# PATCH! Psychometrics-AssisTed BenCHmarking of Large Language models: A Case Study of Proficiency in 8th Grade Mathematics

Qixiang Fang[0000−0003−2689−6653], Daniel L. Oberski[0000−0001−7467−2297], and Dong Nguyen[0000−0002−6062−3117]

Utrecht University, Utrecht, Netherlands
{q.fang,d.l.oberski,d.p.nguyen}@uu.nl

**Abstract.** Many existing benchmarks of large (multimodal) language models (LLMs) focus on measuring LLMs' academic proficiency, often with also an interest in comparing model performance with human test takers'. While such benchmarks have proven key to the development of LLMs, they suffer from several limitations, including questionable measurement quality (e.g., Do they measure what they are supposed to in a reliable way?), lack of quality assessment on the item level (e.g., Are some items more important or difficult than others?) and unclear human population reference (e.g., To whom can the model be compared?). In response to these challenges, we propose leveraging knowledge from psychometrics - a field dedicated to the measurement of latent variables like academic proficiency - into LLM benchmarking. We make three primary contributions. First, we introduce PATCH: a novel framework for **P**sychometrics-**A**ssis**T**ed ben**CH**marking of LLMs. PATCH addresses the aforementioned limitations, presenting a new direction for LLM research. Second, we implement PATCH by measuring several LLMs' proficiency in 8th grade mathematics against 56 human populations. We show that adopting a psychometrics-based approach yields evaluation outcomes that diverge from those based on current benchmarking practices. Third, we release 4 high-quality datasets to support measuring and comparing LLM proficiency in grade school mathematics and science with human populations.

**Keywords:** Item response theory · Vision language models · Evaluation

## 1 Introduction

Large language models (LLMs), including their multimodal variants like vision language models, have witnessed significant advancements in recent years. These models are typically evaluated on established benchmarks that assess their performance across a diverse set of tasks such as *commonsense reasoning* [46,37,8], *coding* [8,17] and *academic proficiency*. Academic proficiency, in particular, has become a crucial part of LLM evaluation, as evidenced by the large number of related benchmarks like MMLU, ARC, GSM8K, DROP and MATH [18,9,10,13,18],

as well as recent model technical reports' increasing focus on them [31,17]. In these benchmarks and reports, LLM performance is also often contrasted with human performance.

Despite the success of existing benchmarks in advancing LLM research, they have limitations. The first concern is measurement quality: *Do these benchmarks measure what they are supposed to in a reliable way?* Many benchmarks are created via crowd-sourced knowledge, by asking a convenience group of individuals (e.g., crowd workers, paper authors) to create new test items (e.g., GSM8K, DROP) or collecting them from (often undocumented) sources (e.g., websites, textbooks, school exams) (e.g., MATH, MMLU, ARC). Without domain expert input and rigorous testing of item quality, undesirable outcomes can occur, including a mismatch between a benchmark and its claimed measurement goal, missing information in a question, wrong answer keys, and low data annotation agreement [30].[1]

Second, current benchmarks do not account for differences across test items, such as item discriminativeness[2] and difficulty (see Section 3.1). For example, consider three items A (easy), B (hard) and C (hard). While answering correctly to A and B would result in the same accuracy score as answering correctly to B and C, the latter (i.e., answering correctly to more difficult items) would imply higher proficiency. Furthermore, items that are too easy or too difficult (i.e., low discriminativeness) will fail to differentiate models of different proficiency levels. Thus, without accounting for item differences, benchmarking results, especially model rankings, can be misleading.

Third, while many benchmarks compare LLMs against humans, the human populations under comparison remain unclear [41]. For instance, human performance in MATH is based on the authors; in MMLU, crowd workers; in MATH, 6 university students. Using such convenience samples (with little information about sample characteristics), the measured human performance cannot be generalised to other human samples or populations beyond that specific sample.

To address these challenges, we propose leveraging insights from psychometrics - a field dedicated to the measurement of *latent* variables like academic proficiency - into LLM benchmarking processes. In particular, we draw on two research areas in psychometrics: *item response theory* (see Section 3.1) and *test development* (see Section 3.2 and 3.3). The former enables more accurate estimation of academic proficiency, compared to common practice in LLM benchmarks (e.g., means, percentages). It can also provide diagnostic information about the quality of each test item. The latter, test development knowledge, can help to build high quality LLM benchmarks where valid comparison to specific human populations can be made.

Our paper makes three primary contributions. First, we present **PATCH**: a novel framework for **P**sychometrics-**A**ssis**T**ed ben**CH**marking of LLMs, which

---

[1] We avoid calling out specific datasets here, but a quick Internet search would reveal many blogs reporting large percentages of errors in existing LLM benchmarks.

[2] In psychometrics, the term "item discrimination" is used. However, given the ambiguity and negative connotation of "discrimination", we adopt "discriminativeness".

addresses the aforementioned limitations of existing benchmarks. Second, we demonstrate the implementation of PATCH by testing several LLMs' proficiency in 8th grade mathematics using the released test items and data from *Trends in International Mathematics and Science Study*[3] (TIMSS) 2011. We show empirically that a psychometrics-based approach can lead to evaluation outcomes that diverge from those obtained through conventional benchmarking practices and that are more informative, underscoring the potential of psychometrics to reshape the LLM benchmarking landscape. Third, we make our evaluation code based on the PATCH framework available[4], along with three other mathematics and science datasets based on TIMSS 2011 and 2008[5].

## 2   Related Work

We are not the first to propose leveraging psychometrics for research on LLMs and other areas in NLP. For instance, psychometric scales have been used to examine the psychological profiles of LLMs such as personality traits and motivations [19,32,11]. The text in these scales can also be used to improve encoding and prediction of personality traits [21,43,45,14]. Psychometrics-based reliability and validity tests have also been proposed or/and used to assess the quality of NLP bias measures [12,44], text embeddings [15], political stance detection [40], annotations [2], user representations [16], and general social science constructs [4].

