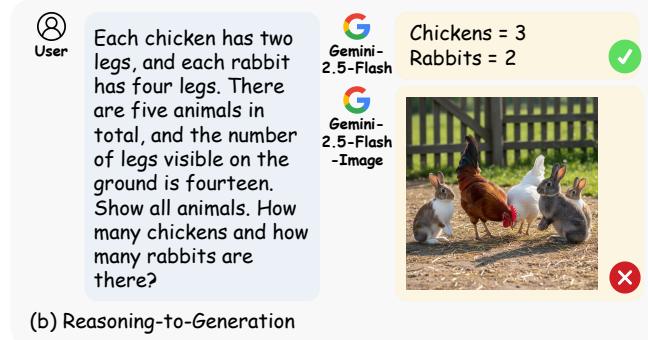
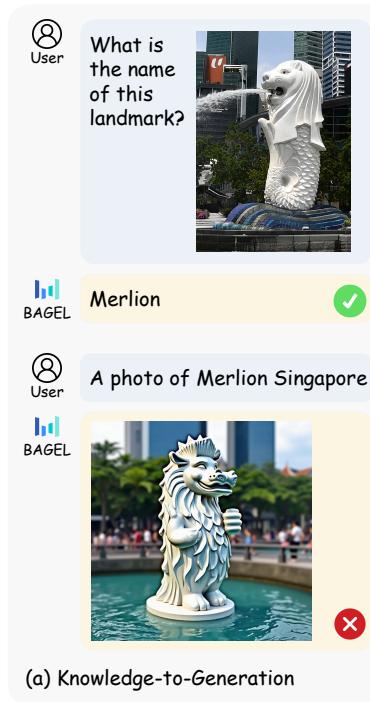
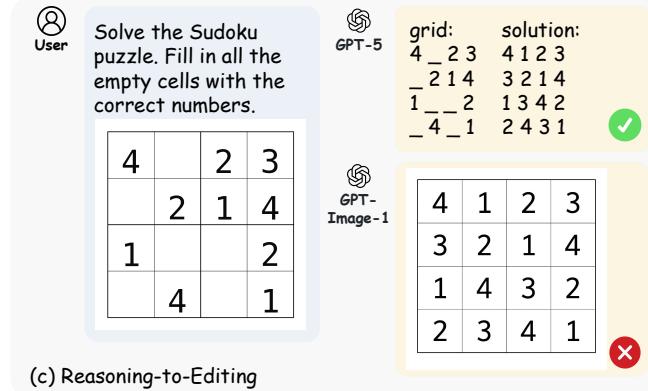


# GIR-BENCH: VERSATILE BENCHMARK FOR GENERATING IMAGES WITH REASONING

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
Anonymous authors  
Paper under double-blind review



(b) Reasoning-to-Generation



(c) Reasoning-to-Editing

Figure 1: Illustration examples of GIR-Bench, which highlight misalignments between the reasoning and generation capabilities of state-of-the-art unified multimodal models.

## ABSTRACT

Unified multimodal models integrate the reasoning capacity of large language models with both image understanding and generation, showing great promise for advanced multimodal intelligence. However, the community still lacks a rigorous reasoning-centric benchmark to systematically evaluate the alignment between understanding and generation, and their generalization potential in complex visual tasks. To this end, we introduce **GIR-Bench**, a comprehensive benchmark that evaluates unified models across three complementary perspectives. Firstly, we explore whether models can consistently leverage the same knowledge for both understanding and generation (GIR-Bench-Uni). Secondly, we investigate whether models can perform reasoning-centric text-to-image generation that requires applying logical constraints and implicit knowledge to generate faithful visual content (GIR-Bench-T2I). Thirdly, we evaluate whether models can handle multi-step reasoning in editing (GIR-Bench-Edit). For each subset, we carefully design different task-specific evaluation pipelines tailored for each task. This enables fine-grained and interpretable evaluation while mitigating biases from the prevalent MLLM-as-a-Judge paradigm. Extensive ablations over various unified models and generation-only systems have shown that: Although unified models are more

054 capable of reasoning-driven visual tasks, they still exhibit a persistent gap between  
 055 understanding and generation. The data and code for GIR-Bench are available at  
 056 <https://anonymous.4open.science/r/GIR-Bench-7E40>.  
 057

## 058 1 INTRODUCTION

061 Image generation and editing techniques (Stability AI, 2024; Labs, 2024; Batifol et al., 2025; Liu  
 062 et al., 2025) have evolved rapidly, demonstrating strong capabilities in producing high-quality vi-  
 063 sual content aligned with explicit prompts. However, existing models still struggle with complex  
 064 visual generation tasks that require multi-step reasoning. This limitation has triggered the research  
 065 focus toward unified multimodal models, where a single model supports both image generation and  
 066 understanding. By leveraging the intrinsic knowledge and reasoning abilities of multimodal large  
 067 language models (MLLM), such unified approaches enable richer expressiveness and more control-  
 068 lable image generation. Recent breakthroughs, such as Gemini-2.5-Flash-Image (DeepMind, 2025)  
 069 and GPT-Image (OpenAI, 2025b), further highlight the transformative potential of this paradigm,  
 070 showing that unified models can fundamentally reshape real-world applications, empowering users  
 071 to accomplish complex visual tasks through natural language interaction. Compared to generation-  
 072 only models, unified models (Deng et al., 2025; Chen et al., 2025a; Xie et al., 2025) promise sub-  
 073 stantial gains. With these advances comes a critical **research question**: *how do we comprehensively  
 074 evaluate the extent to which enhanced multimodal understanding improves generation capabilities?*

075 Earlier image generation benchmarks (Ghosh et al., 2023; Huang et al., 2023) mainly focused on  
 076 object attributes and compositional evaluation, but remained limited to shallow mappings between  
 077 text and visual content. Recent works (Zhao et al., 2025; Niu et al., 2025; Sun et al., 2025) have  
 078 attempted to consider the reasoning capability. However, existing benchmarks remain limited in both  
 079 their evaluation dimensions and protocols, making them inadequate for capturing the full capabilities  
 080 of unified models.

- 081 • **For evaluation dimensions**, they cannot quantify the alignment between reasoning and generation  
 082 within unified models. As shown in Figure 1, we observe consistent misalignments in transferring  
 083 knowledge to generation, reasoning to generation, and reasoning to editing. For instance, a unified  
 084 model may correctly recognize a real-world entity (*e.g.*, the Merlion in Figure 1), but still fail to  
 085 generate it with an explicit description. Revealing and quantifying the gap is crucial, as it not  
 086 only uncovers the inherent limitations of current unified models but also verifies whether they can  
 087 deploy their knowledge and reasoning abilities equally across understanding and generation tasks.
- 088 • **For evaluation protocols**, prevalent benchmarks always propose some challenging tasks such  
 089 as idiom interpretation to investigate unified models, but they cannot decompose these designed  
 090 reasoning-centric tasks into more evaluable and interpretable forms. Instead, they rely heavily  
 091 on the MLLM-as-a-Judge paradigm, where MLLMs serve directly as evaluators and evaluation  
 092 scores are obtained through visual question answering. Such dependence inevitably couples the  
 093 evaluation result with the biases and limitations of multimodal models themselves.

094 To this end, we introduce **GIR-Bench**, a benchmark designed to systematically explore the capabili-  
 095 ty boundaries of unified models in reasoning-driven image generation and editing. The benchmark  
 096 consists of three components: 1) **GIR-Bench-Uni**: We first systematically evaluate the gap between  
 097 understanding and generation for the same real-world entity. As shown in Figure 1(a), we collect  
 098 300 real-world entities across zoology, botany, and geography, design reasoning-oriented implicit  
 099 prompts, and pair them with curated reference images to construct paired tests. This design enables  
 100 us to quantify the gap between understanding and generating the same entity. 2) **GIR-Bench-T2I**:  
 101 We then explore reasoning-centric text-to-image generation (*c.f.*, Figure 1(b)). It requires models not  
 102 only to retrieve relevant knowledge but also to apply precise logical reasoning to faithfully satisfy  
 103 specified constraints. We design 300 carefully crafted prompts spanning three dimensions: numeri-  
 104 cal reasoning, spatial layout, and text rendering. 3) **GIR-Bench-Edit**: We further evaluate whether  
 105 models can perform global planning and reasoning-driven local modifications (*c.f.*, Figure 1(c)).  
 106 We construct 370 cases spanning visual puzzles, visual logic, and reasoning perception. Each case  
 107 consists of an input image and its associated ground-truth image, *i.e.*, reducing evaluation bias.

108 In terms of evaluation, unlike previous works that heavily rely on the MLLM-as-a-Judge paradigm,  
 109 we design a series of task-specific evaluation pipelines tailored for each task. Our results demon-  
 110 strate that these specialized evaluations not only provide fine-grained and interpretable assessments



Figure 2: Examples of leading models on the GIR-Bench. Designed complex and various tasks pose challenges to current models.

but also effectively mitigate the inherent biases of large multimodal models. Within GIR-Bench, we systematically evaluate 21 state-of-the-art models. Our challenging benchmark reveals the limitations of leading models in performing generation tasks that require complex visual reasoning, as shown in Figure 2. Massive results show that integrating understanding with generation enables models to perform more complex visual tasks. However, for unified models, a significant gap remains between understanding and generation. Thus, effective integration of the two is crucial for unlocking the potential of unified multimodal models. In summary, our contributions are threefold:

- We propose **GIR-Bench**, a comprehensive reasoning-centric benchmark that evaluates unified multimodal models across three perspectives.
- We design a diverse suite of tasks with task-specific evaluation pipelines that provide fine-grained and interpretable metrics, moving beyond the limitations of the MLLM-as-a-Judge paradigm.
- Through extensive experiments, we reveal both the performance gap between unified and generation-only models and the internal gap between understanding and generation.

## 2 GIR-BENCH

In this section, we introduce the main components of GIR-Bench. Section 2.1 details the data sources, dataset construction, prompt suites and evaluation pipeline of GIR-Bench-Uni. In Section 2.2, we present the task dimensions and evaluation suite of GIR-Bench-T2I. Section 2.3 elaborates on the evaluation dimensions, image sources and evaluation metrics of GIR-Bench-Edit. Finally, the experiments and the insights derived from GIR-Bench will be discussed in Section 3.

