

Leveraging LLM-Respondents for Item Evaluation: a Psychometric Analysis

Yunting Liu¹ | Shreya Bhandari² | Zachary A. Pardos¹

¹Berkeley School of Education, University of California, Berkeley

²Electrical Engineering and Computer Science, University of California, Berkeley

Correspondence

Email: pardos@berkeley.edu

Effective educational measurement relies heavily on the curation of well-designed item pools (i.e., possessing the right psychometric properties). However, item calibration is time-consuming and costly, requiring a sufficient number of respondents for the response process. We explore using six different LLMs (GPT-3.5, GPT-4, Llama 2, Llama 3, Gemini-Pro, and Cohere Command R Plus) and various combinations of them using sampling methods to produce responses with psychometric properties similar to human answers. Results show that some LLMs have comparable or higher proficiency in College Algebra than college students. No single LLM mimics human respondents due to narrow proficiency distributions, but an ensemble of LLMs can better resemble college students' ability distribution. The item parameters calibrated by LLM-Respondents have high correlations (e.g. > 0.8 for GPT-3.5) compared to their human calibrated counterparts, and closely resemble the parameters of the human subset (e.g. 0.02 Spearman correlation difference). Several augmentation strategies are evaluated for their relative performance, with resampling methods proving most effective, enhancing the Spearman correlation from 0.89 (human only) to 0.93 (augmented human).

KEYWORDS

Large Language Models, Psychometric analysis, Item Response Theory, Simulation, Data augmentation

Practitioner Notes

- What is already known about this topic
 - Collection of human responses to candidate test items is common practice in educational measurement when designing an assessment
 - Large Language Models (LLMs) have been found to rival human abilities in a variety of subject areas
 - Data augmentation using AI has been an effective strategy for enhancing machine learning model performance
- What this paper adds
 - The first psychometric analysis of the ability distribution of a variety of open source and proprietary LLMs as compared to humans
 - Finds that 50 LLM-Respondents produce item parameters similar to 50 undergraduate respondents
 - Using LLM-Respondents to augment human response data gave mixed results
- Implications for practice and/or policy
 - The moderate performance of LLM-Respondents by themselves could provide a low-cost option for curating quality items for low stakes formative or summative assessments
 - This methodology provides a scalable way to evaluate vast amounts of generative AI-produced items

1 | INTRODUCTION

Generating sets of well-functioning items for mathematics assessments often requires multiple iterations of calibration, involving extensive human participation. One of the most common techniques used in building an item pool is Item Response Theory (IRT), where ample response-level data is typically required for accurate item calibration and scaling [1, 2]. This process is time-consuming, costly, and significantly limits the rapid adaptation of educational assessments to different sets of students. For example, the PISA main survey requires between $N = 250$ to $N = 750$ respondents per item per country [3]. Thus, the time and cost involved in obtaining responses from human respondents remains a significant area of resource expenditure.

With the advent of advanced AI technologies, novel ways to address these challenges have arisen. Recent developments in Large Language Models (LLMs) are achieving near-human performance [4, 5, 6, 7], leading to speculation about whether they can competently generate high-fidelity synthetic data without the traditional need for full data collection [8, 9, 10]. In our domain, we explore whether the capabilities of LLMs can be leveraged to provide responses resulting in psychometric properties similar to those derived from human respondents' answers. This research is guided by three critical research questions:

- **RQ1:** Which language model or configuration of language models best mimic human respondent abilities in mathematics, as measured by Item Response Theory (IRT)?
- **RQ2:** How do the psychometric properties of items fit to human responses compare to those fit to LLM-Respondents?
- **RQ3:** Can the augmentation of human respondent data with LLM-Respondent contributions yield item parameters comparable to those obtained from solely increasing human data collection?

If this approach is successful, it would mean that questions, including those produced via generative AI [11], could be tested and evaluated en masse nearly instantly for use in a variety of educational contexts such as computer tutoring systems and other formative and summative assessment scenarios.

In this study, we investigate the capabilities of various Large Language Models (LLMs), including GPT-3.5, GPT-4, Llama 2, Llama 3, Gemini-Pro, and Cohere Command R Plus, to generate assessment responses. Specifically, we prompt each model with 20 items sourced from the OpenStax Creative Commons textbook for College Algebra, producing 150 responses per model. Our analysis focuses on assessing whether these models can effectively replicate the response characteristics of our target population, which consists of undergraduate students in the United States, and on comparing the performance of these models against each other. To this end, we compare the model-generated responses to those obtained from U.S. undergraduates on the popular crowdsourcing platform, Prolific.

