

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 |

065 mantic meaning of the result. This execution-based
 066 verifiability enables a more comprehensive, fine-
 067 grained, and objective evaluation of the perceptual
 068 capabilities of a model.

069 Recognizing this potential, recent benchmarking
 070 efforts begin to explore I2C generation within spe-
 071 cialized domains. For instance, ChartMimic (Yang
 072 et al., 2024) and Plot2Code (Wu et al., 2025) fo-
 073 cus on synthesizing plotting scripts from scientific
 074 charts, while Design2Code (Si et al., 2025) eval-
 075 uates the transformation of web screenshots into
 076 front-end implementations. Although these works
 077 have established a vital foundation, current eval-
 078 uations remain largely restricted in image variety,
 079 programming languages, and task scenarios. To
 080 comprehensively assess the limits of perception,
 081 coding and reasoning, a more diverse and general-
 082 purpose testbed is essential to evaluate whether
 083 LMMs can generalize across varied technical and
 084 professional domains.

085 To bridge this gap, we introduce Omni-I2C, a
 086 comprehensive benchmark for stress-testing the
 087 limits of Image-to-Code (I2C) generation. The
 088 dataset contains 1.1k meticulously curated samples
 089 drawn from real-world user applications, cover-
 090 ing a broad spectrum of technical graphics such
 091 as scientific visualizations, molecular structures,
 092 and complex symbolic notations, and supporting
 093 diverse programming languages. Furthermore, we
 094 propose a multi-level evaluation framework that
 095 decomposes model performance into fine-grained
 096 attributes, enabling a more nuanced understanding
 097 of model strengths and weaknesses across different
 098 dimensions. Unlike previous metrics restricted to
 099 specific domains or coding formats, our evaluation
 100 logic offers comprehensive coverage and demon-
 101 strates superior alignment with human judgment
 102 through extensive empirical validation.

103 Through an extensive evaluation of 13 propri-
 104 etary and open-weight LMMs, we reveal a pro-
 105 found performance gap in high-fidelity image-to-
 106 code generation. Even leading frontier models,
 107 such as Gemini 3 Pro and GPT-5.1, frequently

108 falter in the challenging scenarios presented by
 109 our benchmark. These results highlight substantial
 110 room for improvement and position Omni-I2C as a
 111 challenging benchmark for advancing LMMs. Our
 112 contributions are summarized as follows:

- 113 • We present Omni-I2C, a meticulously curated
 114 dataset of 1.1k items, including 5 program-
 115 ming languages, 8 major subjects, and 43
 116 distinct figure types. It serves as a rigorous
 117 testbed for evaluating the perception and cod-
 118 ing capabilities of LMMs.
- 119 • We propose an evaluation framework that as-
 120 sesses code-level integrity and image-level
 121 perceptual consistency, enabling more diag-
 122 nostic and attributable analyses of model be-
 123 havior than traditional heuristic metrics.
- 124 • Our comprehensive analysis of SOTA LMMs
 125 exposes a significant performance gap in high-
 126 fidelity reconstruction, identifying critical fail-
 127 ure modes and charting a path toward more
 128 precise, trusted multimodal agents.

2 Related Work 129

2.1 Multimodal Code Benchmarks 130

131 The rapid advancement of Large Language Models
 132 (LLMs) has reshaped automated code generation,
 133 , from general-purpose foundation models (OPE-
 134 NAI, 2025a; Deepmind, 2025; Anthropic, 2025)
 135 to specialized architectures (Guo et al., 2024; Hui
 136 et al., 2024) that excel at complex programming
 137 tasks (Jimenez et al., 2023; Zhou et al., 2023; Yao
 138 et al., 2024). Model evaluation, however, has
 139 largely remained single-modal, relying on text-
 140 only benchmarks such as HumanEval (Chen, 2021),
 141 MBPP (Austin et al., 2021), and DS-1000 (Lai
 142 et al., 2023). As models become multimodal, eval-
 143 uation must expand to real-world settings where
 144 visual perception is integral.

145 Recent multimodal benchmarks—MMCode (Li
 146 et al., 2024), Design2Code (Si et al., 2025),

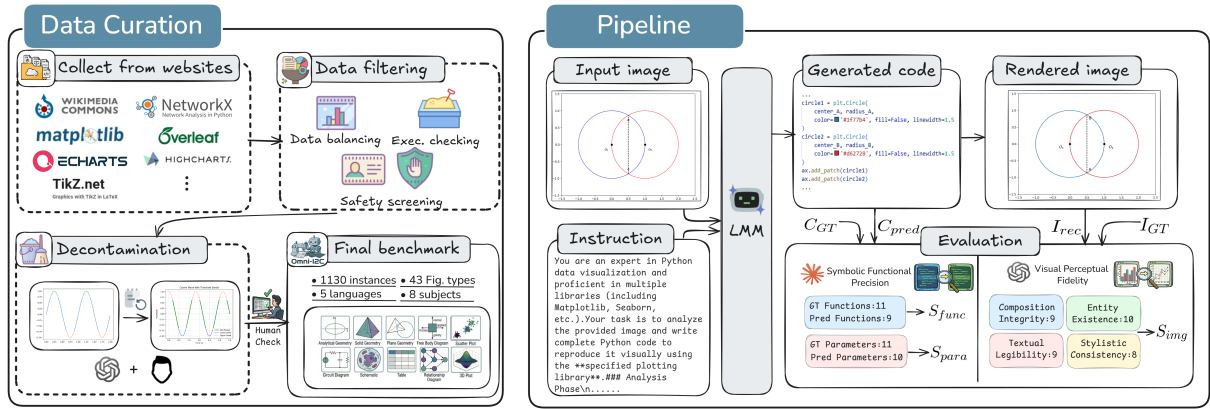


Figure 2: The Data Curation and pipeline of Omni-I2C. **Data Curation** (Left): From raw web collection to a refined benchmark with diverse themes and languages, ensured by strict filtering and human verification. **Pipeline** (Right): The proposed workflow for image-to-code generation and automatic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

Plot2Code (Wu et al., 2025), and Chart-Mimic (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 ex-

ecutable 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. A.

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. A.), 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 exclud-

ing 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. A.

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 the App. B.2.

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. H 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, 2025b), 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. B 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.F, 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 1inspace 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 raw parameter counts; for instance, the 32B Qwen3-VL consistently surpasses the larger 72B Qwen2.5-VL, illustrating a distinct generational leap. Most notably, Qwen3-235B 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.

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**, second best are underlined for proprietary models and open-weight models

| 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> | 99.5 | 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 | 71.9 | <u>94.3</u> | 49.0 | 82.0 |
| Gemini 3 Pro | <u>57.4</u> | 39.7 | <u>71.8</u> | 92.1 | 60.1 | 56.5 | 72.1 | 99.5 | <u>53.5</u> | 31.3 | <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 | 87.0 |
| 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 | <u>80.7</u> | 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-Instruct | 52.7 | 35.2 | 64.0 | 94.9 | 61.0 | 51.6 | 68.1 | 100.0 | 47.8 | 33.5 | 54.9 | <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-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 |
| Gamma3-27B | 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 |

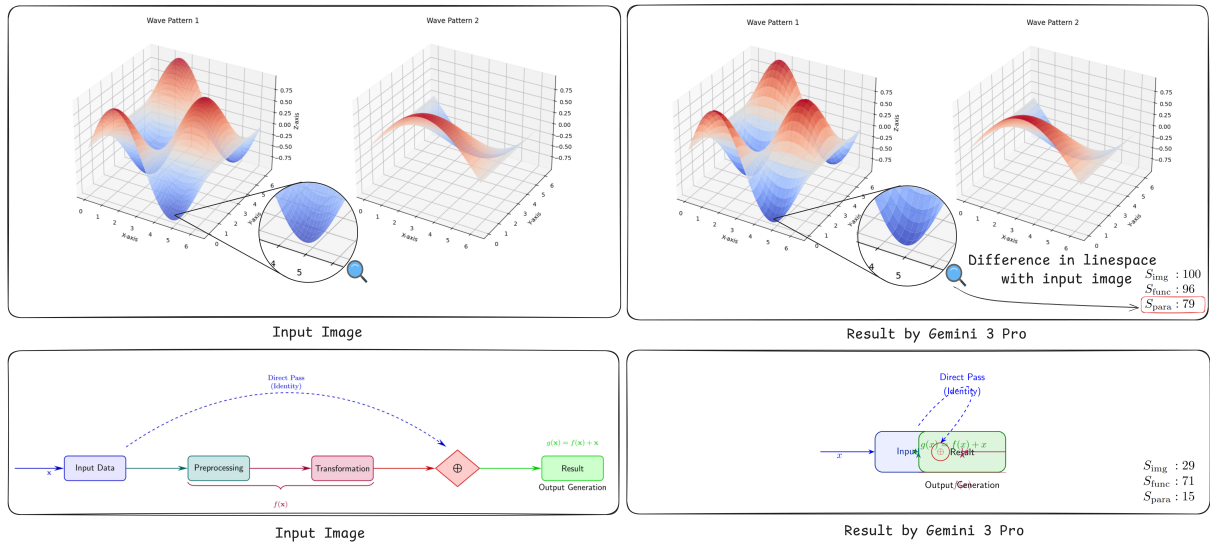


Figure 3: Representative examples from Gemini 3 pro.

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, and Qwen3-VL-235B—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 shown in App. B.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. C).

The results, summarized in Tab. 3, indicate that

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. G.

| 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 |

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 B.

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 2.5 Pro, GPT 4.1 mini, InternVL 3.5-8B, and Qwen3-VL-235B). We then construct a Majority-Vote (MV) ranking as a silver standard for expert consensus. As shown in Tab. 4, 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 re-

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. G.

| 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 | |

stricted 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.

