations in each state. Hovering over a state displays its name and the exact population number of that area. Clicking on a state limits the visualization data to those only living in that state, and all other components update to reflect the distribution of the population only living in that area. Clicking the state again restores the national population statistics. The right part of the map shows the distribution of the population in a selected question. Up to 2 distributions are displayed simultaneously, and each chart can toggle between a bar chart and a pie chart. Hovering over a chart displays the option description and the number of individuals who chose it. Hovering over the question code displays the full question text, and clicking on a code would open a dialog for switching to another question.

**Condition Filter** The Condition Filter allows users to filter the visualized data by selecting specific questions and options, thereby narrowing the visualization population of other components to individuals who have chosen particular options for certain questions. Multiple conditions can be applied simultaneously, the condition filter can also be combined with the geographic filters provided by the Map Filter.

**Individual Information** The Individual Information provides more detail about the individual, it displays random 100 voters from the currently selected population in map view and condition filter view. The top-left section lists the names of the voters, while the top-right section shows detailed information about the currently selected voter, including gender, age, and sample contents. The bottom section features a chat component that allows users to engage in natural language conversations with the selected voter's agent to ask more detailed questions. Whenever the visualization population is changed via the map filter or condition filter, 100 new voters are randomly selected from the updated population and displayed in the list.

**Overview Distribution** This component is divided into two parts. The upper part displays the most voted options in each state for a selected question. The option with the highest number of voters for a particular is filled with a corresponding texture, and the shade of color represents the proportion of votes for that option in the state, a deeper color indicates a higher proportion of votes. When the mouse hovers over a specific region, the proportions of other options are displayed. This view differs from Map Filter, which shows absolute numbers. In contrast, this view displays the relative proportions of options within each state. The lower parts display the participation rates and the trends toward radicalism or conservatism of the selected population across various categories. We have categorized all questions into nine major classes, each represented by a floating bubble. If more voters select the "DK/RF" option, the bubble becomes smaller; otherwise, it becomes larger. Additionally, we have scored the level of radicalism or conservatism for each question option. The higher a bubble floats, the more radical the selected population is on that category of issues; conversely, a lower position indicates a more conservative stance.

**High-Dimensional View** The High-Dimensional View utilizes Sankey diagrams and scatterplot matrices to simultaneously display the distribution of options across multiple questions. Users can select the desired dimensions to visualize using buttons at the top of each interface. The system performs permutations and combinations of the selected dimensions, calculates the number of individuals in the current visualization population for each combination, and displays this information in both the scatterplot matrix and the Sankey diagram. Hovering over the corresponding areas in the charts reveals the exact number of individuals. Due to the large number of voters in some states, to ensure system stability and smooth performance, the high-dimensional views perform statistical analyses based on sampled populations.

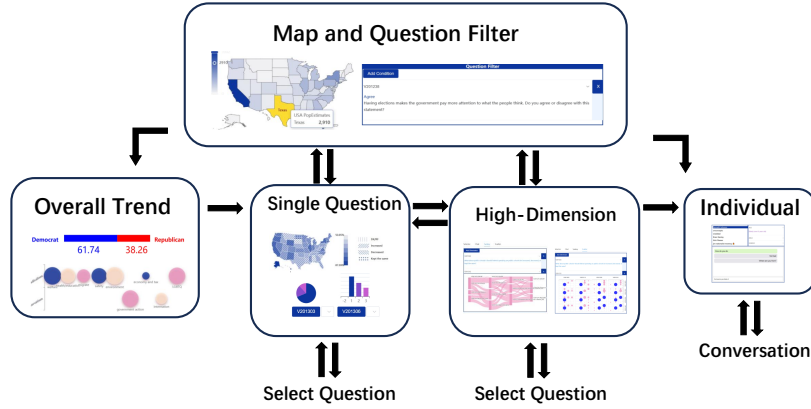


Figure 12: Visualization Workflow

### 6.3 Usage Strategies

Our dataset contains a large volume of discrete question options, so we aim to use visualization to reveal the correlations between these options. The primary approach is to offer multiple visualization methods, enabling users to display overall trends at the national level. As illustrated in Figure 12, users initially view the overall distribution trend at the national level. By applying filters to select different groups, users can refine their view using map filter and question filter to identify specific populations and display their aggregate distribution. Subsequently, users can examine the distribution trends of these specific populations regarding particular questions or categories of questions. This process can be further refined to individual information, where users can engage in dialogues to obtain more detailed insights.

## 7 Related Works

### 7.1 Political Election Research

Traditional election prediction methods mainly rely on opinion polls, expert judgment, and statistical models [20; 6; 10; 22; 19; 8; 11; 7; 64]. Agent-based model (ABM), as an emerging method, provides a more objective and accurate prediction method by simulating individual voter behavior, combining micro-individual characteristics and macro-socioeconomic factors [51; 57]. The ABM method is particularly adept at capturing the diversity and evolving dynamics within actual voting contexts, providing a higher level of detail and adaptability compared to conventional statistical approaches [25]. In recent years, with the rapid development of LLM, researchers have discovered its potential to solve problems in the field of political science [35; 27]. Preliminary research has shown positive outcomes in domains including electoral prediction, policy evaluation, and the simulation of public sentiment [54; 41].

### 7.2 Multi-Agent Simulation by Large Language Models

Agent-based simulations by LLMs have gained wide attention recently for their promising application value and possibility that may shed light on solving general problems paradigm [65; 28; 24]. While individual-level simulation (also known as role-playing agents) focuses on highly reliable and reproducible human-like behavior [56; 63; 66; 58], multi-agent simulation pays more attention to the collaboration and interaction mechanism between agents and the overall achievement of specific tasks and events [21; 30; 50; 23; 34]. Multi-agent simulations also vary depending on different scenarios, wherein general-purpose scenarios highlight the intelligence within LLMs [46; 67; 42] while specific-domain scenarios emphasize the combination between workflows and domain specialization, like journalism [36; 33], economy [31; 70], social media [9; 45; 37; 69], etc.

## 8 Conclusions

In this paper, we introduce the **ElectionSim**, a multi-agent large-scale election simulation framework supported by a million-level voter pool and customized distribution sampling strategy. We apply the ElectionSim to the U.S. presidential election under the poll-based presidential election benchmark, namely **PPE**, and achieve high-accuracy simulation results on both the voter level and the state level. Further analyses, including prompt, ablation, and case study have fully demonstrated the robustness and effectiveness of ElectionSim.