2 | RELATED WORK

2.1 | Simulated Data in Educational Measurement and Educational Data Mining

Analyzing examinee responses to test questions is indispensable in the field of measurement. While gathering real data can be time-consuming, costly, and often incomplete, simulation is a useful and economical technique since it can usually be done on a laptop without additional costs. Therefore, researchers commonly use simulation to validate models [12], compare different models [13], and evaluate estimation methods. In fact, among a random sample of publications in the field, 60% of the studies used simulation, while the ratio for real data is just 41% [14]. Typically, a respondent distribution is specified, and then the item response level data (i.e., dichotomous or polytomous response) is simulated accordingly, with no thought process involved. Now, thanks to the advancement of generative AI, we can simulate some response level data and possibly gain insights into the cognitive structure behind the response process. While most work on simulation in educational settings is based on dialogue [15, 16], there are indeed some researchers conducting item response level data simulation. For example, Xu and Zhang [17] demonstrated the possibility of simulating student behavior based on assessment history. Lu and Wang [18] used insights from teachers to create generative students with various profiles, and then used the generative students' outcomes to guide item development.

In the realm of educational data mining, gathering real learner data can pose privacy concerns [19] and present challenges with the costs of managing logged data [20]. To address these challenges, researchers have at times leveraged simulated data. LLMs are often used to generate the training datasets needed to train and test other models. For example, by using pseudocode to generate synthetic datasets, researchers have been able to develop test cases of teaching activities to inform the development of a Teaching Outcome Model (TOM) [21]. Similarly, researchers have proposed using LLMs as data annotators to create synthetic data that can be used to train other models [8], mimicking the framework of Teacher-Student Learning (TSL) [22]. Simulated data has also been used within the educational data mining community to evaluate latent trait models [23, 24].

2.2 | Data Augmentation

Data augmentation is a method often employed to increase the volume and diversity of data by generating new data from the existing set; it can also be applied to mitigate the 'incomplete data' problem [25]. Having a limited number of data points often leads to weaker generalization capabilities, which can act as an obstacle to the effectiveness of studies [26]. Thus, data augmentation is commonly used to enrich datasets and enhance their suitability for training models. For example, to augment image data, techniques such as resizing, rotating, and shifting images are frequently used [27, 28]. Researchers have also explored introducing noise in LLM training data [27], adding audio tracks or temporal shifts in speech recognition [29, 30], and leveraging Generative Adversarial Networks (GANs) to generate

training data for medical imaging [31], ultimately creating datasets that are more generalizable and effective.

2.3 | OER and automation

In recent years, the field of Open Educational Resources (OER) has seen significant growth and adoption [32], allowing researchers to benefit from a corpus of educational resources at no cost and open materials that can be freely distributed, remixed, and adapted. With the rise of large language models (LLMs), the education sector is experimenting with automating the generation of these resources to reduce costs and enhance efficiency. In particular, there has been an emphasis on automatic item generation, hint generation, and skill tagging. For item generation, much research is focusing on utilizing the capabilities of LLMs to generate math questions either through template-based approaches [33], open-ended generation (Socratic style questions or math word problems) [33, 34, 35, 36], multiple-choice question generation [37, 38, 39], or generation from structured formats (i.e., a bullet-point list) [11]. Hint generation has also been a focus, with researchers examining the effectiveness of LLMs in providing hints (i.e., worked solutions) to support learning in mathematics [40, 41], computer programming [42, 43, 44, 45, 46, 47], and various other STEM subjects. Additionally, studies have investigated human-AI collaboration in skill tagging, assessing its effectiveness across multiple languages and its speed and accuracy [48, 49]. However, unlike these areas, the topic of using LLMs to simulate respondents remains under-researched. Thus, this paper aims to study the feasibility and effectiveness of LLMs in simulating respondents.

3 | METHODS

3.1 | Model Selection

We selected six Large Language Models (LLMs) to generate responses that simulate answers from undergraduate college students in the U.S. to assessment questions. Our selection included GPT-3.5, GPT-4, Llama 2, the newer Llama 3, Gemini-Pro, and Cohere Command R Plus. These models were chosen for their varying levels of sophistication, reported accuracy on mathematics items, and widespread popularity, allowing us to simulate a broad spectrum of student abilities, from lower to higher academic proficiency [50, 51, 52]. For implementation, we utilized APIs for each model. As Llama does not offer a direct API, we accessed Llama 2 and Llama 3 via the Replicate API.

3.2 | Selection of Items and Prompt Engineering

College Algebra was chosen as the subject because pre-authored questions were available under a CC BY license from an open textbook publisher, OpenStax¹. Additionally, we used a dataset from an earlier study that calibrated its item pool and had already collected responses from human participants via Prolific for 20 of the OpenStax College Algebra questions in Lesson 2.2: Linear Equations in One Variable [11]. This prior data collection contained some missingness in the data, so we were able to effectively use $N \geq 99$ for all items.

We distinctively formatted the questions, each prefixed with a label, in the following format: "Q1: <question1>" followed by a double newline, then "Q2: <question2>", and continued this pattern for all 20 questions. Specifically, the prompt was:

Q1: Given $m = 4$, find the equation of the line in slope-intercept form passing through the point $(2, 5)$.