5.2 Different Prompting Methods

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

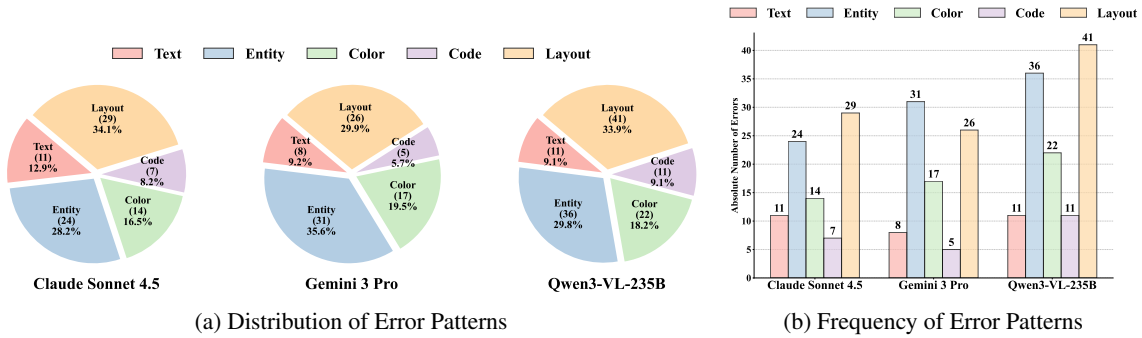


Figure 4: Comparative Analysis of Error Patterns in Claude 4.5 Sonnet, Gemini 3 Pro, and Qwen3-VL-235B-A22B-Instruct. Case studies and evaluations across various languages are detailed in the App. E.

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. D

| 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 | 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 | 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 |

lines 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 variance underscores that achieving "plug-and-play" stability remains an open challenge. An extended analysis of model behaviors is provided in App. D.

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. E.

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 "Grand Challenge" to guide the development of more robust, visually-grounded multimodal 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.

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, Zhifang 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 Schaekermann, 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*.

| | | |
|-----|--|-----|
| 674 | Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. <i>Advances in neural information processing systems</i> , 36:34892–34916. | 728 |
| 675 | | 729 |
| 676 | | 730 |
| 677 | | 731 |
| 678 | 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. <i>arXiv preprint arXiv:2310.02255</i> . | 732 |
| 679 | | 733 |
| 680 | | 734 |
| 681 | | 735 |
| 682 | | 736 |
| 683 | | 737 |
| 684 | 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. <i>CoRR</i> . | 738 |
| 685 | | 739 |
| 686 | | 740 |
| 687 | | 741 |
| 688 | | 742 |
| 689 | | 743 |
| 690 | OPENAI. 2025a. gpt-5. https://openai.com/index/introducing-gpt-5/ . | 744 |
| 691 | | 745 |
| 692 | OPENAI. 2025b. gpt-5.1. https://openai.com/index/gpt-5-1/ . | 746 |
| 693 | | 747 |
| 694 | 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 <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 20023–20070. | 748 |
| 695 | | 749 |
| 696 | | 750 |
| 697 | | 751 |
| 698 | | 752 |
| 699 | | 753 |
| 700 | | 754 |
| 701 | | 755 |
| 702 | Josselin S Roberts, Tony Lee, Chi H Wong, Michihiro Yasunaga, Yifan Mai, and Percy Liang. 2024. Image2struct: Benchmarking structure extraction for vision-language models. <i>Advances in Neural Information Processing Systems</i> , 37:115058–115097. | 756 |
| 703 | | 757 |
| 704 | | 758 |
| 705 | | 759 |
| 706 | | 760 |
| 707 | David Rogers and Mathew Hahn. 2010. Extended-connectivity fingerprints. <i>Journal of Chemical Information and Modeling</i> , 50(5):742–754. | 761 |
| 708 | | 762 |
| 709 | | 763 |
| 710 | 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 <i>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)</i> , pages 3956–3974. | 764 |
| 711 | | 765 |
| 712 | | 766 |
| 713 | | 767 |
| 714 | | 768 |
| 715 | | 769 |
| 716 | | 770 |
| 717 | | 771 |
| 718 | 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. <i>arXiv preprint arXiv:2510.17932</i> . | 772 |
| 719 | | 773 |
| 720 | | 774 |
| 721 | | 775 |
| 722 | | 776 |
| 723 | 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. <i>arXiv preprint arXiv:2503.19786</i> . | 777 |
| 724 | | 778 |
| 725 | | 779 |
| 726 | | 780 |
| 727 | | 781 |
| | | 782 |
| | | 783 |
| | 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. <i>Advances in Neural Information Processing Systems</i> , 37:95095–95169. | 784 |
| | | 785 |
| | | 786 |
| | | 787 |
| | 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. <i>arXiv preprint arXiv:2310.03731</i> . | 788 |
| | | 789 |
| | | 790 |
| | | 791 |
| | Weiyun Wang, Zhangwei Gao, Lixin Gu, Hengjun Pu, Long Cui, Xingguang Wei, Zhaoyang Liu, Linglin Jing, Shenglong Ye, Jie Shao, and 1 others. 2025. Internv1.5: Advancing open-source multimodal models in versatility, reasoning, and efficiency. <i>arXiv preprint arXiv:2508.18265</i> . | 792 |
| | | 793 |
| | | 794 |
| | | 795 |
| | 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 <i>Findings of the Association for Computational Linguistics: NAACL 2025</i> , pages 3006–3028. | 796 |
| | | 797 |
| | | 798 |
| | | 799 |
| | Cheng Yang, Chufan Shi, Yaxin Liu, Bo Shui, Junjie Wang, Mohan Jing, Linran Xu, Xinyu Zhu, Siheng Li, Yuxiang Zhang, and 1 others. 2024. Chart-mimic: Evaluating Imm’s cross-modal reasoning capability via chart-to-code generation. <i>arXiv preprint arXiv:2406.09961</i> . | 800 |
| | | 801 |
| | | 802 |
| | | 803 |
| | 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. <i>arXiv preprint arXiv:2502.14846</i> . | 804 |
| | | 805 |
| | | 806 |
| | | 807 |
| | Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik Narasimhan. 2024. tau-bench: A benchmark for tool-agent-user interaction in real-world domains. <i>arXiv preprint arXiv:2406.12045</i> . | 808 |
| | | 809 |
| | | 810 |
| | 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 <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 9556–9567. | 811 |
| | | 812 |
| | | 813 |
| | | 814 |
| | 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 <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 15134–15186. | 815 |
| | | 816 |
| | | 817 |
| | | 818 |

784 Chi Zhang, Haibo Qiu, Qiming Zhang, Yufei Xu, Zhix-
785 iong Zeng, Siqu Yang, Peng Shi, Lin Ma, and Jing
786 Zhang. 2025a. Perceptual-evidence anchored rein-
787 forced learning for multimodal reasoning. *arXiv*
788 *preprint arXiv:2511.18437*.

789 Chi Zhang, Haibo Qiu, Qiming Zhang, Zhixiong Zeng,
790 Lin Ma, and Jing Zhang. 2025b. Deepsketcher: In-
791 ternalizing visual manipulation for multimodal rea-
792 soning. *arXiv preprint arXiv:2509.25866*.

793 Yi-Fan Zhang, Xingyu Lu, Shukang Yin, Chaoyou
794 Fu, Wei Chen, Xiao Hu, Bin Wen, Kaiyu Jiang,
795 Changyi Liu, Tianke Zhang, and 1 others. 2025c.
796 Thyme: Think beyond images. *arXiv preprint*
797 *arXiv:2508.11630*.

798 Ziwei Zheng, Michael Yang, Jack Hong, Chenxiao
799 Zhao, Guohai Xu, Le Yang, Chao Shen, and Xing
800 Yu. 2025. Deepeyes: Incentivizing" thinking with
801 images" via reinforcement learning. *arXiv preprint*
802 *arXiv:2505.14362*.

803 Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou,
804 Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue
805 Ou, Yonatan Bisk, Daniel Fried, and 1 others.
806 2023. Webarena: A realistic web environment
807 for building autonomous agents. *arXiv preprint*
808 *arXiv:2307.13854*.

A 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.

A.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. 6.

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.

A.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 Idiomaticity. 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"

Table 6: 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/example_s/zh/index.html https://www.highcharts.com/demo |
| SVG | https://commons.wikimedia.org/wiki/Main_Page |

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.

A.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

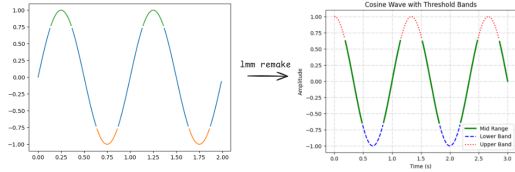


Figure 5: Examples of LLM-based code restructuring with foreground and numerical perturbations.

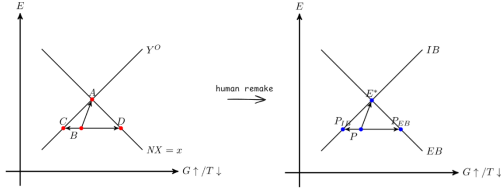


Figure 6: Comparison between original images and human-refined samples with modified key attributes.

887 training session prior to the mass annotation phase.
 888 Annotators are briefed on the “Structural Corre-
 889 spondence” requirement and trained on specific
 890 taxonomic guidelines to minimize ambiguity. A
 891 prime example of our classification protocol in-
 892 volves distinguishing between *Analytic Geometry*
 893 and *Function-related*. We enforced a strict math-
 894 ematical definition: a *Function-related* is defined
 895 by the property that every x value corresponds to a
 896 unique y value (i.e., the vertical line test), whereas
 897 *Analytic Geometry* figures are not bound by this
 898 constraint. Cases with ambiguous classifications
 899 are adjudicated through a majority voting mech-
 900 anism among senior annotators.

901 Following these protocols, annotators utilized
 902 a custom-built data collection platform that syn-
 903 chronizes code editing with real-time rendering
 904 to enforce a unified metadata schema. Upon up-
 905 loading the source code and corresponding visual
 906 artifacts, annotators labeled each entry across three
 907 hierarchical dimensions: *Subject*, *Figure Type*, and
 908 *Code Type*.

