

VCode: A Multimodal Coding Benchmark with SVG as Symbolic Visual Representation

Anonymous CVPR submission

Paper ID 20

Abstract

001 *Code has emerged as a precise, executable medium for*
002 *linguistic-centric tasks, leaving visual-centric coding un-*
003 *derexplored. Conventional image representations rely on*
004 *RGB pixels that capture visual appearance but offer limited*
005 *symbolic abstraction. In this work, we advocate SVG code*
006 *as a compact, interpretable, and executable visual repre-*
007 *sentation. We introduce VCode, a benchmark that reframes*
008 *multimodal understanding as code generation: given an im-*
009 *age, a model must produce SVG that preserves symbolic*
010 *meaning for downstream reasoning. VCode covers general*
011 *commonsense, professional disciplines, and visual-centric*
012 *perception. To assess symbolic fidelity, we propose Code-*
013 *VQA, a novel evaluation protocol where a policy model an-*
014 *swers questions over rendered SVGs; correct answers in-*
015 *dicade faithful symbolic preservation. We also introduce*
016 *VCoder, an agentic framework that augments VLMs via*
017 *test-time revision and visual tool use, yielding substantial*
018 *improvements over strong baselines.*

019 1. Introduction

020 To advance reasoning and agentic intelligence, code has
021 emerged as a powerful medium for interacting with digi-
022 tal environments [12, 13, 23]. However, recent benchmarks
023 predominantly emphasize *linguistic-centric* coding abilities
024 (e.g. program synthesis and debugging) [2, 4, 10, 11, 21].
025 In the multi-modal regime, code is often leveraged to gen-
026 erate *synthetic* visual assets—such as charts [25, 28], dia-
027 grams [5, 15, 16], or websites [3, 17]—which are not di-
028 rectly grounded in the natural visual world.

029 When representing natural images, the dominant prac-
030 tice relies on dense pixels. In contrast, humans often per-
031 ceive and reason through sparse symbolic sketches that em-
032 phasize spatial relationships, object counts, and structural
033 outlines [8]. Building on this intuition, we propose using
034 Scalable Vector Graphics (SVG) code as an alternative vi-
035 sual representation, owing to its compact, interpretable, and

executable nature [6, 15, 26, 29]. This motivates a funda- 036
mental question: **can code serve as a naive visual repre-** 037
sentation of natural images? 038

In this work, we introduce VCode, a multimodal 039
coding benchmark that explores this research problem. 040
VCode is constructed by repurposing existing multimodal 041
understanding benchmarks across three domains: Gen- 042
eral commonsense (MM-Vet [31]), College-level disci- 043
plines (MMMU [32]), and Visual-centric Perception (CV- 044
Bench [22]). VCode reframes these tasks as visual coding: 045
given an image, a model must generate SVG code that faith- 046
fully reconstructs its symbolic representation. To evaluate 047
this transformation, we propose **CodeVQA**, a novel proto- 048
col in which a vision-language model (VLM) must answer 049
core questions about the original image by reasoning *exclu-* 050
sively over the rendered SVG. Experiments on VCode reveal 051
that while coding quality improves with reasoning ability, 052
existing coders still fail to preserve fine-grained visual rela- 053
tions (e.g. near vs. far). 054

To address this persistent gap, we augment existing 055
coders with two complementary capabilities. **(i) Thinking** 056
with Revision: the model compares intermediate render- 057
ings with the original image, explicitly articulates discrep- 058
ancies, and iteratively refines the SVG. **(ii) Acting with Vi-** 059
sual Tools: we equip the coder with external perception 060
toolboxes (e.g. detectors and segmenters [14, 27]) to sup- 061
ply structured cues (objects, shapes, text) as coding context. 062
Together, these strategies yield a +12.3 overall gain over 063
the top-performing Claude-4-Opus. Our contributions are 064
summarized as follows: 065

- **VCode & CodeVQA:** We recast multimodal understand- 066
ing as visual-centric coding. We introduce VCode and 067
the CodeVQA protocol, which tests whether the gener- 068
ated SVG code serves as an adequate and faithful visual 069
representation for downstream reasoning. 070
- **VCoder:** An agentic framework that augments VLMs 071
into strong multimodal coders via *Thinking with Revision* 072
and *Acting with Visual Tools*, achieving significant gains 073
over strong baselines. 074

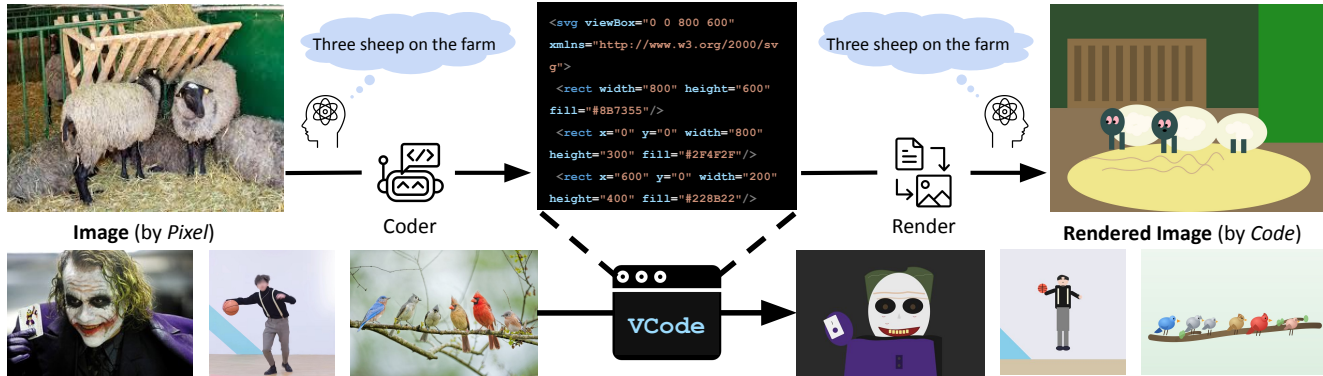


Figure 1. **Illustration of VCode.** An RGB image (left, represented by pixels) is translated into symbolic SVG code (middle) via VLM as Coder and rendered back into an image (right, represented by code), aiming to preserve symbolic meaning (e.g. “three sheep on the farm”). As shown at the bottom, VCode provides a compact, interpretable, and executable representation of original images.