The most closely related work to our paper is the use of IRT models in NLP for constructing more informative test datasets [23], comparison of existing evaluation datasets and instances (e.g., difficulty, discriminativeness) [39,42,35,22,36], as well as identification of difficult instances from training dynamics [25,24]. Our work distinguishes itself from these papers in two aspects. First, we do not apply IRT to *existing* LLM datasets/benchmarks. Instead, we introduce a framework for benchmarking LLMs by leveraging both IRT and test development knowledge from psychometrics. The goal of this framework is to generate new, high-quality benchmarks for LLMs that warrant valid comparison with human populations. Second, we demonstrate our framework with a mathematics proficiency test validated on 56 human populations, and compare LLM performance with human performance. To the best our knowledge, we are the first to apply psychometrically validated (mathematics) proficiency tests to LLMs and make valid model versus human comparison.

## 3   Preliminaries

### 3.1   Item Response Theory

Item response theory (IRT) refers to a family of mathematical models that describe the functional relationship between responses to a test item, the test item's

---

[3] http://timssandpirls.bc.edu/timss2015/encyclopedia/

[4] https://github.com/fqixiang/patch_llm_benchmarking_with_psychometrics

[5] https://zenodo.org/records/12531906

characteristics (e.g., item difficulty and discriminativeness) and test taker's standing on the latent construct being measured (e.g., academic proficiency) [1]. Unlike classical test theory and current LLM benchmarks, which focus on the total or mean score of a test, IRT models takes into account the characteristics of both the items and the individuals being assessed, offering advantages like item quality diagnostics and more accurate estimation of test takers' proficiency. As such, IRT models have gained widespread adoption in various fields, including education, psychology, and healthcare, where trustworthy measurement and assessment are crucial.

We describe below three fundamental IRT models suitable for different types of test items: the 3-parameter logistic (3PL) model for multiple choice items scored as either incorrect or correct, the 2-parameter logistic (2PL) model for open-ended response items scored as either incorrect or correct, as well as the generalised partial credit (GPC) model for open-ended response items scored as either incorrect, partially correct, or correct.

The 3PL model gives the probability that a test taker, whose proficiency is characterised by the latent variable $\theta$, will respond correctly to item $i$:

$$P\left(x_i = 1 \mid \theta, a_i, b_i, c_i\right) = c_i + \frac{1 - c_i}{1 + \exp\left(-1.7 \cdot a_i \cdot (\theta - b_i)\right)} \equiv P_{i,1}\left(\theta\right) \quad (1)$$

where $x_i$ is the scored response to item $i$ (1 if correct and 0 if incorrect); $\theta$ is the proficiency of the test taker, where a higher value implies a greater probability of responding correctly; $a_i$ is the slope parameter of item $i$, characterising its discriminativeness (i.e., how well the item can tell test takers with higher $\theta$ from those with lower $\theta$)[6]; $b_i$ is the location parameter of item $i$, characterising its difficulty; $c_i$ is the lower asymptote parameter of item $i$, reflecting the chances of test takers with very low proficiency selecting the correct answer (i.e., guessing). Correspondingly, the probability of an incorrect response to item $i$ is: $P_{i,0} = P\left(x_i = 0 \mid \theta_k, a_i, b_i, c_i\right) = 1 - P_{i,1}\left(\theta_k\right)$. The 2PL model has the same form as the 3PL model (Equation 1), except that the $c_i$ parameter is fixed at zero (i.e., no guessing).

The GPC model [29] gives the probability that a test taker with proficiency $\theta$ will have, for the $i^{\text{th}}$ item, a response $x_i$ that is scored in the $l^{\text{th}}$ of $m_i$ ordered score categories:

$$P\left(x_i = l \mid \theta, a_i, b_i, d_{i,1}, \cdots, d_{i,m_i-1}\right) = \frac{\exp\left(\sum_{v=0}^{l} 1.7 \cdot a_i \cdot (\theta - b_i + d_{i,v})\right)}{\sum_{g=0}^{m_i-1} \exp\left(\sum_{v=0}^{g} 1.7 \cdot a_i \cdot (\theta - b_i + d_{i,v})\right)}$$
$$\equiv P_{i,l}\left(\theta\right)$$
$$(2)$$

where $m_i$ is the number of response score categories for item $i$; $x_i$ is the response score of item $i$ between 0 and $m_i - 1$ (e.g., 0, 1 and 2, for incorrect,

---

[6] The number 1.7 is a scaling parameter to preserve historical interpretation of parameter $a_i$ on the normal ogive scale [6]. Also applies to 2PL and GPC models.

partially correct, and correct responses); $\theta$, $a_i$, $b_i$ have the same interpretations as in the 3PL and 2PL models; $d_{i,1}$ is the category $l$ threshold parameter. Setting $d_{i,0} = 0$ and $\sum_{j=1}^{m_i-1} d_{i,j} = 0$ resolves the indeterminacy of the model parameters.

Assuming conditional independence, the joint probability of a particular response pattern $x$ across a set of $n$ items is given by:

$$P\left(x \mid \theta, \text{ item parameters }\right) = \prod_{i=1}^{n} \prod_{l=0}^{m_i-1} P_{i,l}\left(\theta\right)^{u_{i,l}} \tag{3}$$

where $P_{i,l}\left(\theta\right)$ is of the form specific to the type of item (i.e., 3PL, 2PL or GPC); $m_i$ equals 2 for dichotomously scored items and 3 for polytomously scored items; $u_{i,l}$ is an indicator defined as:

$$u_{i,l} = \begin{cases} 1 \text{ if response } x_i \text{ is in category } l \\ 0 \text{ otherwise} \end{cases}$$

This function can be viewed as a likelihood function to be maximised by the item parameters. With the estimated item parameters, $\theta$ can then be estimated [34].

