# Uncertainty Quantification for LLM-Based Survey Simulations

Chengpiao Huang [1]    Yuhang Wu [2]    Kaizheng Wang [1 3]

## Abstract

We investigate the use of large language models (LLMs) to simulate human responses to survey questions, and perform uncertainty quantification to assess the fidelity of the simulations. Our approach converts imperfect black-box LLM-simulated responses into confidence sets for population parameters of human responses. A key innovation lies in determining the optimal number of simulated responses: too many produce overly narrow confidence sets with poor coverage, while too few yield excessively loose estimates. Our method adaptively selects the simulation sample size that ensures valid average-case coverage guarantees. The selected sample size itself further provides a quantitative measure of LLM-human misalignment. Experiments on real survey datasets reveal heterogeneous fidelity gaps across different LLMs and domains.

## 1. Introduction

Large language models (LLMs) have demonstrated remarkable capabilities in mimicking human behaviors. Recent studies have leveraged LLMs to simulate human responses in various domains, including economic and social science experiments (Aher et al., 2023; Horton, 2023; Chen et al., 2023; Bisbee et al., 2024; Huang et al., 2024; Yang et al., 2024; Ziems et al., 2024), market research (Brand et al., 2023; Gui & Toubia, 2023; Goli & Singh, 2024; Wang et al., 2024), and education (Zelikman et al., 2023; Lu & Wang, 2024). The typical simulation procedure consists in prompting an LLM with a real or fictional persona as well as a survey question, and collecting the LLM's responses. Compared to traditional survey methods that recruit and query real people, LLM simulations offer significant advantages in terms of time and cost efficiency, enabling the generation of large-scale synthetic responses with minimal effort.

However, a growing body of evidence suggests that LLMs are not perfectly aligned with the human population, and in some cases, the misalignment can be substantial (Aher et al., 2023; Santurkar et al., 2023). This raises critical concerns about the reliability of insights derived from LLM-generated data. It remains a challenge how to properly simulate human responses using LLMs and how to account for their imperfections when using the simulated samples to make inference about the true human population.

We propose to address this challenge through the lens of *uncertainty quantification*. Specifically, we seek to construct confidence sets for population statistics of human responses based on LLM-generated data. The confidence sets will provide a quantitative assessment of the fidelity of the LLM simulations. A central question in confidence set construction is:

*How many synthetic samples should be generated?*

On one hand, generating too many samples risks overfitting the synthetic distribution, which may deviate from the real human population. On the other hand, generating too few samples yields overly large and uninformative confidence sets. The optimal sample size depends on the discrepancy between the synthetic and real populations, which is unknown in practice. This necessitates a data-driven approach to determine the appropriate number of simulated responses.

**Main contributions.** In this paper, we develop a general framework to address these challenges. Our key contributions are as follows:

- (Formulation) We provide a rigorous mathematical framework for uncertainty quantification in LLM-based survey simulations.

- (Methodology) We propose a flexible methodology that transforms simulated responses into valid confidence sets for population parameters of human responses. Our approach adaptively selects the simulation sample size based on the observed misalignment between the LLM and human populations. It is applicable to any LLM, regardless of its fidelity, and can be combined with any method for confidence set construction.

[1]Department of Industrial Engineering and Operations Research, Columbia University [2]Decision, Risk, and Operations Division, Columbia Business School [3]Data Science Institute, Columbia University. Correspondence to: Kaizheng Wang <kaizheng.wang@columbia.edu>.

**Related works.** Our work relates to research on assessing the fidelity of LLM simulations and measuring their alignment with real human populations. Prior studies have explored similarity metrics between synthetic and human distributions (Santurkar et al., 2023; He-Yueya et al., 2024; Dominguez-Olmedo et al., 2024; Durmus et al., 2024; Calderon et al., 2025) and Turing-type tests (Argyle et al., 2023; Mei et al., 2024) to evaluate LLM reliability. While these approaches provide valuable insights into LLM misalignment, they do not offer methods for leveraging imperfect LLM simulations to draw reliable conclusions about human populations. In contrast, our work provides a principled approach for constructing confidence sets that account for the inherent discrepancies between LLM-generated and human responses.

## 2. Problem Setup

To highlight the key challenges, we focus on the simple setting where an LLM simulates binary responses to a survey question.

### 2.1. Motivating Example: Educational Test

Suppose a school wants to estimate the proportion $\mu \in [0, 1]$ of students that can answer a newly designed test question correctly. It will not only provide insights into student progress but also evaluate the question's effectiveness in differentiating among students with varying levels of understanding. Such information can guide the school in tailoring teaching strategies to better address student needs.

The most direct approach is to give the test to $n$ students and collect their results $y_1, ..., y_n \in \{0, 1\}$, where $y_i$ indicates whether student $i$ answers the question correctly. A point estimate for $\mu$ is the sample mean $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$. Given $\alpha \in (0, 1)$, we can construct a confidence interval for $\mu$ based on the Hoeffding's concentration inequality (e.g., Theorem 2.8 in (Boucheron et al., 2013)) to construct a finite-sample confidence interval

$$\left[ \bar{y} - \sqrt{\frac{\log(2/\alpha)}{2n}}, \ \bar{y} + \sqrt{\frac{\log(2/\alpha)}{2n}} \right], \qquad (1)$$

which has at least $(1 - \alpha)$ coverage probability for every $n \in \mathbb{Z}_+$.

Alternatively, the school may use an LLM to simulate students' responses to the question. Compared with directly testing on real students, this approach is more time-efficient and cost-saving. If we prompt the LLM $k$ times with random student profiles, then it generates $k$ synthetic responses, which leads to synthetic outcomes $y_1^{\text{syn}}, ..., y_k^{\text{syn}} \in \{0, 1\}$. We may also compute the sample mean $\bar{y}_k^{\text{syn}} = \frac{1}{k} \sum_{i=1}^k y_i^{\text{syn}}$

and the confidence interval

$$\mathcal{I}^{\text{syn}}(k) = \left[ \bar{y} - c\sqrt{\frac{\log(2/\alpha)}{2n}}, \ \bar{y} + c\sqrt{\frac{\log(2/\alpha)}{2n}} \right], \quad (2)$$

where $c > 1$ is a scaling parameter. Such a dilation by $c$ is necessary; without it, whenever the LLM-generated data deviates from the student population (even by the slightest amount), the interval $\mathcal{I}^{\text{syn}}(k)$ may never achieve $(1 - \alpha)$ coverage regardless of $k$.

Due to the misalignment between the LLM and students, the distribution of the synthetic responses may be very different from the true response distribution. In this case, the sample mean $\bar{y}^{\text{syn}}$ can be a poor estimate of $\mu$, and $\mathcal{I}^{\text{syn}}(k)$ is generally not a valid confidence interval for $\mu$. In particular, as $k \to \infty$, the interval concentrates tightly around the synthetic mean $\mathbb{E}[y_1^{\text{syn}}]$ and fails to cover the true mean $\mu$. When $k$ is small, the interval becomes too wide to be informative, even though it may cover $\mu$ with high probability.

### 2.2. Simulation Sample Size as a Fidelity Measure

In this work, we will develop a principled approach for choosing a good simulation sample size $\widehat{k}$, so that $\mathcal{I}^{\text{syn}}(\widehat{k})$ is a valid confidence interval for $\mu$ while having a modest width. Solving this problem has the following important implications.

1. The choice of $\widehat{k}$ offers valuable information for future simulation tasks on the appropriate number of synthetic samples to generate, so as to produce reliable confidence intervals. It also helps avoid generating excessive samples and improves computational efficiency.

2. The width of $\mathcal{I}^{\text{syn}}(\widehat{k})$ provides an assessment of the alignment between the LLM and the human population. A wide confidence interval indicates high uncertainty of its estimate of the true $\mu$, and thus a large gap between the synthetic data distribution and the true population.

3. The sample size $\widehat{k}$ reflects the size of the target population that the LLM can represent. We make an analogy using the classical theory of parametric bootstrap. Suppose a model is trained via maximum likelihood estimation over $k$ i.i.d. human samples. When performing parametric bootstrap for uncertainty quantification, the bootstrap sample size is usually set to be the training sample size $k$. Thus, our simulation sample size $\widehat{k}$ reveals the LLM as "being made up of" $\widehat{k}$ people from the population. The larger $\widehat{k}$ is, the more diversity that the LLM appears to capture. In contrast, a small $\widehat{k}$ could imply the peculiarity of the LLM compared to the major population.

*Remark* 2.1 (Comparison with existing works). Existing works typically measure LLM misalignment using integral probability metrics and $f$-divergences (Santurkar et al., 2023; Dominguez-Olmedo et al., 2024; Durmus et al., 2024), which do not carry operational meanings themselves and can be hard to interpret. In contrast, our simulation sample size $\widehat{k}$ provides actionable guidance and is easy to understand.

## 3. Methodology

We now introduce our method for choosing a good simulation sample size $\widehat{k}$. We focus on the simple setting in Section 2.1, but our method can be extended to more general settings. Our approach makes use of similar test questions for which real students' results are available. On theses questions, we can compare LLM simulations with real students' results on these questions, and use it to guide the choice of $k$.

Specifically, we assume access to $m$ similar test questions. For example, they can come from previous tests or a question bank. For $j \in [m]$, the $j$-th test question has been tested on $n_j$ real students, with test results $\mathcal{D}_j = \{y_{j,i}\}_{i=1}^{n_j}$. We also simulate LLM responses $\mathcal{D}_j^{\mathsf{syn}} = \{y_{j,i}^{\mathsf{syn}}\}_{i=1}^{K}$ to the $j$-th test question, and $\mathcal{D}^{\mathsf{syn}} = \{y_i^{\mathsf{syn}}\}_{i=1}^{K}$ to the new test question. Here $K \in \mathbb{Z}_+$ is the simulation budget.