Table 1. **Comparison of VCode with existing benchmarks.** VCode differs by generating code directly from natural images without extra query guidance, focusing on broad natural domains (G: General, C: College, P: Perception), and introducing CodeVQA to evaluate symbolic meaning preservation.

Benchmarks	Domain	Inputs	Outputs	Evaluation
<i>Coding</i>				
HumanEval [4]	Algorithm	Text	Code	Execute Pass
ChartMimic [28]	Chart	Text & Img	Code	Similarity
Design2Code [17]	Web UI	Text & Img	Code	Similarity
SWE-Bench [11]	GitHub	Text & Code	Code	Execute Pass
<i>Multi-modal</i>				
MM-Vet [31]	General	Img & Text	Text	OpenQA
MMMU [32]	College	Img & Text	Text	OpenQA/MCQ
CV-Bench [22]	Perception	Img & Text	Text	MCQ
VCode (Ours)	G&C&P	Img.	Code	Render→VQA

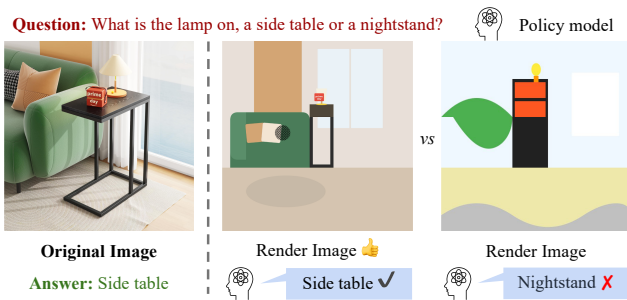


Figure 2. **Illustration of the CodeVQA prototype:** given an image and a question, the policy model answers based on the rendered image.

2.2. Evaluation Metrics

The key to our evaluation lies in defining the correspondence between the input image \mathcal{V} and the rendered SVG image $\tilde{\mathcal{V}}$.

Visual Similarity. An ideal SVG should faithfully preserve semantic content rather than merely matching pixels. We leverage a pretrained visual encoder $f(\cdot)$, such as SigLIP [33], to extract embeddings for both \mathcal{V} and $\tilde{\mathcal{V}}$, and compute their cosine similarity:

$$\mathcal{L}_{sim} = \max \cos(f(\mathcal{V}), f(\tilde{\mathcal{V}})). \quad (2)$$

We complement this with vision-only similarity metrics, including DINO [34], SSIM, and LPIPS.

CodeVQA. A more direct criterion is whether the rendered image $\tilde{\mathcal{V}}$ alone supports correct reasoning. We propose CodeVQA, where the evaluation is not constrained by the original image, but focuses directly on the correctness of answers derived from $\tilde{\mathcal{V}}$. Let ϕ be a policy model that outputs an answer \mathcal{A} given an image and a question \mathcal{Q} . The goal is formulated as:

$$\mathcal{A} = \phi(\tilde{\mathcal{V}}, \mathcal{Q}), \quad \mathcal{L}_{vqa} = \max 1[\text{Evaluator}(\mathcal{A})], \quad (3)$$

where $1[\cdot]$ is the indicator function and $\text{Evaluator}(\cdot)$ is a

- **Evaluation and Insights:** Extensive experiments expose the weaknesses of frontier VLMs in visual-centric coding. Human studies show consistent reasoning patterns between humans and VLMs on rendered SVGs, highlighting the promise of symbolic visual coding.

2. VCode Benchmark

2.1. Task Definitions

As illustrated in Fig. 2, given an input RGB image \mathcal{V} , a vision-language model ψ is tasked with generating SVG code \mathcal{C} that encodes the image. Rendering this code yields a rendered image $\tilde{\mathcal{V}}$. The objective is to minimize the discrepancy between the symbolic information of the original and rendered images:

$$\mathcal{L} = \min \|I(\mathcal{V}) - I(\tilde{\mathcal{V}})\|, \quad (1)$$

where $I(\cdot)$ denotes a representation of symbolic information. The challenge lies in defining an applicable measure for this symbolic information, which we elaborate on below.

rule-based matcher (for multiple-choice) or an LLM-as-Judge (for open-ended tasks). If the answer is correct, the SVG suffices to convey the required semantics; otherwise, it reveals a gap in representational fidelity.

Code Token Length. Beyond faithful representation, an effective coder should be concise. To assess this efficiency, we evaluate the length of the generated SVG in terms of its token count $|\mathcal{C}|$.

2.3. Data Curation

To develop appropriate question sets \mathcal{Q} for each associated image \mathcal{V} , we repurpose existing multimodal visual question answering benchmarks. To ensure diversity in taxonomy and difficulty, we focus on three representative scenarios: **(i) Commonsense Perception:** Assesses the ability to capture everyday semantics and spatial relationships. We incorporate all 218 image-question pairs from MM-Vet [31]; **(ii) Professional Knowledge:** Targets domain-specific, diploma-level tasks demanding deep reasoning. We filter the MMMU [32] development set to retain 146 single-image VQA instances; **(iii) Visual-centric:** Evaluates visually intensive settings (e.g. counting, distance, depth). We create a balanced 100-pair subset from CV-Bench [22] through stratified sampling.

In total, VCode yields 464 image-question pairs. Although smaller than existing text-based benchmarks, VCode emphasizes predicting extended code sequences, which incurs higher computational costs. A dataset of 464 samples represents an affordable yet sufficient scale for evaluating executable visual representations.

