

Omni-I2C: A Holistic Benchmark for High-Fidelity Image-to-Code Generation

Jiawei Zhou^{1,*} Chi Zhang^{1,*} Xiang Feng¹ Qiming Zhang² Haibo Qiu²
Lihuo He^{3,†} Dengpan Ye^{4,†} Xinbo Gao³ Jing Zhang^{1,†}

¹School of Computer Science, Wuhan University, China

²Independent Researcher, China

³Xidian University, China

⁴Cyberspace Institute of Advanced Technology, Guangzhou University, China

†Corresponding author: jingzhang.cv@gmail.com, lhhe@mail.xidian.edu.cn, yedp@gzhu.edu.cn

Abstract

We present Omni-I2C, a comprehensive benchmark designed to evaluate the capability of Large Multimodal Models (LMMs) in converting complex, structured digital graphics into executable code. We argue that this task represents a non-trivial challenge for the current generation of LMMs: it demands an unprecedented synergy between high-fidelity visual perception—to parse intricate spatial hierarchies and symbolic details—and precise generative expression—to synthesize syntactically sound and logically consistent code. Unlike traditional descriptive tasks, Omni-I2C requires a holistic understanding where any minor perceptual hallucination or coding error leads to a complete failure in visual reconstruction. Omni-I2C features 1130 meticulously curated samples, defined by its breadth across subjects, image modalities, and programming languages. By incorporating authentic user-sourced cases, the benchmark spans a vast spectrum of digital content—from scientific visualizations to complex symbolic notations—each paired with executable reference code. To complement this diversity, our evaluation framework provides necessary depth; by decoupling performance into perceptual fidelity and symbolic precision, it transcends surface-level accuracy to expose the granular structural failures and reasoning bottlenecks of current LMMs. Our evaluation reveals a substantial performance gap among leading LMMs; even state-of-the-art models struggle to preserve structural integrity in complex scenarios, underscoring that multimodal code generation remains a formidable challenge. Data and code are available at <https://github.com/MiliLab/Omni-I2C>.

1 Introduction

In professional workflows, programming is rarely a requirement isolated from visual cues. Whether it

*Equal contribution.

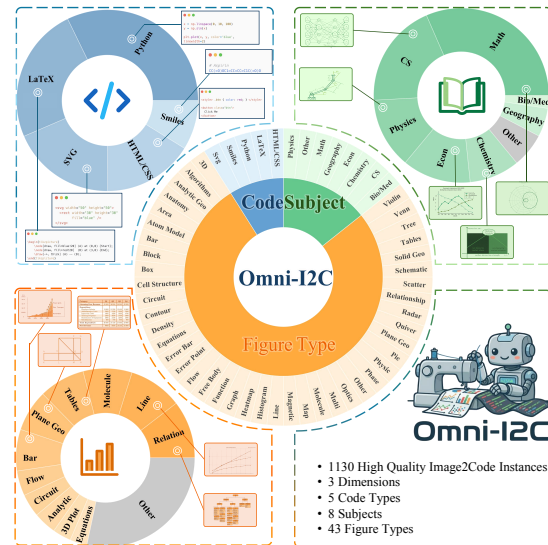


Figure 1: Taxonomy overview of the Omni-I2C benchmark. The dataset features a diverse distribution across 5 code types, 8 major subjects, and 43 distinct figure types.

is an engineer re-implementing complex data interfaces from legacy archives, a developer transforming high-fidelity visual mockups into structured layouts, or a researcher converting technical diagrams and notations into executable logic. The demand for Image-to-Code (I2C) conversion is both ubiquitous and valuable (Yang et al., 2024; Si et al., 2025; Wu et al., 2025; Li et al., 2024). This “pixel-to-program” synthesis represents an advanced application of human-AI collaboration, where LMMs act as intelligent assistants that bridge the gap between visual information and executable functional logic.

Beyond its practical utility, I2C generation offers a unique paradigm for evaluating high-fidelity perception compared to conventional tasks such as Visual Question Answering (VQA) (Yue et al., 2024, 2025; Lu et al., 2023; Antol et al., 2015). While VQA evaluates a model’s ability to extract specific information through targeted queries, I2C generation demands a holistic, structural reconstruction

Table 1: Compared with prior benchmarks, Omni-I2C achieves the greatest breadth in coverage.

| Benchmark | Domain | #Samples | #Subjects | #Figure Types | Output Modalities |
|-------------------------------------|-------------------|--------------|-----------|---------------|----------------------------------|
| ChartMimic (Yang et al., 2024) | Chart | 4,800 | 8 | 22 | Python |
| Chart2Code (Tang et al., 2025) | Chart | 2,023 | 4 | 22 | Python |
| Plot2Code (Wu et al., 2025) | Plot | 368 | - | 6 | Python, R |
| Design2Code (Si et al., 2025) | Web | 484 | - | - | HTML |
| UniSVG (Li et al., 2025) | Icon | 2,850 | - | - | SVG |
| Image2Struct (Roberts et al., 2024) | Doc/Music | 2,400 | 3 | 6 | LaTeX, HTML |
| Omni-I2C | Generalist | 1,130 | 8 | 43 | Python, SVG, HTML, LaTeX, SMILES |

of the entire visual input. It is intrinsically less tolerant of perceptual discrepancies; even a minute error—such as misidentifying a LaTeX subscript or a specific coordinate—will cause the rendered output to deviate, fundamentally altering the semantic meaning of the result. This execution-based verifiability enables a more comprehensive, fine-grained, and objective evaluation of the perceptual capabilities of a model.

Recognizing this potential, recent benchmarking efforts begin to explore I2C generation within specialized domains. For instance, ChartMimic (Yang et al., 2024) and Plot2Code (Wu et al., 2025) focus on synthesizing plotting scripts from scientific charts, while Design2Code (Si et al., 2025) evaluates the transformation of web screenshots into front-end implementations. Although these works have established a vital foundation, current evaluations remain largely restricted in image variety, programming languages, and task scenarios. To comprehensively assess the limits of perception, coding and reasoning, a more diverse and general-purpose testbed is essential to evaluate whether LMMs can generalize across varied technical and professional domains.

To bridge this gap, we introduce Omni-I2C, a comprehensive benchmark for stress-testing the limits of Image-to-Code (I2C) generation. The dataset contains 1130 meticulously curated samples drawn from real-world user applications, covering a broad spectrum of technical graphics such as scientific visualizations, molecular structures, and complex symbolic notations, and supporting diverse programming languages. Furthermore, we propose a multi-level evaluation framework that decomposes model performance into fine-grained attributes, enabling a more nuanced understanding of model strengths and weaknesses across different dimensions. Unlike previous metrics restricted to specific domains or coding formats, our evaluation logic offers comprehensive coverage and demonstrates superior alignment with human judgment through extensive empirical validation.

Through an extensive evaluation of 13 proprietary and open-weight LMMs, we reveal a profound performance gap in high-fidelity image-to-code generation. Even leading frontier models, such as Gemini 3 Pro and GPT-5.1, frequently falter in the challenging scenarios presented by our benchmark. These results highlight substantial room for improvement and position Omni-I2C as a challenging benchmark for advancing LMMs. Our contributions are summarized as follows:

- We present Omni-I2C, a meticulously curated dataset of 1130 items, including 5 programming languages, 8 major subjects, and 43 distinct figure types. It serves as a rigorous testbed for evaluating the perception and coding capabilities of LMMs.
- We propose an evaluation framework that assesses code-level integrity and image-level perceptual consistency, enabling more diagnostic and attributable analyses of model behavior than traditional heuristic metrics.
- Our comprehensive analysis of SOTA LMMs exposes a significant performance gap in high-fidelity reconstruction, identifying critical failure modes and charting a path toward more precise, trusted multimodal agents.

2 Related Work

2.1 Multimodal Code Benchmarks

The rapid advancement of Large Language Models (LLMs) has reshaped automated code generation, from general-purpose foundation models (OpenAI, 2025; Deepmind, 2025; Anthropic, 2025) to specialized architectures (Guo et al., 2024; Hui et al., 2024) that excel at complex programming tasks (Jimenez et al., 2023; Zhou et al., 2023; Yao et al., 2024). Model evaluation, however, has largely remained single-modal, relying on text-only benchmarks such as HumanEval (Chen, 2021), MBPP (Austin et al., 2021), and DS-1000 (Lai

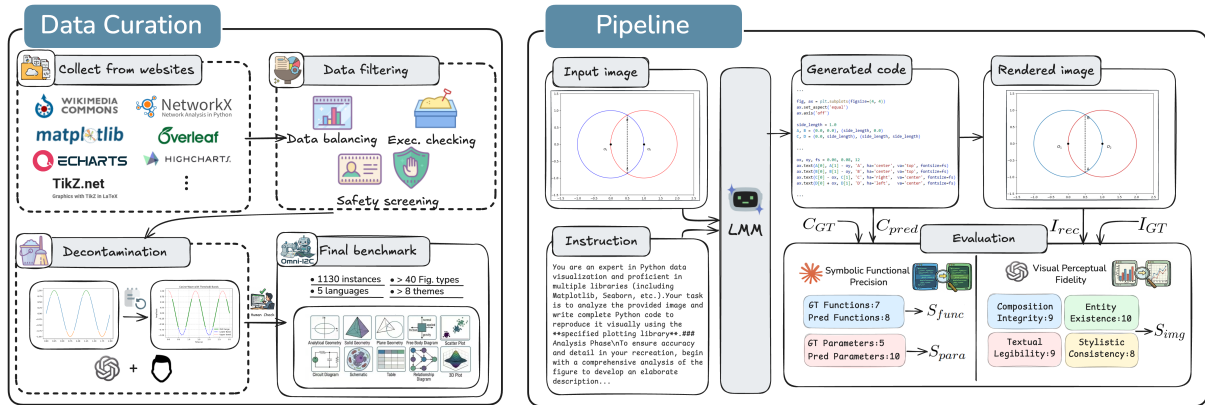


Figure 2: The Data Curation and pipeline of Omni-I2C. **Data Curation** (Left): From raw web collection to a refined benchmark with diverse themes and languages, ensured by strict filtering and human verification. **Pipeline** (Right): The proposed workflow for image-to-code generation and automatic evaluation.

et al., 2023). As models become multimodal, evaluation must expand to real-world settings where visual perception is integral.

Recent multimodal benchmarks—MMCode (Li et al., 2024), Design2Code (Si et al., 2025), Plot2Code (Wu et al., 2025), and ChartMimic (Yang et al., 2024)—take steps in this direction by testing HTML generation and chart-to-code synthesis. Yet they are typically domain-bound, with limited diversity in subjects, visual modalities, and programming languages. To address these limitations, we introduce Omni-I2C, a comprehensive benchmark for high-fidelity visual understanding and image-to-code generation. See Tab. 1 for detailed comparisons with prior benchmarks.

2.2 Visual Perception in LMMs

The evolution of Large Multimodal Models has transitioned from coarse recognition to complex, fine-grained visual reasoning, bolstered by high-quality data and test-time scaling (Alayrac et al., 2022; Liu et al., 2023; Chen et al., 2024; Wang et al., 2024; Lu et al., 2023; Yue et al., 2024; Qiao et al., 2025; Zhang et al., 2025c,b; Zheng et al., 2025; Meng et al., 2025). However, existing VQA benchmarks—often relying on multiple-choice or short-form answers—suffer from sparse or ambiguous supervision (Zhang et al., 2025a). Such formats frequently allow models to exploit language priors or “shortcuts” without true visual grounding, making failures difficult to diagnose. Image-to-code generation mitigates this by introducing dense, execution-based supervision through “pixel-to-program” mapping—where perceptual mistakes translate into functional reconstruction errors—allowing Omni-I2C to assess perceptual fidelity more granularly than traditional VQA.

3 The Omni-I2C Benchmark

3.1 Task Definition

The Omni-I2C benchmark probes cross-modal intelligence by requiring models to synthesize executable code that reconstructs complex digital graphics. Formally, given an input image I and an instruction T specifying the target programming language and reconstruction goal, an LMM f_θ generates code C_{pred} :

$$C_{pred} = f_\theta(I, T).$$

Executing C_{pred} produces a reconstructed image I_{rec} . Omni-I2C then evaluates the fidelity of I_{rec} to I along multiple dimensions, including semantic consistency and structural alignment.

3.2 Benchmark Coverage Analysis

The Omni-I2C benchmark comprises 1,130 high-quality evaluation items curated to span a broad range of programming languages, academic subjects, and visual representations. This coverage provides substantially greater breadth and depth than existing image-to-code benchmarks.

Code Types. Omni-I2C includes five languages: Python, LaTeX, HTML, SVG, and SMILES. They are selected to cover complementary real-world workflows: Python and LaTeX represent widely used academic pipelines for visualization and typesetting, HTML and SVG target broad web and engineering settings for structured layouts and vector graphics, and SMILES brings domain-specific structure from cheminformatics.

Subjects. To reflect practical use cases, Omni-I2C covers over eight subjects, spanning foundational sciences (e.g., Mathematics, Physics, Chemistry) and applied domains (e.g., Computer Science,

Biomedicine, Economics). This diversity pushes models beyond pattern matching toward domain-aware interpretation of field-specific conventions and visual languages.

Image Types. Our benchmark features 43 distinct plot types, ranging from standard statistical charts to complex biological diagrams and geographical maps. This wide coverage ensures a comprehensive evaluation of a model’s visual perception abilities.

By jointly varying code types, domains, and visual formats, Omni-I2C evaluates integrated perception-and-generation ability rather than isolated skills, reducing the impact of single-domain expertise. Fig. 1 summarizes the data distribution; additional taxonomy details and qualitative examples are provided in the App. B.

3.3 Data Curation Process

Omni-I2C is built through a three-stage pipeline consisting of data collection, filtering, and decontamination, as demonstrated in Fig. 2.

Data collection. We adopt a taxonomy-driven strategy. We first establish a structured schema by synthesizing insights from established benchmarks (Yang et al., 2024; Wang et al., 2023) and authoritative reference materials (details in App. B.), which subsequently directs our targeted retrieval process. To ensure high-quality data acquisition, we adhere to three core selection criteria: (i) *Fidelity*, requiring a strict structural alignment between the code and its rendered output; (ii) *Complexity*, prioritizing samples with non-trivial logical depth or intricate spatial arrangements; and (iii) *Authenticity*, focusing on idiomatic code that reflects real-world programming practices. By aligning our efforts with this schema and these criteria, we facilitate a purposeful selection that fills specific data gaps and ensures a balanced distribution by design. We harvest high-fidelity image-code pairs from official galleries, community, tutorial textbooks and document (e.g., *Matplotlib Gallery*, *TikZ.net*) and refine subsets from existing datasets (Yang et al., 2024, 2025) to bolster under-represented categories across the Omni-I2C taxonomy.

Data filtering. To ensure the integrity of the evaluation set, we implement a multi-stage filtering pipeline. First, technical executability is guaranteed via isolated sandbox execution, where any code with compilation or runtime errors is discarded. Second, a manual audit sanitizes the collection by removing private data (PII) and social biases. Third, we employ an LLM-driven semantic

refinement process to develop a granular taxonomy, followed by an expert review to rectify any misalignments. Finally, the dataset is balanced by pruning over-represented classes and excluding instances with extreme sequence lengths. This rigorous curation ensures a high-quality, robust collection of evaluation samples.

Data decontamination. As the raw data originates from public web sources, we develop a refactoring strategy to mitigate data leakage and prevent shortcut learning through memorization. We employ an LLM to systematically refactor the code’s stylistic and aesthetic attributes, such as fonts and layouts. For textual elements within images, we perform semantic sanitization by erasing original text and injecting entirely new, synthetic strings. Furthermore, we manually modify key attributes, including colors, label coordinates, and textual content. To further decouple visual perception from linguistic priors, we occasionally introduce counterfactual “pseudo-labels” by substituting labels or changing the content to nonsense. This design choice ensures that models cannot rely on common-sense guesses or memorized patterns, forcing them to anchor their outputs strictly on the provided visual evidence. Representative visual comparisons of the original and decontaminated samples are provided in App. B.

3.4 Evaluation Metrics

Current I2C evaluation paradigms often fail to reconcile perceptual intuition with objective rigor (Li et al., 2024; Si et al., 2025; Wu et al., 2025; Yang et al., 2024). Methods relying solely on programmatic metrics tend to be library-dependent, while those based purely on LMM-judges may introduce subjective variance. To provide a more comprehensive and robust assessment, Omni-I2C introduces a synergetic perceptual-symbolic evaluation framework that harmonizes qualitative visual perception with quantitative algorithmic rigor.

Under this paradigm, each evaluation instance is characterized by a cross-comparison between the predicted pair (I_{rec}, C_{pred}) and the ground-truth reference (I_{GT}, C_{GT}) . Specifically, this framework bifurcates the assessment into two complementary tracks: Visual Perceptual Fidelity and Symbolic Functional Precision.

Visual Perceptual Fidelity. This track utilizes a frontier LMM to provide a human-aligned score S_{img} of the rendered output I_{rec} . Rather than requiring the model to perform unreliable fine-

grained parsing, this track prioritizes global structural coherence and stylistic essence. Specifically, S_{img} is derived from a multi-dimensional assessment across four axes: *Style*, *Layout*, *Elements*, and *Text*. While the textual axis ensures content accuracy, the former three dimensions are designed to capture holistic, global attributes—such as color harmony and spatial distribution. This serves as a qualitative proxy for human perception, ensuring that the synthesized result is visually plausible and stylistically consistent with the original input.

Symbolic Functional Precision. To complement visual intuition with objective grounding, this track employs an LLM-based evaluator to conduct a systematic audit of the generated code. The evaluator acts as an intelligent parser, first identifying the reference sets of functional logic blocks (\mathcal{F}_{GT}) and fine-grained parameters (\mathcal{P}_{GT}) within the ground-truth code C_{GT} . It then scrutinizes the predicted code C_{pred} to determine the subset of successfully implemented functions (\mathcal{F}_{match}) and accurately aligned parameters (\mathcal{P}_{match}). We quantify the model’s performance via Functional Coverage (S_{func}) and Parametric Accuracy (S_{para}), as:

$$S_{func} = \frac{|\mathcal{F}_{match}|}{|\mathcal{F}_{GT}|}, \quad S_{para} = \frac{|\mathcal{P}_{match}|}{|\mathcal{P}_{GT}|}.$$

Further details regarding the implementation of these metrics are provided in App. C.2. We also provide an updated symbolic evaluation protocol based on checklist-guided verification in App. A.3.

While this pipeline is applied to HTML, LaTeX, Python, and SVG tasks, we adopt a domain-specific approach for SMILES data. We assess reconstruction fidelity by calculating the Tanimoto Similarity (Bajusz et al., 2015) between the Morgan Fingerprints (Rogers and Hahn, 2010) of the ground-truth and generated code. Human inspection confirms that this metric effectively quantifies the accuracy of the reconstructed molecular topology, providing a robust proxy for structural identity that closely aligns with expert perceptual judgment. We provide examples in the App. I for illustration.

4 Experiment

4.1 Baseline Setup

We benchmark 13 LMMs on Omni-I2C, including proprietary frontiers—Claude 4.5 Sonnet (Anthropic, 2025), GPT 5.1 (OPENAI, 2025), Gemini 3 Pro (Deepmind, 2025), and Gemini 2.5

Pro (Comanici et al., 2025)—and a diverse open-weight spectrum (7B–241B) including InternVL 3.5 (Wang et al., 2025), Qwen3-VL (Bai et al., 2025a), Qwen2.5-VL (Bai et al., 2025b), and Gemma3-27B (Team et al., 2025). This setup spans major closed and open ecosystems to ensure a comprehensive comparison across varying model scales and regimes. We employ refined, language-specific prompts that enforce explicit structural constraints (see App. C for details). For the evaluation, we utilize GPT 4.1 to score perceptual fidelity and Claude 4.5 Sonnet to audit code-level functionality.

4.2 Main Results

Tab. 2 presents a comprehensive overview of the benchmarking results across 13 LMMs. Combined with the qualitative analysis in Fig. 3 and App. G, we derive several critical observations regarding the current landscape of image-to-code execution: **The challenging nature of the Omni-I2C.** Omni-I2C reveals that translating complex visuals into executable code remains a significant barrier. Tab. 2 shows that high execution rates do not guarantee high-fidelity reconstruction. Even for frontier models, a performance ceiling exists in semantic and parametric precision, confirming that Omni-I2C provides a rigorous and unsaturated testbed for evaluating the next generation of LMMs.

The Universal Gap Between Logic and Precision. A consistent trend is that Parametric Accuracy (S_{para}) significantly lags behind Functional Coverage (S_{func}). This is because S_{para} identifies fine-grained symbolic discrepancies—such as the subtle `linspace` deviations in Fig. 3—that are difficult to discern via human vision or perceptual metrics (S_{img}). Consequently, S_{para} provides a high-precision audit that visual assessment alone lacks the resolution to capture.

Divergent Expertise among Proprietary Frontiers. Proprietary models maintain a systematic lead, yet they exhibit distinct specialized strengths. Claude 4.5 Sonnet achieves the highest overall execution rate (96.9%), marking it as the most robust model for generating syntactically correct code. GPT 5.1 leads in logical fidelity (S_{func} at 60.1%), demonstrating superior decomposition of complex visual inputs.