909 A.4 Figure Type Taxonomy

910 To rigorously evaluate the model’s code generation
 911 performance across diverse visual scenarios, Omni-
 912 I2C incorporates a granular classification system.
 913 Currently, Omni-I2C encompasses 43 figure types,
 914 covering core visual expressions in scientific re-
 915 search, engineering design, data analysis, and pro-
 916 fessional education.

917 The specific categories include: relationship-
 918 diagram, line-graph, molecular-formula, tables,
 919 bar-chart, plane-geometry, flow-chart, analytical-

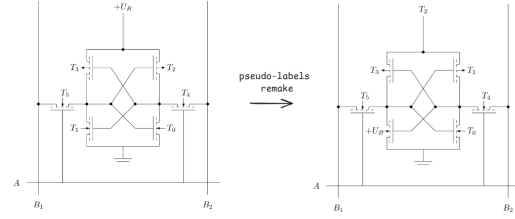


Figure 7: Visualization of counterfactual labels used to decouple perception from language priors.

```

# 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.

920 geometry, circuit, 3d-plot, equations-texts, scatter-
 921 plot, pie-chart, function-related, solid-geometry,
 922 schematic, heatmap, Area, radar-chart, Tree,
 923 Quiver, error-bar, venn-diagram, atom-model, Con-
 924 tour, box-chart, free-body-diagram, error-point,
 925 Density, map, Violin, Histogram, optics-ray-
 926 diagram, cell-structure, Graph, Algorithms,
 927 multi-graph, Phase-Diagram, block-diagram,
 928 physiological-process, magnetic-field-line and
 929 anatomy-diagram. Sparse or ambiguous categories
 930 are aggregated into a miscellaneous group to
 931 maintain statistical rigor.

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

934 A.5 Data Rewriting and Leakage Mitigation

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

938 First, we leverage LLMs to refactor the code’s
939 structural and stylistic attributes, such as layout
940 and aesthetic parameters, ensuring the code devi-
941 ates from its original online form while preserving
942 the target visualization’s core logic. Second, we
943 manually modify key visual attributes, including
944 label coordinates, color schemes, and textual con-
945 tent, to further distance the samples from their raw
946 source. Finally, we deliberately substitute mean-
947 ingful text with nonsensical or counter-factual strings
948 to decouple visual perception from linguistic priors,
949 forcing models to rely strictly on visual grounding
950 rather than common-sense guessing.

951 Illustrative examples of LLM refactoring, man-
952 ual modification, and counter-factual injection are
953 shown in Fig. 5, Fig. 6, and Fig. 7, respectively.
954 The specific prompt employed for LLM-based
955 refactoring is detailed in Fig. 8.

956 B Setups

957 B.1 Inference Setups

958 We evaluate 5 proprietary and 8 open-weight mod-
959 els in our study. We employ greedy decoding and
960 the maximum generation length is 16,384 tokens
961 for all models. For Gemini 2.5 pro, Gemini 3 pro
962 and GPT 5.1, we configure the “thinking effort” to
963 the “minimal”.

964 B.2 Evaluation Setups

965 **Details of Visual Perceptual Fidelity.** To en-
966 sure that the Visual Perceptual Fidelity (S_{img})
967 metric reflects the specific requirements of vari-
968 ous technical domains, we implement a Domain-
969 Aware Weighting Mechanism. Rather than ap-
970 plying a uniform average across all evaluation di-
971 mensions, this approach dynamically adjusts the
972 importance of different visual axes based on the
973 subject matter and figure archetype. The final
974 score S_{img} (on a 100-point scale) is calculated by
975 $S_{img} = (\sum s_i \times w_i) \times 10$, where $s_i \in [0, 10]$
976 represents the raw score for each axis and w_i de-
977 notes the percentage weight. The four assessment
978 axes include: Composition Integrity (s_{comp}), eval-
979 uating global spatial distribution; Textual Legibil-
980 ity (s_{text}), assessing label and symbolic accuracy;
981 Entity Existence (s_{entity}), ensuring critical visual
982 primitives are present; and Stylistic Consistency

(s_{style}), measuring aesthetic attributes like color
and line styles.

983 The assignment of these weighting profiles fol-
984 lows a hierarchical mapping logic. The Charts
985 & Plots profile is applied to archetypes such
986 as line-graphs, bar-charts, pie-charts, scatter-
987 plots, 3d-plots, and function-related visualiza-
988 tions. The Geometry profile encompasses plane-
989 geometry, solid-geometry, and analytical-geometry,
990 where coordinate-based properties are central.
991 Archetypes defined by structural connectivity, in-
992 cluding flow-charts, schematics, circuits, and
993 relationship-diagrams, are mapped to the Diagrams
994 profile. The Chemistry profile is assigned to
995 chemistry-related subjects to prioritize molecular
996 topology; for Math subjects, items containing tabu-
997 lar structures are evaluated under the **Tables** profile,
998 while others are mapped to the Equations profile
999 to emphasize formula rendering. For instances not
1000 covered by these archetypes, they will be pro-
1001 cessed with default profile. The specific weight
1002 configurations are detailed in Table 7. 1003
1004

Table 7: 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% |

1005 The rationale for this design is rooted in human
1006 check, as different image types prioritize different
1007 visual elements—for instance, structural connec-
1008 tivity is paramount for molecular formulas, while
1009 grid alignment is critical for tables. The prompt
1010 template employed to evaluate the visual percep-
1011 tual fidelity (S_{img}) of LMMs within the Omni-I2C
1012 framework is presented in Fig 9.

1013 **Prompts for Image-level Evaluation.** We detail
1014 the template prompts used in Tab. 3 in Fig 11.

1015 **Details of Symbolic Functional Precision.** To im-
1016 plement the Symbolic Functional Precision metrics,
1017 we formulate a structured evaluation prompt that
1018 leverages an LLM to serve as an Expert Code Logic
1019 Analyst. We employ the Python-specific prompt

```

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 Fidelity:** 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 4.

**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

(detailed in Fig. 10) as an example to illustrate the underlying evaluation logic.

C Human Study

To assess the biological plausibility and reliability of our visual perception metrics, we conduct a human alignment study in Tab. B.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 RMSE as meta-evaluation metrics to quantify the

```

# 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_{param}).
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.
* **Precision:** 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

congruence between our automated scoring system and human intuition.

D 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.

E 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 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

Table 8: 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 | 81.99 |
| | 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 | 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 | 83.04 | 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 |

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 constraints, 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 ge-

ometric reasoning capabilities to mitigate layout deviations, while simultaneously improving precise control over code generation.

E.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 in-

voke 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

1334 confirms that while LMMs are approaching syn- 1383
1335 tactic competence, mastering the underlying logic 1384
1336 and environmental constraints of executable code
1337 remains a formidable barrier across all scales and
1338 families.

1339 F Case Study

1340 We conduct a qualitative analysis comparing 13
1341 models on our benchmark. We select **Gemini 3 Pro**
1342 and **Qwen3-VL-235b-a22b-Instruct** for a detailed
1343 comparative case study due to their representative
1344 performance characteristics.

1345 **Case study in Python** We examine a repre- 1385
1346 sentative Python code generation scenario shown 1386
1347 in Fig 12 to analyze the reasoning behind these 1387
1348 scores. In terms of S_{img} (Visual Fidelity), Gemini 1388
1349 3 Pro demonstrates an advantage over Qwen3-VL- 1389
1350 235b-a22b-Instruct by maintaining superior global 1390
1351 visual consistency. The S_{func} metric acutely iden- 1391
1352 tifies critical mathematical errors in frequency esti- 1392
1353 mation within both generated code snippets, while 1393
1354 further highlighting a fundamental gap in logical 1394
1355 depth: Gemini 3 Pro successfully implements com- 1395
1356 plex threshold masking for segmentation, whereas 1396
1357 Qwen3-VL-235b-a22b-Instruct relies on a naive 1397
1358 overlaying of curves, failing to comprehend the 1398
1359 underlying conditional logic of the image. More- 1399
1360 over, fine-grained detection via S_{para} reveals that 1400
1361 Qwen3-VL-235b-a22b-Instruct exhibits lower pre- 1401
1362 cision in parameter control, incurring more errors 1402
1363 in stylistic details such as alpha compared to Gem-
1364 ini 3 Pro.

1365 **Case study in LaTeX** We examine a representa- 1403
1366 tive LaTeX code generation scenario shown in 1404
1367 Fig 13 to analyze the reasoning behind these 1405
1368 scores. In the LaTeX-based plotting task, Gemini 1406
1369 3 Pro achieves an S_{img} score of 81, surpassing 1407
1370 Qwen3-VL-235b-a22b-Instruct (71) in terms of 1408
1371 structural fidelity. The S_{func} evaluation character- 1409
1372 izes a shared limitation in semantic LaTeX imple- 1410
1373 mentation: both models fail to invoke professional 1411
1374 frameworks such as PGFPlots or its associated fill- 1412
1375 between library, defaulting to primitive TikZ op- 1413
1376 erations and manual coordinate systems. Specif- 1414
1377 ically, Gemini 3 Pro employs Bezier curves inte- 1415
1378 grated with region clipping for segmented filling, 1416
1379 whereas Qwen3-VL-235b-a22b-Instruct resorts to 1417
1380 hard-coded vertices, yielding functionally inconsis- 1418
1381 tent trajectories. Moreover, S_{para} metrics indicate 1419
1382 that Qwen3-VL-235b-a22b-Instruct suffers from

significant deviations in numerical parameters, in- 1385
cluding scaling and sampling density. 1386

1387 **Case study in HTML** We examine a repre- 1388
1389 sentative HTML code generation scenario shown 1389
1390 in Fig 14 to analyze the reasoning behind these 1390
1391 scores. For the HTML/ECharts task, Gemini 3 Pro 1391
1392 achieves an S_{func} score of 88, significantly surpass- 1392
1393 ing Qwen3-VL-235b-a22b-Instruct (56). Although 1393
1394 Gemini 3 Pro utilizes static SVG primitives instead 1394
1395 of the standard ECharts framework, it successfully 1395
1396 reconstructs essential functional components, in- 1396
1397 cluding titles, legends, and multi-series markers. 1397
1398 Conversely, Qwen3-VL-235b-a22b-Instruct fails to 1398
1399 include necessary script dependencies and critical 1399
1400 y-axis metadata, resulting in a fragmented struc- 1400
1401 tural representation. Furthermore, fine-grained 1401
1402 S_{para} evaluations highlight a fundamental error 1402
in data scaling by Qwen3-VL-235b-a22b-Instruct, 1403
which fails to align the y-axis range with the ground 1404
truth rainfall metrics. 1405

1406 **Case study in Svg** We examine a representative 1403
1407 Svg code generation scenario shown in Fig 15 to 1404
1408 analyze the reasoning behind these scores. In the 1405
1409 evaluation of SVG vector graphics generation, both 1406
1410 Gemini 3 Pro (S_{img} : 85) and Qwen3-VL-235b- 1407
1411 a22b-Instruct (S_{img} : 80) demonstrate exceptional 1408
1412 capabilities in structural decomposition and primi- 1409
1413 tive reconstruction. Both models achieved an S_{func} 1410
1414 score of 94, indicating a shared proficiency in ac- 1411
1415 curately identifying and mapping core functional 1412
1416 units, including irregular polygon paths, closed 1413
1417 regions, and textual labels. However, in-depth 1414
1418 parametric assessment reveals a significant perfor- 1415
1419 mance divergence in numerical precision: Gemini 1416
1420 3 Pro achieves an S_{para} score of 69, successfully 1417
1421 aligning textual content, ViewBox attributes, and 1418
1422 critical geometric coordinates with high fidelity. 1419
1423 In contrast, Qwen3-VL-235b-a22b-Instruct yields 1420
an S_{para} of 47, exhibiting a substantially lower 1421
success rate in matching complex path details, fill 1422
color values, and specific attribute groups. 1423

1424 Furthermore, we provide qualitative rendering 1424
1425 samples for all models benchmarked in our primary 1425
1426 performance table. As illustrated in Fig 16, Fig 17, 1426
1427 these examples offer a direct comparison of model 1427
1428 outputs on specific data instances. 1428

G 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.

G.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 Table 4, 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 Table 4, 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.

G.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 Table 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.

H 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.

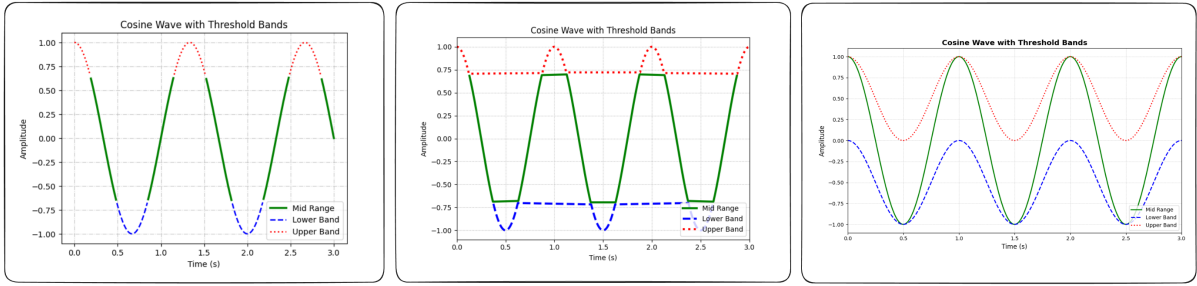
Task:
Compare the Ground Truth Image and the Generated Image. Based ONLY on their visual information, evaluate their similarity.

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.
```