For each question $j$, we form confidence intervals similar to (2) using the synthetic data $\mathcal{D}_j^{\mathsf{syn}}$, aiming to cover the proportion $\mu_j$ of students that answer the $j$-th question correctly:

$$\mathcal{I}_j^{\mathsf{syn}}(k) = \left[ \bar{y}_{j,k}^{\mathsf{syn}} - c\sqrt{\frac{\log(2/\alpha)}{2k}}, \bar{y}_{j,k}^{\mathsf{syn}} + c\sqrt{\frac{\log(2/\alpha)}{2k}} \right], \tag{3}$$

where $\bar{y}_{j,k}^{\mathsf{syn}} = \frac{1}{k}\sum_{i=1}^{k} y_{j,i}^{\mathsf{syn}}$ is the sample mean of the first $k$ samples in $\mathcal{D}_j^{\mathsf{syn}}$. We also set $\mathcal{I}^{\mathsf{syn}}(0) = \mathcal{I}_j^{\mathsf{syn}}(0) = \mathbb{R}$. We will pick $\widehat{k} \in \{0, 1..., K\}$ such that $\mathcal{I}_j^{\mathsf{syn}}(\widehat{k})$ covers $\mu_j$ with high probability. We expect this choice of $\widehat{k}$ to be also good for $\mathcal{I}^{\mathsf{syn}}(k)$, as the test questions are similar.

Ideally, we would like to pick $k$ such that $(1 - \alpha)$-coverage is achieved empirically over the $m$ test questions:

$$\frac{1}{m}\sum_{j=1}^{m} \mathbf{1}\{\mu_j \notin \mathcal{I}_j^{\mathsf{syn}}(k)\} \leq \alpha. \tag{4}$$

As the true $\{\mu_j\}_{j=1}^{m}$ are not available, we use real data $\{\mathcal{D}_j\}_{j=1}^{m}$ to compute the sample means $\bar{y}_j = \frac{1}{n_j}\sum_{i=1}^{n_j} y_{j,i}$ as proxies. We approximate the empirical miscoverage by

$$G(k) = \frac{1}{m}\sum_{j=1}^{m} \mathbf{1}\{\bar{y}_j \notin \mathcal{I}_j^{\mathsf{syn}}(k)\}. \tag{5}$$

Our criterion for selecting $k$ is given by

$$\widehat{k} = \max\left\{0 \leq k \leq K : G(i) \leq \alpha/2 \ \forall i \leq k\right\}. \tag{6}$$

Note that $\widehat{k}$ is well-defined because $G(0) = 0$. Here the threshold is $\alpha/2$ instead of $\alpha$ due to the approximation error of $\bar{y}_j \approx \mu_j$. We provide a detailed explanation in Section A.

In Section B, we give a theoretical analysis of our method, which shows that the selected confidence interval $\mathcal{I}^{\mathsf{syn}}(\widehat{k})$ achieves valid coverage while having a near-optimal width.

## 4. Extension to More General Settings

In this section, we present a more general method that applies to the general setting where survey responses and confidence sets can be multi-dimensional.

### 4.1. General Problem Formulation

Let $\mathcal{Z}$ be a profile space, $\mathcal{P}$ a probability distribution over $\mathcal{Z}$ which represents the true population, and $\mathcal{P}^{\mathsf{syn}}$ a synthetic distribution over $\mathcal{Z}$ used to generate synthetic profiles.

Let $\Psi$ be a collection of survey questions, and $\mathcal{Y}$ be the space of possible responses to the survey questions. When a person with profile $z \in \mathcal{Z}$ is asked a survey question $\psi \in \Psi$, the person gives a response $y$ following a distribution $\mathcal{Q}(\cdot \mid z, \psi)$ over $\mathcal{Y}$. We are interested in the distribution of the population's response to the survey question $\psi$, which is given by $\mathcal{R}(\cdot \mid \psi) = \int_{\mathcal{Z}} \mathcal{Q}(\cdot \mid z, \psi)\,\mathcal{P}(dz)$. In particular, we seek to construct a confidence set for some statistic $\theta(\psi)$ of $\mathcal{R}(\cdot \mid \psi)$, which can be multi-dimensional, say in $\mathbb{R}^d$. We provide examples in Section C.1.

We consider constructing the confidence set by using simulated responses from an LLM. Given a profile $z$, a survey question $\psi$ and a prompt $p$, the LLM simulates a response $y^{\mathsf{syn}}$ from a distribution $\mathcal{Q}^{\mathsf{syn}}(\cdot \mid z, \psi, p)$ which aims to mimic $\mathcal{Q}(\cdot \mid z, \psi)$. We can generate synthetic profiles $\{z_i^{\mathsf{syn}}\}_{i=1}^{K}$ from some distribution $\mathcal{P}^{\mathsf{syn}}$, then feed them into the LLM along with $\psi$ and $p$. The LLM then generates synthetic responses $\{y_i^{\mathsf{syn}}\}_{i=1}^{K}$, where $y_i^{\mathsf{syn}} \sim \mathcal{Q}^{\mathsf{syn}}(\cdot \mid z_i^{\mathsf{syn}}, \psi, p)$. Here $K$ is the simulation budget.

Using the simulated samples $\mathcal{D}^{\mathsf{syn}} = \{y_i^{\mathsf{syn}}\}_{i=1}^{K}$, we can construct a family of candidate confidence sets such as the one-dimensional CLT-based confidence interval (2). More generally, the statistics literature has developed a variety of approaches such as inverting hypothesis tests (Casella & Berger, 2002), the bootstrap (Efron, 1979), and the empirical likelihood ratio function (Owen, 1990). We will assume access to a black-box procedure $\mathcal{C}$ that takes as input a dataset $\mathcal{D}$ and outputs a confidence set $\mathcal{C}(\mathcal{D}) \subseteq \mathbb{R}^d$. Then, we construct a family of confidence sets $\{\mathcal{S}^{\mathsf{syn}}(k)\}_{k=1}^{K}$ by

$$\mathcal{S}^{\mathsf{syn}}(k) = \mathcal{C}\left(\{y_i^{\mathsf{syn}}\}_{i=1}^{k}\right). \tag{7}$$

We also set $\mathcal{S}^{\mathsf{syn}}(0) = \mathbb{R}^d$, so $\theta(\psi) \in \mathcal{S}^{\mathsf{syn}}(0)$ always. We will not impose any assumptions on the quality of the confidence sets produced by $\mathcal{C}$.

As the LLM may not be a faithful reflection of the true human population, we will make use of real data to choose a good confidence set from $\{\mathcal{S}^{\mathsf{syn}}(k)\}_{k=1}^{K}$. We assume that we have collected real human responses from $m$ similar surveys $\psi_1, ..., \psi_m \in \Psi$. For each $j \in [m]$, we have responses $\mathcal{D}_j = \{y_{j,i}\}_{i=1}^{n_j}$ from $n_j$ i.i.d. surveyees $\{z_{j,i}\}_{i=1}^{n_j} \sim \mathcal{P}$, with $y_{j,i} \sim \mathcal{Q}(\,\cdot\mid z_{j,i}, \psi_j)$.

We also simulate LLM responses to these $m$ survey questions. For each $j \in [m]$, we feed synthetic profiles $\{z_{j,i}^{\mathsf{syn}}\}_{i=1}^{K}$, the question $\psi_j$ and the prompt $p$ into the LLM, which simulates responses $\mathcal{D}_j^{\mathsf{syn}} = \{y_{j,i}^{\mathsf{syn}}\}_{i=1}^{K}$ with $y_{j,i}^{\mathsf{syn}} \sim \mathcal{Q}^{\mathsf{syn}}(\,\cdot\mid z_{j,i}^{\mathsf{syn}}, \psi_j, p)$. The datasets $\{\mathcal{D}_j\}_{j=1}^{m}$ and $\{\mathcal{D}_j^{\mathsf{syn}}\}_{j=1}^{m}$ will be used to select a confidence set from $\{\mathcal{S}^{\mathsf{syn}}(k)\}_{k=1}^{K}$.

Our goal is to use $\{\mathcal{D}_j\}_{j=1}^{m}$ and $\{\mathcal{D}_j^{\mathsf{syn}}\}_{j=1}^{m}$ to choose a simulation sample size $\widehat{k} \in [K]$, such that

$$\mathbb{P}\left(\theta(\psi) \in \mathcal{S}^{\mathsf{syn}}(\widehat{k})\right) \approx 1 - \alpha.$$

### 4.2. General Methodology for Sample Size Selection

We now present our general methodology. For each $j \in [m]$, we form confidence sets similar to (7) using the synthetic data $\mathcal{D}_j^{\mathsf{syn}}$:

$$\mathcal{S}_j^{\mathsf{syn}}(k) = \mathcal{C}\left(\{y_{j,i}^{\mathsf{syn}}\}_{i=1}^{k}\right), \quad \forall k \in [K]. \tag{8}$$

We also set $\mathcal{S}_j^{\mathsf{syn}}(0) = \mathbb{R}^d$. We will pick $\widehat{k} \in \{0, 1, ..., K\}$ such that $\mathcal{S}_j^{\mathsf{syn}}(\widehat{k})$ is a good confidence interval for $\theta(\psi_j)$ for each $j \in [m]$. We expect this choice of $\widehat{k}$ to be also good for $\mathcal{S}^{\mathsf{syn}}(k)$, since the survey questions are similar.

Ideally, we would like to pick $k$ such that $(1 - \alpha)$ coverage is achieved empirically over the $m$ survey functions:

$$\frac{1}{m} \sum_{j=1}^{m} \mathbf{1}\{\theta(\psi_j) \notin \mathcal{S}_j^{\mathsf{syn}}(k)\} \leq \alpha. \tag{9}$$

However, the population-level quantities $\{\theta(\psi_j)\}_{j=1}^{m}$ are not available, so we must approximate them by the real data $\{\mathcal{D}_j\}_{j=1}^{m}$. In Section 3, we have taken the approach of constructing unbiased point estimates, but it does not directly extend to the more general case.