3. VCoder Framework

In practice, we find that directly prompting vision-language models to generate SVG code from natural images remains challenging. This difficulty arises from three factors: (i) **Long-Context Code Inputs:** composing thousands of tokens demands strong code reasoning over complex elements; (ii) **Visually-Blind Outputs:** the rendered image is unseen until execution, requiring the model to anticipate visual consequences; and (iii) **Weak Visual Fineness:** language models struggle to capture low-level boundaries and precise coordinates.

To address these intertwined challenges, we propose VCoder, which augments coders at test time with two complementary capabilities: *Thinking with Revision* and *Acting with Visual Tools*.

3.1. Thinking with Revision

Since the initial reconstruction may not always yield a satisfactory result, we design a revision strategy that allows the model to iteratively refine its outputs. This follows a two-step loop:

(i) Comment the Difference. Given an intermediate rendering $\tilde{\mathcal{V}}^{(t)}$ at iteration t , the coder ψ first perceives its deviation from the original image \mathcal{V} . We compute a difference signal quantifying the visual discrepancy:

$$\Delta^{(t)} \leftarrow \psi(\mathcal{V}, \tilde{\mathcal{V}}^{(t)}). \quad (4)$$

(ii) Revise with Render Feedback. The difference signal $\Delta^{(t)}$, together with the current code $\mathcal{C}^{(t)}$ and render $\tilde{\mathcal{V}}^{(t)}$, is provided back to the coder to generate revised code $\mathcal{C}^{(t+1)}$. Executing this code produces an updated reconstruction:

$$\mathcal{C}^{(t+1)} \leftarrow \psi(\mathcal{V}, \tilde{\mathcal{V}}^{(t)}, \Delta^{(t)}, \mathcal{C}^{(t)}), \quad \tilde{\mathcal{V}}^{(t+1)} \leftarrow \text{Render}(\mathcal{C}^{(t+1)}). \quad (5)$$

This revision loop is repeated for $t = 0, 1, \dots, T$, progressively refining the reconstruction until a satisfactory visual outcome is reached.

3.2. Acting with Visual Tools

To overcome inherent limitations in low-level perception, we equip the Coder with external visual tools to extract fine-grained attributes and translate them into structured code signals.

- **Category and Location.** We rely on bounding boxes predicted by Florence-2 [27], expressed as absolute coordinates (x_1, y_1, x_2, y_2) . These semantic labels (e.g. `id='bird'`) and geometry cues allow the Coder to anchor elements accurately, preserving the layout.
- **Shape.** To represent irregular boundaries, we employ SAM-2 [14] to generate segmentation masks. These masks are downsampled into sparse coordinate points via an adaptive simplification strategy. The resulting polygonal paths provide compact yet faithful shape descriptions.
- **Text.** Text carries critical semantic information that cannot be replaced by shapes. We apply OpenOCR to detect and transcribe text regions, encoding them into SVG using the native `<text>` tag, thereby preserving both content and visual attributes without pixel-rendering issues.

4. Experiments

4.1. Baseline and Settings

We evaluate our framework against a comprehensive suite of state-of-the-art models. *Proprietary models* include the Claude and Gemini series [7], the GPT family (GPT-5/5.2, o3, 4o/mini) [1, 9], and Seed-1.6-thinking [18]. These models serve as competitive upper baselines due to their advanced multimodal reasoning. *Open-source* models encompass LLaMA-4-Scout, the Qwen2.5/3-VL series [19], InternVL (3/3.5) [24, 35], Intern-S1, MiniCPM-V-4.5 [30], GLM-4.1/4.5V [20], and specialized models such as OmniSVG and StarVector [15, 29]. This selection covers a diverse range of parameter scales and training paradigms, enabling a robust comparison between leading commercial and open-source approaches.

Table 2. **Main results on VCode** across various top-performing frontier VLM coders. Top half is the proprietary models, while the bottom half is the open-source model. The best scores are in **bold** while the second best are in underline. The Overall score is calculated as an *instance-weighted average* of the three subtasks (MM-Vet, MMMU, and CV-Bench) using their respective question counts.