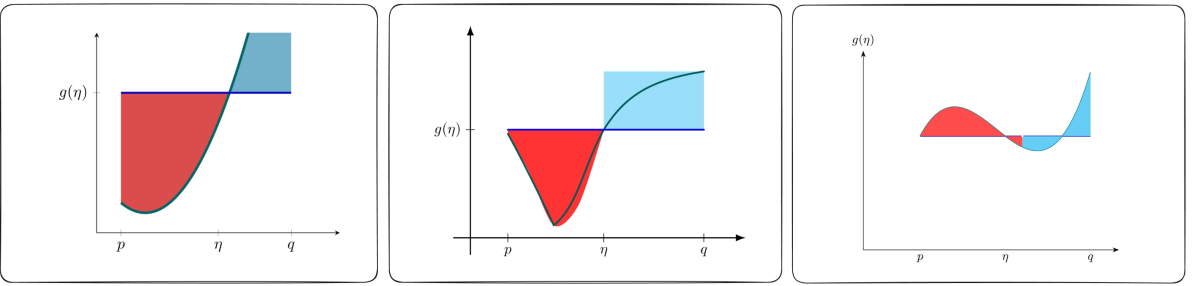
(d) Prompt of Chart2code

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



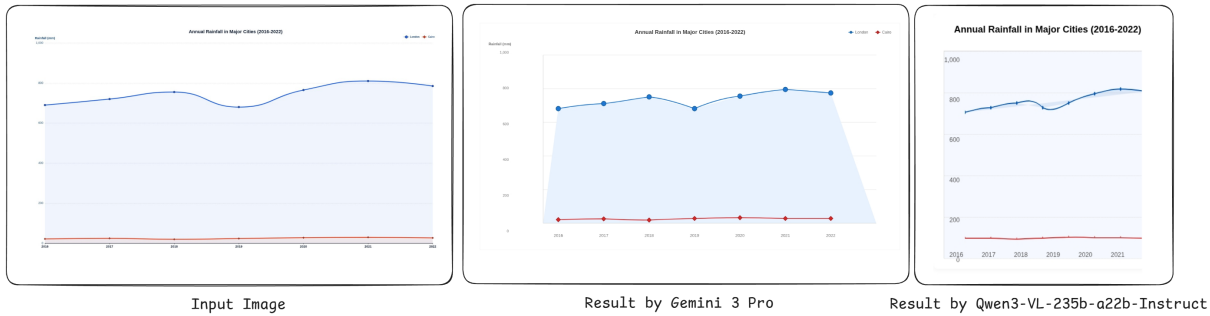
| Gemini 3 Pro | Qwen3-VL-235b-a22b-Instruct |
|--|---|
| <p>S_img: 76</p> <p>S_func: 50</p> <p>S_para: 30</p> <p>GT uses 7 functional units: (1) Generate x with linspace(0, 3.0, 300), (2) Calculate y = cos(1.5*pi*x), (3) Create y_upper using masked_where for values < threshold_upper, (4) Create y_lower 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_img: 66</p> <p>S_func: 21</p> <p>S_para: 27</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 mid-range 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> |

Figure 12: Case of Python



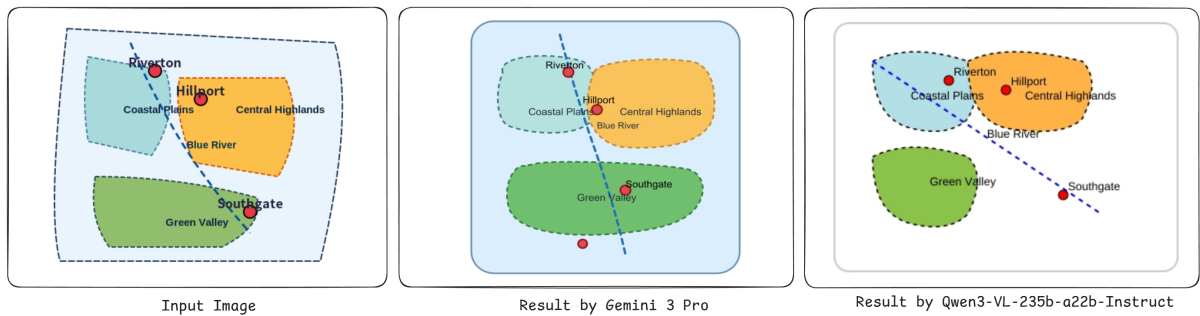
| Gemini 3 Pro | Qwen3-VL-235b-a22b-Instruct |
|---|---|
| <p>S_img: 81</p> <p>S_func: 60</p> <p>S_para: 25</p> <p>The generated code does not use the PGFPlots environment, which is the core structural component of the ground truth. Instead, it uses basic tikz commands and manual coordinate systems. 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_img: 71</p> <p>S_func: 42</p> <p>S_para: 13</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) adplot for the quadratic curve, (4) adplot for the constant line, (5) fillbetween logic with split and segment styling, (6) automatic region detection and coloring. Gen has 4 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> |

Figure 13: Case of LaTeX



| | |
|--|---|
| <p>Gemini 3 Pro</p> <p>S_img: 78 composition: 5 textual: 9 entity: 10 stylistic: 7</p> <p>S_func: 88 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_para: 58 Critical parameters assessed: (1) Title text: 'GT Annual Rainfall in Major Cities (2016-2022)' matches Gen exactly. (2) X-axis years: GT ['2016', '2017', '2018', '2019', '2020', '2021', '2022'] all match Gen. (3) Y-axis label: 'GT Rainfall (mm)' matches Gen. (4) Legend items: GT ['London', 'Cairo'] match Gen. ... Total key parameters: Title (1), X-labels (7), Y-label (1), Legend labels (2), London data (7), Cairo data (7), London line color (1), Cairo line color (1), London marker type (1), Cairo marker type (1), London area color (1), Y-axis scale (1) = 31 parameters. Matched: Title (1), X-labels (7), Y-label (1), Legend labels (2), London data (3 correct: indices 3,5,6), Cairo data (2 correct: indices 1,4), London marker (1), Cairo marker (1) = 18 matched exactly.</p> | <p>Qwen3-VL-235b-a22b-Instruct</p> <p>S_img: 68 composition: 5 textual: 7 entity: 8 stylistic: 7</p> <p>S_func: 56 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 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> <p>S_para: 29 Gen matches the exact years on x-axis and city names in legend. However, y-axis labels are scaled to 0-1000 mm instead of reflecting true data range (London peaks ~810 mm, Cairo ~30 mm), making the rainfall values visually misrepresented. Line and area colors differ (#408080 vs #208040, #d32f2f vs #e84c3d). Tooltip, axis title, and exact data-point values are absent. Thus, most critical data parameters are incorrect or missing.</p> |
|--|---|