Our design of GIR-Bench is guided by three concrete principles intended to keep the benchmark objective, interpretable, and reproducible:

- **Objectivity over Subjectivity:** While many advanced reasoning scenarios (e.g., causal reasoning or open-ended commonsense) typically rely on the “MLLM-as-a-Judge” paradigm, this approach often introduces bias. We prioritized tasks with deterministic solutions (e.g., the unique solution in Sudoku, the original image in Jigsaw Puzzles, or exact answers in arithmetic). This ensures that our evaluation results are reproducible and indisputable.

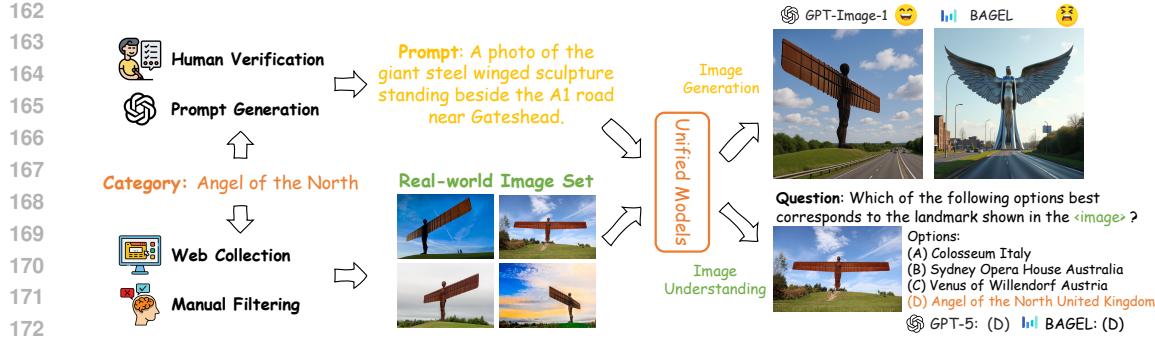


Figure 3: Illustration of GIR-Bench-Uni. For each real-world entity, an implicit prompt drives text-to-image generation, while the corresponding real image is used for image understanding evaluation.

- **Availability of Ground Truth:** We exclusively selected tasks where ground truth can be programmatically generated or strictly verified. For instance, `Spatial Layout` is verified via bounding box coordinates, and `Text Rendering` is checked via OCR. This quantifiable nature is a prerequisite for building a rigorous benchmark.
- **Focus on Implicit Reasoning & Planning:** We explicitly excluded tasks solvable by simple “keyword-to-image” mappings. The selected tasks (e.g., `Visual Puzzle` and `Numerical Reasoning`) compel the model to perform implicit reasoning or global planning before generating visual content that satisfies the constraints. This is key to measuring whether a model truly comprehends the logical constraints within a prompt.

## 2.1 GIR-BENCH-UNI

We collected 300 real-world entities from the Internet and open-source datasets, covering domains such as zoology, botany, and geography. Although previous works have explored the dimension of world knowledge, they typically remain limited to textual interpretation or text-image alignment without a deeper investigation. To evaluate reasoning ability, we utilize GPT-4o to generate implicit descriptions of each entity category, incorporating representative features such as visual appearance, historical context, and distinctive attributes. These prompts are then manually verified and refined to ensure that each uniquely corresponds to the real entity. Meanwhile, we curate a set of high-quality images to serve as reference exemplars for each entity. With entity categories, prompts, and reference image sets in place, we construct paired evaluation sets for both image understanding and image generation as shown in Figure 3. We compute the average DINOv3 (Siméoni et al., 2025) feature similarity between the generated image and the reference image sets as the evaluation metric. For image understanding, reference images are used to formulate visual question-answering cases. This design enables us to systematically investigate whether unified models can leverage the same underlying knowledge and reasoning for both understanding and generating the same entity.

## 2.2 GIR-BENCH-T2I

GIR-Bench-T2I consists of three text-to-image generation tasks that require deep reasoning. We design 300 carefully crafted prompts and construct evaluation pipelines by object detection and text recognition models. Although such problems are trivial for multimodal models, existing image generation models often fail to reason correctly and generate the expected images.

**Numerical Reasoning.** We design prompts that explicitly state mathematical constraints and record the corresponding ground truth, *e.g.*, the chicken-rabbit problem, while models must reason over these constraints to infer the correct objects and their quantities. For example, given the prompt “*A photo of ducks and dogs. There are a total of 10 legs visible, and the total number of animals is 4. Display all the animals clearly.*”, the expected output is 3 ducks and 1 dog. For evaluation, we apply object detection to extract the categories and counts of objects in the generated image and compare them against the ground truth. Notably, a case is counted as correct only when all object counts extracted from the generated image match the expected numbers. This strict criterion is necessary since the quantities of objects specified in the prompts are logically interdependent.

216 Partial correctness indicates a broken reasoning chain, whereas full correctness ensures that the  
 217 model has followed the intended reasoning process.  
 218

219 **Spatial Layout.** We design prompts that specify how objects should be arranged according to  
 220 constraints (e.g., categorical rules or ordered attributes) and record the corresponding ground truth  
 221 layouts. For example, the prompt “*A picture showing a bicycle, a dog, a car, and a cat. Arrange  
 222 the items in a horizontal line, with animals on the left and vehicles on the right.*” requires that the  
 223 bounding boxes of animals appear to the left of those of vehicles. We evaluate results by extract-  
 224 ing bounding boxes from generated images via object detection and verifying whether the spatial  
 225 relations conform to the specified constraints.  
 226

227 **Text Rendering.** While existing models already demonstrate strong text rendering capabilities,  
 228 their ability to reason over textual knowledge remains underexplored. To this end, we collect 60  
 229 short quotes or slogans and design implicit descriptions that correspond to them. For example, the  
 230 target text “*Just do it*” corresponds to the prompt “*A photo of an original 1988 Nike billboard featur-  
 231 ing the brand’s iconic three-word call-to-action slogan in bold capital letters.*” For evaluation, we  
 232 first extract text from generated images. Unlike traditional settings, our prompts are deliberately im-  
 233 plicit, which often leads models to generate the target text along with additional irrelevant content.  
 234 Common metrics (e.g., word accuracy and edit distance) are not suitable, since they wrongly penal-  
 235 ize extra words. Our goal is instead to assess whether the model successfully generates the ground  
 236 truth text while allowing the presence of additional content. To this end, we propose the word-level  
 237 continuous substring score as the primary evaluation metric. It is defined as  $s_{wc}(g, p) = \frac{|\mathcal{W}_{\text{match}}(g, p)|}{|\mathcal{W}(g)|}$ ,  
 238 where  $\mathcal{W}(g)$  denotes the set of words in the ground truth  $g$ , and  $\mathcal{W}_{\text{match}}(g, p)$  counts the number of  
 239 ground truth words that are fully covered by continuous character spans in the predicted text  $p$ .  
 240

### 2.3 GIR-BENCH-EDIT

242 GIR-Bench-Edit evaluates the image editing capabilities of models along three novel dimensions.  
 243 Unlike previous works, each editing case in our benchmark includes both an input image and a  
 244 corresponding ground truth image, thereby mitigating bias in evaluation.  
 245

246 **Visual Puzzle.** We filter real-world images collected in Section 2.1 to retain near-square, high-  
 247 resolution photos (minimum side length  $\geq 1024$  and aspect ratio  $< 1.2$ ), and resize each to a square.  
 248 The processed images are partitioned into grids, and the tiles are randomly permuted such that at  
 249 least half of the tile positions are altered. Given the shuffled image and the accompanying instruc-  
 250 tion, the model is required to reconstruct the original image, thereby evaluating its ability to integrate  
 251 semantic understanding with spatial reasoning in order to restore both the global structure and the  
 252 local coherence of natural images. For evaluation, we compute the Fréchet Inception Distance (FID)  
 253 between generated images and ground truth. To facilitate comparison with other metrics, we further  
 254 normalize FID to the range  $[0, 1]$ , where larger values indicate better reconstruction quality.  
 255

256 **Visual Logic.** We construct a high-quality dataset of Sudoku puzzles using a systematic generation  
 257 pipeline. For ground truth solutions, we employ a constraint-propagation algorithm that maintains  
 258 candidate sets for rows and columns, filling the grid iteratively with randomized choices while en-  
 259 suring the validity of standard Sudoku constraints. For input puzzles, we adopt a deductive removal  
 260 strategy, ensuring that each removed digit preserves the uniqueness of the solution. The puzzles and  
 261 their corresponding solutions are then rendered into images. For evaluation, we apply text detection  
 262 to extract the digits and their positions from generated images, and compute accuracy by comparing  
 263 predictions against the ground truth.  
 264

265 **Reasoning Perception.** We select high-quality images from the LISA (Lai et al., 2024) dataset and  
 266 filter those with small aspect ratios and resize them to a square. Editing instructions are constructed  
 267 from their implicit textual descriptions. The task requires models to edit the target regions into  
 268 green while keeping the background unchanged. Conceptually, this corresponds to segmentation,  
 269 but since directly predicting binary masks is out of scope for editing models, we use this proxy for-  
 270 mulation. Specifically, we instruct the model to segment the target region and render it in solid, fully  
 271 opaque green while leaving the rest of the image unchanged. For evaluation, we extract the edited  
 272 regions from the model outputs and convert them into binary masks. The procedure combines color-  
 273 threshold detection centered on the target green, enhanced, with channel-ratio based seed detection  
 274

270

271

Table 1: Results of various multimodal understanding and generation models on GIR-Bench-Uni.