## References

- [1] Medina, Igielnik and Ulloa. <https://www.nytimes.com/2024/10/13/us/politics/latinos-trump-harris-poll.html>. Accessed: 2024-10-13.
- [2] American National Election Studies. Anes 2020 time series study full release [dataset and documentation], 2021. February 10, 2022 version.
- [3] L. M. Bartels. Uninformed votes: Information effects in presidential elections. *American journal of political science*, pages 194–230, 1996.
- [4] L. M. Bartels. Panel effects in the american national election studies. *Political Analysis*, 8(1):1–20, 1999.
- [5] I. Beltagy, M. E. Peters, and A. Cohan. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*, 2020.
- [6] D. Black. On the rationale of group decision-making. *Journal of political economy*, 56(1):23–34, 1948.
- [7] J. Bohannon. The pulse of the people, 2017.
- [8] P. Burnap, R. Gibson, L. Sloan, R. Southern, and M. Williams. 140 characters to victory?: Using twitter to predict the uk 2015 general election. *Electoral Studies*, 41:230–233, 2016.
- [9] J. Cai, J. Li, M. Zhang, M. Li, C.-S. Wang, and K. Tei. Language evolution for evading social media regulation via llm-based multi-agent simulation. *arXiv preprint arXiv:2405.02858*, 2024.
- [10] J. E. Campbell and T. E. Mann. Forecasting the 1992 presidential election: A user’s guide to the models. *The Brookings Review*, 10(4):22–27, 1992.
- [11] A. Ceron, L. Curini, and S. M. Iacus. *Politics and big data: Nowcasting and forecasting elections with social media*. Routledge, 2016.
- [12] T. Choi and S. Park. Theory building via agent-based modeling in public administration research: Vindications and limitations. *International Journal of Public Sector Management*, 34(6):614–629, 2021.
- [13] A. Chopra, J. Subramanian, B. Krishnamurthy, and R. Raskar. Agenttorch: Agent-based modeling with automatic differentiation. In *Second Agent Learning in Open-Endedness Workshop*.
- [14] A.-A. Choupani and A. R. Mamdoohi. Population synthesis using iterative proportional fitting (ipf): A review and future research. *Transportation Research Procedia*, 17:223–233, 2016.
- [15] Y.-S. Chuang and T. T. Rogers. Computational agent-based models in opinion dynamics: A survey on social simulations and empirical studies. *arXiv preprint arXiv:2306.03446*, 2023.
- [16] M. Cinelli, G. De Francisci Morales, A. Galeazzi, W. Quattrociocchi, and M. Starnini. The echo chamber effect on social media. *Proceedings of the National Academy of Sciences*, 118(9):e2023301118, 2021.
- [17] A. Crooks, C. Castle, and M. Batty. Key challenges in agent-based modelling for geo-spatial simulation. *Computers, Environment and Urban Systems*, 32(6):417–430, 2008.
- [18] M. DeBell, M. Amsbary, T. Brader, S. Brock, C. Good, J. Kamens, N. Maisel, and S. Pinto. Methodology report for the anes 2020 time series study. 2022.
- [19] J. DiGrazia, K. McKelvey, J. Bollen, and F. Rojas. More tweets, more votes: Social media as a quantitative indicator of political behavior. *PloS one*, 8(11):e79449, 2013.
- [20] A. Downs. An economic theory of political action in a democracy. *Journal of political economy*, 65(2):135–150, 1957.
- [21] Y. Du, S. Li, A. Torralba, J. B. Tenenbaum, and I. Mordatch. Improving factuality and reasoning in language models through multiagent debate. *arXiv preprint arXiv:2305.14325*, 2023.
- [22] R. S. Erikson and C. Wlezien. Forecasting us presidential elections using economic and noneconomic fundamentals. *PS: Political Science & Politics*, 47(2):313–316, 2014.
- [23] Z. Fan, J. Tang, W. Chen, S. Wang, Z. Wei, J. Xi, F. Huang, and J. Zhou. Ai hospital: Interactive evaluation and collaboration of llms as intern doctors for clinical diagnosis. *arXiv preprint arXiv:2402.09742*, 2024.
- [24] C. Gao, X. Lan, N. Li, Y. Yuan, J. Ding, Z. Zhou, F. Xu, and Y. Li. Large language models empowered agent-based modeling and simulation: A survey and perspectives. *Humanities and Social Sciences Communications*, 11(1):1–24, 2024.