Figure 14: Case of HTML



| | |
|---|---|
| <p>Gemini 3 Pro</p> <p>S_img: 85 composition: 7 textual: 9 entity: 10 stylistic: 8</p> <p>S_func: 94 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_para: 69 GT parameters identified: 7 text contents ('Riverton', 'Hillport', 'Southgate', 'Coastal Plains', 'Central Highlands', 'Green Valley', 'Blue River'); 3 radii (r=9 each); 3 cx (180,290,320); 3 cy (120,160,320). Total N_vars=16. In Gen: All 7 texts exact. All 3 r=9 exact. cx: 180=180 match, 290=290, 320=320 (3 match). cy: 110≠120, 180≠160, 330≠320 (0 matches). viewBox nearly identical (no scaling, proportional-exact required). Colors semantically consistent across elements. Path d details differ but excluded from params (logic only). M_vars=7 (texts) + 3 (r) + 1 (cx) + 0 (cy) =11.</p> | <p>Qwen3-VL-235b-a22b-Instruct</p> <p>S_img: 80 composition: 6 textual: 9 entity: 10 stylistic: 7</p> <p>S_func: 94 GT decomposition: The GT code contains the following functional units: (1) country-outline path, fill=#fff, stroke=#408080, stroke-width=2, stroke-dasharray=8,4; (2) region-1 fill=#d32f2f, stroke=#d32f2f, stroke-dasharray=8,4; (3) region-2 fill=#e84c3d, stroke=#e84c3d, stroke-dasharray=8,4; (4) region-3 fill=#208040, stroke=#208040, stroke-dasharray=8,4; (5) river path with quadratic curves; (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> <p>S_para: 47 GT Parameter identification: (1) country-outline: d path, fill=#fff, stroke=#408080, stroke-width=2, stroke-dasharray=8,4; (2) region-1 fill=#d32f2f, stroke=#d32f2f, stroke-dasharray=8,4; (3) region-2 fill=#e84c3d, stroke=#e84c3d, stroke-dasharray=8,4; (4) region-3 fill=#208040, stroke=#208040, stroke-dasharray=8,4; (5) river: stroke-dasharray=10,6; (6-8) three circles: cx/cy coordinates, r=9, fill=#e84c3d, stroke=#000000, stroke-width=2; (9-11) city labels: text content (Riverton, Hillport, Southgate), x/y coordinates, font-size=18, fill=#222222; (12-14) region labels: text content (Coastal Plains, Central Highlands, Green Valley), font-size=15, fill=#000000; (15) river label: text='Blue River', font-size=14, fill=#000000. ... Exact matches: 7 text content strings. Mismatches: all fill colors (4), all stroke-dasharray values (4), circle radius (1), all font sizes (6), river stroke color (1), river stroke-width (1). Total matches: 7 out of 15 parameter groups.</p> |
|---|---|

Figure 15: Case of Svg

Case Analysis: 6cf0140514bd4066ae6711234dc518c6

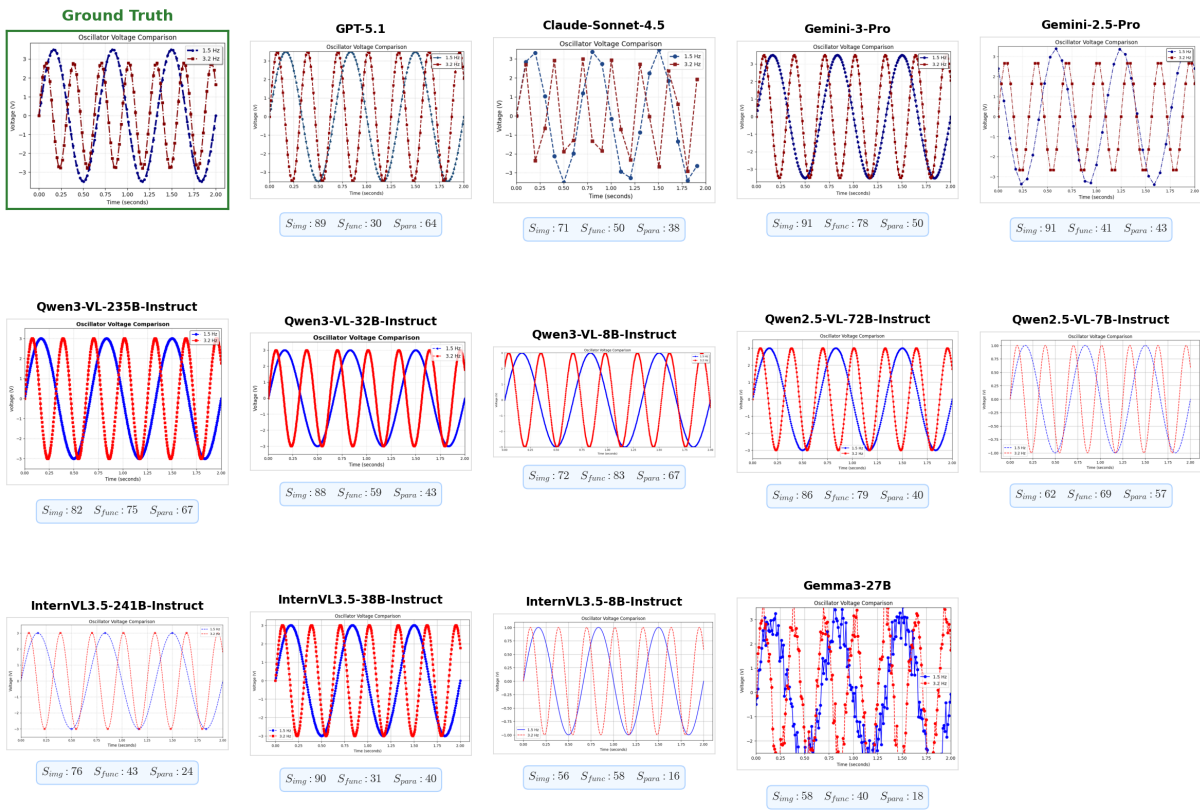


Figure 16: Case 1 of render_samples for all models benchmarked in our Table 2.

Case Analysis: AsymptoteGraphicPipeline_Circle_Geometry_2-88

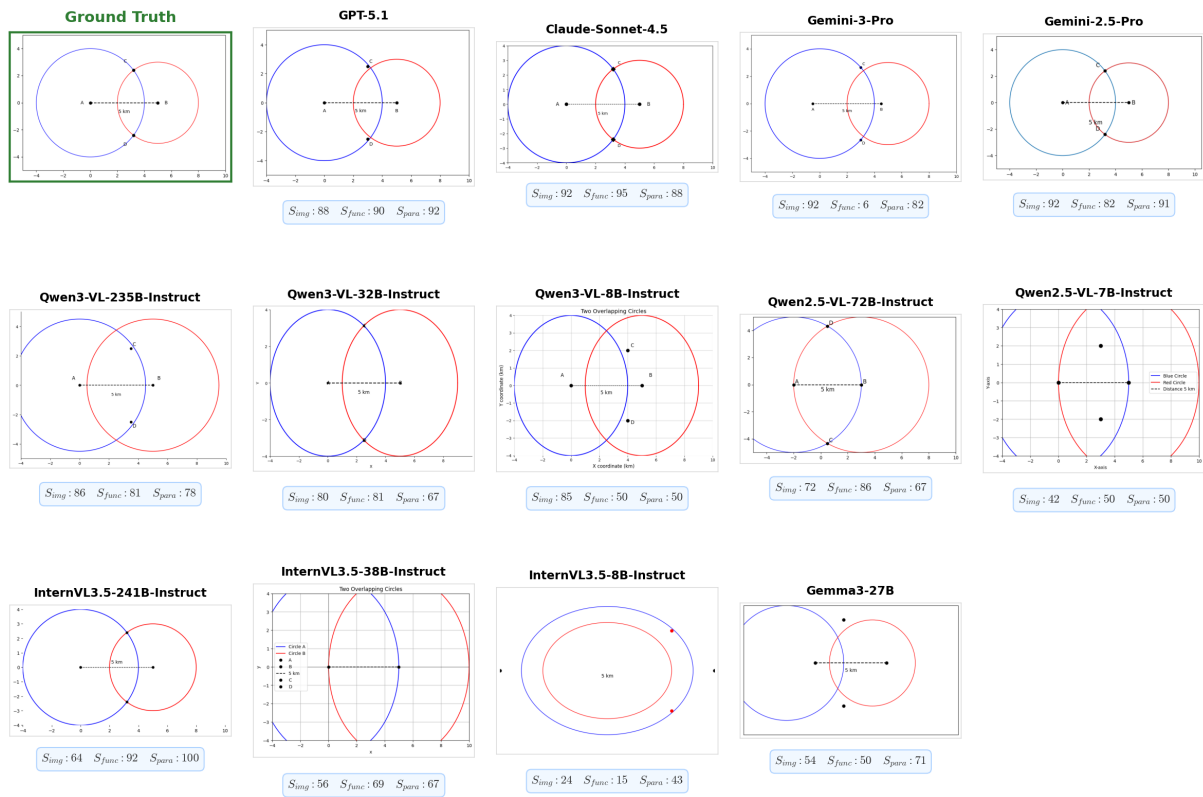


Figure 17: Case 2 of render_samples for all models benchmarked in our Table 2.

