

Benchmarking and Enhancing Text-to-Image Models for Generating Visual Representations in Early Arithmetic Education

Anonymous ACL submission

Abstract

Visual representations are highly effective in early arithmetic education, as they make abstract mathematical symbols more concrete and support the development of numeracy and reasoning skills. However, creating such visuals is labor-intensive for teachers. Thus, we introduce the equation-to-visual generation task and E2V-Bench, a benchmark for generating pedagogically meaningful visuals from arithmetic equations. Developed with insights from primary school math teachers and informed by visual patterns extracted from six educational resources, E2V-Bench comprises 1.5K arithmetic problems spanning four visual types. We also propose new automatic metrics for evaluating generated visuals. A systematic evaluation on E2V-Bench reveals that open-source text-to-image models perform substantially worse than the strongest closed-source models. Building on these findings, we curate a high-quality training dataset and demonstrate that our model adaptation strategies, including rejection sampling fine-tuning, prompt refinement, and regeneration, significantly improve model performance. This work establishes a foundation for studying equation-to-visual generation (a novel reasoning task) and can support teachers in creating visuals for arithmetic education.

1 Introduction

Visuals are essential tools for teaching foundational arithmetic concepts (e.g., addition, subtraction, multiplication and division) in early primary education (Grades 1–3). They help transform abstract mathematical symbols into intuitive representations (Cooper et al., 2018), supporting student understanding and problem solving (Mayer, 2002). For instance, rather than memorizing that “ $7 + 2 = 9$ ”, young learners may benefit from a visual showing seven objects (e.g., apples) combined with two more to make nine, which helps them grasp the meaning of the “+” operation (Mayer, 2002), as illustrated in the bottom right panel of Figure 2.

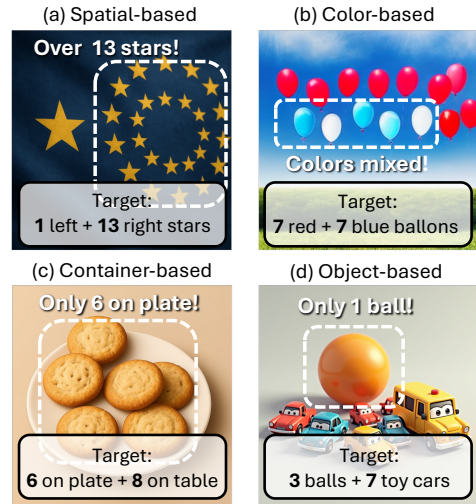


Figure 1: Failure cases of text-to-image models: (a) GPT-Image-1, (b) LMD, (c) Bagel, (d) Show-o2.

However, despite their pedagogical value, producing high-quality instructional visuals remains largely manual and time-consuming (Kaitera and Harmoinen, 2022; Boonen et al., 2016; Xu et al., 2021). Recent work has begun to explore automated generation or retrieval of instructional visuals. Math2Visual (Wang et al., 2025) generates visuals from math word problems, but relies on fixed icon sets and predefined templates, limiting its flexibility in representing standalone arithmetic concepts. Other approaches, such as text-image matching for retrieving visuals for textbook creation (Singh et al., 2023) and text-driven visual question generation (Singh et al., 2019; Luo et al., 2024), focused on retrieval or question construction, rather than generating instructional visuals.

In this work, we focus on primary school level arithmetic education, which deals with learning mathematical operations (addition, subtraction, division, and multiplication) via equations. We formalize the task of equation-to-visual (E2V) and introduce Equation-to-Visual Benchmark (E2V-Bench), a novel benchmark for gener-

ating pedagogically meaningful images from arithmetic equations, together with automatic evaluation metrics for assessing the quality of generated visuals. We manually collected 189 instructional visuals from six educational resources and distilled four recurring visual types that guide our benchmark design. We further conducted interviews with ten primary-school mathematics teachers to refine and validate these design choices.

Using E2V-Bench, we evaluate nine state-of-the-art Text-to-Image (T2I) models across five model families. We show that GPT-Image-1 (OpenAI, 2025), LMD (Lian et al., 2024), Flux.1-dev (Blackforest, 2024), and Bagel (Deng et al., 2025) are top-performing baselines on this task. However, as illustrated in Figure 1, even these strong models exhibit systematic failures in generating pedagogically faithful visuals. Building on these findings, we curate 4,072 high-quality training examples through a generation-and-filtering pipeline from GPT Image-1 and a template-based icon-to-image generation pipeline using the Bagel model. Leveraging this data, we improve the Bagel model through rejection-sampling fine-tuning. In parallel, we adapt prompt-refinement and regeneration techniques for the layout-to-image model LMD, substantially improving its performance. Our contributions are summarized as follows:

- We introduce E2V-Bench, the first exploration on generating pedagogical visuals for arithmetic education, propose automatic metrics and systematically evaluate nine T2I models.
- We release a high-quality training dataset to provide reliable supervision for the E2V task.
- We show that T2I enhancement strategies (i.e., rejection sampling fine-tuning, prompt refinement, and regeneration) yield substantial performance gains on E2V-Bench.

2 Related Work

Visuals in Teaching Primary-School Arithmetic

Visuals are widely recognized as beneficial for learning arithmetic skills (Kaitera and Harmoinen, 2022; Boonen et al., 2016). Well-designed visuals make mathematical ideas more accessible (Mayer, 2002; Evagorou et al., 2015; Small and Lin, 2025), foster student engagement (Cooper et al., 2018), and support more efficient learning (Arcavi, 2003). Various pedagogical designs have been proposed to structure arithmetic relationships. For instance, the

bar model (Hoven and Garelick, 2007) expresses numerical relations through proportional bars, and has been shown to enhance children’s problem-solving ability (Osman et al., 2018) and their use of effective cognitive strategies (Morin et al., 2017). More recently, the Noyon framework depicts mathematical problems with modular iconic components, offering a flexible yet structured approach to visualize mathematical reasoning (Saquib et al., 2021). Collectively, these approaches underscore the importance of carefully designed visuals in arithmetic education.

Automated Visual Generation and Retrieval

Despite their benefits, producing instructional visuals still requires considerable effort from teachers (Xu et al., 2021), motivating automation of visual generation and retrieval for educational use. Prior work has explored aligning images with textbook text (Singh et al., 2023), generating image-based multiple-choice questions from semantic representations (Singh et al., 2019), and combining multimodal information to support question generation (Luo et al., 2024). However, these methods generally focus on broad educational applications rather than making the mathematical concepts explicit. Early efforts such as VisualMath (Dwivedi et al., 2017) rendered simple word problems with pre-existing images, but addressed only basic operations (+, -) and lacked pedagogical validation. More recently, Math2Visual (Wang et al., 2025) generates visuals from math word problems with stronger educational grounding, though its reliance on fixed templates limits flexibility in representing standalone arithmetic concepts within equations. These efforts highlight progress toward automation but underscore the absence of benchmarks for generating pedagogically grounded visuals directly from equations.

Text-to-Image Benchmarks

Benchmarks for T2I models are typically built from curated prompt sets targeting specific capabilities, with general-purpose benchmarks such as DALL-E 3 Eval (Betker et al., 2023), DrawBench (Saharia et al., 2022), PartiPrompts (Yu et al., 2022), and HEIM (Lee et al., 2023) covering challenges such as multi-object generation, spatial and attribute control, and factual consistency. Other datasets focus on narrower skills, such as alignment (e.g., TIFA (Hu et al., 2023), Gecko (Wiles et al., 2025)) or compositional reasoning (e.g., T2I-CompBench (Huang et al., 2023)). More recently,

Kajic et al. (2024) introduced the GECKONUM to study numerical reasoning in T2I models, focusing on counting-based prompts while omitting broader arithmetic operations and educational grounding. To date, no benchmarks systematically evaluate whether T2I models can generate pedagogically meaningful visuals that faithfully represent arithmetic operations for early arithmetic education.

3 Foundations of Equation-to-Visual Generation

3.1 Task Definition

We define the task of E2V generation as producing visuals that faithfully represent the arithmetic concepts encoded in an equation. Rather than reproducing equations as text, generated visuals should depict quantities through concrete objects and spatial arrangements that support learners' intuitive understanding of arithmetic relationships: a form of visual scaffolding commonly used in early mathematics instruction (Kaitera and Harmoinen, 2022; Boonen et al., 2016).

To ground this task in authentic educational practice, we examined how arithmetic equations are visualized in existing instructional materials. Two researchers manually collected 189 visuals from six widely used educational resources, including three textbooks (Ministry of General Education and Instruction, Republic of South Sudan, 2018; Cotton et al., 2021; Moseley and Rees, 2021) and three online platforms (Accessim, 2025; mathematics monster, 2025; fun2dolabs, 2025). After excluding seven decorative visuals that did not correspond to the equations, we conducted a two-stage thematic analysis on the remaining 182 visuals: an exploratory review of 50 visuals to identify recurring patterns, followed by coding the remaining 132 visuals using the identified categories. This analysis revealed four recurring visual types for representing arithmetic equations:

(1) Color-based: Operands are distinguished by assigning distinct colors to otherwise identical objects (e.g., red vs. green apples). Color serves as the primary grouping cue, a perceptual feature that young children readily use for grouping (Harris et al., 1970).

(2) Object-based: Operands are distinguished by using different object types, with object category serving as the primary grouping cue. Prior work shows that young children can readily organize objects by type (Bornstein and Arterberry, 2010).

(3) Spatial-based: Operands are distinguished by placing object groups in different spatial positions (e.g., left vs. right). Spatial separation serves as the primary grouping cue, which children naturally use to organize visual information (Quinn et al., 2008).

(4) Container-based: Operands are distinguished by placing objects into distinct containers, with each container corresponding to one operand. Prior work shows that children readily reason about containment relations (Casasola et al., 2003), making containers a natural basis for grouping.

Based on this analysis, we define the E2V task as generating accurate visuals in each of the four visual types given a single arithmetic equation. Additional details of the thematic analysis are provided in Appendix A.

3.2 Evaluative Study with Teachers

To assess the pedagogical value of the proposed visual types, we conducted an evaluative study with ten experienced primary school mathematics teachers (Grades 1-3; demographics in Table 8). Teachers reviewed visuals following our four visual types, covering both cartoon style and realistic style designs across the four basic arithmetic operations, and rated their pedagogical usefulness on a 7-point Likert scale. Study details are provided in Appendix D.

The four visual types align with current classroom practice. Teachers consistently reported that the four visual types reflect visuals commonly used in Grades 1-3 classrooms (Table 13). All participants indicated prior experience with similar visuals, suggesting that our designs align well with existing instructional practices. Across evaluation criteria, teachers gave uniformly high ratings to visuals generated using the proposed visual types (Table 9), indicating teachers find visuals following our proposed visual types pedagogically valuable and effective for supporting arithmetic learning.

Visual style preferences. Both realistic and cartoon style visuals were rated as pedagogically useful (Table 11), with realistic scoring slightly higher than cartoon. Teachers noted complementary advantages: realistic visuals support real-world grounding, while cartoon visuals reduce visual distraction. These results suggest that supporting both visual styles is beneficial, allowing teachers to select styles based on instructional goals.

Usefulness of Visual Types Across Arithmetic Operations. Teachers' ratings of each visual type's usefulness in teaching vary by operation (Table 12).