¹<https://openstax.org/details/books/college-algebra-2e>

Q2: Find the slope of a line that passes through the points $(2, -1)$ and $(-5, 3)$.

...

Q20: For the following exercises, solve the equation for x . State all x -values that are excluded from the solution set. $2 - 3/(x + 4) = (x + 2)/(x + 4)$. Answer choices: Excluded values: -4 and $x=-3$; Excluded values: 4 and $x=-3$.

This helped to ensure clarity and separation between items. We simulated 150 respondents for each LLM, as this was deemed the right sample size for conducting further analysis. To assess the accuracy of the responses, the second author of the study manually graded them and noted the accuracy for each one.

3.3 | Augmentation Procedure

In a real-world case, sometimes only partial data is gathered. To explore the possibility of augmenting the data, we treated each human respondent in our dataset as a centroid, using only 50 human responses to represent the partial data. We then identified the nearest synthetic respondent from our pool of synthetically generated answers for each human centroid, allowing us to map the AI responses directly onto the characteristics of individual human responses. Next, we conducted a resampling procedure. First, we resampled a subset of 50 synthetic responses, selecting them based on the distribution of the models represented in the original matched 50, aiming to maintain the proportionality observed in this initial sampling. Lastly, we expanded this by resampling a subset of 100 synthetic responses using the same criteria.

3.4 | IRT analysis

Contrary to sum-score analysis or percentage correct metrics, we plan to use Item Response Theory (IRT) to estimate the latent ability of human and LLM-Respondents. This method has several advantages over sum-score analysis [53]. First, IRT assumes a latent trait θ , transforming all estimations onto a logit scale instead of the sum-score scale, which greatly improves measurement precision. Second, IRT provides person-level fit data, which can be done independently of other respondents. Lastly, and most relevant to our purpose, IRT maps both persons and items onto the same scales, enabling equating to be carried out without assuming population score distributions. In fact, IRT equating may be the best method when tests of differing difficulties are given to nonrandom groups of examinees who differ in ability [54].

The simplest IRT model is often called the Rasch model or the one-parameter logistic model (1PL). The probability of individuals responding to a binary item (i.e., True/False) is determined by the individual's trait level and the difficulty of the item, which is often presented as:

$$P(X_{ij} = 1 | \theta_i, \beta_j) = \frac{\exp(\theta_i - \beta_j)}{1 + \exp(\theta_i - \beta_j)}$$

where:

X_{ij} refers to the response made by individual i to item j . If the response is correct/true, then $X_{ij} = 1$.

θ_i refers to the trait level of individual i .

β_j refers to the difficulty of item j .

A significant merit of using the Rasch model is that the estimates of latent trait and difficulty are mapped onto the same logit scale. The logits are interval units, which make the Wright Map (also known as the item-person map) a useful tool for presenting both item difficulties and person abilities arranged along the same logit scale. The location of item difficulty denotes the ability level at which the individual has approximately a 50% probability of answering the item correctly. The results from human calibration using Rasch analysis are shown in Figure 1. Items were ranked by their difficulties in ascending order from top to bottom [55].

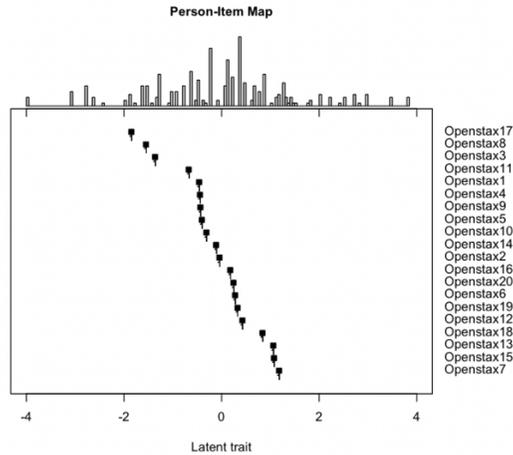


FIGURE 1 Item parameters calibrated by human respondents

IRT analysis will be carried out for two purposes: Firstly, we will use the item parameters calibrated by human respondents as a standard so that the parameters for LLM-Respondents can be estimated. By treating item parameters as fixed, we will be able to compare the proficiency distribution of college students and LLM-Respondents. Within the fixed parameter calibration (FPC) realm, we will choose multiple weights updating and multiple EM cycles (MWU-MEM) as they are the most robust estimation methods [56]. Secondly, separate calibrations will be carried out on human and AI respondent groups, informing us of the practicality of item calibration using AI generated data.