- [25] M. Gao, Z. Wang, K. Wang, C. Liu, and S. Tang. Forecasting elections with agent-based modeling: Two live experiments. *Plos one*, 17(6):e0270194, 2022.
- [26] S. Giorgi, V. E. Lynn, K. Gupta, F. Ahmed, S. Matz, L. H. Ungar, and H. A. Schwartz. Correcting sociodemographic selection biases for population prediction from social media. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 16, pages 228–240, 2022.
- [27] P. Gujral, K. Awaldhi, N. Jain, B. Bhandula, and A. Chakraborty. Can llms help predict elections?(counter) evidence from the world’s largest democracy. *arXiv preprint arXiv:2405.07828*, 2024.
- [28] T. Guo, X. Chen, Y. Wang, R. Chang, S. Pei, N. V. Chawla, O. Wiest, and X. Zhang. Large language model based multi-agents: A survey of progress and challenges. *arXiv preprint arXiv:2402.01680*, 2024.
- [29] J. Hoey, T. Schröder, J. Morgan, K. B. Rogers, D. Rishi, and M. Nagappan. Artificial intelligence and social simulation: Studying group dynamics on a massive scale. *Small Group Research*, 49(6):647–683, 2018.
- [30] S. Hong, X. Zheng, J. Chen, Y. Cheng, J. Wang, C. Zhang, Z. Wang, S. K. S. Yau, Z. Lin, L. Zhou, et al. Metagpt: Meta programming for multi-agent collaborative framework. *arXiv preprint arXiv:2308.00352*, 2023.
- [31] J. J. Horton. Large language models as simulated economic agents: What can we learn from homo silicus? Technical report, National Bureau of Economic Research, 2023.
- [32] M. Jusup, P. Holme, K. Kanazawa, M. Takayasu, I. Romić, Z. Wang, S. Geček, T. Lipić, B. Podobnik, L. Wang, et al. Social physics. *Physics Reports*, 948:1–148, 2022.
- [33] X. Li, Y. Xu, Y. Zhang, and E. C. Malthouse. Large language model-driven multi-agent simulation for news diffusion under different network structures. *arXiv preprint arXiv:2410.13909*, 2024.
- [34] J. Liang, R. Ye, M. Han, R. Lai, X. Zhang, X. Huang, and Z. Wei. Debatrix: Multi-dimensinal debate judge with iterative chronological analysis based on llm. *arXiv preprint arXiv:2403.08010*, 2024.
- [35] M. Linegar, R. Kocielnik, and R. M. Alvarez. Large language models and political science. *Frontiers in Political Science*, 5:1257092, 2023.
- [36] X. Liu, S. Yang, X. Zhang, H. Kuang, L. Sun, Y. Yang, S. Chen, X. Huang, and Z. Wei. Ai-press: A multi-agent news generating and feedback simulation system powered by large language models. *arXiv preprint arXiv:2410.07561*, 2024.
- [37] H. Lyu, J. Huang, D. Zhang, Y. Yu, X. Mou, J. Pan, Z. Yang, Z. Wei, and J. Luo. Gpt-4v (ision) as a social media analysis engine. *arXiv preprint arXiv:2311.07547*, 2023.
- [38] C. M. Macal and M. J. North. Agent-based modeling and simulation. In *Proceedings of the 2009 winter simulation conference (WSC)*, pages 86–98. IEEE, 2009.
- [39] B. Major, A. Blodorn, and G. Major Blascovich. The threat of increasing diversity: Why many white americans support trump in the 2016 presidential election. *Group Processes & Intergroup Relations*, 21(6):931–940, 2018.
- [40] W. E. Miller, S. A. Traugott, and S. Traugott. *American national election studies data sourcebook, 1952-1986*. Harvard University Press, 1989.
- [41] F. Moghimifar, Y.-F. Li, R. Thomson, and G. Haffari. Modelling political coalition negotiations using llm-based agents. *arXiv preprint arXiv:2402.11712*, 2024.
- [42] X. Mou, J. Liang, J. Lin, X. Zhang, X. Liu, S. Yang, R. Ye, L. Chen, H. Kuang, X. Huang, and Z. Wei. Agentsense: Benchmarking social intelligence of language agents through interactive scenarios, 2024.
- [43] X. Mou, Z. Wei, and X. Huang. Unveiling the truth and facilitating change: Towards agent-based large-scale social movement simulation. *arXiv preprint arXiv:2402.16333*, 2024.
- [44] G. Murić, A. Tregubov, J. Blythe, A. Abeliuk, D. Choudhary, K. Lerman, and E. Ferrara. Large-scale agent-based simulations of online social networks. *Autonomous Agents and Multi-Agent Systems*, 36(2):38, 2022.
- [45] M. Papachristou and Y. Yuan. Network formation and dynamics among multi-llms. *arXiv preprint arXiv:2402.10659*, 2024.



- [46] J. S. Park, J. O'Brien, C. J. Cai, M. R. Morris, P. Liang, and M. S. Bernstein. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th annual acm symposium on user interface software and technology*, pages 1–22, 2023.
- [47] E. O. Pérez and M. Tavits. Language shapes people's time perspective and support for future-oriented policies. *American Journal of Political Science*, 61(3):715–727, 2017.
- [48] E. O. Pérez and M. Tavits. Language influences public attitudes toward gender equality. *The Journal of Politics*, 81(1):81–93, 2019.
- [49] W. Qi, H. Lyu, and J. Luo. Representation bias in political sample simulations with large language models. *arXiv preprint arXiv:2407.11409*, 2024.
- [50] C. Qian, W. Liu, H. Liu, N. Chen, Y. Dang, J. Li, C. Yang, W. Chen, Y. Su, X. Cong, et al. Chatdev: Communicative agents for software development. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15174–15186, 2024.
- [51] L. Qiu and R. Phang. Agent-based modeling in political decision making, 02 2020.
- [52] F. Ribeiro, L. Henrique, F. Benevenuto, A. Chakraborty, J. Kulshrestha, M. Babaei, and K. Gummadi. Media bias monitor: Quantifying biases of social media news outlets at large-scale. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 12, 2018.
- [53] S. J. Rosenstone. Forecasting presidential elections. 1981.
- [54] D. Rozado. The political preferences of llms. *arXiv preprint arXiv:2402.01789*, 2024.
- [55] T. C. Schelling. Models of segregation. *The American economic review*, 59(2):488–493, 1969.
- [56] Y. Shao, L. Li, J. Dai, and X. Qiu. Character-llm: A trainable agent for role-playing. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 13153–13187, 2023.
- [57] P. Sobkowicz. Quantitative agent based model of opinion dynamics: Polish elections of 2015. *PloS one*, 11(5):e0155098, 2016.
- [58] L. Sun, S. Wang, X. Huang, and Z. Wei. Identity-driven hierarchical role-playing agents. *arXiv preprint arXiv:2407.19412*, 2024.
- [59] S. Tang. Idea, action, and outcome. *Innovation in the Social Sciences*, 2(2):123–170, 2024.
- [60] R. A. Teixeira. *Red, blue, and purple America: the future of election demographics*. Rowman & Littlefield, 2009.
- [61] E. Vraga. Party differences in political content on social media. *Online Information Review*, 40(5):595–609, 2016.
- [62] W. Wang, D. Rothschild, S. Goel, and A. Gelman. Forecasting elections with non-representative polls. *International Journal of Forecasting*, 31(3):980–991, 2015.
- [63] Z. M. Wang, Z. Peng, H. Que, J. Liu, W. Zhou, Y. Wu, H. Guo, R. Gan, Z. Ni, J. Yang, et al. Rolellm: Benchmarking, eliciting, and enhancing role-playing abilities of large language models. *arXiv preprint arXiv:2310.00746*, 2023.
- [64] Z. Wei, Y. He, W. Gao, B. Li, L. Zhou, and K.-f. Wong. Mainstream media behavior analysis on twitter: a case study on uk general election. In *Proceedings of the 24th ACM Conference on Hypertext and Social Media*, pages 174–178, 2013.
- [65] Z. Xi, W. Chen, X. Guo, W. He, Y. Ding, B. Hong, M. Zhang, J. Wang, S. Jin, E. Zhou, et al. The rise and potential of large language model based agents: A survey. *arXiv preprint arXiv:2309.07864*, 2023.
- [66] C. Xie, C. Chen, F. Jia, Z. Ye, K. Shu, A. Bibi, Z. Hu, P. Torr, B. Ghanem, and G. Li. Can large language model agents simulate human trust behaviors? *arXiv preprint arXiv:2402.04559*, 2024.
- [67] S. Yue, S. Wang, W. Chen, X. Huang, and Z. Wei. Synergistic multi-agent framework with trajectory learning for knowledge-intensive tasks. *arXiv preprint arXiv:2407.09893*, 2024.
- [68] N. Yusuf, N. Al-Banawi, and H. A. R. Al-Imam. The social media as echo chamber: The digital impact. *Journal of Business & Economics Research (Online)*, 12(1):1, 2014.