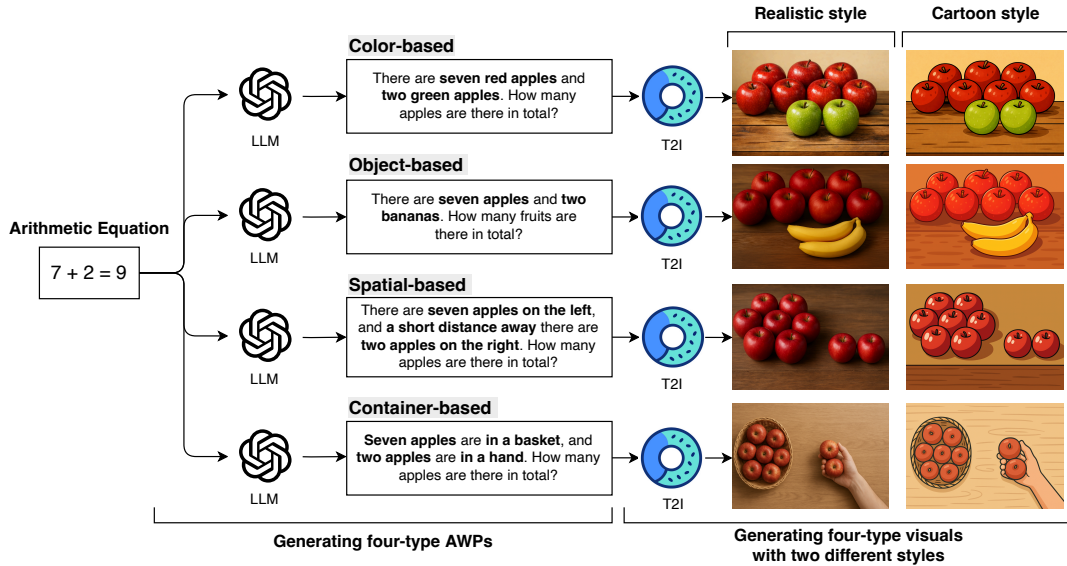


Figure 2: E2V-Bench visual generation pipeline. Equations are systematically generated from four basic arithmetic operations. For each equation, an LLM generates four arithmetic word problems (AWPs), one for each visual type. These AWP are then used as prompts for T2I models to produce visuals in both realistic and cartoon styles.

Container-based visuals were preferred for multiplication and division, where explicit grouping supports understanding of the underlying operations. Color-based visuals were most often associated with addition, as color cues help distinguish addends while conveying their combination. Spatial-based visuals were frequently cited as helpful for subtraction, where spatial separation supports the notion of removing quantities, and were also considered useful for grouping in multiplication and division. Object-based visuals were viewed as a more cognitively demanding option, better suited for challenging students who have already developed a basic understanding of operations.

4 E2V-Bench: Dataset and Metrics

Although the E2V task is defined over mathematical equations, directly prompting existing T2I models with symbolic expressions is ineffective, as models tend to generate stylized renderings of equations rather than object-based visuals that convey quantity and structure (Appendix C.1). This limitation arises because current T2I models are designed for natural language inputs rather than abstract symbolic expressions. To address this modality gap, we adopt a two-step formulation: we first verbalize each equation into an arithmetic word problem (AWP), mirroring how teachers contextualize equations in practice, and then generate visuals from the resulting AWP. This approach enables the generation of pedagogically meaningful visuals

from equations while remaining compatible with standard T2I pipelines.

4.1 Benchmark Construction

We construct E2V-Bench in two phases: equation generation and corresponding AWP generation (Figure 2). We enumerate arithmetic equations over the four basic operations ($+$, $-$, \times , \div), restricting quantities to be below 20 for primary-school suitability, yielding 371 unique equations. Using in-context learning with Gemini-2.5-Flash (Google, 2025), we generate four AWP per equation, one for each visual type (Container, Object, Color, and Spatial-based), guided by eight manually crafted in-context examples per type.

To validate correctness, two researchers reviewed 120 sampled AWP (30 per visual type) and confirmed full alignment with the intended equations and visual types. Pedagogical value was further validated by ten primary-school teachers, who evaluated 32 generated AWP and rated them as suitable for classroom use (Appendix D.3). The final dataset contains 1,484 AWP, with 1,184 for training and 300 for testing (statistics in Table 5).

4.2 Evaluation Metrics

We evaluate model outputs using two criteria informed by prior work and interviews with ten primary-school mathematics teachers.

Quantity Accuracy measures whether the number of objects in the generated visual matches the quan-

Category	Model	Quantity Acc. (%)		Overall Acc. (%)		CLIPScore	
		Realistic	Cartoon	Realistic	Cartoon	Realistic	Cartoon
Closed-source	Recraft-v3	6.00	5.67	4.67	2.67	0.81	0.78
	GPT-Image-1	27.00	22.33	21.33	17.00	0.80	0.81
Diffusion-based	SD-3.5-large	6.67	8.33	3.33	5.00	0.80	0.80
	Flux.1-dev	8.67	9.33	6.33	5.67	0.77	0.75
Layout-to-image	Blueprint	2.67	1.33	1.33	1.00	0.74	0.73
	LMD	22.33	15.33	12.00	9.67	0.75	0.71
Transformer-based	Show-o2	4.33	4.33	2.33	3.00	0.82	0.82
	Bagel	8.00	11.00	5.67	8.00	0.77	0.79
Prompt refinement	PAE	1.00	3.00	1.00	1.00	0.70	0.69

Table 1: Evaluation of T2I models on E2V-Bench. We report quantity accuracy, overall accuracy, and CLIPScore for both realistic-style and cartoon-style outputs. SD-3.5-large denotes the Stable Diffusion-3.5-large.

326 tities specified in the AWP. In practice, we adapt
327 CountGD (Amini-Naieni et al., 2024) to automati-
328 cally count object instances in generated visuals.

329 **Overall Accuracy** evaluates whether a visual
330 aligns with the input AWP. Although recent bench-
331 marks use vision-language models as judges, we
332 explicitly avoid them due to their unreliability in
333 exact object counting (Choudhury et al., 2025). To
334 guarantee precise and reproducible evaluation in
335 early arithmetic settings, we adopt a determinis-
336 tic, rule-based pipeline tailored to each visual type
337 using OpenCV and Scikit-Learn (Appendix C.2).

338 Both criteria are critical in educational settings,
339 where learning depends on the correctness of
340 instructional information (Metzger et al., 2003;
341 Goldin and Shteingold, 2001). To validate align-
342 ment with human judgment, we randomly sam-
343 pled 216 visuals generated by nine models (Ta-
344 ble 1) from the E2V-Bench test set and had two re-
345 searchers independently annotate each visual. They
346 disagreed on only two cases, resolved through dis-
347 cussion. Cohen’s kappa between human annota-
348 tions and the automatic metrics was 0.96 for both
349 criteria, indicating substantial agreement. To fur-
350 ther assess pedagogical validity, we collected ex-
351 pert feedback from ten primary-school teachers
352 during the user study; all teachers agreed that the
353 two criteria adequately capture visual quality for
354 teaching arithmetic skills (details in Appendix D).
355 Following common practice in visual evaluation,
356 we also report CLIPScore (Hessel et al., 2021).

357 5 Benchmarking T2I Models on 358 E2V-Bench

359 We evaluate nine T2I models on E2V-Bench,
360 spanning closed-source (GPT-Image-1 (OpenAI,
361 2025), Recraft-v3 (Recraft, 2025)), diffusion-based
362 (Stable Diffusion-3.5-large (Esser et al., 2024),

363 Flux.1-dev (Blackforest, 2024)), layout-to-image
364 (LMD (Lian et al., 2024), Blueprint (Gani et al.,
365 2024)), transformer-based (Show-o2 (Xie et al.,
366 2025), Bagel (Deng et al., 2025)), and prompt-
367 refinement (PAE (Mo et al., 2024)) models. For
368 each AWP, we generate both cartoon and realistic
369 style visuals. Results are reported in Table 1, with
370 example visuals shown in Figure 13. All results are
371 reported from a single run per model.

372 **Performance Across Model Families.** The
373 closed-source GPT-Image-1 model achieves the
374 best performance among the nine models, with an
375 overall accuracy of 21.33% (realistic) and 17.00%
376 (cartoon). The other closed-source model, Recraft-
377 v3, performs substantially worse, with overall ac-
378 curacy of 4.67% (realistic) and 2.67% (cartoon).

379 Open-source models exhibit substantial variation
380 across families. Diffusion-based models perform
381 poorly overall: Stable Diffusion-3.5-large achieves
382 low overall accuracy (3.33% realistic, 5.00% car-
383 toon), while Flux.1-dev shows modest improve-
384 ment, reaching 6.33% (realistic) and 5.67% (car-
385 toon), but still struggles to reliably control object
386 counts. Layout-to-image models show mixed per-
387 formance. Blueprint performs worst across most
388 metrics, whereas LMD substantially outperforms
389 other open-source models, achieving 12.00% over-
390 all accuracy in the realistic setting and 9.67% in the
391 cartoon setting, alongside relatively strong quan-
392 tity accuracy. These results indicate that explicit
393 layout conditioning provides benefits for numerical
394 grounding. Transformer-based models exhibit vari-
395 able behavior: Bagel surpasses Show-o2 in both
396 quantity and overall accuracy, achieving 5.67%
397 overall accuracy in the realistic setting and 8.00%
398 in the cartoon setting. Finally, PAE records near-
399 zero overall accuracy, confirming that prompt re-
400 finement alone is ineffective at generating images

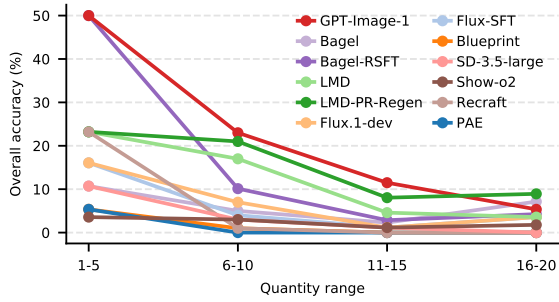


Figure 3: Overall accuracy of realistic-style visuals evaluated across different quantity ranges. Cartoon-style visual results are in Figure 10.

with correct object counts, consistent with prior findings (Cao et al., 2025b). CLIPScores remain comparable across models, indicating similar levels of overall semantic alignment (Hessel et al., 2021).

Variation Across Quantities and Visual Types.

Across models, overall accuracy declines sharply as object counts increase (Figure 3), dropping from around 50% for 1–5 objects to below 15% beyond 11 objects. Performance also varies by visual type (Figure 4): GPT-Image-1 shows the most balanced profile, particularly outperforming other models in the color visual type. LMD performs comparably to GPT-Image-1 in the object, spatial, and container visual types. The Bagel emerges as the second-best performer on container and color-based visuals, indicating its potential for further improvement. The LMD-PR-Regen and Bagel-RSFT are enhanced versions of LMD and Bagel; experimental details and analysis are provided in Section 6.

Key Insights and Open-source Potential. Overall, the closed-source GPT-Image-1 delivers the highest-quality outputs but cannot serve as a reproducible research baseline due to its proprietary nature. Among open-source models, LMD, Flux.1-dev, and Bagel stand out as the strongest representatives and demonstrate potential for further improvement. The next section shows how targeted strategies, including rejection-sampling fine-tuning, prompt refinement, and regeneration, can further enhance model performance.

6 Enhancing Equation-to-Visual Generation

As a foundational study of E2V generation, we explore how existing T2I models, such as LMD, Bagel and Flux.1-dev, can be improved through targeted strategies in the task. Considering costs,

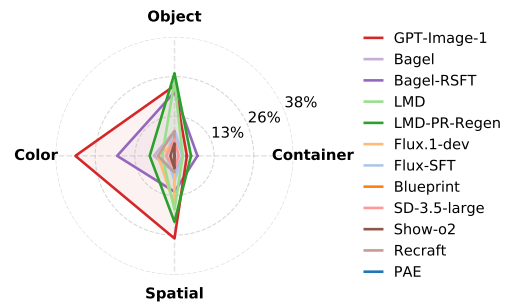


Figure 4: Overall accuracy of realistic-style visuals evaluated across different visual types. Cartoon-style visual results are in Figure 11.

Method	Quant Err. (%)	Bound Err. (%)	Format Err. (%)
LMD	41.67	18.67	2.00
LMD+PR	27.00	22.33	0.00
LMD+PR+Regen	21.00	6.67	0.00

Table 2: Error rates for layout generation.

we report results for the realistic style; however, the same protocols can be applied to cartoon style.