Scaling Laws and Generational Leaps in Open-Weight Models. Performance tracks with both scale and architecture. The InternVL 3.5 series exhibits clear scaling effects up to 241B. Meanwhile, architectural advancements can outweigh

Table 2: Main results on Omni-I2C. Exec. denotes Execution Rate. For **SMILES**, we report Exec. and Tanimoto Similarity (T.S.) score. Best results in each category are **bold**, and second best are underlined within proprietary models and open-weight models, respectively. Updated results under the refined benchmark and evaluation protocol are provided in Appendix A; the corresponding updated main results are shown in Tab. 6.

| Model | Overall(1130) | | | | HTML(182) | | | | LaTeX(274) | | | | Python(364) | | | | SVG(210) | | | | SMILES(100) | |
|--------------------------------|---------------|-------------|-------------|-------------|-------------|-------------|-------------|--------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | S_{func} | S_{para} | S_{img} | Exec | S_{func} | S_{para} | S_{img} | Exec | S_{func} | S_{para} | S_{img} | Exec | S_{func} | S_{para} | S_{img} | Exec | S_{func} | S_{para} | S_{img} | Exec | T.S. | Exec |
| <i>Proprietary Models</i> | | | | | | | | | | | | | | | | | | | | | | |
| Claude Sonnet 4.5 | 56.3 | 34.2 | 74.6 | 96.9 | 61.6 | 56.7 | <u>76.8</u> | <u>99.5</u> | 52.4 | 30.2 | 67.1 | 93.1 | 55.4 | 34.5 | 80.9 | 98.9 | 58.2 | 19.6 | 71.8 | 97.6 | <u>59.4</u> | 94.0 |
| GPT 5.1 | 60.1 | <u>37.9</u> | 71.3 | 92.0 | 64.6 | 61.1 | 73.0 | 100.0 | 56.0 | 33.3 | 55.4 | <u>85.0</u> | 60.6 | 39.2 | 82.1 | 94.8 | 60.7 | 21.7 | <u>71.9</u> | <u>94.3</u> | 49.0 | 82.0 |
| Gemini 3 Pro | <u>57.4</u> | 39.7 | <u>71.8</u> | <u>92.1</u> | 60.1 | 56.5 | 72.1 | <u>99.5</u> | <u>53.5</u> | <u>31.3</u> | <u>56.5</u> | 83.6 | 55.8 | <u>38.6</u> | <u>82.0</u> | <u>95.1</u> | 62.9 | 37.9 | 73.7 | <u>94.3</u> | 48.0 | <u>87.0</u> |
| Gemini 2.5 Pro | 56.7 | 35.2 | 70.6 | 87.8 | <u>63.2</u> | <u>56.9</u> | 78.1 | 98.9 | 48.9 | 27.4 | 52.3 | 72.3 | <u>56.6</u> | 37.1 | 80.7 | 93.4 | <u>61.3</u> | <u>23.1</u> | 70.3 | 90.0 | 75.6 | 85.0 |
| Average | 57.6 | 36.8 | 72.1 | 92.2 | 62.4 | 57.8 | 75.0 | 99.5 | 52.7 | 30.6 | 57.8 | 83.5 | 57.1 | 37.4 | 81.4 | 95.6 | 60.8 | 25.6 | 71.9 | 94.1 | 58.0 | 87.0 |
| <i>Open-Weight Models</i> | | | | | | | | | | | | | | | | | | | | | | |
| InternVL3.5-241B-A28B-Instruct | 52.7 | 35.2 | 64.0 | 94.9 | 61.0 | 51.6 | 68.1 | 100.0 | <u>47.8</u> | 33.5 | <u>54.9</u> | <u>90.9</u> | 50.8 | 30.4 | 73.5 | 94.0 | 54.9 | 31.4 | <u>55.9</u> | 96.7 | 68.9 | 96.0 |
| InternVL3.5-38B-Instruct | 39.3 | 24.9 | 49.6 | 90.7 | 45.6 | 41.9 | 53.1 | 100.0 | 31.8 | 23.2 | 39.0 | 87.6 | 41.1 | 24.6 | 59.3 | 87.4 | 40.4 | 12.8 | 43.5 | <u>91.0</u> | 58.8 | <u>94.0</u> |
| InternVL3.5-8B-Instruct | 27.7 | 17.8 | 36.5 | 81.5 | 41.6 | 36.1 | 45.4 | <u>98.9</u> | 23.8 | 15.3 | 32.9 | 85.0 | 22.9 | 15.7 | 42.2 | 72.0 | 29.2 | 8.9 | 23.8 | 77.1 | 45.3 | 84.0 |
| Qwen3-VL-235B-A22B-Instruct | 58.0 | 39.1 | 64.0 | <u>92.0</u> | 64.9 | 57.2 | <u>61.4</u> | 97.8 | 52.1 | 37.0 | 60.5 | 90.2 | 59.2 | 37.6 | <u>72.7</u> | <u>90.7</u> | 57.7 | 28.5 | 55.9 | 90.5 | <u>66.9</u> | <u>94.0</u> |
| Qwen2.5-VL-72B-Instruct | 48.8 | 37.1 | 55.3 | 90.5 | 53.8 | 52.6 | 48.3 | 97.8 | 45.6 | <u>36.2</u> | 53.3 | 93.1 | 48.0 | 32.7 | 63.7 | 88.2 | 50.1 | <u>32.4</u> | 49.3 | 89.0 | 29.8 | 82.0 |
| Qwen3-VL-32B-Instruct | <u>53.4</u> | <u>37.7</u> | <u>57.7</u> | 84.3 | <u>61.2</u> | <u>55.5</u> | 55.2 | 90.1 | 47.3 | 31.9 | 50.4 | 79.9 | <u>52.4</u> | <u>35.6</u> | 65.2 | 84.9 | <u>56.1</u> | 33.6 | 56.5 | 84.8 | 51.6 | 83.0 |
| Qwen3-VL-8B-Instruct | 47.2 | 31.0 | 49.8 | 82.2 | 55.5 | 47.1 | 53.0 | 86.8 | 39.1 | 28.6 | 38.4 | 75.2 | 46.5 | 28.5 | 59.0 | 82.4 | 51.7 | 24.6 | 46.2 | 85.7 | 53.4 | 85.0 |
| Qwen2.5-VL-7B-Instruct | 31.5 | 23.3 | 36.8 | 73.5 | 38.2 | 39.8 | 43.8 | 83.0 | 24.1 | 18.3 | 29.8 | 66.4 | 33.5 | 22.9 | 45.9 | 78.6 | 32.0 | 16.1 | 24.4 | 76.7 | 15.3 | 51.0 |
| Gemma3-27B | 32.2 | 21.4 | 41.1 | 76.8 | 43.1 | 41.9 | 50.5 | 97.8 | 18.0 | 17.4 | 21.6 | 53.3 | 34.9 | 21.5 | 52.2 | 82.7 | 36.6 | 8.7 | 38.9 | 89.5 | 23.1 | 55.0 |
| Average | 43.4 | 29.7 | 50.5 | 85.2 | 51.7 | 47.1 | 53.2 | 94.7 | 36.6 | 26.8 | 42.3 | 80.2 | 43.3 | 27.7 | 59.3 | 84.5 | 45.4 | 21.9 | 43.8 | 86.8 | 45.9 | 80.4 |

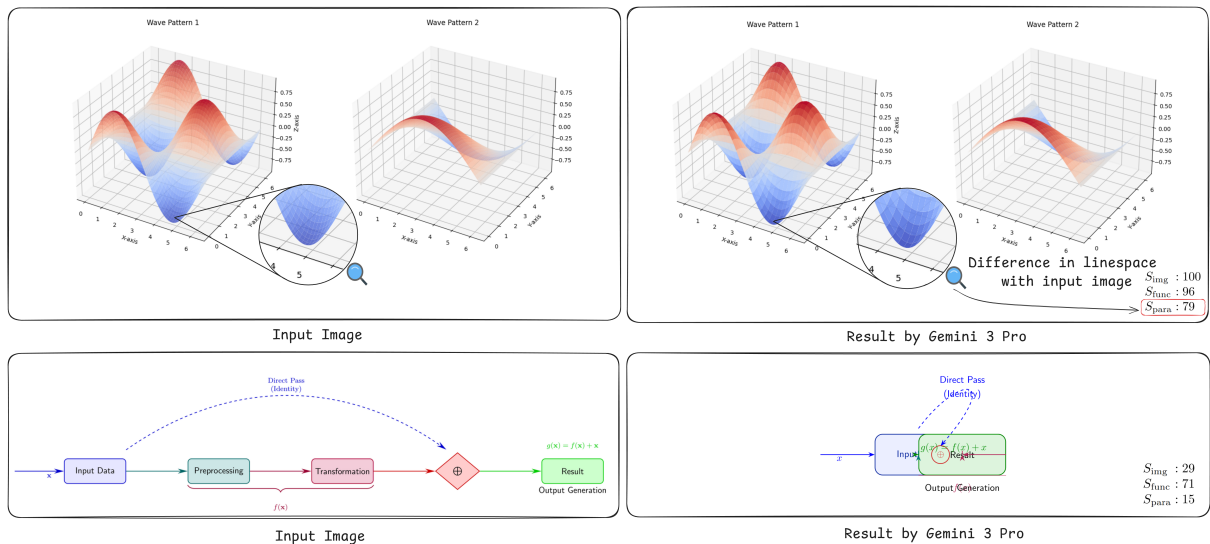


Figure 3: Representative examples from Gemini 3 pro.

raw parameter counts; for instance, the 32B Qwen3-VL consistently surpasses the larger 72B Qwen2.5-VL, illustrating a distinct generational leap. Most notably, Qwen3-VL-235B-A22B-Instruct achieves state-of-the-art results in structured LaTeX tasks, rivaling proprietary frontiers and marking a significant milestone for open-source precision in specialized domains.

SVG and LaTeX Complexity as a Rigorous Geometric Discriminator. SVG and LaTeX serve as the most demanding tasks due to its reliance on intricate geometric topologies rather than hierarchical tags. This poses a significant barrier for small-scale models (7B–8B), which frequently fail to maintain syntactic validity. The 18-point gap in S_{para} of SVG tasks among proprietary models further highlights that translating visual geometry into precise spatial code is a key differentiator for

frontier LMMs.

5 Discussion

5.1 Study on Evaluation Metrics

Image Level Evaluation. Within the Visual Perceptual Fidelity track, we prompt an LMM judge to quantify the visual fidelity (S_{img}) of the reconstructed outputs. To validate the reliability of this metric, we conduct a human-preference alignment study involving three candidate models—Gemini 2.5 Pro, InternVL 3.5-8B-Instruct, and Qwen3-VL-235B-A22B-Instruct—evaluated by three distinct judges: GPT 4.1, Qwen-VL-Max, and Qwen3-VL-Plus. We benchmark our proposed S_{img} metric against four existing image-level evaluation paradigms by measuring their respective alignment with human consensus. The prompt templates of existing image-level evaluation paradigms are

Table 3: **Image-level alignment analysis.** Gray shading highlights the best-performing setting (Omni-I2C with GPT 4.1). A detailed description of the evaluation metrics and the rationale behind their selection are provided in the App. H.

| Prompt | Judge | Ken. τ (\uparrow) | Spr. ρ (\uparrow) | RMSE (\downarrow) |
|-----------------|---------------|----------------------------|----------------------------|-----------------------|
| Chart2Code | GPT 4.1 | 0.7254 | 0.7769 | 0.5474 |
| | Qwen-VL-Max | 0.5971 | 0.6402 | 0.6951 |
| | Qwen3-VL-Plus | 0.5736 | 0.6156 | 0.7293 |
| ChartMimic | GPT 4.1 | 0.7835 | 0.8252 | 0.4934 |
| | Qwen-VL-Max | 0.7241 | 0.7723 | 0.5575 |
| | Qwen3-VL-Plus | 0.7564 | 0.8004 | 0.5265 |
| Plot2Code | GPT 4.1 | 0.7918 | 0.8355 | 0.4780 |
| | Qwen-VL-Max | 0.7800 | 0.8240 | 0.4972 |
| | Qwen3-VL-Plus | 0.7623 | 0.7999 | 0.5229 |
| Self-Reflection | GPT 4.1 | 0.7979 | 0.8417 | 0.4661 |
| | Qwen-VL-Max | 0.7275 | 0.7712 | 0.5508 |
| | Qwen3-VL-Plus | 0.7214 | 0.7700 | 0.5575 |
| Omni-I2C | GPT 4.1 | 0.8304 | 0.8681 | 0.4284 |
| | Qwen-VL-Max | 0.7897 | 0.8323 | 0.4780 |
| | Qwen3-VL-Plus | 0.7931 | 0.8353 | 0.4780 |

shown in App. C.2. For the annotation process, 10 human participants are each assigned 50 randomly sampled instances. For each instance, annotators are presented with the ground-truth image (I_{GT}) alongside three reconstructions (I_{rec}) generated by the contestant models, which they subsequently ranked according to a standardized preference policy (further details are provided in the App. D).

The results, summarized in Tab. 3, indicate that both the judge model and the prompting strategy significantly influence the reliability of perceptual assessment. Notably, the Omni-I2C prompt paired with GPT 4.1 achieves the strongest correlation with human preferences, yielding a Kendall’s τ of 0.8304 and the lowest RMSE of 0.4284. Even when deployed with alternative judges such as Qwen-VL-Max or Qwen3-VL-Plus, our metric consistently outperforms methods from prior work.

Code Level Evaluation. Given the semantic complexity of code generation, we utilize an LLM-based auditor to compute Functional Coverage (S_{func}) and Parametric Accuracy (S_{para}). While this enables scalable and automatic assessment, a single LLM judge is prone to systematic biases and capability limitations. Aggregating multiple judges (e.g., via consensus-based voting (Cobbe et al., 2021)) yields a more stable evaluation signal by mitigating idiosyncratic errors. The prompt templates of code-level evaluation paradigms are shown in App C.

To examine judge reliability, we evaluate three frontier models—Claude 4.5 Sonnet, Gemini 2.5 Pro, and GPT 4.1—by eliciting rankings over 1,130 instances across four target models (Gemini

Table 4: **Judge agreement and alignment to majority vote (MV).** All values are computed over 1,130 indices and reported as mean \pm std unless otherwise noted. *MV align* is the composite alignment to the MV ranking, defined as $\frac{1}{2}(\tau + \rho)$, where τ is Kendall’s τ and ρ is Spearman’s ρ against MV. *Closest to MV (count/%)* reports how often a judge achieves the highest MV align among judges on each item. A detailed description of the evaluation metrics and the rationale behind their selection are provided in the App. H.

| Per-judge vs MV | | |
|--------------------------------|--|------------------------------------|
| Code Judge | MV align \uparrow | Closest to MV (count/%) \uparrow |
| Claude Sonnet 4.5 | 0.8436 ($\tau=0.8546, \rho=0.8325$) | 691 / 61.15% |
| Gemini 2.5 pro | 0.7747 ($\tau=0.8182, \rho=0.7312$) | 222 / 19.65% |
| GPT 4.1 | 0.8419 ($\tau=0.8510, \rho=0.8329$) | 217 / 19.20% |
| Overall (all judges) | | |
| Kendall’s W \uparrow | 0.7405 \pm 0.2270 | |
| Kendall’s τ \uparrow | 0.6517 \pm 0.3378 | |
| Spearman’s ρ \uparrow | 0.6080 \pm 0.3343 | |
| Rank range (mean) \downarrow | 0.91 | |

2.5 Pro, GPT 4.1 mini, InternVL 3.5-8B-Instruct, and Qwen3-VL-235B-A22B-Instruct). We then construct a Majority-Vote (MV) ranking as a silver standard for expert consensus. As shown in Tab. 8, the judges exhibit high agreement overall (Kendall’s $W = 0.74$, mean Kendall’s $\tau = 0.65$, and mean Spearman’s $\rho = 0.61$). Notably, the observed disagreements are small in magnitude, with a mean rank range of 0.91, suggesting that variances are largely restricted to local swaps rather than systematic inversions. We further evaluate each judge’s alignment with the MV: Claude 4.5 Sonnet demonstrates the highest fidelity to the consensus, achieving an MV align score of 0.8436 ($\tau = 0.8546, \rho = 0.8325$). On a per-instance basis, Sonnet 4.5 is the judge most frequently closest to the MV ranking (61.15% of cases), significantly outperforming its peers.

Based on the results, we adopt Claude 4.5 Sonnet as our final code judge. This choice provides a high-fidelity proxy for the consensus ensemble while optimizing for API cost. We note, however, that multi-judge aggregation remains a principled and more robust option for studies where computational budgets permit higher overhead. An updated analysis of code-level judge selection is provided in App. A.4.

5.2 Different Prompting Methods

We evaluate four representative models—Gemini 2.5 Pro, GPT 4.1 mini, Qwen3-VL-235B-A22B-Instruct, and InternVL 3.5-8B-Instruct—under five distinct prompting strategies to assess their impact on Omni-I2C performance. These include: (i) Direct (zero-shot); (ii) Few-Shot (in-context exam-

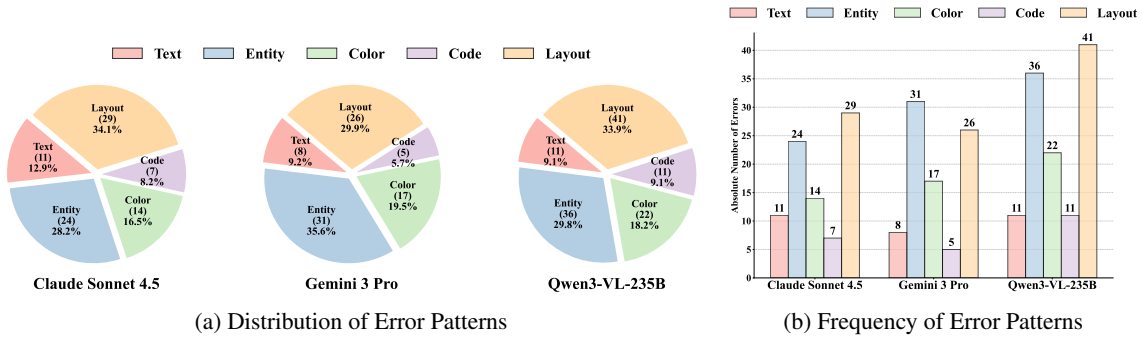


Figure 4: Comparative Analysis of Error Patterns in Claude 4.5 Sonnet, Gemini 3 Pro, and Qwen3-VL-235B-A22B-Instruct. Case studies and evaluations across various languages are detailed in the App. F.

Table 5: Studies on different inference prompting methods. The best results for each model are highlighted in **bold**. Values exceeding and falling below the corresponding average are highlighted in red and blue. Detailed cross-lingual results are provided in App. E

| Model | Method | S_{func} | S_{para} | S_{img} | Exec |
|-----------------------------|-----------------|-------------------|-------------------|------------------|--------------|
| Gemini2-5 Pro | Direct | 66.39 | 56.49 | 69.70 | 88.32 |
| | Few-Shot | 63.80 | 54.10 | 68.74 | 86.73 |
| | Scaffold | 50.62 | 42.57 | 51.65 | 83.01 |
| | Self-Commentary | 67.67 | 57.17 | 69.31 | 92.39 |
| | Hint-Enhanced | 66.16 | 54.68 | 67.81 | 87.79 |
| GPT 4.1 mini | Direct | 45.80 | 35.83 | 53.22 | 81.33 |
| | Few-Shot | 62.83 | 50.91 | 59.14 | 92.03 |
| | Scaffold | 36.20 | 28.97 | 43.50 | 82.83 |
| | Self-Commentary | 48.50 | 37.95 | 57.92 | 84.96 |
| | Hint-Enhanced | 64.19 | 52.87 | 59.50 | 89.29 |
| Qwen3-VL-235B-A22B-Instruct | Direct | 54.53 | 38.98 | 54.93 | 86.46 |
| | Few-Shot | 56.70 | 39.83 | 56.85 | 87.35 |
| | Scaffold | 42.16 | 28.85 | 42.27 | 82.92 |
| | Self-Commentary | 55.15 | 39.61 | 55.90 | 87.26 |
| | Hint-Enhanced | 59.27 | 41.94 | 55.25 | 85.84 |
| InternVL3.5-8B-Instruct | Direct | 23.43 | 20.75 | 28.23 | 71.77 |
| | Few-Shot | 22.02 | 21.56 | 28.08 | 71.15 |
| | Scaffold | 11.64 | 9.65 | 11.88 | 67.70 |
| | Self-Commentary | 25.01 | 21.77 | 29.39 | 76.81 |
| | Hint-Enhanced | 25.21 | 21.36 | 28.52 | 72.65 |
| Average | Direct | 47.54 | 38.02 | 51.52 | 81.97 |
| | Few-Shot | 51.34 | 41.60 | 53.20 | 84.31 |
| | Scaffold | 35.16 | 27.51 | 37.32 | 79.12 |
| | Self-Commentary | 49.08 | 39.12 | 53.13 | 85.35 |
| | Hint-Enhanced | 53.71 | 42.71 | 52.77 | 83.89 |

ples); (iii) Hint-Enhanced, eliciting visual rationales before code generation; (iv) Scaffold, which overlays a coordinate grid for spatial grounding; and (v) Self-Commentary, which interleaves code with explanatory logic. Detailed implementations and full templates for each strategy are provided in App. E. Our results in Tab. 5 reveal a high sensitivity to prompting across all scales. For instance, GPT 4.1 mini’s S_{func} jumps by nearly 20% when transitioning from Direct to Hint-Enhanced prompting. Overall, reasoning-augmented strategies (*Hint-Enhanced*, *Self-Commentary*) consistently outperform baselines by externalizing the image-to-code translation process. Conversely, *Scaffold* prompting often degrades performance due to visual noise, while *Few-Shot* remains a formidable baseline for syntactically rigid tasks like HTML. This high vari-

ance underscores that achieving “plug-and-play” stability remains an open challenge. An extended analysis of model behaviors is provided in App. E.

5.3 Error Analysis

To elucidate the fundamental bottlenecks in current LMMs, we conduct a manual audit of 240 instances (100 per model) on Claude 4.5 Sonnet, Gemini 3 Pro, and Qwen3-VL-235B-A22B-Instruct. We define a five-tier taxonomy to categorize observed failure modes: (1) Code-level Logic, where syntactically valid code fails to map visual requirements to correct programmatic attributes; (2) Textual Precision, involving OCR failures in dense labels or complex symbolic notations; (3) Entity Integrity, characterized by the omission or misinterpretation of visual primitives like data points or legends; (4) Colorimetric Accuracy, where generated color codes deviate from the ground truth; and (5) Spatial Layout, involving distorted aspect ratios or erroneous coordinate transformations.

As illustrated in Fig. 4 (a), the error distributions exhibit a high degree of consistency across different models. Notably, Entity- and Layout-related errors are dominant, collectively accounting for over 60% of the total. This indicates that current models continue to face significant challenges in recognizing visual entities and structures, as well as in implementing precise layout control through code. Furthermore, as depicted in Fig. 4 (b), while the distributional patterns remain similar, the absolute volume of errors varies due to disparities in fundamental capabilities, such as visual perception and code generation. Consequently, Qwen3-VL-235B-A22B-Instruct exhibits a significantly higher error count compared to the other two models. A more granular analysis, including exhaustive definitions and qualitative case studies for each category, is provided in App. F.

6 Conclusion

We introduce Omni-I2C, a benchmark for Image-to-Code generation characterized by its breadth across languages and subjects, and its depth in evaluation. By decoupling visual fidelity from symbolic precision, Omni-I2C exposes critical bottlenecks in I2C mapping. Our results reveal a significant performance gap even among frontier LMMs, establishing a substantial challenge to guide the development of more robust, visually-grounded multi-modal agents.

7 Limitation

Despite its multi-faceted evaluation design, Omni-I2C is subject to several limitations. Omni-I2C currently lacks natural scene images, a common gap in Image-to-Code research that hinders the evaluation of LMMs in generalizing to real-world visual complexities. Also, the dual-dimension evaluation relies on LLM/LMM-as-a-judge, which introduces non-negligible API costs and potential latency, limiting its scalability for rapid iterative testing. Future iterations of Omni-I2C will aim to bridge the “natural image gap” and investigate hybrid evaluation paradigms that balance semantic accuracy with resource efficiency, thereby setting more rigorous benchmarks for the community.

8 Acknowledgments

This work was supported in part by the New Generation Artificial Intelligence-National Science and Technology Major Project (No. 2025ZD0123602), the National Natural Science Foundation of China (Nos. 62276203 and 62472325), and the Hubei Provincial Key R&D Program, China (Grant No. 2025BAB021). The numerical calculations in this paper were carried out in part on the supercomputing system at the Supercomputing Center of Wuhan University.

We also thank the contributors to the data collection and subsequent human alignment experiments. They include Chunqing Du, Qian Cai, Zhuo Song, Xiaoshu Yang, Yuhao Zhang, Junhan Liu, Zheyu Zheng, Jiahui Li, Jiaming Sun, Duojie Chen, Zihan Lou, Zhiwei Wang, Tong Liu, and Zhiyi Yang.

References

Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm

Reynolds, and 1 others. 2022. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems*, 35:23716–23736.

Anthropic. 2025. claude-sonnet-4.5. <https://www.anthropic.com/news/claude-sonnet-4-5>.

Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. 2015. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433.

Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and 1 others. 2021. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*.

Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen, Xionghui Chen, Zesen Cheng, Lianghao Deng, Wei Ding, Chang Gao, Chunjiang Ge, Wenbin Ge, Zhi-fang Guo, Qidong Huang, Jie Huang, Fei Huang, Binyuan Hui, Shutong Jiang, Zhaohai Li, Mingsheng Li, and 45 others. 2025a. Qwen3-vl technical report. *arXiv preprint arXiv:2511.21631*.

Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, and 8 others. 2025b. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*.

Dávid Bajusz, Anita Rácz, and Károly Héberger. 2015. Why is tanimoto index an appropriate choice for fingerprint-based similarity calculations? *Journal of Cheminformatics*, 7(1):1–13.

Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. 2024. Sharegpt4v: Improving large multi-modal models with better captions. In *European Conference on Computer Vision*, pages 370–387. Springer.

Mark Chen. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, and 1 others. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.

Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, and 1 others. 2025. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*.

Deepmind. 2025. gemini-3-pro. <https://deepmind.google/models/gemini/pro/>.

- Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Yu Wu, YK Li, and 1 others. 2024. Deepseek-coder: When the large language model meets programming—the rise of code intelligence. *arXiv preprint arXiv:2401.14196*.
- Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, and 1 others. 2024. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186*.
- Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. 2023. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770*.
- Yuhang Lai, Chengxi Li, Yiming Wang, Tianyi Zhang, Ruiqi Zhong, Luke Zettlemoyer, Wen-tau Yih, Daniel Fried, Sida Wang, and Tao Yu. 2023. Ds-1000: A natural and reliable benchmark for data science code generation. In *International Conference on Machine Learning*, pages 18319–18345. PMLR.
- Jinke Li, Jiarui Yu, Chenxing Wei, Hande Dong, Qiang Lin, Liangjing Yang, Zhicai Wang, and Yanbin Hao. 2025. Unisvg: A unified dataset for vector graphic understanding and generation with multimodal large language models. In *Proceedings of the 33rd ACM International Conference on Multimedia*, pages 13156–13163.
- Kaixin Li, Yuchen Tian, Qisheng Hu, Ziyang Luo, Zhiyong Huang, and Jing Ma. 2024. Mmcode: Benchmarking multimodal large language models for code generation with visually rich programming problems. *arXiv preprint arXiv:2404.09486*.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916.
- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. 2023. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*.
- Fanqing Meng, Lingxiao Du, Zongkai Liu, Zhixiang Zhou, Quanfeng Lu, Daocheng Fu, Botian Shi, Wenhai Wang, Junjun He, Kaipeng Zhang, and 1 others. 2025. Mm-eureka: Exploring visual aha moment with rule-based large-scale reinforcement learning. *CoRR*.
- OpenAI. 2025. gpt-5. <https://openai.com/index/introducing-gpt-5/>.
- OPENAI. 2025. gpt-5.1. <https://openai.com/index/gpt-5-1/>.
- Runqi Qiao, Qiuna Tan, Guanting Dong, MinhuiWu MinhuiWu, Chong Sun, Xiaoshuai Song, Jiapeng Wang, Zhuoma Gongque, Shanglin Lei, Yifan Zhang, and 1 others. 2025. We-math: Does your large multimodal model achieve human-like mathematical reasoning? In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 20023–20070.
- Josselin S Roberts, Tony Lee, Chi H Wong, Michihiro Yasunaga, Yifan Mai, and Percy Liang. 2024. Image2struct: Benchmarking structure extraction for vision-language models. *Advances in Neural Information Processing Systems*, 37:115058–115097.
- David Rogers and Mathew Hahn. 2010. Extended-connectivity fingerprints. *Journal of Chemical Information and Modeling*, 50(5):742–754.
- Chenglei Si, Yanzhe Zhang, Ryan Li, Zhengyuan Yang, Ruibo Liu, and Diyi Yang. 2025. Design2code: Benchmarking multimodal code generation for automated front-end engineering. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 3956–3974.
- Jiahao Tang, Henry Hengyuan Zhao, Lijian Wu, Yifei Tao, Dongxing Mao, Yang Wan, Jingru Tan, Min Zeng, Min Li, and Alex Jinpeng Wang. 2025. From charts to code: A hierarchical benchmark for multimodal models. *arXiv preprint arXiv:2510.17932*.
- Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, and 1 others. 2025. Gemma 3 technical report. *arXiv preprint arXiv:2503.19786*.
- Ke Wang, Junting Pan, Weikang Shi, Zimu Lu, Houxing Ren, Aojun Zhou, Mingjie Zhan, and Hongsheng Li. 2024. Measuring multimodal mathematical reasoning with math-vision dataset. *Advances in Neural Information Processing Systems*, 37:95095–95169.
- Ke Wang, Houxing Ren, Aojun Zhou, Zimu Lu, Sichun Luo, Weikang Shi, Renrui Zhang, Linqi Song, Mingjie Zhan, and Hongsheng Li. 2023. Math-coder: Seamless code integration in llms for enhanced mathematical reasoning. *arXiv preprint arXiv:2310.03731*.
- Weiyun Wang, Zhangwei Gao, Lixin Gu, Hengjun Pu, Long Cui, Xingguang Wei, Zhaoyang Liu, Linglin Jing, Shenglong Ye, Jie Shao, and 1 others. 2025. Internvl3. 5: Advancing open-source multimodal models in versatility, reasoning, and efficiency. *arXiv preprint arXiv:2508.18265*.
- Chengyue Wu, Zhixuan Liang, Yixiao Ge, Qiushan Guo, Zeyu Lu, Jiahao Wang, Ying Shan, and Ping Luo. 2025. Plot2code: A comprehensive benchmark for evaluating multi-modal large language models in code generation from scientific plots. In *Findings of the Association for Computational Linguistics: NAACL 2025*, pages 3006–3028.

- Cheng Yang, Chufan Shi, Yaxin Liu, Bo Shui, Junjie Wang, Mohan Jing, Linran Xu, Xinyu Zhu, Siheng Li, Yuxiang Zhang, and 1 others. 2024. Chartmimic: Evaluating Imm’s cross-modal reasoning capability via chart-to-code generation. *arXiv preprint arXiv:2406.09961*.
- Yue Yang, Ajay Patel, Matt Deitke, Tanmay Gupta, Luca Weihs, Andrew Head, Mark Yatskar, Chris Callison-Burch, Ranjay Krishna, Aniruddha Kembhavi, and 1 others. 2025. Scaling text-rich image understanding via code-guided synthetic multimodal data generation. *arXiv preprint arXiv:2502.14846*.
- Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik Narasimhan. 2024. tau-bench: A benchmark for tool-agent-user interaction in real-world domains. *arXiv preprint arXiv:2406.12045*.
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, and 1 others. 2024. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9556–9567.
- Xiang Yue, Tianyu Zheng, Yuansheng Ni, Yubo Wang, Kai Zhang, Shengbang Tong, Yuxuan Sun, Botao Yu, Ge Zhang, Huan Sun, and 1 others. 2025. Mmmu-pro: A more robust multi-discipline multimodal understanding benchmark. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15134–15186.
- Chi Zhang, Haibo Qiu, Qiming Zhang, Yufei Xu, Zhixiong Zeng, Siqi Yang, Peng Shi, Lin Ma, and Jing Zhang. 2025a. Perceptual-evidence anchored reinforced learning for multimodal reasoning. *arXiv preprint arXiv:2511.18437*.
- Chi Zhang, Haibo Qiu, Qiming Zhang, Zhixiong Zeng, Lin Ma, and Jing Zhang. 2025b. Deepsketcher: Internalizing visual manipulation for multimodal reasoning. *arXiv preprint arXiv:2509.25866*.
- Yi-Fan Zhang, Xingyu Lu, Shukang Yin, Chaoyou Fu, Wei Chen, Xiao Hu, Bin Wen, Kaiyu Jiang, Changyi Liu, Tianke Zhang, and 1 others. 2025c. Thyme: Think beyond images. *arXiv preprint arXiv:2508.11630*.
- Ziwei Zheng, Michael Yang, Jack Hong, Chenxiao Zhao, Guohai Xu, Le Yang, Chao Shen, and Xing Yu. 2025. Deepeyes: Incentivizing" thinking with images" via reinforcement learning. *arXiv preprint arXiv:2505.14362*.
- Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan Bisk, Daniel Fried, and 1 others. 2023. Webarena: A realistic web environment for building autonomous agents. *arXiv preprint arXiv:2307.13854*.

A Updated Benchmark and Evaluation Protocol

With the rapid iteration of frontier proprietary models and our continued investigation, we make a minor update to the benchmark and evaluation pipeline. Specifically, we refine a small portion of the benchmark and revise several components of the evaluation protocol, yielding a more rigorous, fair, and reliable assessment framework.

Under the revised benchmark and evaluation framework, the results in the main table are updated in Tab. 6.

Compared with the original version, the absolute values in the main table change due to the stricter benchmark filtering and the refined modality-aware evaluation setup. More details are shown as follows.

A.1 Benchmark deduplication and minor taxonomy refinement

After a more detailed benchmark-level analysis, we removed 50 potentially redundant instances, reducing the benchmark size from 1,130 to 1,080; this accounts for only a small fraction of the full benchmark. Meanwhile, we slightly refined the taxonomy, increasing the number of figure types from 43 to 45. These changes constitute a minor update and do not alter the underlying task definition or intended benchmark scope.

A.2 Refined Evaluation Protocol for HTML and SVG

We also refine the modality-specific evaluation protocol to improve fairness.

HTML prompting refinement For a subset of HTML instances whose ground-truth implementations are rendered with the ECharts library, we found in subsequent optimization experiments that it is more appropriate to explicitly mention ECharts in the inference prompt. We therefore updated the HTML prompting condition so that the generation setup is better aligned with the underlying ground-truth implementation.

SVG evaluation/rendering refinement For the SVG portion of the benchmark, we adopted a fairer and more reasonable evaluation protocol, for example, by predefining the background color during SVG rendering to ensure a more reliable evaluation. This update yields more faithful rendering behavior during evaluation and better reflects the models'

actual performance on SVG code generation. The SVG task definition itself remains unchanged.

Overall, these changes make the evaluation more consistent across modalities and reduce confounding factors caused by execution environments rather than model capability.

A.3 Checklist-based Symbolic Evaluation

A key challenge in LLM-as-judge evaluation is the inherent stochasticity of generative models, especially when the judge is asked to assess code quality in an open-ended manner. This issue is particularly pronounced in symbolic code-level evaluation, where small variations in judge reasoning can lead to unstable judgments on functional correctness and parameter accuracy. A common solution is the Direct GT-to-Code protocol, where the judge compares the predicted code against the ground-truth implementation holistically, and directly infers whether the functional logic and fine-grained parameters are correct. However, in the absence of explicit instance-level criteria, the judge relies on holistic code-to-code comparison, which may introduce extra ambiguity and reduce consistency across runs.

To address this issue, we adopt a **checklist-based** symbolic evaluation protocol. For each instance, we pre-define a deterministic checklist consisting of functional logic blocks and fine-grained parameters. The reference functional units and parameters are first extracted from the ground-truth code using a strong LLM and then manually verified and modified. During evaluation, the LLM judge no longer performs holistic code-to-code comparison; instead, it serves only as a structured parser that checks whether these fixed requirements are correctly implemented in the predicted code. This design reduces evaluation ambiguity and improves stability, interpretability, and reproducibility.

To quantitatively compare the two protocols, we conduct a 10-run variance analysis on a fixed set of model outputs, evaluating both the original Direct GT-to-Code protocol and our checklist-based protocol. The results are summarized in Tab. 7. Compared with Direct GT-to-Code protocol, the checklist-based protocol consistently yields lower variance on all metrics. In particular, it reduces the variance of Logic Score from 8.288 to 7.414 (**10.6%↓**), and the variance of Parameter Score from 30.605 to 5.544 (**81.9%↓**).

These results show that anchoring symbolic eval-

Table 6: Updated main results on Omni-I2C. Exec. denotes Execution Rate. For **SMILES**, we report Exec. and Tanimoto Similarity (T.S.) score. Best results in each category are **bold**, second best are underlined for proprietary models and open-weight models.

| Model | Overall(1080) | | | | HTML(166) | | | | LaTeX(265) | | | | Python(357) | | | | SVG(192) | | | | SMILES(100) | |
|--------------------------------|-------------------|-------------------|------------------|-------------|-------------------|-------------------|------------------|--------------|-------------------|-------------------|------------------|-------------|-------------------|-------------------|------------------|-------------|-------------------|-------------------|------------------|--------------|-------------|-------------|
| | S_{func} | S_{para} | S_{img} | Exec | S_{func} | S_{para} | S_{img} | Exec | S_{func} | S_{para} | S_{img} | Exec | S_{func} | S_{para} | S_{img} | Exec | S_{func} | S_{para} | S_{img} | Exec | T.S. | Exec |
| <i>Proprietary Models</i> | | | | | | | | | | | | | | | | | | | | | | |
| Claude Sonnet 4.5 | 60.4 | 40.8 | 76.7 | 97.2 | 67.2 | <u>55.3</u> | <u>84.2</u> | 99.4 | 47.7 | 29.4 | 66.8 | 92.8 | 62.5 | 41.0 | 81.0 | 98.9 | 68.1 | <u>44.0</u> | 75.8 | 100.0 | <u>59.4</u> | 94.0 |
| GPT 5.1 | <u>60.7</u> | <u>40.9</u> | 74.1 | 92.7 | <u>67.4</u> | 55.2 | 84.2 | <u>99.4</u> | 46.9 | 28.0 | 55.2 | <u>84.5</u> | 63.3 | 42.9 | 81.9 | 94.7 | <u>69.1</u> | <u>42.5</u> | <u>76.9</u> | <u>100.0</u> | 49.0 | 82.0 |
| Gemini 3 Pro | 62.2 | 41.4 | <u>76.1</u> | <u>92.8</u> | 68.9 | 56.0 | 86.2 | 98.2 | <u>47.2</u> | <u>27.2</u> | <u>56.5</u> | 83.4 | <u>62.6</u> | <u>42.2</u> | <u>81.9</u> | <u>95.0</u> | 76.3 | 47.0 | 83.7 | 100.0 | 48.0 | 87.0 |
| Gemini 2.5 Pro | 58.4 | 39.6 | 71.2 | 89.3 | 61.5 | 54.9 | 80.8 | 98.8 | 45.0 | 26.9 | 52.5 | 72.5 | 62.0 | 40.9 | 80.7 | 93.3 | 67.6 | 41.6 | 71.1 | 99.0 | 75.6 | 85.0 |
| Average | 60.4 | 40.7 | 74.5 | 93.0 | 66.2 | 55.4 | 83.8 | 99.0 | 46.7 | 27.9 | 57.8 | 83.3 | 62.6 | 41.8 | 81.4 | 95.5 | 70.3 | 43.8 | 76.9 | 99.8 | 58.0 | 87.0 |
| <i>Open-Weight Models</i> | | | | | | | | | | | | | | | | | | | | | | |
| InternVL3_5_241B_A28B-Instruct | 50.8 | 36.2 | <u>66.5</u> | 95.3 | 55.6 | 51.9 | <u>76.6</u> | 98.8 | 36.2 | 25.5 | <u>54.5</u> | <u>90.6</u> | <u>55.9</u> | 36.3 | 73.8 | 94.4 | 57.3 | 37.5 | 61.1 | 100.0 | 68.9 | 96.0 |
| InternVL3_5_38B-Instruct | 37.5 | 30.2 | 52.2 | 90.2 | 40.6 | 46.6 | 61.8 | 91.0 | 22.3 | 21.0 | 38.6 | 87.2 | 44.0 | 27.7 | 59.2 | 87.1 | 43.6 | 33.5 | 49.7 | 99.0 | 58.8 | <u>94.0</u> |
| InternVL3_5_8B-Instruct | 28.6 | 24.1 | 42.3 | 88.4 | 36.7 | 40.2 | 56.5 | 97.6 | 14.1 | 16.8 | 34.6 | 89.1 | 34.1 | 20.6 | 45.8 | 77.9 | 31.2 | 26.8 | 34.2 | 99.0 | 48.1 | 87.0 |
| Qwen3-VL-235B-A22B-Instruct | 57.0 | 39.9 | 70.0 | <u>93.8</u> | 65.4 | 54.3 | 79.3 | 100.0 | 46.2 | 30.1 | 60.6 | <u>90.2</u> | 57.1 | 39.4 | <u>72.4</u> | <u>90.5</u> | 64.6 | <u>42.0</u> | 70.4 | 99.5 | <u>66.9</u> | <u>94.0</u> |
| Qwen2.5-VL-72B-Instruct | 46.0 | 34.2 | 59.8 | 92.7 | 52.0 | 51.4 | 72.3 | <u>100.0</u> | 34.4 | <u>25.6</u> | 53.3 | 92.8 | 49.4 | 32.2 | 64.0 | 88.5 | 50.4 | 34.8 | 50.2 | 99.5 | 29.8 | 82.0 |
| Qwen3-VL-32B-Instruct | <u>53.4</u> | <u>37.9</u> | 62.9 | 88.7 | <u>62.0</u> | <u>54.0</u> | 74.0 | <u>100.0</u> | <u>40.6</u> | <u>24.9</u> | 51.0 | 81.1 | 53.1 | <u>37.1</u> | 65.0 | 84.6 | <u>64.2</u> | 43.4 | <u>65.7</u> | <u>100.0</u> | 51.6 | 83.0 |
| Qwen3-VL-8B-Instruct | 46.6 | 33.6 | 53.8 | 85.5 | 59.1 | 51.6 | 70.9 | 97.6 | 33.8 | 22.7 | 38.4 | 75.1 | 46.7 | 31.1 | 58.6 | 82.1 | 53.3 | 37.6 | 51.5 | 99.0 | 53.4 | 85.0 |
| Qwen2.5-VL-7B-Instruct | 30.4 | 25.9 | 41.5 | 78.7 | 39.0 | 45.2 | 59.9 | 94.0 | 17.0 | 14.8 | 29.6 | 65.7 | 34.7 | 22.5 | 46.0 | 78.7 | 33.7 | 30.6 | 33.8 | 99.5 | 15.3 | 51.0 |
| Gemma3-27B-Instruct | 33.5 | 26.5 | 45.1 | 78.3 | 39.8 | 45.0 | 62.1 | 99.4 | 12.0 | 13.7 | 21.8 | 52.8 | 40.9 | 24.5 | 52.3 | 82.6 | 44.1 | 31.7 | 49.1 | 99.5 | 23.1 | 55.0 |
| Average | 42.6 | 32.1 | 54.9 | 88.0 | 50.0 | 48.9 | 68.2 | 97.6 | 28.5 | 21.7 | 42.5 | 80.5 | 46.2 | 30.2 | 59.7 | 85.2 | 49.2 | 35.3 | 51.7 | 99.4 | 46.2 | 80.8 |

uation to fixed, verifiable checklists effectively mitigates the stochasticity of LLM judges. We therefore adopt the checklist-based protocol for the updated evaluation reported in this appendix.

Table 7: Stability analysis: Comparison of Mean Variance (over 10 runs) between Direct GT-to-Code and Checklist-based method.

| Method | Metric | Mean Var. | Improv. |
|----------------------|--------------------|--------------|----------------|
| Direct GT-to-Code | Logic Score | 8.288 | - |
| Our Checklist | Logic Score | 7.414 | 10.6% ↓ |
| Direct GT-to-Code | Param Score | 30.605 | - |
| Our Checklist | Param Score | 5.544 | 81.9% ↓ |

A.4 Judge Model Update for Code-level Evaluation

With the rapid evolution of frontier LLMs, we revisit the choice of the code-level judge in the updated evaluation. In particular, we include the newly available GPT 5.3 codex in the candidate pool and re-evaluate judge reliability to identify a more suitable model for symbolic code-level assessment.

To this end, we compare three frontier models—Claude Sonnet 4.5, Gemini 2.5 Pro, and GPT 5.3 codex—by eliciting rankings over 1,080 instances across three target models (InternVL 3.5-8B-Instruct, Qwen2.5-VL-72B-Instruct, and Qwen3-VL-235B-A22B-Instruct). We then construct a Majority-Vote (MV) ranking as a silver standard for expert consensus. As shown in Tab. 8, the judges exhibit high overall agreement (Kendall’s $W = 0.76$, mean Kendall’s $\tau = 0.63$, and mean Spearman’s $\rho = 0.67$). Moreover, the disagreements remain small in magnitude, with a

Table 8: **Judge agreement and alignment to majority vote (MV)**. All values are computed over 774 indices and reported as mean±std unless otherwise noted. *MV align* is the composite alignment to the MV ranking, defined as $\frac{1}{2}(\tau + \rho)$, where τ is Kendall’s τ and ρ is Spearman’s ρ against MV.

| Per-judge vs MV | | | |
|----------------------|--|-----------------|---------------------------|
| Code Judge | MV align ↑ | | Closest to MV (count/%) ↑ |
| Claude Sonnet 4.5 | 0.8482 ($\tau=0.8335, \rho=0.8629$) | | 276 / 35.62% |
| Gemini 2.5 Pro | 0.7556 ($\tau=0.7368, \rho=0.7743$) | | 215 / 27.80% |
| GPT 5.3 codex | 0.8517 ($\tau=0.8390, \rho=0.8643$) | | 283 / 36.58% |
| Overall (all judges) | | | |
| Kendall’s W ↑ | | 0.7618 ± 0.2482 | |
| Kendall’s τ ↑ | | 0.6315 ± 0.3668 | |
| Spearman’s ρ ↑ | | 0.6781 ± 0.3660 | |
| Rank range (mean) ↓ | | 0.58 | |

mean rank range of 0.58, indicating that most discrepancies are limited to local swaps rather than systematic inversions. We further evaluate each judge’s alignment with the MV ranking. GPT 5.3 codex achieves the highest fidelity to the consensus, with an MV align score of 0.8514 ($\tau = 0.8390, \rho = 0.8643$). On a per-instance basis, it is also the judge most frequently closest to the MV ranking, accounting for 36.58% of all cases and outperforming the other candidates.

Based on these results, we adopt GPT 5.3 codex as the code judge in the updated evaluation. This result suggests that, under the revised symbolic evaluation setting, newer code-specialized frontier models can provide a stronger proxy for consensus-based judgment.

A.5 More Discussions about the Updated Main Results

The updated main results in Tab.6 reveal a more differentiated pattern of model strengths. Among proprietary models, Claude Sonnet 4.5 achieves the highest execution rate (97.2%) and shows particular

strength on LaTeX, Gemini 3 Pro leads on code-level metrics (S_{func} , S_{para}) and performs strongly on HTML and SVG, while GPT 5.1 remains the strongest model for Python generation. The revised results also make LaTeX the clearest bottleneck in Omni-I2C, with all models showing a more pronounced performance drop on this task, especially at smaller scales. By contrast, SVG appears less difficult than in the previous version under the refined protocol, suggesting that the updated rendering and evaluation setup provides a more faithful estimate of model capability. The updated results also echo the findings in open-weight models: although scaling remains beneficial, architectural improvements can outweigh parameter count, as reflected in the consistent advantage of Qwen3-VL-32B over Qwen2.5-VL-72B and the particularly strong LaTeX performance of Qwen3-VL-235B.

Overall, these updates refine the empirical picture without changing the main takeaway, and we recommend the updated benchmark and evaluation pipeline for future use.

B Dataset Details

This part provides exhaustive details regarding our data sourcing, the principles of curation, the human-centric annotation pipeline, figure type taxonomy and data rewriting details.

B.1 Data Composition and Sourcing

The construction of Omni-I2C follows a hybrid approach, integrating high-fidelity real-world samples with logically rigorous synthetic data to balance ecological validity and structural complexity. To ensure ethical and legal rigor, all real-world instances are exclusively curated from publicly accessible platforms under permissive licenses or copyright-free terms, ensuring full compliance for academic research and redistribution.

Real-world Data Sourcing To reflect the diversity of practical applications, we systematically curated samples from 12 representative platforms. These sources span a wide range of subjects and visual modalities, providing a robust foundation for evaluating real-world performance. A comprehensive list of these platforms and their corresponding URLs is provided in Tab. 9.

Synthetic Data Integration and Refinement In addition to original collections, we incorporate selected data from existing benchmarks (Yang et al.,

2024, 2025) to further broaden our scope. Recognizing that raw synthetic data often lacks the necessary complexity or suffers from fidelity issues, we implemented a rigorous re-processing pipeline. This refinement procedure involves code refactoring and visual alignment.

Table 9: Comprehensive list of real-world data sources.

| Category | Source URL |
|----------|---|
| Python | https://matplotlib.org/stable/gallery/index.html https://networkx.org/documentation/stable/auto_examples/index.html https://plotly.com/python/ https://altair-viz.github.io/gallery/index.html https://python-graph-gallery.com/ |
| LaTeX | https://tikz.net/ https://www.overleaf.com/gallery https://texample.net/ https://ctan.mirrors.hoobly.com/graphics/pgf/contrib/pgfplots/doc/pgfplots.pdf |
| HTML | https://echarts.apache.org/examples/zh/index.html https://www.highcharts.com/demo |
| SVG | https://commons.wikimedia.org/wiki/Main_Page |

B.2 Data Collection Principles

To ensure the quality and diagnostic value of Omni-I2C, we adhere to three core curation principles during the real-world data collection process:

Structural Correspondence. We enforce strict congruence between code snippets and their rendered outputs. Every visual component must be directly traceable to specific logical blocks within the code.