- [69] X. Zhang, H. Kuang, X. Mou, H. Lyu, K. Wu, S. Chen, J. Luo, X. Huang, and Z. Wei. SoMeLVLM: A large vision language model for social media processing. In L.-W. Ku, A. Martins, and V. Srikumar, editors, *Findings of the Association for Computational Linguistics ACL 2024*, pages 2366–2389, Bangkok, Thailand and virtual meeting, Aug. 2024. Association for Computational Linguistics.
- [70] Q. Zhao, J. Wang, Y. Zhang, Y. Jin, K. Zhu, H. Chen, and X. Xie. Competeai: Understanding the competition behaviors in large language model-based agents. *arXiv preprint arXiv:2310.17512*, 2023.

## A Simulation Result of the 2024 Presidential Election

Based on the existing data, we conduct a simulation for the 2024 election. Since the ANES survey results for 2024 and the US Census results are not yet public, and considering that there will not be drastic changes in the population structure of each state within 2-4 years, we temporarily fit the demographic distribution information of each state using the ANES survey results from 2020 and the US Census statistics from 2022.

We report the forecast results in Figure 13 and the results of 15 battleground states in Table 14. The simulation results show that the Democratic Party led by Harris will win 8 of the 15 battleground states, and the Republican Party led by Trump will win 7 of them. According to the simulation results, the Democratic Party led by Harris has a certain advantage in the election.

**We emphasize that we have no intention of influencing actual election activities. This simulation is only for academic research and discussion. The predictions and viewpoints included in this study are for informational purposes only and do not represent the position of the authors or the research team. These predictions should not be interpreted as definitive forecasts or guarantees.**

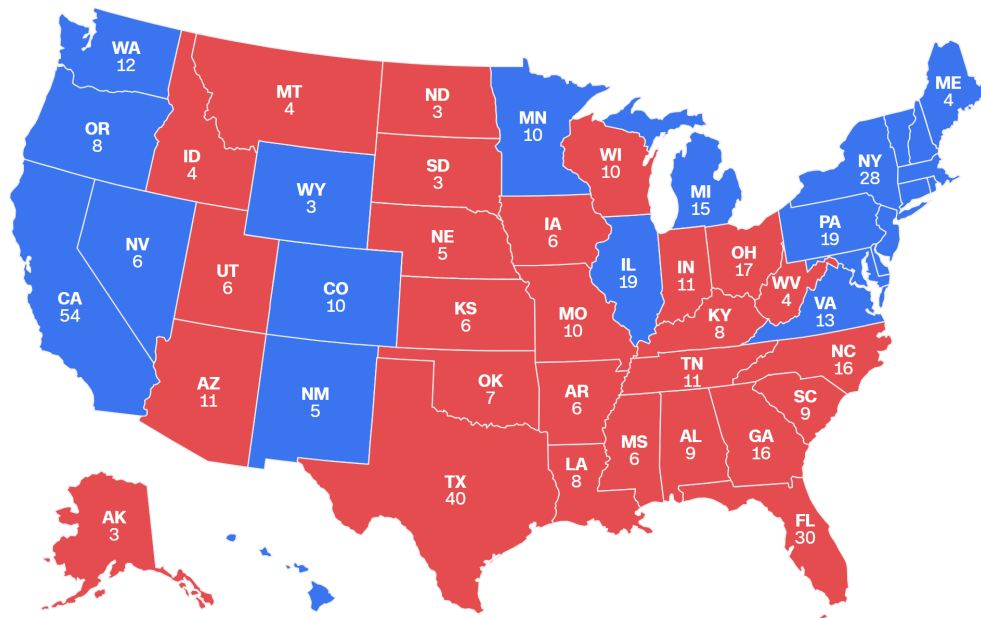


Figure 13: The forecast results for the 2024 U.S. Presidential Election. Red indicates the state won by the Republican Party, and blue indicates the state won by the Democratic Party.

State	Candidates	Prediction Result	
		Relative Vote Share	Winner
Arizona	Harris-Walz	0.3931	
	Trump-Vance	0.6069	★
Colorado	Harris-Walz	0.5518	★
	Trump-Vance	0.4482	
Florida	Harris-Walz	0.4086	
	Trump-Vance	0.5914	★
Georgia	Harris-Walz	0.4518	
	Trump-Vance	0.5482	★
Iowa	Harris-Walz	0.4685	
	Trump-Vance	0.5315	★
Michigan	Harris-Walz	0.5324	★
	Trump-Vance	0.4676	
Minnesota	Harris-Walz	0.5508	★
	Trump-Vance	0.4492	
Nevada	Harris-Walz	0.6156	★
	Trump-Vance	0.3844	
New Hampshire	Harris-Walz	0.6638	★
	Trump-Vance	0.3362	
North Carolina	Harris-Walz	0.4729	
	Trump-Vance	0.5271	★
Ohio	Harris-Walz	0.4163	
	Trump-Vance	0.5837	★
Pennsylvania	Harris-Walz	0.5320	★
	Trump-Vance	0.4680	
Texas	Harris-Walz	0.3683	
	Trump-Vance	0.6317	★
Virginia	Harris-Walz	0.6125	★
	Trump-Vance	0.3875	
Wisconsin	Harris-Walz	0.4563	
	Trump-Vance	0.5437	★

Table 14: Simulation results for the 2024 presidential election in 15 battleground states.