6.1 Layout-to-Image Model

LMD is a representative layout-to-image framework that first generates bounding box layouts in JSON format from text input, and then conditions a diffusion model to produce the final visual (Lian et al., 2024).¹ Given that accurate object counts largely rely on precise layouts (Binyamin et al., 2025), improving the quality of layout generation is essential.

To identify common failure modes in layout generation, two researchers manually reviewed 50 generated layouts and categorized the errors into three recurring types: (1) Quantity error: the number of generated objects does not match the input AWP; (2) Boundary error: bounding boxes extend beyond the image boundary; and (3) Format error: invalid JSON output that cannot be parsed. Grounded on these findings, we developed rule-based scripts to automatically detect each error type. Table 2 reports the frequencies of these errors across 300 samples from the E2V-Bench test set. To address the observed issues of layout generation, we explore two strategies to improve layout generation: prompt refinement and regeneration.

Prompt refinement. We redesigned the prompt instruction to better align the task of generating accurate bounding boxes. Specifically, the revised

¹All LMD experiments use GPT-4o for layout generation, the state-of-the-art multimodal model at the time (Cao et al., 2025a; Ramachandran et al., 2025).

Method	Quant. Acc. (%)	Overall Acc. (%)	CLIPScore
LMD	22.33	12.00	0.75
LMD-PR	26.67	14.67	0.76
LMD-PR-Regen	26.67	15.33	0.76
Flux.1-dev	8.67	6.33	0.77
Flux.1-dev (SFT)	8.00	5.33	0.78
Bagel	8.00	5.67	0.77
Bagel (SFT)	9.00	6.00	0.78

Table 3: Evaluation of refined models on realistic-style visual generation. PR means with prompt refinement and Regen means with regeneration.

prompt moves beyond mechanical box placement and instead frames the model as an expert responsible for designing a coherent scene. The prompt guides the model to: (1) imagine a natural background that accommodates all specified objects, (2) prioritize visibility by making the target objects larger and spatially sensible, (3) group similar objects together where appropriate, and (4) ensure realistic placement (e.g., avoiding floating objects, supporting top-down views for crowded settings). In addition, we manually created 15 high-quality in-context examples tailored to E2V, which served as better guidance during layout generation. Original and enhanced prompts are shown in Appendix C.6.

Regeneration. We further introduce an automatic regeneration loop that detects layout errors and re-prompts the model for correction, improving robustness without architectural modifications. We allow up to two regeneration attempts per layout.

6.1.1 Results

Prompt refinement and regeneration lead to substantial improvements in both layout quality and downstream visual accuracy. As shown in Table 2, prompt refinement reduces major layout errors, and adding regeneration further enhances reliability across all error types. These gains are reflected in model performance on the E2V-Bench test set (Table 3), where overall accuracy increases from 12.00% to 14.67% with prompt refinement, and further to 15.33% with the full regeneration pipeline.

6.2 Diffusion and Transformer-based Models

To quantify practical adaptation gains, we evaluate whether supervised fine-tuning on a curated ground-truth set improves a diffusion model (Flux.1-dev) and a transformer-based model (Bagel) without architectural changes.

Training Dataset Construction. We apply a generation-and-filtering pipeline to curate the ground-truth training dataset. Specifically, we

generate 5,920 AWP pairs from equations in the E2V-Bench training split and generate the corresponding visuals using GPT-Image-1, which achieves the highest accuracy in visual generation (Table 1). We then filter the generated images using our automatic metrics and retain only those passing the overall accuracy check, resulting in 1,055 high-quality image-AWP pairs (Table 6).

Supervised Fine-Tuning. We fine-tuned both Flux.1-dev (12B) and BAGEL-MoT (7B) for 10k optimizer steps using supervised fine-tuning. For Flux.1-dev, we employed LoRA adapters for parameter-efficient adaptation, while for Bagel we enabled only the visual generation pathway and froze the VAE using the official training script (Bagel, 2025). Images were generated at 1024×1024 resolution and optimized with AdamW in bf16 precision. Additional implementation details are provided in Appendix C.7 and C.8.

6.2.1 Results

Table 3 reveals divergent trends: fine-tuning improves the overall accuracy of the transformer-based model Bagel (SFT) from 5.67% to 6.00%, but degrades the overall accuracy of the diffusion-based model Flux.1-dev from 6.33% to 5.33%. We hypothesize that Bagel, an LLM-initialized unified multimodal model trained on interleaved data, retains stronger numerical reasoning capabilities than Flux’s purely generative diffusion architecture, consistent with prior findings on the limitations of diffusion models in numerical reasoning (Kajic et al., 2024). Consequently, Bagel appears more robust for visual reasoning.

6.2.2 Further Enhancements to Bagel

To further improve Bagel’s performance, we conducted follow-up fine-tuning with synthetic data augmentation and iterative refinement (Figure 9). We first constructed a synthetic training dataset using template-based scripts and icon assets, drawing equations exclusively from the E2V-Bench training split to avoid data leakage. Because these synthetic images differ substantially from Bagel’s native generation distribution, directly training on them risks emphasizing stylistic features rather than reasoning. We therefore applied Bagel’s image style conversion module to align the synthetic images with Bagel’s generation style and, after filtering with our metrics, retained 3,017 images as an extended training set (details in Appendix C.5). Following prior work (Adarsh et al., 2025), we mixed the synthetic

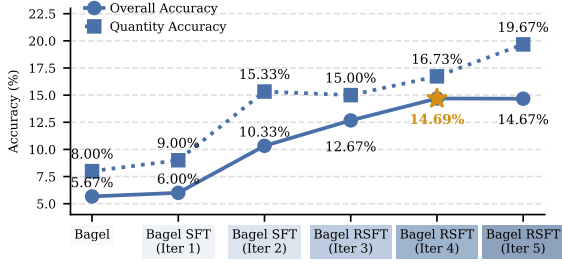


Figure 5: Overall accuracy and quantity accuracy across different training stages of the Bagel model.

and original training data and performed supervised fine-tuning from the previously fine-tuned checkpoint. As shown in Figure 5, this second round of supervised fine-tuning (SFT Iter2) improved overall accuracy from 6.00% to 10.33%.

We further applied rejection-sampling supervised fine-tuning (RSFT), generating ten candidate outputs per training AWP and retaining only samples judged correct by overall accuracy, with at most three visuals preserved per AWP to avoid over-representing easy cases. Each RSFT round was trained for 10k steps, and training stopped once performance no longer improved. As shown in Figure 5, RSFT improves overall accuracy to 12.67% after the first round and peaks at 14.69% after the second. Further RSFT iterations exhibit diminishing returns, potentially due to reduced sample diversity and increased focus on easy cases.

To verify that these gains were not due to additional training iterations, we conducted an ablation where Bagel was only fine-tuned on the original training data without synthetic augmentation. Training to 20k, 30k, and 40k steps yielded overall accuracies of 5.67%, 6.33%, and 7.67%, indicating that the gains stem from synthetic data augmentation and RSFT rather than increased training alone.

6.3 Enhanced Model Analysis and Real-World Validation

Figure 3 and Figure 4 summarize model performance after enhancements. As shown in Figure 3, the enhanced Bagel (Bagel-RSFT) achieves the highest accuracy (50%) in the low-quantity range (1–5 objects), matching GPT-Image-1, but its performance declines rapidly as object counts increase. In contrast, the enhanced LMD (LMD-PR-Regen) attains lower accuracy in low-quantity settings but exhibits more stable performance as quantities grow, suggesting that layout priors help maintain robustness under higher object counts.

Across all visual types (Figure 4), Bagel-RSFT

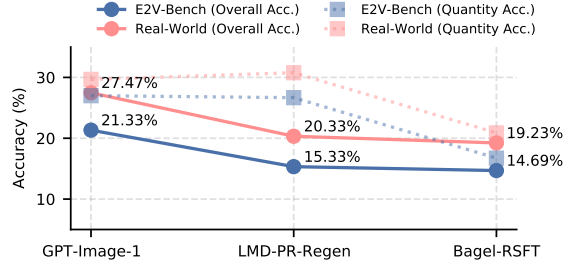


Figure 6: Human evaluation on the real-world educational visual dataset, compared with benchmark results.

improves over the original Bagel, performing best on container-based and second best on color-based visuals. LMD-PR-Regen improves across all visual types, performing best on object-based and second-best on spatial-based visuals, indicating that layout-based priors are effective for generating correct objects and precise spatial arrangements.

To validate performance in real educational settings, we conducted a human evaluation on a manually collected educational dataset of AWP–visual pairs (§ 3.1), sourced from textbooks and educational platforms, comparing GPT-Image-1, LMD-PR-Regen, and Bagel-RSFT (details in App.C.11). Two researchers independently evaluated 546 visuals by comparing them to ground-truth images and assigning binary scores for quantity and overall accuracy. As shown in Figure 6, model performance trends on this real-world dataset mirror those on E2V-Bench, suggesting that the benchmark effectively captures key visual and pedagogical challenges inherent in authentic classroom scenarios.

7 Conclusion

We introduced the equation-to-visual generation task and E2V-Bench, a benchmark to support the study of generating pedagogically meaningful visuals for arithmetic education. Grounded in insights from primary school math teachers and educational resources, E2V-Bench defines four visual types and includes automatic evaluation metrics. Our experiments reveal a substantial performance gap between open and closed-source text-to-image models on this task, highlighting the challenges of educational visual generation. We further show that this gap can be meaningfully reduced through targeted strategies, including rejection sampling fine-tuning, prompt refinement, and regeneration, supported by a curated high-quality training dataset. Together, these contributions establish a foundation for equation-to-visual generation and outline concrete pathways for improving model performance.

635 Limitations

636 **Scope of Representation** Our method is scoped
637 to early arithmetic education (Grades 1–3), focus-
638 ing on the four basic operations and quantities be-
639 low 20. This design choice reflects common class-
640 room practice and enables precise definition of ped-
641 agogical correctness. By constraining the scope,
642 we are able to construct a well-grounded bench-
643 mark and reliable evaluation framework. Nonethe-
644 less, this focus limits direct applicability to more
645 advanced mathematical concepts, such as fractions,
646 ratios, or multi-step equations. Extending equation-
647 to-visual generation beyond early arithmetic re-
648 mains an important direction for future work.

649 **Coverage of Visual Types** Our method mod-
650 els four pedagogically grounded visual types dis-
651 tilled from textbook analysis and validated through
652 teacher feedback. These visual types capture recur-
653 ring and widely used instructional visual patterns in
654 early arithmetic, enabling the task to be formalized
655 and evaluated in a principled manner. However,
656 they do not exhaust the full space of instructional
657 visual representations. Future work could expand
658 the visual taxonomy to cover a broader range of
659 pedagogical designs.

660 **Focus on Foundational Visual Accuracy** Our
661 evaluation focuses on visual correctness, specifi-
662 cally quantity and structural alignment with the in-
663 tended arithmetic concepts, rather than downstream
664 learning outcomes. This emphasis is deliberate: ac-
665 curate visual representation is a prerequisite for
666 any instructional use, and our experiments show
667 that current T2I models still exhibit low overall
668 accuracy on this task (with the strongest models
669 achieving around 20% overall accuracy). Given
670 this gap, establishing reliable generation and eval-
671 uation of correct visuals is a necessary first step
672 before studying higher-level educational effects.
673 Once model accuracy reaches a reliably high level,
674 future work can incorporate human-centered stud-
675 ies to examine learning outcomes and pedagogical
676 impact.

677 References

678 Accessim. 2025. Illustrative Mathematics | v.360 Cur-
679 riculum — accessim.org. <https://accessim.org/>.
680 [Accessed 22-12-2025].