### 3.2 Test Development in Psychometrics

| Psychometrics | LLM Benchmarking |
|---|---|
| 1. Construct and test need specification. | 1. (Construct and) test need specification. |
| 2. Overall planning. | 2. Overall planning. |
| 3. Item development. | 3. Dataset development. |
|   a. Construct refinement. |   a. Existing item collection *OR* |
|   b. Item generation. |    - Quality control. |
|   c. Item review. |   b. Item creation and/or annotation. |
|   d. Piloting of items. |    - Instructions. |
|   e. Psychometric quality analysis. |    - (Pilot) study. |
| 4. Test construction and specification. |    - Agreement analysis. |
| 5. Implementation and testing. |    - Error analysis. |
| 6. Psychometric quality analysis. | 4. Dataset construction. |
| 7. Test scoring and norming. | 5. Model selection and evaluation. |
| 8. Technical Manual. | 6. Benchmark release. |

**Table 1. Contrasting test development between psychometrics and LLM benchmarking.**

Test development in psychometrics concerns the process of developing and implementing a test according to psychometric principles [20]. Table 1 contrasts psychometric test development (based on [20]) with common LLM benchmarking procedures (based on [5,33]). What sets psychometric test development apart

from typical LLM benchmark development is its focus on ensuring that the test matches a well-defined construct via expert-driven item generation, rigorous pilot testing, use of factor analysis and IRT models for item and test diagnostics, establishment of scoring and normalisation standards, and testing on representative samples of intended test takers. The result of this elaborate process is a high-quality test that can assess the construct of interest for the test takers in a valid and reliable way. Many large-scale assessments, such as PISA (Programme for International Student Assessment), TIMSS and PIRLS (Progress in International Reading Literacy Study), conform to such a process.

We will use **P**roficiency in **G**rade **S**chool **M**athematics (PGSM) as the construct of interest to further illustrate this process. In Step 1, the construct of interest and the test need are specified. For instance, how do we define PGSM? Is it based on a specific curriculum? What does existing literature say? Which education levels are we interested in? Is the test meant for comparison between students within a school, or between schools within a country? Such questions help us to clarify what we want to measure and how it can be measured.

In Step 2, we make necessary planning: How many test items? What kind of item format (e.g., multiple choice, short answer questions)? Will the test scores be standardised? How to assess the quality of test items? What are the desired psychometric properties of the test items (e.g., how discriminative and difficult should the items be?) and the test as a whole (e.g., internal consistency)? Will we pilot any test item? Will the test be computer- or paper-based? To sample test takers, what kind of sampling frames and strategies should we use?

In Step 3, we develop test items, which is an iterative procedure involving five steps: (a) construct refinement, where we further clarify the definition of PGSM (e.g., What content domains should be included: number, algebra, and/or probability theory? Is proficiency only about knowing, or also about applying and reasoning?); (b) generate a pool of items with domain experts; (c) review the items for obvious misfit, errors and biases; (d) pilot the items with a representative sample of target test takers; (e) with the responses from the pilot step, we can assess the psychometric properties of the test items with IRT and factor analysis (e.g., item discriminativeness; item difficulty; factor structure[7]). We iterate this procedure until we have a set of test items with acceptable psychometric properties. Then, in Step 4, we construct the PGSM test by specifying, for instance, which items to include (if not all), in which order, how many equivalent test versions, and what scoring instructions to use.

In Step 5, the test gets implemented to the intended test takers, followed by Step 6: another round of quality analysis. If any item displays low quality characteristics (e.g., zero or negative discriminativeness), it will be left out of the final scoring. In Step 7, responses of the test takers are scored for each item, and the resulting item-level scores form the basis for estimating proficiency scores using IRT or simpler procedures like (weighted) sums. It is typical to also normalise the proficiency scores (e.g., with a mean of 500 and a standard

---

[7] Factor structure refers to the correlational relationships between the test items for measuring a construct of interest.

deviation of 100) to facilitate interpretations and comparisons. Finally, in Step 8, a technical manual is compiled, detailing Step 1–7 and corresponding results, to facilitate correct re-use of the response data, the test, as well as interpretation of test scores, among other purposes.

### 3.3   LLM Benchmark Development

Developing LLM benchmarks follows a similar yet different process. Take GSM8K [10] as an example. The authors started by specifying the need for a large, high quality mathematics test at grade school level and of moderate difficulty for LLMs (Step 1). The construct (i.e., PGSM) is not explicitly linked to any specific curriculum. Then, the overall planning is made (Step 2): The number of items should be in the thousands; the items will be curated by crowd workers; agreement and error analysis will be used to investigate the quality of the dataset; GPT-3 will be used to benchmark the dataset and verify dataset difficulty.

In Step 3, where dataset development[8] takes place, often one of the two strategies is used: *either* collect items from existing datasets and other sources and compile them into a new dataset, *or*, like in GSM8K, create own items from scratch (with annotations). The latter is usually an iterative procedure consisting of four parts: creating instructions (and possibly a user interface) for item generation and/or annotation; conducting a (pilot) study to collect the items and/or annotations; check annotator agreement; and assessing errors associated with the items or annotations. This step is iterated until a sufficient number of items and datasets are reached while meeting desired quality standards (e.g., high annotator agreement, low error rate). In total, GSM8K includes 8,500 items with solutions, with identified annotator disagreements resolved and a less than 2% error rate.

In Step 4, the generated items form the final dataset, typically with training, evaluation and testing partitions. In Step 5, selected LLMs are evaluated on the dataset. Finally, in Step 6, the benchmark gets released, which typically consists of the dataset as well as its documentation (often a research paper) and benchmarking results.