Instead, we will use the real data $\{\mathcal{D}_j\}_{j=1}^{m}$ to construct confidence sets for $\{\theta(\psi_j)\}_{j=1}^{m}$. Choose a confidence level $\gamma \in (0, 1)$. For each $j \in [m]$, we use $\mathcal{D}_j$ to construct a confidence set $\mathcal{S}_j$ that satisfies

$$\mathbb{P}\left(\theta(\psi_j) \in \mathcal{S}_j \mid \psi_j\right) \geq \gamma. \tag{10}$$

These confidence sets are easy to construct as the samples in $\mathcal{D}_j$ follow the true response distribution. When $\theta(\psi_j) \in$

$\mathcal{S}_j$, the condition $\mathcal{S}_j \subseteq \mathcal{S}_j^{\mathsf{syn}}(k)$ is sufficient for $\theta(\psi_j) \in \mathcal{S}_j^{\mathsf{syn}}(k)$. Equivalently, when $\theta(\psi_j) \in \mathcal{S}_j$, the condition $\theta(\psi_j) \notin \mathcal{S}_j^{\mathsf{syn}}(k)$ must imply $\mathcal{S}_j \not\subseteq \mathcal{S}_j^{\mathsf{syn}}(k)$. Thus, we take

$$L(k) = \frac{1}{m} \sum_{j=1}^{m} \mathbf{1}\{\mathcal{S}_j \not\subseteq \mathcal{S}_j^{\mathsf{syn}}(k)\} \tag{11}$$

as a proxy for the empirical miscoverage. Since $\theta(\psi_j) \in \mathcal{S}_j$ happens with probability $\gamma$, then the frequency of having $\mathcal{S}_j \not\subseteq \mathcal{S}_j^{\mathsf{syn}}(k)$ is at least $\gamma$ times the frequency of $\theta(\psi_j) \notin \mathcal{S}_j^{\mathsf{syn}}(k)$. Roughly speaking,

$$L(k) \geq \gamma \cdot \left(\frac{1}{m} \sum_{j=1}^{m} \mathbf{1}\{\theta(\psi_j) \notin \mathcal{S}_j^{\mathsf{syn}}(k)\}\right). \tag{12}$$

Combining (9) and (12) leads to the following criterion for selecting $k$:

$$\widehat{k} = \max\{0 \leq k \leq K : L(i) \leq \gamma\alpha, \ \forall i \leq k\}. \tag{13}$$

Note that $\widehat{k}$ is well-defined because $L(0) = 0$. In Section C.2, we give a theoretical analysis of our method, which shows that the selected confidence interval $\mathcal{S}^{\mathsf{syn}}(\widehat{k})$ achieves valid average-case coverage.

## 5. Numerical Experiments

We apply our method in Section 3 to LLMs on real data.

### 5.1. Experiment Setup

**LLMs.** We consider 8 LLMs: GPT-3.5-Turbo (gpt-3.5-turbo), GTP-4o (gpt-4o), and GPT-4o-mini (gpt-4o-mini) (OpenAI, 2022; 2024a;b); Claude 3.5 Haiku (claude-3-5-haiku-20241022) (Anthropic, 2024); Llama 3.1 8B (Llama-3-8B-Instruct-Turbo) and Llama 3.3 70B (Llama-3.3-70B-Instruct-Turbo) (Dubey et al., 2024); Mistral 7B (Mistral-7B-Instruct-v0.3) (Jiang et al., 2023); DeepSeek-V3 (DeepSeek-V3) (Liu et al., 2024).

**Datasets.** We use two datasets for survey questions, each corresponding to one uncertainty quantification task. The first dataset is the OpinionQA dataset created by (Santurkar et al., 2023). It was built from Pew Research's American Trends Panel[1], and contains the general US population's responses to survey questions spanning topics such as science, politics, and health. After pre-processing we have 385 unique questions, each with at least 400 responses. These questions have 5 choices corresponding to ordered sentiments which we map to sentiment scores $-1, -\frac{1}{3}, 0, \frac{1}{3}, 1$.

---

[1] https://www.pewresearch.org/the-american-trends-panel/

For each response, we have information on their political profile, religious affiliation, educational background, socioeconomic status, etc. This information is to generate synthetic profiles. More details about the dataset can be found in Section D.1. We consider the task of constructing a confidence interval for the US population's average sentiment score for a survey question.

The second dataset is the EEDI dataset created by (He-Yueya et al., 2024), which was built upon the NeurIPS 2020 Education Challenge dataset (Wang et al., 2021). It consists of students' responses to mathematics multiple-choice questions on the Eedi online educational platform[2]. All questions have four choices (A, B, C, D). Out of these questions, we use questions that have at least $100$ student responses, which gives a total of $412$ questions. For each student, we have information on their gender, age, and socioeconomic status. This information is used to generate synthetic profiles. More details about the dataset can be found in Section D.2. We consider the task of constructing a confidence interval for the probability of a student answering a question correctly.

**Hyperparameters.** We consider $\alpha \in \{0.05 \cdot \ell : \ell \in [10]\}$ and $c = \sqrt{2}$. For the EEDI dataset, we set the simulation budget $K = 50$ and take $M = 1$. For the OpinionQA dataset, we set $K = 100$ and take $M = 2$.

**Procedure.** We randomly split the survey questions into a training set and a testing set, where the training set is used to select the simulation sample size $\widehat{k}$, and the testing set is used to evaluate the coverage of the constructed confidence interval. More details can be found in Section D.3.

### 5.2. Experiment Results

We now show the experiment results. We omit Llama 3.1 8B for the EEDI dataset experiment because it frequently failed to answer EEDI questions in required formats. As a baseline, we also include a naïve response generator (`random`) that chooses an available answer uniformly at random.

We verify the coverage validity of the selected confidence interval through standard statistical tests. The detailed results are deferred to Section E. Next, we turn to the selected simulation sample size $\widehat{k}$. In Figure 1, we plot the average $\widehat{k}$ over the 100 random splits for various LLMs on the OpinionQA and EEDI datasets, respectively. The error bars represent 95% confidence intervals.

In general, a larger $\widehat{k}$ means that the LLM has stronger simulation power. For the EEDI dataset, DeepSeek-V3 has the best performance, followed by GPT-4o and Claude 3.5 Haiku. For the OpinionQA dataset, GPT-4o has the best performance. Interestingly, while on the EEDI dataset only
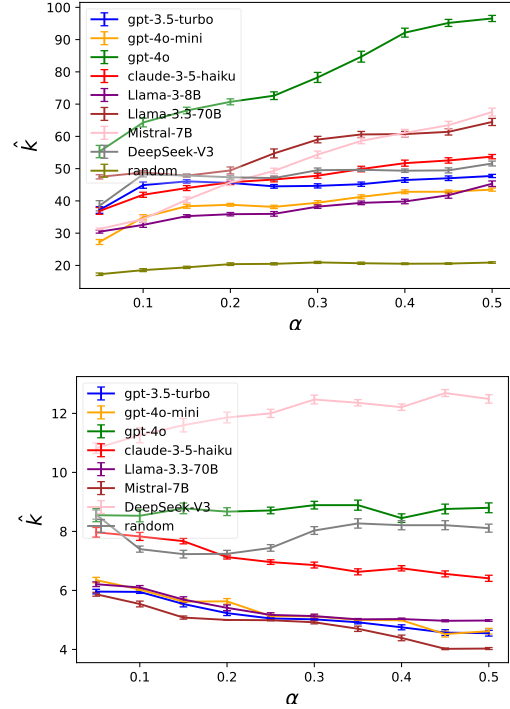
*Figure 1.* Average $\widehat{k}$ for various LLMs and $\alpha$ on the OpinionQA dataset (top) and the EEDI dataset (bottom)

DeepSeek-V3 and GPT-4o seem to outperform the random benchmark, on the OpinionQA dataset all LLMs clearly outperform the random benchmark. Moreover, LLMs exhibit uniformly higher $\widehat{k}$ on the OpinionQA dataset than on the EEDI dataset, suggesting higher fidelity in simulating subjective opinions to social problems than in simulating student answers to mathematics questions.

The experiment results demonstrate the importance of a disciplined approach to using synthetic samples. The ease of LLM-based simulation makes it tempting to generate a large number of responses. Our results show great heterogeneity in the simulation power of different LLMs over different datasets: the largest $\widehat{k}$ is below 100, while the smallest $\widehat{k}$ is less than 10. This means that there is real peril in using excessive synthetic samples.

## 6. Discussions

We developed a general approach for converting imperfect LLM-based survey simulations into statistically valid confidence sets for population statistics of human responses. It identifies a simulation sample size which is useful for future simulation tasks and which reveals the degree of misalignment between the LLM and the target human population.

# References

Aher, G. V., Arriaga, R. I., and Kalai, A. T. Using large language models to simulate multiple humans and replicate human subject studies. In Krause, A., Brunskill, E., Cho, K., Engelhardt, B., Sabato, S., and Scarlett, J. (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 337–371. PMLR, 23–29 Jul 2023. URL https://proceedings.mlr.press/v202/aher23a.html.

Anthropic. Introducing computer use, a new Claude 3.5 Sonnet, and Claude 3.5 Haiku, 2024. URL https://www.anthropic.com/news/3-5-models-and-computer-use.

Argyle, L. P., Busby, E. C., Fulda, N., Gubler, J. R., Rytting, C., and Wingate, D. Out of one, many: Using language models to simulate human samples. *Political Analysis*, 31 (3):337–351, 2023. doi: 10.1017/pan.2023.2.

Bisbee, J., Clinton, J. D., Dorff, C., Kenkel, B., and Larson, J. M. Synthetic replacements for human survey data? The perils of large language models. *Political Analysis*, 32(4): 401–416, 2024. doi: 10.1017/pan.2024.5.

Boucheron, S., Lugosi, G., and Massart, P. *Concentration Inequalities: A Nonasymptotic Theory of Independence*. Oxford University Press, 02 2013. ISBN 9780199535255. doi: 10.1093/acprof:oso/9780199535255.001.0001. URL https://doi.org/10.1093/acprof:oso/9780199535255.001.0001.

Brand, J., Israeli, A., and Ngwe, D. Using LLMs for market research. *Harvard Business School Marketing Unit Working Paper*, (23-062), 2023.

Calderon, N., Reichart, R., and Dror, R. The alternative annotator test for LLM-as-a-judge: How to statistically justify replacing human annotators with LLMs. *arXiv preprint arXiv:2501.10970*, 2025. URL https://arxiv.org/abs/2501.10970.

Casella, G. and Berger, R. *Statistical inference*. CRC press, 2002.

Chen, Y., Liu, T. X., Shan, Y., and Zhong, S. The emergence of economic rationality of GPT. *Proceedings of the National Academy of Sciences*, 120(51):e2316205120, 2023. doi: 10.1073/pnas.2316205120. URL https://www.pnas.org/doi/abs/10.1073/pnas.2316205120.

Dominguez-Olmedo, R., Hardt, M., and Mendler-Dünner, C. Questioning the survey responses of large language models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=Oo7dlLgqQX.

Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Yang, A., Fan, A., et al. The Llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

Durmus, E., Nguyen, K., Liao, T., Schiefer, N., Askell, A., Bakhtin, A., Chen, C., Hatfield-Dodds, Z., Hernandez, D., Joseph, N., Lovitt, L., McCandlish, S., Sikder, O., Tamkin, A., Thamkul, J., Kaplan, J., Clark, J., and Ganguli, D. Towards measuring the representation of subjective global opinions in language models. In *First Conference on Language Modeling*, 2024. URL https://openreview.net/forum?id=zl16jLb91v.

Efron, B. Bootstrap methods: Another look at the jackknife. *The Annals of Statistics*, 7(1):1 – 26, 1979. doi: 10.1214/aos/1176344552. URL https://doi.org/10.1214/aos/1176344552.

Goli, A. and Singh, A. Frontiers: Can large language models capture human preferences? *Marketing Science*, 43(4):709–722, 2024. doi: 10.1287/mksc.2023.0306. URL https://doi.org/10.1287/mksc.2023.0306.

Gui, G. and Toubia, O. The challenge of using LLMs to simulate human behavior: A causal inference perspective. *arXiv preprint arXiv:2312.15524*, 2023.

He-Yueya, J., Ma, W. A., Gandhi, K., Domingue, B. W., Brunskill, E., and Goodman, N. D. Psychometric alignment: Capturing human knowledge distributions via language models. *arXiv preprint arXiv:2407.15645*, 2024.

Horton, J. J. Large language models as simulated economic agents: What can we learn from homo silicus? Technical report, National Bureau of Economic Research, 2023.

Huang, Y., Yuan, Z., Zhou, Y., Guo, K., Wang, X., Zhuang, H., Sun, W., Sun, L., Wang, J., Ye, Y., et al. Social science meets LLMs: How reliable are large language models in social simulations? *arXiv preprint arXiv:2410.23426*, 2024.

Jiang, A. Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D. S., Casas, D. d. l., Bressand, F., Lengyel, G., Lample, G., Saulnier, L., et al. Mistral 7B. *arXiv preprint arXiv:2310.06825*, 2023.

Liu, A., Feng, B., Xue, B., Wang, B., Wu, B., Lu, C., Zhao, C., Deng, C., Zhang, C., Ruan, C., et al. DeepSeek-V3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.

Lu, X. and Wang, X. Generative students: Using LLM-simulated student profiles to support question item evaluation. In *Proceedings of the Eleventh ACM Conference on Learning @ Scale*, L@S '24, pp. 16–27,

New York, NY, USA, 2024. Association for Computing Machinery. ISBN 9798400706332. doi: 10.1145/3657604.3662031. URL https://doi.org/10.1145/3657604.3662031.

Mei, Q., Xie, Y., Yuan, W., and Jackson, M. O. A Turing test of whether AI chatbots are behaviorally similar to humans. *Proceedings of the National Academy of Sciences*, 121(9):e2313925121, 2024. doi: 10.1073/pnas.2313925121. URL https://www.pnas.org/doi/abs/10.1073/pnas.2313925121.

OpenAI. GPT-3.5 Turbo, 2022. URL https://platform.openai.com/docs/models/gpt-3-5#gpt-3-5-turbo.

OpenAI. Hello GPT-4o, 2024a. URL https://openai.com/index/hello-gpt-4o/.

OpenAI. GPT-4o mini: Advancing cost-efficient intelligence, 2024b. URL https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/.

Owen, A. Empirical Likelihood Ratio Confidence Regions. *The Annals of Statistics*, 18(1):90 – 120, 1990. doi: 10.1214/aos/1176347494. URL https://doi.org/10.1214/aos/1176347494.

Santurkar, S., Durmus, E., Ladhak, F., Lee, C., Liang, P., and Hashimoto, T. Whose opinions do language models reflect? In Krause, A., Brunskill, E., Cho, K., Engelhardt, B., Sabato, S., and Scarlett, J. (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 29971–30004. PMLR, 23–29 Jul 2023. URL https://proceedings.mlr.press/v202/santurkar23a.html.

Wang, M., Zhang, D. J., and Zhang, H. Large language models for market research: A data-augmentation approach. *arXiv preprint arXiv:2412.19363*, 2024.

Wang, Z., Lamb, A., Saveliev, E., Cameron, P., Zaykov, J., Hernandez-Lobato, J. M., Turner, R. E., Baraniuk, R. G., Craig Barton, E., Peyton Jones, S., Woodhead, S., and Zhang, C. Results and insights from diagnostic questions: The NeurIPS 2020 education challenge. In Escalante, H. J. and Hofmann, K. (eds.), *Proceedings of the NeurIPS 2020 Competition and Demonstration Track*, volume 133 of *Proceedings of Machine Learning Research*, pp. 191–205. PMLR, 06–12 Dec 2021. URL https://proceedings.mlr.press/v133/wang21a.html.

Yang, K., Li, H., Wen, H., Peng, T.-Q., Tang, J., and Liu, H. Are large language models (LLMs) good social predictors? In Al-Onaizan, Y., Bansal, M., and Chen, Y.-N. (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 2718–2730, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-emnlp.153. URL https://aclanthology.org/2024.findings-emnlp.153/.

Zelikman, E., Ma, W., Tran, J., Yang, D., Yeatman, J., and Haber, N. Generating and evaluating tests for K-12 students with language model simulations: A case study on sentence reading efficiency. In Bouamor, H., Pino, J., and Bali, K. (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 2190–2205, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.135. URL https://aclanthology.org/2023.emnlp-main.135/.

Ziems, C., Held, W., Shaikh, O., Chen, J., Zhang, Z., and Yang, D. Can large language models transform computational social science? *Computational Linguistics*, 50(1):237–291, 03 2024. ISSN 0891-2017. doi: 10.1162/coli_a_00502. URL https://doi.org/10.1162/coli_a_00502.

## A. Choice of Threshold in Sample Size Selection

The choice of the threshold $\alpha/2$ in (6) can be explained as follows. By CLT, when $n_j$ is large, $\mathbb{P}(\mu_j \geq \bar{y}_j) \approx 1/2$. Suppose $\mu_j \notin \mathcal{I}_j^{\mathsf{syn}}(k)$, then $\mu_j$ is either on the left or the right of $\mathcal{I}_j^{\mathsf{syn}}(k)$. In the former case, $\bar{y}_j$ is on the left of $\mu_j$ with probability around $1/2$, which implies $\bar{y}_j \notin \mathcal{I}_j^{\mathsf{syn}}(k)$. Similarly, in the latter case, $\bar{y}_j$ is on the right of $\mu_j$ with probability around $1/2$, and then $\bar{y}_j \notin \mathcal{I}_j^{\mathsf{syn}}(k)$. Roughly speaking, the frequency of having $\bar{y}_j \notin \mathcal{I}_j^{\mathsf{syn}}(k)$ is at least half of the frequency of having $\mu_j \notin \mathcal{I}_j^{\mathsf{syn}}(k)$. In other words, the lower bound

$$G(k) \geq \frac{1}{2} \cdot \frac{1}{m} \sum_{j=1}^{m} \mathbf{1}\{\mu_j \notin \mathcal{I}_j^{\mathsf{syn}}(k)\} \tag{14}$$

approximately holds. Substituting (14) into (4) yields the threshold $\alpha/2$ for choosing $\widehat{k}$.

## B. Theoretical Analysis of the Proposed Method

### B.1. Coverage Guarantee

In this section, we present a theoretical analysis of our proposed method. To do so, we first describe the setup in Section 2.1 in mathematical terms.

The student population can be represented by a distribution $\mathcal{P}$ over a space $\mathcal{Z}$ of possible *student profiles*, say, vectors of background information, classes taken, grades, etc. To simulate student responses from the LLM, synthetic student profiles are generated from a synthetic student population $\mathcal{P}^{\mathsf{syn}}$ over $\mathcal{Z}$, and then fed to the LLM.

We use $\psi$ and $\{\psi_j\}_{j=1}^m$ to refer to the test question of interest and the $m$ similar ones, respectively. Students' performance on test questions are characterized by a *performance function* $F$: a student with profile $z \in \mathcal{Z}$ answers a question $\psi$ correctly with probability $F(z, \psi) \in [0, 1]$. The average student performance on the test questions $\psi$ and $\{\psi_j\}_{j=1}^m$ are then $\mu = \mathbb{E}_{z \sim \mathcal{P}} F(z, \psi)$ and $\mu_j = \mathbb{E}_{z \sim \mathcal{P}} F(z, \psi_j)$, respectively. In addition, the LLM generates synthetic student performance from a *synthetic performance function* $F^{\mathsf{syn}}$: when prompted with a synthetic profile $z^{\mathsf{syn}} \in \mathcal{Z}$, the LLM answers a question $\psi$ correctly with probability $F^{\mathsf{syn}}(z^{\mathsf{syn}}, \psi) \in [0, 1]$.

The collection of the real dataset $\mathcal{D}_j = \{y_{j,i}\}_{i=1}^{n_j}$ can be thought of as drawing $n_j$ i.i.d. student profiles $\{z_{j,i}\}_{i=1}^{n_j} \sim \mathcal{P}$ and then sampling $y_{j,i} \sim \mathrm{Bernoulli}(F(z_{j,i}, \psi_j))$ for each $i \in [n_j]$. Similarly, the generation of the synthetic dataset $\mathcal{D}_j^{\mathsf{syn}} = \{y_{j,i}^{\mathsf{syn}}\}_{i=1}^K$ can be thought of as drawing i.i.d. synthetic profiles $\{z_{j,i}^{\mathsf{syn}}\}_{i=1}^K \sim \mathcal{P}^{\mathsf{syn}}$ and then sampling $y_{j,i}^{\mathsf{syn}} \sim \mathrm{Bernoulli}(F^{\mathsf{syn}}(z_{j,i}^{\mathsf{syn}}, \psi_j))$ for each $i \in [K]$. For $\mathcal{D}^{\mathsf{syn}} = \{y_i^{\mathsf{syn}}\}_{i=1}^K$, we adopt a similar notation $\{z_i^{\mathsf{syn}}\}_{i=1}^K$ for the synthetic profiles. We note that when collecting real or synthetic samples, the performance functions never appear explicitly. They are introduced only to facilitate the problem formulation.

Finally, we assume that the test questions are drawn randomly from a question bank, and that the datasets are independent.