Model name	Success Rate (%)	SigLIP Score	Code Token (K)	CodeVQA												Overall
				Rec	Ocr	Know	MM-Vet			MMMU			CV-Bench			
				Gen	Spat	Math	Avg.	Avg.	Avg.	2D	3D	Avg.				
Orig. Image (4o-mini)	NA	100.0	NA	60.5	78.9	58.5	59.5	70.9	84.2	67.1	50.0	77.4	63.3	70.3	61.7	
Claude-4.5-Sonnet	99.1	66.8	1.9	29.7	57.6	11.9	17.0	57.3	52.7	39.0	42.5	50.4	55.0	52.7	43.1	
Claude-4-Opus	98.2	65.9	1.5	30.4	52.3	13.9	18.5	49.5	50.4	37.5	42.5	41.6	58.3	49.9	41.7	
Claude-4-Sonnet	98.2	65.5	1.6	31.8	51.2	<u>24.9</u>	<u>27.9</u>	44.8	34.6	37.8	39.0	49.0	53.3	51.2	41.1	
GPT-5	100.0	72.3	2.3	<u>33.9</u>	64.9	20.5	23.8	60.5	65.4	<u>43.9</u>	42.5	51.8	66.7	<u>59.2</u>	<u>46.8</u>	
GPT-4o	100.0	60.6	0.6	23.1	58.4	12.7	17.0	51.3	60.4	35.0	44.5	29.3	50.0	39.6	39.0	
GPT-o3	100.0	64.1	1.1	31.3	55.2	17.7	19.7	48.5	61.5	39.8	39.0	47.4	56.7	52.1	42.2	
GPT-4.1	100.0	68.6	1.6	30.8	62.0	15.5	20.4	56.0	55.8	40.9	44.5	48.2	66.7	57.4	45.6	
GPT-4o-mini	100.0	61.1	<u>0.4</u>	20.7	58.4	13.2	18.9	46.8	63.5	33.5	44.5	27.7	48.3	38.0	37.9	
Gemini-2.5-Pro	100.0	66.5	2.4	28.9	57.8	20.0	22.9	47.9	<u>68.5</u>	39.1	<u>45.2</u>	<u>56.1</u>	56.7	56.4	44.7	
Gemini-2.5-Flash	98.0	63.7	1.9	29.3	56.7	17.4	21.1	46.3	53.8	39.1	39.7	48.8	58.3	53.6	42.4	
Seed-1.6-Thinking	100.0	62.8	1.6	18.9	46.5	8.1	11.9	44.1	38.5	28.7	43.2	45.3	51.7	48.5	37.5	
Llama-4-Scout-17B-16E	100.0	55.5	0.7	18.2	44.9	12.4	15.5	32.8	46.2	26.4	42.5	35.0	53.3	44.2	35.3	
Qwen3-VL-235B-A22B	95.1	58.1	<u>1.7</u>	19.3	54.6	8.8	14.5	45.6	53.1	31.1	41.1	22.6	58.3	40.5	36.3	
Qwen2.5-VL-72B	98.7	57.9	0.3	20.6	52.9	14.0	17.3	51.3	43.1	31.8	41.1	21.9	53.3	37.6	36.0	
Qwen2.5-VL-7B	70.6	22.9	0.6	4.9	6.0	3.0	4.0	7.1	3.8	4.8	19.2	17.5	41.7	29.6	14.7	
InternVL3.5-241B-A28B	100.0	60.2	1.0	20.4	52.4	11.9	15.7	39.2	42.3	31.1	43.8	45.3	50.0	47.6	38.7	
Intern-S1	100.0	60.0	1.0	24.7	56.8	12.1	16.0	51.2	41.9	35.2	41.1	46.8	55.0	50.9	40.4	
InternVL3-78B	100.0	57.7	0.7	16.9	52.7	8.3	13.9	40.5	55.0	29.1	41.8	18.3	50.0	34.1	34.2	
MiniCPM-V-4.5	78.9	45.9	0.9	11.8	31.8	4.5	10.8	23.2	26.5	17.7	36.3	23.4	45.0	34.2	27.1	
GLM-4.5V	99.8	63.8	1.6	22.4	54.4	7.1	15.6	46.0	56.9	33.1	40.4	43.1	66.7	54.9	40.1	
GLM-4.1V-Thinking	100.0	61.7	1.2	21.1	52.0	10.4	13.7	44.8	58.8	31.9	43.2	37.9	56.7	47.3	38.8	
OmniSVG	100.0	46.2	5.3	9.2	15.3	3.7	10.4	16.9	11.5	9.4	43.8	24.8	40.0	32.4	25.2	
StarVector	8.3	18.1	1.3	0.0	3.4	0.0	1.6	4.4	0.0	1.5	6.8	0.0	0.0	0.0	2.8	
VCoder (Claude-4-Opus)	99.3	<u>71.0</u>	2.0	46.6	<u>63.4</u>	38.8	41.5	<u>58.1</u>	72.7	54.2 _{+16.7}	48.6 _{+6.2}	57.7	<u>65.0</u>	61.3 _{+11.4}	54.0 _{+12.3}	

210 4.2. Main Results and Analysis

211 Table 2 evaluates full baselines on VCode. We summarize
212 the key observations below:

213 **Stronger reasoning yields better visual coding scores.**
214 Closed-source models consistently outperform open-source
215 counterparts. Gemini-3-Pro sets the strongest standalone
216 baseline with the top SigLIP score (74.4) and the high-
217 est CodeVQA overall (52.1). This pattern indicates that
218 stronger linguistic and multi-modal reasoning ability trans-
219 lates into more faithful symbolic renderings. We also
220 observe a positive correlation between visual similarity
221 (SigLIP) and CodeVQA.

222 **Challenges across different dimensions.** (i) *Best per-*
223 *former still trails the original-image upper bound.* Even
224 Gemini-3-Pro (52.1) remains well below the raw-image up-
225 per bound (61.7), confirming that symbolic representation
226 has ample room for improvement. (ii) *SVG specialists un-*
227 *derperform.* OmniSVG and StarVector rank last due to
228 low success rates for long-context natural images, high-
229 lighting the gap between neatly authored SVG corpora and
230 SVGs derived from real-world images. (iii) *Knowledge and*
231 *Vision-centric perception are tough.* On MMMU, models
232 cluster within a narrow, modest band, failing in demand-
233 ing disciplinary settings. Similarly, CV-Bench scores hover
234 near random guessing on 3D spatial relations, where fine
235 structural abstraction is strictly required.

236 **Absolute gains with VCoder.** Built on Claude-4-Opus,
237 VCoder lifts the Overall score from 41.7 to 54.0 (+12.3) via

test-time revision and vision-tool assistance. It improves
238 across all three domains, demonstrating an effective, agen-
239 tic enhancement for visual-centric coding. 240

241 **Code token length correlates with expressiveness.** Mod-
242 els that emit short SVGs (*e.g.* 0.3K tokens by Qwen-2.5-
243 VL) underperform significantly. By contrast, stronger mod-
244 els (GPT-5, Gemini-2.5-Pro) produce longer sequences (of-
245 ten > 1.5K tokens) and attain higher scores. While length
246 is not sufficient on its own, performance scales with usable
247 context, highlighting long-context generation as a central
248 bottleneck for visual-centric coding.

249 **Qualitative analysis and more visualizations** can be
250 seen in supplementary material.

251 5. Conclusion

252 We introduced VCode, offering a new perspective on
253 visual-centric coding by benchmarking multimodal under-
254 standing with SVG as a naive visual representation. To
255 assess symbolic fidelity, we proposed CodeVQA, a proto-
256 col that evaluates reasoning solely over rendered SVGs.
257 Our extensive study reveals that frontier VLMs struggle
258 to produce faithful SVGs despite strong linguistic reason-
259 ing, exposing a persistent gap between language-driven and
260 vision-driven code generation. To address this, we proposed
261 VCoder, which integrates Test-time Revision and Acting
262 with Visual Tools, yielding substantial absolute improve-
263 ments. We hope this work pathways toward more human-
264 aligned, interpretable multimodal intelligence.