4 | RESULTS

4.1 | LLM-Respondent Simulation

The initial item parameters for the item pool were calibrated on a group of current college students in the United States. Since multiple forms were used, there was missingness at random in the data, effectively resulting in $N \geq 99$ for all items, satisfying the basic requirement for a Rasch analysis. The results are shown in Figure 1. We then fixed the item parameters estimated from the model to obtain the proficiency estimates for the six AI models. We wondered whether the proficiency distribution of synthetic respondents is comparable to that of humans. Results show that most LLM proficiency distributions have a significant overlap with the human respondents. In particular, Llama 3 and GPT-3.5 have the highest mean proficiency distribution, which is higher than the human mean, indicating AI's greater proficiency in College Algebra compared to college students. The mean proficiency of GPT-4 is comparable to humans,

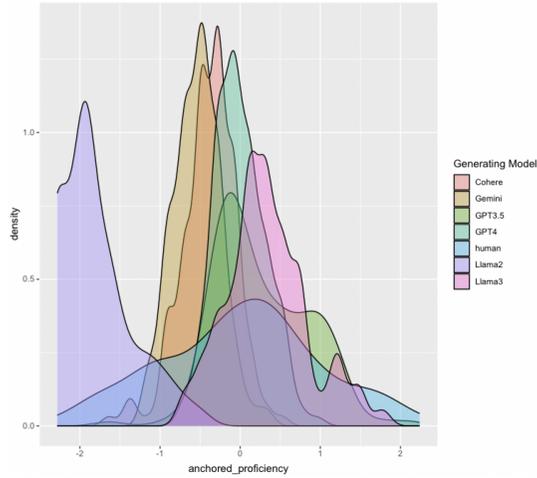


FIGURE 2 Proficiency distribution by Generating Models

while Cohere and Gemini are lower than humans. Llama 2 is the worst among all six models, suggesting its incapability in solving College Algebra problems.

Besides the findings on the average proficiency of each model, an interesting observation is that the variability in proficiency distribution for LLM-Respondents (standard deviation ranging from 0.29 to 0.58) was greatly reduced compared to the human distribution (SD = 0.98). Most AI distributions appear to be sharper than their human counterparts (kurtosis larger than 3). Detailed distributional information can be found in Table 1 and Figure 2. Further analysis might explore ways to increase the variability in proficiency and response patterns.

TABLE 1 Statistical Measures of Response Distributions by Model

Generating Model	Mean	Standard Deviation (SD)
Cohere	-0.40	0.34
GPT3.5	0.27	0.58
GPT4	0.00	0.31
Gemini	-0.54	0.29
Llama2	-1.81	0.44
Llama3	0.37	0.51
human	0.00	0.98

We also performed item calibration on fully simulated data, comparing it against human-calibrated data and AI-calibrated data. Notably, the IRT estimated item difficulties show a significant correlation between AI-generated and human-generated data, especially for GPT-3.5 and GPT-4 (Pearson $\rho > 0.7$), as shown in Table 2. Rank correlation (i.e., Spearman correlation) is also reported since not all distributions have the same mean, and the difficulties scale might not be an interval scale for different estimates. Given that one Gen-AI model might not have the power to represent the human distribution, we decided to use an ensemble method (also denoted as data augmentation) to

create a representative pool of LLM-Respondents whose similarity with human respondents was boosted.

TABLE 2 Evaluation Metrics for Simulation and Augmentation experiments

Generating Source	Pearson ρ	Spearman ρ	RMSE
GPT3.5	0.83	0.87	1.90
GPT4	0.72	0.78	2.27
Cohere	0.65	0.72	2.02
Gemini	0.68	0.61	2.67
Llama3	0.18	0.32	4.03
Llama2	-0.01	-0.07	1.37
Experiment 1	0.92	0.89	0.55
Experiment 2	0.91	0.89	0.54
Experiment 3	0.93	0.93	0.71
Experiment 4	0.75	0.86	1.78

4.2 | Data Augmentation using LLM-Respondent

Given that none of the LLM models have a proficiency distribution resembling that of humans, it is not feasible at this time to fully substitute human respondents with LLM-Respondents from a single LLM. However, LLMs could be used in a hybrid approach where half the respondents are human and half are LLM-Respondents. With this in mind, we propose three hybrid, or data augmentation, strategies listed below:

- An enlarged sample of 50 humans: 50 human respondents (examples) and 50 LLM-Respondents
- A mixture of human respondents and resampled LLM-Respondents using proportions learned from humans
- Fully LLM-Respondents using the mixing proportions learned from humans

The resampling analysis resulted in a set with significant variation in the proportions of each model used. GPT-3.5 was the most prevalent, comprising 36% of the synthetic responses, followed by Llama 2 at 3%. Gemini accounted for 12%, while both Llama 3 and GPT-4 were represented at 8% each. Cohere was the least represented model, constituting only 6% of the responses.

