# Visual Reasoning Benchmark: Evaluating Multimodal LLMs on Classroom-Authentic Visual Problems from Primary Education

**Mohamed Huti**[1]    **Alasdair Mackintosh**[1]    **Amy Waldock**[1]    **Dominic Andrews**[1]

**Maxime Lelièvre**[1]    **Moritz Boos**[1]    **Tobias Murray**[1]

**Paul Atherton**[1]    **Robin A. A. Ince**[1]    **Oliver G. B. Garrod**[1]

[1]Fab AI

## Abstract

AI models have achieved state-of-the-art results in textual reasoning; however, their ability to reason over spatial and relational structures remains a critical bottleneck—particularly in early-grade maths, which relies heavily on visuals. This paper introduces the visual reasoning benchmark (VRB), a novel dataset designed to evaluate Multimodal Large Language Models (MLLMs) on their ability to solve authentic visual problems from classrooms. This benchmark is built on a set of 701 questions sourced from primary school examinations in Zambia and India, which cover a range of tasks such as reasoning by analogy, pattern completion, and spatial matching. We outline the methodology and development of the benchmark which intentionally uses unedited, minimal-text images to test if models can meet realistic needs of primary education. Our findings reveal a "jagged frontier" of capability where models demonstrate better proficiency in static skills such as counting and scaling, but reach a distinct "spatial ceiling" when faced with dynamic operations like folding, reflection, and rotation. These weaknesses pose a risk for classroom use on visual reasoning problems, with the potential for incorrect marking, false scaffolding, and reinforcing student misconceptions. Consequently, education-focused benchmarks like the VRB are essential for determining the functional boundaries of multimodal tools used in classrooms.

## 1   Introduction

AI models have achieved state-of-the-art progress on a wide range of reasoning tasks, from natural language understanding and knowledge retrieval to symbolic mathematics. Benchmarks such as GSM8K [Cobbe et al., 2021], MATH [Hendrycks et al., 2021b] and MMLU [Hendrycks et al., 2021a] now demonstrate above-human performance when reasoning is presented in purely textual form. Yet these advances do not extend to a core dimension of intelligence: the ability to reason visually through patterns, diagrams, and spatial relationships [Wang et al., 2024]. Classic measures such as Raven's Progressive Matrices [Raven, 1936] intentionally minimise language to assess abstract reasoning beyond vocabulary or cultural familiarity, underscoring the importance of non-verbal inference in cognition and assessments. In educational practice, such assessments have been used in school entry decisions, in identifying gifted learners and in the equitable evaluation of children with dyslexia or those learning in a language other than their mother tongue.

Visual reasoning is particularly important in learning foundational mathematics. Effective mathematical problem solving commonly involves constructing and manipulating appropriate visual

representations (e.g., diagrams, number lines, arrays) to support inference [Purcar et al., 2024]. Evidence from early primary classrooms indicates that providing pupils with systematic opportunities to develop such representations enhances arithmetic problem-solving performance [Purcar et al., 2024], and that instruction targeting spatial skills contributes to more advanced computational reasoning [Parkinson and Cutts, 2025]. Maths curricula usually utilise concrete-pictorial-abstract learning pathways, where various visual representations are used to support learning by providing connections between children's experience and abstract conceptions. Consequently, to be effective in supporting student learning in the early years of school, AI models need to be capable of understanding visual reasoning problems well. For example, an AI tutor must be able to recognize student errors in classroom visual problems and provide correct explanations, and for accurate automatic grading, the model needs high accuracy in problem recognition and answering.

Despite rapid progress, contemporary Multimodal Large Language Models (MLLMs) often falter when the decisive information is visual. Recent studies across multimodal and abstract reasoning benchmarks reveal a persistent performance gap between textual and visual modalities, suggesting that models can leverage linguistic cues effectively but still struggle to perceive and reason over the spatial and relational structures encoded in images [Xu et al., 2025].

These findings raise a practical question for education: To what extent do current models possess the kind of visual reasoning capabilities needed to be useful in classrooms? Rather than measuring abstract reasoning ability for complex problems that many adults would find challenging, we focus on whether models can handle the kind of visual problems that primary students would encounter, thus establishing a minimum capability threshold for models to support learners and teachers effectively in visual problem-solving contexts.

The Visual Reasoning Benchmark (VRB) introduced in this paper addresses this question by emphasising minimal-text, classroom-authentic visual problems. In education, visual reasoning (also called non-verbal reasoning) is defined as "using pictures, images or diagrams effectively for solving tasks of higher-order thinking" [Natsheh and Karsenty, 2014]. Here, we use visual reasoning to mean *deriving the correct answer by perceiving and operating on spatial and relational structure present in an image*. The VRB design complements multimodal datasets such as MathVista [Xie et al., 2023] and Math-V [Lu et al., 2024] which include visual problems but retain textual statements and choices, while aligning more closely with VisuLogic [Xu et al., 2025] which deliberately targets non-verbal reasoning. Our contribution is to bring this challenge into authentic educational contexts. By drawing directly on primary-level assessment tasks from Zambia and India, VRB provides the first large-scale, classroom-grounded benchmark for evaluating visual reasoning in low- and middle-income country (LMIC) settings. It uses unedited questions, including slight issues from photocopying or production common in an LMIC setting, revealing how a model might handle genuine student questions (for example, those presented to a chatbot). This is an essential step to determining when such systems could become genuinely useful for learning support. We estimate the minimum capability threshold for classroom usefulness at 94%, the proportion of questions on which three adult annotators independently agreed on the answer (see Section 3.3). In this way, the benchmark is designed to test if models meet grade-level expectations in visual reasoning and to indicate their suitability for use by teachers and students in lesson support and design.

## 2 Related Work

Research on visual reasoning has developed along distinct trajectories, each addressing different aspects of the challenge.

### 2.1 Synthetic Abstract Reasoning

The Abstraction and Reasoning Corpus (ARC) [Chollet, 2019] represents the purest test of abstract visual reasoning, presenting grid-based challenges where participants must uncover hidden transformation rules such as symmetry or repetition. Humans achieve 98–100% accuracy, but frontier models reach 85–95% (ARC-AGI-1) and around 70% (ARC-AGI-2) [Chollet et al., 2025]. MARVEL [Jiang et al., 2024] decomposes this challenge into perception and abstraction components, by testing six core patterns such as 3D-geometry and temporal movement across different visual problems. It reveals that despite high symbolic logic scores, MLLMs often perform near chance specifically on non-verbal reasoning tasks because they struggle with a persistent perceptual bottleneck.

## 2.2 Academic Mathematical Reasoning

A substantial body of work targets mathematical problems based on the interpretation of diagrams, charts, or figures. MathVista [Xie et al., 2023] compiles 6,141 problems from 28 multimodal datasets spanning geometry, algebra, and scientific plots. While frontier models at the time of launch like Gemini 2.0 Flash achieved 73.1%—surpassing the 60.3% human baseline—this success occurs primarily in settings with high textual scaffolding where linguistic cues provide a redundant path to the solution. MATH-V [Lu et al., 2024] extends this approach with 3,040 diagram-augmented contest problems across 16 topics and five difficulty levels, explicitly ensuring visual elements are non-redundant and integral to problem-solving. MathVerse [Zhang et al., 2024] provides the most systematic analysis by using 15,000 questions which include text and diagrams, where some answers depend more on information in the text and others on the diagrams. Their findings reveal a consistent pattern: accuracy declines for questions where the critical information lies in the diagrams. GPT-4V performance drops from 54.7% (Text-Dominant) to 31.6% (Vision-Only), and several models actually improve when images are removed entirely. MV-MATH [Wang et al., 2025] raises the bar with 2,009 problems requiring reasoning across 2–5 images per question, where the best models achieve 34% compared to 75–80% for humans. Across these datasets, multimodal models exploit textual scaffolding effectively, but their diagram-grounded reasoning remains brittle.