Code Idiomatcity. We prioritize implementations that utilize standard and idiomatic libraries (e.g., Matplotlib in Python or standard TikZ libraries). By avoiding obscure or idiosyncratic syntax, we ensure that the evaluation focuses on the model’s fundamental perception, reasoning and coding capabilities rather than its familiarity with “edge-case” language usage, thereby minimizing confounding variables during execution.

Logical Non-triviality. Our selection process favors visualizations with high structural and compositional complexity. We deliberately exclude rudimentary geometric primitives—such as isolated triangles or squares—in favor of multi-layered scientific plots, intricate symbolic notations, and nested

UI components that require deep, multi-step reasoning to synthesize.

For synthetic data integration, our primary objective is to mitigate the data sparsity inherent in natural distributions. We strategically employ synthetic samples to supplement underrepresented or unseen categories identified in the real-world corpus, ensuring that the benchmark remains balanced across diverse subjects and programming paradigms.

B.3 Human Annotation Pipeline

The integrity of our 1,130 benchmark samples is guaranteed by a multi-stage human-in-the-loop process.

Annotator Background We recruit 9 undergraduate students majoring in Computer Science or Artificial Intelligence. All participants had undergone basic research training and possessed prior experience in Python or web development, providing the necessary technical literacy for high-quality curation.

Training and Standardization To ensure inter-annotator consistency, we conducted a rigorous training session prior to the mass annotation phase. Annotators are briefed on the ‘‘Structural Correspondence’’ requirement and trained on specific taxonomic guidelines to minimize ambiguity. A prime example of our classification protocol involves distinguishing between *Analytic Geometry* and *Function-related*. We enforced a strict mathematical definition: a *Function-related* is defined by the property that every x value corresponds to a unique y value (i.e., the vertical line test), whereas *Analytic Geometry* figures are not bound by this constraint. Cases with ambiguous classifications are adjudicated through a majority voting mechanism among senior annotators.

Following these protocols, annotators utilized a custom-built data collection platform that synchronizes code editing with real-time rendering to enforce a unified metadata schema. Upon uploading the source code and corresponding visual artifacts, annotators labeled each entry across three hierarchical dimensions: *Subject*, *Figure Type*, and *Code Type*.

B.4 Figure Type Taxonomy

To rigorously evaluate the model’s code generation performance across diverse visual scenarios, Omni-I2C incorporates a granular classification system. Currently, Omni-I2C encompasses 43 figure types,

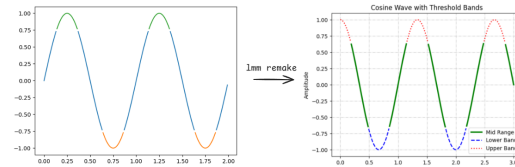


Figure 5: Examples of LLM-based code restructuring with foreground and numerical perturbations.

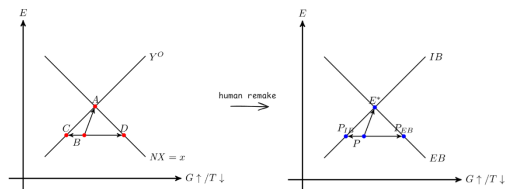


Figure 6: Comparison between original images and human-refined samples with modified key attributes.

covering core visual expressions in scientific research, engineering design, data analysis, and professional education.

The specific categories include: relationship-diagram, line-graph, molecular-formula, tables, bar-chart, plane-geometry, flow-chart, analytical-geometry, circuit, 3d-plot, equations-texts, scatter-plot, pie-chart, function-related, solid-geometry, schematic, heatmap, Area, radar-chart, Tree, Quiver, error-bar, venn-diagram, atom-model, Contour, box-chart, free-body-diagram, error-point, Density, map, Violin, Histogram, optics-ray-diagram, cell-structure, Graph, Algorithms, multi-graph, Phase-Diagram, block-diagram, physiological-process, magnetic-field-line and anatomy-diagram. Sparse or ambiguous categories are aggregated into a miscellaneous group to maintain statistical rigor.

We provide illustrative examples for these image types in Fig. 23, Fig. 24, and Fig. 25.

B.5 Data Rewriting and Leakage Mitigation

To mitigate the risks of data leakage and model memorization, we implement a rigorous data rewriting pipeline consisting of three primary stages:

First, we leverage LLMs to refactor the code’s structural and stylistic attributes, such as layout and aesthetic parameters, ensuring the code deviates from its original online form while preserving the target visualization’s core logic. Second, we manually modify key visual attributes, including label coordinates, color schemes, and textual content, to further distance the samples from their raw source. Finally, we deliberately substitute meaning-

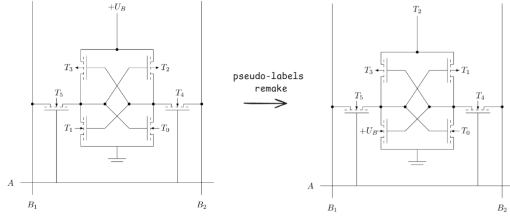


Figure 7: Visualization of counterfactual labels used to decouple perception from language priors.

```

# Role
You are an expert LaTeX and TikZ Visualization Specialist.

# Task
Rewrite the provided LaTeX code to create a NEW dataset sample.
Your goal is to prevent data leakage (memorization) while STRICTLY maintaining
the semantic logic of the '{subject}' domain and the visual structure of a
'{image_type}'.

# Constraints & Rules
1. **Safety Check (Critical):** - Analyze the content inside the '\latex_code'
tags.
- **Check 1 (Image Type):** If the code logic clearly contradicts the Image
Type '{image_type}' (e.g., code draws a generic math function plot but type is
'circuit'), output ONLY: '<TAG_MISMATCH>'.
- **Check 2 (Subject):** If the code content is completely unrelated to or
contradicts the Subject '{subject}' (e.g., Subject is 'Chemistry' but code draws a
computer network topology; Subject is 'Economics' but code draws a molecular
structure), output ONLY: '<TAG_MISMATCH>'.

2. **Rewriting Strategy (TikZ Focus):**
- **Visuals (Foreground Only):** - Change stroke/fill colors (use '\xcolor'
definitions like 'blue!80!black', 'teal'), line widths, arrow styles ('-stealth'),
and node shapes.
- **CRITICAL CONSTRAINT (Background):** **DO NOT CHANGE THE BACKGROUND
COLOR.** Do NOT use '\pagecolor{...}'. Do NOT add a background rectangle fill.
The background must remain default (white/transparent).
- **CRITICAL CONSTRAINT (Contrast):** Ensure all new colors for text,
lines, and borders have **HIGH CONTRAST** against a white background. Avoid
'yellow', 'white', 'lightgray', or extremely pale colors for thin lines or text.
- **Context:** Rename node labels to synonymous terms suitable for
'{subject}'.
- **Math:** Modify coefficients or values in equations (e.g., change
'$y=x^2$' to '$y=2x^2-1$') BUT ensure the equation remains valid LaTeX.

3. **Robustness:**
- Ensure the output is a complete, compilable LaTeX document starting with
'\documentclass' and ending with '\end{document}'.
- Do NOT omit the preamble (usepackage, tikzlibrary).

# Input Data
To ensure you read the full context, note the metadata first:
- **Subject:** {subject}
- **Image Type:** {image_type}

**Source Code:**
<latex_code>
{code}
</latex_code>

# Output
Provide ONLY the rewritten LaTeX code.

```

Figure 8: The systematic prompt utilized for guiding the LLM in code rewriting and consistency checking.

ful text with nonsensical or counter-factual strings to decouple visual perception from linguistic priors, forcing models to rely strictly on visual grounding rather than common-sense guessing.

Illustrative examples of LLM refactoring, manual modification, and counter-factual injection are shown in Fig. 5, Fig. 6, and Fig. 7, respectively. The specific prompt employed for LLM-based refactoring is detailed in Fig. 8.

C Setups

C.1 Inference Setups

We evaluate 5 proprietary and 8 open-weight models in our study. We employ greedy decoding and the maximum generation length is 16,384 tokens for all models. For Gemini 2.5 pro, Gemini 3 pro

and GPT 5.1, we configure the “thinking effort” to the “minimal”.

C.2 Evaluation Setups

Details of Visual Perceptual Fidelity. To ensure that the Visual Perceptual Fidelity (S_{img}) metric reflects the specific requirements of various technical domains, we implement a Domain-Aware Weighting Mechanism. Rather than applying a uniform average across all evaluation dimensions, this approach dynamically adjusts the importance of different visual axes based on the subject matter and figure archetype. The final score S_{img} (on a 100-point scale) is calculated by $S_{img} = (\sum s_i \times w_i) \times 10$, where $s_i \in [0, 10]$ represents the raw score for each axis and w_i denotes the percentage weight. The four assessment axes include: Composition Integrity (s_{comp}), evaluating global spatial distribution; Textual Legibility (s_{text}), assessing label and symbolic accuracy; Entity Existence (s_{entity}), ensuring critical visual primitives are present; and Stylistic Consistency (s_{style}), measuring aesthetic attributes like color and line styles.

The assignment of these weighting profiles follows a hierarchical mapping logic. The Charts & Plots profile is applied to archetypes such as line-graphs, bar-charts, pie-charts, scatter-plots, 3d-plots, and function-related visualizations. The Geometry profile encompasses plane-geometry, solid-geometry, and analytical-geometry, where coordinate-based properties are central. Archetypes defined by structural connectivity, including flow-charts, schematics, circuits, and relationship-diagrams, are mapped to the Diagrams profile. The Chemistry profile is assigned to chemistry-related subjects to prioritize molecular topology; for Math subjects, items containing tabular structures are evaluated under the **Tables** profile, while others are mapped to the Equations profile to emphasize formula rendering. For instances not covered by these archetypes, they will be processed with default profile. The specific weight configurations are detailed in Tab. 10.

The rationale for this design is rooted in human check, as different image types prioritize different visual elements—for instance, structural connectivity is paramount for molecular formulas, while grid alignment is critical for tables. The prompt template employed to evaluate the visual perceptual fidelity (S_{img}) of LLMs within the Omni-I2C framework is presented in Fig. 9.

Table 10: Domain-specific weight configurations (w_i) for S_{img} . Weights are assigned based on a mapping of figure types and subjects to prioritize domain-critical features.

| Profile | Comp. (w_{comp}) | Text (w_{text}) | Ent. (w_{entity}) | Sty. (w_{style}) |
|----------------|-------------------------|------------------------|--------------------------|-------------------------|
| Charts & Plots | 30% | 20% | 35% | 15% |
| Geometry | 30% | 15% | 35% | 20% |
| Diagrams | 30% | 20% | 35% | 15% |
| Tables | 35% | 35% | 15% | 15% |
| Equations | 30% | 35% | 20% | 15% |
| Chemistry | 30% | 20% | 35% | 15% |
| Default | 25% | 25% | 25% | 25% |

```

You are a Senior Visual QA Auditor.
Your goal is to assess the rendering alignment between a Reference image and a Candidate
Image. Focus ONLY on visual fidelity and user perception.

### 1. Context
- Domain: {subject} / {figure_type}
- Rendering Format: {code_type}

### 2. Scoring Instructions (0-10)
Assign raw scores for the following 4 dimensions based on the **Pre-defined Weights**..
Directly evaluate the images based on the following anchors and critical guidelines.

**Important: All scores must be integers (whole numbers) from 0 to 10. Do not use decimal
values.**

**Scoring Anchors:**
- 9-10: Pixel-perfect, virtually indistinguishable.
- 7-8: Minor stylistic flaws but functionally identical.
- 4-6: Significant errors in alignment or text that impact readability.
- 0-3: Structural collapse or missing key data.

**Critical Scoring Guidelines:**
- **Textual Tolerance:** Focus on semantic accuracy. **DO NOT** penalize for logical line
breaks, indentation, or hyphens used for text-wrapping (e.g., "Su-plex" vs "Suplex") unless
they alter meaning.
- **Structural Penalty:** If the Candidate exhibits "Structural Collapse" (e.g., extreme
stretching, unreadable scaling, or complete loss of layout), you **MUST** cap the scores for
Dim 1 and Dim 3 below 40.

**Evaluation Dimensions:**
* Dim 1: Composition & Spatial Integrity (Weight: {w_comp}%)
* Dim 2: Symbolic & Textual Legibility (Weight: {w_text}%)
* Dim 3: Entity Existence (Weight: {w_entity}%)
* Dim 4: Stylistic Consistency (Weight: {w_style}%)

### 3. Output Format (JSON Only)
{
  "evaluation_details": {
    "composition_integrity": { "raw_score": [integer 0-10], "comment": "Reasoning for the
score." },
    "textual_legibility": { "raw_score": [integer 0-10], "comment": "Reasoning for the
score." },
    "entity_existence": { "raw_score": [integer 0-10], "comment": "Reasoning for the
score." },
    "stylistic_consistency": { "raw_score": [integer 0-10], "comment": "Reasoning for the
score." }
  }
}

**Note: All raw_score values must be integers (0, 1, 2, ..., 10). The final weighted score
will be converted to a percentage (0-100).**

```

Figure 9: Prompt of Visual Perceptual Fidelity

Prompts for Image-level Evaluation. We detail the template prompts used in Tab. 3 in Fig. 11.

Comparison with non-LMM Metrics. To further validate the effectiveness of our LMM-based evaluation, we incorporate traditional non-LMM visual similarity metrics as baselines, including **LPIPS**, **SSIM**, and **CLIP Score**. We conduct our experiments across three representative models: Gemini 2.5 Pro, Qwen3-VL-235B-A22B-Instruct, and InternVL3.5-8B-Instruct.

Following the protocol described in the main paper, we conduct a human preference analysis for these metrics, and the results are summarized in Tab 11. As shown, while LMM-based judges may exhibit certain biases—a recognized limitation of current fully automated evaluation

paradigms—their evaluations remain more consistent with human preferences than traditional non-LMM metrics. This suggests that, despite their inherent limitations, our LMM-based metric serves as a reasonably strong and practical proxy for human judgment in the context of fully automated evaluation.

Table 11: Correlation analysis between automated metrics and human preferences. Results show that our LMM-based metric provides a more consistent proxy for human judgment compared to traditional non-LMM baselines.

| Metric | Kendall’s τ \uparrow | Spearman’s ρ \uparrow | RMSE \downarrow |
|-------------|-----------------------------|------------------------------|-------------------|
| LPIPS | 0.2717 | 0.3024 | 0.9401 |
| SSIM | 0.1155 | 0.1298 | 1.0526 |
| CLIP Score | 0.4414 | 0.4831 | 0.7994 |
| Ours | 0.8304 | 0.8681 | 0.4284 |

Details of Symbolic Functional Precision. To implement the Symbolic Functional Precision metrics, we formulate a structured evaluation prompt that leverages an LLM to serve as an Expert Code Logic Analyst. We employ the Python-specific prompt (detailed in Fig. 10) as an example to illustrate the underlying evaluation logic.

D Human Study

To assess the biological plausibility and reliability of our visual perception metrics, we conduct a human alignment study in Tab. C.2. Ten graduate researchers specializing in Computer Science and AI are recruited as annotators. To ensure high inter-rater consistency, all participants undergo a standardized training session prior to the formal annotation process.

From our benchmark, we curate a representative subset of 100 samples; each annotator independently ranked 50 randomly assigned cases. We selected *Gemini 2.5 Pro*, *Qwen3-VL-235B-A22B-Instruct*, and *InternVL3.5-8B-Instruct* for evaluation, as they represent a broad spectrum of model architectures and performance tiers, providing a clear basis for alignment testing. The annotation is performed in a double-blind manner: model identities are strictly anonymized, and the display order is randomized (see Fig. 18).

This process yielded 500 individual judgments, which are subsequently aggregated to form a consensus ground-truth ranking for the 100 samples. We employed Kendall’s τ , Spearman’s ρ , and

```

# Role
You are an expert Code Logic Analyst and Image Reconstruction Auditor. Your task is to compare two Python code snippets:
1. **Ground Truth Code (GT):** The original source code with logic comments.
2. **Generated Code (Gen):** Code generated to reproduce the image.

# Objective
Evaluate the **Functional Consistency** between the two codes across two dimensions: Logic and Parameters.
* **Focus:** Core logic, visual elements (quantity/category), data transformation methods, and key constants.
* **Ignore:** Boilerplate settings like 'figsize', 'dpi', 'savefig' paths, or minor style differences (e.g., exact color hex codes) unless they change the semantic meaning.

# Execution Protocol (Mandatory)
To ensure scoring consistency, you must follow these steps in order:
1. **GT Decomposition:** Analyze the GT code first. Identify and count all distinctive "Functional Units" (Logic) and "Key Parameters" (Data). Establish the total counts (N_{{logic}} and N_{{var}}).
2. **Gen Comparison:** Evaluate the Gen code by cross-referencing it against the specific list derived from the GT.

# Evaluation Criteria

### 1. Logic Evaluation
* **Functional Units:** Distinct logical steps, loop structures, coordinate calculations, or rendering methods.
* **Fidelity:** If Gen hardcodes a value that GT derives through logic/calculation, count it as 0 implementation for that unit.

### 2. Parameter Evaluation
* **Key Parameters:** Specific constants, coefficients, exponents, variable names, and critical coordinates.
* **Semantic Attributes:** Only evaluate attributes (like Color or Labels) if they are essential for identifying data series or categories.

# Scoring Rules
For both streams, identify the total count in GT (N) and the matched count in Gen (M):
1. **Logic Score (N_logic, M_logic):**
* **N_logic:** The total number of Functional Units identified in GT.
* **M_logic:** The count of units successfully implemented in Gen.
* **Precision:** Partial matches (e.g., correct intent but slightly flawed implementation) are allowed and scored as **0.5**. This metric is a **Float**.
2. **Parameter Score (N_var, M_var):**
* **N_var:** The total number of **Essential** Variables/Parameters identified in GT (excluding typography/spacing).
* **M_var:** The count of parameters that match exactly in Gen.
* **Precision:** Parameters must match exactly. Any mismatch results in 0. This metric must be an **Integer**.

# Output Format (JSON Only - CRITICAL)
**IMPORTANT:** You must output ONLY a valid JSON object. Do NOT include any markdown code blocks, explanations, or other text before or after the JSON. Start your response directly with {{ and end with }}.
{{
  "logic_evaluation": {
    "reasoning": "Identify missing environments, structural mismatches, or layout logic errors based on GT intent.",
    "metrics": {
      "gt_total_functions": <Integer: Total structural/logical units (N_logic) in GT>,
      "gen_implemented_functions": <Float: Units matched (M_logic) in Gen; partial match 0.5 allowed>,
      "logic_score": <Float: (gen_implemented_functions / gt_total_functions) * 100>
    }
  },
  "parameter_evaluation": {
    "reasoning": "Identify specific incorrect constants, coefficients, exponents, variable names, or coordinate errors. Note: Color is evaluated only if it carries semantic meaning; ignore minor styling.",
    "metrics": {
      "gt_total_parameters": <Integer: Total key variables/parameters (N_var) in GT>,
      "gen_implemented_parameters": <Integer: Parameters matched exactly (M_var) in Gen; mismatch is 0>,
      "parameter_score": <Float: (gen_implemented_parameters / gt_total_parameters) * 100>
    }
  }
}}
# Ground Truth Code:
{{
  "PYTHON":
  {gt_code}
}}
# Generated Code:
{{
  "PYTHON":
  {pred_code}
}}

```

Figure 10: Example Prompt of Symbolic Functional Precision

RMSE as meta-evaluation metrics to quantify the congruence between our automated scoring system and human intuition.

E Extended Analysis of Prompting Strategies

This section provides the exhaustive analysis and qualitative observations for the five prompting strategies discussed in Sec. 5.2. As established in our benchmarking, the “pixel-to-program” mapping is highly susceptible to the structure of the input prompt.

The Efficacy of Reasoning-Augmented Generation. Strategies that encourage explicit reasoning—*Self-Commentary* and *Hint-*

Enhanced—consistently deliver the highest performance gains. By asking the model to produce a visual rationale before generating code (*Hint-Enhanced*) or to interleave logic with inline comments (*Self-Commentary*), we mitigate translation errors from pixels to syntax. For example, Gemini 2.5 Pro achieves its peak S_{func} (67.67%) with *Self-Commentary*, suggesting that internalizing the “reasoning-then-coding” workflow is vital for complex figure archetypes.

The Failure of Scaffold Prompting. Contrary to its success in simpler spatial tasks, *Scaffold* prompting fails broadly on Omni-I2C. For models like InternVL 3.5-8B, S_{func} drops to nearly half of the Direct baseline. Our analysis suggests that the dot-grid overlay, while intended to aid grounding, introduces significant visual noise. This interference obscures fine-grained geometric textures and symbolic notations, which are critical for the high-precision requirements of our benchmark.

Few-Shot Prompting as a Formidable Baseline. We find that *Few-Shot* prompting remains exceptionally competitive, often rivaling or surpassing complex reasoning prompts for structured outputs like HTML and LaTeX. GPT 4.1 mini exhibits its largest gains here, exceeding the Direct baseline by over 11 points in overall execution. This suggests that for tasks with rigid syntactic patterns, providing in-context examples is more effective than prompting for abstract reasoning.

Robustness and Stability Gaps. The staggering variance in performance—where a single change in strategy can swing scores by 15–20%—highlights a fundamental lack of inherent robustness in current LMMs. This “prompt-dependency” indicates that models are not yet capable of stable, zero-shot code reconstruction for complex visual inputs, necessitating carefully engineered “wrappers” to elicit frontier-level performance. We provide the prompting templates in Fig. 26.

F Extended Error Analysis

This section provides the exhaustive definitions and specific failure manifestations for the five-tier taxonomy introduced in Sec. 5.3. These criteria are strictly followed during our manual audit of 300 sampled instances to ensure a consistent and objective error analysis across all evaluated LMMs.

