

Model-Free Assessment of Simulator Fidelity via Quantile Curves

Garud Iyengar*

Yu-Shiou Willy Lin*

Kaizheng Wang*

This version: December 5, 2025

Abstract

Simulation of complex systems originated in manufacturing and queuing applications. It is now widely used for large-scale, ML-based systems in research, education, and consumer surveys. However, characterizing the discrepancy between simulators and ground truth remains challenging for increasingly complex, machine-learning-based systems. We propose a computationally tractable method to estimate the quantile function of the discrepancy between the simulated and ground-truth outcome distributions. Our approach focuses on output uncertainty and treats the simulator as a black box, imposing no modeling assumptions on its internals, and hence applies broadly across many parameter families, from Bernoulli and multinomial models to continuous, vector-valued settings. The resulting quantile curve supports confidence interval construction for unseen scenarios, risk-aware summaries of sim-to-real discrepancy (e.g., VaR/CVaR), and comparison of simulators' performance. We demonstrate our methodology in an application assessing LLM simulation fidelity on the *WorldValueBench* dataset spanning four LLMs.

Keywords: Simulation, Quantile function estimation, Human-AI alignment, Conformal inference, Output Uncertainty Quantification

1 Introduction

The adoption of simulators across operations and manufacturing, agent-based modeling in social-science research, user surveys, and education/training (Zhang et al., 2019; Macal, 2016; Argyle et al., 2023; Aher et al., 2023) has accelerated with recent advances in artificial intelligence (AI). At the platform level, industry ecosystems such as NVIDIA Omniverse and Earth-2 exemplify the push toward high-fidelity, AI-enabled digital twins. Large language models (LLMs) are being applied to create and improve the ability of generative AI agents that can replicate human responses (Park et al., 2023; Lu et al., 2025). Collectively, these developments fuel the sustained growth in simulation across domains.

Against this backdrop, it becomes increasingly important to understand the discrepancy between simulation outputs and the real-world. Such discrepancies are documented across many domains: the behavior of LLMs, model discrepancy in computing systems, and sim-to-real gap in robotics (Gao et al., 2025; Durmus et al., 2023; Tobin et al., 2017; Peng et al., 2018). In light of these developments, we provide a finite-sample guaranteed distribution of the sim-to-real discrepancy without imposing any structural assumptions on the simulator.

Related Literature The sim-to-real discrepancy has been the focus of the uncertainty quantification (UQ) literature. Figure 1 positions our work within this landscape. Following the taxonomy of

*Department of IEOR and Data Science Institute, Columbia University. Emails: garud@ieor.columbia.edu, y15782@columbia.edu, kaizheng.wang@columbia.edu.

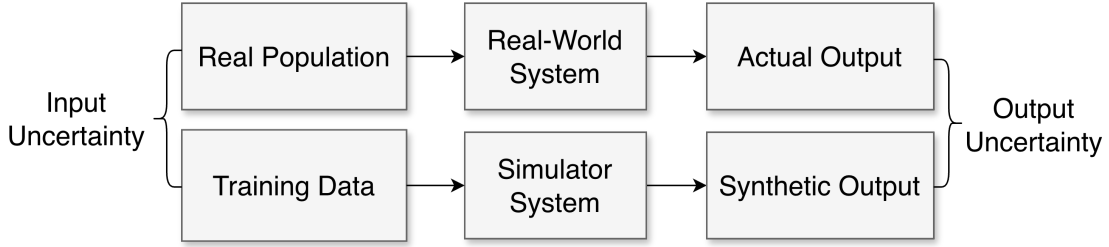


Figure 1: Simulation Uncertainty Quantification.

Roy and Oberkamp (2011), we distinguish *input* uncertainty from *output* uncertainty. The work on input uncertainty (see Chen et al., 2024; Barton et al., 2014; Lam, 2022) assumes that the simulator is accurate in that it creates a faithful representation of the real-world when the inputs follow the correct distribution. The goal here is to characterize how errors in characterizing the input distribution from a finite amount of noisy input-output data propagate to errors in the output. While the initial work on input uncertainty focused on building confidence intervals for functionals of the output, recent work by Chen et al. (2024) constructs Kolmogorov-Smirnov-type uncertainty bands for the CDF of entire output discrepancy. In contrast, modern digital twins rely on complex ML models (e.g., an LLM or a deep neural network), and one does not have access to the model directly, and moreover, it is computationally prohibitive to calibrate its parameters exactly. Consequently, we treat the simulator as a black box and focus on directly characterizing the sim-to-real discrepancy of the output, an approach termed *output uncertainty quantification* by Jeon et al. (2024). Moreover, the classical input uncertainty literature is primarily concerned with Monte Carlo noise, which is the variance induced by stochastic sampling conditional on a fixed input model. Our work, by contrast, explicitly incorporates both variance and bias of the simulator relative to the real world, which is precisely the focus of “sim-to-real” discrepancy. A further distinction of our work is that much of the UQ literature is asymptotic, whereas we provide finite-sample guarantees, which is critical when only a limited number of scenarios are available.

Recent work in LLM literature has also begun to focus on output uncertainty using model-agnostic discrepancy measures, e.g. Santurkar et al. (2023a) aggregate the sim-to-real gap into a single scalar summary of overall error or bias, and Huang et al. (2025) bound the discrepancy at a fixed quantile level. These approaches correspond to evaluating particular functionals of the distribution of the sim-to-real discrepancy. To obtain more flexible and informative diagnostics, we propose a model-free procedure that approximates the entire quantile function of the discrepancy. Our approach is close in spirit to conformal-inference methods (Vovk et al., 2005; Bates et al., 2021) that obtain distribution-free guarantees for a black-box predictor. However, classical conformal methods are designed for pointwise coverage at a given input rather than characterizing the distribution of the sim-to-real discrepancy. Recent conformal variants that target distributional objects (Snell et al., 2022; Budde et al., 2025) still assume more homogeneous data structures and do not address the heterogeneous sample sizes. Our framework fills this gap by providing finite-sample, distribution-level guarantees tailored specifically to assessing black-box simulators against heterogeneous real-world data.

Contributions Ideally, one wants to recover the quantile function for a specified discrepancy measure between the real-world and the simulated distributions. With the quantile function, we can directly compute the confidence intervals for the estimators, and risk summaries, e.g. VaR/CVaR. In

practice, we only observe finite samples of the real-world and simulated outcomes, and therefore, the true quantile function is not computable. We calibrate a model-free conservative estimate for the quantile function, and prove that our estimate comes with finite-sample guarantees for any desired α -quantile.

Building on this, we extend our framework to pairwise comparison of simulators. This allows us to benchmark a candidate simulator against a reference and to make statistically robust claims about which one is closer to reality. For instance, our procedure can certify that a given simulator achieves a smaller sim-to-real gap on at least 90% of questions and pairwise hypothesis testing claims between simulators. We also illustrate our methodology on the *WorldValueBench* (Zhao et al., 2024) dataset, which is constructed from the World Values Survey (Haerpfer et al., 2020). We generate synthetic response profiles and compare the sim-to-real discrepancy of four different LLMs relative to the survey data.

Our main contributions are twofold:

1. Our procedure is *model-free* in that it does not impose any parametric assumptions on either the simulator or the ground truth; therefore, it can be applied to any black-box simulators. Moreover, it provides a novel assessment method to compare between simulators.
2. Since our procedure yields a *quantile function*, it naturally supports a wide range of summaries, from means to tail-risk measures such as Conditional Value at Risk (CVaR). In addition, by leveraging the simulated estimator, we can construct confidence intervals for the real-world parameters at any significance level α .

The remainder of the paper proceeds as follows. In Section 2 we present a motivating example, and formally state the problem. In Section 3 we present the main theoretical results, and in Section 4 we discuss an application of our methodology. We also address the tightness of our methodology to the actual quantile curve in Section 5. In Section 6 we conclude by discussing future directions.

2 Framework and Motivating Example

In this section, we start with a concrete use case that motivates our formulation, and then formally define the quantile estimation problem. Although the terminology below is specific to the example use case, the setting can be generalized. See Appendix D for more examples of simulation systems.

Suppose a media research company plans to survey customers on a particular topic with multinomial outcomes (eg, “Agree”, “Neutral”, “Disagree”). The company wants to estimate the customer opinion before committing resources for an expensive population study. The company has access to a database of past questions and human answers, and has been training an LLM-based “digital twin” for its customer base. By querying this digital twin with the new question, the company can generate an estimate of the mean response for the new question. The problem now facing the company is to characterize how close this estimate is to the true population mean. Can one provide a confidence interval for the true value? Or, better yet, estimate the quantile function for a suitable discrepancy measure between the simulator and population estimates? This work addresses precisely these uncertainty quantification problems.

To formally describe the discrepancy between simulator (LLM) and ground truth (human population), and the theoretical challenges we face, we consider two levels of randomness in this problem, which we term as the scenario uncertainty and the finite-sample uncertainty.

First, scenarios (questions) are drawn as $\psi \sim \Psi$, where Ψ is a scenario pool distribution, and the past dataset consists of a total of m scenarios, $\{\psi_j\}_{j=1}^m \sim \Psi$. A real system (human) is characterized by a latent profile $z \in \mathcal{Z}$ with a population distribution \mathcal{P} over the latent profile

state \mathcal{Z} . For each pair (ψ, z) , the real system (human population) produces a categorical outcome Y^{gt} , with conditional distribution $Q^{\text{gt}}(\cdot \mid z, \psi)$ over a space $\mathcal{Y} := \{y \in \mathbb{N}^d : \sum_{i=1}^d y_i = 1\}$. The simulator (LLM) produces an outcome Y^{sim} with conditional distribution $Q^{\text{sim}}(z^{\text{sim}}, \psi, r)$ over the same outcome space \mathcal{Y} , where the $z^{\text{sim}} \in \mathcal{Z}_{\text{sim}}$ denotes a synthetic profile drawn i.i.d. from a simulator population distribution \mathcal{P}^{sim} on \mathcal{Z}_{sim} , and r denotes LLM settings, including prompting strategy, hyperparameters, and other API settings.¹ In the running example, $Q^{\text{gt}}(\cdot \mid z, \psi) = \text{Categorical}(\Pi^{\text{gt}}(z, \psi))$, where $\Pi^{\text{gt}}(\psi, z) := (Q^{\text{gt}}(\{1\} \mid z, \psi), \dots, Q^{\text{gt}}(\{d\} \mid z, \psi)) \in \mathcal{P}^d$ and \mathcal{P}^d is the $d-1$ -simplex defined $:= \{u \in [0, 1]^d : \sum_{i=1}^d u_i = 1\}$. For any question ψ , we can marginalize the population effect, hence denote by $Q^{\text{gt}}(\cdot \mid \psi) = \mathbb{E}_{z \sim \mathcal{P}}[Q^{\text{gt}}(\cdot \mid z, \psi)]$, the conditional distribution of outcome Y^{gt} given ψ .

Let $p(\psi)$ be a population statistic of interest, which is a functional of the conditional distribution $Q^{\text{gt}}(\cdot \mid \psi)$ and lives in a parameter space Θ . In our running example, $\Theta = \mathcal{P}^d$, $p(\psi)$ is the mean response of survey respondents, and simulator can be defined similarly under the simulator population \mathcal{P}^{sim} . Further examples of $p(\psi)$ are given in Section 3. To simplify notation, we will write $p(\psi)$ and $q(\psi)$ as p_ψ and q_ψ . Summing up:

$$\begin{aligned} p_\psi &:= p(\psi) := \mathbb{E}_{y \sim Q^{\text{gt}}}[y] = \mathbb{E}_{z \sim \mathcal{P}}[\Pi^{\text{gt}}(\psi, z)] \in \Theta, \\ q_\psi &:= q(\psi) := \mathbb{E}_{z \sim \mathcal{P}^{\text{sim}}}[\Pi^{\text{sim}}(\psi, z)] \in \Theta. \end{aligned}$$

Second, we only observe finite samples per scenario. For each $j \in [m]$, we are given n_j i.i.d. profiles $z_{j,1:n_j} \sim \mathcal{P}$, with which we observe n_j ground-truth outcomes $y_{j,i}^{\text{gt}} \sim Q^{\text{gt}}(\cdot \mid z_{j,i}, \psi_j)$. Comparably, we generate k simulator outcomes $y_{j,\ell}^{\text{sim}} \sim Q^{\text{sim}}(\cdot \mid z_{j,\ell}^{\text{sim}}, \psi_j, r)$ using a simulation pool $z_{j,1:k}^{\text{sim}} \sim \mathcal{P}^{\text{sim}}$ with fixed k across j to standardize simulator sampling. For brevity, we will write $p(\psi_j)$ and $q(\psi_j)$ as p_j and q_j , and \hat{p}_j and \hat{q}_j be estimators of p_j and q_j . Note that we set the number of simulated samples k to be fixed across j so that $\{\hat{q}_j\}_{j=1}^m$ are identically distributed. This reflects the fact that simulator sample collections are relatively inexpensive.

The dataset is $\mathcal{D} = \{(\psi_j, \hat{p}_j, \hat{q}_j, n_j, k)\}_{j=1}^m$. For each scenario j , let $\mathcal{D}_j^{\text{gt}} = \{(z_{j,i}, y_{j,i}^{\text{gt}})\}_{i=1}^{n_j}$ denote the human-side data used to construct \hat{p}_j , and $\mathcal{D}_j^{\text{sim}} = \{(z_{j,l}^{\text{sim}}, y_{j,l}^{\text{sim}})\}_{l=1}^k$ denote the simulator outputs used to construct \hat{q}_j . The discrepancy for a specific question ψ between simulated and real output distributions is defined as

$$\Delta_\psi := L(p_\psi, \hat{q}_\psi),$$

where $L : \Theta \times \Theta \rightarrow [0, \infty)$ is a user-chosen discrepancy function. Our method is agnostic to the choice of L and allows practitioners to choose one that suits their application, such as Kullback–Leibler (KL) divergence for categorical outputs or a Wasserstein distance for nonparametric empirical distributions. Let F_Δ denote the cumulative distribution function (CDF) of Δ_ψ when $\psi \sim \Psi$, and define the associated population quantile function $V(\alpha) := \inf\{t \in \mathbb{R} : F_\Delta(t) \geq \alpha\}, \forall \alpha \in [0, 1]$.

Ideally, our target is this function $V(\alpha)$, which describes the entire distribution of sim-to-real discrepancies across scenarios. However, this function is intractable, so we state our main goal:

Construct a calibrated function $\hat{V}(\cdot, \mathcal{D}) : [0, 1] \rightarrow \mathbb{R}$ such that, for a new $\psi \sim \Psi$ and all $\alpha \in [0, 1]$,

$$\mathbb{P}_{\psi \sim \Psi}(\Delta_\psi \leq \hat{V}(\alpha, \mathcal{D}) \mid \mathcal{D}) \approx \alpha - \varepsilon_m,$$

holds with high probability over the draw of \mathcal{D} . The quantity ε_m should vanish as $m \rightarrow \infty$.