## 2.3 Research-Focused Visual Reasoning

Recent work has moved toward minimising textual scaffolding to isolate genuine visual reasoning capabilities. VisuLogic [Xu et al., 2025] introduces 1,000 fully non-verbal, human-verified puzzles across six categories, with humans scoring 51.4% while leading MLLMs remain below 30%. MM-IQ [Cai et al., 2025] extends this approach to a larger scale (2,700 test items) across eight reasoning paradigms, yielding similar outcomes (33% for models vs 51% for humans). BabyVision [Chen et al., 2026] targets an even more fundamental level, testing core visual abilities that humans acquire before language. Across 388 items spanning four categories, the best-performing MLLM scores 49.7%, falling below the level of six-year-old children and far short of the adult baseline (94.1%). VERIFY [Liu et al., 2025] emphasises fidelity by pairing 600 diagrammatic items with human-annotated reasoning steps. Top models achieve 21.7% accuracy and frequently generate correct answers with flawed reasoning. ME2 [Park et al., 2025] shifts focus from accuracy to explanatory fidelity, requiring solution steps that reference diagram annotations; even strong models struggle to align explanations with visuals. These efforts mark a shift from outcome-only scoring toward process-aware evaluation, revealing that models frequently show incorrect reasoning even when answers are correct.

## 2.4 The Educational Authenticity Gap

Collectively, these benchmarks reveal a consistent pattern: models that excel at text-based reasoning often collapse when equivalent reasoning must be conducted visually. This reflects dual challenges of perception (detecting and segmenting visual elements) and abstraction (inferring relational rules and generalizing across contexts). While contemporary vision encoders provide strong perceptual grounding, abstraction remains elusive.

In educational practice, another critical gap emerges in this landscape: **none of these core visual reasoning benchmarks are grounded in authentic primary-education assessment contexts**. Synthetic puzzles like ARC or MM-IQ probe abstraction in isolation but lack pedagogical relevance. Contest-style datasets like MATH-V target advanced mathematical reasoning but represent specialised academic contexts rather than foundational learning. Research-focused benchmarks often yield low-to-moderate human performance, suggesting a disconnection from natural learning tasks and classroom relevance.

Our benchmark addresses this gap by drawing directly from primary-level educational assessments used in real classrooms. We focus on visual reasoning tasks that millions of children successfully navigate as part of their learning journeys. This approach not only provides a real-world test of reasoning capabilities but also establishes clear expectations based on age-appropriate educational standards. In doing so, we create the first large-scale evaluation that is simultaneously rigorous for AI systems and directly relevant to educational contexts where such reasoning is integral to effective learning.

# 3 Methodology

## 3.1 Source of Questions

We extracted multiple-choice questions (MCQs) from Zambia's National End of Primary Exams Special Paper 2 (Non-Verbal Reasoning) and the Jawahar Navodaya Vidyalaya's Selection Test Class 6 (JNVST Class 6) from India.

Special Paper 2 from Zambia is administered as part of Grade 7's End of Primary Leaving Exams. Together with Special Paper 1 (verbal reasoning), it serves as an aptitude test. The aim of the test is to assess cognitive skills beyond curriculum subject areas and to provide information to support secondary school enrolment. We sourced questions from Special Paper 2 Exam Papers administered in 2018, 2019, 2021, and 2022.

JNVST is the entrance exam for the Jawahar Navodaya Vidyalayas (JNV) schools, a network of co-educational, residential schools fully financed and administered by the Government of India. It is a national-level examination based on exams designed by the Central Board of Secondary Education (CBSE) and conducted independently in each state by the Navodaya Vidyalaya Samiti, an autonomous organisation[1]. We sourced questions from JNVST administered in 2014–2020, 2022, and 2024.

## 3.2 Question Processing

**Question Extraction and Pre-processing**

The MCQ questions were extracted from the exam PDFs. Items were extracted at the question level by manual region-of-interest cropping (figures + options) using a PDF viewer (Preview) with occasional readjustment. Crops were exported as black-and-white JPEGs using default viewer settings. Note that all the original source questions are also black-and-white prints. To reflect how a model might be used in practice by students and teachers, we applied no normalisation or sharpening. We did not remove on-page text (e.g. option letters) if they were a part of the cropped region. Each MCQ received a deterministic Question ID keyed to its source. A lightweight index links Question ID to image path and the source's question. Models were evaluated in an image-first, minimal-text method—mirroring recent research that minimises linguistic scaffolding to isolate vision-centric reasoning. Prompts standardise response format only; they never include hints, examples, or domain content. Some of the exams also had linked answers, which were provided in separate documents. These were parsed and matched to the extracted questions, although additional verification took place as described in the next section. To prevent benchmark leakage, we release all evaluation prompts and code needed to reproduce our pipeline, but we do not publicly distribute the question images or answer keys. We will work directly with model developers to facilitate evaluation of new models against the benchmark.

## 3.3 Human Review, Marking and Annotation

To ensure the quality of the MCQs and for selection and annotation, each MCQ was independently checked, answered, and categorised by three human markers (two education experts and one non-expert). The markers' answers were used to provide the verified answer which models could be scored against. Where markers felt there was no correct answer, they recorded a 'refuse' rather than a random selection. In order to be retained in the dataset, we selected those MCQs with an answer with a majority consensus between the three human markers. 34 questions were excluded as there was no majority consensus answer or 2 or more markers gave 'refuse'. A further 10 questions were removed from the benchmark due to fatal errors (9 with corrupted or incorrect images and 1 with a question presentation error).

In total, 658 questions were retained with full consensus answers across all markers, and an additional 43 MCQs where 2 out of 3 markers agreed. We chose to retain these questions with a 2 out of 3 answer consensus, because there was still a clear answer, but these were likely to be more challenging as one human had failed to answer correctly. Our aim was to ensure the benchmark would test even the most capable models rigorously. This gave a final dataset containing 701 questions.