681 Shivam Adarsh, Kumar Shridhar, Caglar Gulcehre,
682 Nicholas Monath, and Mrinmaya Sachan. 2025.

SIKeD: Self-guided iterative knowledge distillation
for mathematical reasoning. In *Findings of the As-
sociation for Computational Linguistics: ACL 2025*,
pages 9868–9880, Vienna, Austria. Association for
Computational Linguistics.

N. Amini-Naieni, T. Han, and A. Zisserman. 2024.
Countgd: Multi-modal open-world counting. In *Ad-
vances in Neural Information Processing Systems
(NeurIPS)*.

Abraham Arcavi. 2003. The role of visual representa-
tions in the learning of mathematics. *Educational
studies in mathematics*, 52(3):215–241.

Bagel. 2025. GitHub - ByteDance-Seed/Bagel: Open-
source unified multimodal model — github.com.
<https://github.com/ByteDance-Seed/Bagel>.
git. [Accessed 20-09-2025].

James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jian-
feng Wang, Linjie Li, Long Ouyang, Juntang Zhuang,
Joyce Lee, Yufei Guo, et al. 2023. Improving image
generation with better captions. *Computer Science*.
<https://cdn.openai.com/papers/dall-e-3.pdf>, 2(3):8.

Lital Binyamin, Yoad Tewel, Hilit Segev, Eran Hirsch,
Royi Rassin, and Gal Chechik. 2025. Make it count:
Text-to-image generation with an accurate number of
objects. In *Proceedings of the Computer Vision
and Pattern Recognition Conference*, pages 13242–
13251.

Blackforest. 2024. black-forest-labs/FLUX.1-
dev · Hugging Face — huggingface.co.
[https://huggingface.co/black-forest-labs/
FLUX.1-dev](https://huggingface.co/black-forest-labs/FLUX.1-dev). [Accessed 12-02-2025].

Anton JH Boonen, Helen C Reed, Judith Schoonen-
boom, and Jelle Jolles. 2016. It’s not a math lesson–
we’re learning to draw! teachers’ use of visual rep-
resentations in instructing word problem solving in
sixth grade of elementary school. *Frontline Learning
Research*, 4(5):55–82.

Marc H Bornstein and Martha E Arterberry. 2010. The
development of object categorization in young chil-
dren: hierarchical inclusiveness, age, perceptual at-
tribute, and group versus individual analyses. *Devel-
opmental psychology*, 46(2):350.

Pu Cao, Feng Zhou, Junyi Ji, Qingye Kong, Zhixiang
Lv, Mingjian Zhang, Xuekun Zhao, Siqu Wu, Yinghui
Lin, Qing Song, et al. 2025a. Preliminary explo-
rations with gpt-4o (mni) native image generation.
arXiv preprint arXiv:2505.05501.

Yuefan Cao, Xuyang Guo, Jiayan Huo, Yingyu Liang,
Zhenmei Shi, Zhao Song, Jiahao Zhang, and Zhen
Zhuang. 2025b. Text-to-image diffusion models can-
not count, and prompt refinement cannot help. *arXiv
preprint arXiv:2503.06884*.

Marianella Casasola, Leslie B Cohen, and Elizabeth
Chiarello. 2003. Six-month-old infants’ categoriza-
tion of containment spatial relations. *Child develop-
ment*, 74(3):679–693.

850	I Loshchilov and F Hutter. 2019. "decoupled weight decay regularization", 7th international conference on learning representations, iclr. <i>New Orleans, LA, USA, May</i> , (6-9):2019.	906
851		907
852		908
853		909
854	Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. 2024. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. In <i>International Conference on Learning Representations (ICLR)</i> .	910
855		911
856		912
857		
858		
859		
860		
861	Haohao Luo, Yang Deng, Ying Shen, See-Kiong Ng, and Tat-Seng Chua. 2024. Chain-of-exemplar: Enhancing distractor generation for multimodal educational question generation . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 7978–7993, Bangkok, Thailand. Association for Computational Linguistics.	916
862		917
863		918
864		919
865		
866		
867		
868		
869	mathematics monster. 2025. Mathematics Monster — mathematics-monster.com. https://www.mathematics-monster.com/ . [Accessed 22-12-2025].	920
870		921
871		922
872		923
873	Richard E. Mayer. 2002. Multimedia learning . volume 41 of <i>Psychology of Learning and Motivation</i> , pages 85–139. Academic Press.	924
874		925
875		
876	Miriam J Metzger, Andrew J Flanagin, and Lara Zwarun. 2003. College student web use, perceptions of information credibility, and verification behavior. <i>Computers & Education</i> , 41(3):271–290.	926
877		927
878		
879		
880	Ministry of General Education and Instruction, Republic of South Sudan. 2018. <i>Primary Mathematics: Pupil's Book 2</i> . Mountain Top Publishers Ltd., Nairobi, Kenya. Funded by the Global Partnership for Education.	928
881		929
882		930
883		931
884		932
885	Wenyi Mo, Tianyu Zhang, Yalong Bai, Bing Su, Ji-Rong Wen, and Qing Yang. 2024. Dynamic prompt optimizing for text-to-image generation. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 26627–26636.	933
886		934
887		
888		
889		
890	Lisa L Morin, Silvana MR Watson, Peggy Hester, and Sharon Raver. 2017. The use of a bar model drawing to teach word problem solving to students with mathematics difficulties. <i>Learning Disability Quarterly</i> , 40(2):91–104.	935
891		936
892		937
893		938
894		939
895	Cherri Moseley and Janet Rees. 2021. <i>Cambridge Primary Mathematics Digital Learner's Book 2</i> , 2nd edition edition. Cambridge University Press, Cambridge, United Kingdom. Digital format.	940
896		941
897		942
898		943
899	OpenAI. 2025. OpenAI Platform — platform.openai.com. https://platform.openai.com/docs/models/gpt-image-1 . [Accessed 24-09-2025].	944
900		945
901		946
902		947
903	opencv. 2025. opencv-python — pypi.org. https://pypi.org/project/opencv-python/ . [Accessed 22-12-2025].	948
904		949
905		
	Sharifah Osman, Che Nurul Azieana Che Yang, Mohd Salleh Abu, Norulhuda Ismail, Hanifah Jambari, and Jeya Amantha Kumar. 2018. Enhancing students' mathematical problem-solving skills through bar model visualisation technique. <i>International Electronic Journal of Mathematics Education</i> , 13(3):273–279.	950
		951
		952
		953
		954
		955
	Prolific. 2025. Prolific Easily collect high-quality data from real people — prolific.com. https://www.prolific.com/ . [Accessed 13-02-2025].	956
		957
		958
	Paul C Quinn, Ramesh S Bhatt, and Angela Hayden. 2008. Young infants readily use proximity to organize visual pattern information. <i>Acta Psychologica</i> , 127(2):289–298.	
	Rahul Ramachandran, Ali Garjani, Roman Bachmann, Andrei Atanov, Oguzhan Fatih Kar, and Amir Zamir. 2025. How well does gpt-4o understand vision? evaluating multimodal foundation models on standard computer vision tasks. <i>arXiv preprint arXiv:2507.01955</i> .	
	Recraft. 2025. Recraft v3. https://www.recraft.ai/docs#models . [Accessed 29-01-2025].	
	Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. 2022. Photorealistic text-to-image diffusion models with deep language understanding. <i>Advances in neural information processing systems</i> , 35:36479–36494.	
	Nazmus Saquib, Rubaiat Habib Kazi, Li-yi Wei, Gloria Mark, and Deb Roy. 2021. Constructing embodied algebra by sketching. In <i>Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems</i> , pages 1–16.	
	scikit learn. 2025. scikit-learn: machine learning in Python &x2014; scikit-learn 0.16.1 documentation — scikit-learn.org. https://scikit-learn.org/ . [Accessed 22-12-2025].	
	Anjali Singh, Ruhi Sharma Mittal, Shubham Atreja, Mourvi Sharma, Seema Nagar, Prasenjit Dey, and Mohit Jain. 2019. Automatic generation of leveled visual assessments for young learners. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 33, pages 9713–9720.	
	Janvijay Singh, Vilém Zouhar, and Mrinmaya Sachan. 2023. Enhancing textbooks with visuals from the web for improved learning . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 11931–11944, Singapore. Association for Computational Linguistics.	
	Marian Small and Amy Lin. 2025. <i>Eyes on math: A visual approach to teaching math concepts</i> . Teachers College Press.	

959 Junling Wang, Anna Rutkiewicz, April Wang, and Mrin- 1012
960 maya Sachan. 2025. [Generating pedagogically mean- 1013
961 ingful visuals for math word problems: A new bench- 1014
962 mark and analysis of text-to-image models](#). In *Find- 1015
963 ings of the Association for Computational Linguistics: 1016
964 ACL 2025*, pages 11229–11257, Vienna, Austria. As- 1017
965 sociation for Computational Linguistics. 1018

966 Ke Wang, Junting Pan, Weikang Shi, Zimu Lu, Mingjie 1019
967 Zhan, and Hongsheng Li. 2024. Measuring mul-
968 timodal mathematical reasoning with math-vision
969 dataset. *arXiv preprint arXiv:2402.14804*.

970 Olivia Wiles, Chuhan Zhang, Isabela Albuquerque,
971 Ivana Kajic, Su Wang, Emanuele Bugliarello, Ya-
972 sumasa Onoe, Pinelopi Papalampidi, Ira Ktena,
973 Christopher Knutsen, Cyrus Rashtchian, Anant
974 Nawalgaria, Jordi Pont-Tuset, and Aida Nematzadeh.
975 2025. [Revisiting text-to-image evaluation with
976 gecko: on metrics, prompts, and human rating](#). In
977 *The Thirteenth International Conference on Learning
978 Representations*.

979 Jinheng Xie, Zhenheng Yang, and Mike Zheng Shou.
980 2025. Show-o2: Improved native unified multimodal
981 models. *arXiv preprint arXiv:2506.15564*.

982 Yi Xu, Roger Smeets, and Rafael Bidarra. 2021. Pro-
983 cedural generation of problems for elementary math
984 education. *International Journal of Serious Games*,
985 8(2):49–66.

986 Shih-Ying Yeh, Yu-Guan Hsieh, Zhidong Gao,
987 Bernard BW Yang, Giyeong Oh, and Yanmin Gong.
988 2023. Navigating text-to-image customization: From
989 lycoris fine-tuning to model evaluation. In *The
990 Twelfth International Conference on Learning Repre-
991 sentations*.

992 Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Lu-
993 ong, Gunjan Baid, Zirui Wang, Vijay Vasudevan,
994 Alexander Ku, Yinfei Yang, Burcu Karagol Ayan,
995 Ben Hutchinson, Wei Han, Zarana Parekh, Xin
996 Li, Han Zhang, Jason Baldridge, and Yonghui Wu.
997 2022. [Scaling autoregressive models for content-rich
998 text-to-image generation](#). *Transactions on Machine
999 Learning Research*. Featured Certification.

1000 A Details of Thematic Analysis of Visual 1001 Types

1002 A.1 Procedure

1003 We conducted a thematic analysis to identify recur-
1004 ring visual types from visuals collected across six
1005 educational sources. In total, we collected 189 visu-
1006 als; seven were decorative and did not correspond
1007 to their associated equations, and were therefore
1008 excluded, leaving 182 visuals. In the initial phase,
1009 we sampled 50 visuals to identify recurring themes.
1010 Through iterative discussion, two researchers iden-
1011 tified and consolidated four distinct visual types.

Coding was performed collaboratively, with a fo-
cus on refining the categories through close reading
and comparison.

In the systematic evaluation phase, the same two
researchers manually analyzed the remaining 132
visuals. This broader analysis allowed us to as-
sess the distribution of problems across the four
identified visual types.

1020 A.2 Results

The results of the systematic evaluation are pre-
sented in Table 4. The distribution of examples
is relatively balanced across spatial-based (54),
object-based (50), container-based (49), and color-
based (29) visual types, suggesting that the identi-
fied visual types can effectively categorize visuals
from educational sources.