Note that \hat{V} is a calibrated function that is indexed both by the prescribed level α and by the calibration sample \mathcal{D} . There are two distinct sources of randomness we must control: (i) *scenario*

¹We keep r fixed during calibration so that variation in \hat{q}_j reflects scenario differences rather than encoding drift or model choice; choosing a different r defines a different simulator.

uncertainty, coming from drawing finitely many scenarios $\psi_j \sim \Psi$, and (ii) *finite-sample uncertainty* within each scenario, since p_j and q_j are only observed through finite samples that produce \hat{p}_j and \hat{q}_j . Our construction of $\hat{V}(\alpha, \mathcal{D})$ explicitly addresses both layers, and the main theorem will show that we obtain non-vacuous, finite-sample coverage guarantees.

3 Theoretical Result

We introduce our method for constructing \hat{V} and then provide theoretical guarantees, following up with an extension to directly compare simulators' performance.

3.1 Methodology

Recall that our data \mathcal{D} consists of m scenarios; the j -th scenario has population parameters $p_j := p(\psi_j)$ and $q_j := q(\psi_j)$ in the real and simulator; \hat{p}_j and \hat{q}_j are their point estimates.

Our procedure has two steps:

1. For each j , compute a confidence set $\mathcal{C}_j(\hat{p}_j)$ on p_j , and calculate the pseudo-discrepancy $\hat{\Delta}_j := \sup_{u \in \mathcal{C}_j(\hat{p}_j)} L(u, \hat{q}_j)$.
2. Denote $\hat{V}_m(\alpha) :=$ the α -quantile of $\{\hat{\Delta}_j\}_{j=1}^m$.

Intuitively, $\hat{\Delta}_j$ is a worst-case discrepancy for scenario j over plausible values of p_j , and $\hat{V}_m(\alpha)$ is the empirical α -quantile of these worst-case discrepancies. We will show in the next section that $\hat{V}_m(\alpha)$ is exactly the quantile function we are aiming for. Before that, we comment that the key design choice in our methodology is the construction of confidence set $\mathcal{C}_j(\cdot)$. Note that $\mathcal{C}_j(\cdot)$ is equivalent to the uncertainty set in robust optimization, which ensures that the true $p_j \in \mathcal{C}_j(\hat{p}_j)$ with a prescribed probability. The reason for introducing the uncertainty set will be discussed in the proof sketch. For now we illustrate how to construct these sets using off-the-shelf concentration bounds. In what follows we fix a coverage level $\gamma \in (0, 1)$ (in Theorem 3.7 we take $\gamma = \frac{1}{2}$) and suppress γ in the notation, writing $\mathcal{C}_j(\hat{p}_j) := \mathcal{C}_j(\hat{p}_j, \gamma)$. Consistent with our motivating example, we start with a multinomial setting.

Example 3.1 (Multinomial Confidence Set). *Suppose we have multinomial outcomes with $\Theta = \mathcal{P}^d := \{u \in [0, 1]^d : \sum_{i=1}^d u_i = 1\}$ and denote KL divergence as $\text{KL}(\cdot \| \cdot)$. Then*

$$\begin{aligned} \mathbb{P}(p_j \in \mathcal{C}_j(\hat{p}_j)) &\geq \gamma, \text{ where} \\ \mathcal{C}_j(\hat{p}_j) &:= \left\{ u \in \mathcal{P}^d : \text{KL}(\hat{p}_j \| u) \leq \frac{d-1}{n_j} \log\left(\frac{2(d-1)}{\gamma}\right) \right\}. \end{aligned} \tag{3.1}$$

Bound (3.1) is a variant of Chernoff-Hoeffding Bound, see Lemma A.5 (Mardia et al., 2019). We also present examples for bounded outcomes, Bernoulli outcomes and general non-parametric distributions.

Example 3.2 (Bounded Outcomes). *Suppose outcome Y is bounded between $[a, b]$ and $\hat{p}_j := \frac{1}{n_j} \sum_{i=1}^{n_j} Y_{j,i}$ denote the sample mean for scenario j ., then the set*

$$\begin{aligned} \mathbb{P}(p_j \in \mathcal{C}_j(\hat{p}_j)) &\geq \gamma, \text{ where} \\ \mathcal{C}_j(\hat{p}_j) &:= \left\{ u \in [a, b] : |u - \hat{p}_j| \leq (b-a) \sqrt{\frac{\log(2/\gamma)}{2n_j}} \right\}. \end{aligned}$$

Example 3.3 (Bernoulli Confidence Set). Suppose the outcome Y is distributed according to a Bernoulli Distribution $Ber(p_j)$. Then the set

$$\mathbb{P}(p_j \in \mathcal{C}_j(\hat{p}_j)) \geq \gamma, \text{ where}$$

$$\mathcal{C}_j(\hat{p}_j) := \left\{ u \in \mathbb{R} : \text{KL}(\hat{p}_j \| u) \leq \frac{1}{n_j} \log\left(\frac{2}{\gamma}\right) \right\}.$$

Example 3.4 (Nonparametric W_1 Confidence Set). For question j , let \hat{P}_j denote the empirical distribution from n_j samples and assume that the true outcome Y is σ -sub-Gaussian. Define the confidence set

$$\mathcal{C}_j^{W_1}(\hat{P}_j) := \left\{ Q : W_1(\hat{P}_j, Q) \leq r_j(n_j, \gamma, \sigma) \right\},$$

where,

$$r_j = \frac{512\sigma}{\sqrt{n_j}} + \sigma \sqrt{\frac{256e}{n_j} \log \frac{1}{1-\gamma}}.$$

Then $\mathcal{C}_j^{W_1}(\hat{P}_j)$ covers the true distribution with probability at least γ . The above form is justified via concentration inequality proved in [L.A. and Bhat 2022](#).

Next, we show that \mathcal{C}_j can be efficiently computed for a given scenario j :

- *KL-Divergence*: With the KL-based set $\mathcal{C}_j(\hat{p}_j) = \{u \in \mathcal{P}^d : \text{KL}(\hat{p}_j \| u) \leq r_j\}$ and loss $L(u, \hat{q}_j) = \text{KL}(u \| \hat{q}_j)$, the pseudo-discrepancy is

$$\hat{\Delta}_j = \sup_{u \in \mathcal{C}_j(\hat{p}_j)} \text{KL}(u \| \hat{q}_j).$$

The maximizer for Bregman divergence lies on the boundary $\text{KL}(\hat{p}_j \| u) = r_j$ and can be computed via a one-dimensional dual search for a Lagrange multiplier or other maximization techniques.

- *Wasserstein-1*: For the W_1 -based set $\mathcal{C}_j^{W_1}(\hat{P}_j) = \{u : W_1(u, \hat{P}_j) \leq r_j\}$ and loss $L(u, \hat{Q}_j) = W_1(u, \hat{Q}_j)$. Then, by the triangle inequality,

$$\hat{\Delta}_j = \sup_{u \in \mathcal{C}_j^{W_1}(\hat{P}_j)} W_1(u, \hat{Q}_j) \leq W_1(\hat{P}_j, \hat{Q}_j) + r_j.$$

In both cases, $\hat{\Delta}_j$ reduces to a tractable, real-valued number per question, and other loss functions can be computed using similar techniques.

3.2 Calibrated Quantile Curve Theory

With the above methodology, we introduce two assumptions on which our theoretical guarantee relies and state the main theorem.

Assumption 3.1 (Independent data). Scenarios ψ_j are drawn i.i.d. from Ψ . In addition, given (ψ_1, \dots, ψ_m) , the pairs $\mathcal{D} = \{(\mathcal{D}_j^{\text{gt}}, \mathcal{D}_j^{\text{sim}})\}_{j=1}^m$ are independent with $\mathcal{D}_j^{\text{gt}} \perp\!\!\!\perp \mathcal{D}_j^{\text{sim}}$ conditional on ψ_j , and the new $(\psi, \mathcal{D}^{\text{sim}})$ is independent of \mathcal{D} .

Assumption 3.2 (Regular Discrepancy). The discrepancy $L : \Theta \times \Theta \rightarrow [0, \infty)$ is jointly continuous on $\Theta \times \Theta$ and satisfies $L(u, u) = 0$ for all $u \in \Theta$.

Theorem 3.1. Suppose Assumption 3.1 and 3.2 hold. For any simulation sample size $k \in \mathbb{N}$, define the per-scenario discrepancy (unobservable) $\Delta_j^{(k)}$ and the pseudo-discrepancy (observable) $\hat{\Delta}_j^{(k)}$ by

$$\Delta_j^{(k)} := L(p_j, \hat{q}_j), \quad \hat{\Delta}_j^{(k)} := \sup_{u \in \mathcal{C}_j(\hat{p}_j)} L(u, \hat{q}_j),$$

where $\mathcal{C}_j(\hat{p}_j) \subset \Theta$ are data-driven compact confidence sets satisfying $\mathbb{P}(p_j \in \mathcal{C}_j(\hat{p}_j) \mid \psi_j, n_j) \geq \frac{1}{2}$. Then, for any $\alpha \in (0, 1)$, with probability at least $1 - \eta$ over \mathcal{D} , we have the following guarantee:

$$\mathbb{P}_{\psi \sim \Psi} \left(\Delta_\psi^{(k)} \leq \hat{V}_m \left(1 - \frac{\alpha}{2} \right) \mid \mathcal{D} \right) \geq 1 - \alpha - \frac{\varepsilon(\alpha, m, \eta)}{\sqrt{m}}, \quad (3.2)$$

where

$$\begin{aligned} \varepsilon(\alpha, m, \eta) = & \sqrt{2\alpha \log \frac{2m}{\eta} + \frac{(\log \frac{2m}{\eta})^2 + 4 \log \frac{2m}{\eta}}{m}} \\ & + \frac{\log \frac{2m}{\eta} + 2}{\sqrt{m}} + \sqrt{\frac{\log(4/\eta)}{2}}. \end{aligned}$$

In particular, the remainder is $O(\sqrt{(\log m)/m})$, which vanishes as $m \rightarrow \infty$ for any α .

We can rewrite the guarantee in Theorem 3.1 by the coverage level α rather than the tail parameter $1 - \alpha$. In particular, for each fixed interior α and large m this behaves as

$$\mathbb{P}_{\psi \sim \Psi} \left(\Delta_\psi^{(k)} \leq \hat{V}_m \left(\frac{1+\alpha}{2} \right) \mid \mathcal{D} \right) \approx \alpha - o_m(1). \quad (3.3)$$

Heuristically, we can interpret the above formulation as follows: to get α -level coverage for the true discrepancy distribution, we evaluate the empirical pseudo-quantile at the slightly more central index $(\alpha + 1)/2$ (the finite-sample uncertainty), and we pay a vanishing $o_m(1)$ loss in coverage coming from having only m scenarios (the scenario uncertainty). The complete proof is in Appendix B. We present the intuition of the proof below.

Proof Sketch:

For simplicity we ignore the dependence of Δ on k , as it will not affect the final theoretical guarantee. Suppose we have access to an oracle that returns the true p_j for each scenario j , then this guarantees we have i.i.d. $\{\Delta_j\}_{j=1}^m$. Therefore, by adopting the well-known Dvoretzky-Kiefer-Wolfowitz (DKW) inequality on the m scenarios, we can use the empirical quantile to upper-bound the true quantile function with probability $1 - \eta$. However, we face two difficulties: (i) We do not have access to this oracle, (ii) each scenario has heterogeneous n_j sample size.

To solve this problem, we wrap the uncertainty set \mathcal{C}_j around the estimated parameter \hat{p}_j , which guarantees $\mathbb{P}(p_j \in \mathcal{C}_j) \geq \frac{1}{2}$. Recall that we want to know the α quantile level of the actual quantile curve of Δ_ψ for any α , and by taking sup over \mathcal{C}_j , we ensure that the pseudo-discrepancy exceeds the true discrepancy for that scenario with at least probability $\frac{1}{2}$, independently across scenarios. This allows us to infer the magnitude of the true discrepancy quantile from the pseudo-discrepancy quantile curve, since with probability $\geq \frac{1}{2}$ it provides an upper bound. To account for the part where the empirical quantile of pseudo-discrepancies doesn't provide an upper bound to the true discrepancies, we shift to a more conservative tail index to align with the real sim-to-real tail. This accounts for one source of the correction terms, and along with the DKW correction terms as the other source, we obtain our theorem statement. \square

Remark 1. Note that we set $\gamma = \frac{1}{2}$ in Theorem 3.1, which can actually be changed to any $\gamma \in (0, 1)$. The choice of γ would directly affect the size of \mathcal{C}_j and the tightness of the calibrated quantile curve, but will not affect the theoretical guarantee provided above.

A subtle but important modeling choice is that Theorem 3.1 characterizes the distribution of $\Delta^{(k)} = L(p, \hat{q})$, where \hat{q} is estimated using k samples, instead of the distribution of $\Delta = L(p, q)$. This why our guarantee holds *uniformly* for all sample sizes $k \in \mathbb{N}$. A simple adjustment to the proof yields an approximation to the quantile function of Δ . See Corollary B.1 in Appendix B for details. Conceptually, $\Delta^{(k)}$ is the correct target when one wants to characterize performance for a fixed query budget. In contrast, Δ is more natural when simulator queries are cheap, and the goal is to quantify the *model's* inherent bias. Our framework supports both these approaches. For brevity, we will, henceforth, ignore the k superscript in following theoretical discussions.

We now turn to the benefits of obtaining a distributional level characterization of sim-to-real discrepancy. With this calibrated quantile, for an unseen scenario ψ we can derive a confidence set for p_ψ given a simulated \hat{q}_ψ . For a target confidence level $\bar{\alpha} \in (0, 1)$, define

$$S_{\bar{\alpha}} := \{u \in \Theta : L(u, \hat{q}_\psi) \leq \tau_{\bar{\alpha}}\}, \quad \tau_{\bar{\alpha}} := \hat{V}_m(1 - \frac{\bar{\alpha}}{2}).$$

By Theorem 3.1, $p_\psi \in S_{\bar{\alpha}}$ with probability at least $1 - \bar{\alpha}$ up to an $o_m(1)$ remainder.

In addition, the calibrated quantile curve can be compressed to provide summary statistics. Returning to the motivating example, a media research firm can use a calibrated AUC to assess average simulator bias, while a pollster or risk-averse product team may track a calibrated CVaR $_\alpha$ to bound rare but consequential misreads of public sentiment.

Concretely, define the index-adjusted (calibrated) curve

$$\hat{V}_m^{\text{cal}}(\tau) := \hat{V}_m\left(\frac{1+\tau}{2}\right), \quad \tau \in [0, 1],$$

which matches the coverage guarantee in Theorem 3.1. We then summarize overall average sim-to-real discrepancy via the calibrated AUC

$$\text{AUC}_{\text{cal}} := \int_0^1 \hat{V}_m^{\text{cal}}(\tau) d\tau,$$

and tail risk via a calibrated right-tail CVaR: for $\alpha \in (0, 1)$,

$$\text{CVaR}_\alpha^{\text{cal}} := \frac{1}{\alpha} \int_{1-\alpha}^1 \hat{V}_m^{\text{cal}}(u) du.$$

In words, AUC_{cal} aggregates the entire calibrated curve into a single average-bias summary, while $\text{CVaR}_\alpha^{\text{cal}}$ averages over the worst α fraction of scenarios under the same finite-sample-corrected quantile indexing.