To support analysis of this final dataset, we annotated each question by 'task', and by 'skill'. The task considers how the question was structured (e.g. odd one out, pattern completion, matching). The
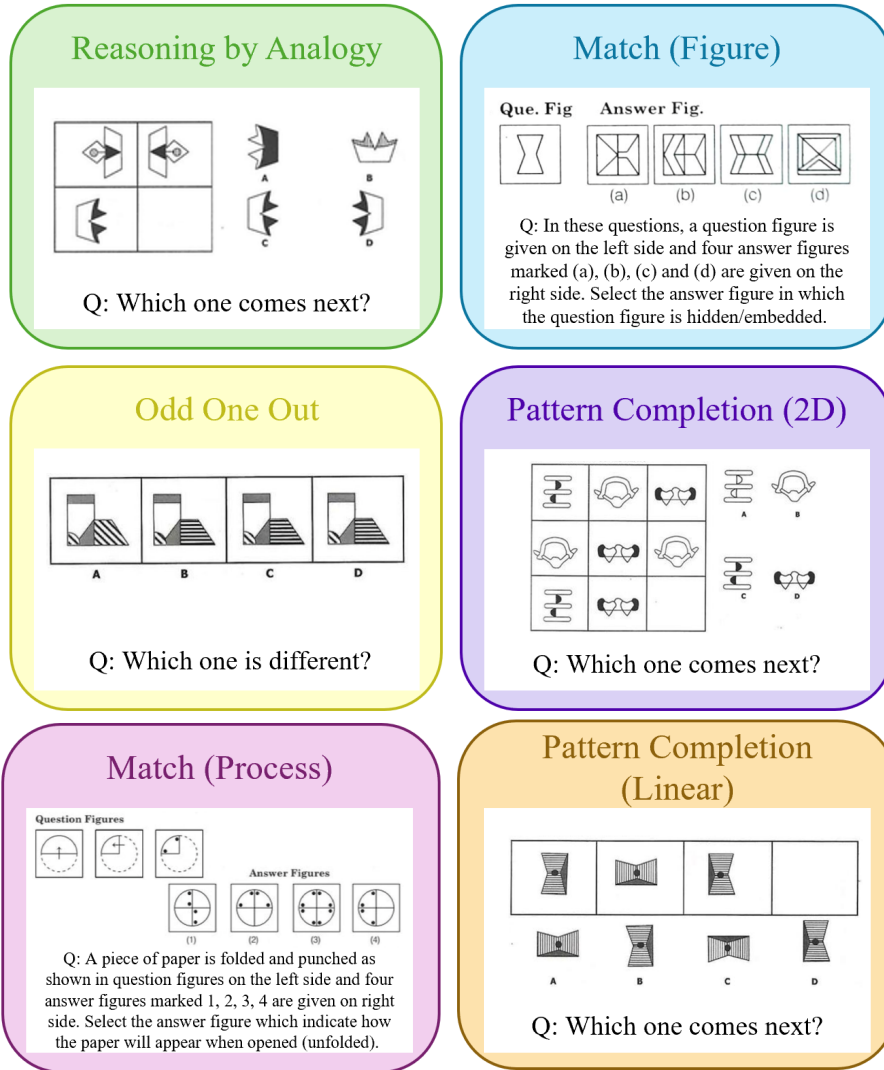
---

[1]https://examcart.in/community/blogs/jnvst-class-6-exam-date-syllabus-eligibility-pattern-pyqs

skill considers the underlying visual operation(s) (e.g. folding, counting, rotating) needed to solve the questions. To identify task and skill category names, a review was done of sources of visual reasoning questions. The categories we selected were mostly derived from guidance material used to support the UK '11+' Secondary Entrance Exams [Broadbent et al., 2018]. Whilst the '11+' does not have official question categories, this guidance material (for a well-recognised exam taken at the same age as our question set) offered a useful starting point. Where the existing categories did not match our question set well, we adapted. For skills, we added Scaling and also combined two others to create Reconfiguring Shapes and Layering. For tasks, we added Match (process) and Match (figure). In total, we used 6 tasks and 10 skills. See Figures 1, 2 and 3 for question samples and Appendix A for definitions of categories and tags.

For tasks, only one label was assigned. All task labels had 2 out of 3 or full marker agreement, with the majority label assigned used in the final annotations. For skills, flexibility was given to assign more than one skill tag because some questions involved observation of two or more metrics, for example, a shape being both rotated and its shading changing. Similarly, the final skill tags were assigned where there was majority agreement between markers (i.e. if 2 out of 3 markers or more assigned a tag, it was used in the final annotations).

Where appropriate and necessary, a pedagogy expert made minor edits to MCQs. Changes focused on removing irrelevant information, such as instructions for how to fill in answer booklets. However, MCQs were intentionally left as close to the original as possible to determine how models performed with genuine student questions.

Finally, to evaluate the robustness of multimodal models, we also annotated the quality of the question image. We categorised items based on their potential impact on student comprehension: No Error, Minor (Slight artefacts), and Moderate (clear visual noise that requires more cognitive effort to interpret). (See Appendix B for full definitions).

## Reasoning by Analogy

Q: Which one comes next?

## Match (Figure)

Q: In these questions, a question figure is given on the left side and four answer figures marked (a), (b), (c) and (d) are given on the right side. Select the answer figure in which the question figure is hidden/embedded.

## Odd One Out

Q: Which one is different?

## Pattern Completion (2D)

Q: Which one comes next?

## Match (Process)

Q: A piece of paper is folded and punched as shown in question figures on the left side and four answer figures marked 1, 2, 3, 4 are given on right side. Select the answer figure which indicate how the paper will appear when opened (unfolded).

## Pattern Completion (Linear)

Q: Which one comes next?

Figure 1: Question examples showing the six task categories in our Visual Reasoning Benchmark.

## Shape and Form

**Que. Fig**   **Answer Fig.**

(a) (b) (c) (d)

Q: A question figure is given on the left side and four answer figures marked (a), (b), (c) and (d) are given on the right side. Select the answer figure which is exactly the same as the question figure.

## Folding

Q: Which one comes next?

## Rotation

Q: Which one comes next?

## Reconfiguring Shapes

**Que. Fig**   **Answer Fig.**

(a) (b) (c) (d)

Q: A question figure is given on the left side and four answer figures marked (a), (b), (c) and (d) are given on the right side. Select the answer figure which can be formed from the cut-out pieces given in the question figure.

## Reflection

Q: Which one comes next?

## Shading and Line Type

Q: Which one comes next?

## Scaling

Q: Which one comes next?

## Counting

Q: Which one comes next?

## Layering

Q: Which one comes next?

## Other

Problem Figures

Answer Figures

(a) (b) (c) (d)

Q: Which one comes next?

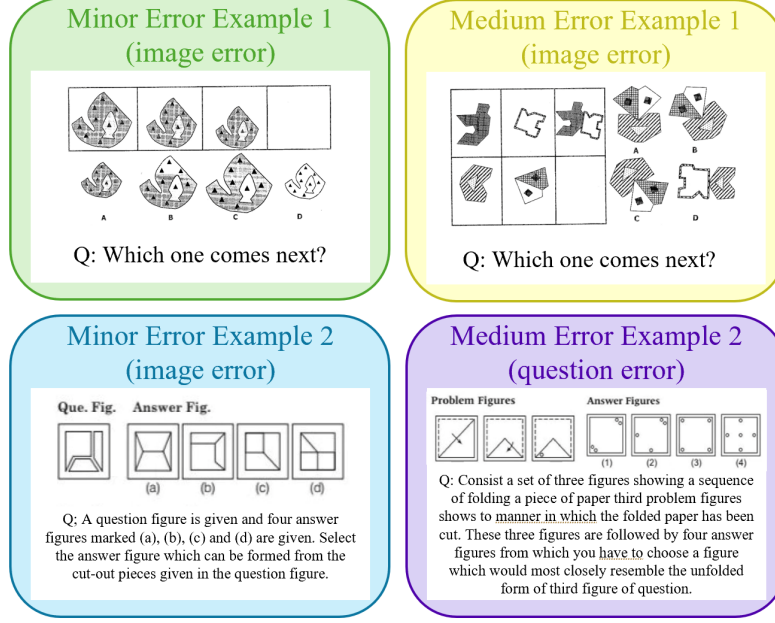Figure 2: Question examples showing the ten skill tags in our Visual Reasoning Benchmark.

Figure 3: Question examples showing the two types of error in our Visual Reasoning Benchmark.