References

- 265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
- Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th annual acm symposium on user interface software and technology*, pages 1–22, 2023. 1 322 323 324
- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 3
- [2] Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*, 2021. 1
- [3] Tony Beltramelli. pix2code: Generating code from a graphical user interface screenshot. In *Proceedings of the ACM SIGCHI Symposium on Engineering Interactive Computing Systems*, New York, NY, USA, 2018. Association for Computing Machinery. 1
- [4] Mark Chen and Jerry Tworek et.al. Evaluating large language models trained on code. 2021. 1, 2
- [5] Siqi Chen, Xinyu Dong, Haolei Xu, Xingyu Wu, Fei Tang, Hang Zhang, Yuchen Yan, Linjuan Wu, Wenqi Zhang, Guiyang Hou, Yongliang Shen, Weiming Lu, and Yueting Zhuang. Svcgenius: Benchmarking llms in svg understanding, editing and generation, 2025. 1
- [6] Yamei Chen, Haoquan Zhang, Yangyi Huang, Zeju Qiu, Kaipeng Zhang, Yandong Wen, and Weiyang Liu. Symbolic graphics programming with large language models. *arXiv preprint arXiv:2509.05208*, 2025. 1
- [7] Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blisstein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025. 3
- [8] Yushi Hu, Weijia Shi, Xingyu Fu, Dan Roth, Mari Ostendorf, Luke Zettlemoyer, Noah A Smith, and Ranjay Krishna. Visual sketchpad: Sketching as a visual chain of thought for multimodal language models. *Advances in Neural Information Processing Systems*, 37:139348–139379, 2024. 1
- [9] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024. 3
- [10] Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*, 2024. 1
- [11] Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. Swebench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770*, 2023. 1, 2
- [12] Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, et al. Agentbench: Evaluating llms as agents. *arXiv preprint arXiv:2308.03688*, 2023. 1
- [13] Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 322 323 324
- [14] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, Eric Mintun, Junting Pan, Kalyan Vasudev Alwala, Nicolas Carion, Chaoyuan Wu, Ross Girshick, Piotr Dollár, and Christoph Feichtenhofer. Sam 2: Segment anything in images and videos, 2024. 1, 3 325 326 327 328 329 330 331
- [15] Juan A. Rodriguez, Abhay Puri, Shubham Agarwal, Issam H. Laradji, Pau Rodriguez, Sai Rajeswar, David Vazquez, Christopher Pal, and Marco Pedersoli. Starvector: Generating scalable vector graphics code from images and text. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16175–16186, 2025. 1, 3 332 333 334 335 336 337 338
- [16] Pratyusha Sharma, Tamar Rott Shaham, Manel Baradad, Stephanie Fu, Adrian Rodriguez-Munoz, Shivam Duggal, Phillip Isola, and Antonio Torralba. A vision check-up for language models. In *arXiv preprint*, 2024. 1 339 340 341 342
- [17] Chenglei Si, Yanzhe Zhang, Zhengyuan Yang, Ruibo Liu, and Diyi Yang. Design2code: How far are we from automating front-end engineering?, 2024. 1, 2 343 344 345
- [18] ByteDance Seed Team. Seed-oss open-source models. <https://github.com/ByteDance-Seed/seed-oss>, 2025. 3 346 347 348
- [19] Qwen Team. Qwen2.5-vl, 2025. 3 349
- [20] V Team. Glm-4.5v and glm-4.1v-thinking: Towards versatile multimodal reasoning with scalable reinforcement learning, 2025. 3 350 351 352
- [21] Minyang Tian, Luyu Gao, Shizhuo Dylan Zhang, Xinan Chen, Cunwei Fan, Xuefei Guo, Roland Haas, Pan Ji, Kittithat Krongchon, Yao Li, Shengyan Liu, Di Luo, Yutao Ma, Hao Tong, Kha Trinh, Chenyu Tian, Zihan Wang, Bohao Wu, Yanyu Xiong, Shengzhu Yin, Minhui Zhu, Kilian Lieret, Yanxin Lu, Genglin Liu, Yufeng Du, Tianhua Tao, Ofir Press, Jamie Callan, Eliu Huerta, and Hao Peng. Scicode: A research coding benchmark curated by scientists, 2024. 1 353 354 355 356 357 358 359 360
- [22] Shengbang Tong, Ellis Brown, Penghao Wu, Sanghyun Woo, Manoj Middepogu, Sai Charitha Akula, Jihan Yang, Shusheng Yang, Adithya Iyer, Xichen Pan, Austin Wang, Rob Fergus, Yann LeCun, and Saining Xie. Cambrian-1: A fully open, vision-centric exploration of multimodal llms, 2024. 1, 2, 3 361 362 363 364 365 366
- [23] Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv:2305.16291*, 2023. 1 367 368 369 370
- [24] Weiyun Wang, Zhangwei Gao, Lixin Gu, Hengjun Pu, Long Cui, Xingguang Wei, Zhaoyang Liu, Linglin Jing, Shenglong Ye, Jie Shao, et al. Internvl3.5: Advancing open-source multimodal models in versatility, reasoning, and efficiency. *arXiv preprint arXiv:2508.18265*, 2025. 3 371 372 373 374 375
- [25] Chengyue Wu, Yixiao Ge, Qishan Guo, Jiahao Wang, Zhixuan Liang, Zeyu Lu, Ying Shan, and Ping Luo. Plot2code: A comprehensive benchmark for evaluating multi-modal large 376 377 378