3.3 Pairwise Simulator Comparison

In Section 3.2, we have developed a distributional characterization of a single simulator's sim-to-real discrepancy performance. A natural next step is to ask: *Can we extend this viewpoint to compare two simulators, i.e., obtain a distributional characterization for pairwise performance?*

To that end, we mirror the construction in Section 3.2. Consider two simulators, denoted S_1, S_2 and their corresponding point estimates as $\hat{q}^{(1)}, \hat{q}^{(2)}$. For a given scenario ψ , define the performance discrepancy as

$$\delta_\psi := L(p_\psi, \hat{q}_\psi^{(1)}) - L(p_\psi, \hat{q}_\psi^{(2)}).$$

We are interested in making the inferential claim: “ S_1 is at least as good as S_2 ”, where $\delta_\psi < 0$ means S_1 is closer to the ground truth than S_2 on scenario ψ .

The challenge, as in Section 3.2, is that p_ψ is unobserved and, with heterogeneous n_j , the direct plug-in discrepancies $L(\hat{p}_\psi, \hat{q}_\psi^{(1)}) - L(\hat{p}_\psi, \hat{q}_\psi^{(2)})$ are non-i.i.d., making it harder to establish uniform claims across quantile levels. We therefore follow the same idea and work with a conservative, data-driven surrogate that constructs a confidence set, enabling a valid distributional comparison between S_1 and S_2 . The methodology is as follows:

1. For each j , compute a confidence set $\mathcal{C}_j(\hat{p}_j)$ on p_j , and calculate the pseudo-performance discrepancy $\hat{\delta}_j := \sup_{u \in \mathcal{C}_j(\hat{p}_j)} \{L(u, \hat{q}_j^{(1)}) - L(u, \hat{q}_j^{(2)})\}$.
2. Denote $\hat{U}_m(\alpha) :=$ The α – quantile of $\{\hat{\delta}_j\}_{j=1}^m$.

With these, we present the following theorem:

Theorem 3.2. *Suppose Assumption 3.1 and 3.2 hold. For any simulation sample size $k \in \mathbb{N}$, denote $\hat{q}_\psi^{(1)}, \hat{q}_\psi^{(2)}$ as point estimates of the two simulators for scenario ψ . Define the per-scenario performance discrepancy δ_j and the pseudo-performance discrepancy $\hat{\delta}_j$ by*

$$\begin{aligned}\delta_j &:= L(p_j, \hat{q}_j^{(1)}) - L(p_j, \hat{q}_j^{(2)}) \\ \hat{\delta}_j &:= \sup_{u \in \mathcal{C}_j(\hat{p}_j, \gamma)} \left[L(u, \hat{q}_j^{(1)}) - L(u, \hat{q}_j^{(2)}) \right],\end{aligned}$$

where confidence set $\mathcal{C}_j(\hat{p}_j, \gamma)$ satisfies $\mathbb{P}(p_j \in \mathcal{C}_j(\hat{p}_j, \gamma) \mid \psi_j, n_j) \geq \gamma$ for a fixed $\gamma \in (0, 1)$.

Then, for any $\alpha \in (0, 1)$ and $\eta \in (0, 1)$ such that with probability at least $1 - \eta$ over the calibration data \mathcal{D} ,

$$\mathbb{P}_{\psi \sim \Psi} \left(\delta_\psi \leq \hat{U}_m(1 - \frac{\alpha}{2}) \mid \mathcal{D} \right) \geq 1 - \alpha - \frac{\varepsilon(\alpha, m, \eta)}{\sqrt{m}},$$

where

$$\begin{aligned}\varepsilon(\alpha, m, \eta) &:= \sqrt{2\alpha \log \frac{2m}{\eta} + \frac{(\log \frac{2m}{\eta})^2 + 4 \log \frac{2m}{\eta}}{m}} \\ &\quad + \frac{\log \frac{2m}{\eta} + 2}{\sqrt{m}} + \sqrt{\frac{\log(4/\eta)}{2}}.\end{aligned}$$

In particular, the remainder is $O(\sqrt{(\log m)/m})$, which vanishes as $m \rightarrow \infty$ for any α .

The proof of Theorem 3.2 exactly mirrors the argument for Theorem 3.1, where we construct analogous pseudo-discrepancies via confidence sets, and then apply the same quantile-bounding machinery. Intuitively speaking, Theorem 3.2 says that: If the data-driven performance discrepancy quantile $\hat{U}_m(1 - \frac{\alpha}{2}) \leq 0$ for some $\bar{\alpha}$, then – up to an $O(\sqrt{(\log m)/m})$ adjustment term – S_1 is at least as good as S_2 on at least a $(1 - \bar{\alpha})$ fraction of scenarios with high probability.

4 Application: LLM Fidelity Profiling

4.1 Dataset and Methodology

We evaluate our procedure on the *WorldValueBench* dataset, curated by Zhao et al. (2024) and constructed from the World Values Survey (Haerpfer et al., 2020). The data comprise survey

questions asked across 64 countries that elicit attitudes spanning 12 categories (e.g., social values, security, migration). Individual-level covariates are available for each respondent, including gender, age, migration status, education, marital status, etc.

Q199. How interested would you say you are in politics? Are you

- 1 Very interested
- 2 Somewhat interested
- 3 Not very interested
- 4 Not at all interested

Figure 2: Example of World Value Questions. Retrieved from [Haerpfer et al. \(2020\)](#).

Preprocessing. After data cleaning, we retain 235 distinct questions and responses from 96,220 individuals. Each question offers a categorical set of possible answers, as illustrated in Figure 2. To place heterogeneous answer sets on a common scale, we map each question’s categories to the interval $[-1, 1]$, producing a bounded real-valued outcome for every scenario. Individual-level covariates are then used to construct synthetic profiles and corresponding prompts for simulation. Additional details on preprocessing and dataset characteristics are provided in Appendix C.

Methodology We generate synthetic responses and estimate $\{(\hat{q}_j)\}_{j=1}^{235}$ for four LLMs: GPT-4o (gpt-4o), GPT-5 MINI (gpt-5-mini), LLAMA 3.3 70B (Llama-3.3-70B-Instruct-Turbo), and QWEN 3 235B (Qwen3 235B A22B Thinking 2507 FP8). As a benchmark, we also construct a *Uniform* baseline: for each question, this generator samples an answer uniformly at random from the available choices for each individual. We set the simulation budget to $k = 500$ for each question and then randomly select 200 of these simulated responses to estimate \hat{q}_j . The same 500 human participants used to generate the synthetic profiles are also used to calculate $\{\hat{p}_j\}_{j=1}^{235}$ and construct the corresponding confidence sets \mathcal{C}_j . Because some participants decline or are unable to answer certain questions, the effective human sample size n_j to construct \hat{p}_j varies slightly across questions, typically ranging from 450 to 500.

Finally, since outcomes are bounded within $[-1, 1]$, we construct per-question confidence sets \mathcal{C}_j following Example 3.2, and compute the corresponding pseudo-gaps under our general framework. We then apply the procedure described in Section 3 to obtain a fidelity profile for each candidate LLM. Throughout, we use the squared-error discrepancy $L(p, q) = (p - q)^2$, and set the confidence-level parameter $\delta = 0.05$. For completeness, we also apply the same workflow to a Bernoulli setting (EEDI) and a multinomial setting (OpinionQA); see Appendix E.

4.2 Fidelity Profiling

The calibrated $V(\alpha)$ are shown in Figure 3, which compares models by how tightly their synthetic outcomes track the human distribution across items. We plot $\hat{V}_\ell(\tilde{\alpha})$ against α under the adjusted index $\tilde{\alpha}$ provided in Equation 3.3. Lower-flatter curves indicate uniformly small discrepancies, while elbows reveal rare but severe misses. Our result suggests that, relative to the uniform baseline, all simulators outperform it on over 70% of questions, but fail to capture the outliers. Another observation is that GPT-4o lies lowest across all quantiles, indicating the most reliable alignment, with GPT-5-MINI close behind yet failing to capture some outlier questions. LLAMA 3.3 70B and QWEN-3-235B clearly performing worse, suggesting a dominance in performance among simulators.

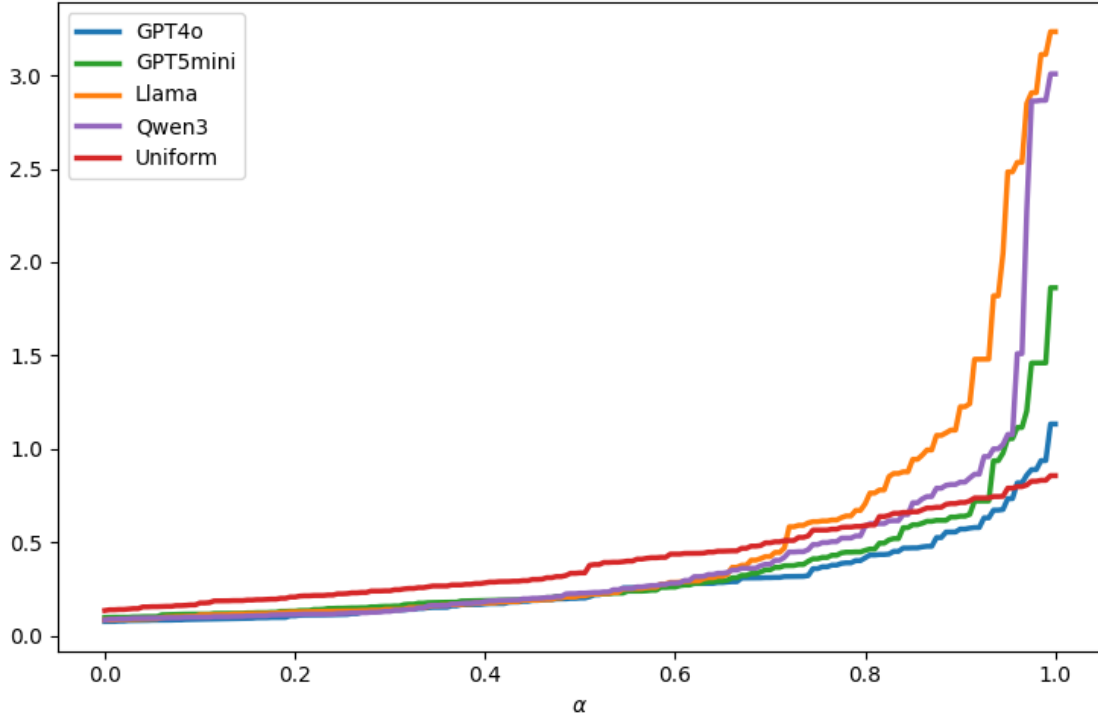


Figure 3: Calibrated $V(\alpha)$ across LLMs.

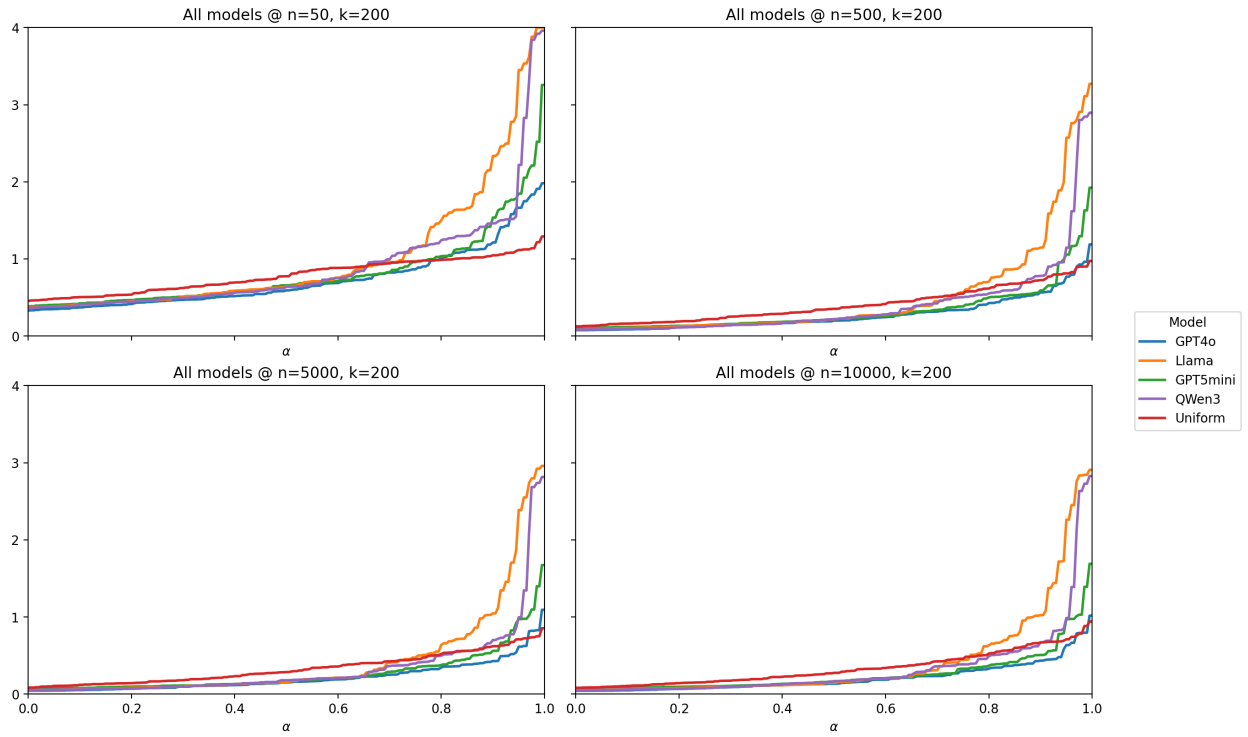


Figure 4: Robustness check of simulator performance under different n -levels.

We also provide a robustness check on whether this dominance of performance is an anomaly for this particular set of numbers of n_j . Figure 4 shows performance across LLMs given a fixed n_j , where we consider $n_j = \{50, 500, 5000, 10000\}$. The results indicate that the dominance of GPT-4O over GPT-5-MINI, LLAMA 3.3 70B, and QWEN-3-235B is consistent across n_j , hence validating our method’s result.

5 Tightness Analysis of Calibrated Quantiles

A natural concern with our approach is the tightness of the calibrated quantile curve from Theorem 3.1. Since we rely on off-the-shelf concentration inequalities to construct pseudo-gaps and empirical distribution gaps, it might be the case that our method returns a very loose quantile that is of little practical use. In order to address this concern, we conduct numerical analysis to isolate the effect of such finite-sample calibration and order-statistic adjustments from the simulation misalignment. We run the following analysis:

Since we have 96,220 respondents per question in the World Value Bench dataset, we treat the estimate \hat{p}_j based on all samples as the underlying true p_j . Therefore, we can calculate the oracle gap $\Delta_j^* := L(p_j, \hat{q}_j)$ for each question j , and hence construct the true quantile function. We set $k = 200$ and randomly sample $n_j = 100, 200, 500, 1000$ to calculate $\{\hat{\Delta}_j\}$ under different finite-sample calibration and order-statistic adjustments to evaluate the tightness of our proposed calibrated quantiles. We present our result in Figure 5.