# 4 Results

## 4.1 Visual Reasoning Benchmark (VRB) Performance

At the time of writing, 45 multimodal models have been evaluated covering both proprietary and open weighted models from major providers and spanning a wide range of sizes and price points. All models answer 4-option multiple-choice questions using a single image, minimal-text prompt. We report item-level accuracy as the primary metric, together with 95% bootstrap confidence intervals over 1000 resamples (Figure 4). Chance performance on VRB is 25%, so scores near this level are equivalent to random guessing over the options.



Figure 4: Accuracy on VRB for a subset of models. Errors show 95% bootstrap confidence intervals.

Across the full panel, accuracies range from 23% up to 78%. Gemini-3.0 Flash and Gemini-3.0 Pro attain the highest scores at 78% and 76% respectively. However, given that these questions are targeted at end-of-primary students and do not require advanced subject knowledge, this means that even the currently highest scores are relatively weak, and suggests the benchmark is currently far from saturation. This is echoed by the overall wide variation, and the long tail of models close to the 25% chance level, highlighting the need for improvements in this area by model developers.

8

Taken together, these results show that primary-level visual reasoning remains challenging for current multimodal LLMs. Even the best available models only solve a subset of items reliably and performance varies substantially across the model landscape. In the following sections, we unpack this variation by model type, reasoning capability, cost and later by skill profiles.

### 4.1.1 Model groups: weights and reasoning capability

Figure 5 compares performance on VRB for proprietary (closed-weight, API-access) models and open-weight models. Proprietary models largely dominate the top of the leaderboard; Kimi-K2.5 is the main exception at 60% and then GLM-4.6V and QWEN3-VL are the only other open-weight models appearing in the upper half, reaching an accuracy ≈45%. Small open models such as Gemma-3 4B and Mistral Small 2.1 3B remain around the mid 20% range effectively at or near the chance baseline. It is also worth noting that many of the strongest open-weight reasoning models currently available are text-only and therefore do not appear in this multimodal evaluation.

These results indicate that, as of January 2026, open-weight models mostly lag behind proprietary systems in visual reasoning on VRB, though the 60% performance of Kimi-K2.5 represents a significant step forward in closing this gap. However, in LMIC contexts where on-device or locally hosted models are preferred for privacy, connectivity, and cost—this progress meets a practical limit. Because Kimi-K2.5 is a trillion-parameter model requiring massive hardware resources, it remains unsuitable for low resource local hosting. Consequently, our results highlight a sharp trade-off: the scale required for open-weight models to achieve frontier-level visual reasoning creates an infrastructure barrier that negates the cost and accessibility benefits of local hosting.
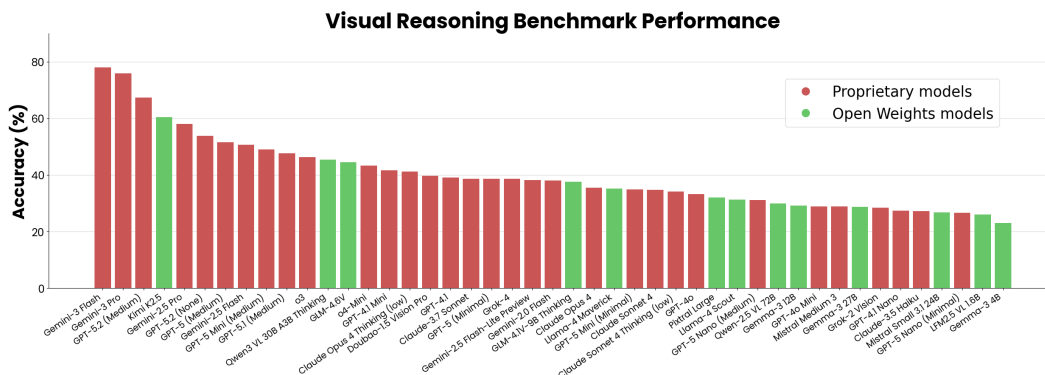


Figure 5: Accuracy on the Visual Reasoning Benchmark (VRB) by Weights Availability.

To examine the role of explicit reasoning, we group models into reasoning variants and non-reasoning variants (Figure 6). Reasoning models are those that expose dedicated "thinking" or chain-of-thought modes (e.g. Gemini 3.0, GPT-5 (medium), Claude Opus 4 Thinking), while their corresponding base chat models and other standard instruction-following models are treated as non-reasoning variants.
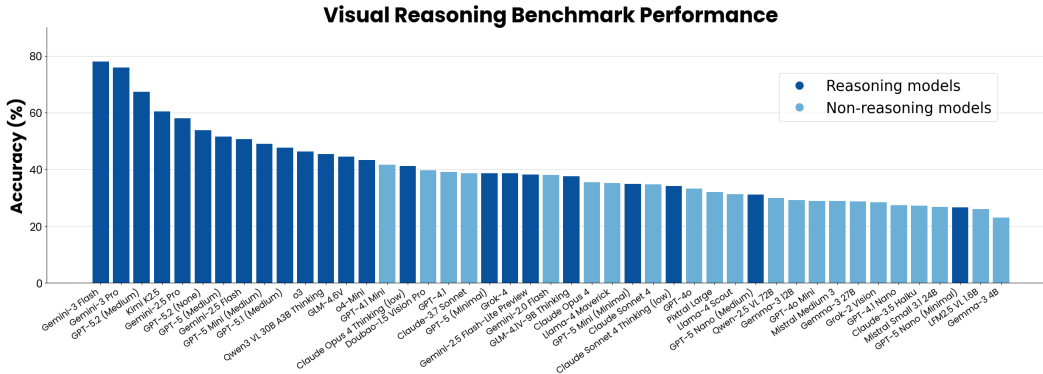
Figure 6: Accuracy on the Visual Reasoning Benchmark (VRB) by Reasoning Capability.

### 4.1.2 Accuracy vs Cost: The Visual Value Frontier

For models with published pricing, we record the input token cost in USD per million tokens and examine the trade-off between VRB accuracy and inference cost (Figure 7). Token prices span nearly four orders of magnitude, from around **$0.01** to close to **$100** per million input tokens. At the very low end of this range, models achieve only ≈**23%** accuracy and are unlikely to be useful for most classroom tasks involving visual reasoning. Around **$0.10–$0.20** per million tokens, the best models reach accuracies of roughly **38–40%**. The current value frontier is dominated at the high end by a notable anomaly: Gemini-3 Flash, achieving an accuracy of 78% at a price point of only **$0.50**, exceeding the performance of models that cost ten times as much.

For any given cost band there is typically a spread of around 15–20 percentage points between the lowest- and highest-performing models, indicating substantial variation in "visual value for money".
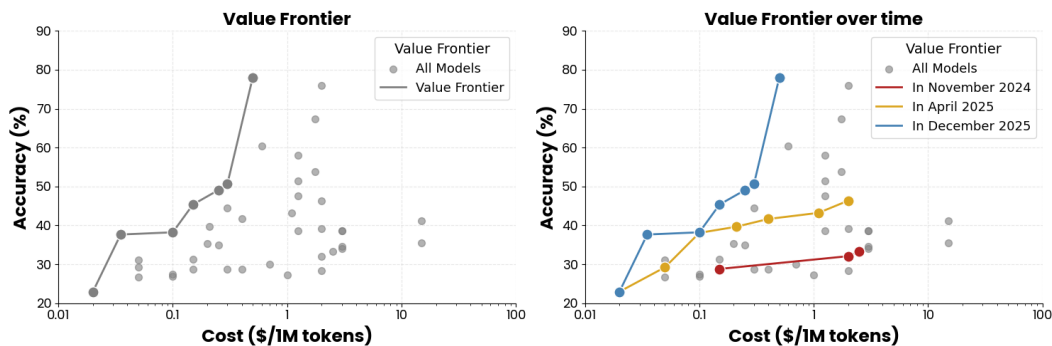


Figure 7: Left: **Cost vs accuracy,** the best models being on the value frontier (grey line). Right: **Value Frontier over time** on VRB.