To evaluate the relative performance of these three strategies, we proposed four experiments to test the effectiveness of different strategies on the item calibration process. The benchmark performance of the calibration is set by the human respondents; namely, we use all available data from human respondents to calibrate the item pool and gather item parameter estimates. Each experiment is designed to explore a different augmentation strategy. The proposed experiments are as follows:

Experiment 1 We use only half the number of the respondent pool and do the calibration ($N = 50$), representing a real-world scenario where there is a limitation in the budget, so only part of the intended respondents were collected.

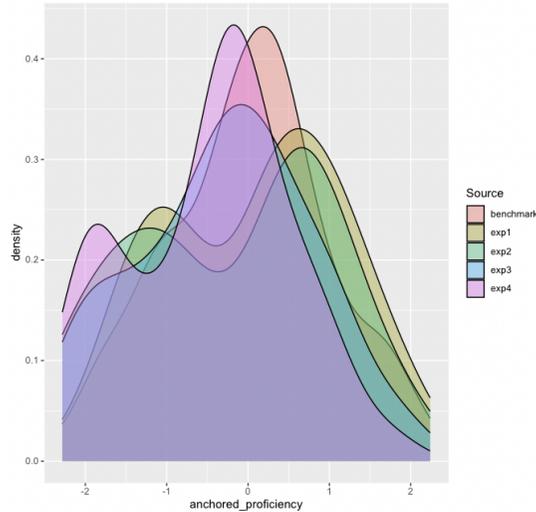


FIGURE 3 Proficiency distribution by Augmentation Experiments

Experiment 2 In addition to Experiment 1, we enlarge the data size by twice using augmentation strategies ($N = 100$).

Experiment 3 We use a mixture of human respondents and fully resampled data in a 1:1 ratio ($N = 100$).

Experiment 4 We use fully resampled data, with a size equal to the number of the benchmark dataset ($N = 100$).

In terms of evaluation criteria, Pearson correlation and Spearman correlation with the benchmark condition are reported. We also use Root Mean Square Error (RMSE) to evaluate how accurate the new estimates are. RMSE measures the average difference between values predicted by a model and the actual values. $RMSE = \sqrt{\frac{\sum_{i=1}^N \|y(i) - \hat{y}(i)\|^2}{N}}$. RMSE is always non-negative, and a value of 0 (almost never achieved in practice) would indicate a perfect recovery of the true/benchmark data.

Results show that when there aren't a sufficient number of respondents (Experiment 1), some difficulty estimates from the IRT analysis are not trustworthy, especially at the very ends of the respondent distribution (too-easy or too-hard items), resulting in a relatively high but not perfect correlation (Pearson $\rho = 0.92$, Spearman $\rho = 0.89$), and $RMSE = 0.55$. Among all augmentation strategies, Experiment 3 has the best result; it raised the Spearman ρ from 0.89 to 0.93, indicating it is an effective strategy for recovering the order of item difficulties. However, since the mean center of the respondent distribution in Experiment 3 is tilted to the left ($\mu = -0.29$) compared to Experiment 1 ($\mu = 0.08$), the RMSE is larger as a universal shift is applied to all item parameters. Experiment 2 yields results comparable to Experiment 1, suggesting the strategy might not be beneficial in the current settings; practical reasons will be discussed in the Discussion section. As a validation, we also calibrated the respondent proficiency distribution using the fixed item parameter methods, and the results are shown in Figure 3.

5 | DISCUSSION AND CONCLUSIONS

In this study, we explored six different LLMs (GPT-3.5, GPT-4, Llama 2, Llama 3, Gemini-Pro, and Cohere Command R Plus) and various combinations using sampling methods to achieve psychometric properties similar to those from

human respondents' answers. Our findings are structured around three key conclusions: the proficiency of LLMs in approximating ability distributions (RQ1), item parameter correlation (RQ2), and the effectiveness of data augmentation (RQ3). Firstly, when comparing the proficiency between LLM responses and human responses, the results show that although some LLMs have comparable or even higher abilities in College Algebra, their distributions alone cannot fully represent the human distribution due to their narrow spread. Interestingly, this novel application of Item Response Theory (IRT) to LLM abilities reveals a first-of-its-kind distribution spread of abilities from multiple promptings, as opposed to the point estimates or percent-correct scores reported in other studies. Secondly, we compared the item parameters calibrated by AI responses and human responses, finding a relatively high Spearman correlation of 0.87 for GPT-3.5 and a lower 0.78 for GPT-4. Notably, GPT-3.5 emerged as the most human-like AI respondent, exhibiting Spearman correlations within 0.02 of those from 50 human respondents. Finally, since no single LLM currently has the capability to represent humans by itself, we explored ensembling approaches using three strategies. Among these, a mixture of human respondents and fully resampled data in a 1:1 ratio (Experiment 3) provided the best result, raising the Spearman correlation from 0.89 to 0.93. However, these augmentation results were mixed; while Spearman and Pearson correlations improved by 0.04 and 0.01, respectively, this approach substantially increased the RMSE.