## B Prompt Lib

### B.1 Prompt for Demographic Annotation

**Instruction:** You are a professional annotator tasked with evaluating the attributes of a person based on their entire history of speeches. Your role is to assess the person holistically, taking into account all the provided speeches together, rather than evaluating each speech individually. Below is some historical speech information about this person:

**Text:** {text}

Now, please classify the following attributes of the person:

1. Age Group
  - (a) Youth (18-35 years old)
  - (b) Middle-aged (36-65 years old)
  - (c) Elderly (over 65 years old)
2. Gender
  - (a) Male
  - (b) Female
3. Race
  - (a) White
  - (b) Black
  - (c) Asian
  - (d) Hispanic
4. Party Affiliation
  - (a) Democratic Party
  - (b) Republican Party
  - (c) Other Party
  - (d) Independent
5. Ideology
  - (a) Liberal
  - (b) Moderate
  - (c) Conservative

Please provide your answers in the following JSON format for each attribute:

```
““json
{ "AGE": "A", "GENDER": "B", "RACE": "C", "PARTY": "B", "IDEOLOGY": "C" }
““
```

### B.2 Prompt for Voter-wise Simulation

**Instruction:** It's 2020, and you're being surveyed for the 2020 American National Election Studies. You are a real person living in {state} with the following personal information. Please answer the following question as best as you can. You should act consistently with the role you are playing. Do not select the option to refuse to answer.

**Personal information:** {personal info}

**Question:** {question}

**Options:** {options}

You should give your answer (you only need to answer the option letter number) in JSON format as example below:

```
““json
{ "answer": "xxx" }
““
```

### B.3 Prompt for State-wise Simulation

**Instruction:** It's 2020, and you're being surveyed for the 2020 American National Election Studies. You are a real person living in {state} with the following personal information. Please answer the following question as best as you can. You should act consistently with the role you are playing. Do not select the option to refuse to answer.

Some of your historical comments on social media platforms: {historical comments}

**Personal information:** {personal info}

**Candidates Information:** In the 2020 United States presidential election, the Republican ticket is led by incumbent President Donald Trump, who is known for his assertive communication style and strict immigration policies. Trump is focusing on economic management and a tough stance on law and order, reflecting his commitment to his "America First" approach. His running mate is Vice President Mike Pence. On the Democratic side, former Vice President Joe Biden is the nominee, with Senator Kamala Harris from California as his running mate. Harris is the first African-American, first Asian-American, and third female vice presidential nominee on a major party ticket. Biden's campaign emphasizes unity and healing, with a focus on addressing the public health and economic impacts of the ongoing COVID-19 pandemic, civil unrest following the killing of George Floyd, the future of the Affordable Care Act, and the composition of the U.S. Supreme Court.

**Question:** {question}

**Options:** {options}

You should give your answer (you only need to answer the option letter number) in JSON format as example below:

```
“{
  "answer": "xxx"
}”
```

### B.4 Prompt for Generating Personal Biography

**Instruction:** You are a very outstanding biographer. Now there is some information about a person. Please generate a description of his past experiences based on this information. Please return to this biography in the second person, with the sentence structure of "You are xxx".

**Personal information:** {personal info}

You should give your answer and reason in JSON format as below:

```
“{
  "answer": "xxx"
}”
```

## C Demographic Feature Annotation

### C.1 Language Filtering

We utilize *langid*<sup>‡‡</sup> to detect the languages of user posts. However, it’s important to note that the English posts may still include views from users in other English-speaking regions, like the UK.

### C.2 Repeatability Calculation

We calculate the repeatability score with the following steps:

1. **Sampling:** Let  $P = \{p_1, p_2, p_3, p_4, p_5\}$  be a set of five sampled posts from a user’s historical tweets.
2. **Jaccard Similarity:** For any two posts  $p_i$  and  $p_j$ , define the Jaccard similarity  $J(p_i, p_j)$  as:

$$J(p_i, p_j) = \frac{|A_i \cap A_j|}{|A_i \cup A_j|}$$

where  $A_i$  and  $A_j$  are the sets of unique words in posts  $p_i$  and  $p_j$ , respectively. The numerator  $|A_i \cap A_j|$  represents the size of the intersection of the two sets, and the denominator  $|A_i \cup A_j|$  represents the size of their union.

3. **Repeatability Score:** The user post repeatability rate, represented by the mean Jaccard similarity, can be expressed as:

$$\mu_J = \frac{1}{25} \sum_{i=1}^5 \sum_{\substack{j=1 \\ j \neq i}}^5 J(p_i, p_j)$$

We set a threshold of 0.28 to filter out users exceeding this limit, effectively removing spam accounts and advertisements to ensure our analysis reflects authentic interactions.

### C.3 Test Set Construction

#### C.3.1 Implementation Setting

For all commercial APIs, we apply the prompt template B.1 to prompt the model for demographic annotation and clean the annotation results to match predefined answers. We set the temperature as 0 for all commercial APIs.

#### C.3.2 Manual Verification

We conduct manual verification on samples with inconsistent results from the three commercial API models used. We recruit five annotators to follow the instruction B.1 and provide Twitter homepage links corresponding to the inconsistent samples as additional information for further annotation and to establish ground truth results. For each sample, at least two different annotators are responsible for the annotation. If inconsistencies still occur after manual annotation, the sample is deemed unable to obtain the corresponding attribute and is marked as a null value in the relevant dimension. We calculate consistency scores for the annotation results from the five annotators, as shown in Figure 14.