We also track how this value frontier evolves over time by comparing snapshots from late 2024, mid-2025 and late 2025. Across almost all price points, the frontier shifts upwards: for example, at around **$0.10/M** the best available accuracy increases from ≈30% to ≈38–39% over this period, while at higher costs the top models move from the low-30s to around 60%. This pattern mirrors trends in text-based benchmarks, but here it shows that non-verbal visual reasoning capabilities are also improving rapidly. VRB is thus sensitive enough to register incremental gains in visual reasoning and to highlight models that offer the best trade-off between performance, openness and cost.

## 4.2 Task and Skill Level Performance

We analysed performance across two dimensions to investigate which types of questions MLLMs are performing better on, and where improvements are most needed. As highlighted in Section 3.3[2]:

- **Task** considers how the question was structured (e.g. odd one out, pattern completion, matching).
- **Skill** considers the underlying visual operation(s) (e.g. folding, counting, rotating) needed to solve the questions.

These dimensions help us to understand and distinguish between where the models fail (the visual layout) and why they fail (the cognitive bottleneck).

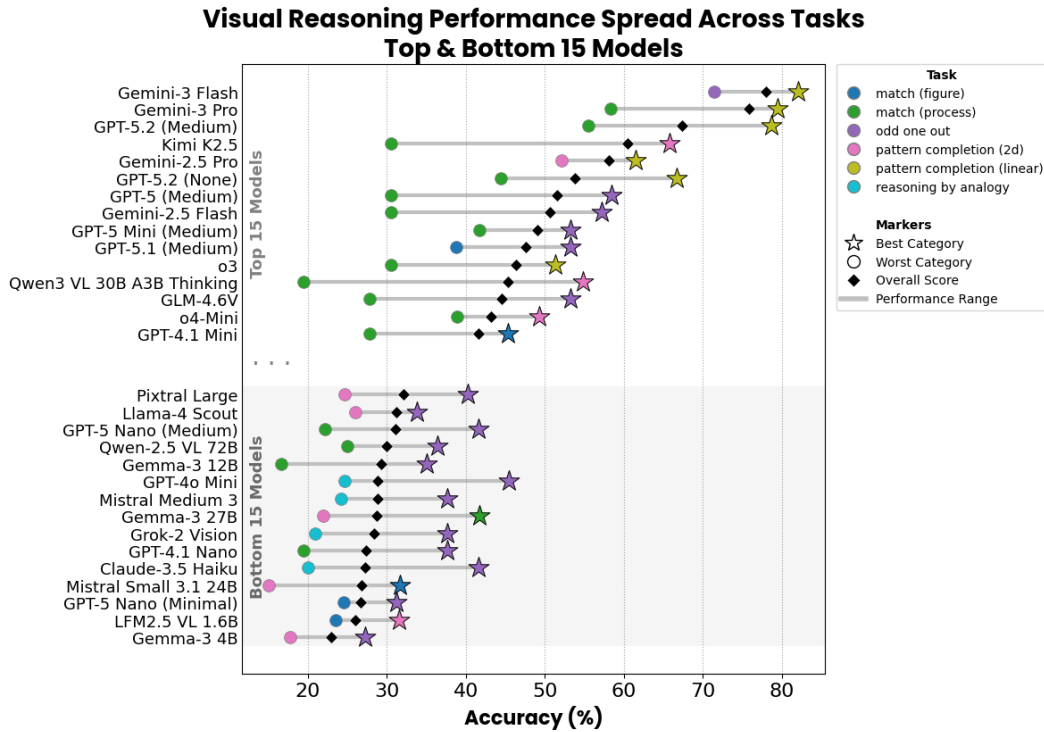### 4.2.1 Task formats and performance



Figure 8: The spread of performance over tasks for the top and bottom 15 models.

To uncover the mask of overall performance accuracy we analysed the accuracy by looking at the task format. Figure 8 shows the performance spread across task categories for the top and bottom 15 models, highlighting each model's overall score (diamond), best task (star), worst task (circle), and the range between the best and the worst tasks (grey line).

For the Top 15 models, we can see that there is a large volatility within models. While many top-tier models show a 10–25 percentage point gap between their best and worst task, this gap is even wider in newer open-source models like Kimi-K2.5 and Qwen3-VL, where the difference is nearly 35 points. We observe that even as open-source models become more capable on average, their performance becomes more inconsistent across different types of visual logic. The strongest models achieve the highest accuracies on *pattern completion (linear)* tasks which suggests that models find it easier to extend sequences once a rule is inferred. On the other hand, the worst-performing task for most

---

[2]See Figures 1, 2 and 3 for question samples and Appendix A for definitions of categories and tags.

frontier models is *Match (process)*, which could be due to the multi-step transformations required rather than a simple visual comparison.

In the bottom 15 models, overall accuracy clusters much lower, around the 20%–40% range. These models perform better on *Odd One Out* tasks, which focus on direct comparisons rather than learning and applying a transformational rule.

Task formats give a useful view on which question types models can struggle on. However, a single task category can use multiple visual operations e.g. *Match (Figure)* might only require *counting* in one question and *rotation* or *folding* in another. Because skill tags are unevenly distributed across tasks and can co-occur within the same item, task-level results do not show us the full picture. The next section shifts from presentation to the source of difficulty by estimating how each visual skill affects accuracy.

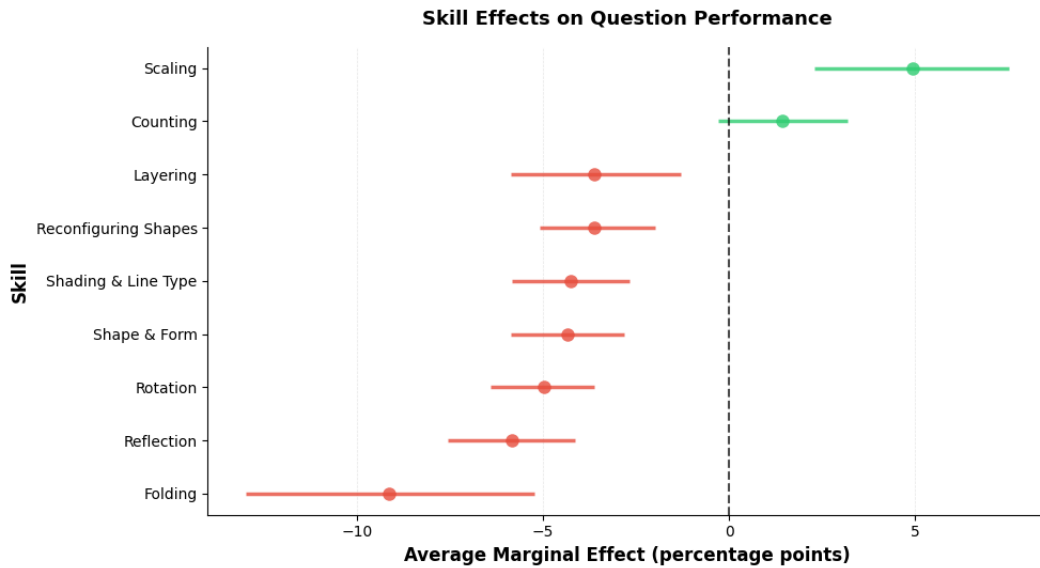### 4.2.2 Skill effects on question performance



Figure 9: **Average marginal effects of specific skills on question performance.** This plot shows the predicted change in accuracy (percentage points) associated with each skill tag, accounting for model identity, dataset, and task category. Uncertainty is represented by 95% bootstrap confidence intervals.