Category	Spatial	Object	Container	Color
Count	54	50	49	29

Table 4: Distribution of examples across visual types.

1028 B Dataset Details

1029 B.1 Statistics of E2V-Bench Dataset and 1030 Training Dataset

The statistics of the E2V-Bench dataset are shown
in Table 5, while those of our curated ground-truth
training dataset are presented in Table 6.

Category	Sub-category	Train	Test	Total
Visual Type	Container-based	296	75	371
	Spatial-based	296	75	371
	Object-based	296	75	371
	Color-based	296	75	371
Operation	Addition	320	80	400
	Subtraction	608	152	760
	Multiplication	112	28	140
	Division	144	40	184
Quantity Range	1–5	352	56	408
	6–10	300	100	400
	11–15	264	87	353
	16–20	268	57	325
Total		1184	300	1484

Table 5: Statistics of E2V-Bench dataset.

1034 B.2 Comparison of E2V-Bench with Existing 1035 Benchmarks

We present the comparison of E2V-Bench with re-
lated multimodal benchmarks in Table 7.

Category	Sub-category	Count
Visual Type	Container-based	111
	Spatial-based	271
	Object-based	353
	Color-based	431
Operation	Addition	211
	Subtraction	737
	Multiplication	132
	Division	86
Quantity Range	1-5	628
	6-10	293
	11-15	140
	16-20	105
Total		1166

Table 6: Statistics of curated training dataset.

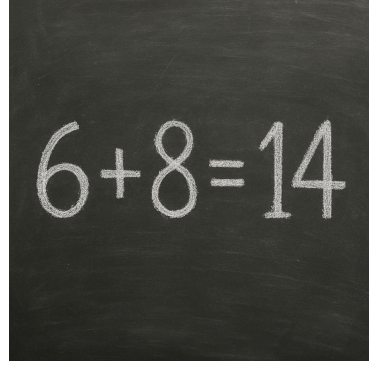


Figure 8: Example of a visual generated by GPT-Image-1 with the prompt: “Create a realistic style image to visualize this equation:6+8=14”

C Experiments Details

C.1 Naive Equation-to-Visual Generation

We first explored a straightforward approach by directly prompting a T2I model to visualize equations. Example outputs are shown in Figure 7 and Figure 8. As illustrated, the generated visuals mainly decorate the equation text with colors or stylistic effects, rather than producing meaningful representations.



Figure 7: Example of a visual generated by GPT-Image-1 with the prompt: “Create a cartoon style image to visualize this equation:3+4=7”

C.2 Automatic Metrics

C.2.1 Quantity Accuracy

For quantity accuracy, we use CountGD (Amini-Naieni et al., 2024) to automatically count object instances in generated visuals.

C.2.2 Overall Accuracy

For overall accuracy, we adopt a rule-based evaluation framework tailored to each visual type, implemented using OpenCV (opencv, 2025) and Scikit-Learn (scikit learn, 2025). Overall accuracy is defined at the image level: a generated visual is

marked as correct only if it satisfies all task-specific structural and semantic constraints implied by the corresponding equation and visual type.

Across all visual types, we first detect objects using a counting-oriented detection model (CountGD (Amini-Naieni et al., 2024)), producing object counts and bounding boxes that serve as shared inputs to downstream evaluations. Based on the annotated visual type, we then apply specialized verification rules, as described below.

Object-based visuals. For object-based visuals, overall accuracy is equivalent to quantity accuracy. A visual is considered correct if the detected number of each object type exactly matches the ground-truth counts derived from the equation.

Color-based visuals. For color-based visuals, we evaluate both count and color assignment. Using object bounding boxes from the counting stage, we extract object crops and estimate dominant colors in HSV space with Gaussian-weighted aggregation. A visual is marked correct only if (1) the number of detected objects matches the ground truth and (2) the multiset of predicted colors matches the expected color distribution.

Spatial-based visuals. For spatial-based visuals, we verify whether objects are correctly grouped according to spatial relations (e.g., left/right or separated groups). We apply a dual clustering strategy that combines spatial proximity and geometric similarity. Specifically, we first perform agglomerative hierarchical clustering with single linkage on object center coordinates to capture spatial adjacency, which is effective for separated or chain-like layouts. If this strategy fails to recover the expected grouping, we additionally apply k-means cluster-

Benchmark	#Entries	Domain	Task Type	Suitability for Arithmetic Education
Evaluating Numerical Reasoning in T2I (Kajic et al., 2024)	1,386	Numbers, objects	T2I generation	Broad coverage of numerical reasoning, but not designed for early arithmetic education.
Math2Visual (Wang et al., 2025)	1,903	Math word problems	Template-based visual generation	Targets math word problem education but lacks visual diversity due to its heavy reliance on templates and icon datasets.
MathVista (Lu et al., 2024)	6,141	Math + vision	Multimodal VQA	Evaluation on reasoning, not visual generation.
MATHVision (Wang et al., 2024)	3,040	General mathematics, competition-level problems	Visual math problem-solving; multimodal reasoning evaluation	Targets middle–high school competition-level math, not early arithmetic education.
E2V-Bench (Ours)	5,738 (1,484 AWP + 4,072 AWP–visual pairs + 182 real educational visuals)	Arithmetic equations, arithmetic word problems	Equation-to-Visual generation	Explicitly targets early arithmetic education with diverse, non-template-based visual generation.

Table 7: Comparison of E2V-Bench with related multimodal benchmarks.

1093	ing over object bounding box geometry (width and	Here are some examples:	1133
1094	height) to distinguish groups based on size and po-	...	1134
1095	sition. A visual is considered correct only if the	Now, for the following equation, generate the	1135
1096	predicted group sizes exactly match the ground-	outputs following the same style:	1136
1097	truth grouping specified by the equation.	Equation: {equation}	1137
1098	Container-based visuals. For container-based	Please respond in the following format:	1138
1099	visuals, we jointly evaluate object–container assign-	AWP: <natural language word problem>	1140
1100	ment and count consistency. Container regions are	object_count: { "<object_name>": <count>, ... }	1141
1101	detected using a grounding-based detector Ground-	The prompt for generating color-based AWP:	1142
1102	ingDINO (Liu et al., 2025), and objects are as-	You are a primary school math teacher. Your task	1144
1103	signed to containers based on spatial overlap and	is to create a simple arithmetic word	1145
1104	relative position. We use GroundingDINO for con-	problem that corresponds to a given	1146
1105	tainer detection because it generalizes better to	arithmetic equation. The arithmetic word	1147
1106	container-like objects and achieves higher accu-	problem should use everyday objects that	1148
1107	racy than CountGD in our preliminary experiments.	young children (Grades 1--3) can easily	1149
1108	A visual is marked correct only if each container	understand.	1150
1109	contains the exact number of objects specified by	Each operand should be represented by identical	1151
1110	the equation.	objects, differentiated only by their colors	1152
1111	C.3 Prompts for Generating AWP	. . .	1153
1112	The prompt for generating object-based AWP:	In addition to the word problem, you must also	1154
1113	You are a primary school math teacher. Your task	provide:	1155
1114	is to create a simple arithmetic word	(1) the total count of the objects, and	1156
1115	problem that corresponds to a given	(2) a structured description of how quantities	1157
1116	arithmetic equation. The arithmetic word	are grouped by color.	1158
1117	problem should use everyday objects that	Here are some examples:	1159
1118	young children (Grades 1--3) can easily	...	1160
1119	understand.	Now, for the following equation, generate the	1161
1120	Each operand in the equation should be	outputs following the same style:	1162
1121	represented by a different type of object to	Equation: {equation}	1163
1122	clearly distinguish the groups of numbers.	Please respond in the following format:	1164
1123	Do not use colors or containers as	AWP: <natural language word problem>	1165
1124	distinguishing features; rely only on object	object_count: { "<object_name>": <total_count> }	1166
1125	types.	visual_structure: {	1167
1126	In addition to the word problem, you must also	" <object_name>": {	1168
1127	provide:	"total": <total_count> ,	1169
1128	(1) the total count of each object type	" <color_1>": <count> ,	1170
1129	mentioned	" <color_2>": <count>	1171
1130		}	1172
1131		}	1173
1132		}	1174
		}	1175
		}	1176
		}	1177
		}	1178
		}	1179
		}	1180

1181	The prompt for generating spatial-based AWP:	"<container1_name>": { "total": <count>, "<object_name>": <count> },	1248
1182	You are a primary school math teacher. Your task	"<container2_name>": { "total": <count>, "<object_name>": <count> }	1249
1183	is to create a simple arithmetic word	}	1250
1184	problem that corresponds to a given		1251
1185	arithmetic equation. The arithmetic word		1252
1186	problem should use everyday objects that		
1187	young children (Grades 1--3) can easily		
1188	understand.		
1189			
1190	Each operand should be represented by placing	C.4 Prompts for Generating Visuals	1253
1191	groups of objects in distinct spatial	The prompt for generating realistic style visuals:	1254
1192	positions (e.g., left/right, near/far).	Create a realistic style image of: <AWP>	1255
1193		The prompt for generating cartoon style visuals:	1256
1194	In addition to the word problem, you must also	Create a cartoon style image of: <AWP>	1257
1195	provide:		
1196	(1) the total count of the objects, and	C.5 Synthetic Data Generation and Style	1258
1197	(2) a structured description of how quantities	Conversion	1259
1198	are grouped spatially.		
1199		Synthetic data generation. We constructed a	1260
1200	Here are some examples:	synthetic dataset using a template-based genera-	1261
1201	...	tion pipeline grounded in arithmetic equations. We	1262
1202		exclusively used equations from the E2V-Bench	1263
1203	Now, for the following equation, generate the	<i>training split</i> to avoid any data leakage.	1264
1204	outputs following the same style:	Given an equation and its annotated visual type,	1265
1205		the pipeline deterministically parses the arithmetic	1266
1206	Equation: {equation}	operands and operator, and instantiates a corre-	1267
1207		sponding visual scene using a library of icon assets.	1268
1208	Please respond in the following format:	This icon library was collected from flaticon (flati-	1269
1209	AWP: <natural language word problem>	con, 2025) and consists of 20 high-quality PNG	1270
1210	object_count: { "<object_name>": <total_count> }	icons representing common objects appearing in	1271
1211	visual_structure: {	E2V-Bench. Each visual is composed of one or	1272
1212	"group1": { "total": <count>, "<object_name>":	more "groups", where each group specifies (i) the	1273
1213	<count> },	object type (icon), (ii) the number of instances,	1274
1214	"group2": { "total": <count>, "<object_name>":	and (iii) optional visual attributes such as color	1275
1215	<count> }	or container membership. A dynamic layout al-	1276
1216	}	gorithm places these groups onto a fixed-size can-	1277
1217		vas, automatically selecting grid dimensions, icon	1278
1218	The prompt for generating container-based	size, padding, and inter-group spacing to ensure	1279
1219	AWP:	balanced composition across a wide range of quan-	1280
1220	You are a primary school math teacher. Your task	tities.	1281
1221	is to create a simple arithmetic word	In parallel, we generate a paired AWP using lan-	1282
1222	problem that corresponds to a given	guage templates aligned with the underlying arith-	1283
1223	arithmetic equation. The arithmetic word	metic operation. For each sample, we also record	1284
1224	problem should use everyday objects that	structured metadata: including object counts, color	1285
1225	young children (Grades 1--3) can easily	assignments, container structure, and spatial group-	1286
1226	understand.	ing, which is serialized into the accompanying CSV	1287
1227		file. To increase diversity, each equation is rendered	1288
1228	Each operand should be represented by placing	multiple times with randomized icon choices, color	1289
1229	objects into one or more containers.	assignments, and layout configurations.	1290
1230		Motivation for style conversion. The result-	1291
1231	In addition to the word problem, you must also	ing synthetic images are visually clean and icon-	1292
1232	provide:	centric, which differs substantially from Bagel's	1293
1233	(1) the total count of each object type, and	native generation distribution. Directly fine-tuning	1294
1234	(2) a structured description of how objects are	Bagel on these images risks overfitting to super-	1295
1235	distributed across containers.	ficial stylistic artifacts rather than improving the	1296
1236	Here are some examples:		
1237	...		
1238			
1239	Now, for the following equation, generate the		
1240	outputs following the same style:		
1241			
1242	Equation: {equation}		
1243			
1244	Please respond in the following format:		
1245	AWP: <natural language word problem>		
1246	object_count: { "<object_name>": <total_count> }		
1247	visual_structure: {		