### C.4 Consistency Calculation

For the test set, we obtain ground truth results through API annotation and manual verification. We compare the annotation results of each commercial API with the ground truth results. Let  $N$  represent the total number of samples in the test set, and let  $C$  denote the number of samples with consistent annotations between the API results and the ground truth. We calculate the consistency score  $S$  for each commercial API using the formula:

$$S = \frac{C}{N}$$

---

<sup>‡‡</sup><https://github.com/saffsd/langid.py>

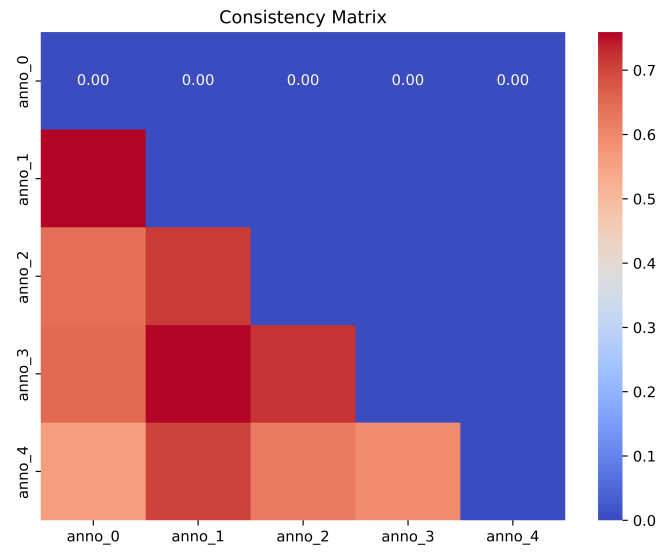


Figure 14: Agreement score among annotators, with an average of 67.23%

This score quantifies the proportion of samples where the API annotations match the established ground truth, providing an assessment of the API's reliability.



## D Supplementary Case Study

Here are a few examples of interactions with agents with specific profiles.

### D.1 Case 1

**Agent A** is a middle-aged white male who lives in an urban area, identifies as conservative, and is affiliated with the Democratic Party.

USER: Who did you vote for?

AGENT A: I voted for Joe Biden.

USER: Why did you vote for Biden?

AGENT A: **I voted for Biden** because I believe in his vision for the country, his experience in government, and his commitment to addressing issues like **healthcare, climate change, and economic inequality**. I also felt that he was the best candidate to bring the country together during a time of significant division.

USER: It is now 2024, Biden has withdrawn from the race, Harris now represents the Democratic Party, and Trump now represents the Republican Party. Who will you vote for now? Why?

AGENT A: **I will vote for Kamala Harris** because I believe in her leadership, her progressive policies, and her commitment to social justice and equality. I also trust her ability to continue the work started by the Biden administration and to address the pressing issues facing our country, such as **healthcare, climate change, and economic recovery**.

USER: If something happened right now that would change your vote, what would it be?

AGENT A: If something happened right now that would change my vote, it would be **a significant policy shift or a major scandal involving Kamala Harris** that I found deeply concerning and incompatible with my values.

From the interview transcript, it is evident that Agent A's voting choices are significantly influenced by their alignment with the candidates' policies and values. Initially, Agent A supported Joe Biden due to his belief in Biden's vision for the country, his experience in governance, and his commitment to tackling issues such as healthcare, climate change, and economic inequality. Additionally, Agent A valued Biden's potential to unite the nation during a period of significant division. Following Biden's withdrawal from the race, Agent A shifted their support to Kamala Harris, citing her leadership, progressive policies, commitment to social justice, and potential to continue the work initiated by the Biden administration while addressing urgent national issues. Notably, Agent A's voting intention is not static; they indicated that a significant policy shift or a major scandal involving Kamala Harris that conflicts with their values could alter their vote. This reflects an attentiveness to the candidates' performance and a strong emphasis on personal values.

## D.2 Case 2

**Agent B** is a middle-aged white suburban female, a staunch conservative, and a supporter of the Republican Party.

USER: Who did you vote for? Why?

AGENT B: **I voted for Donald Trump** because I believe in his conservative values and his commitment to making America strong economically and on the global stage. I also appreciate his stance on reducing regulations and his support for pro-life policies.

USER: Do you think Trump is more conservative than Biden? Why?

AGENT B: Yes. **I believe Trump is more conservative than Biden** because Trump's policies and public stances align more closely with conservative values, such as **lower taxes, reduced government regulation, strong national defense, and support for the Second Amendment**. Biden, on the other hand, tends to support more progressive policies, such as **higher taxes on the wealthy, increased government spending on social programs, and stronger environmental regulations**.

USER: OK. Now it is 2024, and the two candidates representing the two parties are Trump and Harris. Who will you vote for? Why?

AGENT B: I will vote for Donald Trump. I believe he aligns more closely with my conservative values and has a proven track record of supporting policies that benefit the economy and national security. I also trust his leadership and vision for the country.

USER: If you knew that Trump would not make the issues you care about a top political priority, but Harris would, would that change your voting intention?

AGENT B: No, it would not change my voting intention. **I still believe that Trump's overall platform and conservative values align more closely** with my beliefs and the best interests of the country.

Agent B is a staunch conservative supporter who supports the Republican Party and Donald Trump because his policies and positions align with her conservative values, including low taxes, reduced government regulation, a strong national defense, and support for the Second Amendment. Although she knows that in a hypothetical election in 2024, Trump may not make the issues she cares about a priority, while Kamala Harris will, she still chooses to vote for Trump because she believes that Trump's overall political stance and conservative values are more in line with her beliefs and the best interests of the country, and she trusts Trump's leadership and vision for the country's future. This also reflects to some extent that Trump's supporters are more determined.

## E Questionnaire

Here we show the full questionnaire for this experiment.

<b>Q01</b>	<b>Voting Behavior</b>
Question	ORDER OF MAJOR PARTY CANDIDATE NAMES
Value Labels	1. Democrat first / Republican second 2. Republican first / Democrat second
<b>Q02</b>	<b>Social Security</b>
Question	Next I am going to read you a list of federal programs. For each one, I would like you to tell me whether you would like to see spending increased, decreased, or kept the same. What about Social Security? Should federal spending on Social Security be increased, decreased, or kept the same?
Value Labels	-2. DK/RF 1. Increased 2. Decreased 3. Kept the same
<b>Q03</b>	<b>Education</b>
Question	What about public schools? Should federal spending on public schools be increased, decreased, or kept the same?
Value Labels	-2. DK/RF 1. Increased 2. Decreased 3. Kept the same
<b>Q04</b>	<b>Immigration</b>
Question	What about tightening border security to prevent illegal immigration? Should federal spending on tightening border security to prevent illegal immigration be increased, decreased, or kept the same?
Value Labels	-2. DK/RF 1. Increased 2. Decreased 3. Kept the same
<b>Q05</b>	<b>Criminal Justice</b>
Question	What about dealing with crime? Should federal spending on dealing with crime be increased, decreased, or kept the same?
Value Labels	-2. DK/RF 1. Increased 2. Decreased 3. Kept the same
<b>Q06</b>	<b>Social Welfare</b>
Question	What about welfare programs? Should federal spending on welfare programs be increased, decreased, or kept the same?
Value Labels	-2. DK/RF 1. Increased 2. Decreased 3. Kept the same
<b>Q07</b>	<b>Infrastructure</b>
Question	What about building and repairing highways? Should federal spending on building and repairing highways be increased, decreased, or kept the same?