272

273

| Type    | Model                  | Image Understanding |        |           |         | Image Generation |        |           |         |
|---------|------------------------|---------------------|--------|-----------|---------|------------------|--------|-----------|---------|
|         |                        | Zoology             | Botany | Geography | Overall | Zoology          | Botany | Geography | Overall |
| Und     | Qwen2.5VL-7B           | 0.943               | 0.990  | 1.000     | 0.978   | -                | -      | -         | -       |
|         | Qwen2.5-VL-32B         | 0.951               | 0.990  | 1.000     | 0.976   | -                | -      | -         | -       |
|         | GPT-5                  | 0.983               | 1.000  | 1.000     | 0.994   | -                | -      | -         | -       |
|         | Gemini-2.5-Flash       | 0.991               | 1.000  | 1.000     | 0.997   | -                | -      | -         | -       |
| Gen     | SD-3.5-Large           | -                   | -      | -         | -       | 0.263            | 0.163  | 0.437     | 0.288   |
|         | HiDream-II-Full        | -                   | -      | -         | -       | 0.298            | 0.218  | 0.617     | 0.378   |
|         | FLUX.1-schnell         | -                   | -      | -         | -       | 0.239            | 0.197  | 0.440     | 0.292   |
| Unified | Show-o2-7B             | 0.894               | 0.910  | 1.000     | 0.935   | 0.200            | 0.128  | 0.265     | 0.198   |
|         | Janus-Pro-7b           | 0.813               | 0.810  | 1.000     | 0.874   | 0.201            | 0.111  | 0.321     | 0.211   |
|         | BLIP3o-NEXT-SFT-3B     | 0.951               | 0.970  | 1.000     | 0.974   | 0.260            | 0.169  | 0.360     | 0.263   |
|         | Ovis-U1-3B             | 0.878               | 0.850  | 1.000     | 0.909   | 0.225            | 0.130  | 0.377     | 0.244   |
|         | OmniGen2               | 0.932               | 0.940  | 0.984     | 0.952   | 0.214            | 0.218  | 0.451     | 0.294   |
|         | UniPic2-Metaquery-9B   | -                   | -      | -         | -       | 0.269            | 0.195  | 0.440     | 0.301   |
|         | UniWorld-V1            | -                   | -      | -         | -       | 0.236            | 0.220  | 0.451     | 0.302   |
|         | BAGEL-7B               | 0.911               | 0.900  | 1.000     | 0.937   | 0.242            | 0.200  | 0.445     | 0.295   |
|         | BAGEL-7B w/ CoT        | 0.935               | 0.970  | 1.000     | 0.968   | 0.256            | 0.243  | 0.525     | 0.341   |
|         | Qwen-Image             | -                   | -      | -         | -       | 0.293            | 0.319  | 0.677     | 0.429   |
|         | Gemini-2.5-Flash-Image | -                   | -      | -         | -       | 0.448            | 0.565  | 0.772     | 0.595   |
|         | GPT-Image-1            | -                   | -      | -         | -       | 0.568            | 0.700  | 0.800     | 0.689   |

290

291

obtained by comparing the input and generated images. The resulting binary mask is then compared with the ground-truth mask using Intersection-over-Union (IoU) as the evaluation metric.

292

293

### 3 EXPERIMENTAL RESULTS AND INSIGHTS

294

#### 3.1 EXPERIMENT SETTINGS

295

**Evaluated Models.** We evaluate 21 representative models, including: **Multimodal understanding models:** Qwen2.5-VL-7B (Bai et al., 2025), Qwen2.5-VL-32B (Bai et al., 2025), GPT-5 (OpenAI, 2025a), Gemini-2.5-Flash (Comanici et al., 2025); **Image generation models:** SD-3.5-Large (Stability AI, 2024), HiDream-II-Full (Cai et al., 2025), FLUX.1-schnell (Labs, 2024); **Image editing models:** FLUX.1-Kontext (Batifol et al., 2025), ICEdit (Zhang et al., 2025), Step1X-Edit (Liu et al., 2025); **Unified multimodal models:** Uniworld (Lin et al., 2025), UniPic2 (Wei et al., 2025), BAGEL (Deng et al., 2025), Ovis-U1 (Wang et al., 2025a), OmniGen2 (Wu et al., 2025b), Show-o2 (Xie et al., 2025), Janus-Pro (Chen et al., 2025b), BLIP3o (Chen et al., 2025a), Qwen-Image (Wu et al., 2025a), GPT-Image-1 (OpenAI, 2025b), Gemini-2.5-Flash-Image (DeepMind, 2025).

296

297

**Implementation Details.** For object detection, we employ the grounding capability of InternVL3.5-38B (Wang et al., 2025b) to detect object categories and bounding boxes from generated images. For text detection, we use PPOCR v5 (Cui et al., 2025) to detect and recognize textual content, retaining only segments with confidence scores greater than 0.5.

298

#### 3.2 EVALUATION ON REAL-WORLD UNDERSTANDING AND GENERATION

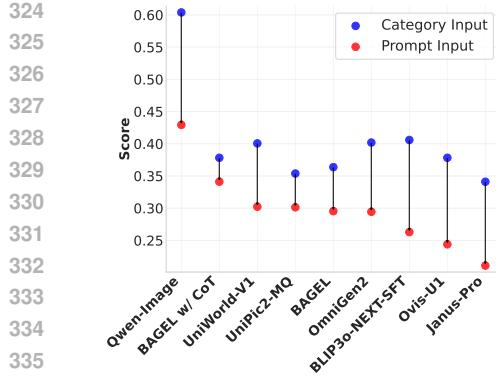
299

Table 1 reports results on GIR-Bench-Uni, which jointly evaluates real-world entity understanding and generation. The left panel shows understanding accuracy, while the right panel presents generation performance under implicit prompts. Overall, unified models outperform generation-only systems on reasoning-centric generation tasks, indicating that joint training across understanding and generation yields tangible benefits. GPT-Image-1 achieves the best performance and significantly outperforms other models, providing a strong upper bound. Unexpectedly, however, open-source unified models do not show a clear advantage in generation compared to strong generation-only models. To further probe this phenomenon, we conducted two complementary analyses:

300

301

**Knowledge and Reasoning Capacity.** We further probed whether generation failures stem from missing world knowledge or from difficulty in reasoning. To this end, we evaluated image understanding and reported the results in Table 1. Unified models exhibit consistently strong understanding: BLIP3o attains the highest overall score of 0.974, while the lowest, Janus-Pro, still achieves



324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377

Figure 4: Performance decline from category inputs to implicit prompts.

Prompt: A photo of the sculpture cast by Arturo Di Modica in 1989 as an impromptu gift symbolizing market resilience after the 1987 crash, now standing beside New York's primary stock exchange.

Category: Charging Bull United States



Figure 5: Qualitative cases in GIR-Bench-Uni, showing both direct category inputs and implicit prompt inputs.

Table 2: Results of various multimodal models on GIR-Bench-T2I and GIR-Bench-Edit.

| Type    | Model                  | GIR-Bench-T2I       |                |                |         | GIR-Bench-Edit |              |                      |         |
|---------|------------------------|---------------------|----------------|----------------|---------|----------------|--------------|----------------------|---------|
|         |                        | Numerical Reasoning | Spatial Layout | Text Rendering | Overall | Visual Puzzle  | Visual Logic | Reasoning Perception | Overall |
| Gen     | SD-3.5-Large           | 0.107               | 0.069          | 0.227          | 0.134   | -              | -            | -                    | -       |
|         | HiDream-II-Full        | 0.062               | 0.218          | 0.180          | 0.153   | -              | -            | -                    | -       |
|         | FLUX.1-schnell         | 0.045               | 0.195          | 0.238          | 0.159   | -              | -            | -                    | -       |
| Edit    | FLUX.1-Kontext-dev     | -                   | -              | -              | -       | 0.044          | 0.000        | 0.271                | 0.105   |
|         | ICEdit                 | -                   | -              | -              | -       | 0.023          | 0.030        | 0.233                | 0.095   |
|         | Step1X-Edit            | -                   | -              | -              | -       | 0.037          | 0.041        | 0.132                | 0.070   |
| Unified | UniWorld-V1            | 0.062               | 0.207          | 0.146          | 0.139   | 0.038          | 0.013        | 0.111                | 0.054   |
|         | UniPic2-Metaquery-9B   | 0.107               | 0.172          | 0.126          | 0.135   | 0.107          | 0.030        | 0.261                | 0.133   |
|         | Ovis-U1-3B             | 0.051               | 0.356          | 0.082          | 0.163   | 0.203          | 0.033        | 0.049                | 0.095   |
|         | OmniGen2               | 0.057               | 0.195          | 0.167          | 0.140   | 0.029          | 0.051        | 0.139                | 0.073   |
|         | Qwen-Image             | 0.153               | 0.368          | 0.150          | 0.224   | -              | -            | -                    | -       |
|         | Qwen-Image-Edit        | -                   | -              | -              | -       | 0.230          | 0.012        | 0.233                | 0.158   |
|         | BAGEL-7B               | 0.057               | 0.287          | 0.163          | 0.169   | 0.131          | 0.058        | 0.105                | 0.098   |
|         | BAGEL-7B w/ CoT        | 0.249               | 0.448          | 0.120          | 0.272   | 0.224          | 0.050        | 0.147                | 0.140   |
|         | Gemini-2.5-Flash-Image | 0.362               | 0.759          | 0.806          | 0.642   | 0.375          | 0.249        | 0.396                | 0.340   |
|         | GPT-Image-1            | 0.328               | 0.690          | 0.813          | 0.610   | 0.403          | 0.202        | 0.444                | 0.350   |

0.874. Such uniformly high scores suggest that entity recognition and basic reasoning are not the main bottlenecks. Moreover, with respect to implementation, it is worth noting that UniWorld and UniPic2 do not release official image understanding scripts. Both are developed upon a frozen Qwen2.5-VL-7B, and thus their understanding capability can be reasonably approximated by the performance of Qwen2.5-VL-7B. In a similar vein, Qwen-Image is built on Qwen2.5-VL (variant undisclosed); here we use Qwen2.5-VL-32B as a reference.