*Comparison with Psychometrics* While sharing similarity with test development in psychometrics, benchmark development for LLMs falls short on four aspects. First, the construct of interest is often under-specified, leading to a mismatch between the intended construct and what the dataset actually measures. Take GSM8K as an example: While the dataset is intended to measure proficiency in grade school mathematics, the target grade level(s) are unclear and it only focuses on one content domain (algebra), missing other relevant ones like geometry and data. This is likely the result of not using established mathematics curricula and domain experts to develop test items.

---

[8] Note that we use the term "dataset development" here, contrasting "item development" in psychometrics, because of LLM benchmarks' typical emphasis on large and multiple datasets rather than concrete test items.

Second, despite researchers' interest in comparing LLM performance with human test takers (e.g., the GSM8K paper claims that "a bright middle school student should be able to solve every problem"), such comparisons usually cannot be made because the test has not been designed with humans in mind or validated on any representative samples of the test's target user populations.

Third, besides agreement and error analysis, LLM benchmarks can benefit from psychometric analysis of test items, (i.e., checking item discriminativeness and difficulty, as well as the factor structure of the items). While this is not yet the norm, there have been promising attempts (see Section 2).

Lastly, the released benchmark often does not contain sufficient details about the process of benchmark creation. For instance, the GSM8K paper does not report instructions for item generation and annotation, results of the pilot study, agreement statistics, or annotator characteristics, all of which are important for external researchers to independently validate the quality of the benchmark.

## 4   PATCH: Psychometrics-AssisTed benCHmarking of LLMs

Figure 1 illustrates PATCH, our conceptualisation of a Psychometrics-AssisTed framework for benCHmarking LLMs.[9] Under PATCH, the first step is to define the construct of interest (e.g., proficiency in 8th grade mathematics). The second step is to find an existing validated psychometric test measuring this property; alternatively, a test can be developed from scratch, following the procedures described in Section 3.2, which likely requires collaboration with experienced psychometricians. The term "validated" means that the test has been tested on a representative sample of the target population of (human) test takers and fulfils psychometric quality requirements (e.g., sufficiently many discriminative items well distributed across different difficulty levels; showing high reliability (e.g., high internal consistency) and validity (e.g., the test's factor structure matches the construct definition)).

Next (Step 3→4), we use the items from the validated psychometric test to construct prompts for the LLMs under evaluation and then sample responses. A response typically consists of a task description, an explanation and an answer (key). Therefore, in Step 4→5, we extract the answer (key) for each item's response, then grade it to obtain item scores (Step 5→6).

For Step 2→7, the responses of human test takers (and of LLMs, if a sufficient number of LLMs are involved) can be used to estimate IRT item parameters and subsequently the latent proficiency scores for each test taker (human or LLM) with uncertainty estimates. Multiple IRT models are often used because of the adoption of different types of test items. These latent proficiency scores are typically standardised $z$-scores (i.e., mean of 0 and standard deviation of 1), which can optionally go through further normalisation (e.g., re-scaling to a

---

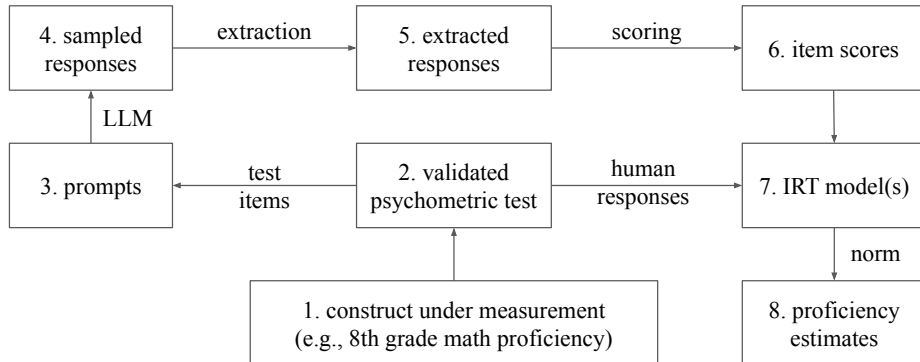[9] PATCH is partly inspired by the Hexagon Framework of scientific measurements proposed by [27].

**Fig. 1. PATCH: A *P*sychometrics-*AssisT*ed framework for ben*CH*marking LLMs.**

mean of 500 and a standard deviation of 100) (Step 6→7). These final proficiency scores enable comparison with other models and populations.

At the heart of PATCH lies a validated psychometric test, which not only provides the basis for accurate measurement of the capability of interest but also facilitates comparison between LLMs and human test takers. Unfortunately, developing such a test can be a long and expensive process; utilising existing tests can be a shortcut, which, however, should satisfy three conditions: clear human population reference; test items available; human responses and/or item parameter estimates available. The second and third are in practice difficult to meet, as many test institutes do not make their test items public due to commercial interests (e.g., SAT) or the need to measure trends over time (e.g., PISA). Collaboration with test institutes would alleviate this problem.

To the best of our knowledge, among academic proficiency tests, only TIMSS and PIRLS tests from certain years can be readily used for PATCH-based LLM benchmarking. TIMSS measures proficiency in grade school mathematics and science (4th grade, 8th grade, and final year of secondary school), while PIRLS assesses reading comprehension in 9/10-year-olds. Both TIMSS and PIRLS are administered in a large number of countries and regions with representative student samples, enabling country/region-level comparisons. In the following section, we demonstrate PATCH by measuring several LLMs' proficiency in 8th grade mathematics, using the latest available data from TIMSS 2011.