Value Labels	-2. DK/RF 1. Increased 2. Decreased 3. Kept the same
<b>Q08</b>	<b>Aid to Poor</b>
Question	What about aid to the poor? Should federal spending on aid to the poor be increased, decreased, or kept the same?
Value Labels	-2. DK/RF 1. Increased 2. Decreased 3. Kept the same
<b>Q09</b>	<b>Environment</b>
Question	What about protecting the environment? Should federal spending on protecting the environment be increased, decreased, or kept the same?
Value Labels	-2. DK/RF 1. Increased 2. Decreased 3. Kept the same
<b>Q10</b>	<b>Government</b>
Question	How much do you feel that having elections makes the government pay attention to what the people think?
Value Labels	-2. DK/RF 1. A good deal 2. Some 3. Not much
<b>Q11</b>	<b>Economy</b>
Question	Which party do you think would do a better job of handling the nation's economy?
Value Labels	-2. DK/RF 1. Democrats would do a better job 2. Not much difference between them 3. Republicans would do a better job
<b>Q12</b>	<b>Health Care</b>
Question	Which party do you think would do a better job of handling health care?
Value Labels	-2. DK/RF 1. Democrats would do a better job 2. Not much difference between them 3. Republicans would do a better job
<b>Q13</b>	<b>Immigration</b>
Question	Which party do you think would do a better job of handling immigration?
Value Labels	-2. DK/RF 1. Democrats would do a better job 2. Not much difference between them 3. Republicans would do a better job
<b>Q14</b>	<b>Taxes</b>
Question	Which party do you think would do a better job of handling taxes?
Value Labels	-2. DK/RF 1. Democrats would do a better job 2. Not much difference between them 3. Republicans would do a better job
<b>Q15</b>	<b>Environment</b>

Question	Which party do you think would do a better job of handling the environment?
Value Labels	-2. DK/RF 1. Democrats would do a better job 2. Not much difference between them 3. Republicans would do a better job
<b>Q16</b>	<b>Education</b>
Question	Some people think the government should provide fewer services even in areas such as health and education in order to reduce spending. Other people feel it is important for the government to provide many more services even if it means an increase in spending. And, of course, some people have a neutral position. Which of the following best describes your view?
Value Labels	-2. DK/RF 1. Government should provide fewer services 2. Neutral 3. Government should provide more services
<b>Q17</b>	<b>Defense</b>
Question	Some people believe that we should spend less money for defense. Others feel that defense spending should be increased. And, of course, some people have a neutral position. Which of the following best describes your view?
Value Labels	-2. DK/RF 1. Decrease defense spending 2. Neutral 3. Increase defense spending
<b>Q18</b>	<b>Health Care</b>
Question	There is much concern about the rapid rise in medical and hospital costs. Some people feel there should be a government insurance plan which would cover all medical and hospital expenses for everyone. Others feel that all medical expenses should be paid by individuals through private insurance plans like Blue Cross or other company paid plans. And, of course, some people have a neutral position. Which of the following best describes your view?
Value Labels	-2. DK/RF 1. Government insurance plan 2. Neutral 3. Private insurance plan
<b>Q19</b>	<b>Social Welfare</b>
Question	Some people feel the government in Washington should see to it that every person has a job and a good standard of living. Others think the government should just let each person get ahead on their own. And, of course, some people have a neutral position. Which of the following best describes your view?
Value Labels	-2. DK/RF 1. Government should see to jobs and standard of living 2. Neutral 3. Government should let each person get ahead on own
<b>Q20</b>	<b>Aid to Blacks</b>

Question	Some people feel that the government in Washington should make every effort to improve the social and economic position of blacks. Others feel that the government should not make any special effort to help blacks because they should help themselves. And, of course, some people have a neutral position. Which of the following best describes your view?
Value Labels	-2. DK/RF 1. Government should help blacks 2. Neutral 3. Blacks should help themselves
<b>Q21</b>	<b>Environment</b>
Question	Some people think we need much tougher government regulations on business in order to protect the environment. Others think that current regulations to protect the environment are already too much of a burden on business. And, of course, some people have a neutral position. Which of the following best describes your view?
Value Labels	-2. DK/RF 1. Tougher regulations on business needed to protect environment 2. Neutral 3. Regulations to protect environment already too much a burden on business
<b>Q22</b>	<b>Abortion</b>
Question	Would you be pleased, upset, or neither pleased nor upset if the Supreme Court reduced abortion rights?
Value Labels	-2. DK/RF 1. Pleased 2. Upset 3. Neither pleased nor upset
<b>Q23</b>	<b>Criminal Justice</b>
Question	Do you favor or oppose the death penalty for persons convicted of murder?
Value Labels	-2. DK/RF 1. Favor 2. Oppose
<b>Q24</b>	<b>US Position in World</b>
Question	Do you agree or disagree with this statement: 'This country would be better off if we just stayed home and did not concern ourselves with problems in other parts of the world.'
Value Labels	-2. DK/RF 1. Agree 2. Disagree
<b>Q25</b>	<b>US Position in World</b>
Question	How willing should the United States be to use military force to solve international problems?
Value Labels	-2. DK/RF 1. Willing 2. Moderately willing 3. Not willing
<b>Q26</b>	<b>Inequality</b>
Question	Do you think the difference in incomes between rich people and poor people in the United States today is larger, smaller, or about the same as it was 20 years ago?