**Bridging Reasoning and Generation.** To disentangle the effect of reasoning, we compare two types of inputs: (i) a *category input*, which directly specifies the target entity using the template “*a photo of {category}*”, and (ii) the original *prompt input*, which requires reasoning to infer the entity. Figure 4 presents the per-model performance under both conditions. Across all models, scores with prompt input exhibit a marked reduction compared to category input, often by a significant margin. This pattern indicates that while models are capable of rendering entities when explicitly named, their performance deteriorates when the entity must be reasoned. The limitation therefore lies not in the ability to render the object itself, but in transferring reasoned constraints into the generative process. We further illustrate this gap through qualitative case studies in Section 4.

### 3.3 EVALUATION ON REASONING-CENTRIC TEXT-TO-IMAGE GENERATION

We report results on GIR-Bench-T2I in Table 2 (left). Overall, unified models outperform generation-only models, indicating that coupling understanding and generation yields measurable gains on reasoning-driven text-to-image generation. Nevertheless, even the strongest proprietary



Figure 6: Illustrative examples from GIR-Bench-Uni and GIR-Bench-T2I.

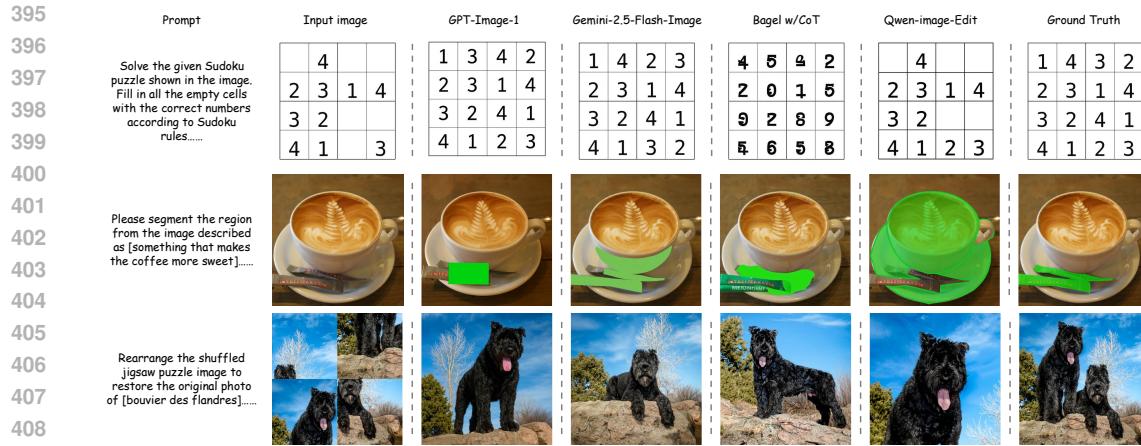


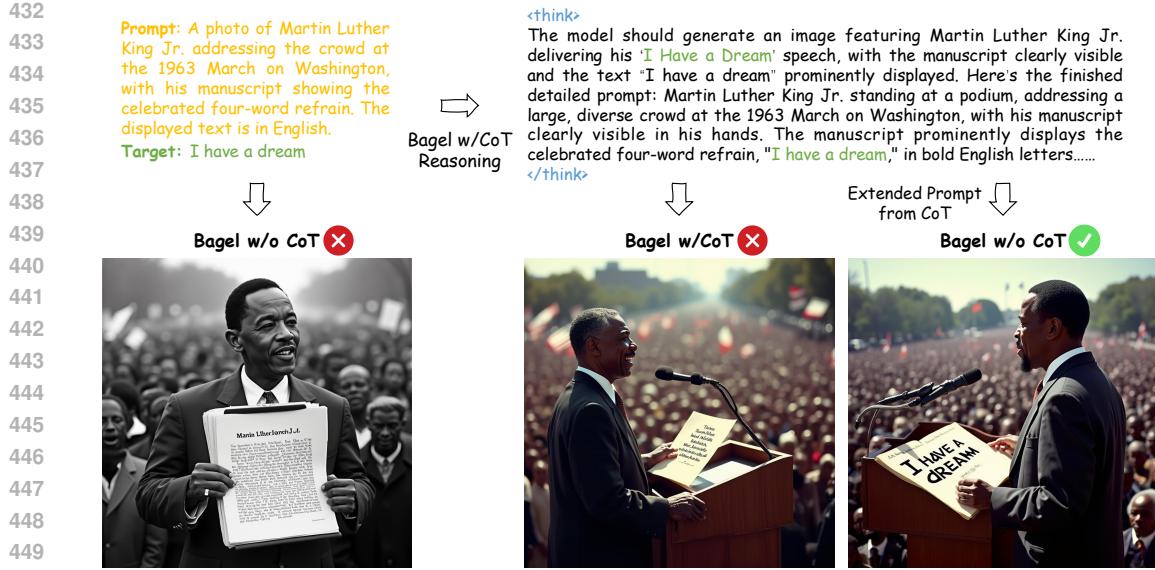
Figure 7: Representative outputs from GIR-Bench-Edit.

models (*i.e.*, GPT-Image-1 and Gemini-2.5-Flash-Image) are far from perfect, highlighting the limitations of current unified models in visual generation with reasoning.

Specifically, for numerical reasoning, the highest score is just 0.362 with Gemini-2.5-Flash-Image, underscoring that models are largely unable to generate correct quantities when reasoning over implicit prompts. Among open-source unified models, BAGEL w/ CoT exhibits a substantial gain over its base counterpart (0.249 vs. 0.057), suggesting that explicit chain-of-thought helps transfer arithmetic constraints into the generative process. In spatial Layout, proprietary models again dominate, with Gemini-2.5-Flash-Image at 0.759. Within open-source unified models, BAGEL w/ CoT (0.448) outperforms Qwen-Image (0.368) and Ovis-U1 (0.356), indicating that explicit planning improves global arrangement beyond local object placement. For text rendering, under implicit descriptions, the performance gap between proprietary and open-source models is the most pronounced: GPT-Image-1 (0.813) and Gemini-2.5-Flash-Image (0.806) substantially outperform all others. While many models can render text reliably with explicit input text, they struggle to produce correct outputs when the target text must first be reasoned and then generated. Notably, BAGEL w/CoT does not improve over BAGEL (0.120 vs. 0.163), suggesting that current reasoning traces are not yet effectively grounded in the generation process, thereby reflecting the broader gap between reasoning and generation. We further provide qualitative analysis in Section 4.

### 3.4 EVALUATION ON REASONING-CENTRIC IMAGE EDITING

Table 2 (right) presents results on GIR-Bench-Edit, which covers three editing tasks: visual puzzle, visual logic, and reasoning perception. In contrast to the trends observed in Section 2.1 and



486 **Qualitative Comparison with MLLM-as-a-Judge.** Figure 9 contrasts MLLM-as-a-Judge with our  
 487 explicit, task-specific metrics. MLLM-as-a-Judge sometimes assigns high scores based on superfi-  
 488 cial resemblance or biased reasoning, whereas our metrics correctly penalize mismatches with the  
 489 ground truth. For example, in the Lotus Temple case, judge-based scoring mislabels the Qwen-  
 490 Image output as correct despite clear structural inaccuracies. In the pens-and-pencils case, it fails  
 491 to reward Gemini-2.5-Flash-Image for producing the exact target. These cases highlight that our  
 492 metrics yield more consistent and interpretable evaluations.

493 **Misalignment Understanding-Generation** Based on our quantitative ablations and qualitative  
 494 case studies presented, we identify three primary factors driving the misalignment between the un-  
 495 derstanding and generation capabilities of unified models:

- 497 • **Asymmetry of Reasoning Capabilities:** Our experiments reveal a fundamental asymmetry:  
 498 models often possess the knowledge but fail to “activate” it during the generation process. As  
 499 illustrated in our case studies and experiments (in Figure 4 and Figure 5), the model successfully  
 500 generates the correct entity when explicitly named (*e.g.*, “Charging Bull United States”) but fails  
 501 when the same entity must be reasoned from an implicit prompt. This suggests that the reasoning  
 502 capability is localized within the LLM component. While the model correctly solves the inter-  
 503 mediate logic (*e.g.*, identifying the entity), the image generator—which is optimized for visual  
 504 fidelity rather than logical adherence—fails to receive or adhere to this reasoned state, leading to  
 505 generation failures.
- 506 • **Information Bottleneck in Heterogeneous Architectures:** Many unified models evaluated in  
 507 GIR-Bench (*e.g.*, BAGEL, Qwen-Image) employ heterogeneous architectures that couple a ca-  
 508 pable LLM with a separate generation head via a lightweight interface. This design creates a  
 509 significant information bottleneck. While the LLM performs deep reasoning, the interface often  
 510 compresses this rich semantic state into a limited number of condition tokens. Consequently, fine-  
 511 grained logical constraints (*e.g.*, exact counts in Numerical Reasoning) are often diluted or lost  
 512 before reaching the pixel generation stage.
- 513 • **Lack of Process-Oriented Training Data:** Current multimodal pre-training relies heavily on  
 514 static image-text pairs. These datasets map a final text description directly to an image but fail  
 515 to capture the *process* of generation. Existing models lack exposure to interleaved “reasoning-  
 516 trace” data (*e.g.*, [Reasoning Step] → [Intermediate Visual State] →  
 517 [Refined Image]). Without such data, the model struggles to learn how to decompose a  
 518 complex abstract instruction into a sequential plan for visual execution.

## 5 CONCLUSION

521 In this work, we presented GIR-Bench, a reasoning-centric benchmark that systematically evaluates  
 522 unified multimodal models across understanding, generation, and editing. By grounding evaluation  
 523 in explicit, task-specific metrics rather than the MLLM-as-Judge paradigm, GIR-Bench exposes  
 524 fine-grained weaknesses that are otherwise hidden behind holistic scores. Extensive experiments and  
 525 analyses show that while unified models consistently surpass generation-only models on reasoning-  
 526 intensive tasks, they still struggle with reliably transferring reasoning into faithful visual outputs.  
 527 These findings point to the importance of advancing unified models capable of seamlessly integrat-  
 528 ing reasoning and generation.

## 530 REFERENCES

531 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang,  
 532 Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan,  
 533 Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng,  
 534 Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report. *arXiv*  
 535 *preprint arXiv:2502.13923*, 2025.