- 379 language models in code generation from scientific plots,
380 2024. 1
- 381 [26] Ronghuan Wu, Wanchao Su, Kede Ma, and Jing Liao. Icon-
382 shop: Text-guided vector icon synthesis with autoregressive
383 transformers. *ACM Trans. Graph.*, 42(6), 2023. 1
- 384 [27] Bin Xiao, Haiping Wu, Weijian Xu, Xiyang Dai, Houdong
385 Hu, Yumao Lu, Michael Zeng, Ce Liu, and Lu Yuan.
386 Florence-2: Advancing a unified representation for a vari-
387 ety of vision tasks. *arXiv preprint arXiv:2311.06242*, 2023.
388 1, 3
- 389 [28] Cheng Yang, Chufan Shi, Yaxin Liu, Bo Shui, Junjie Wang,
390 Mohan Jing, Linran Xu, Xinyu Zhu, Siheng Li, Yuxiang
391 Zhang, et al. Chartmimic: Evaluating lmm’s cross-modal
392 reasoning capability via chart-to-code generation. *arXiv*
393 *preprint arXiv:2406.09961*, 2024. 1, 2
- 394 [29] Yiyang Yang, Wei Cheng, Sijin Chen, Xianfang Zeng, Ji-
395 axu Zhang, Liao Wang, Gang Yu, Xinjun Ma, and Yu-Gang
396 Jiang. Omnisvg: A unified scalable vector graphics genera-
397 tion model. *arXiv preprint arxiv:2504.06263*, 2025. 1, 3
- 398 [30] Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui,
399 Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui He,
400 et al. Minicpm-v: A gpt-4v level mllm on your phone. *Nat*
401 *Commun* 16, 5509 (2025), 2025. 3
- 402 [31] Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang,
403 Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang.
404 Mm-vet: Evaluating large multimodal models for integrated
405 capabilities. In *International conference on machine learn-*
406 *ing*. PMLR, 2024. 1, 2, 3
- 407 [32] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi
408 Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming
409 Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Ren-
410 liang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo
411 Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhua Chen.
412 Mmmu: A massive multi-discipline multimodal understand-
413 ing and reasoning benchmark for expert agi. In *Proceedings*
414 *of CVPR*, 2024. 1, 2, 3
- 415 [33] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and
416 Lucas Beyer. Sigmoid loss for language image pre-training.
417 In *Proceedings of the IEEE/CVF international conference on*
418 *computer vision*, pages 11975–11986, 2023. 2
- 419 [34] Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun
420 Zhu, Lionel M. Ni, and Heung-Yeung Shum. Dino: Detr
421 with improved denoising anchor boxes for end-to-end object
422 detection, 2022. 2
- 423 [35] Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shen-
424 glong Ye, Lixin Gu, Hao Tian, Yuchen Duan, Weijie Su,
425 Jie Shao, et al. Internvl3: Exploring advanced training and
426 test-time recipes for open-source multimodal models. *arXiv*
427 *preprint arXiv:2504.10479*, 2025. 3

VCode: A Multimodal Coding Benchmark with SVG as Symbolic Visual Representation

Supplementary Material

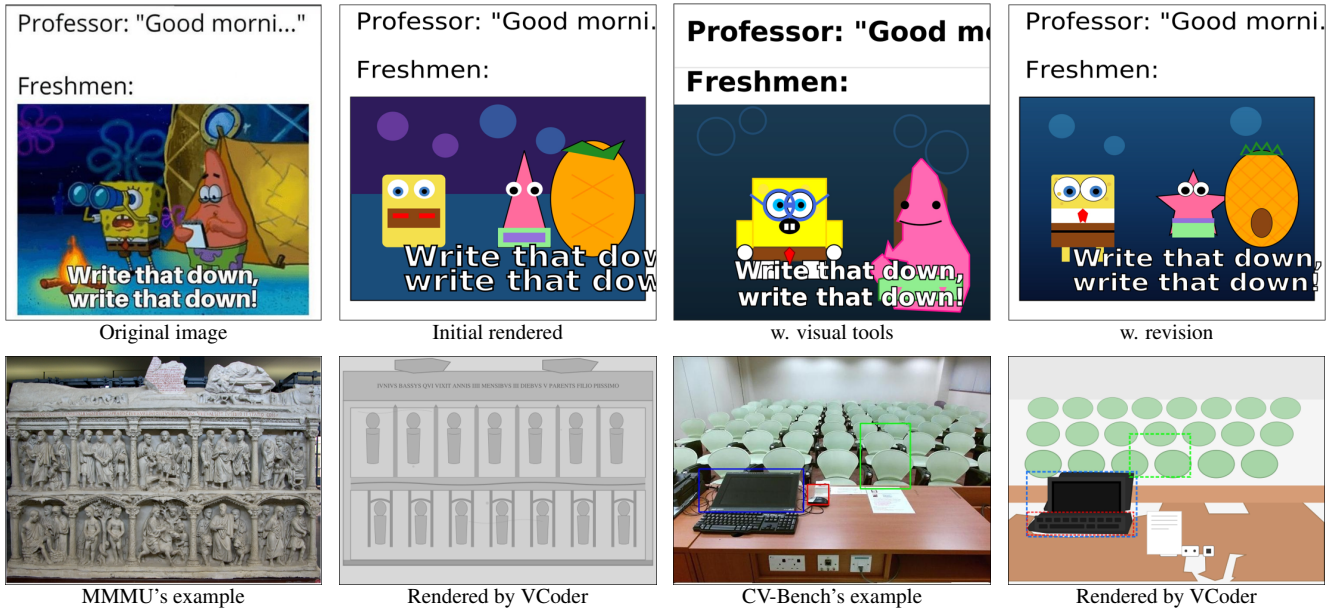


Figure 3. **Qualitative examples from VCode.** **Top row (a–d):** an internet meme rendered progressively by initial decoding, visual-tool assistance, and revision. **Bottom row:** challenge samples from MMMU (Art-Theory) and CV-Bench (3D), alongside their SVG renderings by VCoder.