Value Labels	-2. DK/RF 1. Larger 2. Smaller 3. About the same
<b>Q27</b>	<b>Environment</b>
Question	Do you think the federal government should be doing more about rising temperatures, should be doing less, or is it currently doing the right amount?
Value Labels	-2. DK/RF 1. Should be doing more 2. Should be doing less 3. Is currently doing the right amount
<b>Q28</b>	<b>Parental Leave</b>
Question	Do you favor, oppose, or neither favor nor oppose requiring employers to offer paid leave to parents of new children?
Value Labels	-2. DK/RF 1. Favor 2. Oppose 3. Neither favor nor oppose
<b>Q29</b>	<b>LGBTQ+ Rights</b>
Question	Do you think business owners who provide wedding-related services should be allowed to refuse services to same-sex couples if same-sex marriage violates their religious beliefs, or do you think business owners should be required to provide services regardless of a couple's sexual orientation?
Value Labels	-2. DK/RF 1. Should be allowed to refuse 2. Should be required to provide services
<b>Q30</b>	<b>LGBTQ+ Rights</b>
Question	Should transgender people - that is, people who identify themselves as the sex or gender different from the one they were born as - have to use the bathrooms of the gender they were born as, or should they be allowed to use the bathrooms of their identified gender?
Value Labels	-2. DK/RF 1. Have to use the bathrooms of the gender they were born as 2. Be allowed to use the bathrooms of their identified gender
<b>Q31</b>	<b>LGBTQ+ Rights</b>
Question	Do you favor or oppose laws to protect gays and lesbians against job discrimination?
Value Labels	-2. DK/RF 1. Favor 2. Oppose
<b>Q32</b>	<b>LGBTQ+ Rights</b>
Question	Do you think gay or lesbian couples should be legally permitted to adopt children?
Value Labels	-2. DK/RF 1. Yes 2. No
<b>Q33</b>	<b>LGBTQ+ Rights</b>
Question	Which comes closest to your view? You can just tell me the number of your choice.
Value Labels	-2. DK/RF 1. Gay and lesbian couples should be allowed to legally marry 2. Gay and lesbian couples should be allowed to form civil unions but not legally marry 3. There should be no legal recognition of gay or lesbian couples' relationship
<b>Q34</b>	<b>Immigration</b>

Question	Some people have proposed that the U.S. Constitution should be changed so that the children of unauthorized immigrants do not automatically get citizenship if they are born in this country. Do you favor, oppose, or neither favor nor oppose this proposal?
Value Labels	-2. DK/RF 1. Favor 2. Oppose 3. Neither favor nor oppose
<b>Q35</b>	<b>Immigration</b>
Question	What should happen to immigrants who were brought to the U.S. illegally as children and have lived here for at least 10 years and graduated high school here? Should they be sent back where they came from, or should they be allowed to live and work in the United States?
Value Labels	-2. DK/RF 1. Should be sent back where they came from 2. Should be allowed to live and work in the US
<b>Q36</b>	<b>Immigration</b>
Question	Do you favor, oppose, or neither favor nor oppose building a wall on the U.S. border with Mexico?
Value Labels	-2. DK/RF 1. Favor 2. Oppose 3. Neither favor nor oppose
<b>Q37</b>	<b>Unrest</b>
Question	During the past few months, would you say that most of the actions taken by protestors to get the things they want have been violent, or have most of these actions by protesters been peaceful, or have these actions been equally violent and peaceful?
Value Labels	-2. DK/RF 1. Mostly violent 2. Mostly peaceful 3. Equally violent and peaceful
<b>Q38</b>	<b>Government</b>
Question	Do you think it is better when one party controls both the presidency and Congress, better when control is split between the Democrats and Republicans, or doesn't it matter?
Value Labels	-2. DK/RF 1. Better when one party controls both 2. Better when control is split 3. It doesn't matter
<b>Q39</b>	<b>Government</b>
Question	Would you say the government is pretty much run by a few big interests looking out for themselves or that it is run for the benefit of all the people?
Value Labels	-2. DK/RF 1. Run by a few big interests 2. For the benefit of all the people
<b>Q40</b>	<b>Government</b>
Question	Do you think that people in government waste a lot of the money we pay in taxes, waste some of it, or don't waste very much of it?



Value Labels	-2. DK/RF 1. Waste a lot 2. Waste some 3. Don't waste very much
<b>Q41</b>	<b>Election Integrity</b>
Question	Do you favor, oppose, or neither favor nor oppose allowing convicted felons to vote once they complete their sentence?
Value Labels	-2. DK/RF 1. Favor 2. Oppose 3. Neither favor nor oppose
<b>Q42</b>	<b>Democratic Norms</b>
Question	How important is it that news organizations are free to criticize political leaders?
Value Labels	-2. DK/RF 1. Not important 2. Moderately important 3. Important
<b>Q43</b>	<b>Democratic Norms</b>
Question	How important is it that the executive, legislative, and judicial branches of government keep one another from having too much power?
Value Labels	-2. DK/RF 1. Not important 2. Moderately important 3. Important
<b>Q44</b>	<b>Democratic Norms</b>
Question	How important is it that elected officials face serious consequences if they engage in misconduct?
Value Labels	-2. DK/RF 1. Not important 2. Moderately important 3. Important
<b>Q45</b>	<b>Democratic Norms</b>
Question	How important is it that people agree on basic facts even if they disagree politically?
Value Labels	-2. DK/RF 1. Not important 2. Moderately important 3. Important
<b>Q46</b>	<b>Democratic Norms</b>
Question	Would it be helpful, harmful, or neither helpful nor harmful if U.S. presidents could work on the country's problems without paying attention to what Congress and the courts say?
Value Labels	-2. DK/RF 1. Helpful 2. Harmful 3. Neither helpful nor harmful
<b>Q47</b>	<b>Democratic Norms</b>
Question	Do you favor, oppose, or neither favor nor oppose elected officials restricting journalists' access to information about government decision-making?

Value Labels	-2. DK/RF 1. Favor 2. Oppose 3. Neither favor nor oppose
<hr/>	
<b>Q48</b>	<b>Gender Resentment</b>
Question	'Many women interpret innocent remarks or acts as being sexist.' Do you agree, neither agree nor disagree, or disagree with this statement?
Value Labels	-2. DK/RF/technical error 1. Agree 2. Neither agree nor disagree 3. Disagree
<hr/>	
<b>Q49</b>	<b>Gender Resentment</b>
Question	'Women seek to gain power by getting control over men.' Do you agree, neither agree nor disagree, or disagree with this statement?
Value Labels	-2. DK/RF/technical error 1. Agree 2. Neither agree nor disagree 3. Disagree
<hr/>	