536 Stephen Batifol, Andreas Blattmann, Frederic Boesel, Saksham Consul, Cyril Diagne, Tim Dock-  
 537 horn, Jack English, Zion English, Patrick Esser, Sumith Kulal, et al. Flux. 1 kontext: Flow match-  
 538 ing for in-context image generation and editing in latent space. *arXiv e-prints*, pp. arXiv–2506,  
 539 2025.

540 Qi Cai, Jingwen Chen, Yang Chen, Yehao Li, Fuchen Long, Yingwei Pan, Zhaofan Qiu, Yiheng  
 541 Zhang, Fengbin Gao, Peihan Xu, et al. Hidream-i1: A high-efficient image generative foundation  
 542 model with sparse diffusion transformer. *arXiv preprint arXiv:2505.22705*, 2025.

543 Juhai Chen, Zhiyang Xu, Xichen Pan, Yushi Hu, Can Qin, Tom Goldstein, Lifu Huang, Tianyi  
 544 Zhou, Saining Xie, Silvio Savarese, et al. Blip3-o: A family of fully open unified multimodal  
 545 models-architecture, training and dataset. *arXiv preprint arXiv:2505.09568*, 2025a.

546 Xiaokang Chen, Zhiyu Wu, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai Yu, and  
 547 Chong Ruan. Janus-pro: Unified multimodal understanding and generation with data and model  
 548 scaling. *arXiv preprint arXiv:2501.17811*, 2025b.

549 G. Comanici et al. Gemini 2.5: Pushing the frontier with advanced reasoning. *arXiv preprint*, 2025.  
 550 URL <https://arxiv.org/abs/2507.06261>.

551 Cheng Cui, Ting Sun, Manhui Lin, Tingquan Gao, Yubo Zhang, Jiaxuan Liu, Xueqing Wang, Zelun  
 552 Zhang, Changda Zhou, Hongen Liu, Yue Zhang, Wenyu Lv, Kui Huang, Yichao Zhang, Jing  
 553 Zhang, Jun Zhang, Yi Liu, Dianhai Yu, and Yanjun Ma. Paddleocr 3.0 technical report, 2025.  
 554 URL <https://arxiv.org/abs/2507.05595>.

555 Google DeepMind, 2025. URL <https://developers.googleblog.com/en/introducing-gemini-2-5-flash-image/>.

556 Chaorui Deng, Deyao Zhu, Kunchang Li, Chenhui Gou, Feng Li, Zeyu Wang, Shu Zhong, Wei-  
 557 hao Yu, Xiaonan Nie, Ziang Song, Guang Shi, and Haoqi Fan. Emerging properties in unified  
 558 multimodal pretraining. *arXiv preprint arXiv:2505.14683*, 2025.

559 Dhruba Ghosh, Hannaneh Hajishirzi, and Ludwig Schmidt. Geneval: An object-focused framework  
 560 for evaluating text-to-image alignment. *Advances in Neural Information Processing Systems*, 36:  
 561 52132–52152, 2023.

562 Kaiyi Huang, Kaiyue Sun, Enze Xie, Zhenguo Li, and Xihui Liu. T2i-compbench: A compre-  
 563 hensive benchmark for open-world compositional text-to-image generation. *Advances in Neural  
 564 Information Processing Systems*, 36:78723–78747, 2023.

565 Black Forest Labs. Flux. <https://github.com/black-forest-labs/flux>, 2024.

566 Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. Lisa: Rea-  
 567 soning segmentation via large language model. In *Proceedings of the IEEE/CVF Conference on  
 568 Computer Vision and Pattern Recognition (CVPR)*, pp. 9579–9589, June 2024.

569 Bin Lin, Zongjian Li, Xinhua Cheng, Yuwei Niu, Yang Ye, Xianyi He, Shenghai Yuan, Wangbo Yu,  
 570 Shaodong Wang, Yunyang Ge, et al. Uniworld: High-resolution semantic encoders for unified  
 571 visual understanding and generation. *arXiv preprint arXiv:2506.03147*, 2025.

572 Shiyu Liu, Yucheng Han, Peng Xing, Fukun Yin, Rui Wang, Wei Cheng, Jiaqi Liao, Yingming  
 573 Wang, Honghao Fu, Chunrui Han, Guopeng Li, Yuang Peng, Quan Sun, Jingwei Wu, Yan Cai,  
 574 Zheng Ge, Ranchen Ming, Lei Xia, Xianfang Zeng, Yibo Zhu, Binxing Jiao, Xiangyu Zhang,  
 575 Gang Yu, and Dixin Jiang. Step1x-edit: A practical framework for general image editing. *arXiv  
 576 preprint arXiv:2504.17761*, 2025.

577 Yiyang Ma, Xingchao Liu, Xiaokang Chen, Wen Liu, Chengyue Wu, Zhiyu Wu, Zizheng Pan,  
 578 Zhenda Xie, Haowei Zhang, Xingkai Yu, Liang Zhao, Yisong Wang, Jiaying Liu, and Chong  
 579 Ruan. Janusflow: Harmonizing autoregression and rectified flow for unified multimodal under-  
 580 standing and generation, 2024.

581 Yuwei Niu, Munan Ning, Mengren Zheng, Weiyang Jin, Bin Lin, Peng Jin, Jiaqi Liao, Kunpeng  
 582 Ning, Chaoran Feng, Bin Zhu, and Li Yuan. Wise: A world knowledge-informed semantic eval-  
 583 uation for text-to-image generation. *arXiv preprint arXiv:2503.07265*, 2025.

584 OpenAI. Gpt-5 system card. Technical report, OpenAI, 2025a. URL <https://cdn.openai.com/gpt-5-system-card.pdf>.

585 OpenAI, 2025b. URL <https://openai.com/index/image-generation-api/>.

594 Xichen Pan, Satya Narayan Shukla, Aashu Singh, Zhuokai Zhao, Shlok Kumar Mishra, Jialiang  
 595 Wang, Zhiyang Xu, Juhai Chen, Kunpeng Li, Felix Juefei-Xu, Ji Hou, and Saining Xie. Transfer  
 596 between modalities with metaqueries. *arXiv preprint arXiv:2504.06256*, 2025.

597

598 Oriane Siméoni, Huy V. Vo, Maximilian Seitzer, Federico Baldassarre, Maxime Oquab, Cijo Jose,  
 599 Vasil Khalidov, Marc Szafraniec, Seungeun Yi, Michaël Ramamonjisoa, Francisco Massa, Daniel  
 600 Haziza, Luca Wehrstedt, Jianyuan Wang, Timothée Darct, Théo Moutakanni, Leonel Sentana,  
 601 Claire Roberts, Andrea Vedaldi, Jamie Tolan, John Brandt, Camille Couprie, Julien Mairal, Hervé  
 602 Jégou, Patrick Labatut, and Piotr Bojanowski. DINOv3, 2025. URL <https://arxiv.org/abs/2508.10104>.

603

604 Stability AI. <https://stability.ai/news/introducing-stable-diffusion-3-5>,  
 605 2024.

606

607 Kaiyue Sun, Rongyao Fang, Chengqi Duan, Xian Liu, and Xihui Liu. T2i-reasonbench: Bench-  
 608 marking reasoning-informed text-to-image generation. *arXiv preprint arXiv:2508.17472*, 2025.

609

610 Chameleon Team. Chameleon: Mixed-modal early-fusion foundation models. *arXiv preprint  
 arXiv:2405.09818*, 2024.

611

612 Guo-Hua Wang, Shanshan Zhao, Xinjie Zhang, Liangfu Cao, Pengxin Zhan, Lunhao Duan, Shiyin  
 613 Lu, Minghao Fu, Jianshan Zhao, Yang Li, and Qing-Guo Chen. Ovis-u1 technical report. *arXiv  
 preprint arXiv:2506.23044*, 2025a.

614

615 Weiyun Wang, Zhangwei Gao, Lixin Gu, Hengjun Pu, Long Cui, Xingguang Wei, Zhaoyang Liu,  
 616 Linglin Jing, Shenglong Ye, Jie Shao, et al. Internvl3.5: Advancing open-source multimodal  
 617 models in versatility, reasoning, and efficiency. *arXiv preprint arXiv:2508.18265*, 2025b.

618

619 Xinlong Wang, Xiaosong Zhang, Zhengxiong Luo, Quan Sun, Yufeng Cui, Jinsheng Wang, Fan  
 620 Zhang, Yueze Wang, Zhen Li, Qiying Yu, et al. Emu3: Next-token prediction is all you need.  
 621 *arXiv preprint arXiv:2409.18869*, 2024.

622

623 Hongyang Wei, Baixin Xu, Hongbo Liu, Cyrus Wu, Jie Liu, Yi Peng, Peiyu Wang, Zexiang Liu,  
 624 Jingwen He, Yidan Xietian, Chuanxin Tang, Zidong Wang, Yichen Wei, Liang Hu, Boyi Jiang,  
 625 William Li, Ying He, Yang Liu, Xuchen Song, Eric Li, and Yahui Zhou. Skywork unipic 2.0:  
 626 Building kontext model with online rl for unified multimodal model, 2025. URL <https://arxiv.org/abs/2509.04548>.

627

628 Chenfei Wu, Jiahao Li, Jingren Zhou, Junyang Lin, Kaiyuan Gao, Kun Yan, Sheng ming Yin, Shuai  
 629 Bai, Xiao Xu, Yilei Chen, Yuxiang Chen, Zecheng Tang, Zekai Zhang, Zhengyi Wang, An Yang,  
 630 Bowen Yu, Chen Cheng, Dayiheng Liu, Deqing Li, Hang Zhang, Hao Meng, Hu Wei, Jingyuan  
 631 Ni, Kai Chen, Kuan Cao, Liang Peng, Lin Qu, Minggang Wu, Peng Wang, Shuting Yu, Tingkun  
 632 Wen, Wensen Feng, Xiaoxiao Xu, Yi Wang, Yichang Zhang, Yongqiang Zhu, Yujia Wu, Yuxuan  
 633 Cai, and Zenan Liu. Qwen-image technical report, 2025a. URL <https://arxiv.org/abs/2508.02324>.