**Assumption B.1** (Randomly sampled questions)**.** The questions $\psi, \psi_1, ..., \psi_m$ are independently sampled from a distribution over a space $\Psi$.

**Assumption B.2** (Independent data)**.** For each $j \in [m]$, conditioned on $\psi_j$, the datasets $\mathcal{D}_j$ and $\mathcal{D}_j^{\mathsf{syn}}$ are independent. Conditioned on $\psi_1, ..., \psi_m$, the dataset tuples $(\mathcal{D}_1, \mathcal{D}_1^{\mathsf{syn}}), ..., (\mathcal{D}_m, \mathcal{D}_m^{\mathsf{syn}})$ are independent. Finally, $(\psi, \mathcal{D}^{\mathsf{syn}})$ is independent of $\{(\psi_j, \mathcal{D}_j, \mathcal{D}_j^{\mathsf{syn}})\}_{j=1}^m$.

We are now ready to state the theoretical guarantee of our approach. Its proof is deferred to Section B.3. We note that the assumption $\mathbb{P}(\mu_j \geq \bar{y}_j \mid \psi_j) = 1/2$ is a CLT approximation and is for mathematical convenience only.

**Theorem B.3** (Coverage guarantee)**.** *Let Assumptions B.1 and B.2 hold. Assume that $\mathbb{P}(\bar{y}_j \leq \mu_j \mid \psi_j) = 1/2$ for each $j \in [m]$. Fix $\alpha \in (0, 1)$. Then the simulation sample size $\widehat{k}$ defined by (6) satisfies*

$$\mathbb{P}\left(\mu \in \mathcal{I}^{\mathsf{syn}}(\widehat{k})\right) \geq 1 - \alpha - \sqrt{\frac{2}{m}}.$$

*The probability is taken with respect to randomness of $\{(\psi_j, \mathcal{D}_j, \mathcal{D}_j^{\mathsf{syn}})\}_{j=1}^m$, $\psi$ and $\mathcal{D}^{\mathsf{syn}}$.*

On average, the chosen simulation sample size $\widehat{k}$ leads to a confidence interval $\mathcal{I}^{\mathsf{syn}}(\widehat{k})$ that covers the true mean $\mu$ with probability at least $1 - \alpha - O(\sqrt{1/m})$. As $m \to \infty$, the aforementioned lower bound converges to $1 - \alpha$.

## B.2. Sharpness of Sample Size Selection

We have seen that the chosen interval $\mathcal{I}^{\text{syn}}(\widehat{k})$ has good coverage properties. In this section, we complement this result by showing that the interval is not overly conservative. To simplify computation, we slightly modify the setting.

**Example B.4** (Gaussian performance score). Consider the setting in Section B.1 with the following modifications. On a test question $\psi \in \Psi$, the performance (e.g., score) of a real student follows a Gaussian distribution with mean $\mathbb{E}_{z \sim \mathcal{P}} F(z, \psi)$ and variance 1, instead of a Bernoulli distribution with mean $\mathbb{E}_{z \sim \mathcal{P}} F(z, \psi)$. Similarly, the performance of the LLM follows a Gaussian distribution with mean $\mathbb{E}_{z^{\text{syn}} \sim \mathcal{P}^{\text{syn}}} F^{\text{syn}}(z^{\text{syn}}, \psi)$ and variance 1. Moreover, the confidence intervals $\mathcal{I}^{\text{syn}}(k)$ defined in (2) and $\mathcal{I}_j^{\text{syn}}(k)$ defined in (3) are changed to

$$\mathcal{I}^{\text{syn}}(k) = \left[ \bar{y}_k^{\text{syn}} - \frac{C}{\sqrt{k}}, \ \bar{y}_k^{\text{syn}} + \frac{C}{\sqrt{k}} \right],$$

$$\mathcal{I}_j^{\text{syn}}(k) = \left[ \bar{y}_{j,k}^{\text{syn}} - \frac{C}{\sqrt{k}}, \ \bar{y}_{j,k}^{\text{syn}} + \frac{C}{\sqrt{k}} \right],$$

respectively, where $C = 2\Phi^{-1}(1 - \alpha/4)$. For simplicity, we suppose that the real datasets have the same size: $n_j = n$ for all $j \in [m]$. Finally, we define

$$\Delta = \sup_{\psi \in \Psi} \left| \mathbb{E}_{z \sim \mathcal{P}} F(z, \psi) - \mathbb{E}_{z^{\text{syn}} \sim \mathcal{P}^{\text{syn}}} F^{\text{syn}}(z^{\text{syn}}, \psi) \right|.$$

In Theorem B.4, the quantity $\Delta$ measures the discrepancy between the distributions of the real students' performance and of the simulated students' performance. The following theorem presents a lower bound on the chosen simulation sample size $\widehat{k}$. Its proof can be found in Section B.4.

**Theorem B.5** (Sharpness of chosen sample size). *Consider the setting of Theorem B.4. Let $\widehat{k}$ be chosen by the procedure (6). Choose $\delta \in (0, 1)$. There exists a constant $C' > 0$ determined by $\alpha$ such that when $m > C' \log(n/\delta)$, the following holds with probability at least $1 - \delta$:*

$$\widehat{k} \geq \min \left\{ K, \ n, \ \left( \frac{C}{5\Delta} \right)^2 \right\}.$$

*When this happens, the selected confidence interval $\mathcal{I}^{\text{syn}}(\widehat{k})$ has width $O\left(\max\{K^{-1/2}, n^{-1/2}, \Delta\}\right)$.*

Theorem B.5 implies that the interval $\mathcal{I}^{\text{syn}}(\widehat{k})$ is the shortest possible. To see this, suppose that the simulation budget $K$ is large, then Theorem B.5 states that with high probability, $\mathcal{I}^{\text{syn}}(\widehat{k})$ has width $O(\max\{\Delta, n^{-1/2}\})$. This is the optimal width because of the following reasons. First, in the worst case, any $\mathcal{I}^{\text{syn}}(k)$ that covers the true mean with high probability must have width $\Omega(\Delta)$ in order to address the distribution shift between the real and simulated responses. Second, as $n$ real human responses can identify the true mean up to an error of $O(n^{-1/2})$, then any valid $\mathcal{I}^{\text{syn}}(k)$ must also have width $\Omega(n^{-1/2})$. This shows the sharpness of the chosen sample size $\widehat{k}$ and the confidence interval $\mathcal{I}^{\text{syn}}(\widehat{k})$.

## B.3. Proof of Theorem B.3

We will prove the following stronger guarantee.

**Lemma B.6** (Conditional coverage). *Consider the setting of Theorem B.3. Let $\delta \in (0, 1)$. With probability at least $1 - \delta$,*

$$\mathbb{P}\left( \mu \in \mathcal{I}^{\text{syn}}(\widehat{k}) \,\Big|\, \widehat{k} \right) \geq 1 - \alpha - \sqrt{\frac{2 \log(1/\delta)}{m}}. \tag{15}$$

By Theorem B.6, we obtain

$$\mathbb{P}\Big(\mu \in \mathcal{I}^{\mathsf{syn}}(\widehat{k})\Big) = \mathbb{E}\left[\mathbb{P}\Big(\mu \in \mathcal{I}^{\mathsf{syn}}(\widehat{k}) \mid \widehat{k}\Big)\right]$$

$$= \int_0^\infty \mathbb{P}\left(\mathbb{P}\Big(\mu \in \mathcal{I}^{\mathsf{syn}}(\widehat{k}) \mid \widehat{k}\Big) > t\right) dt$$

$$\geq \int_0^{1-\alpha} \left[1 - \exp\left(-\frac{m}{2}(t - (1-\alpha))^2\right)\right] dt$$

$$\geq 1 - \alpha - \int_{-\infty}^{1-\alpha} \exp\left(-\frac{m}{2}(t - (1-\alpha))^2\right) dt$$

$$\geq 1 - \alpha - \sqrt{\frac{2}{m}}.$$

We will now prove Theorem B.6. Define $\varepsilon = \sqrt{2\log(1/\delta)/m}$ and a deterministic oracle sample size

$$\bar{k} = \inf\left\{k \in [K] : \mathbb{P}\Big(\mu \notin \mathcal{I}^{\mathsf{syn}}(k)\Big) > \alpha + \varepsilon\right\}. \tag{16}$$

If $\bar{k} = \inf \emptyset$ does not exist, then there is nothing to prove. Now suppose that $\bar{k} \in [K]$ exists. We will prove that with probability at least $1 - \delta$, it holds that $G(\bar{k}) > \alpha/2$. When this event happens, we have $\widehat{k} < \bar{k}$, which implies $\mathbb{P}\big(\mu \notin \mathcal{I}^{\mathsf{syn}}(\widehat{k}) \mid \widehat{k}\big) \leq \alpha + \varepsilon$ and thus (15), thanks to the independence of $\widehat{k}$ and $(\psi, \mathcal{D}^{\mathsf{syn}})$.