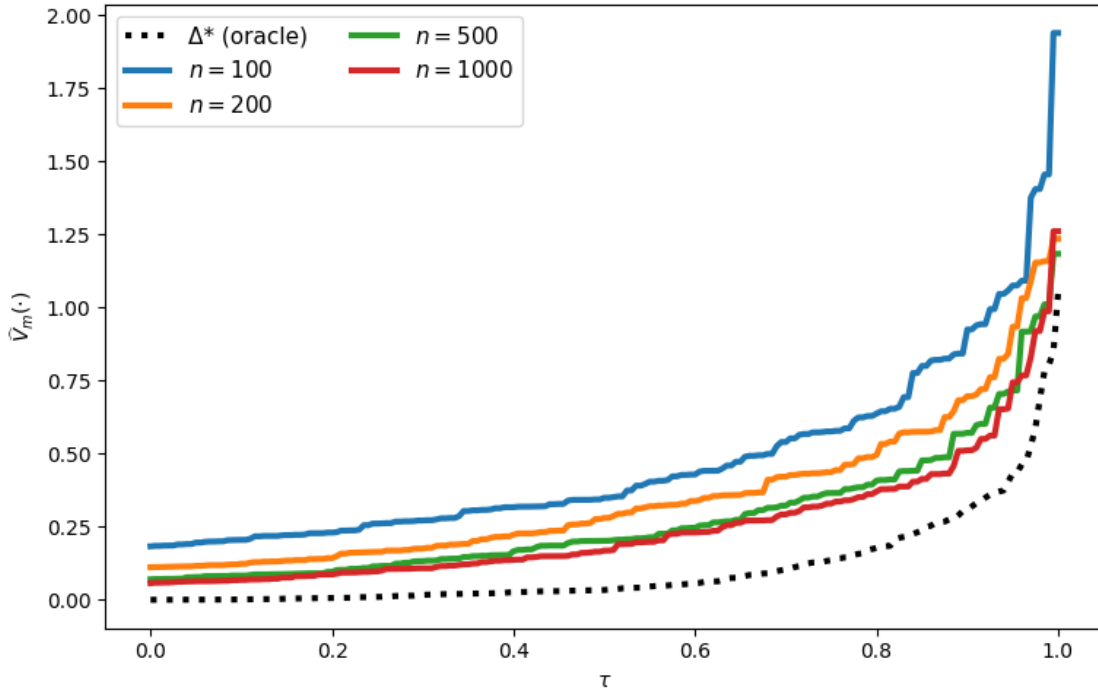


Figure 5: Tightness analysis of different n_j under GPT-4o.

As Figure 5 shows, the envelope quantile is loose when we only have $n_j = 100$. Nonetheless, it converges to the true quantile $\{\Delta^*\}$ when n grows to 200, 500, 1,000, indicating that our method is tight given sufficient real-world sample size per scenario.

However, in most cases we will not have such a large dataset as benchmark to test tightness. Subsequently, we will provide a data-driven tightness guarantee.

Theorem 5.1. *Suppose the setup of Theorem 3.1 and Assumptions 3.1–3.2 hold. Let $\gamma \in (1/2, 1]$ and define lower pseudo-discrepancies*

$$\Delta_j^- := \inf_{u \in C_j(\hat{p}_j)} L(u, \hat{q}_j),$$

where $C_j(\hat{p}_j) \subset \Theta$ are data-driven confidence sets satisfying $\mathbb{P}(p_j \in C_j(\hat{p}_j) \mid \psi_j, n_j) \geq \gamma$. Let $\hat{V}_m^-(\alpha)$ denote the empirical α -quantile of $\{\Delta_j^-\}_{j=1}^m$, and $\hat{V}_m(\alpha)$ the empirical α -quantile of the upper pseudo-discrepancies from Theorem 3.1 (under the same coverage level γ).

For any $\eta \in (0, 1)$, the following uniform lower bound holds: with probability at least $1 - \eta$ over the calibration data \mathcal{D} , for all $\alpha \in (0, 1)$,

$$\mathbb{P}_{\psi \sim \Psi} \left(\Delta(\psi) \geq \hat{V}_m^-(\gamma\alpha) \mid \mathcal{D} \right) \geq 1 - \alpha - \frac{\varepsilon_-^{(\gamma)}(\alpha, m, \eta)}{\sqrt{m}}, \quad (5.1)$$

where $\varepsilon_-^{(\gamma)}(\alpha, m, \eta)$ is of order $O(\sqrt{\frac{\log m}{m}})$ uniformly in α .

Let $V : [0, 1] \rightarrow \mathbb{R}$ denote the quantile function of $\Delta(\psi)$. Then, for any fixed $\tau \in (0, 1)$, as $m \rightarrow \infty$ we obtain the asymptotic band

$$\hat{V}_m^-(\gamma\tau) \leq V(\tau) \leq \hat{V}_m(\gamma + (1 - \gamma)\tau) + o_m(1), \quad (5.2)$$

where $o_m(1) \rightarrow 0$ as $m \rightarrow \infty$.

Proof is deferred to Appendix B. The main takeaway is that by adopting Theorem 5.1 and 3.1, we can calibrate a confidence band for the actual quantile curve.

Remark 2. As discussed in Section 3.1, the factor γ in (5.1) and (5.2) directly reflects the per-scenario coverage level of the confidence sets $C_j(\hat{p}_j)$. For the tail behavior (e.g., small τ), the upper side of the band is intrinsically conservative: we only guarantee

$$V(\tau) \leq \hat{V}_m(\gamma + (1 - \gamma)\tau) + o_m(1),$$

so when τ is close to 0 the index $\gamma + (1 - \gamma)\tau$ is bounded away from 1. In other words, the right-hand side controls $V(\tau)$ using a substantially more *central* empirical quantile, and the resulting upper band can be quite loose in the extreme left tail. The same phenomenon appears symmetrically in the right tail when $\tau \approx 1$. For any fixed $\gamma < 1$, this conservativeness near very small or very large τ is unavoidable: the band must widen in the extremes to pay for the finite per-scenario coverage level.

6 Discussion

We present a model-free estimator for the quantile function that makes no parametric assumptions on either the simulator or the ground truth, delivers finite-sample guarantees for any desired α -quantile, and supports flexible summaries from means to tail-risk measures. We further develop a pairwise comparison framework to directly compare simulators' ability to replicate real-world outcomes. Finally, we apply our methods to a real-world dataset, assessing multiple LLMs and demonstrating their relative effectiveness at simulating human responses.

Despite its broad applicability, the method we introduce leaves several natural extensions. First, our proof relies on DKW-type concentration – often conservative for small m – and a grid-uniform step that further loosens constants; tightening these bounds is an immediate target. Second, many applications involve temporally dependent, dynamic simulation processes; extending our static framework to dynamic settings would broaden applicability. Third, our analysis assumes i.i.d. scenarios, whereas covariate shift or endogenous sampling may invalidate marginal guarantees; addressing such distribution shifts is an important avenue for future work.

Acknowledgement

Kaizheng Wang’s research is supported by NSF grants DMS-2210907 and DMS-2515679 and a Data Science Institute seed grant SF-181 at Columbia University.

References

- AHER, G., ARRIAGA, R. I. and KALAI, A. T. (2023). Using large language models to simulate multiple humans and replicate human subject studies. In *Proceedings of the 40th International Conference on Machine Learning*. ICML’23, JMLR.org.
- ARGYLE, L. P., BUSBY, E. C., FULDA, N., GUBLER, J. R., RYTTHING, C. and WINGATE, D. (2023). Out of one, many: Using language models to simulate human samples. *Political Analysis* **31** 337–351.
- BARTON, R. R., NELSON, B. L. and XIE, W. (2014). Quantifying input uncertainty via simulation confidence intervals. *INFORMS Journal on Computing* **26** 74–87.
- BATES, S., ANGELOPOULOS, A., LEI, L., MALIK, J. and JORDAN, M. (2021). Distribution-free, risk-controlling prediction sets. *Journal of the ACM (JACM)* **68** 1–34.
- BUDDE, C. E., HARTMANN, A., MEGGENDORFER, T., WEININGER, M. and WIENHÖFT, P. (2025). Statistical model checking beyond means: Quantiles, cvar, and the dkw inequality (extended version). *arXiv preprint arXiv:2509.11859* .
- CHEN, M., LAM, H. and LIU, Z. (2024). Quantifying distributional input uncertainty via inflated kolmogorov-smirnov confidence band. *arXiv preprint arXiv:2403.09877* .
- DURMUS, E., NGUYEN, K., LIAO, T. I., SCHIEFER, N., ASKELL, A., BAKHTIN, A., CHEN, C., HATFIELD-DODDS, Z., HERNANDEZ, D., JOSEPH, N. ET AL. (2023). Towards measuring the representation of subjective global opinions in language models. *arXiv preprint arXiv:2306.16388* .
- GAO, Y., LEE, D., BURTCH, G. and FAZELPOUR, S. (2025). Take caution in using llms as human surrogates. *Proceedings of the National Academy of Sciences* **122** e2501660122.
URL <https://www.pnas.org/doi/abs/10.1073/pnas.2501660122>
- HAERPFER, C., INGLEHART, R., MORENO, A., WELZEL, C., KIZILOVA, K., DIEZ-MEDRANO, J., LAGOS, M., NORRIS, P., PONARIN, E., PURANEN, B. ET AL. (2020). World values survey: Round seven – country-pooled datafile (2017–2020).
URL <https://doi.org/10.14281/18241.1>
- HE-YUEYA, J., MA, W. A., GANDHI, K., DOMINGUE, B. W., BRUNSKILL, E. and GOODMAN, N. D. (2024). Psychometric alignment: Capturing human knowledge distributions via language models. *arXiv preprint arXiv:2407.15645* .
- HOEFFDING, W. (1963). Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association* **58** 13–30.
URL <https://www.tandfonline.com/doi/abs/10.1080/01621459.1963.10500830>
- HUANG, C., WU, Y. and WANG, K. (2025). Uncertainty quantification for LLM-based survey simulations. In *Forty-second International Conference on Machine Learning*.
URL <https://openreview.net/forum?id=nY1Ge2wxtP>
- JEON, Y., CHU, Y., PASUPATHY, R. and SHASHAANI, S. (2024). Uncertainty quantification using simulation output: Batching as an inferential device.
URL <https://arxiv.org/abs/2311.04159>

- L.A., P. and BHAT, S. P. (2022). A wasserstein distance approach for concentration of empirical risk estimates. *Journal of Machine Learning Research* **23** 1–61.
URL <http://jmlr.org/papers/v23/20-965.html>
- LAM, H. (2022). Cheap bootstrap for input uncertainty quantification. In *Proceedings of the 2022 Winter Simulation Conference*. IEEE.
- LU, Y., HUANG, J., HAN, Y., YAO, B., BEI, S., GESI, J., XIE, Y., HE, Q., WANG, D. ET AL. (2025). Prompting is not all you need! evaluating llm agent simulation methodologies with real-world online customer behavior data. *arXiv preprint arXiv:2503.20749* .
- MACAL, C. (2016). Everything you need to know about agent-based modelling and simulation. *Journal of Simulation* **10** 144–156.
- MARDIA, J., JIAO, J., TÁNCZOS, E., NOWAK, R. D. and WEISSMAN, T. (2019). Concentration inequalities for the empirical distribution of discrete distributions: beyond the method of types. *Information and Inference: A Journal of the IMA* **9** 813–850.
URL <https://doi.org/10.1093/imaiai/iaz025>
- MASSART, P. (1990). The tight constant in the dvoretzky-kiefer-wolfowitz inequality. *The annals of Probability* 1269–1283.
- PARK, J. S., O’BRIEN, J., CAI, C. J., MORRIS, M. R., LIANG, P. and BERNSTEIN, M. S. (2023). Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th annual acm symposium on user interface software and technology*.
- PENG, X. B., ANDRYCHOWICZ, M., ZAREMBA, W. and ABBEEL, P. (2018). Sim-to-real transfer of robotic control with dynamics randomization. In *2018 IEEE international conference on robotics and automation (ICRA)*. IEEE.
- ROY, C. J. and OBERKAMPF, W. L. (2011). A comprehensive framework for verification, validation, and uncertainty quantification in scientific computing. *Computer Methods in Applied Mechanics and Engineering* **200** 2131–2144.
URL <https://www.sciencedirect.com/science/article/pii/S0045782511001290>
- SANTURKAR, S., DURMUS, E., LADHAK, F., LEE, C., LIANG, P. and HASHIMOTO, T. (2023a). Whose opinions do language models reflect? In *Proceedings of the 40th International Conference on Machine Learning*. ICML’23, JMLR.org.
- SANTURKAR, S., DURMUS, E., LADHAK, F., LEE, C., LIANG, P. and HASHIMOTO, T. (2023b). Whose opinions do language models reflect? In *Proceedings of the 40th International Conference on Machine Learning*. ICML’23, JMLR.org.
- SNELL, J. C., ZOLLO, T. P., DENG, Z., PITASSI, T. and ZEMEL, R. (2022). Quantile risk control: A flexible framework for bounding the probability of high-loss predictions. *arXiv preprint arXiv:2212.13629* .
- TOBIN, J., FONG, R., RAY, A., SCHNEIDER, J., ZAREMBA, W. and ABBEEL, P. (2017). Domain randomization for transferring deep neural networks from simulation to the real world. In *2017 IEEE/RSJ international conference on intelligent robots and systems (IROS)*. IEEE.
- VOVK, V., GAMMERMAN, A. and SHAFER, G. (2005). *Algorithmic Learning in a Random World*. Springer-Verlag, Berlin, Heidelberg.

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. (2021). Results and insights from diagnostic questions: The neurips 2020 education challenge. In *Proceedings of the NeurIPS 2020 Competition and Demonstration Track* (H. J. Escalante and K. Hofmann, eds.), vol. 133 of *Proceedings of Machine Learning Research*. PMLR.

URL <https://proceedings.mlr.press/v133/wang21a.html>

ZHANG, L., ZHOU, L., REN, L. and LAILI, Y. (2019). Modeling and simulation in intelligent manufacturing. *Computers in Industry* **112** 103123.

URL <https://www.sciencedirect.com/science/article/pii/S0166361519303239>

ZHAO, W., MONDAL, D., TANDON, N., DILLION, D., GRAY, K. and GU, Y. (2024). Worldvalues-bench: A large-scale benchmark dataset for multi-cultural value awareness of language models.