Our findings hold much promise for the automatic curation of items for tutoring systems. Simply put, it seems plausible to leverage AI respondents to curate an item pool that has a desirable spread of difficulty. Given the performance of 150 AI respondents from GPT-3.5 closely mirroring that of 50 humans (Experiment 1), AI respondents could be used as an initial filtering phase to reliably narrow down a larger item pool, and then have human respondents further refine the selection using the more manageable subset of items. This would significantly help optimize human resources, saving both time and money. For classroom environments, this research allows nimble testing of new questions, enabling selection of only quality assessments to present in the classroom.

6 | LIMITATIONS AND FUTURE WORK

Our study still has limitations. Firstly, we only utilized a single College Algebra lesson, which makes it difficult to generalize the results to other lessons within the same subject or to different subject areas. Additionally, our OpenStax item pool consists only of questions without images, figures, or tables because not all the LLMs support multimodal capabilities. Furthermore, in our experiment, the original human dataset displays bimodality. Therefore, when augmenting the data without resampling (Experiment 2), bimodality was also exhibited. While Experiment 3 mitigates this impact by using an effective resampling strategy, it would be beneficial to utilize a pool of human respondents that portray a normal distribution from the start. Due to leveraging prior data collection, we used the respondent sample size from their study, rather than calculating what the effective size should actually be.

In the future, analyses should be conducted to determine how valid the estimated proficiencies are when derived from a measurement tool calibrated by augmentation experiments versus those from benchmark data. To address this, we should investigate whether the structure of the tool with augmented respondents is the same as it is with the original human population [57]. Typically, this is done by investigating the reliability of the measurement tool and conducting statistical analyses such as Explanatory Factor Analysis (EFA) or Confirmatory Factor Analysis (CFA) [58].

The field should also aim to refine and expand methodologies to significantly improve the accuracy of responses generated by the models. In our study, we presented all 20 questions in a single prompt. However, it may be beneficial to question the model independently for each item. Although our initial trials with this technique for GPT-3.5 resulted in degenerate outcomes, this method could be extended to other models to assess its effectiveness more comprehensively. Additionally, more sophisticated prompt engineering techniques should be explored. Our prompt was simple

and uniform across all models. However, it may be advantageous to customize prompts based on the model to better suit the specific strengths and design of each model. Incorporating few-shot learning by including one or two example responses in the prompt may also be effective. Finally, employing hallucination mitigation techniques such as self-consistency [59], in which the model is prompted multiple times for the same question and the modal response is selected, may help minimize erroneous outputs and thus increase accuracy. This technique has proven effective in reducing hallucination rates to near zero for answering mathematics questions [41]. Thus, expanding this technique, along with exploring other hallucination mitigation methods, may prove beneficial.

references

- [1] König C, Khorrarnel L, Yamamoto K, Frey A. The benefits of fixed item parameter calibration for parameter accuracy in small sample situations in large-scale assessments. *Educational Measurement: Issues and Practice* 2021;40(1):17–27.
- [2] Kim S, Kolen MJ. Application of IRT fixed parameter calibration to multiple-group test data. *Applied Measurement in Education* 2019;32(4):310–324.
- [3] Mazzeo J, von Davier M. Linking scales in international large-scale assessments. *Handbook of international large-scale assessment: Background, technical issues, and methods of data analysis* 2014;p. 229–258.
- [4] Katz DM, Bommarito MJ, Gao S, Arredondo P. Gpt-4 passes the bar exam. *Philosophical Transactions of the Royal Society A* 2024;382(2270):20230254.
- [5] Chang Y, Wang X, Wang J, Wu Y, Yang L, Zhu K, et al. A survey on evaluation of large language models. *ACM Transactions on Intelligent Systems and Technology* 2023;.
- [6] Liu X, Yu H, Zhang H, Xu Y, Lei X, Lai H, et al. Agentbench: Evaluating llms as agents. *arXiv preprint arXiv:230803688* 2023;.
- [7] Achiam J, Adler S, Agarwal S, Ahmad L, Akkaya I, Aleman FL, et al. Gpt-4 technical report. *arXiv preprint arXiv:230308774* 2023;.
- [8] Ding B, Qin C, Zhao R, Luo T, Li X, Chen G, et al. Data augmentation using llms: Data perspectives, learning paradigms and challenges. *arXiv preprint arXiv:240302990* 2024;.
- [9] Ye J, Xu N, Wang Y, Zhou J, Zhang Q, Gui T, et al. LLM-DA: Data Augmentation via Large Language Models for Few-Shot Named Entity Recognition. *arXiv preprint arXiv:240214568* 2024;.
- [10] Whitehouse C, Choudhury M, Aji AF. Llm-powered data augmentation for enhanced crosslingual performance. *arXiv preprint arXiv:230514288* 2023;.
- [11] Bhandari S, Liu Y, Pardos ZA. Evaluating ChatGPT-generated Textbook Questions using IRT. In: *Proceedings of the Generative AI for Education Workshop (GAIED) at the Thirty-seventh Conference on Neural Information Processing Systems (NeurIPS)*. New Orleans, LA; 2023. .
- [12] Swaminathan H, Hambleton RK, Rogers HJ. 21 Assessing the Fit of Item Response Theory Models. *Handbook of statistics* 2006;26:683–718.
- [13] Harwell M, Stone CA, Hsu TC, Kirisci L. Monte Carlo studies in item response theory. *Applied psychological measurement* 1996;20(2):101–125.
- [14] Feinberg RA, Rubright JD. Conducting simulation studies in psychometrics. *Educational Measurement: Issues and Practice* 2016;35(2):36–49.