## 5    Demonstration: Measuring LLM Proficiency in 8th Grade Mathematics

### 5.1    Data: TIMSS 2011 8th Grade Mathematics

56 countries/regions participated in TIMSS 2011, with typically a random sample of about 150 schools in each country/region and a random sample of about

4,000 students from these schools. These sample sizes are determined on the basis of a $\leq .035$ standard error for each country's mean proficiency estimate. The use of random sampling makes unbiased proficiency estimates possible at the population level. TIMSS 2011 has released a publicly available database[10], of which three components are relevant to our study:

*Test Items* The TIMSS 2011 study has released 88 mathematics test items, 48 of which are multiple choice, 30 open-ended items scored as either incorrect or correct, and 10 open-ended items scored as either incorrect, partially correct, or correct. These items assess four content domains representative of 8th grade mathematics curriculum (agreed upon by experts from participating countries/regions): number, algebra, geometry, data and chance. Within each domain, items are designed to cover various subtopics (e.g., decimals, functions, patterns) and three cognitive domains: knowing, applying and reasoning. These test items are only available in a PDF file that can be downloaded from the NCES website, which includes also scoring instructions.[11] To extract them into a format compatible with LLMs, we used OCR tools to extract as much textual information as possible, converted mathematical objects (e.g., numbers, symbols, equations, tables) into LaTeX format (following earlier benchmarks like MATH) [18] and figures into JPEG format. See Appendix A.1 for examples. We have released this LLM-compatible version of test items, as well as an eighth grade science test dataset from TIMSS 2011, an advanced secondary school mathematics test dataset from TIMSS 2008, and an advanced secondary school physics test dataset from TIMSS 2008.

*IRT and Item Parameters* The dataset also specifies the IRT model used for each test item and contains the item parameter estimates (e.g., discriminativeness, difficulty), which we use to reconstruct the final IRT model for proficiency estimation.

*Student Responses and Proficiency Estimates* Lastly, responses of the sampled students to each test item and their proficiency estimates are also available, allowing us to construct proficiency score distributions for each country/region.

### 5.2   LLMs: GPT-4, Gemini-Pro and Qwen with Vision Capability

Considering that more than $1/3$ of the test items contain visual elements, we selected four competitive vision language models: GPT-4 with Vision (GPT-4V), Gemini-Pro-Vision, as well as the open-source Qwen-VL-Plus and Qwen-VL-Max [3]. There are more LLMs with vision capability. However, our goal is to showcase PATCH, not to benchmark as many LLMs as possible.

A major concern in using these LLMs is data contamination, which is difficulty to check due to inaccessible (information about) training data. However,

---

[10] https://timssandpirls.bc.edu/timss2011/international-database.html
[11] https://nces.ed.gov/timss/pdf/TIMSS2011_G8_Math.pdf

as our focus is on demonstrating the PATCH framework, data contamination is less worrying. Furthermore, data contamination is still unlikely for four reasons. First, these test items are copyrighted, forbidding commercial use. Second, the test items are hard to extract from the source PDF. Third, to the best of our knowledge, these test items do not exist in current LLM mathematics benchmarks. Fourth, we prompted the selected LLMs to explain or provide solutions to the test items' IDs (available in the source PDF). All failed to recognise these specific test IDs.

### 5.3   Prompts and Temperature

We design two separate prompts for each test item: the system message and the user message. We design the system message according to the prompt engineering guide by OpenAI, utilising chain-of-thought and step-by-step instructions on how to respond to the user message (i.e., with a classification of question type, an explanation and an answer (key)).[12] The system message is the same for all test items (see Appendix A.2). Furthermore, to account for LLMs' sensitivity to slight variations in prompts [38,26], we generate 10 additional variants of the system prompt with slight perturbations (e.g., lowercase a heading, vary the order of unordered bullet points).

The user message is item-specific, containing both the item's textual description and the associated image(s) in base 64 encoded format. See Appendix A.1 for examples.[13]

Following [31]'s technical report, we set the temperature parameter at 0.3 for multiple choice items and 0.6 for the others. See Appendix B for example responses.

### 5.4   Scoring and Proficiency Estimation

We manually scored the sampled responses from the LLMs following the official scoring rubrics of TIMSS 2011. Then, for multiple choice items, we apply the 3PL model (Equation 1); for open-ended items, we apply the GPC model (Equation 2) if partially correct response is admissible, otherwise the 2PL model. We use maximum likelihood to obtain unbiased estimates of model proficiency scores ($\theta$) with the `mirt` package in R [7]. This results in 11 $\theta$ estimates per model corresponding to 11 system message variants. We then use inverse variance weighting [28] to combine these estimates. Inverse variance weighting gives more weight to estimates that are more precise (i.e., having lower variance) and less weight to those that are less precise (i.e., having higher variance). This way, we obtain a more accurate *overall* $\theta$ estimate and its 95% confidence interval (CI) for each model.

---

[12] https://platform.openai.com/docs/guides/prompt-engineering

[13] We are aware of other prompt engineering techniques like few-shot prompting and self-consistency. We did not experiment with them, as our focus is on demonstrating PATCH.
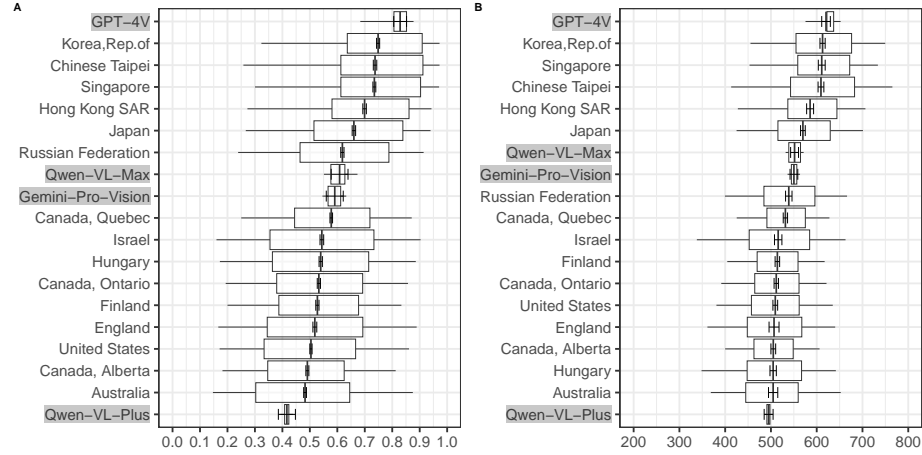
## 5.5   Results



**Fig. 2. Distribution of proficiency estimates for GPT-4V, Gemini-Vision-Pro, Qwen-VL-Plus, Qwen-VL-Max and selected participating countries/regions of the TIMSS 2011 8th grade mathematics test.** Left figure (A) shows the proficiency estimates based on the percentages of correct responses. Right figure (B) shows the IRT-based proficiency estimates. The middle vertical line in each box plot represents the weighted mean proficiency score, with the error bars indicating its 95% confidence interval. The borders of each box indicate the range of the middle 50% of all values, with the two whiskers indicating the 5th and 95th percentiles.