Omni-I2C Annotation Re-evaluation (50 Samples)

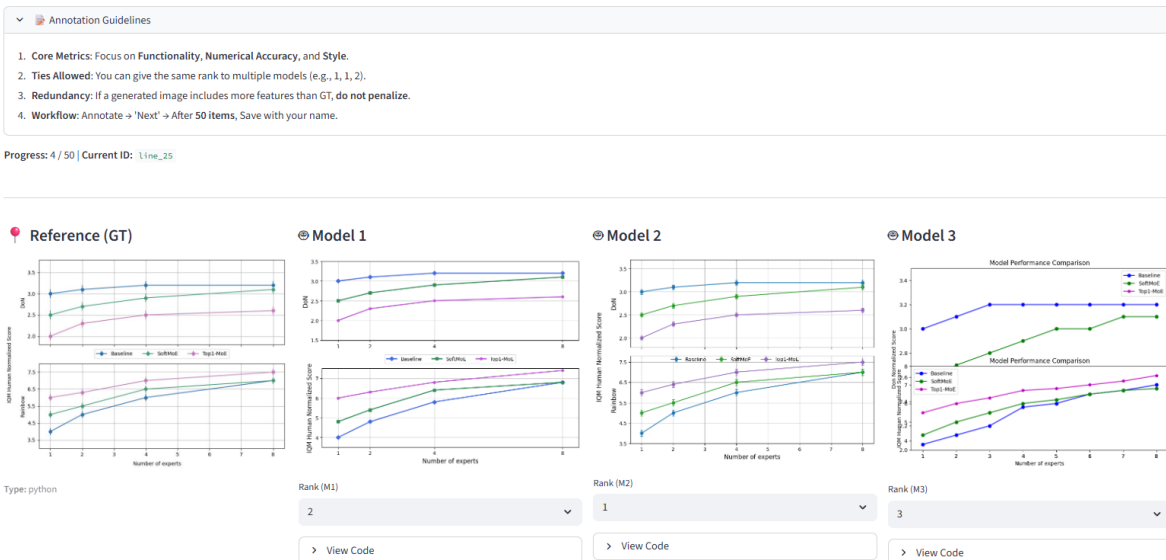


Figure 18: A screenshot of the human evaluation questionnaire.

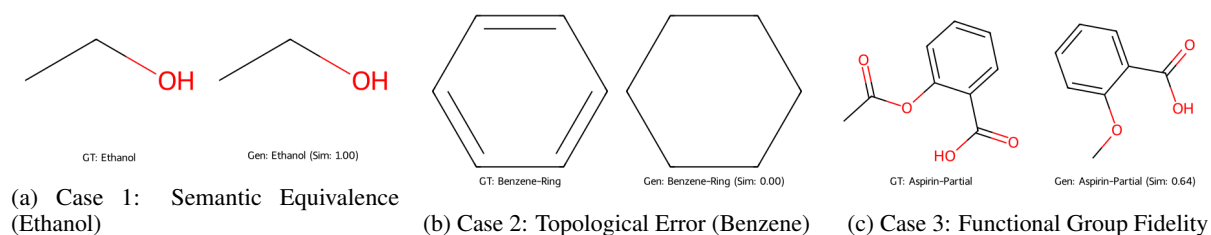


Figure 19: Visual comparison of SMILES reconstruction cases. (a) Illustrates identical topology despite different string sequences. (b) Highlights critical perception errors in aromaticity captured by Tanimoto similarity. (c) Demonstrates the metric's sensitivity to localized functional group deletions.

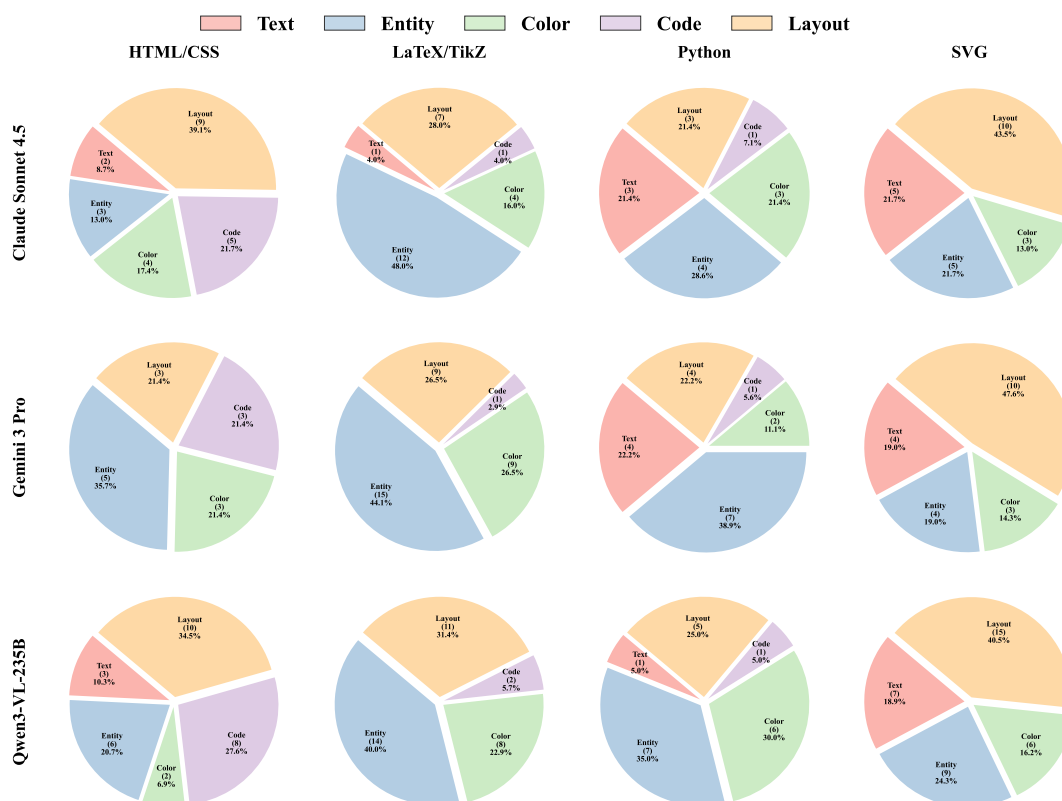
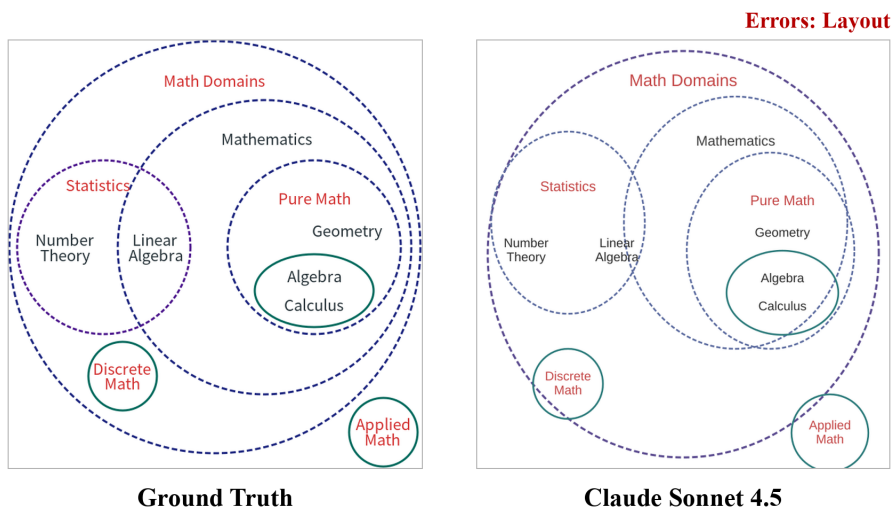
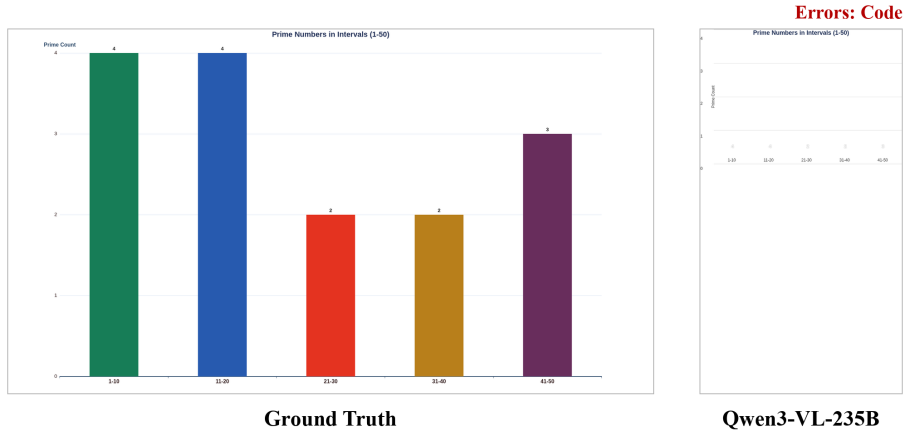


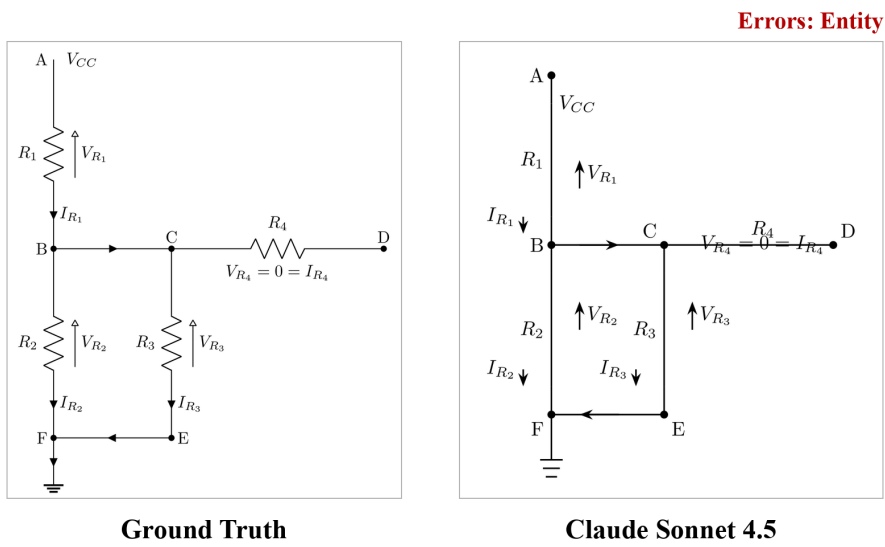
Figure 20: Fine-grained error distribution across languages



(a) Layout Error



(b) Code Error



(c) Entity Error

Figure 21: Representative error cases produced by different models.