1297	model’s visual reasoning over quantities and group	Caption: A realistic image of landscape scene	1352
1298	structure. To mitigate this mismatch, we apply	depicting a green car parking on the left of	1353
1299	Bagel’s image-to-image generation module to con-	a blue truck, with a red air balloon and a	1354
1300	vert the synthetic images into Bagel-aligned visual	bird in the sky	1355
1301	styles before using them for training.	Objects: [('a green car', [21, 281, 211, 159]),	1356
		('a blue truck', [269, 283, 209, 160]), ('a	1357
		red air balloon', [66, 8, 145, 135]), ('a	1358
		bird', [296, 42, 143, 100])]	1359
1302	Style conversion with Bagel. We perform style	Background prompt: A realistic landscape scene	1360
1303	conversion using the original Bagel model in an	Negative prompt:	1361
1304	image-conditioned generation setting. Given a		
1305	synthetic icon-based image and its corresponding	Enhanced prompt:	1362
1306	AWP, we prompt the model to regenerate the im-	You are an expert in drawing bounding box for	1363
1307	age in a realistic Bagel style while preserving the	layout-to-image generation.	1364
1308	underlying semantic structure. Specifically, each	When drawing bounding box, you first	1365
1309	input image is passed through Bagel’s vision en-	think about a reasonable scene	1366
1310	coder and VAE, and generation is guided by a text	that can incorporate all objects,	1367
1311	prompt of the form: “Change this image to realistic	describe the scene in text. Then	1368
1312	style: [AWP]”.	draw the bounding box for each	1369
		object, output bounding box and	1370
		the scene prompt.	1371
			1372
1313	Filtering and final dataset. After style conver-	Note:	1373
1314	sion, we filtered the generated images with our	1. You can draw limited amount of	1374
1315	automatic metrics. This process resulted in 3,017	other objects to make the whole	1375
1316	high-quality images, which form the final synthetic	image realistic, but the quantity	1376
1317	extension used in our follow-up fine-tuning exper-	of objects specified in the	1377
1318	iments. By separating structural generation from	prompt should be accurate.	1378
1319	stylistic alignment, this pipeline allows us to inject	2. Bounding box should reflect the	1379
1320	large-scale, logically controlled supervision while	shape of the object, and the	1380
1321	maintaining compatibility with Bagel’s visual gen-	object mentioned in the prompt	1381
1322	eration space.	should be the focus of the image	1382
		and their bounding box should be	1383
1323	C.6 Prompt Refinement Details	BIG for visualization.	1384
1324	We present the original prompt and the enhanced	3. If not specified in the prompt,	1385
1325	prompt for layout generation in LMD.	make sure same type of objects	1386
1326	Original prompt:	are grouping together. If the	1387
1327	You are an intelligent bounding box generator. I	prompt involve same type of	1388
1328	will provide you with a caption for a photo	objects in different color, group	1389
1329	, image, or painting. Your task is to	objects of the same color	1390
1330	generate the bounding boxes for the objects	together.	1391
1331	mentioned in the caption, along with a	4. Please place the bounding boxes in	1392
1332	background prompt describing the scene. The	a natural and spatially sensible	1393
1333	images are of size 512x512. The top-left	way --- for example, objects	1394
1334	corner has coordinate [0, 0]. The bottom-	should not be floating in the air.	1395
1335	right corner has coordinate [512, 512]. The	If there are too many objects,	1396
1336	bounding boxes should not overlap or go	you can use a top-down view as	1397
1337	beyond the image boundaries. Each bounding	indicated in the Background	1398
1338	box should be in the format of (object name,	prompt. Similarly, if the objects	1399
1339	[top-left x coordinate, top-left y	are inside a container, you may	1400
1340	coordinate, box width, box height]) and	also use a top view to make both	1401
1341	should not include more than one object. Do	the container and the objects	1402
1342	not put objects that are already provided in	visible.	1403
1343	the bounding boxes into the background	5. Make sure no bounding box exceeds	1404
1344	prompt. Do not include non-existing or	the image boundary.	1405
1345	excluded objects in the background prompt.		1406
1346	Use “A realistic scene” as the background	Example:	1407
1347	prompt if no background is given in the	1. Input prompt: There are five	1408
1348	prompt. If needed, you can make reasonable	balloons floating in the air. A	1409
1349	guesses. Please refer to the example below	short distance away, there are	1410
1350	for the desired format.	eight green balloons also	1411
1351		floating.	1412
		Output Bounding Box:[('balloon', [8,	1413
		62, 95, 100]), ('balloon', [115,	1414
		62, 96, 102]), ('balloon', [9,	1415
		177, 93, 98]), ('balloon', [118,	1416
		176, 96, 101]), ('balloon', [14,	1417
		293, 97, 97]), ('balloon', [294,	1418
		27, 97, 103]), ('balloon', [403,	1419
		24, 100, 107]), ('balloon', [291,	1420

1421 139, 102, 98]), ('balloon',
 1422 [410, 137, 95, 102]), ('balloon',
 1423 [285, 244, 107, 104]), ('balloon'
 1424 ', [405, 244, 97, 105]), ('
 1425 balloon', [284, 361, 110, 93]),
 1426 ('balloon', [407, 357, 100, 96])]
 1427 Output Background Prompt: A realistic
 1428 scene
 1429 Negative prompt:

1430 C.7 Details of SFT on Flux.1-dev

1431 We fine-tuned Flux.1-dev (12B parameters) for 10
 1432 epochs with 10 repetitions, using a batch size of 5.
 1433 This batch size balanced GPU memory constraints
 1434 (single RTX 4090) and batch diversity, which we
 1435 found to improve convergence stability in prelim-
 1436 inary runs. Gradient checkpointing was enabled
 1437 to reduce memory consumption and allow deeper
 1438 model tuning without sacrificing input resolution.

1439 For optimization, we used the AdamW_BF16
 1440 optimizer (Loshchilov and Hutter, 2019) with an
 1441 initial learning rate of 1e-5. Learning rate was
 1442 decayed with a polynomial scheduler and no warmup,
 1443 which yielded smoother training dynamics com-
 1444 pared to linear decay or cosine schedules. All fine-
 1445 tuning was conducted with BF16 precision.

1446 To enable parameter-efficient fine-tuning, we
 1447 adopted LoRA adapters via the Lycoris frame-
 1448 work (Yeh et al., 2023), updating only the atten-
 1449 tion layers instead of full model weights. Images
 1450 were generated at 1024×1024 resolution to ensure high-
 1451 quality visuals suitable for educational use cases
 1452 while remaining computationally feasible. Fine-
 1453 tuning required approximately 48 hours per run on
 1454 a single NVIDIA RTX 4090 GPU.

1455 C.8 Details of SFT on Bagel

1456 We fine-tuned BAGEL-MoT (7B parameters) for
 1457 10,000 optimizer steps. We used a packed-batch
 1458 scheme of one sequence per rank targeting 10,240
 1459 tokens (flush at 11,520; per-sample cap 10,240)
 1460 to balance GPU memory limits with sequence di-
 1461 versity observed to aid convergence in preliminary
 1462 runs. FSDP and non-reentrant activation check-
 1463 pointing were enabled to reduce memory footprint
 1464 and permit longer effective contexts without lower-
 1465 ing input fidelity.

1466 For optimization, we used AdamW with an initial
 1467 learning rate of 2×10^{-5} ($\beta_1 = 0.9$, $\beta_2 = 0.95$,
 1468 $\epsilon = 1 \times 10^{-15}$) and gradient clipping at 1.0.
 1469 The learning rate followed a constant schedule
 1470 with 2000 warm-up steps. All fine-tuning was
 1471 conducted in bfloat16 precision (bf16) using
 1472 torch.amp.autocast.

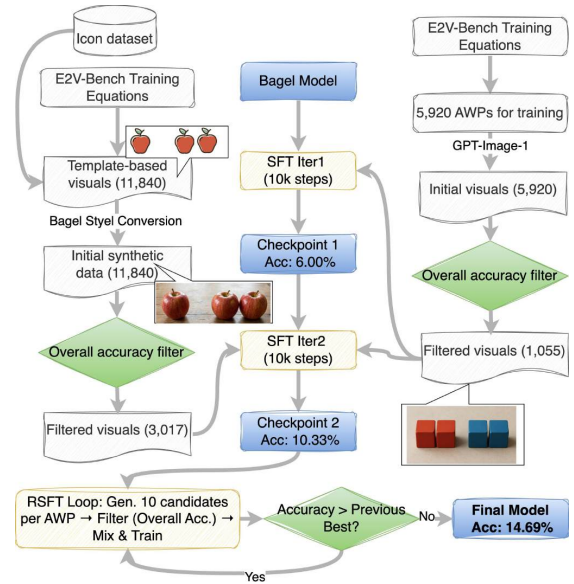


Figure 9: Bagel training pipeline. The process begins with supervised fine-tuning (SFT Iter1) on filtered GPT-Image-1 generated data. Performance is further enhanced through synthetic data augmentation (SFT Iter2) and subsequently refined through iterative rejection-sampling supervised fine-tuning (RSFT).

1473 To limit trainable state, we froze the VAE and
 1474 enabled only the visual generation pathway (i.e.,
 1475 visual understanding disabled). We fine-tuned from
 1476 the HuggingFace BAGEL checkpoint², loading
 1477 EMA weights only before training. Fine-tuning
 1478 required approximately 12 hours per run on a 4
 1479 NVIDIA GH200 GPUs.

1480 C.9 Details of Bagel Training

1481 We present the training pipeline of Bagel model in
 1482 Figure 9.

1483 C.10 Additional Results

1484 Figures 10 and 11 present the overall accuracy of
 1485 cartoon style visuals generated by different models,
 1486 evaluated across varying quantity ranges and visual
 1487 types, respectively.

²<https://huggingface.co/ByteDance-Seed/BAGEL-7B-MoT>

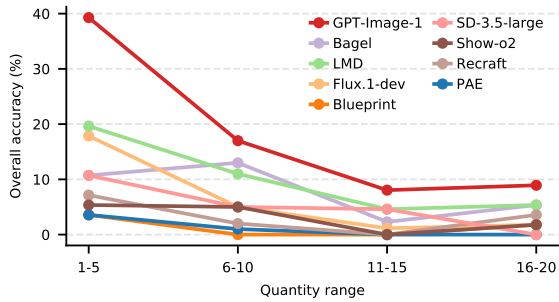


Figure 10: Overall accuracy of cartoon style visuals evaluated across different quantity ranges. Corresponding figure for realistic style visuals are in Figure 3

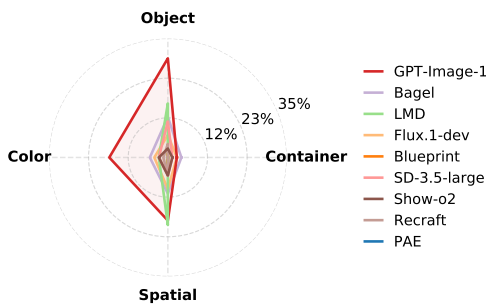


Figure 11: Overall accuracy of cartoon style visuals evaluated across different visual types. Corresponding figure for realistic style visuals are in Figure 4.