Figure 2 shows the proficiency score distribution and ranking of 15 selected participating countries and regions, GPT-4V and Gemini-Pro-Vision. Only 15 countries are shown here to save space. The complete figures can be found in Appendix C. The proficiency scores (x-axis) on the left panel are percentages of correct responses, which is the default approach in current LLM benchmarking; the proficiency estimates on the right panel are based on IRT. We make three observations. First, regardless of the method of proficiency estimation, GPT-4V has the overall best performance relative to Gemini-Pro-Vision and the average proficiency of 8th grade students of each participating country/region. Second, the method of proficiency estimation affects the ranking results. For instance, while Chinese Taipei is ranked 3rd on the left, it is ranked 4th on the right; Gemini-Pro-Vision is ranked 8th on the left, but ranked 7th on the right. Similarly, while Hungary is ranked 11th on the left, it drops to the 16th place on the right. Third, the method of proficiency estimation affects the estimated 95% CIs, which are usually wider when IRT is used (as it accounts for both item and test taker variances). Notably, while on the left panel the CI of GPT-4V does not overlap with the second best, South Korea, indicating a statistically significant difference, they overlap on the right panel, suggesting otherwise. This finding

shows that the adoption of PATCH is likely going to make a difference to LLM benchmark results.

## 6  Conclusion

In this paper, we propose PATCH, a psychometrics-inspired framework to address current limitations of LLM benchmarks, including questionable measurement quality, lack of quality assessment on the item level and unwarranted comparison between humans and LLMs. We demonstrate PATCH with an 8th grade mathematics proficiency test, where PATCH yields evaluation outcomes that diverge from those based on existing benchmarking practices. This underscores the potential of PATCH to reshape the LLM benchmarking landscape.

## 7  Limitations

Our paper has the following limitations, among others. First, PATCH requires validated tests, which can be resource-intensive if tests need to be developed from scratch. However, this also opens up opportunities for collaboration between LLM researchers, psychometricians and test institutes. Second, the validity, reliability, and fairness of using tests validated solely on humans for LLM benchmarking are debatable due to possibly differing notions of proficiency and cognitive processes between LLMs and humans. Nonetheless, such tests are still better than non-validated benchmarks, particularly for comparison of model and human performance. Advancing LLM benchmarking further requires tests validated on LLMs (and humans for model-human comparisons), necessitating theoretical work on LLM-specific constructs and the development of LLM-specific IRT models and testing procedures. Third, our experiment only includes two proprietary LLMs and one proficiency test. We consider this sufficient for demonstrating PATCH, but not enough if the goal is to benchmark as many LLMs as possible across different tests.

**Disclosure of Interests.** The authors have no competing interests to declare that are relevant to the content of this article.

## References

1. AERA, APA, NCME: The Standards for Educational and Psychological Testing. American Educational Research Association (2014)

2. Amidei, J., Piwek, P., Willis, A.: Identifying annotator bias: A new IRT-based method for bias identification. In: Proceedings of the 28th International Conference on Computational Linguistics. pp. 4787–4797 (2020). https://doi.org/10.18653/V1/2020.COLING-MAIN.421

3. Bai, J., Bai, S., Yang, S., Wang, S., Tan, S., Wang, P., Lin, J., Zhou, C., Zhou, J.: Qwen-vl: A frontier large vision-language model with versatile abilities. ArXiv **abs/2308.12966** (2023)

4. Birkenmaier, L., Lechner, C., Wagner, C.: ValiTex - A uniform validation framework for computational text-based measures of social science constructs. arXiv (2023), https://arxiv.org/abs/2307.02863

5. Bowman, S.R., Angeli, G., Potts, C., Manning, C.D.: A large annotated corpus for learning natural language inference. In: Màrquez, L., Callison-Burch, C., Su, J. (eds.) Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing. pp. 632–642. Association for Computational Linguistics, Lisbon, Portugal (2015). https://doi.org/10.18653/v1/D15-1075

6. Camilli, G.: Teacher's corner: origin of the scaling constant d= 1.7 in item response theory. Journal of Educational Statistics **19**(3), 293–295 (1994). https://doi.org/10.3102/10769986019003293

7. Chalmers, R.P.: mirt: A multidimensional item response theory package for the r environment. Journal of Statistical Software **48**, 1–29 (2012). https://doi.org/10.18637/jss.v048.i06

8. Chen, M., Tworek, J., Jun, H., Yuan, Q., Ponde, H., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., Ray, A., Puri, R., Krueger, G., Petrov, M., Khlaaf, H., Sastry, G., Mishkin, P., Chan, B., Gray, S., Ryder, N., Pavlov, M., Power, A., Kaiser, L., Bavarian, M., Winter, C., Tillet, P., Such, F.P., Cummings, D.W., Plappert, M., Chantzis, F., Barnes, E., Herbert-Voss, A., Guss, W.H., Nichol, A., Babuschkin, I., Balaji, S., Jain, S., Carr, A., Leike, J., Achiam, J., Misra, V., Morikawa, E., Radford, A., Knight, M.M., Brundage, M., Murati, M., Mayer, K., Welinder, P., McGrew, B., Amodei, D., McCandlish, S., Sutskever, I., Zaremba, W.: Evaluating large language models trained on code. arXiv (2021), https://arxiv.org/abs/2107.03374