-
- [15] Shaikh O, Chai VE, Gelfand M, Yang D, Bernstein MS. Rehearsal: Simulating conflict to teach conflict resolution. In: Proceedings of the CHI Conference on Human Factors in Computing Systems; 2024. p. 1–20.
- [16] Markel JM, Opferman SG, Landay JA, Piech C. Gpteach: Interactive ta training with gpt-based students. In: Proceedings of the tenth acm conference on learning@ scale; 2023. p. 226–236.
- [17] Xu S, Zhang X. Leveraging generative artificial intelligence to simulate student learning behavior. arXiv preprint arXiv:231019206 2023;.
- [18] Lu X, Wang X. Generative Students: Using LLM-Simulated Student Profiles to Support Question Item Evaluation. arXiv preprint arXiv:240511591 2024;.
- [19] Hutt S, Das S, Baker RS. The Right to Be Forgotten and Educational Data Mining: Challenges and Paths Forward. International Educational Data Mining Society 2023;.
- [20] Jacob J, Jha K, Kotak P, Puthran S. Educational data mining techniques and their applications. In: 2015 International Conference on Green Computing and Internet of Things (ICGCIoT) IEEE; 2015. p. 1344–1348.
- [21] Ndukwe IG, Daniel BK, Butson RJ. Data science approach for simulating educational data: Towards the development of teaching outcome model (TOM). Big Data and Cognitive Computing 2018;2(3):24.
- [22] Hu C, Li X, Liu D, Chen X, Wang J, Liu X. Teacher-student architecture for knowledge learning: A survey. arXiv preprint arXiv:221017332 2022;.
- [23] Badrinath A, Wang F, Pardos Z. pyBKT: An Accessible Python Library of Bayesian Knowledge Tracing Models. In: Proceedings of the 14th International Conference on Educational Data Mining; 2021. p. 468–474.
- [24] Beheshti B, Desmarais MC, Naceur R. Methods to Find the Number of Latent Skills. International Educational Data Mining Society 2012;.
- [25] Seltzer MH. The use of data augmentation in fitting hierarchical models to educational data. PhD thesis, The University of Chicago; 1991.
- [26] Kieser F, Wulff P, Kuhn J, Küchemann S. Educational data augmentation in physics education research using ChatGPT. Physical Review Physics Education Research 2023;19(2):020150.
- [27] Xie Z, Wang SI, Li J, Lévy D, Nie A, Jurafsky D, et al. Data noising as smoothing in neural network language models. arXiv preprint arXiv:170302573 2017;.
- [28] Shorten C, Khoshgoftaar TM. A survey on image data augmentation for deep learning. Journal of big data 2019;6(1):1–48.
- [29] Deng L, Acero A, Plumpe M, Huang X. Large-vocabulary speech recognition under adverse acoustic environments. In: INTERSPEECH Citeseer; 2000. p. 806–809.
- [30] Hannun A, Case C, Casper J, Catanzaro B, Diamos G, Elsen E, et al. Deep speech: Scaling up end-to-end speech recognition. arXiv 2014. arXiv preprint arXiv:14125567 2014;.
- [31] Frid-Adar M, Diamant I, Klang E, Amitai M, Goldberger J, Greenspan H. GAN-based synthetic medical image augmentation for increased CNN performance in liver lesion classification. Neurocomputing 2018;321:321–331.
- [32] Henderson S, Ostaszewski N. Barriers, incentives, and benefits of the open educational resources (OER) movement: An exploration into instructor perspectives. First Monday 2018;.
- [33] Shridhar K, Macina J, El-Assady M, Sinha T, Kapur M, Sachan M, Automatic Generation of Socratic Subquestions for Teaching Math Word Problems; 2022.