A. Qualitative Analysis

Fig. 3 presents qualitative results by comparing original image and the rendered image by VCoder. **Top row.** Across four stages, the initial decoding misses layout and semantics. Adding *visual tools* recovers key geometry (e.g. the starfish character’s triangular body and facial features), while *revision* corrects fine details (character proportions, text alignment, spacing), yielding a rendering that closely matches the meme’s structure and intent. **Bottom row.** VCoder produces SVGs that are both more faithful to the source and more interpretable for downstream reasoning. The left example (MMMU) is knowledge-intensive: accurately depicting a multi-panel historical relief requires domain cues and fine structural abstraction, where base models often lose detail. The right example (CV-Bench) is vision-centric: success hinges on *visually grounded prompts* that localize and size objects correctly (e.g. monitor in front of keyboard, receding rows of chairs), after which revision tightens residual misalignments. These examples underscore the challenges posed by VCode.

B. Implement Details

We implement our model using the PyTorch framework on an NVIDIA RTX 4090 GPU with 24GB of memory. The maximum output length is set to 16,384 tokens, while for the Qwen2.5-VL models we use 8,192 tokens.

For evaluation, different protocols are used depending on the benchmark. In MM-Vet, we employ gpt-4-0613 as the evaluator to score model responses. In CV-Bench and MMMU experiments, we adopt a rule-based string matching parser to determine correctness.

For SigLip2, we use the siglip2-so400m-patch14-384. The token cost reported in our tables is measured using the tiktoken library with the cl100k_base encoding.

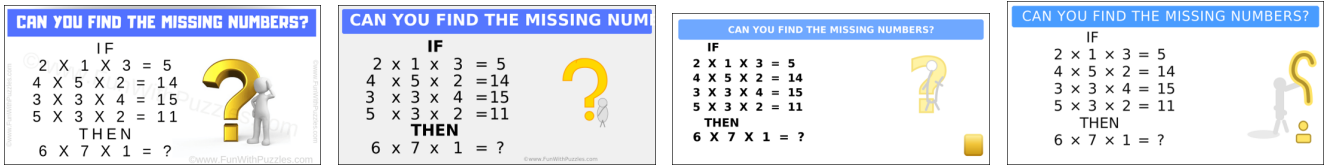
It is worth noting that in the *img2svg* experiments, StarVector cannot take textual prompts as input. It directly performs image-to-SVG generation.

451 **C. More Visualizations**

452 **C.1. VCoder vs. Baselines**

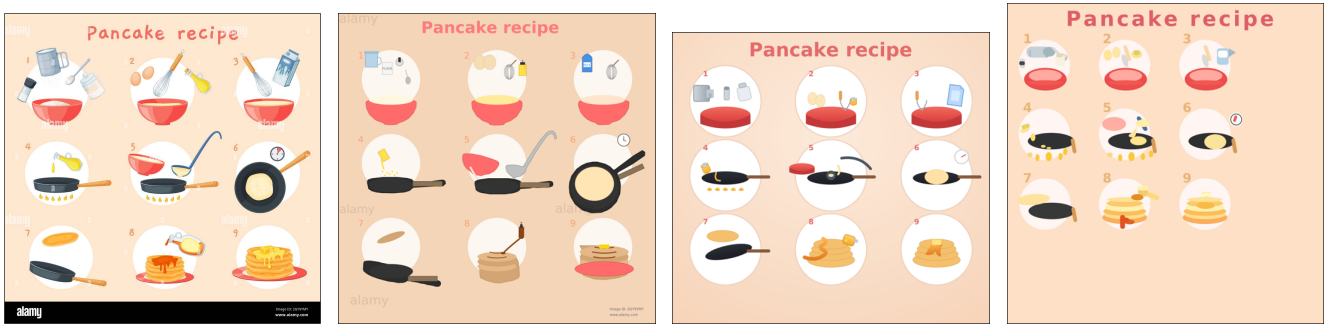
453 In this section, we present qualitative comparisons between VCoder and baseline models on representative examples from
 454 three benchmarks. For each case, we display: (a) the original reference image, (b) the output generated by VCoder, and (c–d)
 455 the visual results produced by the two strongest baseline models. These comparisons clearly illustrate VCoder’s superior
 456 ability to faithfully interpret and reconstruct visual content while preserving semantic consistency with the reference images.

457 **C.1.1. MM-VET**



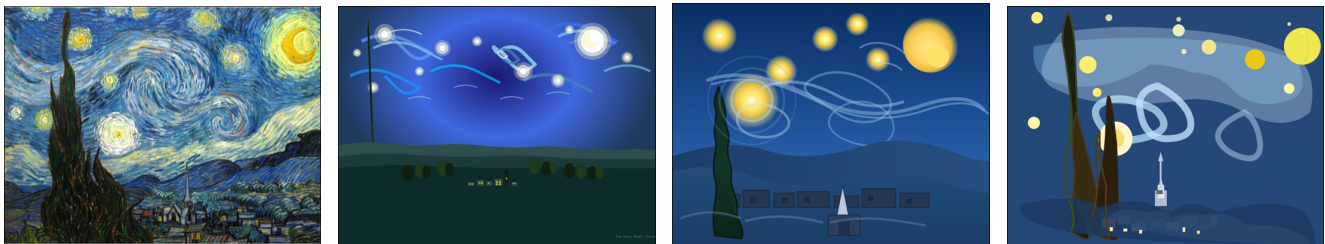
458 (a) Original image (b) VCoder (c) GPT-5 (d) GPT-4.1

459 **Question:** Find the pattern of how the "X" operator is redefined, and answer the given equation in the image. **Answer:** 13



460 (a) Original image (b) VCoder (c) GPT-5 (d) GPT-4.1

461 **Question:** What should we add in the third step? **Answer:** milk



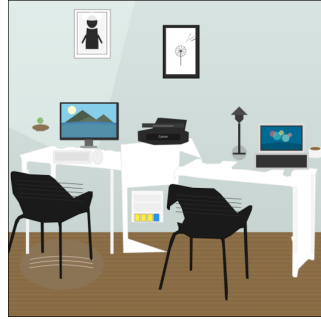
462 (a) Original image (b) VCoder (c) GPT-5 (d) GPT-4.1

463 **Question:** Can you give a short introduction to this painting?