634

635 Chengyue Wu, Xiaokang Chen, Zhiyu Wu, Yiyang Ma, Xingchao Liu, Zizheng Pan, Wen Liu,  
 636 Zhenda Xie, Xingkai Yu, Chong Ruan, et al. Janus: Decoupling visual encoding for unified  
 637 multimodal understanding and generation. *arXiv preprint arXiv:2410.13848*, 2024a.

638

639 Chenyuan Wu, Pengfei Zheng, Ruiran Yan, Shitao Xiao, Xin Luo, Yueze Wang, Wanli Li, Xiyan  
 640 Jiang, Yixin Liu, Junjie Zhou, Ze Liu, Ziyi Xia, Chaofan Li, Haoge Deng, Jiahao Wang, Kun  
 641 Luo, Bo Zhang, Defu Lian, Xinlong Wang, Zhongyuan Wang, Tiejun Huang, and Zheng Liu.  
 642 Omnigen2: Exploration to advanced multimodal generation. *arXiv preprint arXiv:2506.18871*,  
 2025b.

643

644 Junfeng Wu, Yi Jiang, Chuofan Ma, Yuliang Liu, Hengshuang Zhao, Zehuan Yuan, Song Bai,  
 645 and Xiang Bai. Liquid: Language models are scalable multi-modal generators. *arXiv preprint  
 arXiv:2412.04332*, 2024b.

646

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

648 Yang Ye, Xianyi He, Zongjian Li, Bin Lin, Shenghai Yuan, Zhiyuan Yan, Bohan Hou, and Li Yuan.  
 649 Imgedit: A unified image editing dataset and benchmark. *arXiv preprint arXiv:2505.20275*, 2025.  
 650

651 Zechuan Zhang, Ji Xie, Yu Lu, Zongxin Yang, and Yi Yang. In-context edit: Enabling instructional  
 652 image editing with in-context generation in large scale diffusion transformer. *arXiv preprint*  
 653 *arXiv:2504.20690*, 2025.

654 Xiangyu Zhao, Peiyuan Zhang, Kexian Tang, Hao Li, Zicheng Zhang, Guangtao Zhai, Junchi  
 655 Yan, Hua Yang, Xue Yang, and Haodong Duan. Envisioning beyond the pixels: Benchmark-  
 656 ing reasoning-informed visual editing. *arXiv preprint arXiv:2504.02826*, 2025.  
 657

## 659 A ETHICS STATEMENT

661 This work introduces GIR-Bench, a benchmark designed to evaluate unified multimodal models on  
 662 reasoning-driven image generation and editing tasks. Our study does not involve human subjects,  
 663 personally identifiable information, or sensitive private data. All images used in the benchmark  
 664 are either collected from publicly available open-source datasets or generated through automated  
 665 pipelines, and we have carefully filtered the data to avoid inappropriate or harmful content.  
 666

667 We acknowledge that generative models, when misused, could produce misleading or harmful out-  
 668 puts. To mitigate such risks, GIR-Bench is constructed solely for research purposes, with a focus  
 669 on advancing transparent and interpretable evaluation of multimodal reasoning. We explicitly avoid  
 670 releasing prompts or data that could be exploited for malicious generation, and we emphasize that  
 671 our benchmark should not be deployed in downstream applications without appropriate safeguards.  
 672

673 The release of GIR-Bench, along with its evaluation protocols, aims to foster reproducibility, fair-  
 674 ness, and accountability in multimodal research. By replacing opaque “MLLM-as-a-Judge” eval-  
 675 uation with task-specific, interpretable metrics, we encourage ethical research practices and reduce  
 676 risks of hidden biases. Our study complies with the ICLR Code of Ethics, and we affirm our com-  
 677 mitment to responsible AI research.  
 678

## 679 B REPRODUCIBILITY STATEMENT

680 We have taken extensive measures to ensure the reproducibility of GIR-Bench. All dataset con-  
 681 struction procedures, including the collection of real-world entities, prompt generation, and editing  
 682 case design, are described in detail in Section 2.1, 2.2 and 2.3, with additional implementation de-  
 683 tails provided in the Appendix. Evaluation metrics and pipelines are explicitly defined for each  
 684 task, covering object detection, text recognition, and segmentation-based analysis, enabling oth-  
 685 others to replicate our results without relying on subjective judgments. Comprehensive experimental  
 686 results across 21 models are reported in Table 2 and Table 1, and we also provide qualitative ex-  
 687 amples in Figures 10–12 to illustrate typical outcomes. To further support reproducibility, we re-  
 688 lease all data, prompts, and evaluation code at [https://anonymous.4open.science/r/  
 689 GIR-Bench-7E40](https://anonymous.4open.science/r/GIR-Bench-7E40), ensuring that our benchmark can be independently verified and extended by  
 690 the research community.  
 691

## 692 C RELATED WORK

### 693 C.1 UNIFIED MULTIMODAL MODELS

695 Research in multimodal generation has shifted from modality-specific architectures to unified mod-  
 696 els that handle both understanding and generation across modalities (Lin et al., 2025; Wu et al.,  
 697 2024b;a; Ma et al., 2024; Wei et al., 2025; Wu et al., 2025b; Chen et al., 2025b; Wang et al., 2025a).  
 698 Chameleon (Team, 2024) pioneered an early-fusion, token-based transformer that can generate and  
 699 interpret text and images interchangeably, matching or surpassing proprietary systems. Emu3 (Wang  
 700 et al., 2024) shows that pure autoregressive modeling can achieve general multimodal intelligence  
 701 without relying on diffusion. MetaQueries (Pan et al., 2025) introduces learnable query tokens to  
 702 interface a frozen multimodal LLM with a diffusion decoder, enabling knowledge-augmented image

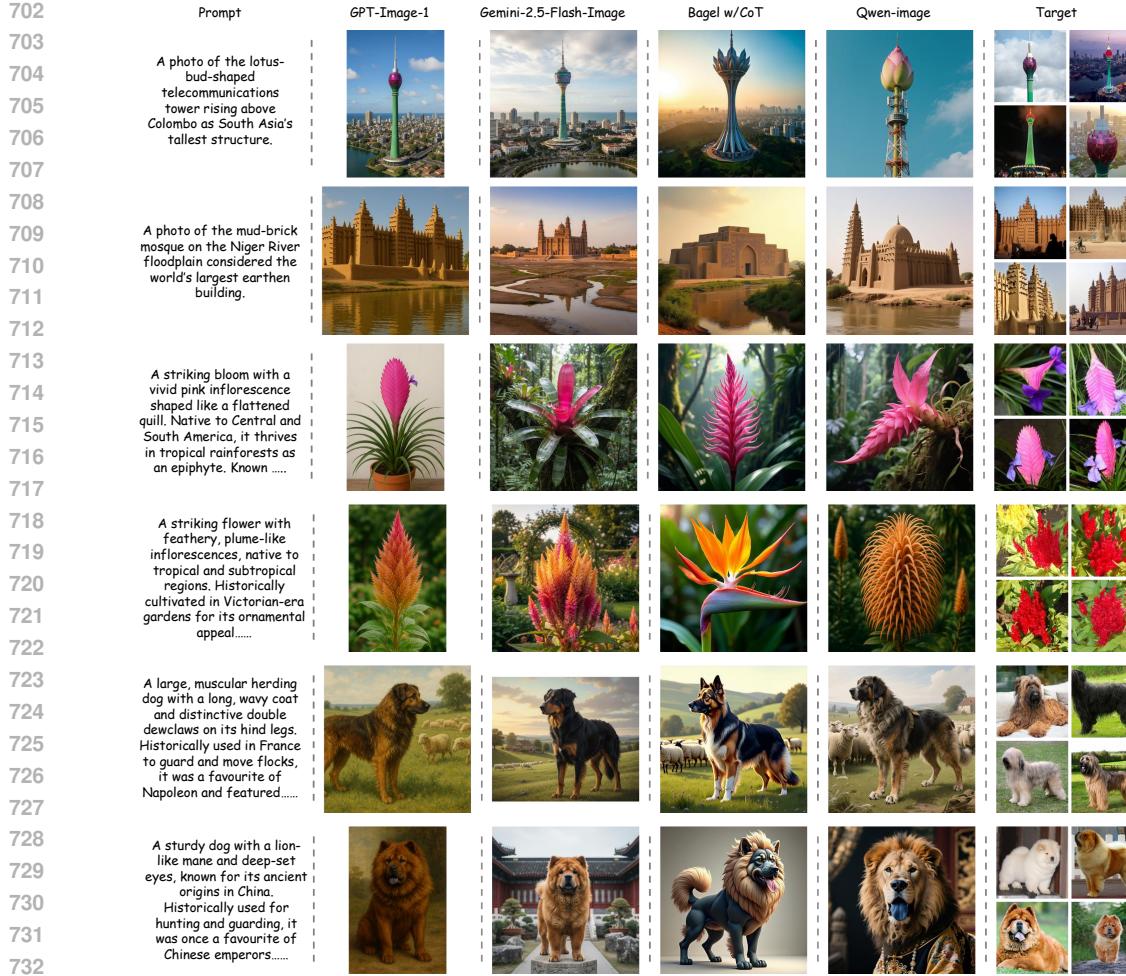


Figure 10: Examples from GIR-Bench-Uni.

generation without degrading understanding performance. Some works attempt to integrate diffusion into unified frameworks. BLIP3-o (Chen et al., 2025a) uses a diffusion transformer to generate semantically rich CLIP image features and employs sequential pretraining (understanding followed by generation). BAGEL (Deng et al., 2025) is a decoder-only foundation model pretrained on trillions of interleaved multimodal tokens, exhibiting emergent capabilities in complex multimodal reasoning such as free-form image manipulation. Proprietary systems (e.g., GPT-Image-1 (OpenAI, 2025b) and Gemini 2.5 Flash Image (DeepMind, 2025)) further showcase the powerful capabilities of unified multimodal models and have given rise to a wide range of real-world applications. These advancements underscore the need for comprehensive and reasoning-centric benchmarks.