Code-level Logic Errors. This category encompasses instances where the generated code is syntactically valid and executable but fails to implement

Table 12: Studies on different prompting methods. The best results for each model are highlighted in **bold**. Values exceeding and falling below the corresponding average are highlighted in red and blue, respectively.

| Model | Method | Overall(1130) | | | | HTML(182) | | | | LaTeX(274) | | | | Python(364) | | | | SVG(210) | | | | SMILES(100) | |
|-----------------------------|-----------------|---------------|-------------------|-------------------|------------------|---------------|-------------------|-------------------|------------------|--------------|-------------------|-------------------|------------------|--------------|-------------------|-------------------|------------------|--------------|-------------------|-------------------|------------------|--------------|--------------|
| | | Exec | S _{func} | S _{para} | S _{img} | Exec | S _{func} | S _{para} | S _{img} | Exec | S _{func} | S _{para} | S _{img} | Exec | S _{func} | S _{para} | S _{img} | Exec | S _{func} | S _{para} | S _{img} | Exec | T.S. |
| Gemini2-5 Pro | Direct | 88.32 | 66.39 | 56.49 | 69.70 | 99.45 | 78.78 | 85.09 | 74.57 | 69.34 | 53.46 | 46.12 | 52.40 | 95.05 | 71.69 | 59.19 | 80.91 | 88.10 | 63.34 | 40.57 | 68.62 | 96.00 | 81.09 |
| | Few-Shot | 86.73 | 63.80 | 54.10 | 68.74 | 98.90 | 77.76 | 81.48 | 70.29 | 73.36 | 54.95 | 46.44 | 56.81 | 92.86 | 68.17 | 56.01 | 80.44 | 78.10 | 55.67 | 37.05 | 62.69 | 97.00 | 80.93 |
| | Scaffold | 83.01 | 50.62 | 42.57 | 51.65 | 99.45 | 72.62 | 77.47 | 63.78 | 62.04 | 39.17 | 34.65 | 35.43 | 86.81 | 46.19 | 38.13 | 54.46 | 91.43 | 54.19 | 30.33 | 57.43 | 79.00 | 52.41 |
| | Self-Commentary | 92.39 | 67.67 | 57.17 | 69.31 | 92.31 | 68.58 | 75.41 | 61.99 | 86.13 | 62.28 | 53.76 | 57.91 | 95.33 | 70.72 | 58.53 | 79.93 | 93.81 | 68.62 | 43.43 | 72.14 | 96.00 | 77.65 |
| | Hint-Enhanced | 87.79 | 66.16 | 54.68 | 67.81 | 98.90 | 79.91 | 83.78 | 72.26 | 72.26 | 53.95 | 45.12 | 50.53 | 93.41 | 69.64 | 56.60 | 78.82 | 90.00 | 64.16 | 38.60 | 67.41 | 85.00 | 75.61 |
| GPT 4.1 mini | Direct | 81.33 | 45.80 | 35.83 | 53.22 | 100.00 | 60.35 | 65.86 | 61.61 | 43.07 | 22.97 | 17.25 | 26.72 | 92.86 | 53.42 | 42.70 | 70.88 | 91.43 | 49.76 | 22.16 | 49.92 | 89.00 | 53.03 |
| | Few-Shot | 92.03 | 62.83 | 50.91 | 59.14 | 99.45 | 80.50 | 83.91 | 63.39 | 86.13 | 55.67 | 45.28 | 50.86 | 93.68 | 65.17 | 53.48 | 72.29 | 93.81 | 52.82 | 25.20 | 43.47 | 85.00 | 52.16 |
| | Scaffold | 82.83 | 36.20 | 28.97 | 43.50 | 98.35 | 55.08 | 55.96 | 54.27 | 45.99 | 18.68 | 14.25 | 21.97 | 93.96 | 37.93 | 33.84 | 55.49 | 94.76 | 39.70 | 16.35 | 41.47 | 90.00 | 29.50 |
| | Self-Commentary | 84.96 | 48.50 | 37.95 | 57.92 | 100.00 | 63.51 | 67.25 | 63.21 | 54.74 | 27.11 | 21.20 | 34.63 | 94.51 | 55.69 | 46.09 | 73.12 | 91.90 | 50.92 | 20.31 | 57.40 | 91.00 | 54.84 |
| | Hint-Enhanced | 89.29 | 64.19 | 52.87 | 59.50 | 100.00 | 84.17 | 84.74 | 65.80 | 78.10 | 54.47 | 45.04 | 42.96 | 93.68 | 66.59 | 55.07 | 73.52 | 87.62 | 55.42 | 31.65 | 51.31 | 88.00 | 57.34 |
| Qwen3-VL-235B-A22B-Instruct | Direct | 86.46 | 54.53 | 38.98 | 54.93 | 93.41 | 58.02 | 60.16 | 50.59 | 71.17 | 42.04 | 38.55 | 45.23 | 89.56 | 66.41 | 33.77 | 64.40 | 90.95 | 47.19 | 30.23 | 54.93 | 95.00 | 67.99 |
| | Few-Shot | 87.35 | 56.70 | 39.83 | 56.85 | 98.90 | 65.57 | 63.18 | 55.85 | 75.55 | 45.51 | 40.54 | 49.93 | 91.76 | 68.55 | 34.50 | 66.62 | 82.86 | 43.07 | 27.91 | 49.81 | 92.00 | 67.45 |
| | Scaffold | 82.92 | 42.16 | 28.85 | 42.27 | 97.80 | 58.28 | 57.63 | 46.16 | 62.41 | 30.15 | 27.08 | 31.03 | 88.74 | 47.43 | 21.40 | 48.52 | 85.24 | 34.71 | 19.15 | 42.72 | 86.00 | 37.89 |
| | Self-Commentary | 87.26 | 55.15 | 39.61 | 55.90 | 96.15 | 60.31 | 59.82 | 52.59 | 71.90 | 41.37 | 37.73 | 46.35 | 92.86 | 68.16 | 36.75 | 66.76 | 86.67 | 46.09 | 29.50 | 52.39 | 94.00 | 68.25 |
| | Hint-Enhanced | 85.84 | 59.27 | 41.94 | 55.25 | 97.80 | 67.40 | 67.17 | 55.50 | 64.96 | 43.49 | 38.82 | 42.26 | 90.66 | 72.10 | 37.11 | 65.62 | 90.48 | 50.57 | 32.52 | 54.00 | 94.00 | 66.86 |
| InternVL3.5-8B-Instruct | Direct | 71.77 | 23.43 | 20.75 | 28.23 | 97.25 | 40.01 | 42.21 | 38.87 | 40.88 | 11.49 | 11.35 | 17.52 | 71.98 | 22.57 | 23.05 | 34.22 | 80.95 | 26.13 | 10.44 | 22.60 | 90.00 | 50.63 |
| | Few-Shot | 71.15 | 22.02 | 21.56 | 28.08 | 100.00 | 43.99 | 44.18 | 42.53 | 54.01 | 16.93 | 17.61 | 25.09 | 70.60 | 23.05 | 23.60 | 34.98 | 60.00 | 7.82 | 3.56 | 7.51 | 91.00 | 50.00 |
| | Scaffold | 67.70 | 11.64 | 9.65 | 11.88 | 96.15 | 18.36 | 16.68 | 14.38 | 40.15 | 5.62 | 4.44 | 6.23 | 62.64 | 11.47 | 12.64 | 16.53 | 80.00 | 13.96 | 5.17 | 9.02 | 84.00 | 13.02 |
| | Self-Commentary | 76.81 | 25.01 | 21.77 | 29.39 | 97.25 | 39.57 | 42.76 | 37.96 | 47.08 | 12.36 | 12.47 | 18.32 | 76.65 | 24.12 | 24.29 | 35.72 | 91.43 | 30.44 | 11.34 | 25.43 | 91.00 | 49.75 |
| | Hint-Enhanced | 72.65 | 25.21 | 21.36 | 28.52 | 98.90 | 41.97 | 42.21 | 38.53 | 48.54 | 12.52 | 13.01 | 19.55 | 71.98 | 24.54 | 23.05 | 34.81 | 77.14 | 28.42 | 11.24 | 20.67 | 84.00 | 45.25 |
| Average | Direct | 81.97 | 47.54 | 38.02 | 51.52 | 97.53 | 59.29 | 63.33 | 56.41 | 56.11 | 32.49 | 28.32 | 35.47 | 87.36 | 53.52 | 39.68 | 62.60 | 87.86 | 46.60 | 25.85 | 49.02 | 92.50 | 63.41 |
| | Few-Shot | 84.31 | 51.34 | 41.60 | 53.20 | 99.31 | 66.95 | 68.19 | 58.02 | 72.26 | 43.27 | 37.47 | 45.67 | 87.22 | 56.23 | 41.90 | 63.58 | 78.69 | 39.84 | 23.43 | 40.87 | 91.25 | 62.64 |
| | Scaffold | 79.12 | 35.16 | 27.51 | 37.32 | 97.94 | 51.09 | 51.94 | 44.65 | 52.65 | 23.41 | 20.10 | 23.67 | 83.04 | 35.76 | 26.50 | 43.75 | 87.86 | 35.64 | 17.75 | 37.66 | 84.75 | 33.20 |
| | Self-Commentary | 85.35 | 49.08 | 39.12 | 53.13 | 96.43 | 57.99 | 61.31 | 53.94 | 64.96 | 35.78 | 31.29 | 39.30 | 89.84 | 54.67 | 41.41 | 63.88 | 90.95 | 49.02 | 26.14 | 51.84 | 93.00 | 62.62 |
| | Hint-Enhanced | 83.89 | 53.71 | 42.71 | 52.77 | 98.90 | 68.36 | 69.47 | 58.02 | 65.97 | 41.11 | 35.50 | 38.83 | 87.43 | 58.22 | 42.96 | 63.19 | 86.31 | 49.64 | 28.50 | 48.35 | 87.75 | 61.27 |

the intended visual logic. Common manifestations include the misconfiguration of library-specific parameters (e.g., incorrect Matplotlib axis scales or plot types) and the inclusion of non-functional code blocks that, while error-free, result in blank or semantically distorted renderings. This failure mode highlights a “logic-intent” gap, where the model maintains syntactic fluency but cannot map visual requirements to the appropriate programmatic attributes.

Textual Precision Errors. These failures center on Optical Character Recognition (OCR) and string formatting accuracy within the code. Typical issues include the misidentification of alphanumeric characters—particularly in high-density labels—and the misparsing of complex symbolic notations such as LaTeX subscripts, superscripts, and Greek letters. Such errors lead to illegible or factually incorrect text in the final rendered output.

Entity Integrity Errors. This pertains to the identification and reconstruction of discrete visual primitives. Models frequently omit critical components (e.g., specific data points, legend keys, or axis ticks) or misinterpret the nature of a visual object (e.g., treating a data line as a grid line). These failures result in an incomplete or structurally altered representation of the source data within the synthesized code.

Colorimetric Accuracy Errors. This category assesses the model’s ability to extract color data and map it to its corresponding semantic entities. Errors manifest as generated color codes (Hex or RGB) that deviate significantly from the ground truth. This reflects a fundamental limitation in the

model’s color perception or its ability to precisely translate perceived palettes into programmatic constants.

Spatial Layout Errors. This category addresses the global topological structure and coordinate-based positioning defined by the code. Failures such as distorted aspect ratios, erroneous coordinate transformations, or incorrect entity arrangement reveal a deficit in translating complex visual hierarchies into precise spatial constraints. These errors often represent a failure in high-level spatial reasoning rather than low-level code syntax.

As illustrated in Fig. 20, despite variations in the fundamental capabilities of different models, the error distribution within specific languages exhibits cross-model consistency. This suggests that error patterns in code generation are intrinsically linked to the target language.

Specifically, for low-level descriptive languages such as SVG, the frequency of Layout- and Text-related errors is notably higher compared to other categories. This stems primarily from SVG’s lack of high-level semantic encapsulation; when processing long sequences of numerical coordinates, the model often loses its grasp of global spatial structure. As illustrated in Fig. 21a, although Claude Sonnet 4.5 correctly generated all entities and text with accurate colors, it failed to properly reconstruct the spatial arrangement and relational logic between these entities.

Conversely, in HTML tasks, the language’s high syntactic tolerance often causes errors to manifest as implicit rendering anomalies. Consequently, compared to languages with stricter syntax con-

straints, HTML tasks exhibit a statistically higher incidence of Code-related errors. For instance, as shown in Fig. 21b, the omission of an explicit height definition for the parent container in the output code resulted in the height collapse of internal color blocks. While the image rendered successfully, this logical coding error caused the final visualization to diverge significantly from the ground truth.

Simultaneously, the dominance of Entity-related errors in LaTeX and Python underscores the challenges models face in defining entities when utilizing complex macros and functional tools. As demonstrated in Fig. 21c, while the model successfully generated the majority of the content, it omitted specific components, such as the resistor present in the original diagram.

In conclusion, future advancements in MLLMs must prioritize enhancing fine-grained spatial geometric reasoning capabilities to mitigate layout deviations, while simultaneously improving precise control over code generation.

F.1 Code-related Error

To facilitate a more precise analysis of the execution failures encountered during our evaluation, we conduct a comprehensive diagnostic study on cases from the Omni-I2C benchmark where models failed to produce valid renderings. For each programming language, we developed a specialized taxonomy to categorize the root causes of these errors. By decoupling surface-level failures into granular dimensions, we expose the specific structural and logic-synthesis bottlenecks of current LMMs.

For Python-based visualizations, we categorize exceptions into four hierarchical levels. At the structural level, we identify Syntax & Parsing Failures and scope issues like Name Resolution & Import Failure. Interaction with library interfaces is monitored through Attribute Access Violations and Argument Specification Errors. Within the data logic itself, we distinguish between Type Semantics Violations, Value Domain Violations, and the Shape & Dimensionality Mismatches prevalent in tensor operations. Finally, execution context issues are captured by Resource & I/O Failures, Backend & Environment Limitations, and miscellaneous Undefined / Other Runtime Errors.

In the domain of SVG vector graphics, our taxonomy addresses the dual nature of the format as both a structured document and a graphical description. Failures in the XML layer are defined by

Well-formedness Violations and Entity & Encoding Errors. Semantic issues within the graphic description include Geometric & Viewport Mismatches, where elements fall outside the canvas, and Attribute & Value Invalidity, such as malformed path strings. We also account for rendering constraints including Renderer Capability Limits, Typography & Asset Failures, and other Undefined/Other SVG Errors. This systematic attribution allows us to isolate whether a “blank output” stems from a parsing collapse or a subtle geometric miscalculation.

The compilation of LaTeX graphics involves a sensitive interplay between macro expansion and coordinate geometry, leading to eight diagnostic categories. Syntactic failures are identified as Token & Delimiter Errors or Command & Macro Definition issues. Environmental conflicts often arise from Dependency & Package omissions or Engine-Specific Conflicts between compilers like pdfLaTeX and XeLaTeX. Regarding drawing logic, we track Coordinate & Unit Violations and Visual/Logical Nullity—cases where valid code produces no visual output. Toolchain-related failures are captured through Externalization & Conversion errors and Miscellaneous Runtime exceptions.

In contrast, for HTML/CSS and SMILES, we employ a flattened error attribution logic due to their distinct execution characteristics. For HTML/CSS, the high fault tolerance of web browsers often results in “silent failures” (e.g., blank pages or invisible elements) that do not trigger explicit exceptions, making granular attribution less reliable. Similarly, failures in SMILES are predominantly monolithic, centered almost exclusively on Invalid SMILES Syntax. For these domains, we primarily report the binary success of the rendering or parsing process rather than a detailed sub-categorization, ensuring that our qualitative analysis remains grounded in observable execution data.

The aggregated distribution of rendering failures across all evaluated models, as illustrated in Fig. 22, reveals consistent diagnostic patterns despite individual variations in model behavior. Globally, LaTeX-TikZ stands out as the most error-prone language for all tested models, highlighting the significant challenge of synthesizing code that requires both strict syntactic precision and complex macro expansion. Within the LaTeX-TikZ domain, *Token & Delimiter Errors*, such as unbalanced braces or missing math-mode delimiters, constitute the majority of failures. This is frequently followed by

Dependency & Package issues, where models invoke specialized commands without including the necessary macro packages, suggesting that LMMs struggle with the brittle, state-dependent nature of TeX compilation.

In the Python domain, the error distribution highlights a clear gap between syntactic knowledge and logical reasoning. *Shape & Dimensionality Mismatch* is the most prevalent failure mode, followed by *Argument Specification Error*, indicating that while models can often generate valid Python syntax, they frequently falter when reasoning about the underlying tensor dimensions or specific function signatures required for scientific visualization. This phenomenon is further influenced by model scale; by bifurcating the evaluated models into two cohorts—those above and below the 38B parameter threshold—we observe that smaller models exhibit a markedly higher frequency of *Shape & Dimensionality Mismatches* and *Value Domain Violations* compared to their larger counterparts. This suggests that increased parameter scale is a critical factor in developing the internal logic necessary to handle complex data structures and numerical constraints.

The failure modes in SVG generation provide further insight into the divergent capabilities of different model architectures and training paradigms. For closed-source models, *Geometric & Viewport Mismatch* is the leading cause of failure, implying that while these models successfully generate well-formed XML, they struggle to align graphical elements accurately within defined canvas boundaries. In contrast, open-source models are primarily hindered by *Well-formedness Violations*, frequently failing at the basic structural level of the XML document. This distinction is particularly evident when comparing specific model families: within the SVG tasks, *Well-formedness Violations* are the dominant failure mode for the Qwen3-VL family, whereas the InternVL3.5 family aligns more closely with the behavior of closed-source models, where *Geometric & Viewport Mismatch* represents the primary bottleneck.

Ultimately, these findings underscore a broader hierarchical trend in model maturity. For more accessible formats like SVG and Python, the challenge for advanced models has shifted from basic syntactic integrity to complex spatial reasoning and structural alignment. However, for languages with high-entry barriers like LaTeX-TikZ, even the most capable models remain susceptible to fundamental

parsing collapses. This taxonomy-based analysis confirms that while LMMs are approaching syntactic competence, mastering the underlying logic and environmental constraints of executable code remains a formidable barrier across all scales and families.

G Case Study

We conduct a qualitative analysis comparing 13 models on our benchmark. We select **Gemini 3 Pro** and **Qwen3-VL-235b-a22b-Instruct** for a detailed comparative case study due to their representative performance characteristics.

Case study in Python We examine a representative Python code generation scenario shown in Fig. 12 to analyze the reasoning behind these scores. In terms of S_{img} (Visual Fidelity), Gemini 3 Pro demonstrates an advantage over Qwen3-VL-235b-a22b-Instruct by maintaining superior global visual consistency. The S_{func} metric acutely identifies critical mathematical errors in frequency estimation within both generated code snippets, while further highlighting a fundamental gap in logical depth: Gemini 3 Pro successfully implements complex threshold masking for segmentation, whereas Qwen3-VL-235b-a22b-Instruct relies on a naive overlaying of curves, failing to comprehend the underlying conditional logic of the image. Moreover, fine-grained detection via S_{para} reveals that Qwen3-VL-235b-a22b-Instruct exhibits lower precision in parameter control, incurring more errors in stylistic details such as alpha compared to Gemini 3 Pro.

Case study in LaTeX We examine a representative LaTeX code generation scenario shown in Fig. 13 to analyze the reasoning behind these scores. In the LaTeX-based plotting task, Gemini 3 Pro achieves an S_{img} score of 81, surpassing Qwen3-VL-235b-a22b-Instruct (71) in terms of structural fidelity. The S_{func} evaluation characterizes a shared limitation in semantic LaTeX implementation: both models fail to invoke professional frameworks such as PGFPlots or its associated fill-between library, defaulting to primitive TikZ operations and manual coordinate systems. Specifically, Gemini 3 Pro employs Bezier curves integrated with region clipping for segmented filling, whereas Qwen3-VL-235b-a22b-Instruct resorts to hard-coded vertices, yielding functionally inconsistent trajectories. Moreover, S_{para} metrics indicate

that Qwen3-VL-235b-a22b-Instruct suffers from significant deviations in numerical parameters, including scaling and sampling density.

Case study in HTML We examine a representative HTML code generation scenario shown in Fig. 14 to analyze the reasoning behind these scores. For the HTML/ECharts task, Gemini 3 Pro achieves an S_{func} score of 88, significantly surpassing Qwen3-VL-235b-a22b-Instruct (56). Although Gemini 3 Pro utilizes static SVG primitives instead of the standard ECharts framework, it successfully reconstructs essential functional components, including titles, legends, and multi-series markers. Conversely, Qwen3-VL-235b-a22b-Instruct fails to include necessary script dependencies and critical y-axis metadata, resulting in a fragmented structural representation. Furthermore, fine-grained S_{para} evaluations highlight a fundamental error in data scaling by Qwen3-VL-235b-a22b-Instruct, which fails to align the y-axis range with the ground truth rainfall metrics.

Case study in SVG We examine a representative SVG code generation scenario shown in Fig. 15 to analyze the reasoning behind these scores. In the evaluation of SVG vector graphics generation, both Gemini 3 Pro (S_{img} : 85) and Qwen3-VL-235b-a22b-Instruct (S_{img} : 80) demonstrate exceptional capabilities in structural decomposition and primitive reconstruction. Both models achieved an S_{func} score of 94, indicating a shared proficiency in accurately identifying and mapping core functional units, including irregular polygon paths, closed regions, and textual labels. However, in-depth parametric assessment reveals a significant performance divergence in numerical precision: Gemini 3 Pro achieves an S_{para} score of 69, successfully aligning textual content, ViewBox attributes, and critical geometric coordinates with high fidelity. In contrast, Qwen3-VL-235b-a22b-Instruct yields an S_{para} of 47, exhibiting a substantially lower success rate in matching complex path details, fill color values, and specific attribute groups.

Furthermore, we provide qualitative rendering samples for all models benchmarked in our primary performance table. As illustrated in Fig. 16, Fig. 17, these examples offer a direct comparison of model outputs on specific data instances.

H Judge Alignment and Metric Validation

This section details the methodologies and rationales used to ensure our automated evaluation tracks—both symbolic and visual—align with objective truth and human judgment.

H.1 Symbolic Track: Code-level Reliability

For the symbolic track, we evaluate the reliability of our code judges (Claude 4.5 Sonnet, Gemini 2.5 Pro, and GPT 4.1) by measuring their alignment with the Majority Vote (MV) ranking. As shown in Tab. 8, we employ a suite of metametrics to quantify both individual judge accuracy and group consensus. Kendall’s τ and Spearman’s ρ : These are non-parametric rank correlation coefficients. τ measures the proportion of “concordant pairs” (pairs of items ranked in the same order by both the judge and MV), making it highly sensitive to even minor swaps in ranking. ρ measures the strength of the monotonic relationship between ranks. MV align: Defined as $\frac{1}{2}(\tau + \rho)$, this composite metric serves as our primary indicator of alignment. We adopt this average to balance the strict ordinal consistency of τ with the rank-intensity correlation of ρ . As indicated in Tab. 8, Claude 4.5 Sonnet achieves the highest MV align (0.8436), justifying its selection as the official symbolic evaluator. Closest to MV (count/%) : This represents the frequency with which a specific judge’s ranking is the most similar to the majority consensus on an instance-by-instance basis. This metric exposes judge “outliers” and confirms that Claude 4.5 Sonnet is the most consistently representative of the collective “wisdom of the crowd” (61.15% frequency). Kendall’s W (Coefficient of Concordance): This measures the agreement among all judges simultaneously. A value of 0.7405 indicates a strong consensus across the evaluated frontier models, suggesting that “quality” in Image-to-Code tasks is an objective property that these models can identify consistently.

H.2 Perceptual Track: Image-level Alignment

To validate the S_{img} metric, we compare the LMM judge outputs against a “ground truth” established by 10 human annotators. The results in Tab. 3 demonstrate the superiority of the Omni-I2C prompting strategy. Ordinal Alignment (τ and ρ): Since S_{img} is intended to act as a proxy for “at-a-glance” human preference, it is critical

that the model-generated scores result in the same ranking order as human consensus. Our results show that the GPT 4.1 + Omni-I2C configuration reaches a peak Kendall’s τ of 0.8304, significantly outperforming deduction-based alternatives like Chart2Code. Root Mean Square Error (RMSE): Unlike the symbolic track, which is purely rank-based, S_{img} is a scalar value. We use RMSE to quantify the absolute deviation between the normalized LMM scores and human-assigned preference scores. A lower RMSE (0.4284 for Omni-I2C) indicates that the judge is well-calibrated—meaning its “score magnitude” matches human perception of quality, rather than just the “relative order.” Rationale for Metric Selection: Robustness to Diverse Modalities: By using rank correlation (τ , ρ), we ensure the metric remains valid across different image types (e.g., SVG vs. LaTeX) where absolute visual differences might vary in scale. Sensitivity to Structural Failures: The inclusion of RMSE ensures that when a model produces a catastrophic layout failure, the judge penalizes it with a score magnitude that reflects the severity perceived by a human, rather than just ranking it “last”. The consistent performance of our multi-dimensional prompt across different judges (Qwen-VL-Max and Qwen3-VL-Plus) confirms that a structured rubric focusing on Style, Layout, Elements, and Text is the most reliable way to stabilize perceptual evaluation in LMMs.

To further investigate the reliability of the structured rubric across different visual aspects, we conducted a fine-grained alignment experiment targeting the four specific sub-dimensions: *Composition*, *Text*, *Entity*, and *Style*. We recruit 16 professional annotators, each assigned to evaluate and rank 90 items across the four specific dimensions. We then calculate the correlation between the LMM-based scores and human rankings. This analysis is based on the samples generated by Gemini 2.5 Pro, Qwen3-VL-235B-A22B-Instruct, and InternVL3.5-8B-Instruct.

We then calculated the correlation between the LMM-based scores and human rankings using Kendall’s τ , Spearman’s ρ , and RMSE. As summarized in Tab. 13, the LMM judge demonstrates a remarkably strong correlation with human judgment in Composition ($\tau = 0.7837$) and Style ($\tau = 0.7367$). These results suggest that for macroscopic visual features—such as global spatial arrangements and overall aesthetic coherence—the LMM acts as a highly reliable judge that closely

aligns with human visual perception.