By Hoeffding's inequality (e.g., Theorem 2.8 in (Boucheron et al., 2013)) and the conditional independence of $(\mathcal{D}_1, \mathcal{D}_1^{\mathsf{syn}}), ..., (\mathcal{D}_m, \mathcal{D}_m^{\mathsf{syn}})$ given $(\psi_1, ..., \psi_m)$,

$$\mathbb{P}\left(G(\bar{k}) \geq \frac{1}{m}\sum_{j=1}^m \mathbb{P}\Big(\bar{y}_j \notin \mathcal{I}_j^{\mathsf{syn}}(\bar{k})\Big) - \sqrt{\frac{\log(1/\delta)}{2m}}\right) \geq 1 - \delta. \tag{17}$$

We now bound $\mathbb{P}\big(\bar{y}_j \notin \mathcal{I}_j^{\mathsf{syn}}(\bar{k})\big)$. For each $j \in [m]$ and $k \in [K]$,

$$\mathbf{1}\left\{\bar{y}_j \notin \mathcal{I}_j^{\mathsf{syn}}(k)\right\} \geq \mathbf{1}\left\{\bar{y}_j < \mu_j \text{ and } \mu_j < \min \mathcal{I}_j^{\mathsf{syn}}(k)\right\} + \mathbf{1}\left\{\bar{y}_j \geq \mu_j \text{ and } \mu_j > \max \mathcal{I}_j^{\mathsf{syn}}(k)\right\}.$$

By the conditional independence of $\mathcal{D}_j$ and $\mathcal{D}_j^{\mathsf{syn}}$ given $\psi_j$,

$$\mathbb{P}\Big(\bar{y}_j < \mu_j \text{ and } \mu_j < \min \mathcal{I}_j^{\mathsf{syn}}(k)\Big)$$

$$= \mathbb{E}\left[\mathbb{P}\Big(\bar{y}_j < \mu_j \mid \psi_j\Big) \cdot \mathbb{P}\Big(\mu_j < \min \mathcal{I}_j^{\mathsf{syn}}(k) \mid \psi_j\Big)\right]$$

$$= \mathbb{E}\left[\frac{1}{2} \cdot \mathbb{P}\Big(\mu_j < \min \mathcal{I}_j^{\mathsf{syn}}(k) \mid \psi_j\Big)\right]$$

$$= \frac{1}{2}\mathbb{P}\Big(\mu_j < \min \mathcal{I}_j^{\mathsf{syn}}(k)\Big).$$

Similarly,

$$\mathbb{P}\Big(\bar{y}_j \geq \mu_j \text{ and } \mu_j > \max \mathcal{I}_j^{\mathsf{syn}}(k)\Big) = \frac{1}{2}\mathbb{P}\Big(\mu_j > \max \mathcal{I}_j^{\mathsf{syn}}(k)\Big).$$

Therefore,

$$\mathbb{P}\Big(\bar{y}_j \notin \mathcal{I}_j^{\mathsf{syn}}(k)\Big) \geq \frac{1}{2}\left[\mathbb{P}\Big(\mu_j < \min \mathcal{I}_j^{\mathsf{syn}}(k)\Big) + \mathbb{P}\Big(\mu_j > \max \mathcal{I}_j^{\mathsf{syn}}(k)\Big)\right]$$

$$= \frac{1}{2}\mathbb{P}\Big(\mu_j \notin \mathcal{I}_j^{\mathsf{syn}}(k)\Big) = \frac{1}{2}\mathbb{P}\Big(\mu \notin \mathcal{I}^{\mathsf{syn}}(k)\Big), \tag{18}$$

10

where the last equality is due to Assumption B.1. When the event in (17) happens,

$$
G(\bar{k}) \geq \frac{1}{m} \sum_{j=1}^{m} \mathbb{P}\Big(\bar{y}_j \notin \mathcal{I}_j^{\mathsf{syn}}(\bar{k})\Big) - \sqrt{\frac{\log(1/\delta)}{2m}}
$$

$$
\geq \frac{1}{2}\mathbb{P}\Big(\mu \notin \mathcal{I}^{\mathsf{syn}}(\bar{k})\Big) - \sqrt{\frac{\log(1/\delta)}{2m}} \qquad \text{(by (18))}
$$

$$
> \frac{\alpha}{2}. \qquad \text{(by definition of } \bar{k})
$$

This completes the proof.

### B.4. Proof of Theorem B.5

By Hoeffding's inequality (e.g., Theorem 2.8 in (Boucheron et al., 2013)) and a union bound, the following happens with probability at least $1 - \delta$:

$$
G(k) \leq \frac{1}{m} \sum_{j=1}^{m} \mathbb{P}\Big(\bar{y}_j \notin \mathcal{I}_j^{\mathsf{syn}}(k)\Big) + \sqrt{\frac{\log(n/\delta)}{2m}}, \quad \forall k \leq \min\left\{ n,\ K,\ \left(\frac{C}{5\Delta}\right)^2 \right\} \tag{19}
$$

We now show that the right hand side of (19) is at most $\alpha/2$ for $m$ large. For all $j \in [m]$ and $k \in [K]$,

$$
\mathbb{P}\Big(y_j^{\mathsf{syn}} \notin \mathcal{I}_j^{\mathsf{syn}}(k)\Big) = \mathbb{P}\left( |y_j^{\mathsf{syn}} - \bar{y}_{j,k}^{\mathsf{syn}}| > \frac{C}{\sqrt{k}} \right).
$$

Since $\bar{y}_{j,k}^{\mathsf{syn}} \sim N(\mu_j^{\mathsf{syn}}, 1/k)$ and $\bar{y}_j \sim N(\mu_j, 1/n)$, then

$$
y_j^{\mathsf{syn}} - \bar{y}_{j,k}^{\mathsf{syn}} \sim N\left( \mu_j - \mu_j^{\mathsf{syn}},\ \frac{1}{k} + \frac{1}{n} \right).
$$

When $k \leq \min\{n, K,\ (\frac{C}{5\Delta})^2\}$, we have

$$
\mathbb{P}\Big(\bar{y}_j \notin \mathcal{I}_j^{\mathsf{syn}}(k)\Big) \leq \mathbb{P}\left( |(\bar{y}_j - \bar{y}_{j,k}^{\mathsf{syn}}) - (\mu_j - \mu_j^{\mathsf{syn}})| + \Delta > \frac{C}{\sqrt{k}} \right)
$$

$$
= 2\Phi\left( -\frac{C/\sqrt{k} - \Delta}{\sqrt{k^{-1} + n^{-1}}} \right) \leq 2\Phi\left( -\frac{4C/(5\sqrt{k})}{\sqrt{k^{-1} + k^{-1}}} \right) \qquad (\Delta \leq C/(5\sqrt{k}))
$$

$$
= 2\Phi\left( -\frac{2\sqrt{2}C}{5} \right) = 2\Phi\left( -\frac{4\sqrt{2}\Phi^{-1}(1 - \alpha/4)}{5} \right) = 2\Phi\left( \frac{4\sqrt{2}\Phi^{-1}(\alpha/4)}{5} \right) < \frac{\alpha}{2}.
$$

Let

$$
\xi = \frac{\alpha}{2} - 2\Phi\left( \frac{4\sqrt{2}\Phi^{-1}(\alpha/4)}{5} \right),
$$

then $\xi > 0$ and

$$
\mathbb{P}\Big(\bar{y}_j \notin \mathcal{I}_j^{\mathsf{syn}}(k)\Big) \leq \frac{\alpha}{2} - \xi. \tag{20}
$$

When $m > \log(n/\delta)/(2\xi^2)$, substituting (20) into (19) yields that for all $k \leq \min\{n,\ K,\ (\frac{C}{5\Delta})^2\}$,

$$
G(k) \leq \frac{1}{m} \sum_{j=1}^{m} \mathbb{P}\Big(\bar{y}_j \notin \mathcal{I}_j^{\mathsf{syn}}(k)\Big) + \sqrt{\frac{\log(n/\delta)}{2m}}
$$

$$
< \left( \frac{\alpha}{2} - \xi \right) + \xi = \frac{\alpha}{2}.
$$

When this happens, we have $\widehat{k} \geq \min\{n,\ K,\ (\frac{C}{5\Delta})^2\}$.

11

# C. More Details of the General Method

In this section, we provide more details of the general method in Section 4.

## C.1. Examples of the General Problem Formulation

**Example C.1** (Educational test evaluation). In the educational test evaluation example in Section 2.1, each $z \in \mathcal{Z}$ is a student profile, each $\psi \in \Psi$ is a test question, the response space is $\mathcal{Y} = \{0, 1\}$, and $\mathcal{Q}(\cdot \mid z, \psi) = \text{Bernoulli}(F(z, \psi))$. The statistic $\theta(\psi)$ is the probability of a student answering the question correctly: $\mathbb{E}_{y \sim \mathcal{R}(\cdot|\psi)}[y] = \mathbb{E}_{z \sim \mathcal{P}}[F(z, \psi)]$.

**Example C.2** (Market research). Suppose a company is interested in learning its customers' willingness-to-pay (WTP) for a new product, which is the highest price a customer is willing to pay for the product. Then, each $z \in \mathcal{Z}$ can represent a customer profile (e.g., age, gender, occupation), each survey question $\psi$ is about a certain product, and a customer's response $y$ is a noisy observation of the customer's WTP. Then $\mathcal{R}(\cdot \mid \psi)$ is the distribution of the customer population's WTP. We may take $\theta(\psi)$ as the $\tau$-quantile of the WTP distribution $\mathcal{R}(\cdot \mid \psi)$, for some $\tau \in (0, 1)$:

$$\theta(\psi) = \inf \left\{ q \in [0, \infty) : \mathbb{P}_{y \sim \mathcal{R}(\cdot|\psi)}(y \leq q) \geq \tau \right\}.$$

An LLM can be used to simulate customers' WTP for the product.

**Example C.3** (Public survey). Suppose an organization is interested in performing a public survey in a city. Each survey question $\psi$ is a multiple-choice question with 5 options. An example is "How often do you talk to your neighbors?", with 5 choices "Basically every day", "A few times a week", "A few times a month", "Once a month", and "Less than once a month". Every $z \in \mathcal{Z}$ is a person's profile (e.g., age, gender, occupation), the response space $\mathcal{Y}$ is the standard orthonormal basis $\{e_i\}_{i=1}^5$ in $\mathbb{R}^5$, where $y = e_i$ indicates that a person chooses the $i$-th option. We can take $\theta(\psi) = \mathbb{E}_{y \sim \mathcal{Q}(\cdot|\psi)}[y] \in \mathbb{R}^5$, which summarizes the proportion of people that choose the $i$-th option. An LLM can be used to simulate people's answers to the survey question.

## C.2. Theoretical Analysis of the General Method

We now present the coverage guarantee for our method, which shows that the chosen confidence set $\mathcal{S}^{\text{syn}}(\widehat{k})$ has coverage probability at least $1 - \alpha - O(1/\sqrt{m})$. Its proof is deferred to Section C.3.

**Theorem C.4.** *Let Assumptions B.1 and B.2 hold. Fix $\alpha \in (0, 1)$. The sample size $\widehat{k}$ defined by* (13) *satisfies*

$$\mathbb{P}\left( \theta(\psi) \in \mathcal{S}^{\text{syn}}(\widehat{k}) \right) \geq 1 - \alpha - \gamma^{-1} \sqrt{\frac{1}{2m}}.$$

It is worth noting that our method achieves this coverage without any assumptions on the qualities of the LLM and the procedure $\mathcal{C}$ for confidence set construction. Nevertheless, the size of the chosen confidence set $\mathcal{S}^{\text{syn}}(\widehat{k})$, in terms of the true coverage rate and size, depend on these factors. If there is a large alignment gap between the LLM and the human population, then $\mathcal{S}^{\text{syn}}(\widehat{k})$ will inevitably be large.