## C.2 BENCHMARKS FOR UNIFIED MULTIMODAL MODELS

Existing benchmarks (Huang et al., 2023; Niu et al., 2025; Zhao et al., 2025) for unified multimodal models cover various aspects but still fall short of probing reasoning. GenEval (Ghosh et al., 2023) uses object detection to check whether generated images match explicit constraints on object co-occurrence, position, count and colour, limiting its scope to compositional alignment. ImgEdit (Ye et al., 2025) introduces a large-scale, high-quality dataset of 1.2 million edit pairs covering both single-turn and multi-turn image editing tasks, and proposes ImgEdit-Bench, a benchmark evaluating instruction adherence, editing quality and detail preservation; however, its focus remains on surface-level editing fidelity rather than deeper reasoning. WISE (Niu et al., 2025) provides 1,000 prompts spanning cultural commonsense, spatio-temporal reasoning and natural science, and measures knowledge–image alignment with its WiScore metric, but its evaluation depends on the

|     | Prompt   | GPT-Image-1 | Gemini-2.5-Flash-Image | Bagel w/CoT | Qwen-image | Target                          |
|-----|--|-------------|------------------------|-------------|------------|---------------------------------|
| 756 |  |             |                        |             |            |                                 |
| 757 | A photo of keys, lamps, and blankets. The number of keys is equal to twice the number of lamps, and the number of blankets is equal to the number of lamps. Altogether, there are five items. Display all the objects. |             |                        |             |            | key: 2<br>lamp: 1<br>blanket: 2 |
| 758 |  |             |                        |             |            |                                 |
| 759 |  |             |                        |             |            |                                 |
| 760 |  |             |                        |             |            |                                 |
| 761 |  |             |                        |             |            |                                 |
| 762 | A photo of cups, plates, and bottles. The number of cups is one more than the number of plates. The number of bottles equals the number of plates. In total, there are four items. Show all the objects.               |             |                        |             |            | cup: 2<br>plate: 1<br>bottle: 1 |
| 763 |  |             |                        |             |            |                                 |
| 764 |  |             |                        |             |            |                                 |
| 765 |  |             |                        |             |            |                                 |
| 766 |  |             |                        |             |            |                                 |
| 767 | A photo of an ancient Roman marble tablet in the Capitoline Museums bearing Caesar's succinct three-word report of victory in Latin. The displayed text is in English.   |             |                        |             |            | I came, I saw, I conquered      |
| 768 |  |             |                        |             |            |                                 |
| 769 |  |             |                        |             |            |                                 |
| 770 |  |             |                        |             |            |                                 |
| 771 |  |             |                        |             |            |                                 |
| 772 | A photo of Benjamin Franklin's Poor Richard Almanack opened to the page where punctuality is directly linked with profit, the sentence underlined. The displayed text is in English.                                   |             |                        |             |            | Time is money                   |
| 773 |  |             |                        |             |            |                                 |
| 774 |  |             |                        |             |            |                                 |
| 775 |  |             |                        |             |            |                                 |
| 776 |  |             |                        |             |            |                                 |
| 777 |  |             |                        |             |            |                                 |
| 778 |  |             |                        |             |            |                                 |
| 779 |  |             |                        |             |            |                                 |

Figure 11: Examples from GIR-Bench-T2I.

inherent knowledge of multimodal LLMs, making it prone to bias. T2I-ReasonBench (Sun et al., 2025) introduces 800 prompts across idiom interpretation, textual design, entity and scientific reasoning, but its two-stage evaluation (LLM-generated questions and LMM-based scoring) cannot disentangle understanding from generation. RISEBench (Zhao et al., 2025) extends evaluation to visual editing tasks involving temporal, causal, spatial and logical reasoning and scores instruction adherence, appearance consistency and visual plausibility using human and multimodal judges, yet it again depends on large models as judges rather than explicit, interpretable metrics. Current benchmarks emphasize explicit alignment and surface-level fidelity, lean heavily on multimodal models as judges, and lack a framework to decouple comprehension from generation.

## D IMAGE EDITING INSTRUCTIONS

To ensure consistency and clarity in evaluating editing capabilities, we design explicit instructions for the three tasks in GIR-Bench-Edit. Each instruction specifies the editing objective, the constraints to be preserved, and the expected output format. Placeholders such as `{caption}` and `{category}` are replaced with task-specific descriptions or categories for each case.

**Reasoning Perception:** Please segment the region from the image described as `{caption}`. Output the image with a solid, fully opaque green (#00FF00) segmentation mask overlay, and leave the rest of the image unchanged.

**Visual Logic:** Solve the given Sudoku puzzle shown in the image. Fill in all the empty cells with the correct numbers according to Sudoku rules. Keep the existing numbers and the grid lines unchanged. Write the solution numbers in a consistent font and size, aligned within each cell, so that the final image shows the completed Sudoku solution.

**Visual Puzzle:** Rearrange the shuffled jigsaw puzzle image to restore the original photo of `{caption}`.

## E MORE QUALITATIVE RESULTS

In this section, we provide representative qualitative results from the three components of GIR-Bench. Figure 10 shows examples from GIR-Bench-Uni, Figure 11 from GIR-Bench-T2I, and Fig-

| 810 | Prompt   | Input image | GPT-Image-1 | Gemini-2.5-Flash-Image | Bagel w/CoT | Qwen-image-Edit | Ground Truth |
|-----|--|-------------|-------------|------------------------|-------------|-----------------|--------------|
| 811 | Solve the given Sudoku puzzle shown in the image. Fill in all the empty cells with the correct numbers according to Sudoku rules.....  |             |             |                        |             |                 |              |
| 812 | Solve the given Sudoku puzzle shown in the image. Fill in all the empty cells with the correct numbers according to Sudoku rules.....  |             |             |                        |             |                 |              |
| 813 | Solve the given Sudoku puzzle shown in the image. Fill in all the empty cells with the correct numbers according to Sudoku rules.....  |             |             |                        |             |                 |              |
| 814 | Solve the given Sudoku puzzle shown in the image. Fill in all the empty cells with the correct numbers according to Sudoku rules.....  |             |             |                        |             |                 |              |
| 815 | Solve the given Sudoku puzzle shown in the image. Fill in all the empty cells with the correct numbers according to Sudoku rules.....  |             |             |                        |             |                 |              |
| 816 | Solve the given Sudoku puzzle shown in the image. Fill in all the empty cells with the correct numbers according to Sudoku rules.....  |             |             |                        |             |                 |              |
| 817 | Solve the given Sudoku puzzle shown in the image. Fill in all the empty cells with the correct numbers according to Sudoku rules.....  |             |             |                        |             |                 |              |
| 818 | Solve the given Sudoku puzzle shown in the image. Fill in all the empty cells with the correct numbers according to Sudoku rules.....  |             |             |                        |             |                 |              |
| 819 | Please segment the region from the image described as [if pigs are shot and killed, they usually fall to the ground. Which pigs in the picture are most likely to still be alive?.....]                                |             |             |                        |             |                 |              |
| 820 | Please segment the region from the image described as [What are some of the main features of the rainforest environment that support the growth and development of the diverse collection of plants found there?.....] |             |             |                        |             |                 |              |
| 821 | Please segment the region from the image described as [What are some of the main features of the rainforest environment that support the growth and development of the diverse collection of plants found there?.....] |             |             |                        |             |                 |              |
| 822 | Please segment the region from the image described as [What are some of the main features of the rainforest environment that support the growth and development of the diverse collection of plants found there?.....] |             |             |                        |             |                 |              |
| 823 | Rearrange the shuffled jigsaw puzzle image to restore the original photo of [scotch terrier].  |             |             |                        |             |                 |              |
| 824 | Rearrange the shuffled jigsaw puzzle image to restore the original photo of [scotch terrier].  |             |             |                        |             |                 |              |
| 825 | Rearrange the shuffled jigsaw puzzle image to restore the original photo of [scotch terrier].  |             |             |                        |             |                 |              |
| 826 | Rearrange the shuffled jigsaw puzzle image to restore the original photo of [scotch terrier].  |             |             |                        |             |                 |              |
| 827 | Rearrange the shuffled jigsaw puzzle image to restore the original photo of [scotch terrier].  |             |             |                        |             |                 |              |
| 828 | Rearrange the shuffled jigsaw puzzle image to restore the original photo of [scotch terrier].  |             |             |                        |             |                 |              |
| 829 | Rearrange the shuffled jigsaw puzzle image to restore the original photo of [scotch terrier].  |             |             |                        |             |                 |              |
| 830 | Rearrange the shuffled jigsaw puzzle image to restore the original photo of [scotch terrier].  |             |             |                        |             |                 |              |
| 831 | Rearrange the shuffled jigsaw puzzle image to restore the original photo of [scotch terrier].  |             |             |                        |             |                 |              |
| 832 | Rearrange the shuffled jigsaw puzzle image to restore the original photo of [scotch terrier].  |             |             |                        |             |                 |              |
| 833 | Rearrange the shuffled jigsaw puzzle image to restore the original photo of [scotch terrier].  |             |             |                        |             |                 |              |
| 834 | Rearrange the shuffled jigsaw puzzle image to restore the original photo of [scotch terrier].  |             |             |                        |             |                 |              |
| 835 | Rearrange the shuffled jigsaw puzzle image to restore the original photo of [scotch terrier].  |             |             |                        |             |                 |              |
| 836 | Rearrange the shuffled jigsaw puzzle image to restore the original photo of [scotch terrier].  |             |             |                        |             |                 |              |
| 837 | Rearrange the shuffled jigsaw puzzle image to restore the original photo of [scotch terrier].  |             |             |                        |             |                 |              |
| 838 | Rearrange the shuffled jigsaw puzzle image to restore the original photo of [scotch terrier].  |             |             |                        |             |                 |              |
| 839 | Rearrange the shuffled jigsaw puzzle image to restore the original photo of [scotch terrier].  |             |             |                        |             |                 |              |
| 840 | Rearrange the shuffled jigsaw puzzle image to restore the original photo of [scotch terrier].  |             |             |                        |             |                 |              |

Figure 12: Examples from GIR-Bench-Edit.

ure 12 from GIR-Bench-Edit. These cases illustrate typical challenges faced by current models and highlight the reasoning–generation misalignment observed across different tasks.