To find out which visual operations AI models are finding the most challenging, we calculated the marginal effect of each skill on item performance while controlling for model, dataset, and task. The resulting marginal effects represent the average change in accuracy (percentage points) that is associated with the presence of a specific skill.

There are two skills that stand out to facilitate correct responses. Items that are tagged with *Counting* and *Scaling* are strong predictors for success. *Scaling* provides a boost of approximately +5 percentage points while *Counting* has an increase of +1.4 points relative to the absence of these skills. This suggests that multimodal LLMs are relatively more proficient at detecting numerosity and recognising monotonic size changes.

All other skills correlate with a decrease in performance. *Layering* and *Reconfiguring Shapes* result in modest decreases of approximately $-3.5$ percentage points. The most significant drop happens in tasks requiring complex mental manipulation. *Shading and Line Type* and *Shape and Form* both lead to a decrease of $-4.5$ points, whilst *Rotation* ($-5$ points) and *Reflection* ($-5.8$ points) represent substantial hurdles. The most difficult skill is *Folding* which is associated with a decrease of $-9$ points. These findings indicate that while models can "see" and "count" discrete elements, they

12

fundamentally struggle with spatial transformations—the mental gymnastics of rotating, reflecting, or folding objects in 3D space.

## 4.3 Error Analysis



Figure 10: **Model accuracy across three error-severity levels** (No error, Minor, Medium). Thin lines show individual model trajectories; the thick red line shows the mean across all models.

A distinct contribution of VRB is the inclusion of classroom-authentic questions that contain noise such as photocopying artefacts (e.g., faded line segments and inconsistent shading) or minor production errors that are likely to be present in real-world use cases, particularly in LMIC settings.

As shown in Figure 10, we see a decline in overall accuracy as the visual quality decreases. The mean accuracy across all evaluated models drops from 38% on pristine items to approximately 30% on items with moderate error severity.

The steepest rate of decline was for the highest performing models. This suggests that high-performing frontier models achieve these scores through precise pixel matching—relying on low-level visual patterns that lack structural integrity, which is easily disrupted by the presence of artefacts. Lower-performing models nearer the chance threshold of 25% experienced less variance due to the image quality, indicating that these architectures fail at extracting meaningful structural information regardless of the quality of the input.

## 5 Discussion

VRB was designed to investigate a practical question: Are current multimodal LLMs equipped to navigate the non-verbal reasoning challenges that are typically expected of primary school students. Our results demonstrate a "jagged frontier" of capability. While some models achieve high accuracies on certain visual reasoning skills, their ability is unevenly distributed and improvements are most needed for "dynamic" visual reasoning such as *folding, reflection, and rotation*.

### 5.1 AI models demonstrate a "jagged frontier" where their low performance on visual reasoning has issues for supporting children's learning

The VRB results highlight the jagged frontier of modern AI where models that achieve state-of-the-art results on textual reasoning benchmarks falter on primary-level visual logic. To support learning effectively, AI must meet a minimum capability threshold in visual reasoning. An AI tutor that can solve complex textual mathematics but then fails to recognise a student's error in a

simple pattern-completion task is not yet trustworthy for unsupervised classroom use. Similarly, an MLLM for grading tasks that cannot provide accurate results or feedback on a subset of visually grounded questions risks reinforcing or even creating student misconceptions. A further challenge is the 'invisible' nature of such failures, where trust has been built around a tool due to its strong performance across certain tasks, reducing the likelihood that teachers or learners will monitor outputs for hallucinations or subtle errors.

The lack of robustness in model accuracy on questions that contain minor errors typical of learning materials raises further concerns about classroom suitability, especially in settings where photocopying is the norm. This limitation may also erode potential time-saving benefits of AI if teachers must spend additional time improving material quality before use. Taken together, the VRB findings suggest that classroom applications involving visual reasoning must be implemented with explicit human oversight and quality-assurance checks.

## 5.2 AI models show evidence of being constrained by "the spatial ceiling" and improvements are most needed on "dynamic" visual reasoning

The failure patterns that we see in the VRB appear similar to the established cognitive frameworks in human spatial reasoning. This distinguishes between "static" or object-based skills and "dynamic" or spatial transformation skills [Newcombe and Shipley, 2014]. In humans, tasks involving dynamic tasks such as mental rotations or recognising reflections tend to be more difficult to solve than those requiring static skills such as recognising structural features [Caissie et al., 2009, Green and Kluever, 1992]. For example, studies of children aged 9–12 have recorded accuracies between 73% and 81% for object-based skills while scores for spatial tasks involving visualisation and spatial relations dropped to between 55% and 72% [Soluki et al., 2021].

The VRB reveals a similar hierarchy of difficulty for AI models. While contemporary models perform better on static skills such as *counting* and *scaling*, their performance reaches a spatial ceiling when faced with dynamic operations such as *folding, rotation*, and *reflection*. This indicates that while models have acquired an object-level grasp of images, they lack the fundamental internal mechanisms to simulate spatial movements.

This finding aligns with the "spatial gap" found in vision language research [Wang et al., 2024], where vision is treated as a supporting cue and secondary to linguistic priors. This leads to a collapse in relational logic when a task requires the model to navigate spatial transformations even if the object is perceived correctly [Wang et al., 2024, Chen et al., 2026].

Furthermore, the performance drop in mental manipulation tasks suggests a deeper architectural limitation. Current MLLMs seem to struggle to maintain structural integrity and geometric identity through transformations. This deficit is characterized in recent literature as a fundamental failure of "spatial imagination" [Carpenter et al., 1990], confirming that even frontier models still fall short of the visual competencies that emerge in early human development [Chen et al., 2026].

This architectural bottleneck is made worse by the cognitive load of rule induction [Carpenter et al., 1990]. In humans, performance is higher on simple pairs (analogies) than on complex matrix patterns, primarily because the former requires integrating fewer visual features [Soluki et al., 2021, Carpenter et al., 1990]. MLLMs face a parallel bottleneck in which, if the initial visual encoding fails to spot the relational features, the model is effectively reasoning from an inaccurate version of the image [Chen et al., 2026]. Without a solid grasp of what it is seeing, the model's reasoning becomes decoupled from the visual ground truth and leads to misgrounded explanations that lack reasoning fidelity [Wang et al., 2024].

## 5.3 Limitations and Future Directions

Our study is bounded by several limitations, including a focus on high-stakes aptitude exams that represent the upper end of difficulty for 11-year-olds rather than routine early primary tasks that are the core of foundational numeracy. Additionally, the VRB utilises static, single-turn interactions which do not evaluate a model's ability to provide incremental hints or respond to student diagrams in a dialogue. There is also a fundamental difference in processing: humans solve these problems using internal imagery, whereas models operate through visually conditioned language generation, meaning a correct answer does not necessarily imply human-like reasoning.