9. Clark, P., Cowhey, I., Etzioni, O., Khot, T., Sabharwal, A., Schoenick, C., Tafjord, O.: Think you have solved question answering? Try ARC, the AI2 Reasoning Challenge. arXiv (2018), https://arxiv.org/abs/1803.05457

10. Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., Hesse, C., Schulman, J.: Training verifiers to solve math word problems. arXiv (2021), https://arxiv.org/abs/2110.14168

11. Dillion, D., Tandon, N., Gu, Y., Gray, K.: Can AI language models replace human participants? Trends in Cognitive Sciences (2023). https://doi.org/10.1016/j.tics.2023.04.008

12. Du, Y., Fang, Q., Nguyen, D.: Assessing the reliability of word embedding gender bias measures. In: Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing. pp. 10012–10034. Association for Computational Linguistics, Online and Punta Cana, Dominican Republic (2021). https://doi.org/10.18653/v1/2021.emnlp-main.785

13. Dua, D., Wang, Y., Dasigi, P., Stanovsky, G., Singh, S., Gardner, M.: DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs. In: Burstein, J., Doran, C., Solorio, T. (eds.) Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers). pp. 2368–

2378. Association for Computational Linguistics, Minneapolis, Minnesota (2019). https://doi.org/10.18653/v1/N19-1246

14. Fang, Q., Giachanou, A., Bagheri, A., Boeschoten, L., van Kesteren, E.J., Kamalabad, M.S., Oberski, D.: On text-based personality computing: Challenges and future directions. In: Findings of the Association for Computational Linguistics: ACL 2023. pp. 10861–10879 (2023). https://doi.org/10.18653/v1/2023.findings-acl.691

15. Fang, Q., Nguyen, D., Oberski, D.L.: Evaluating the construct validity of text embeddings with application to survey questions. EPJ Data Science **11**(1), 1–31 (Dec 2022). https://doi.org/10.1140/epjds/s13688-022-00353-7

16. Fang, Q., Zhou, Z., Barbieri, F., Liu, Y., Neves, L., Nguyen, D., Oberski, D.L., Bos, M.W., Dotsch, R.: Designing and evaluating general-purpose user representations based on behavioural logs from a measurement process perspective: A case study with snapchat. arXiv (2023), https://arxiv.org/abs/2312.12111

17. Google, G.T.: Gemini: a family of highly capable multimodal models. arXiv (2023), https://arxiv.org/abs/2312.11805

18. Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., Steinhardt, J.: Measuring massive multitask language understanding. In: International Conference on Learning Representations (2021), https://openreview.net/forum?id=d7KBjmI3GmQ

19. Huang, J., Wang, W., Li, E.J., LAM, M.H., Ren, S., Yuan, Y., Jiao, W., Tu, Z., Lyu, M.: On the humanity of conversational AI: Evaluating the psychological portrayal of LLMs. In: The Twelfth International Conference on Learning Representations (2024), https://openreview.net/forum?id=H3UayAQWoE

20. Irwing, P., Hughes, D.J.: Test development. In: The Wiley Handbook of Psychometric Testing, pp. 1–47. John Wiley & Sons, Ltd (2018). https://doi.org/10.1002/9781118489772.ch1

21. Kreuter, A., Sassenberg, K., Klinger, R.: Items from psychometric tests as training data for personality profiling models of twitter users. In: Proceedings of the 12th Workshop on Computational Approaches to Subjectivity, Sentiment & Social Media Analysis. pp. 315–323. Association for Computational Linguistics (2022). https://doi.org/10.18653/v1/2022.wassa-1.35

22. Lalor, J.P., Wu, H., Munkhdalai, T., Yu, H.: Understanding deep learning performance through an examination of test set difficulty: A psychometric case study. In: Riloff, E., Chiang, D., Hockenmaier, J., Tsujii, J. (eds.) Proceedings of the Conference on Empirical Methods in Natural Language Processing. Conference on Empirical Methods in Natural Language Processing. pp. 4711–4716. Association for Computational Linguistics (2018). https://doi.org/10.18653/v1/D18-1500

23. Lalor, J.P., Wu, H., Yu, H.: Building an evaluation scale using item response theory. In: Su, J., Duh, K., Carreras, X. (eds.) Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing. pp. 648–657. Association for Computational Linguistics, Austin, Texas (2016). https://doi.org/10.18653/v1/D16-1062

24. Lalor, J.P., Wu, H., Yu, H.: Learning latent parameters without human response patterns: Item response theory with artificial crowds. In: Inui, K., Jiang, J., Ng, V., Wan, X. (eds.) Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP). pp. 4249–4259. Association for Computational Linguistics, Hong Kong, China (2019). https://doi.org/10.18653/v1/D19-1434

25. Lalor, J.P., Yu, H.: Dynamic data selection for curriculum learning via ability estimation. In: Cohn, T., He, Y., Liu, Y. (eds.) Findings of the Association for Computational Linguistics: EMNLP 2020. pp. 545–555. Association for Computational Linguistics, Online (2020). https://doi.org/10.18653/v1/2020.findings-emnlp.48

26. Loya, M., Sinha, D., Futrell, R.: Exploring the sensitivity of LLMs' decision-making capabilities: Insights from prompt variations and hyperparameters. In: Findings of the Association for Computational Linguistics: EMNLP 2023. pp. 3711–3716 (2023). https://doi.org/10.18653/v1/2023.findings-emnlp.241

27. Mari, L., Wilson, M., Maul, A.: Measurement across the sciences: Developing a shared concept system for measurement. Springer Nature (2023). https://doi.org/10.1007/978-3-031-22448-5