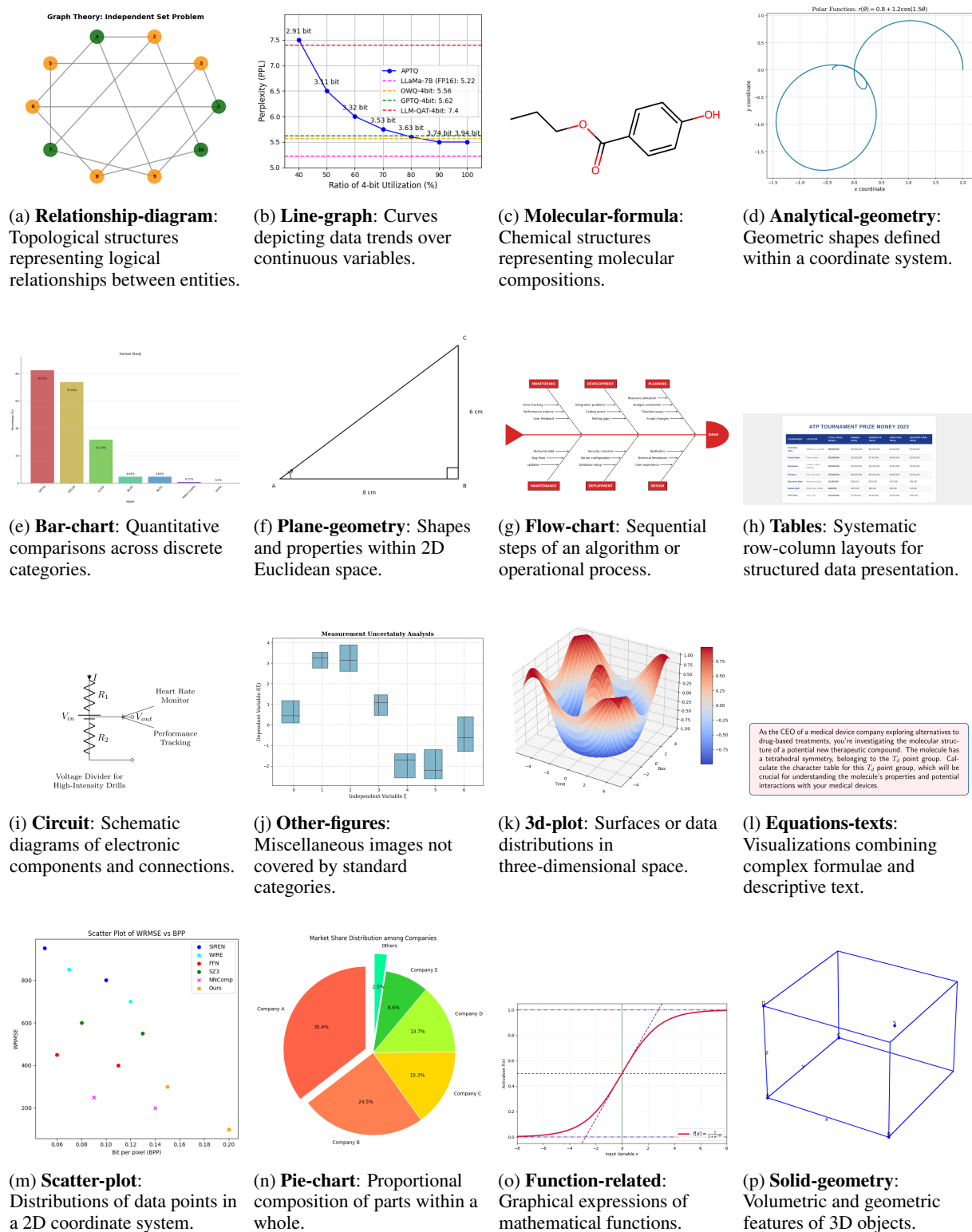
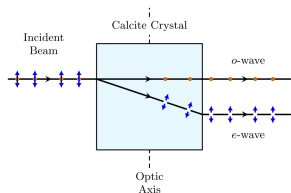
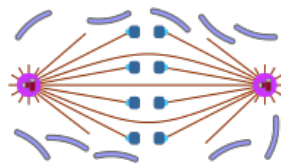


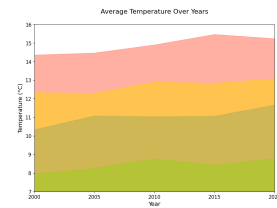
Figure 23: Omni-I2C dataset: Figure Type Taxonomy gallery (Part 1).



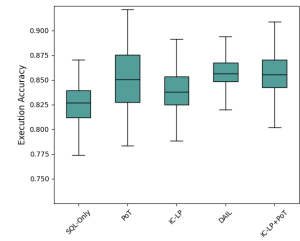
(a) **Schematic:** Simplified diagrams illustrating systems or structures.



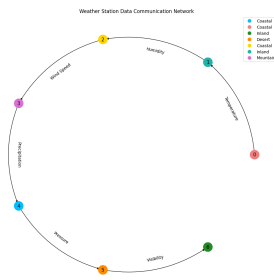
(b) **Physiological-process:** Dynamic representations of biological activities.



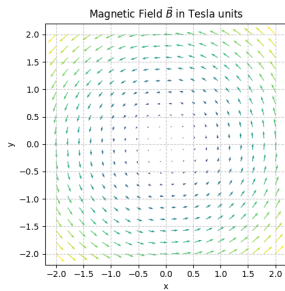
(c) **Area:** Trend charts with filled area sections.



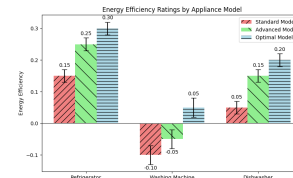
(d) **Box-chart:** Statistical summaries showing quartiles and outliers.



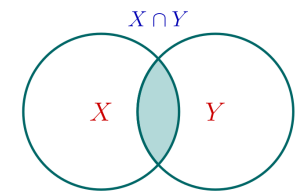
(e) **Graph:** General node-edge connectivity representations.



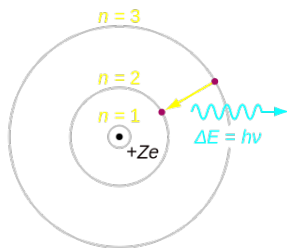
(f) **Quiver:** Distribution of vector fields in a spatial domain.



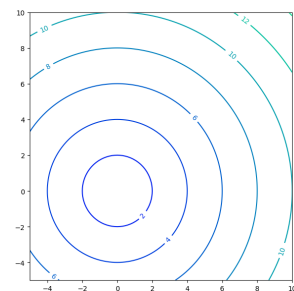
(g) **Error-bar:** Statistical charts indicating data uncertainty ranges.



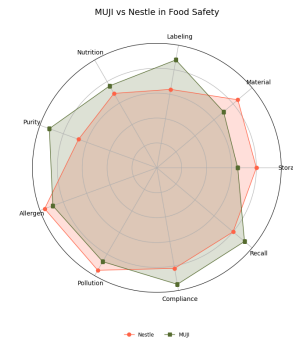
(h) **Venn-diagram:** Overlapping sets showing logical relations.



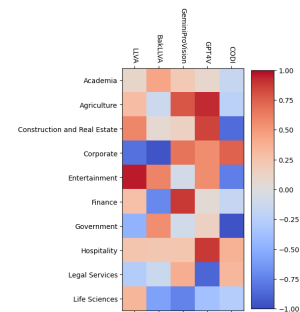
(i) **Atom-model:** Physical models of atomic structures and electrons.



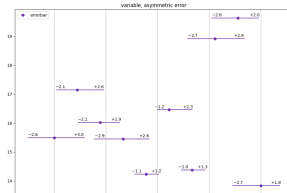
(j) **Contour:** Numerical distributions with isoline features.



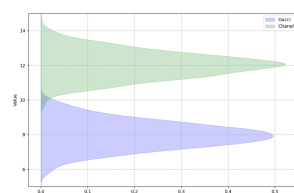
(k) **Radar-chart:** Comparative visualizations of multi-dimensional data.



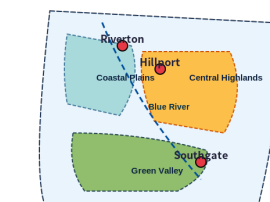
(l) **Heatmap:** Matrix intensities represented by color gradients.



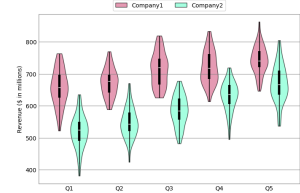
(m) **Error-point:** Data points with designated error offsets.



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

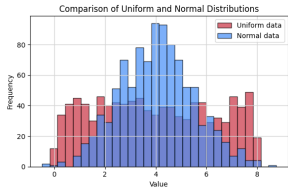


(o) **Map:** Visualizations of geographic locations or regions.

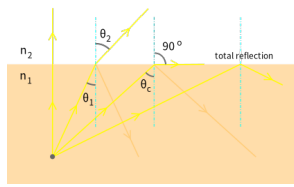


(p) **Violin:** Distribution plots combining density and box plot features.

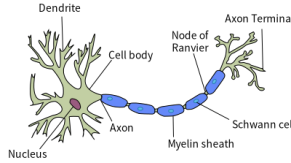
Figure 24: Omni-I2C dataset: Figure Type Taxonomy gallery (Part 2).