URL <https://arxiv.org/abs/2404.16308>

A Technical Lemmas

Lemma A.1. (*Dvoretzky-Kiefer-Wolfowitz (DKW) Inequality via [Massart \(1990\)](#)*)

Let X_1, X_2, \dots, X_n be i.i.d. real-valued random variables with cumulative distribution function (CDF) F^* , and let \hat{F}_n be the empirical distribution function defined by

$$\hat{F}_n(x) := \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{X_i \leq x\}.$$

Then, for any $\varepsilon > 0$,

$$\mathbb{P} \left(\sup_{x \in \mathbb{R}} |\hat{F}_n(x) - F^*(x)| > \varepsilon \right) \leq 2e^{-2n\varepsilon^2}.$$

Equivalently, for any confidence level $\delta \in (0, 1)$, with probability at least $1 - \delta$,

$$\sup_{x \in \mathbb{R}} |\hat{F}_n(x) - F^*(x)| \leq \sqrt{\frac{1}{2n} \log \left(\frac{2}{\delta} \right)}.$$

Lemma A.2. [*Hoeffding (additive) via [Hoeffding \(1963\)](#)*]

Let $Z_1, \dots, Z_n \in [0, 1]$ be independent, $T = \sum_{i=1}^n Z_i$, and $\mu = \mathbb{E}[T]$. For any $t \in [0, \mu]$,

$$\mathbb{P}(T \leq \mu - t) \leq \exp \left(-\frac{2t^2}{\sum_{i=1}^n (1 - 0)^2} \right) = \exp \left(-\frac{2t^2}{n} \right).$$

Lemma A.3. [*Chernoff-Hoeffding for one-parameter exponential family*]

Let X_1, \dots, X_n be i.i.d. with density (or mass) in the one-parameter canonical exponential family

$$p_\theta(x) = \exp\{\theta T(x) - A(\theta)\} h(x), \quad \theta \in \Theta,$$

where $T(x)$ is the sufficient statistic, $A(\theta)$ is the log-partition function (convex, differentiable on Θ), and the mean map is $\mu(\theta) := \mathbb{E}_\theta[T(X)] = A'(\theta)$. Assume Θ is an open interval and all quantities below are finite.

Define the empirical mean of the sufficient statistic

$$\bar{T}_n = \frac{1}{n} \sum_{i=1}^n T(X_i).$$

For each t in the range of \bar{T}_n let θ_t be the (unique) canonical parameter satisfying $\mu(\theta_t) = t$. Then for any $\varepsilon > 0$,

$$\mathbb{P}(D(p_{\theta_{\bar{T}_n}} \| p_\theta) > \varepsilon) \leq 2e^{-n\varepsilon}$$

Proof:

Define the shifted log-MGF under p_θ ,

$$\psi_\theta(\lambda) := \log \mathbb{E}_\theta[e^{\lambda T(X)}] = A(\theta + \lambda) - A(\theta),$$

where the displayed equality follows from the exponential-family form (for λ in the domain where the expectation is finite). For an attainable mean m denote θ_m as the unique solution of $\mu(\theta_m) = m$.

By adopting the one-sided Chernoff bound, for any real λ such that expectations exist and any $m \in \mathbb{R}$,

$$\begin{aligned} \mathbb{P}_\theta(\bar{T}_n \geq m) &= \mathbb{P}(e^{\lambda n \bar{T}_n} \geq e^{\lambda n m}) \\ &\leq e^{-\lambda n m} \mathbb{E}_\theta[e^{\lambda n \bar{T}_n}] \\ &= \exp(-n(\lambda m - \psi_\theta(\lambda))). \end{aligned}$$

Optimizing over λ gives the Chernoff bound

$$\begin{aligned} \mathbb{P}_\theta(\bar{T}_n \geq m) &\leq \exp(-n\psi_\theta^*(m)), \\ \implies \psi_\theta^*(m) &:= \sup_{\lambda \in \mathbb{R}} \{\lambda m - \psi_\theta(\lambda)\}, \end{aligned}$$

where $\psi_\theta^*(m)$ is the Fenchel-Legendre transform of log-MGF. A symmetric argument with $\lambda < 0$ yields the lower-tail bound

$$\mathbb{P}_\theta(\bar{T}_n \leq m) \leq \exp(-n\psi_\theta^*(m)).$$

We next link the Fenchel-Legendre transform to KL-Divergence using exponential tilting. First, adopting change of variables $\eta = \theta + \lambda$. Then

$$\begin{aligned} \psi_\theta^*(m) &= \sup_{\eta} \{\langle \eta - \theta, m \rangle - (A(\eta) - A(\theta))\} \\ &= A(\theta) + A^*(m) - \langle \theta, m \rangle, \end{aligned}$$

where $A^*(m) = \sup_{\eta} \{\langle \eta, m \rangle - A(\eta)\}$ is the convex conjugate of A . When m is attainable, the supremum is achieved at $\eta = \theta_m$, and therefore

$$\psi_\theta^*(m) = \langle \theta_m - \theta, m \rangle - (A(\theta_m) - A(\theta)).$$

But for exponential-family densities one has the following direct algebraic identity for the KL:

$$\begin{aligned} D(p_{\theta_1} \| p_{\theta_2}) &= \mathbb{E}_{\theta_1} \left[\log \frac{p_{\theta_1}(X)}{p_{\theta_2}(X)} \right] \\ &= \mathbb{E}_{\theta_1} [(\theta_1 - \theta_2) T(X) - (A(\theta_1) - A(\theta_2))] \\ &= (\theta_1 - \theta_2) \mathbb{E}_{\theta_1} [T(X)] - (A(\theta_1) - A(\theta_2)). \end{aligned}$$

Taking $\theta_1 = \theta_m$ and $\theta_2 = \theta$ (so $\mathbb{E}_{\theta_1}[T] = m$) yields

$$\psi_{\theta}^*(m) = D(p_{\theta_m} \| p_{\theta}).$$

Combining this with the one-sided Chernoff-bound above yields the one-sided KL-form Chernoff bounds

$$\begin{aligned}\mathbb{P}_{\theta}(\bar{T}_n \geq m) &\leq e^{-nD(p_{\theta_m} \| p_{\theta})}, \\ \mathbb{P}_{\theta}(\bar{T}_n \leq m) &\leq e^{-nD(p_{\theta_m} \| p_{\theta})}.\end{aligned}$$

Fix $\varepsilon > 0$. Because $A'' > 0$ the function $m \mapsto D(p_{\theta_m} \| p_{\theta})$ is continuous, strictly convex and has a unique minimum 0 at $m = A'(\theta)$. Thus the sublevel set $\{m : D(p_{\theta_m} \| p_{\theta}) < \varepsilon\}$ is an open interval (m_-, m_+) ; equivalently

$$\{m : D(p_{\theta_m} \| p_{\theta}) \geq \varepsilon\} = (-\infty, m_-] \cup [m_+, \infty).$$

Hence

$$\{D(p_{\theta_{\bar{T}_n}} \| p_{\theta}) \geq \varepsilon\} \subseteq \{\bar{T}_n \leq m_-\} \cup \{\bar{T}_n \geq m_+\}.$$

Applying the one-sided KL bounds at m_{\pm} (each equals ε) and using the union bound gives

$$\mathbb{P}(D(p_{\theta_{\bar{T}_n}} \| p_{\theta}) \geq \varepsilon) \leq e^{-n\varepsilon} + e^{-n\varepsilon} = 2e^{-n\varepsilon},$$

which proves the lemma.

Lemma A.4. (*Chernoff-Hoeffding Inequality*)

Let $X_1, \dots, X_n \sim \text{Ber}(\tilde{p})$ be i.i.d. Bernoulli random variables with unknown mean \tilde{p} , and define the empirical mean as

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i.$$

Then for any $\varepsilon > 0$,

$$\mathbb{P}(D(\bar{X}_n \| \tilde{p}) > \varepsilon) \leq 2e^{-n\varepsilon},$$

where $D(p \| q)$ is the Kullback-Leibler divergence between Bernoulli distributions with parameters p and q .

Proof: Via Lemma A.3.

Lemma A.5 (Multinomial Chernoff-Hoeffding Bound via [Mardia et al. \(2019\)](#)). For all $d \leq (\frac{nC_0}{4})^{\frac{1}{3}}$ and $P \in \mathcal{M}_k$, the following holds with the universal constants $C_0 = \frac{e^3}{2\pi} \approx 3.1967$, for any $\epsilon > 0$,

$$\Pr(D(\hat{P}_{n,d} \| P) \geq \epsilon) \leq 2(d-1)e^{-\frac{n\epsilon}{d-1}}.$$

Lemma A.6 (Wasserstein distance bound via [L.A. and Bhat \(2022\)](#)). Let X be a sub-Gaussian r.v. with parameter σ . Let F denote the CDF of X . Then, for every $n \geq 1$ and ε such that $\frac{512\sigma}{\sqrt{n}} < \varepsilon < \frac{512\sigma}{\sqrt{n}} + 16\sigma\sqrt{e}$, we have

$$\mathbb{P}(W_1(F_n, F) > \varepsilon) \leq \exp\left(-\frac{n}{256\sigma^2 e} \left(\varepsilon - \frac{512\sigma}{\sqrt{n}}\right)^2\right),$$

where e is Euler's number.

Lemma A.7 (Order-statistic thresholding). *Let $x_1, \dots, x_m \in \mathbb{R}$ and let $x_{(1)} \leq \dots \leq x_{(m)}$ be their order statistics. Fix a threshold $\tau \in \mathbb{R}$ and an integer $N \in \{1, \dots, m\}$. If at least N of the sample values are at least τ , i.e. $|\{j : x_j \geq \tau\}| \geq N$, then*

$$x_{(m-N+1)} \geq \tau.$$

(Equivalently, if at least N of the x_j are strictly larger than τ , the same conclusion holds.)

Proof:

Suppose, for contradiction, that $x_{(m-N+1)} < \tau$. Then all of the first $m - N + 1$ order statistics are strictly less than τ , so there are at most $m - (m - N + 1) = N - 1$ indices j with $x_j \geq \tau$, contradicting $|\{j : x_j \geq \tau\}| \geq N$. Hence $x_{(m-N+1)} \geq \tau$.

B Proof to Theorem 3.1 and 5.1

B.1 Proof to Theorem 3.1:

As a setup, we work conditionally on $\mathcal{G}_j := \sigma(\psi_j, \hat{q}_j, n_j)$, $\mathcal{G} := \sigma(\{\mathcal{G}_j\}_{j=1}^m)$. Here \mathcal{G}_j is the σ -field collecting all “human-side” information for scenario j . Thus $\mathcal{G} := \sigma(\{\mathcal{G}_j\}_{j=1}^m)$ represents the joint information from all scenarios that we treat as fixed when analyzing the remaining randomness coming from the simulator estimates \hat{p}_j . For brevity, we abuse notation and write $p(\psi_j), q(\psi_j)$ as p_j, q_j , and similarly \hat{p}, \hat{q} . We also ignore the superscript k to simplify the proof, and discussions on the dependence on k can be found in Section 3. In addition, for any quantities $\{\Delta_j\}_{j=1}^m$, we denote the sorted version as $\{\Delta_{(i)}\}_{i=1}^m$, ie. $\Delta_{(1)} \leq \dots \leq \Delta_{(m)}$. For any sequence $\{\Delta_j\}_{j=1}^m$, let $\Delta_{(1)} \leq \dots \leq \Delta_{(m)}$ denote its order statistics.

We first consider a fixed level of confidence $\bar{\alpha}$, then extend to $\bar{\alpha}$ holding uniformly across $(0, 1)$. Throughout the proof, $\alpha \in (0, 1)$ denotes a generic quantile index (a function argument) and is distinct from the target coverage level $\bar{\alpha}$.

We seek the quantile function of $\Delta_j = L(p_j, \hat{q}_j)$, but only observe the estimators (\hat{p}_j, \hat{q}_j) . Note that the sequence $\{\Delta_j\}_{j=1}^m$ is i.i.d., since the simulator budget k is fixed and each scenarios $\{\psi_j\}$ are i.i.d. as assumed in Section 2.

Let $\bar{V}_m(\alpha) = \inf\{t : \bar{F}_m(t) \geq \alpha\} = \Delta_{(\lceil m\alpha \rceil)}$ be the empirical α quantile of $\{\Delta_j\}_{j=1}^m$, where $\bar{F}_m(x) = \frac{1}{m} \sum_{j=1}^m \mathbb{I}(\Delta_j \leq x)$. By Lemma A.1, with probability $1 - \delta$:

$$\begin{aligned} \mathbb{P}_{\psi \sim \Psi}(L(p(\psi), \hat{q}(\psi)) \leq \bar{V}_m(1 - \frac{\bar{\alpha}}{2}) | \mathcal{D}) \\ \geq 1 - \frac{\bar{\alpha}}{2} - \sqrt{\frac{\log(2/\delta)}{2m}} \end{aligned} \tag{B.1}$$

and we would have our desired bound.

However, for the above argument to hold, we would need to know p_j , instead we only have the estimated \hat{p}_j . Therefore, we instead construct a randomized pseudo-gaps. Specifically, we set a level of coverage $\gamma = \frac{1}{2}$, and define the two gap terms as

$$\begin{cases} \Delta_j = L(p_j, \hat{q}_j) & , \text{ i.i.d., yet unobservable} \\ \hat{\Delta}_j := \sup_{u \in \mathcal{C}_j(\hat{p}_j, \gamma)} L(u, \hat{q}_j) & , \text{ Not i.i.d., yet observable,} \end{cases}$$

where $\mathcal{C}_j(\hat{p}_j, \gamma) \subset \Theta$ are data-driven confidence sets satisfying $\mathbb{P}(p_j \in \mathcal{C}_j(\hat{p}_j, \gamma) | \psi_j, n_j) \geq \frac{1}{2}$.

By Assumption 3.2 and the compactness of the confidence set $C_j(\hat{p}_j, \gamma) \subset \Theta$, Berge's maximum theorem guarantees that the supremum in the definition of $\hat{\Delta}_j$ is attained, hence $\hat{\Delta}_j$ is well defined. Therefore, by the coverage property of C_j , we have

$$\mathbb{P}(\hat{\Delta}_j \geq \Delta_j | \mathcal{G}_j \cup \sigma(\gamma)) \geq \gamma,$$

and by $\gamma \perp\!\!\!\perp \mathcal{G}_j$ we get $\forall j$

$$\begin{aligned} \mathbb{P}(\hat{\Delta}_j \geq \Delta_j | \mathcal{G}_j) &= \mathbb{E}[\mathbb{P}(\hat{\Delta}_j \geq \Delta_j | \mathcal{G}_j \cup \sigma(\gamma) | \mathcal{G}_j)] \\ &\geq \frac{1}{2}. \end{aligned} \tag{B.2}$$

Furthermore, by the tower property and $\gamma \perp\!\!\!\perp \mathcal{G}_j$,

$$\mathbb{P}(\hat{\Delta}_j \geq \Delta_j) = \mathbb{E}[\mathbb{P}(\hat{\Delta}_j \geq \Delta_j | \mathcal{G}_j)] \geq \frac{1}{2}.$$

Moreover, by the above discussion, conditional on \mathcal{G} , the indicators $Y_j = \mathbf{1}\{\hat{\Delta}_j \geq \Delta_j\}$ are independent across j .

We will use $\{\hat{\Delta}_j\}_{j=1}^m$ to create an upper bound of $\bar{V}_m(1 - \alpha)$, which along side (B.1) will give us our desired envelope. Intuitively, we want to find a larger quantile of $\hat{\Delta}_j$ and with (B.2), we can claim an upper bound of $\bar{V}_m(1 - \alpha)$ with high probability.