C.11 Human Evaluation Details

We conducted a human evaluation on a manually curated educational dataset of AWP–visual pairs, as described in Section 3.1. The dataset was sourced from primary-school mathematics textbooks and online educational platforms.

During manual curation, we collected visuals from arithmetic education sections covering the four basic operations that were accompanied by either an AWP or a symbolic equation. For visuals that were associated only with equations and lacked an explicit AWP, we manually authored a descriptive AWP consistent with the visual content. Among the 182 collected visuals, 52 had corresponding AWPs, while the remaining visuals were originally paired only with equations.

During human evaluation, AWPs were used as input prompts to the models. Two researchers independently evaluated a total of 546 generated visuals by comparing them against the ground-truth images and assigning binary judgments for Quantity Accuracy and Overall Accuracy. A total of three annotation disagreements were identified and resolved through discussion. A screenshot of the annotation interface used in the human evaluation is provided in Figure 12.

C.12 Example Visuals from T2I Models

We present example visuals generated by T2I models on E2V-Bench in Figure 13.

D Teacher Interview Study Details

D.1 Teacher Demographics

We recruited ten primary school math teachers through Prolific (Prolific, 2025) and paid them 16.63 USD per hour, which is adequate given the participants’ country of residence. We present the participants’ demographics in Table 8. All participating teachers use English as the language of instruction.

PID	Teaching Experience	Age	Gender
1	More than 10 years	46	Male
2	6–10 years	36	Female
3	1–2 years	28	Male
4	More than 10 years	54	Female
5	More than 10 years	34	Female
6	More than 10 years	51	Female
7	More than 10 years	53	Male
8	More than 10 years	42	Female
9	6–10 years	27	Female
10	More than 10 years	53	Female

Table 8: Participants’ demographics including teaching experience, age, and gender.

D.2 Study Procedure

Our study obtained ethical approval, and written consent was collected from each participant. The study lasted between 45 minutes and one hour. Participants were first introduced to the background of the study and then completed four sessions, as described below.

In the first session, participants were introduced to the four visual types, accompanied by example visuals for the four arithmetic operations (+, −, ×, ÷). After reviewing each design, they completed a questionnaire to provide feedback.

In the second session, participants were presented with 32 arithmetic problems randomly selected from the E2V-Bench test set. They then completed a questionnaire rating the usefulness of these problems for teaching arithmetic skills.

In the third session, participants were shown examples of realistic and cartoon style visuals. They were asked to indicate their preference between the two styles and explain their reasons.

In the fourth session, participants were introduced to the definitions of our proposed evaluation criteria and asked to provide feedback on them.

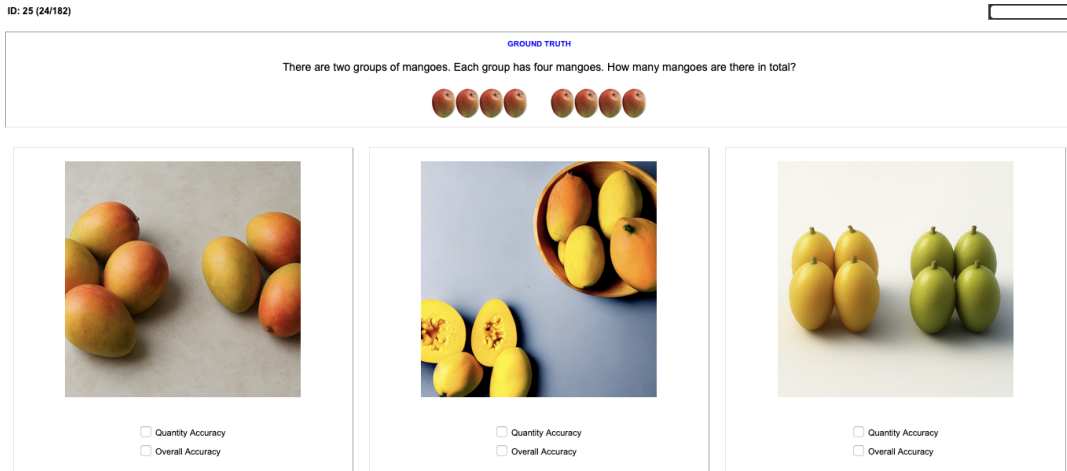


Figure 12: Screenshot of the system used for human evaluation. Annotators assign binary scores for quantity accuracy and overall accuracy: clicking the button indicates a correct judgment, while leaving it unclicked indicates an incorrect one. Visuals are shown in random order for each question.

Finally, after completing all four sessions, participants were asked about their general experiences using visuals to teach arithmetic skills and the types of visuals they typically employ in their classrooms.

D.3 Additional Results

D.3.1 Teachers' Evaluation of Visual Usefulness for Teaching

We present teachers' ratings of the usefulness of images following our proposed four visual types for teaching math equations in Table 9. Scores were collected on a seven-point Likert scale, with seven representing the most positive rating. Teachers consistently gave high ratings across all three statements, with averages close to seven (6.7-6.8) and low variance (0.18-0.46). These results indicate that, from the teachers' perspective, visuals generated with our proposed visual types are useful for teaching and support students' understanding of math equations. Teachers also indicated a strong willingness to use them in their own instruction.

Statement	Average	Variance
Useful for teaching	6.80	0.18
Helpful for student understanding	6.70	0.23
Would like to use in teaching	6.70	0.46

Table 9: Teacher ratings of the usefulness of visuals following the four visual types for teaching arithmetic skills. Scores are on a Likert scale from 1 to 7, where 7 indicates strongly agree.

D.3.2 Teachers' Evaluation of Generated Arithmetic Problems for Teaching

We present teachers' ratings of generated arithmetic problems for teaching in Table 10. Scores were collected on a seven-point Likert scale, with seven representing the most positive rating. Teachers consistently gave high ratings across all three statements, with averages close to seven (6.7-6.9) and low variance (0.10-0.23). These results indicate that, from the teachers' perspective, our generated arithmetic problems are useful for teaching and support students' understanding of math equations. Teachers also indicated a strong willingness to use them in their own instruction.

Statement	Average	Variance
Useful for teaching	6.70	0.23
Helpful for student understanding	6.70	0.23
Would like to use in teaching	6.90	0.10

Table 10: Teacher ratings of the usefulness of generated arithmetic problems following the four visual types for teaching arithmetic skills. Scores are on a Likert scale from 1 to 7, where 7 indicates strongly agree.

D.3.3 Teachers' Preference for Realistic and Cartoon Style Visuals

We present teachers' ratings of generated visuals in realistic and cartoon styles for teaching arithmetic skills in Table 11. Scores were collected on a seven-point Likert scale, with seven representing the most positive rating. Teachers gave high ratings for both styles, with the realistic style (6.3) rated slightly higher than the cartoon style (6.2), and variance

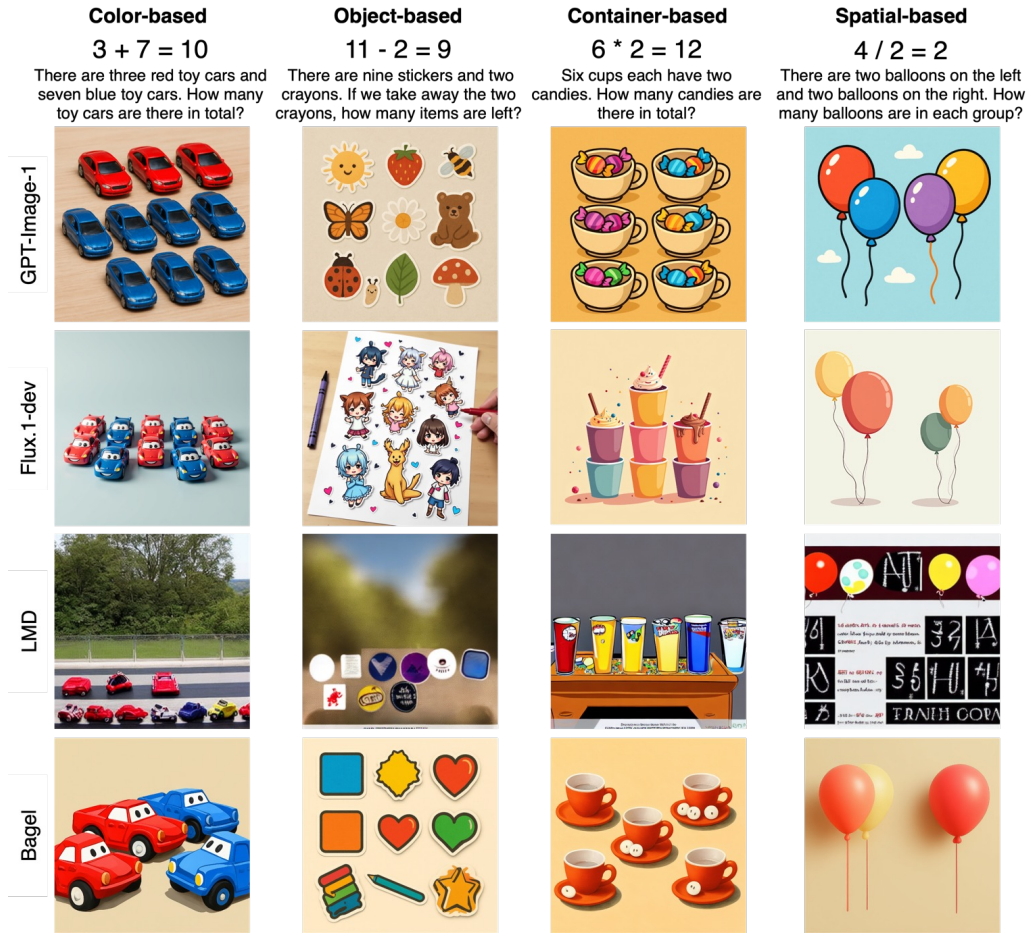


Figure 13: Example visuals from four T2I models on the E2V-Bench task. Each column contains one of the four identified structures (Color-based, Object-based, Container-based, Spatial-based), with the input equation and arithmetic word problem (AWP) description at the top. Each row shows outputs from a different model.

1594 remaining relatively low (0.84–0.90). Out of ten
 1595 teachers, four gave higher scores to the realistic
 1596 style. Six teachers also commented that realistic
 1597 visuals are useful because students can more easily
 1598 connect them with familiar objects from daily life
 1599 (P1–4, P6–7), which helps students better under-
 1600 stand the context of the problems. Three teachers
 1601 gave higher scores to the cartoon style. Two teach-
 1602 ers noted that cartoon visuals are commonly used
 1603 in their teaching (P5, P10), while two teachers em-
 1604 phasized that the simplified design of cartoon style
 1605 can reduce student distractions by using plain back-
 1606 grounds (P10) and identical objects (P7). These
 1607 results suggest that both realistic and cartoon visu-
 1608 als are perceived as useful for teaching arithmetic
 1609 skills, with realistic visuals valued for their con-
 1610 nection to students’ familiar real-world objects and
 1611 cartoon visuals appreciated for their clarity and
 1612 familiarity in educational contexts.

Statement	Average	Variance
Realistic-style visuals are useful for teaching arithmetic skills.	6.30	0.90
Cartoon-style visuals are useful for teaching arithmetic skills.	6.20	0.84

Table 11: Teacher ratings of the usefulness of generated visuals in realistic and cartoon styles for teaching arithmetic skills. Scores are on a Likert scale from 1 to 7, where 7 indicate most useful.

D.3.4 Teacher’s Evaluation of Visuals Across Visual Types

We present teachers’ ratings of the usefulness of four visual types for teaching arithmetic operations in Table 12. Scores were collected on a seven-point Likert scale, with seven representing the most positive rating.

For addition, teachers rated color-based (6.6) and container-based (6.3) visuals as most useful. Seven teachers noted that color-based visuals effectively illustrate the idea of adding new objects into an

existing group, while seven teachers emphasized that container-based visuals help clearly separate different groups, making it easier to highlight the part being added.