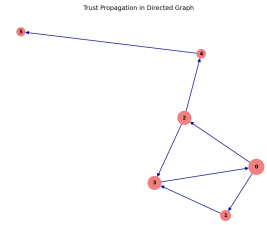
(a) **Histogram:** Statistical bars for frequency or count distributions.



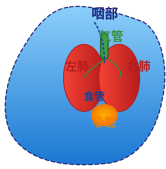
(b) **Optics-ray-diagram:** Propagation paths of light in optical systems.



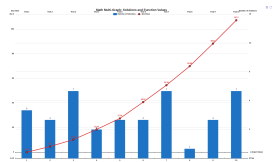
(c) **Cell-structure:** Diagrams illustrating internal biological cell parts.



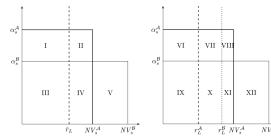
(d) **Algorithms:** Visualization of algorithmic logic or pseudocode.



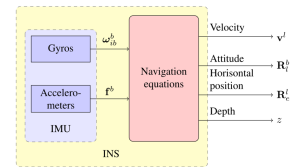
(e) **Anatomy-diagram:** Detailed structural diagrams of biological anatomy.



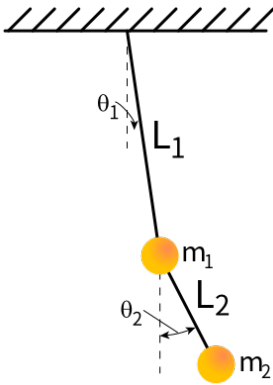
(f) **Multi-graph:** Composite images containing multiple sub-figures.



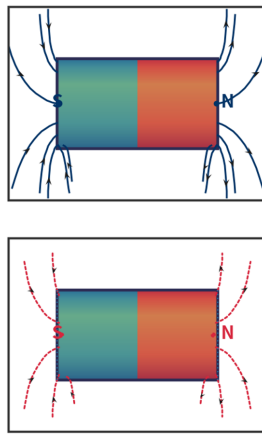
(g) **Phase-Diagram:** State charts describing system phase transitions.



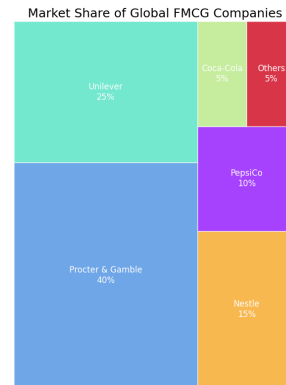
(h) **Block-diagram:** Functional relationships between system modules.



(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
import matplotlib.pyplot as plt
import networkx as nx

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).

## I Datasheet for Omni-I2C

### I.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.

### I.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.

### I.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,

|      |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
|------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|--|--|
| 1667 | TikZ.net, TeXample.net) to capture genuine practical utility.                                                                                                                                                                                                                                                                                                                                                                                              |  |  |
| 1668 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1669 | <b>2. What mechanisms or procedures were used to collect the data?</b>                                                                                                                                                                                                                                                                                                                                                                                     |  |  |
| 1670 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1671 | <b>A2:</b> We utilized a taxonomy-driven approach, pre-defining a comprehensive schema to guide the search and ensure broad coverage across subjects and languages.                                                                                                                                                                                                                                                                                        |  |  |
| 1672 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1673 | <b>3. If the dataset is a sample from a larger set, what was the sampling strategy?</b>                                                                                                                                                                                                                                                                                                                                                                    |  |  |
| 1674 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1675 | <b>A3:</b> Purposeful sampling was used to address gaps in existing benchmarks, with 5 languages and 43 distinct figure types.                                                                                                                                                                                                                                                                                                                             |  |  |
| 1676 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1677 | <b>4. Who was involved in the data collection process?</b>                                                                                                                                                                                                                                                                                                                                                                                                 |  |  |
| 1678 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1679 | <b>A4:</b> 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. |  |  |
| 1680 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1681 | <b>5. Over what timeframe was the data collected?</b>                                                                                                                                                                                                                                                                                                                                                                                                      |  |  |
| 1682 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1683 | <b>A5:</b> The collection, manual auditing, and decontamination process took approximately one month.                                                                                                                                                                                                                                                                                                                                                      |  |  |
| 1684 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1685 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1686 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1687 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1688 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1689 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1690 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1691 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1692 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1693 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1694 | <b>I.4 Preprocessing/cleaning/labeling</b>                                                                                                                                                                                                                                                                                                                                                                                                                 |  |  |
| 1695 | <b>1. Was any preprocessing/cleaning/labeling of the data done?</b>                                                                                                                                                                                                                                                                                                                                                                                        |  |  |
| 1696 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1697 | <b>A1:</b> 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.                                                                                                                                                                                                          |  |  |
| 1698 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1699 | <b>2. Was a data decontamination strategy employed?</b>                                                                                                                                                                                                                                                                                                                                                                                                    |  |  |
| 1700 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1701 | <b>A2:</b> 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.                                                                                                                                                                                               |  |  |
| 1702 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1703 | <b>3. Is the software used to preprocess/clean/label the instances available?</b>                                                                                                                                                                                                                                                                                                                                                                          |  |  |
| 1704 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1705 | <b>A3:</b> The evaluation pipeline and sandbox configurations will be provided upon publication.                                                                                                                                                                                                                                                                                                                                                           |  |  |
| 1706 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1707 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1708 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1709 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1710 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1711 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1712 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1713 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1714 | <b>I.5 Uses</b>                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1715 | <b>1. Has the dataset been used for any tasks already? If so, please provide a description.</b>                                                                                                                                                                                                                                                                                                                                                            |  |  |
| 1716 |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
|      | <b>A1:</b> 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.                                                                                                                                                                                                                                                        |  |  |
|      |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
|      | <b>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.</b>                                                                                                                                                                                                                                                                                                       |  |  |
|      |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
|      | <b>A2:</b> N/A.                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
|      | <b>3. What (other) tasks could the dataset be used for?</b>                                                                                                                                                                                                                                                                                                                                                                                                |  |  |
|      |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
|      | <b>A3:</b> Omni-I2C can be used for evaluating multimodal reasoning, precise coordinate modeling in SVG/HTML, and specialized chemical informatics knowledge via SMILES.                                                                                                                                                                                                                                                                                   |  |  |
|      |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
|      | <b>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?</b>                                                                                                                                                                                                                                                  |  |  |
|      |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
|      | <b>A4:</b> 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.                                                                                                                                                                                                                   |  |  |
|      |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
|      | <b>5. Are there tasks for which the dataset should not be used? If so, please provide a description.</b>                                                                                                                                                                                                                                                                                                                                                   |  |  |
|      |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
|      | <b>A5:</b> N/A.                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
|      |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
|      | <b>I.6 Distribution</b>                                                                                                                                                                                                                                                                                                                                                                                                                                    |  |  |
|      | <b>1. Will the dataset be distributed to third parties outside of the entity? If so, please provide a description.</b>                                                                                                                                                                                                                                                                                                                                     |  |  |
|      |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
|      | <b>A1:</b> Yes, the dataset will be publicly available for the research community.                                                                                                                                                                                                                                                                                                                                                                         |  |  |
|      |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
|      | <b>2. How will the dataset will be distributed? Does the dataset have a digital object identifier (DOI)?</b>                                                                                                                                                                                                                                                                                                                                               |  |  |
|      |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
|      | <b>A2:</b> Omni-I2C will be distributed via the official project website and GitHub.                                                                                                                                                                                                                                                                                                                                                                       |  |  |
|      |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
|      | <b>3. When will the dataset be distributed?</b>                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
|      |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
|      | <b>A3:</b> The dataset will be distributed once the paper is accepted after peer review.                                                                                                                                                                                                                                                                                                                                                                   |  |  |
|      |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
|      | <b>4. Will the dataset be distributed under a copyright or other license?</b>                                                                                                                                                                                                                                                                                                                                                                              |  |  |
|      |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
|      | <b>A4:</b> It will be distributed under the Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International (CC BY-NC-SA 4.0) License.                                                                                                                                                                                                                                                                                                             |  |  |
|      |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
|      | <b>5. Have any third parties imposed IP-based or other restrictions on the data associated with the instances?</b>                                                                                                                                                                                                                                                                                                                                         |  |  |
|      |                                                                                                                                                                                                                                                                                                                                                                                                                                                            |  |  |
|      | <b>A5:</b> No.                                                                                                                                                                                                                                                                                                                                                                                                                                             |  |  |

1767 **6. Do any export controls or other regulatory**  
1768 **restrictions apply to the dataset or to individual**  
1769 **instances?**

1770 **A6:** No.

## 1771 **I.7 Maintenance**

1772 **1. Who will be supporting/hosting/maintaining**  
1773 **the dataset?**

1774 **A1:** The authors.

1775 **2. How can the owner/curator/manager of the**  
1776 **dataset be contacted (e.g., email address)?**

1777 **A2:** Email addresses will be provided on the  
1778 project homepage post-publication.

1779 **3. Is there an erratum? If so, please provide a**  
1780 **link or other access point.**

1781 **A3:** Any errata will be posted on the project  
1782 GitHub repository.

1783 **4. Will the dataset be updated? If so, please**  
1784 **describe how often, by whom, and how updates**  
1785 **will be communicated to users?**

1786 **A4:** Yes. The authors will periodically update  
1787 the benchmark with new instances and categories  
1788 based on community feedback.

1789 **5. Will older versions of the dataset continue to**  
1790 **be supported/hosted/maintained? If so, please**  
1791 **describe how.**

1792 **A5:** Older versions will be archived on GitHub  
1793 to ensure researchers can reproduce results from  
1794 previous studies.

1795 **6. If others want to extend/augment/build**  
1796 **on/contribute to the dataset, is there a mech-**  
1797 **anism for them to do so?**

1798 **A6:** N/A.