First, define $S_\alpha := \{j \in [m] : \Delta_j \geq \bar{V}_m(1 - \alpha)\}$ and $s := |S_\alpha|$. By definition of the empirical $(1 - \alpha)$ -quantile $\bar{V}_m(1 - \alpha)$, at most a fraction $(1 - \alpha)$ of the sample can lie strictly below it. More precisely,

$$|\{j \in [m] : \Delta_j < \bar{V}_m(1 - \alpha)\}| \leq \lfloor m(1 - \alpha) \rfloor.$$

Therefore the number of indices with $\Delta_j \geq \bar{V}_m(1 - \alpha)$ is at least

$$s \geq m - \lfloor m(1 - \alpha) \rfloor = \lceil m\alpha \rceil \geq \lfloor m\alpha \rfloor.$$

We next define $Y_j = \mathbb{I}(\hat{\Delta}_j \geq \Delta_j)$. By (B.2): $\mathbb{P}(Y_j = 1 | \mathcal{G}) \geq \frac{1}{2}$. Next, define $T_\alpha = \{j \in S_\alpha : \hat{\Delta}_j \geq \Delta_j\} = \{j \in S_\alpha : Y_j = 1\} = \{j : \hat{\Delta}_j \geq \Delta_j \geq \bar{V}_m(1 - \alpha)\}$. First, we calculate $\mathbb{E}[T_\alpha]$:

$$\mathbb{E}[T_\alpha | \mathcal{G}] = \mathbb{E}\left[\sum_{j \in S_\alpha} Y_j | \mathcal{G}\right] = \sum_{j \in S_\alpha} \mathbb{P}(Y_j | \mathcal{G}) \geq \frac{1}{2}s.$$

Lemma A.2 implies that for any $\delta \in [0, \frac{1}{2}s]$, we have:

$$\begin{aligned} \mathbb{P}\left(|T_\alpha| \leq \frac{1}{2}s - t | \mathcal{G}\right) &\leq \mathbb{P}\left(|T_\alpha| \leq \mathbb{E}[T_\alpha] - t | \mathcal{G}\right) \\ &\leq \exp\left(-\frac{2t^2}{s}\right), \end{aligned}$$

where we applied $\mathbb{E}[T_\alpha | \mathcal{G}] \geq \frac{1}{2}s$ in the first inequality. By setting $t = c\sqrt{s}$:

$$\mathbb{P}\left(|T_\alpha| \leq \frac{1}{2}s - t | \mathcal{G}\right) \leq \exp(-2c^2) \tag{B.3}$$

With this bound, we can link the actual set of indices we have interest, ie. $U_\alpha = \{j : \hat{\Delta}_j \geq \bar{V}_m(1 - \alpha)\}$ to T_α . By construction, $T_\alpha \subseteq U_\alpha$, hence by (B.3), with probability greater than $1 - \exp(-2c^2)$:

$$|U_\alpha| \geq |T_\alpha| \geq \frac{1}{2}s - c\sqrt{s},$$

which implies at least $\frac{1}{2}s - c\sqrt{s}$ of the $\hat{\Delta}_j$'s are larger than $\bar{V}_m(1 - \alpha)$ with high probability.

We now analyze what coverage guarantee can we get for the inner probability via order statistics for any α . Set $N := \lfloor \frac{1}{2}s - c\sqrt{s} \rfloor$ and define $\hat{V}_m(\alpha) := \inf\{t : \hat{F}_m(t) \geq \alpha\} = \Delta_{(\lceil m\alpha \rceil)}$, where $\hat{F}(x) = \frac{1}{m} \sum_{j=1}^m \mathbb{I}(\hat{\Delta}_j \leq x)$. If at least N sample values exceed $\bar{V}_m(1 - \alpha)$, then by order-statistics calculus (Lemma A.7)

$$\hat{V}_m(1 - \alpha) = \hat{\Delta}_{(m - \lfloor m\alpha \rfloor)} \geq \bar{V}_m(1 - \alpha_{\text{eff}}),$$

whenever α_{eff} is chosen so that $N \geq \lfloor m\alpha \rfloor + 1$ holds when $s \geq \lfloor m\alpha_{\text{eff}} \rfloor + 1$.

A sufficient condition for the above to satisfy is

$$\frac{1}{2}m\alpha_{\text{eff}} - c\sqrt{m\alpha_{\text{eff}}} - 1 \geq m\alpha.$$

Define $\alpha_{\text{eff}}(\alpha, c, m) := \inf\{x \in (0, 1) : \frac{1}{2}x - c\sqrt{\frac{x}{m}} - \frac{1}{m} \geq \alpha\}$. Writing $y = \sqrt{x}$, this is equivalent to $y^2 - \frac{2c}{\sqrt{m}}y - (\frac{2}{m} + 2\alpha) \geq 0$, so the minimal admissible y is

$$y^* = \frac{2c/\sqrt{m} + \sqrt{4c^2/m + 8\alpha + 8/m}}{2},$$

$$\alpha_{\text{eff}}(\alpha, c, m) = (y^*)^2.$$

Thus, we define

$$\mathcal{E}_\alpha := \left\{ \hat{V}_m(1 - \alpha) \geq \bar{V}_m(1 - \alpha_{\text{eff}}(\alpha, c, m)) \right\}, \quad (\text{B.4})$$

$$\text{then } \mathbb{P}(\mathcal{E}_\alpha) \geq 1 - e^{-2c^2}.$$

Apply (B.1) at level $1 - \alpha_{\text{eff}}(\bar{\alpha}, c, m)$:

$$\mathbb{P}_{\psi \sim \Psi} \left(L(p(\psi), \hat{q}(\psi)) \leq \bar{V}_m(1 - \alpha_{\text{eff}}) \middle| \mathcal{D} \right) \geq 1 - \alpha_{\text{eff}}(\bar{\alpha}, c, m) - \varepsilon_m(\delta).$$

On $\mathcal{E}_{\bar{\alpha}}$ in (B.4), $\bar{V}_m(1 - \alpha_{\text{eff}}) \leq \hat{V}_m(1 - \bar{\alpha})$. Hence,

$$\mathbb{P}_{\psi \sim \Psi} \left(L(p(\psi), \hat{q}(\psi)) \leq \hat{V}_m(1 - \bar{\alpha}) \middle| \mathcal{D} \right) \geq 1 - \alpha_{\text{eff}}(\bar{\alpha}, c, m) - \varepsilon_m(\delta),$$

with outer probability at least $1 - \delta - e^{-2c^2}$.

The exact algebraic form is

$$\alpha_{\text{eff}}(\bar{\alpha}, c, m) = \frac{\left(\frac{2c}{\sqrt{m}} + \sqrt{\frac{4c^2}{m} + 8\bar{\alpha} + \frac{8}{m}} \right)^2}{4} =$$

$$2\bar{\alpha} + \frac{c}{\sqrt{m}} \sqrt{8\bar{\alpha} + \frac{4c^2+8}{m}} + \frac{2c^2+2}{m}.$$

Therefore, as $m \rightarrow \infty$, $\alpha_{\text{eff}}(\alpha, c, m) = 2\bar{\alpha} + c\sqrt{8\bar{\alpha}/m} + O(m^{-1})$.

We have shown that for any target level α and choice of $c > 0$, the preceding argument yields a high-probability concentration bound based on $\{\hat{\Delta}_j\}_{j=1}^m$. We now extend the guarantee to hold *uniformly* over all α . Fix $c > 0$ and $\delta \in (0, 1)$. For the grid $\alpha_r := r/m$ ($r = 1, \dots, m$), let

$$\mathcal{E}_{\text{DKW}} := \left\{ \sup_x |\hat{F}_m(x) - F^*(x)| \leq \varepsilon_m(\delta) \right\}$$

$$\mathcal{E}_r := \left\{ (3.2) \text{ holds with } \alpha = \alpha_r \right\}.$$

By DKW, $\Pr(\mathcal{E}_{\text{DKW}}) \geq 1 - \delta$, and by the fixed-level argument, $\Pr(\mathcal{E}_r) \geq 1 - e^{-2c^2}$ for each r . Hence, by a union bound,

$$\Pr\left(\mathcal{E}_{\text{DKW}} \cap \bigcap_{r=1}^m \mathcal{E}_r\right) \geq 1 - \delta - me^{-2c^2}.$$

For any $\alpha \in (0, 1)$ let $r = \lceil m\alpha \rceil$ and denote $\alpha_+ := \alpha_r = r/m \in [\alpha, \alpha + 1/m]$. Since the empirical quantile is piecewise constant on the m -grid,

$$\hat{V}_m(1 - \alpha) = \hat{V}_m(1 - \alpha_+).$$

Applying (3.2) at level α_+ yields

$$\begin{aligned} & \mathbb{P}_{\psi \sim \Psi}\left(\Delta(\psi) \leq \hat{V}_m(1 - \alpha) \mid \mathcal{D}\right) \\ & \geq 1 - 2\alpha_+ - \frac{c}{\sqrt{m}} \sqrt{8\alpha_+ + \frac{4c^2 + 8}{m}} - \frac{2c^2 + 2}{m} - \varepsilon_m(\delta). \end{aligned}$$

Since $\alpha_+ \in [\alpha, \alpha + 1/m]$ and the right-hand side is nonincreasing in α , the same bound holds with α replaced by α_+ , and (optionally) one may absorb the rounding slack $\alpha_+ - \alpha \leq 1/m$ into the $O(m^{-1})$ term by a crude inequality

$$2\alpha_+ \leq 2\alpha + \frac{2}{m}, \quad \sqrt{8\alpha_+ + \frac{4c^2 + 8}{m}} \leq \sqrt{8\alpha + \frac{4c^2 + 16}{m}}.$$

Therefore, with probability at least $1 - \delta - me^{-2c^2}$ over \mathcal{D} , the guarantee (3.2) (at α replaced by α_+) holds *uniformly* for all $\alpha \in (0, 1)$. The form in the main theorem is a simplification with respect to c and δ .

B.2 Corollary to Theorem 3.1:

Corollary B.1. *Suppose Assumptions 3.1 and 3.2 hold. For any simulation sample size $k_j \in \mathbb{N}$, define the per-scenario true sim-to-real discrepancy and its pseudo-discrepancy by*

$$\Delta_j^{*,(k)} := L(p_j, q_j), \quad \hat{\Delta}_j^{*,(k)} := \sup_{\substack{u \in C_j^p(\hat{p}_j) \\ v \in C_j^q(\hat{q}_j)}} L(u, v),$$

where $C_j^p(\hat{p}_j) \subset \Theta_p$, $C_j^q(\hat{q}_j) \subset \Theta_q$ are data-driven compact confidence sets such that

$$\mathbb{P}(p_j \in C_j^p(\hat{p}_j), q_j \in C_j^q(\hat{q}_j) \mid \psi_j, n_j, k_j) \geq \frac{1}{2}.$$

Let $\hat{V}_m^*(\alpha)$ denote the empirical α -quantile of $\{\hat{\Delta}_j^{*,(k)}\}_{j=1}^m$.

Then, for any $\alpha \in (0, 1)$ and any $\eta \in (0, 1)$, with probability at least $1 - \eta$ over \mathcal{D} , we have

$$\mathbb{P}_{\psi \sim \Psi}\left(\Delta_\psi^{*,(k)} \leq \hat{V}_m^*\left(1 - \frac{\alpha}{2}\right) \mid \mathcal{D}\right) \geq 1 - \alpha - \frac{\varepsilon(\alpha, m, \eta)}{\sqrt{m}}, \quad (\text{B.5})$$

where the remainder $\varepsilon(\alpha, m, \eta)$ is the same as in Theorem 3.1, in particular $\varepsilon(\alpha, m, \eta) = O(\sqrt{\log m})$ and $\varepsilon(\alpha, m, \eta)/\sqrt{m} = O(\sqrt{(\log m)/m})$.

Proof. The argument is identical to the proof of Theorem 3.1, except that we now construct separate confidence sets for p_j and q_j with marginal coverages γ_p and γ_q chosen so that $\gamma_p \gamma_q = \frac{1}{2}$. For example, we can set $\gamma_p = \gamma_q = \frac{1}{\sqrt{2}}$. This ensures that the joint event $\{p_j \in C_j^p(\hat{p}_j), q_j \in C_j^q(\hat{q}_j)\}$ plays exactly the same role as the univariate coverage event in Theorem 3.1, i.e., $\mathbb{P}(\hat{\Delta}_j^{*,(k)} \geq \Delta_j^{*,(k)} \mid \mathcal{G}_j) \geq \frac{1}{2}$. All subsequent steps then carry over verbatim, yielding the stated bound. \square

B.3 Proof to Theorem 5.1:

We follow the setup of Theorem 3.1. Recall that we work conditionally on $\mathcal{G}_j := \sigma(\psi_j, \hat{q}_j, n_j)$, $\mathcal{G} := \sigma(\{\mathcal{G}_j\}_{j=1}^m)$. For brevity, write $p_j := p(\psi_j)$, $q_j := q(\psi_j)$, and similarly \hat{p}_j, \hat{q}_j . The true per-scenario discrepancy is $\Delta_j := L(p_j, \hat{q}_j)$, $j = 1, \dots, m$, and by the i.i.d. scenario assumption, $\{\Delta_j\}_{j=1}^m$ are i.i.d. draws from the distribution of $\Delta(\psi)$ under $\psi \sim \Psi$.

By Assumption 3.2 and compactness of $C_j(\hat{p}_j)$, Berge's maximum theorem implies that the infimum is attained on the confidence set, hence Δ_j^- is well defined.

Also by definition, on the event $\{p_j \in C_j(\hat{p}_j)\}$ we clearly have $\Delta_j^- \leq \Delta_j$. Thus, for each j ,

$$\mathbb{P}(\Delta_j^- \leq \Delta_j \mid \mathcal{G}_j) = \mathbb{P}(p_j \in C_j(\hat{p}_j) \mid \mathcal{G}_j) \geq \frac{1}{2}.$$

Define the indicators $Y_j^- := \mathbf{1}\{\Delta_j^- \leq \Delta_j\}$, $j = 1, \dots, m$. Conditionally on \mathcal{G} , the variables $\{Y_j^-\}_{j=1}^m$ are independent and satisfy

$$\mathbb{P}(Y_j^- = 1 \mid \mathcal{G}) \geq \frac{1}{2}, \quad \forall j. \quad (\text{B.6})$$

Fix $\alpha \in (0, 1)$ for the moment. We again denote $\bar{V}_m(\alpha)$ as the empirical α -quantile of the (unobservable) $\{\Delta_j\}_{j=1}^m$, i.e., $\bar{V}_m(\alpha) := \Delta_{(\lceil m\alpha \rceil)}$, where $\Delta_{(1)} \leq \dots \leq \Delta_{(m)}$ are the order statistics, and let $\bar{F}_m(t) := \frac{1}{m} \sum_{j=1}^m \mathbf{1}\{\Delta_j \leq t\}$ be their empirical CDF.

Define the index set of "lower-tail" scenarios

$$S_\alpha^- := \{j \in [m] : \Delta_j \leq \bar{V}_m(\alpha)\}, \quad s^- := |S_\alpha^-|.$$

By definition of $\bar{V}_m(\alpha)$ we have $s^- \geq m\alpha$, and these are the true Δ quantile that we wish to capture. Among these, consider those for which the lower pseudo-gap is valid:

$$\begin{aligned} T_\alpha^- &:= \{j \in S_\alpha^- : Y_j^- = 1\} \\ &= \{j \in S_\alpha^- : \Delta_j^- \leq \Delta_j \leq \bar{V}_m(\alpha)\}. \end{aligned}$$

By (B.6), $\mathbb{E}[|T_\alpha^-| \mid \mathcal{G}] = \sum_{j \in S_\alpha^-} \mathbb{P}(Y_j^- = 1 \mid \mathcal{G}) \geq \frac{1}{2} s^-$.