For subtraction, container-based visuals were rated highest (6.1), followed by spatial-based (5.7). Five teachers mentioned that container-based visuals make the subtraction process easier for students to understand. Three teachers also noted that combining visuals with narrative would be helpful for teaching the subtraction process.

For multiplication and division, container-based visuals received the highest ratings (7.0 for multiplication and 6.9 for division), with spatial-based visuals also considered useful (6.1 for multiplication and 6.4 for division). Seven teachers commented that containers are the most effective way to represent groups in multiplication and division, while four teachers mentioned spatial-based visuals are helpful for showing how items can be grouped or partitioned.

Overall, these results suggest that teachers' evaluations of visual types vary by operation. However, container-based visuals were consistently rated highly across all four operations, with particular strength for multiplication and division.

Operation	Color	Object	Spatial	Container
Addition	6.6	5.3	5.3	6.3
Subtraction	5.5	5.2	5.7	6.1
Multiplication	5.5	4.1	6.1	7.0
Division	5.6	5.1	6.4	6.9

Table 12: Teacher ratings of the usefulness of four visual types (color-based, object-based, spatial-based, container-based) for teaching different arithmetic operations. Scores are on a Likert scale from 1 to 7, where 7 indicates most useful.

D.3.5 Coverage of Visual Types and Teachers' Current Teaching Visuals

As shown in Table 13, teachers expressed very strong agreement with the statement that the four visual types capture the kinds of visuals they currently use when teaching arithmetic skills in Grades 1–3 (average = 6.90). The low variance (0.10) indicates a high level of consistency among teachers, suggesting that the proposed visual types closely align with existing classroom practices.

All ten teachers reported that they have used visuals similar to our designs following the four visual types. Six teachers mentioned using cartoon-style visuals similar to our designs, while the other

four indicated that they have used realistic-style visuals similar to our designs in their classroom teaching.

Statement	Average	Variance
The four visual types cover the kinds of visuals used to teach arithmetic skills in Grades 1–3.	6.90	0.10

Table 13: Teachers' ratings of the extent to which the four visual types cover visuals used in teaching arithmetic skills in Grades 1–3. Scores are on a Likert scale from 1 to 7, where 7 indicates strongly agree.

E Ethical Consideration and Applications

E.1 Potential Risks

One potential risk is that generated visuals may be misinterpreted if they do not accurately capture the intended mathematical relationships, potentially leading to confusion among students and educators. To mitigate this risk, we collaborated closely with primary school mathematics teachers to develop a structured design space aligned with pedagogical standards. The automatic metrics we propose further help ensure the correctness of the curated dataset. However, this dataset should not be used directly for educational purposes without the supervision of educators.

E.2 Terms of Use

This section outlines the terms and conditions for the use of E2V-Bench. By using the code and datasets in this project, users agree to the following terms:

Prohibited Use The code and datasets shall not be used for commercial purposes without prior written consent from the authors.

Attribution When using or referencing the code and datasets, users must provide proper attribution to the original authors.

No Warranty This project is provided as is without any warranties of any kind, either expressed or implied, including but not limited to fitness for a particular purpose. The authors are not responsible for any damage or loss resulting from the use of this project.

Liability The authors shall not be held liable for any direct, indirect, incidental, special, exemplary,

1700	or consequential damages arising in any way out	Users are encouraged to preserve the pedagogical	1749
1701	of the use of the E2V-Bench.	intent of the visual types and evaluation criteria	1750
1702	Updates and Changes The authors reserve the	when extending or applying the benchmark.	1751
1703	right to make changes to the terms of this license	Curated Training Dataset for Model Enhance-	1752
1704	or the E2V-Bench itself at any time.	ment	1753
1705	E.3 Compliance with Artifact Usage and	Intended Use: The curated training dataset is in-	1754
1706	Intended Use Specifications	tended to support research on improving T2I mod-	1755
1707	E.3.1 Compliance with Existing Artifact	els for equation-to-visual generation. It provides	1756
1708	Usage	high-quality supervision for studying model adap-	1757
1709	In our study, we utilized a range of existing arti-	tation techniques.	1758
1710	facts, such as icons from Flaticon (flaticon, 2025)	Restrictions: The dataset consists of automat-	1759
1711	and visuals from six educational sources (Ministry	ically generated and filtered image–text pairs and	1760
1712	of General Education and Instruction, Republic	may still contain inaccuracies or biases inherent	1761
1713	of South Sudan, 2018; Cotton et al., 2021; Mose-	to current generative models. It is therefore not	1762
1714	ley and Rees, 2021; Accessim, 2025; mathematics	recommended for direct instructional use or de-	1763
1715	monster, 2025; fun2dolabs, 2025), to develop our	ployment in educational products without expert	1764
1716	datasets. We rigorously ensured that our usage of	oversight.	1765
1717	these materials was in strict accordance with their	Data Ethics: All data are synthetically gener-	1766
1718	intended purposes. Additionally, we employed var-	ated or derived from open educational resources	1767
1719	ious computational tools within their prescribed	and do not contain personally identifiable informa-	1768
1720	licensing terms, thus adhering to ethical and legal	tion. The dataset is released for research purposes,	1769
1721	standards.	and users are encouraged to adhere to responsible	1770
1722	E.3.2 Specification of Intended Use for	data usage practices consistent with educational	1771
1723	Created Artifacts	and academic norms.	1772
1724	Our research resulted in two primary artifacts: (1)	E.4 Data Collection and Anonymization	1773
1725	the E2V-Bench benchmark and evaluation frame-	Procedures	1774
1726	work, and (2) a curated training dataset for T2I	The benchmark and training datasets do not con-	1775
1727	model enhancement.	tain personal data. Arithmetic equations and corre-	1776
1728	E2V-Bench: Benchmark and Evaluation Frame-	sponding word problems were generated by mod-	1777
1729	work	els, and all visuals were produced using T2I mod-	1778
1730	Intended Use: E2V-Bench is intended for aca-	els.	1779
1731	demetic research on multimodal learning, text-to-	For components involving human participation,	1780
1732	image generation, and educational AI. It supports	including teacher interviews and evaluations, no	1781
1733	systematic evaluation of models’ ability to gener-	personally identifiable information was collected	1782
1734	ate pedagogically meaningful visual representations	or retained. All feedback was anonymized prior to	1783
1735	from arithmetic equations, and facilitates the de-	analysis. In addition, we conducted manual screen-	1784
1736	velopment and comparison of model adaptation	ing to exclude offensive, sensitive, or inappropriate	1785
1737	strategies.	content. These procedures were adopted to ensure	1786
1738	Restrictions: E2V-Bench is designed as a re-	ethical data handling, participant privacy, and re-	1787
1739	search benchmark and should not be used as a stan-	sponsible research practice.	1788
1740	dalone instructional system or deployed directly	E.5 Artifact Documentation	1789
1741	in classroom settings without additional validation.	E.5.1 E2V-Bench Benchmark	1790
1742	High-stakes educational or commercial use is dis-	Domain Coverage E2V-Bench targets early	1791
1743	couraged unless supported by further empirical	arithmetic education (Grades 1–3), focusing on	1792
1744	studies and ethical review.	equation-to-visual generation for foundational	1793
1745	Ethical Considerations: The benchmark de-	arithmetic concepts.	1794
1746	sign is grounded in analyses of publicly available	Operation Coverage The benchmark covers	1795
1747	educational materials and validated through consul-	four basic arithmetic operations: addition, subtrac-	1796
1748	tations with primary school mathematics teachers.	tion, multiplication, and division, with quantities	1797

1798	restricted to values suitable for primary school in-		
1799	struction.		
1800	Visual Types Each equation is associated		
1801	with four pedagogically grounded visual types:		
1802	container-based, object-based, color-based, and		
1803	spatial-based representations.		
1804	E.5.2 Curated Training Dataset		
1805	Visual Style The dataset includes realistic style		
1806	visuals.		
1807	Content Scope All examples are derived from		
1808	arithmetic equations and corresponding arithmetic		
1809	word problems generated to support equation-to-		
1810	visual learning.		
1811	Demographic Representation The dataset does		
1812	not encode demographic attributes. Its educational		
1813	scope reflects arithmetic instruction practices rather		
1814	than representations of specific populations.		
1815	E.6 Use of AI Assistants in Research		
1816	In our study, AI assistants were used sparingly and		
1817	in accordance with ACL’s Policy on AI Writing		
1818	Assistance. We utilized ChatGPT and Grammarly		
1819	for basic paraphrasing and grammar checks, re-		
1820	spectively. These tools were applied minimally		
1821	to ensure the authenticity of our work and to ad-		
1822	here strictly to the regulatory standards set by ACL.		
1823	Our use of these AI tools was focused, responsible,		
1824	and aimed at supplementing rather than replacing		
1825	human input and expertise in our research process.		
1826	E.7 Instructions Given To Participants		
1827	E.7.1 Disclaimer for Annotators		
1828	Thank you for participating in our evaluation pro-		
1829	cess. Please read the following important points		
1830	before you begin:		
1831	• Voluntary Participation: Your participation		
1832	is completely voluntary. You have the free-		
1833	dom to withdraw from the task at any time		
1834	without any consequences.		
1835	• Confidentiality: All data you will be work-		
1836	ing with is anonymized and does not contain		
1837	any personal information. Your responses and		
1838	scores will also be kept confidential.		
1839	• Risk Disclaimer: This task does not involve		
1840	any significant risks. It primarily consists of		
1841	reading and scoring generated visuals.		
	• Queries: If you have any questions or con-		1842
	cerns during the task, please feel free to reach		1843
	out to us.		1844
	E.7.2 Instructions for Experiments		1845
	Thank you for participating in our study. This re-		1846
	search received ethical approval, and informed con-		1847
	sent was obtained from all participants. The study		1848
	took approximately one hour and consisted of four		1849
	tasks. Please read the instructions below carefully.		1850
	Task 1: Visual Type Evaluation You will be pre-		1851
	sented with visuals representing arithmetic equa-		1852
	tions across four visual types (container-based,		1853
	object-based, color-based, and spatial-based) and		1854
	the four basic arithmetic operations. For each set of		1855
	visuals, please complete a questionnaire evaluating		1856
	its usefulness for teaching arithmetic concepts in		1857
	the classroom.		1858
	Task 2: Arithmetic Word Problem Evaluation		1859
	You will review 32 automatically generated arith-		1860
	metic word problems corresponding to equations		1861
	in our benchmark. Please rate whether each word		1862
	problem is suitable for classroom use and align-		1863
	ment with instructional practice.		1864
	Task 3: Visual Style Comparison You will eval-		1865
	uate two visual styles: realistic and cartoon, for		1866
	visualizing arithmetic equations. Please provide		1867
	feedback on how effective each style may be for		1868
	educational use, including perceived benefits and		1869
	potential drawbacks in classroom settings.		1870
	Task 4: Evaluation Criteria Feedback We will		1871
	present the definitions of our proposed automatic		1872
	evaluation metrics for generated visuals. Please		1873
	comment on whether these criteria adequately cap-		1874
	ture visual usefulness for teaching arithmetic and		1875
	provide any suggestions for improvement.		1876
	Please answer all questions honestly and feel		1877
	free to share additional feedback throughout the		1878
	study. Your responses will help validate our design		1879
	choices and evaluation framework. Thank you for		1880
	your time and valuable input.		1881
	E.7.3 Data Consent		1882
	The data you provide during this study will be used		1883
	solely for academic research purposes. All informa-		1884
	tion will be anonymized and securely stored, and		1885
	any published or shared data will be aggregated to		1886
	ensure your privacy. By participating, you agree to		1887
	the use of your data as described, but you retain the		1888

1889 right to withdraw your consent at any time with-
1890 out penalty. If you have any questions about how
1891 your data will be used, please feel free to ask the
1892 research team.