## F ROBUSTNESS AND VALIDATION OF AUTOMATED EVALUATION METRICS

To ensure the reliability of GIR-Bench, we conducted a comprehensive analysis of the robustness of our automated evaluation metrics compared to human judgment and MLLM-based judges. As illustrated in Figure 13, our automated metrics demonstrate a high Pearson correlation of  $\rho \approx 0.96$  with human annotation globally. In contrast, MLLM-as-a-Judge approaches exhibit lower consistency and a notable tendency towards score inflation, particularly in lower-performance regimes (visible as points floating significantly above the diagonal in the lower-left quadrant). This quantitative evidence supports our choice of using deterministic, verifiable metrics over stochastic MLLM judgments.

We specifically investigated edge cases where automated tools might misjudge correct outputs. Our observations indicate:

- **Rare False Negatives:** Mismatches are infrequent and typically occur only when objects are rendered in highly abstract, stylized forms, are heavily occluded, or when text is severely distorted beyond standard legibility. In such cases, human annotators often also find the content ambiguous.
- **Ranking Stability:** Crucially, these edge cases are uniformly distributed and do not disproportionately affect specific models. Consequently, the relative ranking of models remains consistent. The significant performance gap is robust to these minor metric fluctuations.

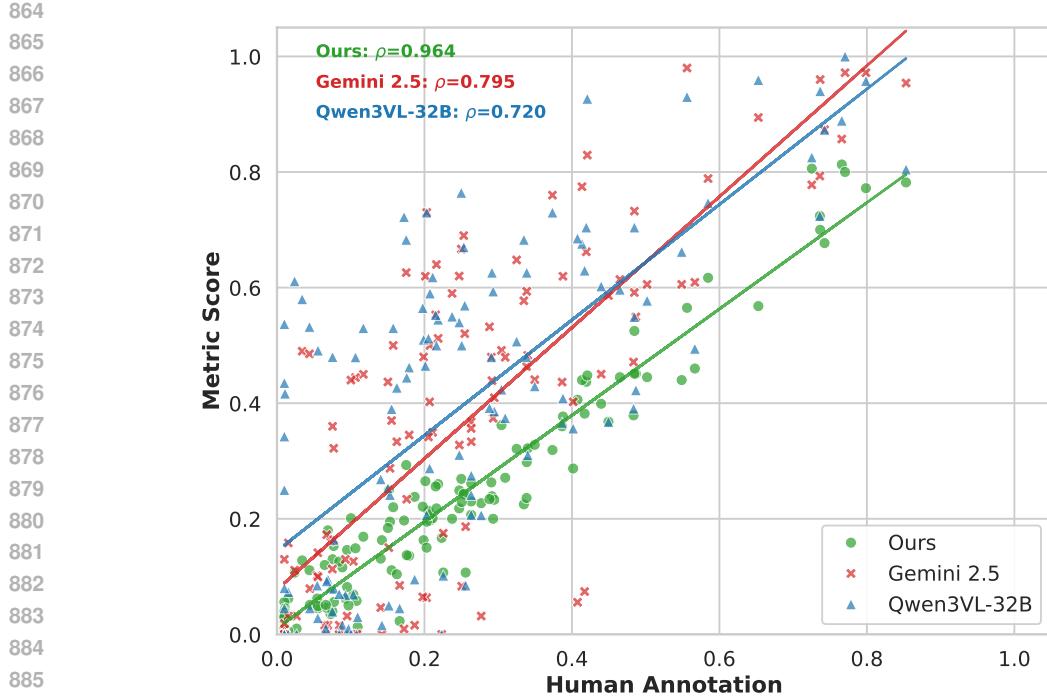


Figure 13: Correlation Analysis of Evaluation Protocols. We aggregated evaluation scores across all tasks in GIR-Bench to compare our metrics against MLLM-as-a-Judge methods. Our task-specific metrics (Green) show a strong linear correlation ( $\rho \approx 0.96$ ) with human judgment, demonstrating high robustness. In contrast, MLLM judges (Red/Blue) exhibit higher variance and inflation bias.

## G EXTENDED DISCUSSION: THE GAP BETWEEN UNDERSTANDING AND GENERATION

In this section, we provide an in-depth analysis of the performance gap observed between the understanding and generation capabilities of unified multimodal models. Based on our quantitative ablations in Figure 4) and qualitative case studies in Figure 5), we identify the root causes of this misalignment and propose potential research directions to bridge this gap.

### G.1 ROOT CAUSES OF THE PERFORMANCE GAP

We identify three primary factors driving the misalignment where models fail to transfer their reasoning capabilities into faithful visual generation:

- **Asymmetry of Reasoning Capabilities:** Our experiments reveal a fundamental asymmetry: models often possess the requisite knowledge but fail to “activate” it during the generation process. As illustrated in our case study on the “Charging Bull” (as shown in Figure 5), the model successfully generates the correct entity when explicitly named (*e.g.*, “Charging Bull United States”) but fails when the same entity must be reasoned from an implicit prompt. This suggests that the reasoning capability is primarily localized within the Large Language Model (LLM) component. While the model correctly solves the intermediate logic (*e.g.*, identifying the entity), the image generator—which is optimized for visual fidelity rather than logical adherence—fails to receive or adhere to this reasoned state, leading to generation failures.
- **Information Bottleneck in Heterogeneous Architectures:** Many unified models evaluated in GIR-Bench (*e.g.*, BAGEL (Deng et al., 2025), Qwen-Image (Wu et al., 2025a)) employ heterogeneous architectures that couple a highly capable autoregressive LLM with a separate generation head via a lightweight interface. This design creates a significant

918 information bottleneck. While the LLM performs deep reasoning, the interface often compresses this rich semantic state into a limited number of condition tokens or a concise embedding. Consequently, fine-grained logical constraints (e.g., exact counts in Numerical Reasoning tasks) are often diluted or lost before reaching the pixel generation stage.

- 919 • **Lack of Process-Oriented Training Data:** Current multimodal pre-training relies heavily  
920 on static image-text pairs. These datasets map a final text description directly to an image  
921 but fail to capture the *process* of generation. Existing models lack exposure to interleaved  
922 “reasoning-trace” data (e.g., [Reasoning Step] → [Intermediate Visual State] →  
923 [Refined Image]). Without such data, the model struggles to learn how to decompose a  
924 complex abstract instruction into a sequential execution plan for visual generation.

## 929 G.2 FUTURE DIRECTIONS

930 To narrow the gap between understanding and generation, we propose three concrete research directions supported by our empirical findings:

- 931 • **Explicit Chain-of-Thought (CoT):** Our experimental results with **BAGEL w/ CoT** (Ta-  
932 ble 2) demonstrate that externalizing the reasoning process significantly boosts per-  
933 formance. For instance, accuracy on the Numerical Reasoning task improves from **0.057**  
934 to **0.249**. This confirms that forcing the model to output an explicit textual plan acts as  
935 a crucial bridge, converting implicit logic into explicit constraints that the generator can  
936 better attend to and execute.
- 937 • **Native Unified Architectures:** We argue that future work should move towards natively  
938 unified paradigms where text and image tokens are treated equally within a single back-  
939 bone. This would effectively eliminate the heterogeneous interface bottleneck, allowing  
940 the visual generation process to attend directly to the full, uncompressed reasoning states  
941 of the MLLM.
- 942 • **Dataset Evolution:** We advocate for the construction of reasoning-trace datasets. Future  
943 benchmarks and training sets should provide not only the final ground-truth image but also  
944 the logical execution trace. This would enable models to be supervised directly on the  
945 reasoning-to-generation mapping, rather than relying solely on end-to-end alignment.

## 946 H CLARIFICATION FOR THE TEXT RENDERING METRIC ( $s_{wc}$ )

947 We illustrate the robustness of our proposed metric ( $s_{wc}$ ) compared to standard accuracy ( $s_{acc}$ ) using  
948 the Ground Truth: “**Make It Happen**”.

949  
950 Table 3: Comparison of metric behaviors under different generation scenarios. Ground Truth:  
951 “**Make It Happen**”.

| 952 Case            | 953 Generated Text (OCR Output)               | 954 $s_{wc}$ (Ours) | 955 $s_{acc}$      |
|---------------------|---|---------------------|--------------------|
| 956 A (Perfect)     | 957 “Make It Happen”                          | 958 <b>1.00</b>     | 959 1.00           |
| 960 B (Extra Words) | 961 “Poster says xxxx <b>Make It Happen</b> ” | 962 <b>1.00</b>     | 963 $\approx 0.00$ |
| 964 C (Merged)      | 965 “ <b>MakeItHappen</b> ”                   | 966 <b>1.00</b>     | 967 0.86           |
| 968 D (Missing)     | 969 “Make It”                                 | 970 <b>0.67</b>     | 971 0.50           |

972 The cases in Table 3 demonstrate the specific advantages of  $s_{wc}$ :

- 973 • **Case B (Robustness to Context):** Unified models often generate descriptive text or arti-  
974 facts (e.g., “Poster says xxxx”) alongside the target. This effectively handles the “chatty”  
975 nature of these models.  $s_{acc}$  fails due to the length penalty imposed by the extra characters,  
976 whereas  $s_{wc}$  correctly credits the presence of the valid content.
- 977 • **Case C (Robustness to Layout):** Tight visual layouts or artistic fonts often lead to merged  
978 OCR outputs.  $s_{wc}$  successfully identifies the legible words within the merged string, while  
979  $s_{acc}$  penalizes the missing whitespace.

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

- **Case D (Penalty for Missing):** Both metrics correctly penalize missing content. It is important to note that  $s_{wc}$  acts as a strict filter: partial matches (e.g., “Make It Hap”) are not credited, ensuring that the metric does not reward incomplete generations.