Applying Lemma A.2, for any $t \in [0, s^-/2]$,

$$\mathbb{P}\left(|T_\alpha^-| \leq \frac{1}{2} s^- - t \mid \mathcal{G}\right) \leq \exp\left(-\frac{2t^2}{s^-}\right).$$

Setting $t = c\sqrt{s^-}$ for some $c > 0$ yields

$$\mathbb{P}\left(|T_\alpha^-| \leq \frac{1}{2} s^- - c\sqrt{s^-} \mid \mathcal{G}\right) \leq e^{-2c^2}. \quad (\text{B.7})$$

Define also $U_\alpha^- := \{j \in [m] : \Delta_j^- \leq \bar{V}_m(\alpha)\}$. Then $T_\alpha^- \subseteq U_\alpha^-$, so on the complement of the event in (B.7),

$$|U_\alpha^-| \geq |T_\alpha^-| \geq \frac{1}{2} s^- - c\sqrt{s^-} \geq \frac{1}{2} m\alpha - c\sqrt{m\alpha}.$$

Let $\hat{F}_m^-(t)$ denote the empirical CDF of $\{\Delta_j^-\}_{j=1}^m$, ie. $\hat{F}_m^-(t) := \frac{1}{m} \sum_{j=1}^m \mathbf{1}\{\Delta_j^- \leq t\}$, and let $\hat{V}_m^-(\beta)$ be its empirical β -quantile.

From the bound on $|U_\alpha^-|$ we obtain

$$\hat{F}_m^-(\bar{V}_m(\alpha)) = \frac{|U_\alpha^-|}{m} \geq \frac{1}{2}\alpha - c\sqrt{\frac{\alpha}{m}} - \frac{1}{m}. \quad (\text{B.8})$$

Thus, for any

$$\beta \leq \beta_{\text{eff}}^-(\alpha, c, m) := \frac{1}{2}\alpha - c\sqrt{\frac{\alpha}{m}} - \frac{1}{m},$$

we have $\hat{F}_m^-(\bar{V}_m(\alpha)) \geq \beta$, and hence, by the definition of the empirical quantile,

$$\hat{V}_m^-(\beta) \leq \bar{V}_m(\alpha), \quad \forall \beta \leq \beta_{\text{eff}}^-(\alpha, c, m). \quad (\text{B.9})$$

Let \mathcal{E}_α^- denote the event that both (B.7) (complement) and (B.9) hold; then $\mathbb{P}(\mathcal{E}_\alpha^- \mid \mathcal{G}) \geq 1 - e^{-2c^2}$.

Finally, we try to bridge $\bar{F}_m(t)$ to the true CDF of $\Delta(\psi)$ under $\psi \sim \Psi$, which we denote as F_Δ . By DKW inequality (Lemma A.1), for any $\delta \in (0, 1)$, with probability at least $1 - \delta$, we have the following event:

$$\begin{aligned} \mathcal{E}_{\text{DKW}} &:= \left\{ \sup_{t \in \mathbb{R}} |\bar{F}_m(t) - F_\Delta(t)| \leq \varepsilon_m(\delta) \right\}, \\ \varepsilon_m(\delta) &:= \sqrt{\frac{\log(2/\delta)}{2m}}. \end{aligned}$$

Work on the intersection $\mathcal{E}_{\text{DKW}} \cap \mathcal{E}_\alpha^-$. Fix $\beta \leq \beta_{\text{eff}}^-(\alpha, c, m)$ and set $t^- := \hat{V}_m^-(\beta)$. By (B.9), we have $t^- \leq \bar{V}_m(\alpha)$. Using this and the fact that F_Δ is nondecreasing, we have

$$F_\Delta(t^-) \leq F_\Delta(\bar{V}_m(\alpha)).$$

On \mathcal{E}_{DKW} , the DKW bound yields

$$F_\Delta(\bar{V}_m(\alpha)) \leq \bar{F}_m(\bar{V}_m(\alpha)) + \varepsilon_m(\delta) \leq \alpha + \varepsilon_m(\delta).$$

Combining the last two equations,

$$F_\Delta(t^-) \leq \alpha + \varepsilon_m(\delta).$$

Equivalently,

$$F_\Delta(\hat{V}_m^-(\beta)) = \mathbb{P}_{\psi \sim \Psi}(\Delta(\psi) \leq \hat{V}_m^-(\beta) \mid \mathcal{D}) \leq \alpha + \varepsilon_m(\delta),$$

hence on the event $\mathcal{E}_{\text{DKW}} \cap \mathcal{E}_\alpha^-$ we have,

$$\begin{aligned} \mathbb{P}_{\psi \sim \Psi}(\Delta(\psi) \geq \hat{V}_m^-(\beta) \mid \mathcal{D}) &\geq 1 - \alpha - \varepsilon_m(\delta), \\ \forall \beta &\leq \beta_{\text{eff}}^-(\alpha, c, m). \end{aligned} \quad (\text{B.10})$$

Finally, we specify the actual adjustment term β for a fixed α . Fix $\alpha \in (0, 1)$. We wish to apply (B.10) with $\beta = \alpha/2$ to keep the argument of \hat{V}_m^- simple. Recall that (B.10) holds whenever $\beta \leq \beta_{\text{eff}}^-(\alpha', c, m)$ for some α' . Thus we need to choose α' such that

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

Writing $y := \sqrt{\alpha'}$, we obtain the quadratic inequality $y^2 - \frac{2c}{\sqrt{m}}y - \left(\alpha + \frac{2}{m}\right) \geq 0$. Solving the minimal admissible solution we get:

$$y^*(\alpha, c, m) := \frac{\frac{2c}{\sqrt{m}} + \sqrt{\frac{4c^2}{m} + 4\alpha + \frac{8}{m}}}{2},$$

$$\alpha_{\text{eff}}^-(\alpha, c, m) := (y^*(\alpha, c, m))^2.$$

For this choice of α' we have $\alpha/2 \leq \beta_{\text{eff}}^-(\alpha', c, m)$, so (B.10) with $\beta = \alpha/2$ gives

$$\mathbb{P}_{\psi \sim \Psi}(\Delta(\psi) \geq \hat{V}_m^-(\frac{\alpha}{2}) | \mathcal{D}) \geq 1 - \alpha_{\text{eff}}^-(\alpha, c, m) - \varepsilon_m(\delta).$$

Explicitly solving quadratic solution, one can derive

$$\alpha_{\text{eff}}^-(\alpha, c, m) - \alpha \leq \frac{c}{\sqrt{m}} \sqrt{4\alpha + \frac{4c^2+8}{m}} + \frac{2c^2+2}{m},$$

hence implies

$$\mathbb{P}_{\psi \sim \Psi}(\Delta(\psi) \geq \hat{V}_m^-(\frac{\alpha}{2}) | \mathcal{D}) \geq$$

$$1 - \alpha - \frac{c}{\sqrt{m}} \sqrt{4\alpha + \frac{4c^2+8}{m}} - \frac{2c^2+2}{m} - \varepsilon_m(\delta). \quad (\text{B.11})$$

Again, we now make the bound uniform over $\alpha \in (0, 1)$ and tie (c, δ) to a target failure probability $\eta \in (0, 1)$.

Consider the grid $\alpha_r := r/m$ for $r = 1, \dots, m$. For each r , the fixed-level argument above yields an event

$$\mathcal{E}_r^- := \left\{ (\text{B.11}) \text{ holds with } \alpha = \alpha_r \right\}.$$

From the Chernoff step we have, for each r , $\mathbb{P}(\mathcal{E}_r^-) \geq 1 - e^{-2c^2}$. Hence, by a union bound,

$$\mathbb{P}\left(\mathcal{E}_{\text{DKW}} \cap \bigcap_{r=1}^m \mathcal{E}_r^-\right) \geq 1 - \delta - me^{-2c^2}.$$

We now choose

$$\delta := \frac{\eta}{2}, \quad c^2 := \frac{1}{2} \log \frac{2m}{\eta},$$

so that

$$\mathbb{P}(\mathcal{E}_{\text{DKW}}) \geq 1 - \frac{\eta}{2}, \quad e^{-2c^2} = \frac{\eta}{2m} \implies me^{-2c^2} = \frac{\eta}{2},$$

and therefore

$$\mathbb{P}\left(\mathcal{E}_{\text{DKW}} \cap \bigcap_{r=1}^m \mathcal{E}_r^-\right) \geq 1 - \eta.$$

On this high-probability event, plugging $\alpha = \alpha_r$ and the choices of (c, δ) into (B.11) gives, for each r ,

$$\mathbb{P}_{\psi \sim \Psi}(\Delta(\psi) \geq \hat{V}_m^-(\frac{\alpha_r}{2}) | \mathcal{D}) \geq$$

$$1 - \alpha_r - \frac{c}{\sqrt{m}} \sqrt{4\alpha_r + \frac{4c^2+8}{m}} - \frac{2c^2+2}{m} - \varepsilon_m(\delta).$$

A direct algebraic simplification, using $c^2 = \frac{1}{2} \log(2m/\eta)$ and $\delta = \eta/2$, shows that the right-hand side is:

$$1 - \alpha_r - \frac{1}{\sqrt{m}} \varepsilon_-(\alpha, m, \eta),$$

where

$$\begin{aligned} \varepsilon_-(\alpha, m, \eta) := & \sqrt{2\alpha \log \frac{2m}{\eta} + \frac{(\log \frac{2m}{\eta})^2 + 4 \log \frac{2m}{\eta}}{m}} \\ & + \frac{\log \frac{2m}{\eta} + 2}{\sqrt{m}} + \sqrt{\frac{\log(4/\eta)}{2}}, \end{aligned}$$

so that for each grid point α_r we can write

$$\mathbb{P}_{\psi \sim \Psi}(\Delta(\psi) \geq \hat{V}_m^-(\frac{\alpha_r}{2}) | \mathcal{D}) \geq 1 - \alpha_r - \frac{\varepsilon_-(\alpha_r, m, \eta)}{\sqrt{m}}.$$

Finally, for an arbitrary $\alpha \in (0, 1)$, let $r := \lceil m\alpha \rceil$ and $\alpha_+ := \alpha_r = r/m \in [\alpha, \alpha + 1/m]$. Since $\hat{V}_m^-(\cdot)$ is nondecreasing,

$$\hat{V}_m^-(\frac{\alpha}{2}) \leq \hat{V}_m^-(\frac{\alpha_+}{2}),$$

and since $\alpha_+ \leq \alpha + 1/m$ and $\varepsilon_-(\alpha_+, m, \eta)$ differs from $\varepsilon_-(\alpha, m, \eta)$ by at most an $O(1)$ amount, we can enlarge ε_- slightly (absorbing the $1/m$ rounding slack) so that

$$\mathbb{P}_{\psi \sim \Psi}(\Delta(\psi) \geq \hat{V}_m^-(\frac{\alpha}{2}) | \mathcal{D}) \geq 1 - \alpha - \frac{\varepsilon_-(\alpha, m, \eta)}{\sqrt{m}},$$

for all $\alpha \in (0, 1)$, on an event of probability at least $1 - \eta$. This is exactly (5.1).

Finally, we derive the confidence band guarantee. Let V denote the (left-continuous) quantile function of $\Delta(\psi)$,

$$V(u) := \inf\{t \in \mathbb{R} : F_\Delta(t) \geq u\},$$

where F_Δ is the CDF of $\Delta(\psi)$ under $\psi \sim \Psi$.

From Theorem 3.1, we already know that on the same high-probability event,

$$\mathbb{P}_{\psi \sim \Psi}(\Delta(\psi) \leq \hat{V}_m(1 - \frac{\alpha}{2}) | \mathcal{D}) \geq 1 - \alpha - \frac{\varepsilon(\alpha, m, \eta)}{\sqrt{m}}, \quad (\text{B.12})$$

for some remainder $\varepsilon(\alpha, m, \eta)$ of the similar algebraic form as $\varepsilon_-(\alpha, m, \eta)$. Rewriting (5.1) and (B.12) in terms of F_Δ yields

$$\begin{aligned} F_\Delta(\hat{V}_m^-(\frac{\alpha}{2})) &\leq \alpha + \frac{\varepsilon_-(\alpha, m, \eta)}{\sqrt{m}}, \\ F_\Delta(\hat{V}_m(1 - \frac{\alpha}{2})) &\geq 1 - \alpha - \frac{\varepsilon(\alpha, m, \eta)}{\sqrt{m}} \end{aligned}$$

By the definition of V as the left-continuous inverse of F_Δ , these inequalities are equivalent to

$$\begin{aligned} \hat{V}_m^-(\frac{\alpha}{2}) &\leq V\left(\alpha + \frac{\varepsilon_-(\alpha, m, \eta)}{\sqrt{m}}\right) \\ V\left(1 - \alpha - \frac{\varepsilon(\alpha, m, \eta)}{\sqrt{m}}\right) &\leq \hat{V}_m\left(1 - \frac{\alpha}{2}\right). \end{aligned}$$

Combining these displays yields the desired quantile sandwich (5.2) for all $\alpha \in (0, 1)$, on an event of probability at least $1 - \eta$. This completes the proof.

C World Value Bench Methodology

C.1 Selection of Survey Questions

The World Values Survey contains 259 questions, grouped into categories such as social values, well-being, economic values, and security. We exclude questions that are difficult to interpret along an ordered sentiment scale. For example, in Figure 6, Question 223 is highly country- and time-specific and does not admit a natural ordering of sentiment, making it hard to align with other questions.

Q223. If there were a national election tomorrow, for which party on this list would you vote? Just call out the number on this card. If DON'T KNOW: Which party appeals to you most?

1. Party 1
2. Party 2
3. Party 3
4. etc.

[INSERT COUNTRY-SPECIFIC LIST OF PARTIES]

Figure 6: Text of Question 223.

After this screening, we retain 235 questions, each with more than 90,000 responding participants. Ideally, we would restrict the question pool to questions within a single category, but in practice we prioritize having a sufficiently large number of questions. Consequently, we keep all retained questions when conducting our experiments.

C.2 Example Question and Preprocess

Below we list three example questions from the dataset.

In your view, how often do the following things occur in this country’s elections?

| | Very often | Fairly often | Not often | Not at all often |
|--|------------|--------------|-----------|------------------|
| Q224. Votes are counted fairly | 1 | 2 | 3 | 4 |
| Q225. Opposition candidates are prevented from running | 1 | 2 | 3 | 4 |
| Q226. TV news favors the governing party | 1 | 2 | 3 | 4 |
| Q227. Voters are bribed | 1 | 2 | 3 | 4 |
| Q228. Journalists provide fair coverage of elections | 1 | 2 | 3 | 4 |
| Q229. Election officials are fair | 1 | 2 | 3 | 4 |
| Q230. Rich people buy elections | 1 | 2 | 3 | 4 |
| Q231. Voters are threatened with violence at the polls | 1 | 2 | 3 | 4 |
| Q232. Voters are offered a genuine choice in the elections | 1 | 2 | 3 | 4 |
| Q233. Women have equal opportunities to run the office | 1 | 2 | 3 | 4 |

Figure 7: Example questions from the Political Interest category.