Table 13: Human alignment analysis across the four sub-dimensions of S_{img} . Results are based on evaluations by 16 professional annotators.

| Dimension | Kendall’s | Spearman’s | RMSE |
|-------------|-----------------|-----------------|--------------|
| | $\tau \uparrow$ | $\rho \uparrow$ | \downarrow |
| Composition | 0.7837 | 0.8222 | 0.4110 |
| Text | 0.6889 | 0.7225 | 0.5267 |
| Entity | 0.6531 | 0.6986 | 0.5167 |
| Style | 0.7367 | 0.7791 | 0.4491 |

I Evaluation Metrics for Molecular Structures

To ensure a robust evaluation of chemical structures, we convert SMILES strings into molecular objects using the RDKit library. We specifically employ *Morgan Fingerprints* with a radius of 2 and a bit-vector size of 2048. The similarity between the ground-truth (G) and predicted (P) molecules is then quantified via the Tanimoto coefficient:

$$\text{Tanimoto}(G, P) = \frac{\mathbf{v}_G \cdot \mathbf{v}_P}{\|\mathbf{v}_G\|^2 + \|\mathbf{v}_P\|^2 - \mathbf{v}_G \cdot \mathbf{v}_P} \quad (1)$$

where \mathbf{v}_G and \mathbf{v}_P represent the fingerprint vectors.

To further elucidate the superiority of Morgan Fingerprint-based Tanimoto similarity over traditional string-based metrics, we present three representative cases encountered during our evaluation. These examples highlight the necessity of capturing chemical topology rather than mere syntactical sequences. The detailed examples show in Fig. 19

Case 1: Semantic Equivalence under Canonicalization Variance The inherent non-uniqueness of SMILES notation often leads to multiple valid strings for a single molecular structure. As shown in our first example (Ethanol, GT: CCO, Gen: OCC), string-based metrics like BLEU or Edit Distance would penalize the model for reversing the atom indexing order. In contrast, our metric yields a perfect Tanimoto score of 1.00, as both strings map to the identical molecular fingerprint. This confirms the metric’s ability to recognize functional identity regardless of the generation sequence.

Case 2: Topological Sensitivity to Atomic Hybridization Small character-level differences in SMILES strings can signify profound changes in chemical properties that are easily overlooked by superficial text metrics. In the second case, we

compare a benzene ring (c1ccccc1) with its saturated counterpart, cyclohexane (C1CCCCC1). While the text-level structure appears nearly identical (differing only by character casing), the Tanimoto similarity based on Morgan Fingerprints drops to **0.00**. This total divergence correctly reflects the fundamental loss of aromaticity and the shift in carbon hybridization from sp^2 to sp^3 —a critical structural shift that renders the reconstructed molecule chemically distinct from the ground truth. String-based metrics, such as Edit Distance or standard token matching, often fail to penalize these “semantic collapses” heavily enough, potentially masking the severity of the prediction error.

Case 3: Fidelity of Complex Functional Groups

In more complex reconstructions, such as Aspirin (GT: CC(=O)Oc1ccccc1C(=O)O), the model might correctly generate the scaffold but fail on localized functional groups. When the acetoxy group is incorrectly reduced to a methoxy group (Gen: COc1ccccc1C(=O)O), the Tanimoto metric provides a nuanced penalty (Score: 0.58), reflecting the partial structural loss. This demonstrates that the metric serves as a robust proxy for human expert judgment, as it effectively quantifies the impact of topological errors on the overall molecular integrity.

```
You are an expert judge at evaluating the visual fidelity and technical consistency of diverse image types. The first image (reference image) is rendered from ground truth code (such as Python, LaTeX, SVG, HTML, or SMILES), and the second image (AI-generated image) is rendered from code generated by an AI assistant. Your task is to score how well the AI-generated output reproduces the ground truth reference.

### Scoring Methodology:
The AI-generated image's score is based on the following criteria, totaling a score out of 100 points:
1. Visual Components and Elements (20 points)
2. Layout and Structure (10 points)
3. Text and Semantic Content (20 points)
4. Data and Technical Accuracy (20 points)
5. Style and Aesthetics (20 points)
6. Clarity and Rendering Quality (10 points)

### Evaluation:
Compare the two images head to head and provide a detailed assessment. Use the following format for your response:

Comments:
- Visual Components: $your comment and subscore}
- Layout and Structure: $your comment and subscore}
- Text and Semantic Content: $your comment and subscore}
- Data and Technical Accuracy: $your comment and subscore}
- Style and Aesthetics: $your comment and subscore}
- Clarity and Rendering Quality: $your comment and subscore}

Score: $your final score out of 100}
```

(a) Prompt of Chartmimic

```
You are a helpful assistant. Please evaluate the similarity between the **first image (the Ground Truth/GT image)** and the **second image (the Generated image)**. The **GT image** is rendered from original reference code, while the **generated image** is rendered from code provided by an AI assistant (covering formats like Python, LaTeX, SVG, HTML, or SMILES). Consider factors such as the overall appearance, technical/structural accuracy (including numerical values, mathematical notations, or chemical connectivity), colors, shapes, positions, and other visual elements.

Begin your evaluation by providing a short explanation assessing how precisely the **Generated image** reproduces both the aesthetic and functional details of the **GT image**. Be as objective as possible. After providing your explanation, you must rate the similarity on a scale of 1 to 10 by strictly following this format: "Rating: [[rating]]", for example: "Rating: [[5]]".
```

(b) Prompt of Plot2code

```
You are an exceptionally strict and meticulous image auditor. Your task is to evaluate the fidelity of a 'Generated Image' (the second image) against a 'Ground Truth Image' (the first image).

Please process the evaluation following these specific Chain of Thought (CoT) and Self-Reflection steps:
1. Initial Component Deconstruction
2. Deep Numerical & Logical Comparison
3. Visual & Style Audit
4. Mandatory Self-Reflection: Challenge your own initial observations.
5. Final Evidence Synthesis

Return ONLY a single JSON object with:
- "thought": A detailed internal monologue.
- "score": An integer from 0 to 100.
- "reason": A concise expert summary.
```

(c) Prompt of Self-Reflexion

```
You are an exceptionally strict and meticulous image analyst. Your task is to evaluate the visual similarity of two images. You must be extremely critical. Any deviation, no matter how small, must be penalized heavily. A perfect score is reserved only for images that are visually indistinguishable to the human eye. Your analysis must be based solely on the visual information in the images provided.

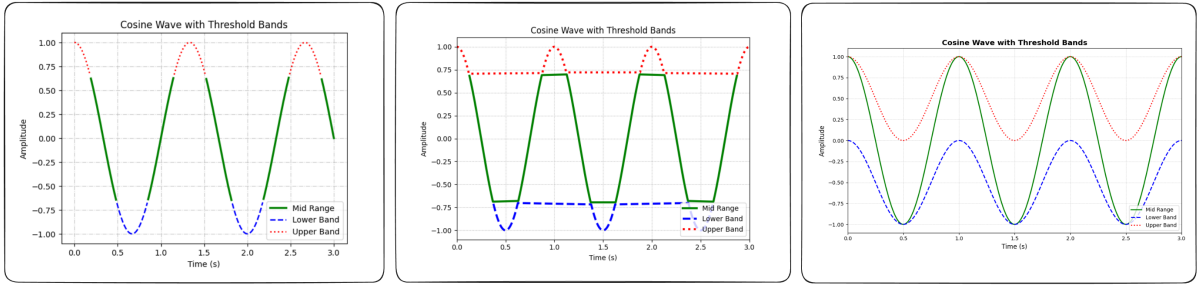
Task:
Compare the Ground Truth Image and the Generated Image. Based ONLY on their visual information, evaluate their similarity.

Evaluation Rules:
1. Strictness is Key: Start with a perfect score of 100 and deduct points for EVERY visual difference, including but not limited to: chart type, data points, colors, line styles, markers, labels (content, font, and position), titles, legends, axes (limits, ticks, scaling), layout, aspect ratio, and any other visual element.
2. Identical Means Identical: A score of 100 is ONLY for images that are pixel-perfect or visually indistinguishable. Even a tiny difference in line thickness or a single different pixel color must result in a lower score.
3. Heavy Penalties: Apply significant penalties for noticeable differences. For example, a different color map or a missing legend should lead to a large deduction.

Output Format:
Return ONLY a single JSON object with two keys: "score" ( an integer from 0 to 100) and "reason" (a concise, expert analysis in English, detailing every detected difference that justifies the score deduction). Do not include any other text, markdown, or explanations outside the JSON object.
```

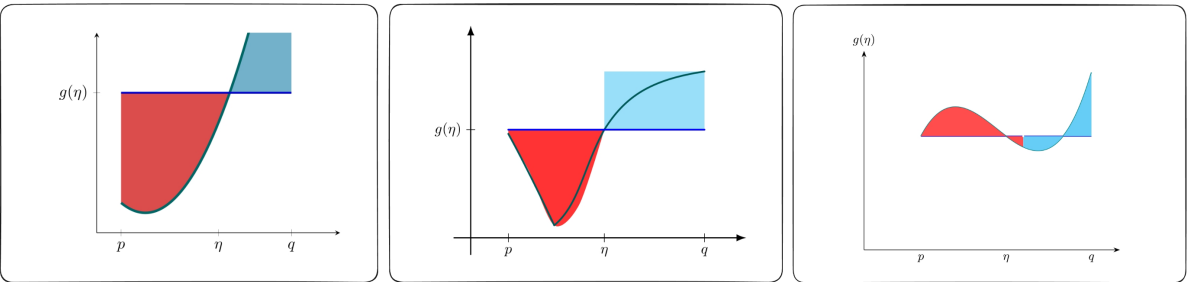
(d) Prompt of Chart2code

Figure 11: Comparison of prompts used in different methods. Since prompt images contain text, a vertical layout ensures readability within a single column.



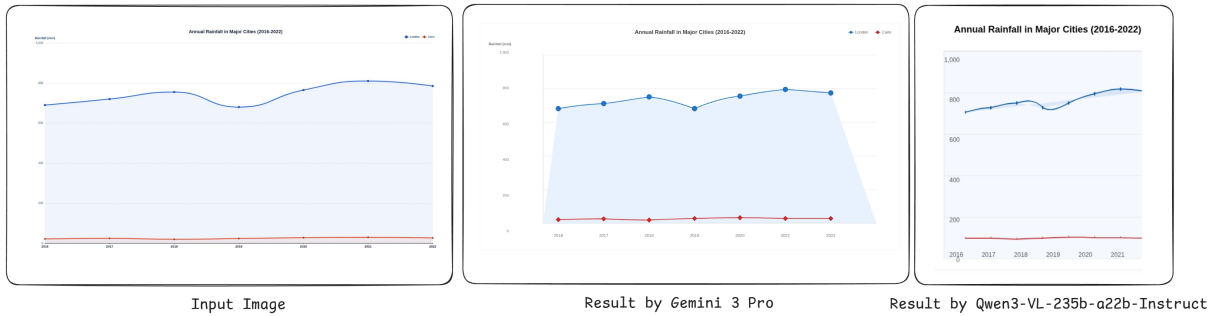
| Gemini 3 Pro | Qwen3-VL-235b-a22b-Instruct |
|---|--|
| <p>S_img: 76</p> <p>composition: 6</p> <p>textual: 10</p> <p>entity: 8</p> <p>stylistic: 7</p> | <p>S_img: 66</p> <p>composition: 5</p> <p>textual: 10</p> <p>entity: 6</p> <p>stylistic: 7</p> |
| <p>S_func: 50</p> <p>GT uses 7 functional units: (1) Generate x with linspace(0, 3, 300), (2) Calculate y = cos(1.5*pi*x), (3) Create y_upper using masked_where for values < threshold_upper, (4) Create y_lower using masked_where for values > threshold_lower, (5) Create y_mid using masked_where for values outside both thresholds, (6) Plot three masked arrays directly, (7) Apply grid with specific linestyle='-' . Gen implements: (1) linspace correctly but different sample count (minor), (2) INCOMPLETE: uses cos(2*pi*x) instead of cos(1.5*pi*x) - this fundamentally changes the wave period and is a critical logic error (0 points), (3-5) Uses boolean masking instead of np.ma.masked_where - different approach but achieves similar segmentation (0.5 points each = 1.5 total), (6) Plots with indexing x[mask], y[mask] instead of masked arrays - functionally similar but different method (0.5), (7) Grid uses '-' instead of '-' (0.5).</p> | <p>S_func: 21</p> <p>GT implements threshold-based masking logic to segment a single cosine wave into three conditional ranges (upper, lower, mid) using np.ma.masked_where. Gen instead generates three independent mathematical functions (mid_range, lower_band, upper_band) without any masking or conditional logic. The core functional approach is fundamentally different: GT uses 1 base function + 3 masking operations = 3 plot calls for masked segments, while Gen uses 3 separate mathematical functions + 3 plot calls. GT has 7 key functional units: (1) base cosine generation, (2) upper threshold masking, (3) lower threshold masking, (4) mid-range masking with compound condition, (5) plot masked, (6) plot lower masked, (7) plot upper masked. Gen only matches the base cosine generation partially (wrong frequency: 2π vs 1.5π) and the three plot calls (but plotting wrong data), giving partial credit for plotting structure but no credit for masking logic.</p> |
| <p>S_para: 30</p> <p>GT key parameters: (1) x range start=0, (2) x range end=3.0, (3) x samples=300, (4) cosine coefficient=1.5, (5) threshold_upper=0.65, (6) threshold_lower=-0.65, (7) y_mid linewidth=2.5, (8) y_lower linewidth=1.8, (9) y_upper linewidth=1.8, (10) grid alpha=0.7. Gen matches: (1) 0 correct, (2) 3 correct (treating 3.0 and 3 as equivalent), (3) 300 vs 600 - mismatch, (4) 1.5 vs 2 - CRITICAL: mismatch affecting wave period, (5) 0.65 vs 0.7 - mismatch, (6) -0.65 vs -0.7 - mismatch, (7) 2.5 vs 3 - mismatch, (8) 1.8 vs 3 - mismatch, (9) 1.8 vs 3 - mismatch, (10) 0.7 correct. Matches: parameters 1, 2, 10 = 3 out of 10.</p> | <p>S_para: 27</p> <p>GT key parameters: (1) x range start=0, (2) x range end=3.0, (3) x samples=300, (4) cosine frequency coefficient=1.5, (5) cosine frequency multiplier=-, (6) threshold_upper=0.65, (7) threshold_lower=-0.65, (8) mid linewidth=2.5, (9) lower/upper linewidth=1.8, (10) grid linestyle='-', (11) grid alpha=0.7. Gen mismatches: frequency is 2π instead of 1.5π (parameters 4 and 5 both wrong as they combine to different value), samples=500 vs 300, threshold values completely offset (replaced by 0.5 offsets in separate functions), mid linewidth=2 vs 2.5, lower/upper linewidth=2 vs 1.8, grid linestyle='-' vs '-'. Matches: x start=0, x end=3, grid alpha=0.7. Total GT parameters: 11. Matched: 3 (x start, x end, grid alpha).</p> |

Figure 12: Case of Python



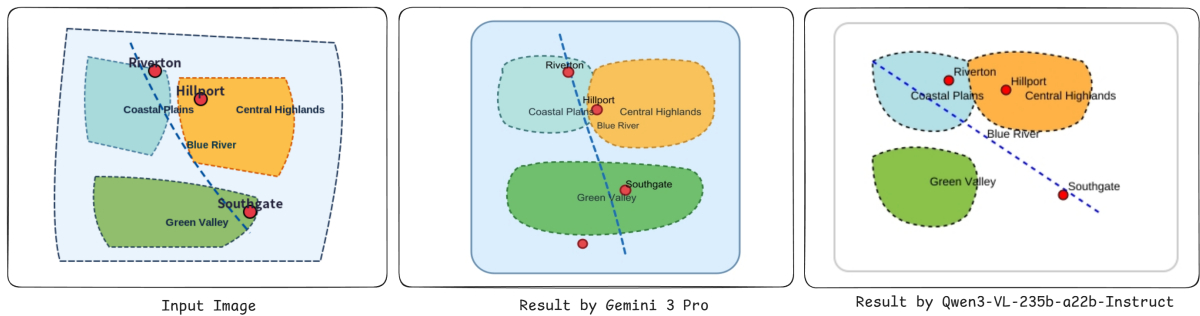
| Gemini 3 Pro | Qwen3-VL-235b-a22b-Instruct |
|--|--|
| <p>S_img: 81</p> <p>composition: 6</p> <p>textual: 10</p> <p>entity: 8</p> <p>stylistic: 8</p> | <p>S_img: 71</p> <p>composition: 4</p> <p>textual: 10</p> <p>entity: 6</p> <p>stylistic: 8</p> |
| <p>S_func: 60</p> <p>The generated code does not use the PGFPlots environment, which is the core structural component of the ground truth. Instead, it uses basic tikz commands and manual coordinate systems. The mathematical quadratic function is replaced by a manual Bezier curve, which lacks precision. Furthermore, the generated code fails to implement the automatic 'fill between' logic with intersection splitting; instead, it manually clips a region for the first area and uses a simple rectangle for the second area, which is geometrically incorrect compared to the function-bounded area in the ground truth. Tick marks and labels are also manually drawn rather than using axis properties.</p> | <p>S_func: 42</p> <p>The GT code uses pgfplots with a custom function declaration (quad), axis environment with tick labels, and the fillbetween library to automatically compute and fill regions between two curves with alternating colors. The Gen code abandons pgfplots entirely, uses raw tikz drawing commands, manually approximates the quadratic function with a cubic spline, and manually specifies fill regions with hardcoded intersection points. GT has 6 functional units: (1) custom function declaration via pgfmathdeclarefunction, (2) axis environment with configured ticks and labels, (3) addplot for the quadratic curve, (4) addplot for the constant line, (5) fillbetween logic with split and segment styling, (6) automatic region detection and coloring. Gen has 5 functional units: (1) manual axis drawing, (2) manual tick label placement, (3) manual curve plotting with approximation, (4) manual fill regions with hardcoded coordinates. The fillbetween library's automatic split logic and segment alternation are completely absent in Gen, replaced by manual region specification.</p> |
| <p>S_para: 25</p> <p>The numerical values in the generated code are entirely different from the ground truth. Specifically, the x-coordinates for p, eta, and q (0.5, 2.5, 4.0 in GT) are defined as 1.3, 3.6, and 6.0 in the generated code. The y-level for the horizontal line (qeta) = 5 in GT is set to 3.0. The axis limits and the specific mathematical coefficients of the quadratic equation are not respected. The color scheme (red, blue, red, cyan) and the label strings themselves are the only parameters that conceptually align with the ground truth.</p> | <p>S_para: 13</p> <p>Gen parameters: (1) scale: 1.5, (2) x-axis endpoint: 4.5, (3) y-axis endpoint: 3.5, (4) x-tick positions: 1, 2.5, 4 (incorrect: should be 0.5, 2.5, 4), (5) tick labels: p, eta, q (correct), (6) horizontal line y-value: 2 (incorrect: should be 5), (7) horizontal line x-range: 1 to 4 (incorrect: should span domain 0.5 to 4), (8) curve approximation formula: 0.5*(x-2.5)*(x-3.5) + 2 (completely different from 2*x^2 - 4*x + 1), (9) samples: 100 (incorrect: should be 160), (10) red fill region endpoint: 2.8 (hardcoded intersection, not computed), (11) red fill color: red!70 (close to red!80!black but not exact), (12) cyan fill color: cyan!60 (incorrect: should be cyan!70!black), (13) no opacity specification (should be 0.7). Gen matches only 2 parameters exactly (tick labels p, eta, q) and partially matches 1 (samples count). All numerical values for function, domain, axis limits, colors, and line properties are incorrect or missing.</p> |

Figure 13: Case of LaTeX



| Gemini 3 Pro | Qwen3-VL-235b-a22b-Instruct |
|--|--|
| <p>S_img: 78</p> <p>S_func: 88</p> <p>S_para: 58</p> <p>GT uses ECharts library to render an interactive line chart with tooltip, legend, grid, axes, and two series (London and Cairo). Gen implements a custom SVG-based static chart with manual DOM structure (wrapper, header, legend, y-labels, SVG elements). Both achieve the same visual goal: a dual-line chart with area fill, markers, axes, grid, legend, and title. ... Gen lacks this interactivity. The absence of interactive tooltip in Gen is a missing functional unit. All other structural elements are implemented, though with different tags (SVG primitives vs. ECharts config). Total GT functional units: 8 (title, legend, x-axis, y-axis, grid, London series with area/markers, Cairo series with markers, tooltip). Gen implements 7 (missing interactive tooltip).</p> | <p>S_img: 68</p> <p>S_func: 56</p> <p>S_para: 29</p> <p>GT uses ECharts library to render an interactive line chart with two series (London and Cairo), while Gen manually draws a static SVG chart. Gen lacks the ECharts script inclusion, dynamic chart initialization, and interactive tooltip. Gen also omits the y-axis title 'Rainfall (mm)' and uses a different layout structure (centered container vs full-height). However, Gen preserves the overall chart concept, line/area rendering, legend, and axis labels, achieving partial structural intent.</p> |

Figure 14: Case of HTML



| Gemini 3 Pro | Qwen3-VL-235b-a22b-Instruct |
|--|---|
| <p>S_img: 85</p> <p>S_func: 94</p> <p>S_para: 69</p> <p>GT units identified: 1. Country-outline path (irregular closed polygon). 2. Region-1 path (Coastal Plains, closed). 3. Region-2 path (Central Highlands, closed). 4. Region-3 path (Green Valley, closed). 5. River path (open curved). 6. Riverton city circle. 7. Hillport city circle. 8. Southgate city circle. In Gen: Rounded rect background partially matches country-outline intent/topology (simplified enclosure, 0.5). All 3 region paths present with equivalent closed bezier topology (1.8 each). River path present with equivalent open curve (1.8). All 3 city circles present with correct primitive type (1.8 each). Text elements all present matching intent/content.</p> | <p>S_img: 80</p> <p>S_func: 94</p> <p>S_para: 47</p> <p>GT decomposition: The GT code contains the following functional units: (1) country-outline path, fill:#fff, stroke:#333, stroke-width=2, stroke-dasharray=8,4; (2) region-1 fill:#d8d0dc, stroke:#407094, stroke-dasharray=8,3; (3) region-2 fill:#e6b4f0, stroke:#403040, stroke-dasharray=6,3; (4) region-3 fill:#90b0e0, stroke:#204060, stroke-dasharray=6,3; (5) river: stroke:#008080, stroke-width=3; (6) three city marker circles; (7) three city label text elements; (8) three region label text elements; (9) one river label text element. Total: 9 distinct logical units. Gen Comparison: ... The Gen code implements: river path (0.5 credit for topological equivalence). Total implemented with credit: 13.5 out of 14 expected structural elements (counting country-outline as mandatory).</p> |

Figure 15: Case of SVG

Case Analysis: 6cf0140514bd4066ae6711234dc518c6

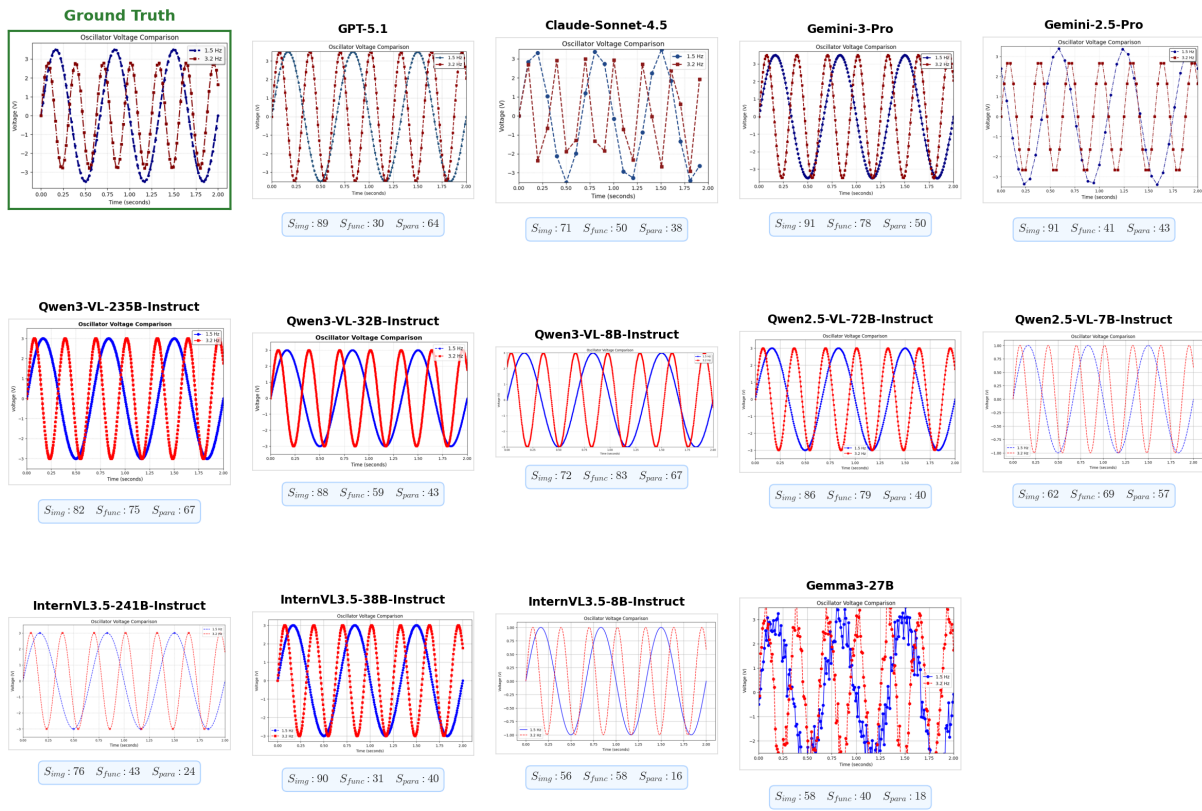


Figure 16: Case 1 of render_samples for all models benchmarked in our Tab. 2.