Future research can take this further through process-aware evaluation to ensure that models demonstrate reasoning fidelity and avoid shortcuts. By drawing on work such as ConceptARC [Beger et al., 2025], future work could also require models to generate natural language explanations to distinguish between correct intended abstractions and correct unintended rules [Beger et al., 2025]. This would allow the differentiation of false patterns that work for specific demonstrations but fail to generalise to broader classroom contexts. Additionally, this would also help determine if failures on classroom tasks stem from a lack of high-order logic or a perceptual bottleneck. Further research could mitigate these perceptual hurdles by evaluating the impact of external tool use by granting models access to Python libraries for precise geometric processing [Beger et al., 2025].

Finally, moving toward multi-turn benchmarking will better simulate real AI tutor use cases and determine whether models can identify why a student's visual reasoning is flawed and provide a more effective corrective scaffold.

## 6 Conclusion

The Visual Reasoning Benchmark (VRB) evaluates multimodal LLMs on classroom-authentic, minimal-text visual reasoning items from primary assessments in Zambia and India. Across a broad set of proprietary and open models, accuracy ranges from near-chance to 78%, showing that primary-level non-verbal reasoning remains unsolved for current systems.

Performance is defined by two consistent limits. Firstly, accuracy can swing sharply across task format, producing a "jagged frontier" where strong performance on linear pattern completion coexists with failures on process matching and 2D matrix completion. Secondly, skill effects reveal a clear spatial ceiling: controlling for task category, dataset, and model identity, one finds that static skills (e.g. *counting, scaling*) have slightly higher accuracy, while *folding, reflection*, and *rotation* pose the greatest challenges.

In realistic LMIC classroom conditions, this brittleness is a deployment risk, and implementation of tasks that require visual reasoning should be done with human oversight.

## Acknowledgments

## References

C. Beger, R. Yi, S. Fu, A. Moskvichev, S. W. Tsai, S. Rajamanickam, et al. Do AI models perform human-like abstract reasoning across modalities? *arXiv preprint arXiv:2510.02125v3*, 2025.

D. Broadbent, C. Hayden, R. Tate, and L. von Kotze, editors. *11+ Non-Verbal Reasoning For GL Assessment*. Coordination Group Publications Ltd. (CGP), 2018.

H. Cai, Y. Yang, and W. Hu. MM-IQ: Benchmarking human-like abstraction and reasoning in multimodal models. *arXiv preprint arXiv:2502.00698*, 2025.

A. F. Caissie, F. Vigneau, and D. A. Bors. What does the mental rotation test measure? An analysis of item difficulty and item characteristics. 2009.

P. A. Carpenter, M. A. Just, and P. Shell. What one intelligence test measures: a theoretical account of the processing in the Raven Progressive Matrices Test. *Psychological Review*, 97(3):404, 1990.

L. Chen, W. Xie, Y. Liang, H. He, H. Zhao, Z. Yang, et al. BabyVision: Visual reasoning beyond language. *arXiv preprint arXiv:2601.06521*, 2026.

F. Chollet. On the measure of intelligence. *arXiv preprint arXiv:1911.01547*, 2019.

F. Chollet, M. Knoop, G. Kamradt, B. Landers, and H. Pinkard. ARC-AGI-2: A new challenge for frontier AI reasoning systems. *arXiv preprint arXiv:2505.11831*, 2025.

D. Cobbe, C. Kosinski, and M. Bavarian. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.

K. E. Green and R. C. Kluever. Components of item difficulty of Raven's matrices. *The Journal of General Psychology*, 119(2):189–199, 1992.

D. Hendrycks, C. Burns, S. Basart, et al. Measuring massive multitask language understanding. In *International Conference on Learning Representations (ICLR)*, 2021a.

D. Hendrycks, C. Burns, C. P. Zou, et al. Measuring mathematical problem solving with the MATH dataset. In *NeurIPS Datasets and Benchmarks*, 2021b.

Y. Jiang, J. Zhang, K. Sun, Z. Sourati, K. Ahrabian, K. Ma, F. Ilievski, and J. Pujara. MARVEL: Multidimensional abstraction and reasoning through visual evaluation and learning. *arXiv preprint arXiv:2404.13591*, 2024.

Y. Liu, Z. Chen, T. Zhang, and X. Wang. VERIFY: A benchmark of visual explanation and reasoning fidelity for MLLMs. *arXiv preprint arXiv:2503.11557*, 2025.

P. Lu, K. Wang, J. Pan, et al. MATH-V: Measuring mathematical reasoning in diagram-augmented contexts. *arXiv preprint arXiv:2402.14804*, 2024.

I. Natsheh and R. Karsenty. Exploring the potential role of visual reasoning tasks among inexperienced solvers. *ZDM*, 46(1):109–122, 2014.

N. S. Newcombe and T. F. Shipley. Thinking about spatial thinking: New typology, new assessments. In *Studying visual and spatial reasoning for design creativity*, pages 179–192. Springer Netherlands, 2014.

J. Park, D. Lee, and K. Cho. Explain with visual keypoints like a real teacher: Multimodal solution explanation with ME2. *arXiv preprint arXiv:2504.03197*, 2025.

J. Parkinson and Q. Cutts. Improving primary school pupils' spatial skills leads to computational thinking gains. In *Proceedings of the 30th ACM Conference on Innovation and Technology in Computer Science Education V. 1*, pages 646–652, 2025.

A. M. Purcar, M. Bocos, A. L. Pop, A. Roman, D. Rad, D. Mara, C. Crisan, R. Răduţ-Taciu, E. L. Mara, I. Todor, and D. G. Triff. The effect of visual reasoning on arithmetic word problem solving. *Education Sciences*, 14(3):278, 2024.

J. C. Raven. Mental tests used in genetic studies: The performances of related individuals on tests mainly educative and mainly reproductive. Master's thesis, University of London, 1936.

S. Soluki, S. Yazdani, A. A. Arjmandnia, J. Fathabadi, S. Hassanzadeh, and V. Nejati. Comprehensive assessment of spatial ability in children: A computerized tasks battery. *Advances in Cognitive Psychology*, 17(1):38, 2021.

J. Wang, Y. Ming, Z. Shi, V. Vineet, X. Wang, Y. Li, and N. Joshi. Is a picture worth a thousand words? Delving into spatial reasoning for vision language models. *arXiv preprint arXiv:2406.14852*, 2024.

P. Wang, Z. Li, F. Yin, D. Ran, and C.-L. Liu. MV-MATH: Evaluating multimodal math reasoning in multi-visual contexts. In *CVPR 2025*, 2025.

Y. Xie, P. Lu, Y. Li, et al. MathVista: Evaluating mathematical reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*, 2023.

W. Xu, J. Wang, W. Wang, et al. VisuLogic: A benchmark for evaluating visual reasoning in multi-modal large language models. *arXiv preprint arXiv:2504.15279*, 2025.

R. Zhang, J. Liu, T. Wang, et al. MathVerse: Does your multi-modal LLM truly see the diagrams? *arXiv preprint arXiv:2403.14624*, 2024.