Q163. All things considered, would you say that the world is better off, or worse off, because of science and technology? Please tell me which comes closest to your view on this scale: 1 means that “the world is a lot worse off,” and 10 means that “the world is a lot better off.” (Code one number):

| | | | | | | | | | | |
|-----------------|---|---|---|---|---|---|---|---|----|------------------|
| A lot worse off | | | | | | | | | | A lot better off |
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | |

Figure 8: Example question from the Science and Technology category.

As discussed in the main text, to obtain a unified scale across heterogeneous categorical response formats, we map all answers to numerical values in $[-1, 1]$. Specifically, we query GPT-5 to determine

Q121. Now we would like to know your opinion about the people from other countries who come to live in [your country] - the immigrants. How would you evaluate the impact of these people on the development of [your country]?

| | | | | |
|-----------|------------|--------------------------|-----------|----------|
| Very good | Quite good | Neither good, nor bad | Quite bad | Very bad |
| 5 | 4 | 3 | 2 | 1 |

Figure 9: Example question from the Migration category.

the direction of this mapping according to its assessment of the “idealness” of each response in terms of public opinion. For example, in Figure 8, GPT-5 maps option 1 (“A lot worse off”) to -1 and option 10 (“A lot better off”) to 1, with intermediate choices linearly spaced between these endpoints.²

C.3 Synthetic Profile Generation

In the original dataset, some respondent IDs appear irregularly or are duplicated. We exclude these records and are left with 96,220 unique participants. The dataset contains demographic information for each participant, summarized in Table 1.

Table 1: Master list of demographic variables and categories.

| Demographic | Categories |
|------------------------------|---|
| Sex | Male; Female |
| Country | 64 countries/societies |
| Age | Numeric (years) |
| Language | Text |
| Migration status | Native-born; First-generation immigrant; Second-generation; Other/missing |
| Marital status | Married; Living together as married; Divorced; Separated; Widowed; Never married/Single |
| Children status | Number of children (numeric) |
| Education level | ISCED levels 0–8 ³ |
| Employment status | Employed full-time; Employed part-time; Self-employed; Unemployed; Student; Homemaker; Retired; Other |
| Household income decile | 1 (lowest) ; 2; 3; 4; 5; 6; 7; 8; 9; 10 (highest) |
| Religious affiliation | Christian; Muslim; Hindu; Buddhist; Jewish; No religion; Other religion |
| Self-identified social class | Lower class; Working class; Lower middle class; Upper middle class; Upper class |
| Settlement Type | Urban; Rural (with Population Size) |

We select a subset of these demographic attributes to define our synthetic profiles, including Country, Sex, Age, and Settlement Type, among others. The LLMs are asked to produce a categorical response encoded as [[1]], [[2]], etc. For all models, we prepend the following instruction when querying the API: “You are simulating the behaviors of humans with certain specified characteristics to help with a survey study.” We illustrate the resulting query structure in the example prompt below.

²This conversion does not represent the authors’ values; it reflects GPT-5’s inferred public opinion, which may be influenced by RLHF or other factors unknown to the authors.

Example prompt:

Pretend that you reside in Mexico. You live in a capital city urban with a town size of 100,000-500,000 area. You are female, your age is between 35-44. and you are Married. You normally speak Spanish; Castilian at home. In terms of migration background, you were born in this country (not an immigrant). In terms of education, you attained Upper secondary education (ISCED 3). Your current employment status is: Housewife not otherwise employed. You work in Not applicable; Never had a job. You belong to Roman Catholic; Latin Church;.

On a scale of 1 to 10, 1 meaning 'A lot worse off' and 10 meaning 'A lot better off', do you think science and technology have improved or worsened the world? Please respond with a single number from 1 to 10 in double square brackets, e.g., [[1]].

D Simulation System Examples

Manufacturing: Factory Production (discrete-event simulation; cycle time).

- Outcome space: $\mathcal{X} = \mathbb{R}_+$ (cycle time or throughput).
- Scenarios: ψ = product mix + scheduling policy.
- Profiles: \mathcal{Z} = machine/operator states, shift team, lot sizes; \mathcal{P} = plant variability.
- Laws: $Q^{\text{gt}}(\cdot \mid z, \psi)$ = empirical cycle-time distribution on the floor; $Q^{\text{sim}}(\cdot \mid a(z), \psi, r)$ = DES output under mirrored inputs.
- Parameters: $\Theta = \mathbb{R}$ (mean cycle time) or \mathbb{R}^2 (mean, variance).
- Discrepancy: L = difference of means, Gaussian KL.
- Sampling: n_j production runs logged; k simulated replications per scenario.

Environment: Urban decarbonization (technology choice; multinomial).

- Outcome space: $\mathcal{X} = \{1, \dots, K\}$ with mean $p(\psi) \in \Delta^{K-1}$ (for example, gas furnace, heat pump, variable refrigerant flow, other).
- Scenarios: ψ consists of city, season, rebate level, carbon price path, and policy bundle.
- Profiles: \mathcal{Z} contains household and building attributes such as income, occupants, roof area, and baseline electricity use, or exogenous drivers including weather and demand shocks.
- Laws: $Q^{\text{gt}}(\cdot \mid z, \psi)$ denotes the empirical technology-choice distribution, and $Q^{\text{sim}}(\cdot \mid a(z), \psi, r)$ denotes the simulator output under mirrored inputs.
- Parameters: $\Theta = \Delta^{K-1}$ for category probabilities, or a low-dimensional reparameterization such as multinomial logistic parameters.
- Discrepancy: L on the simplex, for example the total-variation distance $\frac{1}{2}\|p - q\|_1$, the multiclass Kullback–Leibler divergence $\sum_{c=1}^K p_c \log(p_c/q_c)$.
- Sampling: n_j human records per scenario and k synthetic replications per scenario.

E Applications

E.1 EEDI Dataset

Our dataset is EEDI [He-Yueya et al. \(2024\)](#), built on the NeurIPS 2020 Education Challenge [Wang et al. \(2021\)](#), which consists of student responses to mathematics multiple-choice questions collected on the Eedi online education platform. The full corpus includes 573 distinct questions and 443,433 responses from 2,287 students, and each question has four options A-D that we binarize as “correct/incorrect” based on the student’s or simulator’s choice, consistent with Lemma [A.4](#). We adopt the preprocessed version curated by [Huang et al. \(2025\)](#), which retains questions with at least 100 student responses and excludes items with graphs or diagrams, yielding 412 questions. EEDI also provides individual-level covariates such as gender, age, and socioeconomic status, which the authors of [Huang et al. \(2025\)](#) use to construct synthetic profiles. Under the same problem formulation, they compute $\{\hat{p}_j, \hat{q}_j\}_{j=1}^{412}$ for seven LLMs: GPT-3.5-TURBO (`gpt-3.5-turbo`), GPT-4O (`gpt-4o`), and GPT-4O-MINI (`gpt-4o-mini`); CLAUDE 3.5 HAIKU (`claude-3-5-haiku-20241022`); LLAMA 3.3 70B (`Llama-3.3-70B-Instruct-Turbo`); MISTRAL 7B (`Mistral-7B-Instruct-v0.3`); DEEPSEEK-V3 (`DeepSeek-V3`), and constructed a benchmark random simulator that selects uniformly among the available answer choices. A more detailed exploration into the EEDI dataset and the simulation procedure can be found in [Huang et al. \(2025\)](#).

We apply our methodology to produce a fidelity profile for each candidate LLM ℓ . We use absolute error as the loss, $L(p, q) = |p - q|$. We set $\gamma = 0.5$ uniformly and the DKW failure probability to $\delta = 0.1$, which determines the curve’s effective width at $\alpha \rightarrow 1$ in Figure [10](#). In addition, we set the simulation budget $k = 50$.

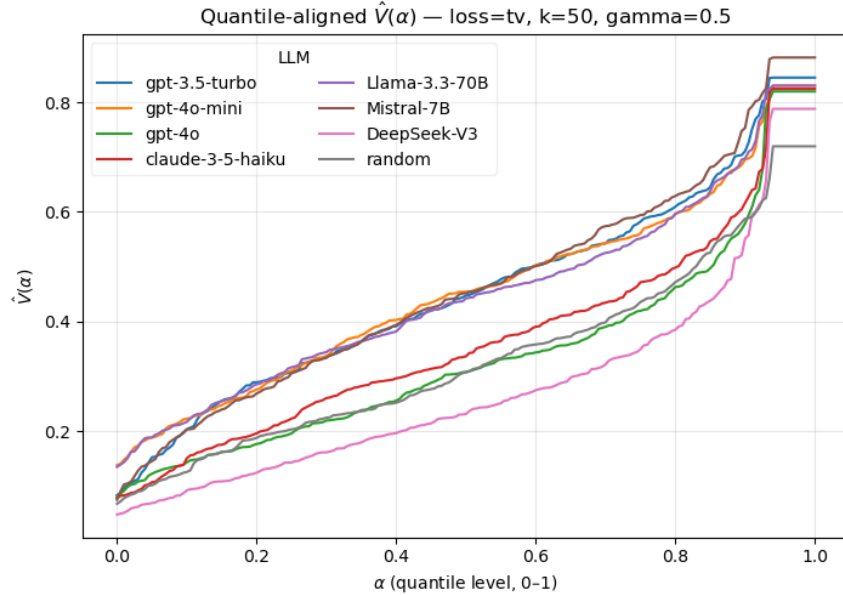


Figure 10: Quantile fidelity profiles $\hat{V}(\alpha)$ across LLMs (Discrepancy: Absolute loss, $k = 50$, $\beta = 0.5$, $\delta = 0.1$).

Figure [10](#) compares models by how tightly their synthetic outcomes track the human distribution across items. We plot $\hat{V}_\ell(\alpha)$ against α , where lower-flatter curves indicate uniformly small discrepancies, while elbows reveal rare but severe misses. DEEPSEEK-V3 lies lowest across most quantiles, indicating the most reliable alignment, with the random benchmark and GPT-4O close behind.

Notably, several models do not outperform the random baseline, suggesting they may be ill-suited for agent-based simulation under this discrepancy function.

E.2 OpinionQA

Our dataset is OpinionQA Santurkar et al. (2023b), built from the Pew Research’s American Trends Panel, which consists of the US population’s responses to survey questions spanning topics such as racial equity, security, and technology. We adopt the preprocessed version curated by Huang et al. (2025), which includes 385 distinct questions and 1,476,868 responses from at least 32,864 people. Each question has 5 choices, corresponding to the order sentiments, which is a multinomial setting. We can construct confidence sets \mathcal{C}_j for multinomial vectors by adopting Example 3.1 with $d = 5$. OpinionQA also provides individual-level covariates such as gender, age, socioeconomic status, religious affiliation, and marital status, and more, which are used to construct synthetic profiles. Under the same problem formulation, the authors of Huang et al. (2025) compute $\{\hat{p}_j, \hat{q}_j\}_{j=1}^{385}$ for seven LLMs: GPT-3.5-TURBO (gpt-3.5-turbo), GPT-4o (gpt-4o), and GPT-4o-MINI (gpt-4o-mini); CLAUDE 3.5 HAIKU (claude-3-5-haiku-20241022); LLAMA 3.3 70B (Llama-3.3-70B-Instruct-Turbo); MISTRAL 7B (Mistral-7B-Instruct-v0.3); DEEPSEEK-V3 (DeepSeek-V3), and constructed a baseline random simulator that selects uniformly among the available answer choices. A more detailed exploration into the OpinionQA dataset and the simulation procedure can be found in Huang et al. (2025). We also provide the same procedure onto a Bernoulli setting using the EEDI dataset, details can be found in Appendix E.

We apply our methodology to produce a fidelity profile for each candidate LLM ℓ . We use total variation as the discrepancy measure, $L(p, q) = \frac{1}{2} \|p - q\|_1$, and set $\delta = 0.05$. In addition, we set the simulation budget $k = 100$, the result is presented in Figure 11.

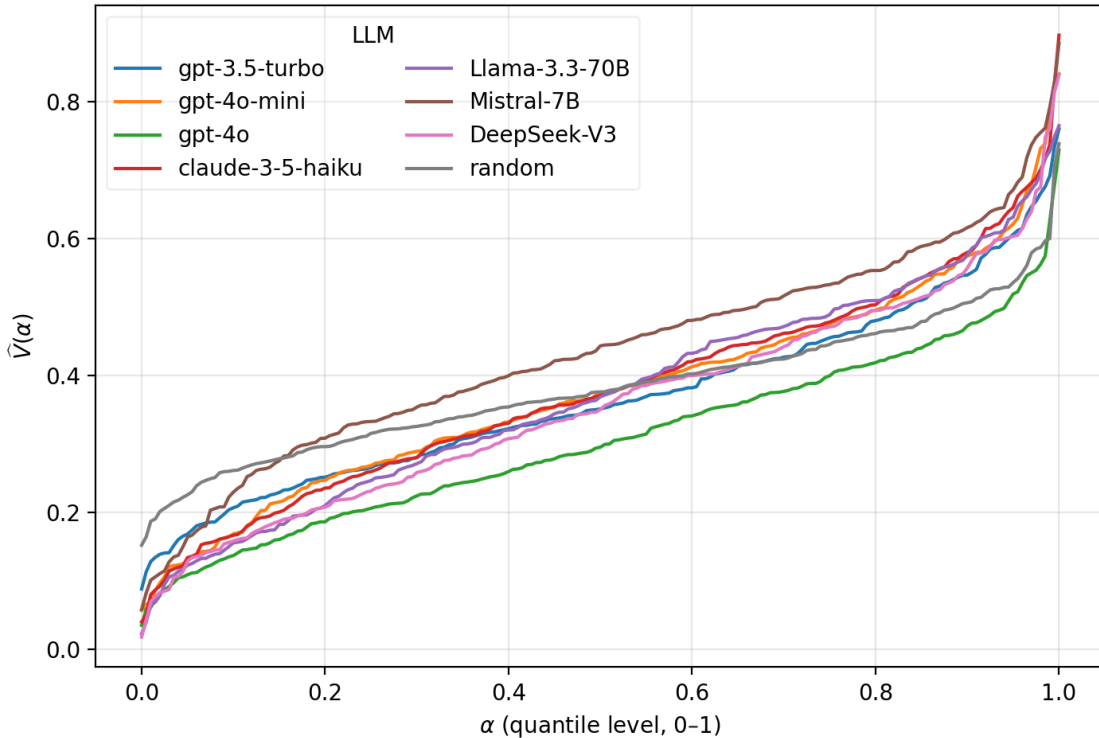


Figure 11: Quantile fidelity profiles $\hat{V}(\alpha)$ across LLMs.

Figure 11 compares models by how tightly their synthetic outcomes track the human distribution across items. We plot $\hat{V}_\ell(\alpha)$ against α , where lower-flatter curves indicate uniformly small discrepancies, while elbows reveal rare but severe misses. GPT-4O lies lowest across most quantiles, indicating the most reliable alignment, with MISTRAL 7B clearly performing worse. Notably, the simulator curves are steeper than the random benchmark, indicating question-dependent alignment and less uniform discrepancies. This suggests the simulators may require further fine-tuning to achieve more uniform discrepancy levels across this set of questions.