Case Analysis: AsymptoteGraphicPipeline_Circle_Geometry_2-88

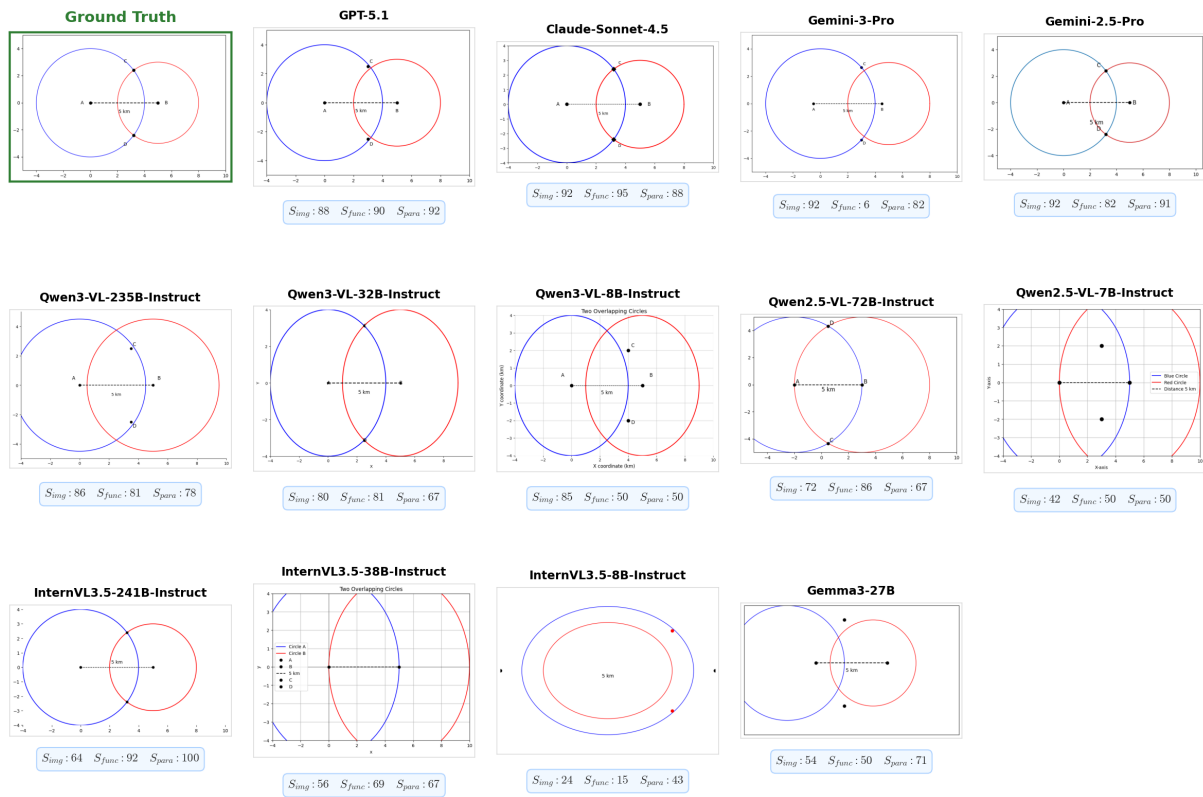


Figure 17: Case 2 of render_samples for all models benchmarked in our Tab. 2.

Omni-I2C Annotation Re-evaluation (50 Samples)

Annotation Guidelines

- Core Metrics: Focus on Functionality, Numerical Accuracy, and Style.
- Ties Allowed: You can give the same rank to multiple models (e.g., 1, 1, 2).
- Redundancy: If a generated image includes more features than GT, do not penalize.
- Workflow: Annotate → Next → After 50 Items, Save with your name.

Progress: 4 / 50 | Current ID: line_25

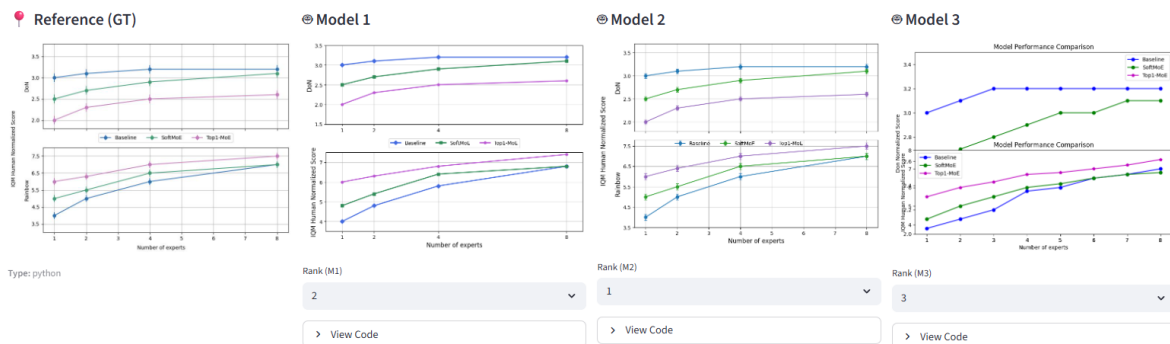


Figure 18: A screenshot of the human evaluation questionnaire.

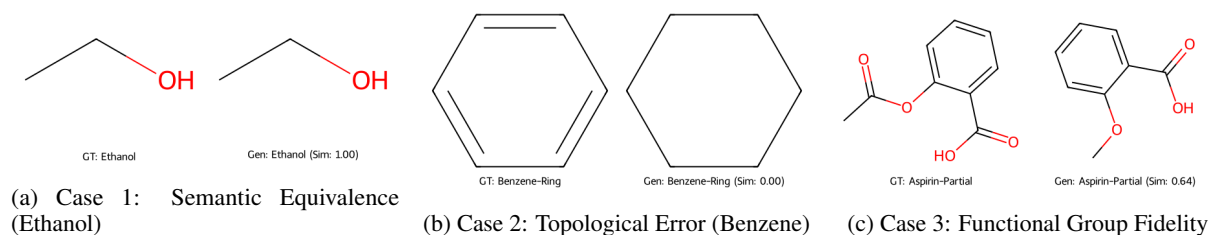


Figure 19: Visual comparison of SMILES reconstruction cases. (a) Illustrates identical topology despite different string sequences. (b) Highlights critical perception errors in aromaticity captured by Tanimoto similarity. (c) Demonstrates the metric's sensitivity to localized functional group deletions.

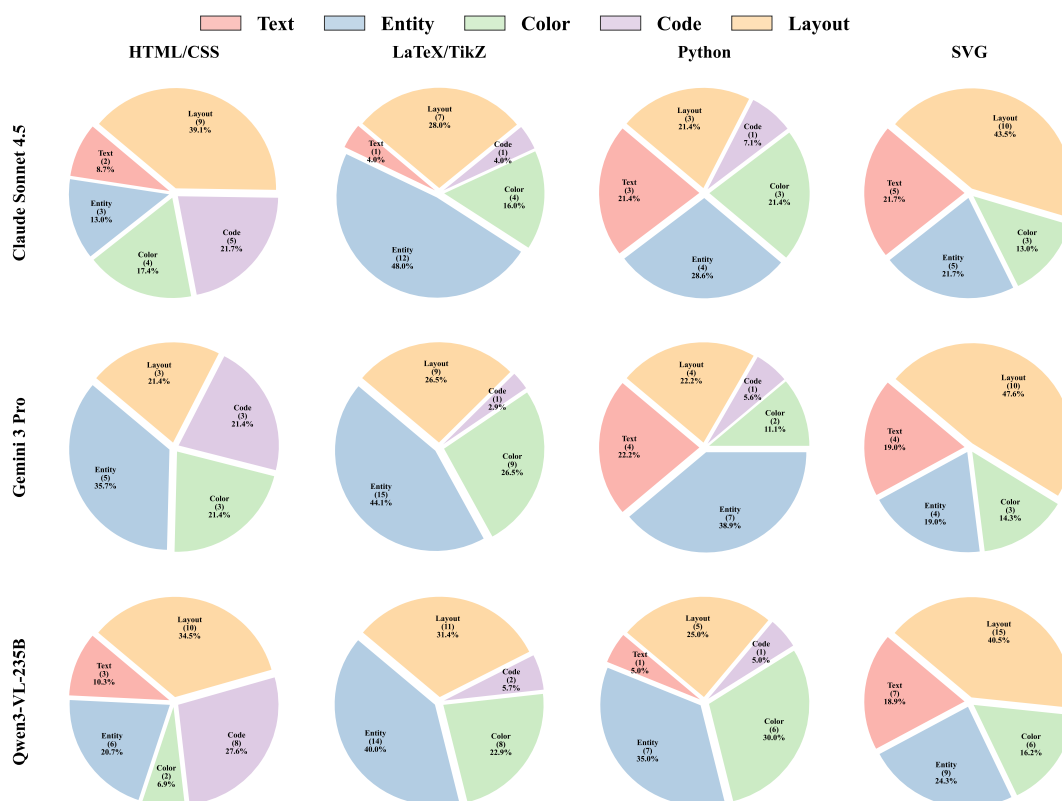
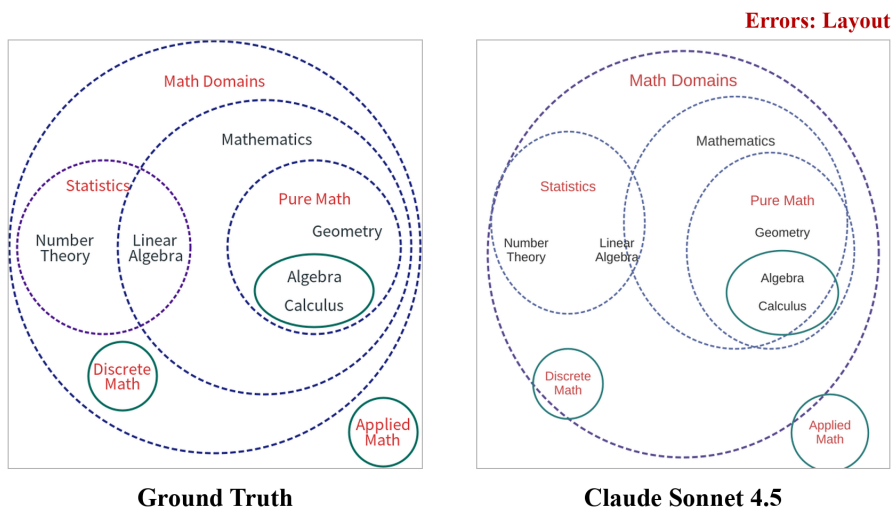
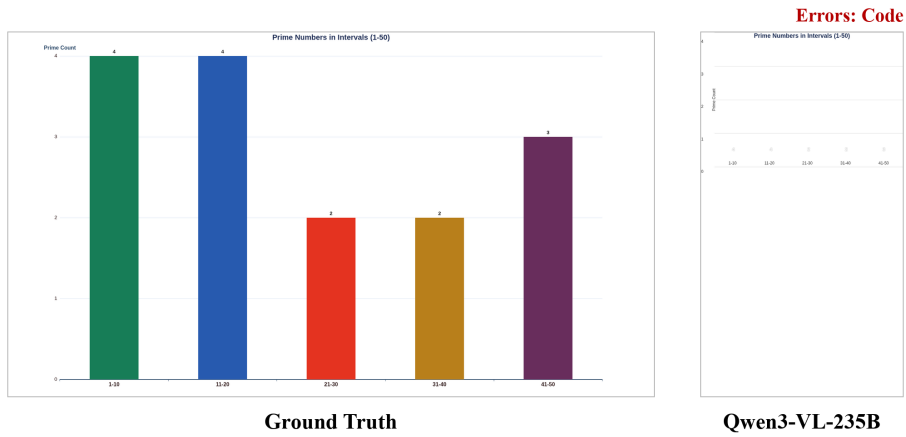


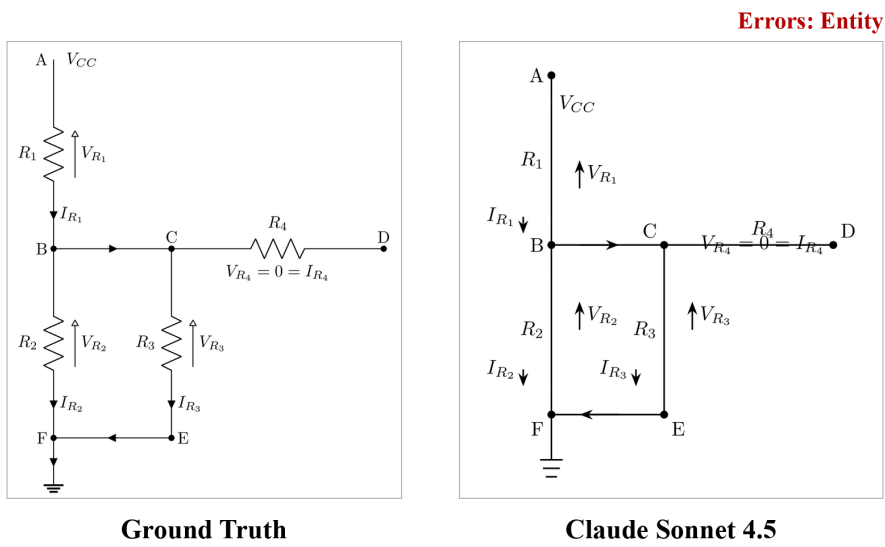
Figure 20: Fine-grained error distribution across languages



(a) Layout Error



(b) Code Error



(c) Entity Error

Figure 21: Representative error cases produced by different models.

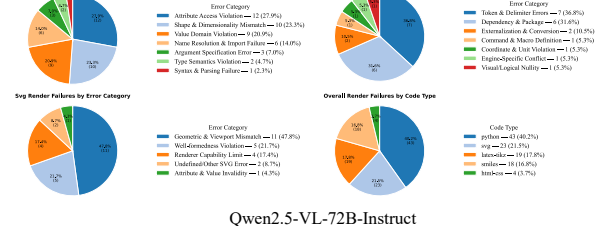
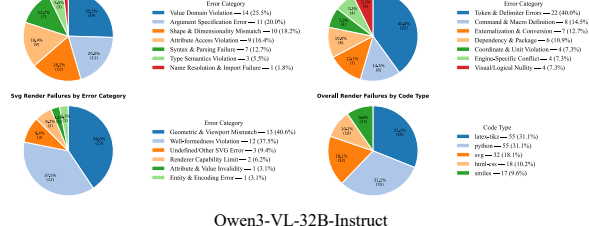
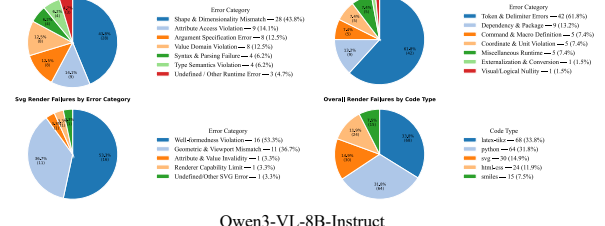
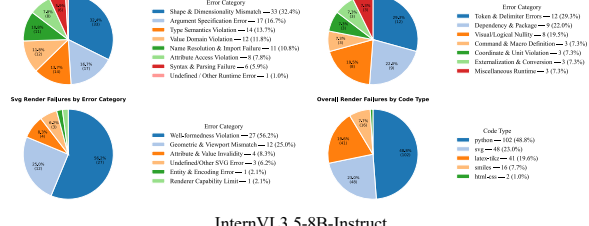
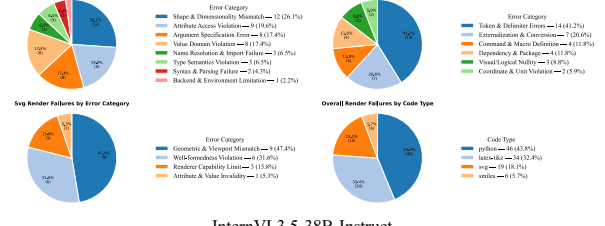
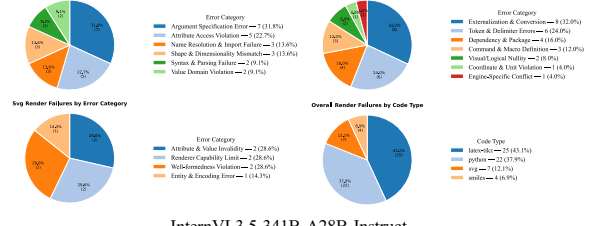
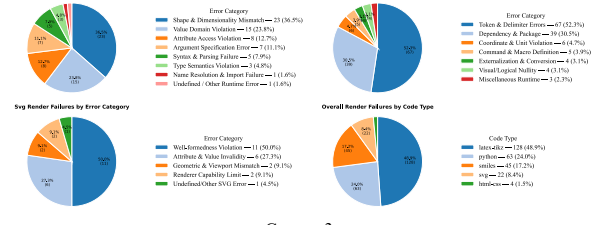
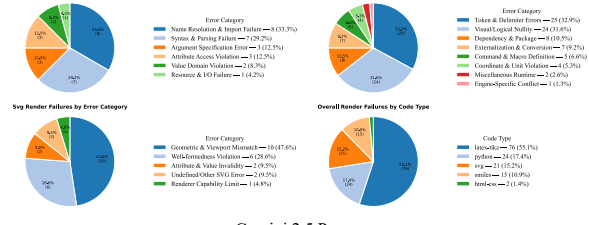
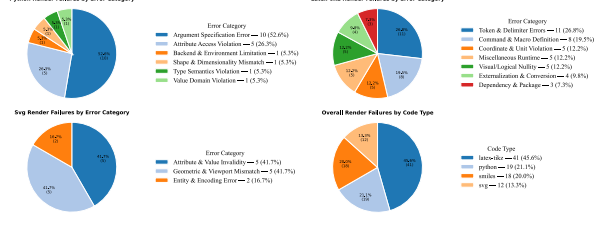
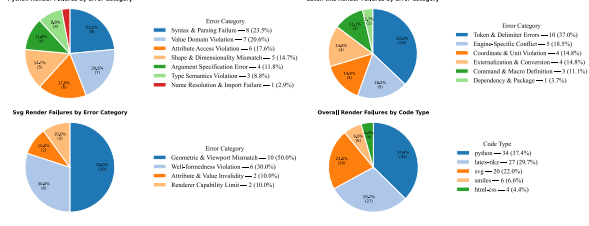
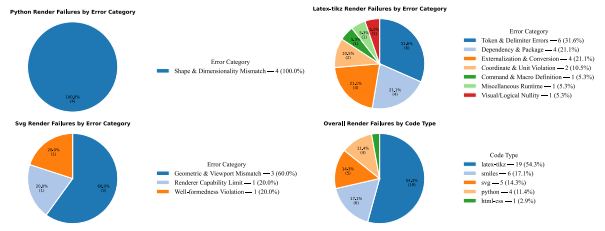
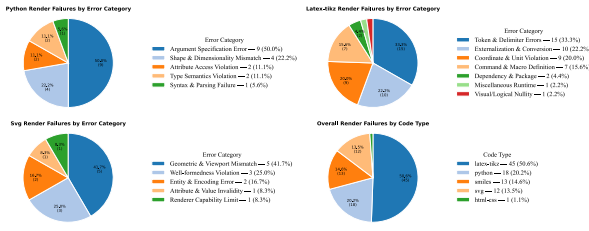


Figure 22: Diagnostic breakdown of failure modes across evaluated models.

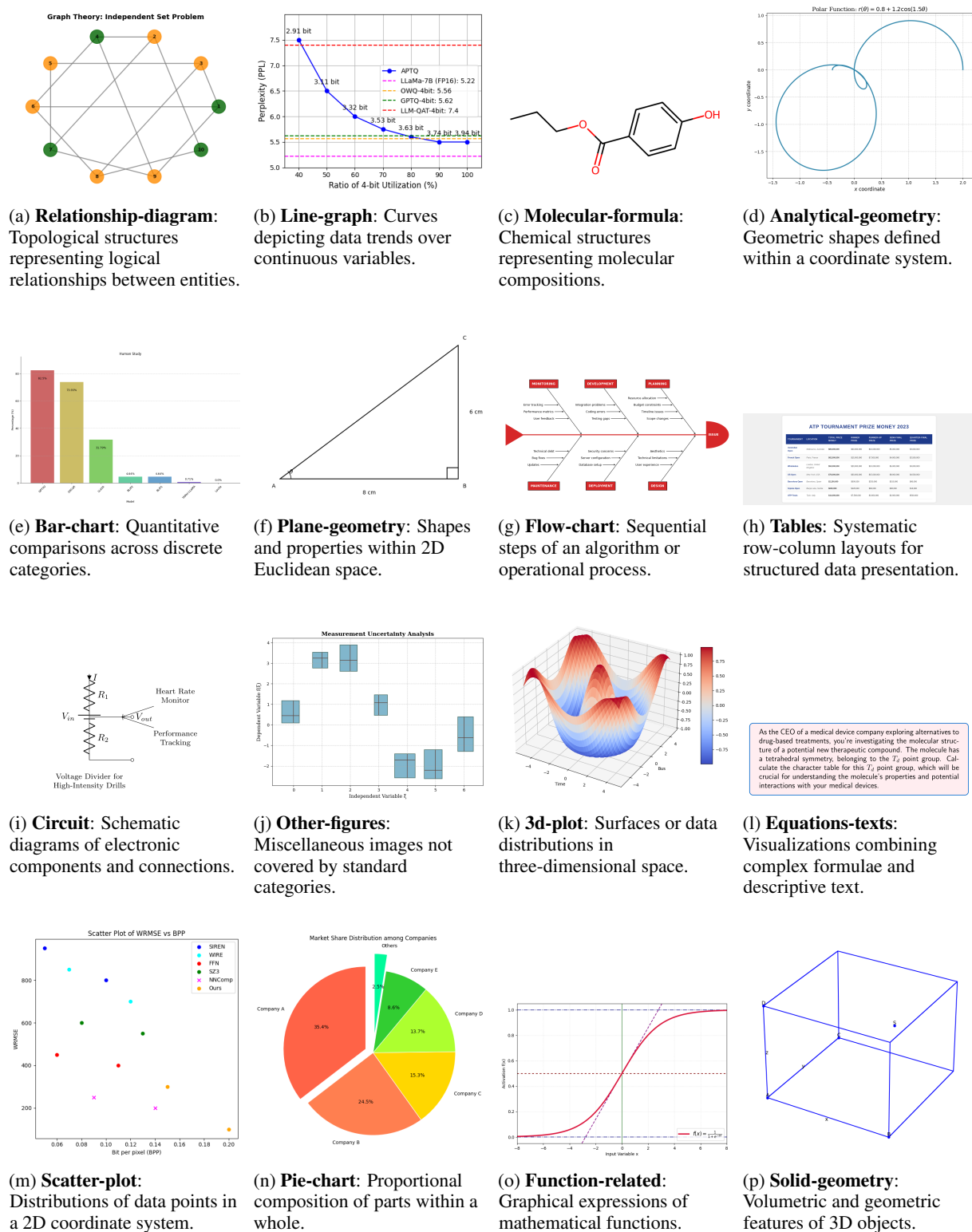
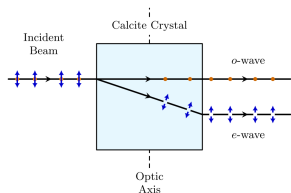
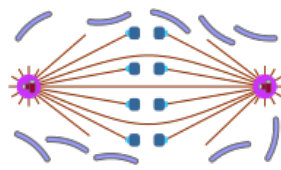


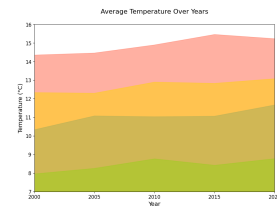
Figure 23: Omni-I2C dataset: Figure Type Taxonomy gallery (Part 1).