- [34] Zhou Z, Ning M, Wang Q, Yao J, Wang W, Huang X, et al., Learning by Analogy: Diverse Questions Generation in Math Word Problem; 2023.
- [35] Keller SU, Automatic Generation of Word Problems for Academic Education via Natural Language Processing (NLP); 2021.
- [36] Onal S, Kulavuz-Onal D. A Cross-Disciplinary Examination of the Instructional Uses of ChatGPT in Higher Education. *Journal of Educational Technology Systems* 0;0(0):00472395231196532. <https://doi.org/10.1177/00472395231196532>.
- [37] CH DR, Saha SK. Generation of Multiple-Choice Questions From Textbook Contents of School-Level Subjects. *IEEE Transactions on Learning Technologies* 2023;16(1):40–52.
- [38] Nagasaka K. Multiple-choice questions in mathematics: Automatic generation, revisited. In: The 25th Asian technology conference in mathematics, virtual format, Radford University, Virginia, USA and Suan Sunandha Rajabhat University, Thailand; 2020. .
- [39] Lee J, Smith D, Woodhead S, Lan A. Math Multiple Choice Question Generation via Human-Large Language Model Collaboration. *arXiv preprint arXiv:240500864* 2024;.
- [40] Pardos ZA, Bhandari S. Learning gain differences between ChatGPT and human tutor generated algebra hints. *arXiv preprint arXiv:230206871* 2023;.
- [41] Pardos ZA, Bhandari S. ChatGPT-generated help produces learning gains equivalent to human tutor-authored help on mathematics skills. *Plos one* 2024;19(5):e0304013.
- [42] Price TW, Dong Y, Zhi R, Paaßen B, Lytle N, Cateté V, et al. A comparison of the quality of data-driven programming hint generation algorithms. *International Journal of Artificial Intelligence in Education* 2019;29:368–395.
- [43] Rivers K, Koedinger KR. Automating hint generation with solution space path construction. In: *Intelligent Tutoring Systems: 12th International Conference, ITS 2014, Honolulu, HI, USA, June 5-9, 2014. Proceedings* 12 Springer; 2014. p. 329–339.
- [44] Piech C, Sahami M, Huang J, Guibas L. Autonomously generating hints by inferring problem solving policies. In: *Proceedings of the second (2015) acm conference on learning@ scale*; 2015. p. 195–204.
- [45] Buwalda M, Jeuring J, Naus N. Use Expert Knowledge Instead of Data: Generating Hints for Hour of Code Exercises. In: *Proceedings of the Fifth Annual ACM Conference on Learning at Scale L@S '18, New York, NY, USA: Association for Computing Machinery*; 2018. <https://doi.org/10.1145/3231644.3231690>.
- [46] Price TW, Dong Y, Barnes T. Generating data-driven hints for open-ended programming. *International Educational Data Mining Society* 2016;.
- [47] Roy Choudhury R, Yin H, Fox A. Scale-driven automatic hint generation for coding style. In: *Intelligent Tutoring Systems: 13th International Conference, ITS 2016, Zagreb, Croatia, June 7-10, 2016. Proceedings* 13 Springer; 2016. p. 122–132.
- [48] Ren C, Pardos Z, Li Z. Human-AI Collaboration Increases Skill Tagging Speed but Degrades Accuracy. *arXiv preprint arXiv:240302259* 2024;.
- [49] Kwak Y, Pardos ZA. Bridging large language model disparities: Skill tagging of multilingual educational content. *British Journal of Educational Technology* 2024;.
- [50] Touvron H, Lavril T, Izacard G, Martinet X, Lachaux MA, Lacroix T, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:230213971* 2023;.
- [51] Plevris V, Papazafeiropoulos G, Jiménez Rios A. Chatbots Put to the Test in Math and Logic Problems: A Comparison and Assessment of ChatGPT-3.5, ChatGPT-4, and Google Bard. *AI* 2023;4(4):949–969.

-
- [52] Team G, Anil R, Borgeaud S, Wu Y, Alayrac JB, Yu J, et al. Gemini: a family of highly capable multimodal models. arXiv preprint arXiv:2312.11805 2023;.
- [53] DeMars C. Item response theory. Oxford University Press; 2010.
- [54] Cook LL, Eignor DR. IRT equating methods. *Educational measurement: Issues and practice* 1991;10(3):37–45.
- [55] Wilson M. Constructing measures: An item response modeling approach. Taylor & Francis; 2023.
- [56] Kim S. A comparative study of IRT fixed parameter calibration methods. *Journal of educational measurement* 2006;43(4):355–381.
- [57] Ahmad S, Zulkurnain N, Khairushalimi F. Assessing the validity and reliability of a measurement model in Structural Equation Modeling (SEM). *British Journal of Mathematics & Computer Science* 2016;15(3):1–8.
- [58] Arafat S, Chowdhury HR, Qusar M, Hafez M. Cross cultural adaptation and psychometric validation of research instruments: a methodological review. *Journal of Behavioral Health* 2016;5(3):129–136.
- [59] Wang X, Wei J, Schuurmans D, Le Q, Chi E, Zhou D. Self-consistency improves chain of thought reasoning in language models. arXiv preprint arXiv:2203.11171 2022;.