## C.3. Proof of Theorem C.4

We will prove the following stronger guarantee.

**Lemma C.5** (Conditional coverage). *Consider the setting of Theorem C.4. Let $\delta \in (0, 1)$. With probability at least $1 - \delta$,*

$$\mathbb{P}\left( \theta(\psi) \in \mathcal{S}^{\text{syn}}(\widehat{k}) \,\Big|\, \widehat{k} \right) \geq 1 - \alpha - \gamma^{-1} \sqrt{\frac{\log(1/\delta)}{2m}}. \tag{21}$$

By Theorem B.6, we obtain

$$\mathbb{P}\Big(\theta(\psi) \in \mathcal{S}^{\mathsf{syn}}(\widehat{k})\Big) = \mathbb{E}\left[\mathbb{P}\Big(\theta(\psi) \in \mathcal{S}^{\mathsf{syn}}(\widehat{k}) \,\Big|\, \widehat{k}\Big)\right]$$

$$= \int_0^\infty \mathbb{P}\left(\mathbb{P}\Big(\theta(\psi) \in \mathcal{S}^{\mathsf{syn}}(\widehat{k}) \,\Big|\, \widehat{k}\Big) > t\right) dt$$

$$\geq \int_0^{1-\alpha} \left[1 - \exp\left(-2m\gamma^2(t - (1-\alpha))^2\right)\right] dt$$

$$\geq 1 - \alpha - \int_{-\infty}^{1-\alpha} \exp\left(-2m\gamma^2(t - (1-\alpha))^2\right) dt$$

$$= 1 - \alpha - \sqrt{\frac{\pi}{8}} \cdot \gamma^{-1} \sqrt{\frac{1}{m}}$$

$$\geq 1 - \alpha - \gamma^{-1}\sqrt{\frac{1}{2m}}.$$

We now prove Theorem C.5. Define $\varepsilon = \gamma^{-1}\sqrt{\frac{\log(1/\delta)}{2m}}$ and a deterministic oracle sample size

$$\bar{k} = \inf\left\{k \in [K] : \mathbb{P}\Big(\theta(\psi) \notin \mathcal{S}^{\mathsf{syn}}(k)\Big) > \alpha + \varepsilon\right\}. \tag{22}$$

If $\bar{k} = \inf \emptyset$ does not exist, then there is nothing to prove. Now suppose $\bar{k} \in [K]$ exists. We will prove that with probability at least $1 - \delta$, it holds that $L(\bar{k}) > \gamma\alpha$. When this event happens, we have $\widehat{k} < \bar{k}$, which implies $\mathbb{P}\big(\theta(\psi) \notin \mathcal{S}^{\mathsf{syn}}(\widehat{k}) \mid \widehat{k}\big) \leq \alpha + \varepsilon$ and thus (21), thanks to the independence of $\widehat{k}$ and $(\psi, \mathcal{D}^{\mathsf{syn}})$.

By Hoeffding's inequality (e.g., Theorem 2.8 in (Boucheron et al., 2013)) and the conditional independence of $(\mathcal{D}_1, \mathcal{D}_1^{\mathsf{syn}}), ..., (\mathcal{D}_m, \mathcal{D}_m^{\mathsf{syn}})$ given $(\psi_1, ..., \psi_m)$,

$$\mathbb{P}\left(L(\bar{k}) \geq \frac{1}{m}\sum_{j=1}^m \mathbb{P}\Big(\mathcal{S}_j \not\subseteq \mathcal{S}_j^{\mathsf{syn}}(\bar{k})\Big) - \sqrt{\frac{\log(1/\delta)}{2m}}\right) \geq 1 - \delta. \tag{23}$$

We now bound $\mathbb{P}\big(\mathcal{S}_j \not\subseteq \mathcal{S}_j^{\mathsf{syn}}(\bar{k})\big)$. For each $j \in [m]$ and $j \in [K]$,

$$\mathbf{1}\left\{\mathcal{S}_j \not\subseteq \mathcal{S}_j^{\mathsf{syn}}(k)\right\} \geq \mathbf{1}\left\{\theta(\psi_j) \in \mathcal{S}_j \text{ and } \theta(\psi_j) \notin \mathcal{S}_j^{\mathsf{syn}}(k)\right\}.$$

By the conditional independence of $\mathcal{D}_j$ and $\mathcal{D}_j^{\mathsf{syn}}$ given $\psi_j$,

$$\mathbb{P}\Big(\mathcal{S}_j \not\subseteq \mathcal{S}_j^{\mathsf{syn}}(k)\Big) \geq \mathbb{E}\left[\mathbb{P}\Big(\theta(\psi_j) \in \mathcal{S}_j \text{ and } \theta(\psi_j) \notin \mathcal{S}_j^{\mathsf{syn}}(k) \,\Big|\, \psi_j\Big)\right]$$

$$= \mathbb{E}\left[\mathbb{P}\Big(\theta(\psi_j) \in \mathcal{S}_j \,\Big|\, \psi_j\Big) \cdot \mathbb{P}\Big(\theta(\psi_j) \notin \mathcal{S}_j^{\mathsf{syn}}(k) \,\Big|\, \psi_j\Big)\right]$$

$$\geq \mathbb{E}\left[\gamma \cdot \mathbb{P}\Big(\theta(\psi_j) \notin \mathcal{S}_j^{\mathsf{syn}}(k) \,\Big|\, \psi_j\Big)\right]$$

$$= \gamma \cdot \mathbb{P}\Big(\theta(\psi_j) \notin \mathcal{S}_j^{\mathsf{syn}}(k)\Big)$$

$$= \gamma \cdot \mathbb{P}\Big(\theta(\psi) \notin \mathcal{S}^{\mathsf{syn}}(k)\Big), \tag{24}$$

where the last equality is due to Assumption B.1. Therefore, when the event in (23) happens,

$$L(\bar{k}) \geq \frac{1}{m} \sum_{j=1}^{m} \mathbb{P}\Big(\mathcal{S}_j \not\subseteq \mathcal{S}_j^{\mathsf{syn}}(\bar{k})\Big) - \sqrt{\frac{\log(1/\delta)}{2m}}$$

$$\geq \gamma \cdot \mathbb{P}\Big(\theta(\psi) \notin \mathcal{S}^{\mathsf{syn}}(\bar{k})\Big) - \sqrt{\frac{\log(1/\delta)}{2m}} \qquad \text{(by (24))}$$

$$> \gamma\alpha. \qquad \text{(by definition of } \bar{k})$$

This completes the proof.

## D. Details of Numerical Experiments

### D.1. The OpinionQA Dataset

**Selection of survey questions.** The original dataset is categorized into topics such as health, crime/security, and political issues. Ideally, we would want to consider questions from the same category to ensure that they are similar enough. However, the category with the most questions has fewer than 200 questions, which is not suited for demonstration. We thus consider all questions. In total, the dataset has 1,442 survey questions, which is too large for our computational resources. We thus selected a subset of questions. First, while the number of choices ranges from 2 to 19, most questions have 5 choices. To give a fair comparison and for simplicity, we therefore only consider questions with 5 choices. Second, not all questions have choices that can be clearly ordered in sentiments; the following question is an example:

*Who do you think has the most responsibility to reduce the amount of made-up news and information? 1. The government, 2. Technology companies, 3. The public, 4. The news media, 5. None of these, 6. Refused.*

In contrast, all example questions in the next section have choices that can be clearly ordered in sentiments. To streamline the process, we ask GPT-4o to determine if a question's choices can be ordered in sentiments and we keep those that have GPT-4o's affirmative answer. This leaves us with 546 questions. To compensate for the loss of similarity by pooling questions across various topics and to further reduce our computational cost, we selected 400 questions that are "most similar" to each other by embedding the question statements using OpenAI's `text-embedding-3-small`, calculating the mean, and selecting the 400 questions with the smallest Euclidean distance to the mean. Out of these 400 questions, 15 questions have various issues with their choices by manual inspection, so we exclude them. This leaves us with 385 questions. All these questions happen to have at least 400 responses.

**Example questions.** The questions in the OpinionQA dataset span a wide range of topics, including health, crime/security, and political issues. Some example questions are as follows:

- *How much, if at all, do you think wages and incomes are contributing to your opinion about how the economy is doing?*

  *1. A great deal    2. A fair amount    3. Not too much    4. Not at all    5. Refused*

- *Regardless of whether you would want to move, how likely is it that you will move to a different community at some point in the future?*

  *1. Very likely    2. Somewhat likely    3. Not too likely    4. Not at all likely    5. Refused*

- *How much, if anything, would you be willing to change about how you live and work to help reduce the effects of global climate change? Would you be willing to make:*

  *1. A lot of changes  2. Some changes  3. Only a few changes  4. No changes at all  5. Refused*

**Profiles.** Excluding surveyees with missing information, each of the 385 questions we consider has at least 400 responses. Since there was no information on the surveyees' identification, by dropping repeated profiles we can only say that there are at least 32,864 surveyees. Each surveyee is described by 12 features. Their corresponding categories are listed in Table 1.

*Table 1.* Categories of surveyees' features in the OpinionQA dataset.

| Feature | Options |
|---|---|
| US citizenship | 'Yes', 'No' |
| Region | 'Northeast', 'Midwest', 'South', 'West' |
| Sex | 'Male', 'Female' |
| Age | '18-29', '30-49', '50-64', '65+' |
| Marital status | 'Married', 'Divorced', 'Separated', 'Widowed', 'Never been married' |
| Race | 'White', 'Black', 'Asian', 'Hispanic', 'Other' |
| Educational background | 'Less than high school', 'High school graduate', 'Some college, no degree', 'Associate's degree', 'College graduate/some postgrad', 'Postgraduate' |
| Income | 'Less than $30,000', '$30,000-$50,000', '$50,000-$75,000', '$75,000-$100,000', '$100,000 or more' |
| Religious affiliation | 'Protestant', 'Roman Catholic', 'Mormon', 'Orthodox', 'Jewish', 'Muslim', 'Buddhist', 'Hindu', 'Atheist', 'Agnostic', 'Other', 'Nothing in particular' |
| Religious attendance | 'More than once a week', 'Once a week', 'Once or twice a month', 'A few times a year', 'Seldom', 'Never' |
| Political party | 'Republican', 'Democrat', 'Independent', 'Other' |
| Political ideology | 'Very conservative', 'Conservative', 'Moderate', 'Liberal', 'Very liberal' |

**Synthetic response generation.** We generate synthetic profiles by bootstrapping the 32,864 unique real profiles. We then generate synthetic answers by prompting LLMs to pretend that they are a surveyee with the synthetic profile and answer the question. An example prompt is as follows:

> *Pretend that you reside in the US and you are a US citizen from the West region of the country. You are female, your age is between 18 and 29, and you are single. In terms of race, you are white. In terms of education, you attended college but did not graduate. Your annual income is less than $30,000. Religion-wise, you do not belong to any particular religion, and you never attend religious services. Politically, you are affiliated with a political party that is not Democratic or Republican, and you consider your political ideology to be liberal. Please answer the following question:*
>
> *How much, if at all, do you think what happens to black people in the country overall affects what happens in your own life? ['1. A lot', '2. Some', '3. Not much', '4. Not at all', '5. Refused'].*
>
> *Please provide your answer choice (a single number from 1 to 5) in double square brackets.*

While LLMs generally provide explanations for the math questions in EEDI dataset, they usually directly provide answers for the OpinionQA dataset; e.g., '[[2]]'.