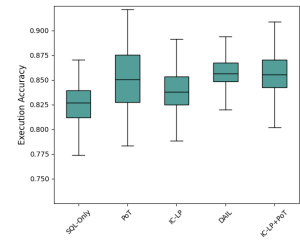
(a) **Schematic:** Simplified diagrams illustrating systems or structures.



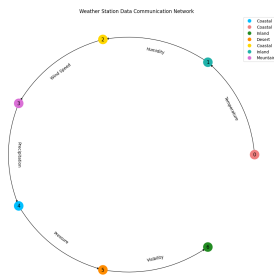
(b) **Physiological-process:** Dynamic representations of biological activities.



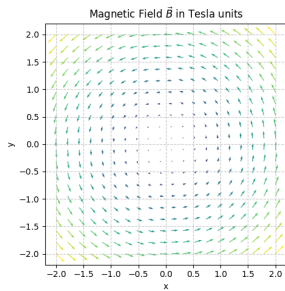
(c) **Area:** Trend charts with filled area sections.



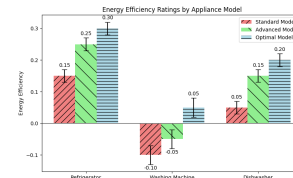
(d) **Box-chart:** Statistical summaries showing quartiles and outliers.



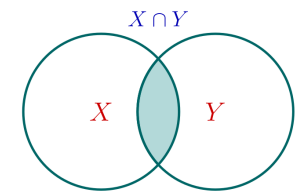
(e) **Graph:** General node-edge connectivity representations.



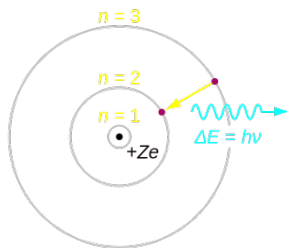
(f) **Quiver:** Distribution of vector fields in a spatial domain.



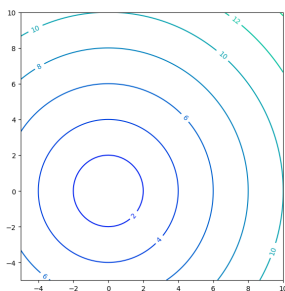
(g) **Error-bar:** Statistical charts indicating data uncertainty ranges.



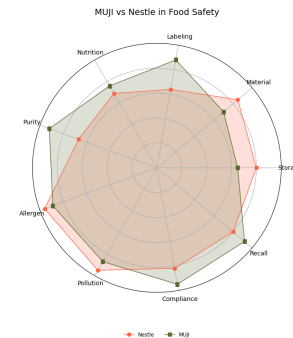
(h) **Venn-diagram:** Overlapping sets showing logical relations.



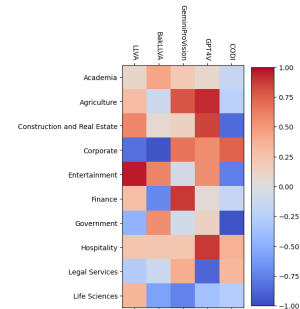
(i) **Atom-model:** Physical models of atomic structures and electrons.



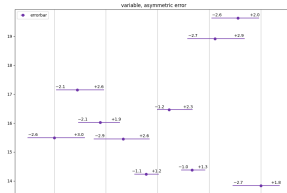
(j) **Contour:** Numerical distributions with isoline features.



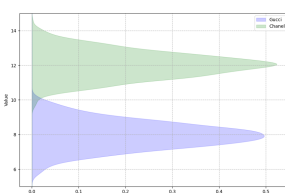
(k) **Radar-chart:** Comparative visualizations of multi-dimensional data.



(l) **Heatmap:** Matrix intensities represented by color gradients.



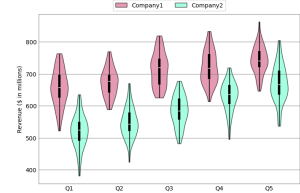
(m) **Error-point:** Data points with designated error offsets.



(n) **Density:** Spatial distributions of probability densities.

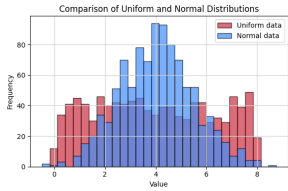


(o) **Map:** Visualizations of geographic locations or regions.

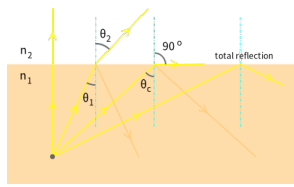


(p) **Violin:** Distribution plots combining density and box plot features.

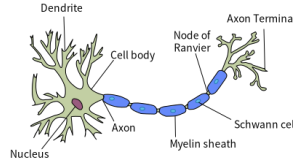
Figure 24: Omni-I2C dataset: Figure Type Taxonomy gallery (Part 2).



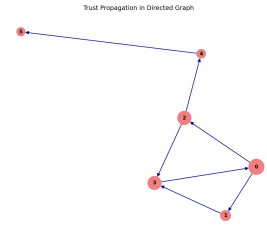
(a) **Histogram:** Statistical bars for frequency or count distributions.



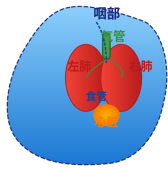
(b) **Optics-ray-diagram:** Propagation paths of light in optical systems.



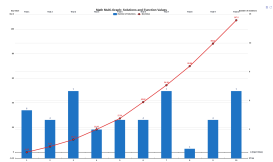
(c) **Cell-structure:** Diagrams illustrating internal biological cell parts.



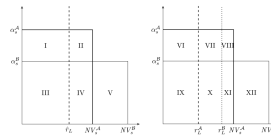
(d) **Algorithms:** Visualization of algorithmic logic or pseudocode.



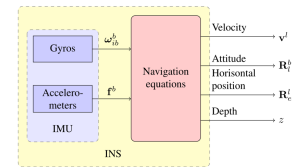
(e) **Anatomy-diagram:** Detailed structural diagrams of biological anatomy.



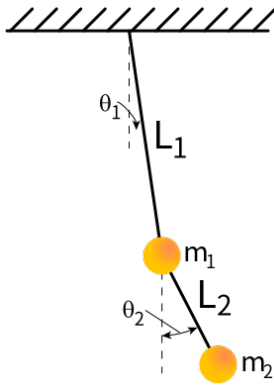
(f) **Multi-graph:** Composite images containing multiple sub-figures.



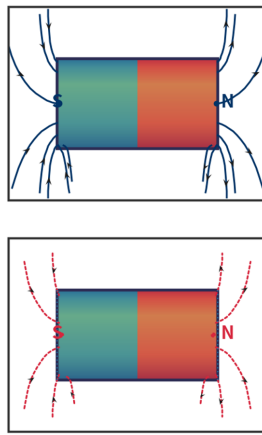
(g) **Phase-Diagram:** State charts describing system phase transitions.



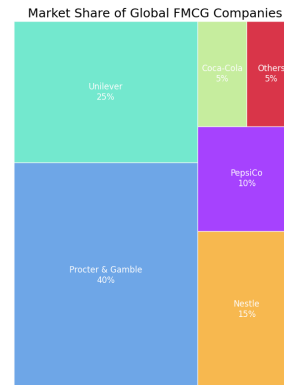
(h) **Block-diagram:** Functional relationships between system modules.



(i) **Free-body-diagram:** Classical mechanics models for force analysis.



(j) **Magnetic-field-line:** Orientations and distributions of magnetic fields.



(k) **Tree:** Hierarchical structures with parent-child relationships.

Figure 25: Omni-I2C dataset: Figure Type Taxonomy gallery (Part 3).

You are an expert in Python data visualization and proficient in multiple libraries (including Matplotlib, Seaborn, Plotly, Bokeh, Altair, etc.).

(a) Direct (Zero-shot) Prompting

You are an expert in Python data visualization and proficient in multiple libraries (including Matplotlib, Seaborn, Plotly, Bokeh, Altair, etc.). Your task is to analyze the provided image and write complete Python code to reproduce it visually using the specified plotting library.

```

Requirements:
1. Target Library: Use the library requested by the user. If no specific library is mentioned, default to Matplotlib.
2. Data: Create reasonable dummy data based on visual trends observed in the image.
3. Dimensions: You MUST set the figure dimensions to width: {width}, height: {height}. Use the syntax strictly appropriate for the chosen library.
4. Visuals: Match colors, styles, limits, and annotations exactly. Adapt the implementation (e.g., interactive tooltips vs. static text) to the strengths of the chosen library.

*** EXAMPLES OF EXPECTED OUTPUT STRUCTURE ***

Example 1: Complex Network Graph
```python
Part 1: Importing Libraries
Part 2: Data Preparation
G = nx.DiGraph()
modules = {0: "Sales", 1: "Marketing", 2: "Product Development",
 3: "R&D", 4: "Finance", 5: "Human Resources",
 6: "IT Support", 7: "Customer Service", 8: "Operations",
 9: "Legal", 10: "Compliance", 11: "Supply Chain"}
G.add_nodes_from(modules.keys())
edges = [(0, 1, 3), (1, 2, 4), (1, 9, 2), (2, 3, 4), (3, 4, 3),
 (3, 5, 2), (4, 7, 2), (5, 6, 1), (6, 10, 4), (7, 8, 3),
 (8, 10, 4), (10, 11, 2)]
G.add_weighted_edges_from(edges)
pos = nx.spring_layout(G, seed=42)
node_labels = modules
edge_labels = {(start, end): f'{weight}' for start, end, weight in edges}

Part 3: Plot Configuration and Rendering
plt.figure(figsize=(width, height))
node_size = [800 + 200 * (i % 3) for i in range(len(modules))]
node_color = ["lightblue" if i % 3 == 0 else "green" if i % 3 == 1 else "red" for i in range(len(modules))]
nx.draw(G, pos, node_size=node_size, node_color=node_color, with_labels=False, arrows=True)
nx.draw_networkx_labels(G, pos, labels=node_labels, font_size=9)
nx.draw_networkx_edge_labels(G, pos, edge_labels=edge_labels, font_color='red')
plt.title("Marketing units and their collaboration relationships", fontsize=15)
plt.axis('off')

Part 4: Saving Output
plt.tight_layout()
plt.savefig("output_graph.pdf", bbox_inches='tight')
...
*** IMPORTANT ***
* Output ONLY the Python code inside a single Markdown block ('python ...').
* NO conversational text.

```

### (b) Few-Shot Prompting

You are an expert in Python data visualization and proficient in multiple libraries (including Matplotlib, Seaborn, Plotly, Bokeh, Altair, etc.). Your task is to analyze the provided image and write complete Python code to reproduce it visually using the specified plotting library.

```

Requirements:
1. Target Library: Use the library requested by the user. If no specific library is mentioned, default to Matplotlib.
2. Data: Create reasonable dummy data based on visual trends observed in the image.
3. Dimensions: You MUST set the figure dimensions to width: {width}, height: {height}. Use the syntax strictly appropriate for the chosen library.
4. Visuals: Match colors, styles, limits, and annotations exactly. Adapt the implementation (e.g., interactive tooltips vs. static text) to the strengths of the chosen library.
5. Comments: Append a high-quality comment to EVERY code line explaining the visual intent.

*** IMPORTANT ***
- Output ONLY the Python code inside a single Markdown block ('python ...').
- NO conversational text.

```

### (c) Self-Commentary Prompting

You are an expert in Python data visualization and proficient in multiple libraries (including Matplotlib, Seaborn, Plotly, Bokeh, Altair, etc.). Your task is to analyze the provided image and write complete Python code to reproduce it visually using the specified plotting library.

```

Anchor-Based Image Analysis
You will be provided with two images. The first image is the original picture. The second image is the same picture overlaid with a 6x6 dot matrix to serve as reference anchors.
- Each dot is labeled with two-dimensional coordinates (x, y).
- Within each column, the x-coordinate increases from top to bottom.
- Within each row, the y-coordinate increases from left to right.
- Your Task: First, use these dots as reference anchors to generate a detailed description of the picture (e.g., "Between dot (1,1) and (1,2), there is a legend..."). Then, based on this analysis, generate the code to reproduce the image.

```

```

Requirements:
1. Target Library: Use the library requested by the user. If no specific library is mentioned, default to Matplotlib.
2. Data: Create reasonable dummy data based on visual trends observed in the image.
3. Dimensions: You MUST set the figure dimensions to width: {width}, height: {height}. Use the syntax strictly appropriate for the chosen library.
4. Visuals: Match colors, styles, limits, and annotations exactly. Adapt the implementation (e.g., interactive tooltips vs. static text) to the strengths of the chosen library.

```

```

*** IMPORTANT ***
- First, provide the anchor-based description of the image.
- Then, output the Python code inside a single Markdown block ('python ...').
- NO conversational text.

```

### (d) Scaffold Prompting

You are an expert in Python data visualization and proficient in multiple libraries (including Matplotlib, Seaborn, Plotly, Bokeh, Altair, etc.). Your task is to analyze the provided image and write complete Python code to reproduce it visually using the specified plotting library.

```

Analysis Phase
To ensure accuracy and detail in your recreation, begin with a comprehensive analysis of the figure to develop an elaborate description. This description should cover, but not be limited to, the following aspects:
1. Layout and Structural Analysis: Identify the overall composition and arrangement of elements. Note the presence of subplots, grids, columns, or nested structures and how they are aligned.
2. Element Identification: Determine the nature of the components present—such as chart types, text blocks, mathematical notations, geometric shapes, or UI components. Identify if elements are independent or share common structural frameworks (like axes or containers).
3. Content and Data Interpretation: Summarize the information or data patterns being conveyed. For visualizations, identify trends; for documents or formulas, identify the semantic content and hierarchy to generate accurate dummy data or text.
4. Visual and Stylistic Features: Identify supplementary elements that contribute to clarity and aesthetics, including color palettes, font styles, line weights, markers, legends, transparency, and specific rendering artifacts.

```

```

Requirements:
1. Target Library: Use the library requested by the user. If no specific library is mentioned, default to Matplotlib.
2. Data: Create reasonable dummy data based on visual trends observed in the image.
3. Dimensions: You MUST set the figure dimensions to width: {width}, height: {height}. Use the syntax strictly appropriate for the chosen library.
4. Visuals: Match colors, styles, limits, and annotations exactly. Adapt the implementation (e.g., interactive tooltips vs. static text) to the strengths of the chosen library.

```

```

*** IMPORTANT ***
- First, provide the hint-enhanced description of the image.
- Then, output the Python code inside a single Markdown block ('python ...').
- NO conversational text.

```

### (e) Hint-Enhanced Prompting

Figure 26: Detailed visualization of the evaluated prompting strategies (Part I).

Figure 26: Detailed visualization of the evaluated prompting strategies (Part II).

## J Datasheet for Omni-I2C

### J.1 Motivation

**1. For what purpose was the dataset created? Was there a specific task in mind? Was there a specific gap that needed to be filled? Please provide a description.**

**A1:** Omni-I2C is designed to evaluate the capability of Large Multimodal Models (LMMs) in converting complex, structured digital graphics into executable code. Existing benchmarks are often restricted in image variety and programming languages, failing to capture the structural complexity found in authentic, user-sourced applications. This “Grand Challenge” requires a synergy between high-fidelity visual perception and precise generative expression to achieve “pixel-to-program” alignment.

**2. Who created this dataset (e.g., which team, research group) and on behalf of which entity (e.g., company, institution, organization)?**

**A2:** This dataset is created by the authors of this paper.

**3. Who funded the creation of the dataset? If there is an associated grant, please provide the name of the grantor and the grant name and number.**

**A3:** N/A.

### J.2 Composition

**1. What do the instances that comprise the dataset represent? Please provide a description.**

**A1:** Omni-I2C consists of 1,130 high-quality Image2Code instances. Each instance comprises an input digital graphic (image) and its corresponding ground-truth implementation in executable code, spanning technical domains such as scientific visualizations, molecular structures, and complex symbolic notations.

**2. How many instances are there in total (of each type, if appropriate)?**

**A2:** There are 1,130 evaluation instances in total, distributed across five programming and markup languages: Python (364), LaTeX (274), HTML (182), SVG (210), and SMILES (100).

**3. Does the dataset contain all possible instances or is it a sample? If the dataset is a sample, then what is the larger set?**

**A3:** It is a curated sample of real-world items. To ensure practical utility, we harvested image-code pairs from official documentation, specialized

tutorials, and community galleries, representing a wide spectrum of digital content.

**4. What data does each instance consist of? “Raw” data or features?**

**A4:** Each instance consists of the “raw” input image, the corresponding textual instruction, and the ground-truth executable code.

**5. Is there a label or target associated with each instance? If so, please provide a description.**

**A5:** Yes. Each target is the executable code that accurately reconstructs the visual input. The fidelity of the reconstructed image is evaluated across semantic and structural dimensions.

**6. Is any information missing from individual instances?**

**A6:** No.

**7. Are relationships between individual instances made explicit?**

**A7:** Yes. Instances are explicitly categorized by their programming language, academic discipline (8 subjects), and figure archetype (43 types).

**8. Are there recommended data splits (e.g., training, development/validation, testing)?**

**A8:** Omni-I2C is primarily intended as a benchmark for evaluation.

**9. Are there any errors, sources of noise, or redundancies in the dataset?**

**A9:** No. All instances were verified through an isolated sandbox environment to ensure the code compiles and produces the correct visual output.

**10. Is the dataset self-contained, or does it link to or otherwise rely on external resources?**

**A10:** The dataset is self-contained. While raw data was collected from public sources, the final curated benchmark includes all necessary code and images.

**11. Does the dataset contain data that might be considered confidential?**

**A11:** No. A manual audit was performed to prune any Personally Identifiable Information (PII).

**12. Does the dataset contain data that might be offensive?**

**A12:** No.

### J.3 Collection Process

**1. How was the data associated with each instance acquired?**

**A1:** Data was acquired through a systematic three-stage pipeline: collection, filtering, and decontamination. We targeted official documentation and specialized galleries (e.g., Matplotlib Gallery,

TikZ.net, TeXample.net) to capture genuine practical utility.

**2. What mechanisms or procedures were used to collect the data?**

**A2:** We utilized a taxonomy-driven approach, pre-defining a comprehensive schema to guide the search and ensure broad coverage across subjects and languages.

**3. If the dataset is a sample from a larger set, what was the sampling strategy?**

**A3:** Purposeful sampling was used to address gaps in existing benchmarks, with 5 languages and 43 distinct figure types.

**4. Who was involved in the data collection process?**

**A4:** Beyond the primary authors, nine internal annotators from the same research institution were engaged for data collection. These individuals were compensated at a rate of eight times the local minimum hourly wage. To prevent interference with their regular schedules, the tasks were conducted on a flexible, part-time basis without fixed daily hour requirements. The average time commitment per annotator was approximately two days.

**5. Over what timeframe was the data collected?**

**A5:** The collection, manual auditing, and decontamination process took approximately one month.

#### J.4 Preprocessing/cleaning/labeling

**1. Was any preprocessing/cleaning/labeling of the data done?**

**A1:** Yes. Candidates were executed in isolated sandboxes to discard instances with compilation errors. Manual auditing was used to prune social biases and PII. Finally, we used an LLM to refine the categorization of the remaining data.

**2. Was a data decontamination strategy employed?**

**A2:** Yes. We employed an LLM to reformulate code styles and aesthetic parameters (fonts, layouts). For images, we performed “content erasure” and generated new synthetic strings to prevent models from relying on linguistic priors or memorization.

**3. Is the software used to preprocess/clean/label the instances available?**

**A3:** The evaluation pipeline and sandbox configurations will be provided upon publication.

#### J.5 Uses

**1. Has the dataset been used for any tasks already? If so, please provide a description.**

**A1:** Yes. We benchmarked 13 state-of-the-art LMMs, including GPT-5.1, Claude 4.5 Sonnet, and Gemini 3 Pro, evaluating their ability to reconstruct structural integrity in complex scenarios.

**2. Is there a repository that links to any or all papers or systems that use the dataset? If so, please provide a link or other access point.**

**A2:** N/A.

**3. What (other) tasks could the dataset be used for?**

**A3:** Omni-I2C can be used for evaluating multimodal reasoning, precise coordinate modeling in SVG/HTML, and specialized chemical informatics knowledge via SMILES.

**4. Is there anything about the composition of the dataset or the way it was collected that might impact future uses? Is there anything a future user could do to mitigate these undesirable harms?**

**A4:** No. The dataset was designed with counter-intuitive or “nonsense” content swaps to force models to rely on visual perception rather than linguistic priors, which actually enhances its robustness as a perceptual benchmark.

**5. Are there tasks for which the dataset should not be used? If so, please provide a description.**

**A5:** N/A.

#### J.6 Distribution

**1. Will the dataset be distributed to third parties outside of the entity? If so, please provide a description.**

**A1:** Yes, the dataset will be publicly available for the research community.

**2. How will the dataset will be distributed? Does the dataset have a digital object identifier (DOI)?**

**A2:** Omni-I2C will be distributed via the official project website and GitHub.

**3. When will the dataset be distributed?**

**A3:** The dataset will be distributed once the paper is accepted after peer review.

**4. Will the dataset be distributed under a copyright or other license?**

**A4:** It will be distributed under the Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International (CC BY-NC-SA 4.0) License.

**5. Have any third parties imposed IP-based or other restrictions on the data associated with the instances?**

**A5:** No.

**6. Do any export controls or other regulatory restrictions apply to the dataset or to individual instances?**

**A6:** No.

## **J.7 Maintenance**

**1. Who will be supporting/hosting/maintaining the dataset?**

**A1:** The authors.

**2. How can the owner/curator/manager of the dataset be contacted (e.g., email address)?**

**A2:** Email addresses will be provided on the project homepage post-publication.

**3. Is there an erratum? If so, please provide a link or other access point.**

**A3:** Any errata will be posted on the project GitHub repository.

**4. Will the dataset be updated? If so, please describe how often, by whom, and how updates will be communicated to users?**

**A4:** Yes. The authors will periodically update the benchmark with new instances and categories based on community feedback.

**5. Will older versions of the dataset continue to be supported/hosted/maintained? If so, please describe how.**

**A5:** Older versions will be archived on GitHub to ensure researchers can reproduce results from previous studies.

**6. If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for them to do so?**

**A6:** N/A.