28. Marín-Martínez, F., Sánchez-Meca, J.: Weighting by inverse variance or by sample size in random-effects meta-analysis. Educational and Psychological Measurement **70**(1), 56–73 (2010). https://doi.org/10.1177/0013164409344534

29. Muraki, E.: A generalised partial credit model: Application of an EM algorithm. Applied Psychological Measurement **16**(2), 159–176 (1992). https://doi.org/10.1002/j.2333-8504.1992.tb01436.x

30. Nie, Y., Zhou, X., Bansal, M.: What can we learn from collective human opinions on natural language inference data? In: Webber, B., Cohn, T., He, Y., Liu, Y. (eds.) Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP). pp. 9131–9143. Association for Computational Linguistics, Online (2020). https://doi.org/10.18653/v1/2020.emnlp-main.734

31. OpenAI: GPT-4 technical report. arXiv (2023), https://arxiv.org/abs/2303.08774

32. Pellert, M., Lechner, C.M., Wagner, C., Rammstedt, B., Strohmaier, M.: AI Psychometrics: Assessing the psychological profiles of large language models through psychometric inventories. Perspectives on Psychological Science (2023). https://doi.org/10.1177/17456916231214460

33. Raji, D., Denton, E., Bender, E.M., Hanna, A., Paullada, A.: AI and the everything in the Whole Wide World Benchmark. In: Vanschoren, J., Yeung, S. (eds.) Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks. vol. 1 (2021), https://datasets-benchmarks-proceedings.neurips.cc/paper_files/paper/2021/file/084b6fbb10729ed4da8c3d3f5a3ae7c9-Paper-round2.pdf

34. Reise, S.P., Revicki, D.A.: Handbook of item response theory modeling. Taylor & Francis New York, NY (2014)

35. Rodriguez, P., Barrow, J., Hoyle, A.M., Lalor, J.P., Jia, R., Boyd-Graber, J.: Evaluation examples are not equally informative: How should that change NLP leaderboards? In: Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers). pp. 4486–4503 (2021). https://doi.org/10.18653/v1/2021.acl-long.346

36. Rodriguez, P., Htut, P.M., Lalor, J.P., Sedoc, J.: Clustering examples in multi-dataset benchmarks with item response theory. In: Proceedings of the Third Workshop on Insights from Negative Results in NLP. pp. 100–112 (2022). https://doi.org/10.18653/v1/2022.insights-1.14

37. Sakaguchi, K., Bras, R.L., Bhagavatula, C., Choi, Y.: WinoGrande: An adversarial winograd schema challenge at scale. Commun. ACM **64**(9), 99–106 (2021). https://doi.org/10.1145/3474381

38. Sclar, M., Choi, Y., Tsvetkov, Y., Suhr, A.: Quantifying language models' sensitivity to spurious features in prompt design or: How I learned to start worrying about prompt formatting. In: The Twelfth International Conference on Learning Representations (2024), https://openreview.net/forum?id=RIu5lyNXjT

39. Sedoc, J., Ungar, L.: Item response theory for efficient human evaluation of chatbots. In: Proceedings of the First Workshop on Evaluation and Comparison of NLP Systems. pp. 21–33 (2020). https://doi.org/10.18653/v1/2020.eval4nlp-1.3

40. Sen, I., Flöck, F., Wagner, C.: On the reliability and validity of detecting approval of political actors in tweets. In: Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP). pp. 1413–1426. Association for Computational Linguistics, Online (2020). https://doi.org/10.18653/v1/2020.emnlp-main.110

41. Tedeschi, S., Bos, J., Declerck, T., Hajič, J., Hershcovich, D., Hovy, E., Koller, A., Krek, S., Schockaert, S., Sennrich, R., Shutova, E., Navigli, R.: What's the meaning of superhuman performance in today's NLU? In: Rogers, A., Boyd-Graber, J., Okazaki, N. (eds.) Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). pp. 12471–12491. Association for Computational Linguistics, Toronto, Canada (2023). https://doi.org/10.18653/v1/2023.acl-long.697

42. Vania, C., Htut, P.M., Huang, W., Mungra, D., Pang, R.Y., Phang, J., Liu, H., Cho, K., Bowman, S.: Comparing test sets with item response theory. In: Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers). pp. 1141–1158 (2021). https://doi.org/10.18653/v1/2021.acl-long.92

43. Vu, H., Abdurahman, S., Bhatia, S., Ungar, L.: Predicting responses to psychological questionnaires from participants' social media posts and question text embeddings. In: Findings of the Association for Computational Linguistics: EMNLP 2020. pp. 1512–1524. Association for Computational Linguistics, Online (2020). https://doi.org/10.18653/v1/2020.findings-emnlp.137

44. van der Wal, O., Bachmann, D., Leidinger, A., van Maanen, L., Zuidema, W., Schulz, K.: Undesirable biases in NLP: Addressing challenges of measurement. Journal of Artificial Intelligence Research **79**, 1–40 (2024). https://doi.org/10.1613/jair.1.15195

45. Yang, F., Yang, T., Quan, X., Su, Q.: Learning to answer psychological questionnaire for personality detection. In: Findings of the Association for Computational Linguistics: EMNLP 2021. pp. 1131–1142. Association for Computational Linguistics (2021). https://doi.org/10.18653/v1/2021.findings-emnlp.98

46. Zellers, R., Holtzman, A., Bisk, Y., Farhadi, A., Choi, Y.: HellaSwag: Can a machine really finish your sentence? In: Korhonen, A., Traum, D., Màrquez, L. (eds.) Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics. pp. 4791–4800. Association for Computational Linguistics, Florence, Italy (2019). https://doi.org/10.18653/v1/P19-1472

## Appendices

Due to page limit, please find the full appendices on https://github.com/fqixiang/patch_llm_benchmarking_with_psychometrics/blob/main/bnaic_paper.pdf.