### D.2. The EEDI Dataset

**Example questions.** Some example questions from the EEDI dataset are as follows:

- *What number belongs in the box?* $\Box + 7 = 2$

  *A) 9    B) -5    C) -6    D) 5*

- *If you multiply a square number by* 9*, you get a square number. Is this statement:*

  *A) always true    B) sometimes true    C) never true    D) impossible to say*

- *Which calculation is equal to* $-20$*?*

  *A)* $2 \times (-2) - (-4) \times 4$    *B)* $-28 - (-4) \times 2$    *C)* $(-5)\top 2 + 5$    *D)* $(-42) \div (-2) + 1$

**Profile distribution.** Excluding students with missing information which take up less than 10% of the total population, there are 2,111 students who answered at least one of the 412 questions. Each student is described by three features: gender, age, and whether or not they are eligible for free school meals or premium pupil. Gender is represented by 1 or 2, where 1 corresponds to female and 2 corresponds to male. The students' ages are rounded to integers from 11 and 14. Whether or not a student is eligible for free school meals is represented by 0 or 1, where 0 corresponds to not eligible and 1 corresponds to eligible. The distribution of these students' features is presented in Table 2.

*Table 2.* Summary statistics of students' features in the EEDI dataset.

|  | min | max | mean | median | standard deviation |
|---|---|---|---|---|---|
| Gender | 1 | 2 | 1.4988 | 1 | 0.5001 |
| Age | 11 | 14 | 11.2776 | 11 | 0.4696 |
| Premium Pupil | 0 | 1 | 0.2842 | 0 | 0.4512 |

**Synthetic response generation.** For each question, we generate synthetic profiles by sampling with replacement from the real profiles. We then generate synthetic answers by prompting LLMs to pretend that they are a student with the synthetic profile and answer the question. We adapted the prompt from (He-Yueya et al., 2024) with slight modifications to reduce computational cost. An example prompt featuring an 11-year-old boy who is not eligible for free school meals is as follows:

> *Pretend that you are an 11-year-old student. Your gender is male. You are not eligible for free school meals or pupil premium due to being relatively financially advantaged. Given your characteristics, is it likely that you would be able to solve the following problem?*
>
> *Problem: [Insert question here]*
>
> *If yes, put the final answer choice (a single letter) in double square brackets. If you are likely to struggle with this problem, put a plausible incorrect answer choice (a single letter) in double square brackets.*

An example answer from GTP-4o when given the second example question above is as follows:

> *As an 11-year-old student, I might have learned about square numbers and multiplication in school. However, the problem may be a bit tricky if I haven't thought about how multiplying square numbers by other numbers can also result in square numbers. I might not immediately realize that 9 is actually a square number itself (3 squared), which makes this property more evident.*
>
> *Considering this, I could find the reasoning challenging and decide based on a misconception. I might go with a plausible incorrect answer choice like [[B]] because I might think that it's only sometimes possible without realizing the full mathematical principle involved.*

### D.3. Experiment Procedure

We now describe more formally our experiment procedure for applying our method to each dataset. Denote the dataset by $\{(\psi_j, \mathcal{D}_j)\}_{j=1}^J$, where $\psi_j$ is a survey question and $\mathcal{D}_j = \{y_{j,i}\}_{i=1}^{n_j}$ is a collection of human responses. For each $j \in [J]$, we simulate $K$ responses $\mathcal{D}_j^{\mathsf{syn}}$ from an LLM. We then randomly split $\mathscr{D} = \{(\mathcal{D}_j, \mathcal{D}_j^{\mathsf{syn}})\}_{j=1}^J$ into a training set $\mathscr{D}^{\mathrm{tr}} = \{(\mathcal{D}_j, \mathcal{D}_j^{\mathsf{syn}})\}_{j \in \mathcal{J}_{\mathrm{tr}}}$ and a testing set $\mathscr{D}^{\mathrm{te}} = \{(\mathcal{D}_j, \mathcal{D}_j^{\mathsf{syn}})\}_{j \in \mathcal{J}_{\mathrm{te}}}$, with $|\mathscr{D}^{\mathrm{tr}}| : |\mathscr{D}^{\mathrm{te}}| = 3 : 2$.

**Selection of simulation sample size.** We apply our method in (6) with the training set $\mathscr{D}^{\mathrm{tr}}$ to select a simulation sample size $\widehat{k}$.

**Evaluation of selected sample size.** We use $\mathscr{D}^{\mathrm{te}}$ to evaluate the quality of the chosen simulation sample size $\widehat{k}$. As the true population mean $\mu$ is unavailable, the true coverage probability $\mathbb{P}\big(\mu \in \mathcal{I}^{\mathsf{syn}}(\widehat{k})\big)$ cannot be computed. However, we can apply the same idea as (5) in Section 3 to compute a proxy for the miscoverage level. For each survey question $j \in \mathcal{J}_2$, the selected sample size $\widehat{k}$ leads to the synthetic confidence set $\mathcal{I}_j^{\mathsf{syn}}(\widehat{k}) = \mathcal{C}\big(\{y_{j,i}^{\mathsf{syn}}\}_{i=1}^{\widehat{k}}\big)$. We form the sample mean $\bar{y}_j$ from real data $\mathcal{D}_j$ and define

$$\widetilde{G}(k) = \frac{2}{|\mathcal{J}_{\mathrm{te}}|} \sum_{j \in \mathcal{J}_{\mathrm{te}}} \mathbf{1}\big\{\bar{y}_j \notin \mathcal{I}_j^{\mathsf{syn}}(k)\big\}. \tag{25}$$

16

The proof of Theorem B.3 shows that, for every $k \in [K]$ and survey question $j$,

$$\frac{1}{2}\mathbb{P}\big(\mu \notin \mathcal{I}^{\mathsf{syn}}(k)\big) \leq \mathbb{P}\big(\bar{y}_j \notin \mathcal{I}_j^{\mathsf{syn}}(k)\big) = \frac{1}{2}\mathbb{E}\big[\widetilde{G}(k)\big].$$

Thus, if $\mathbb{E}\big[\widetilde{G}(\widehat{k})\big] \leq \alpha$, then $\mathbb{P}\big(\mu \notin \mathcal{I}^{\mathsf{syn}}(\widehat{k})\big) \leq \alpha$ must hold. To this end, we will test a hypothesis $H_0 : \mathbb{E}\big[\widetilde{G}(\widehat{k})\big] \leq \alpha$ against its alternative $H_1 : \mathbb{E}\big[\widetilde{G}(\widehat{k})\big] > \alpha$.

## E. Experiment Results on Coverage Validity

In Figure 2, we present histograms of $p$-values for the hypothesis test $\mathbb{E}\big[\widetilde{G}(\widehat{k})\big] \leq \alpha$ against $\mathbb{E}\big[\widetilde{G}(\widehat{k})\big] > \alpha$ across various LLMs and $\alpha$'s over the OpinionQA and EEDI datasets. The $p$-values are computed using a one-sided $z$-test over the 100 random splits. As can be seen from the histograms, all $p$-values are reasonably large, indicating that the hypothesis $\mathbb{E}\big[\widetilde{L}(\widehat{k})\big] \leq \alpha$ cannot be rejected (e.g., at the 0.05 significance level) for any LLM and $\alpha$ across both datasets. These experiment results verify the theoretical guarantees in Theorem B.3, showing that the miscoverage rate is effectively controlled by our method.
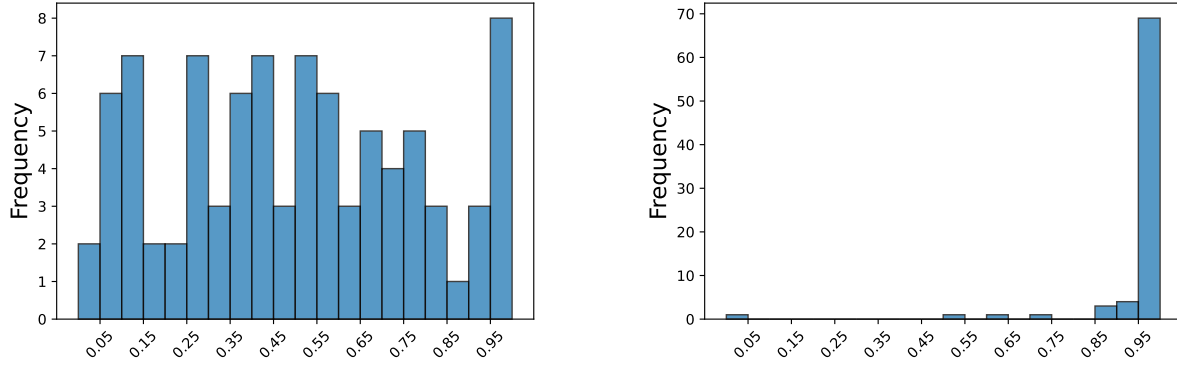


*Figure 2.* Histograms of $p$-values for the hypothesis test $\mathbb{E}\big[\widetilde{G}(\widehat{k})\big] \leq \alpha$ against $\mathbb{E}\big[\widetilde{G}(\widehat{k})\big] > \alpha$ across various LLMs and $\alpha$'s over the OpinionQA dataset (left) and the EEDI dataset (right).