# A  Definition of Skill Tags and Task Categories

| Skill Tag | Definition |
| --- | --- |
| Shape and Form | The question involves shape (regular or irregular) as a key aspect and the ability to distinguish subtle differences in shape or form, including shape changes and shapes being added in different configurations. Where shape change is not the primary process, for example if it occurs as a result of reflection, no tag is given. |
| Counting | The question requires counting the number of objects or parts of objects to determine the solution. |
| Shading and Line Type | The question involves a change in shading, fill, or line style, including type (curved/straight, etc.). Shading may be a complete fill with lines, dots or similar, or another pattern of objects contained in a shape. |
| Rotation | The question involves an object or parts of an object which have been rotated. This is commonly by 45°, 90°, or 180° but could also be other turns. Where either reflection or rotation could occur, both should be tagged. |
| Reflection | The question involves an object or parts of an object which have been reflected. This could be horizontal, vertical, or diagonal symmetry. Where either reflection or rotation could occur, both should be tagged. |
| Scaling | The question involves objects or parts of objects which increase or decrease in size. This may be over a sequence of shapes or just 2 shapes getting smaller or larger. |
| Folding | This question involves a shape being folded and punched or cut before the pattern on the unfolded shape is determined, or questions where folds are made in an asymmetric way. |
| Reconfiguring Shapes | This question involves objects or parts of objects changing position, being brought together or separated. They may or may not be touching but should not overlap. |
| Layering | This question involves identifying shapes hidden within larger figures and visualising objects which are placed one on top of the other, or which were on top of each other and are taken apart. |

| Task Category | Definition |
| --- | --- |
| Reasoning by Analogy | The question involves a comparison where the relationship between the first set of images needs to be matched by the relationship between the second image set. |
| Odd One Out | The question involves a group of images where most are identical and the different one must be selected, or where most images share a common feature and the one without this feature is identified. |
| Pattern Completion (Linear) | The question involves a sequence. The last image should be chosen to continue the series in the same way. |
| Pattern Completion (2D) | The question involves a grid structure with one part missing and needs to be filled in to complete the grid following the existing pattern. |
| Match (Process) | The question involves a group of question figures which show a process. The answer is chosen to match the output of this process. |
| Match (Figure) | The question involves a single problem figure and several answer figures. The answer figure is chosen which matches the problem figure in the way described in the question. |

# B   Definitions of Error Types and Severity

| Error Type | Definition |
|---|---|
| Image error | This is an issue with the images that make up the question (for example, inconsistent shading or inaccurate sizing). |
| Question presentation error | This is an issue with the text in the question or another issue with the production of the question, such as a mouse cursor overlaying an image. |

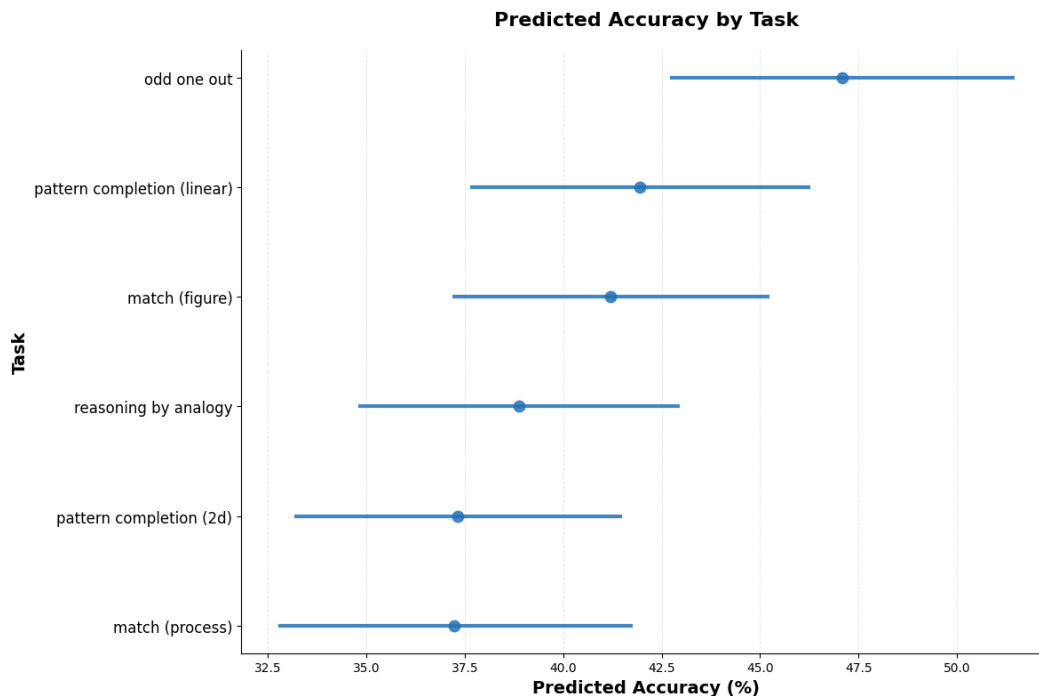| Severity | Definition |
|---|---|
| Minor | A small error which is not likely to affect a student's ability to understand and answer the question. The correct answer is clear. It is likely to be either easily recognisable (such as a misprint with a mouse cursor over the image), not essential for answering the problem (such as an absent dotted line for a fold), or does not have a significant impact on meaning (such as minor grammar issues). |
| Medium | A more substantial error which could affect a student's ability to distinguish between the options, but the correct answer remains clear. For example, an issue with degrees of rotation where one correct answer still stands out or misleading wording in the written question where the visual image is still clear. |

# C   Additional Figures



Figure 11: Task effects on item difficulty ranked from most difficult (bottom) to easiest (top).
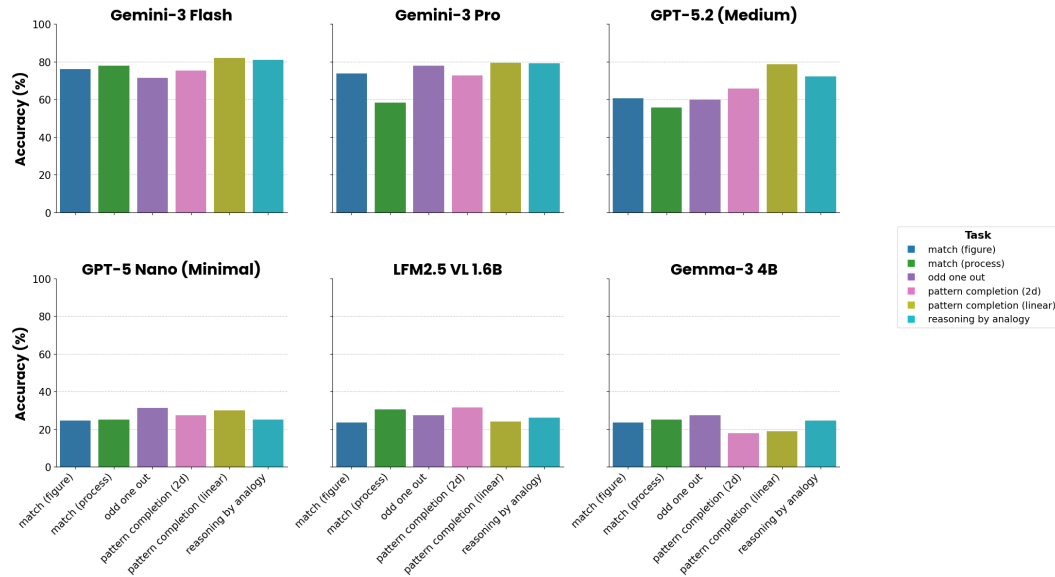
Figure 12: Per-model category bar charts (Top and bottom 3 models).

## Extended Experimental Details

## Data / Code Availability

The evaluation prompts, scoring code, and analysis pipeline are publicly available at `https://github.com/AI-for-Education/vrb-benchmark`. The question images and answer keys are held privately to preserve benchmark integrity; access can be arranged by contacting the authors.