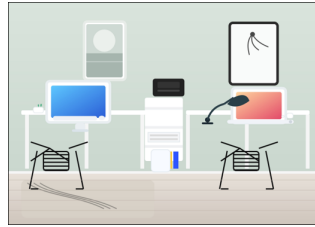
464 **Answer:** The Starry Night is an oil-on-canvas painting by the Dutch Post-Impressionist painter Vincent van Gogh. Painted
 465 in June 1889, it depicts the view from the east-facing window of his asylum room at Saint-Rémy-de-Provence, just before
 466 sunrise, with the addition of an imaginary village. It has been in the permanent collection of the Museum of Modern Art in
 467 New York City since 1941, acquired through the Lillie P. Bliss Bequest. Widely regarded as Van Gogh’s magnum opus, The
 468 Starry Night is one of the most recognizable paintings in Western art.



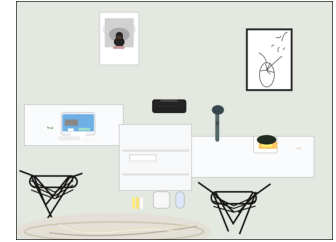
(a) Original image



(b) VCoder



(c) GPT-5



(d) GPT-4.1

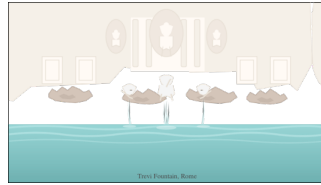
469

Question: On the right desk, what is to the left of the laptop? **Answer:** table lamp/desk lamp

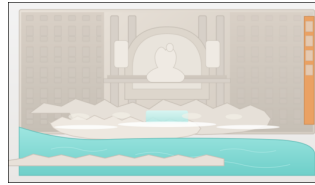
470



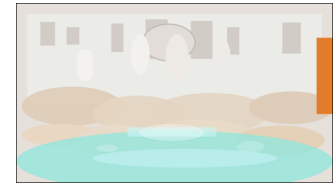
(a) Original image



(b) VCoder



(c) GPT-5



(d) GPT-4.1

471

Question: What is the name of this landmark? **Answer:** Trevi Fountain

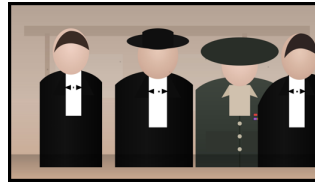
472



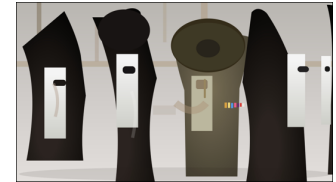
(a) Original image



(b) VCoder



(c) GPT-5



(d) GPT-4.1

473

Question: Can you give a short introduction to this movie?

474

Answer: The Godfather is a 1972 American crime film[2] directed by Francis Ford Coppola, who co-wrote the screenplay with Mario Puzo, based on Puzo's best-selling 1969 novel of the same title. The film stars Marlon Brando, Al Pacino, James Caan, Richard Castellano, Robert Duvall, Sterling Hayden, John Marley, Richard Conte, and Diane Keaton. It is the first installment in The Godfather trilogy, chronicling the Corleone family under patriarch Vito Corleone (Brando) from 1945 to 1955. It focuses on the transformation of his youngest son, Michael Corleone (Pacino), from reluctant family outsider to ruthless mafia boss.

475

476

477

478

479

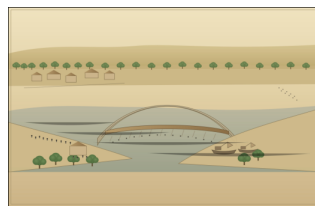
480



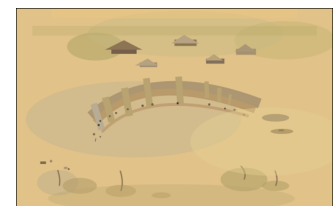
(a) Original image



(b) VCoder



(c) GPT-5



(d) GPT-4.1

481

Question: Can you give a short introduction to this painting?

482

Answer: Along the River During the Qingming Festival (Qingming Shanghe Tu) is a handscroll painting by the Song dynasty painter Zhang Zeduan (1085–1145) and copied many times in the following centuries. It captures the daily life of people and the landscape of the capital, Bianjing (present-day Kaifeng) during the Northern Song. The theme is often said to celebrate the festive spirit and worldly commotion at the Qingming Festival, rather than the holiday's ceremonial aspects,

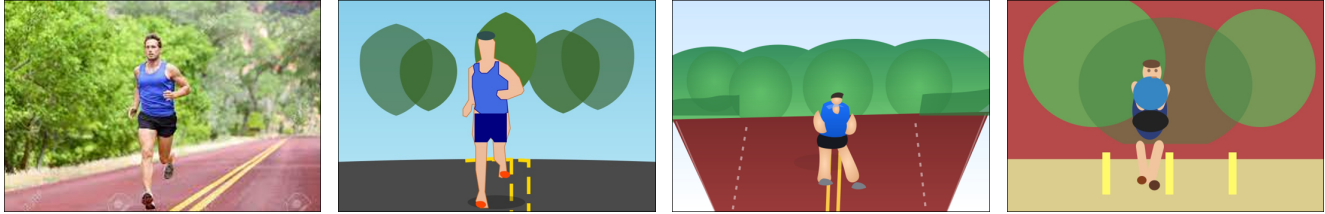
483

484

485

486

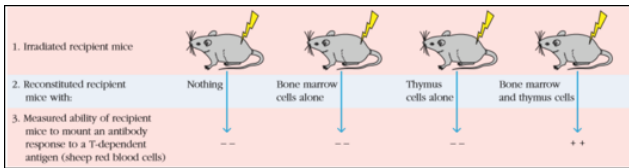
487 such as tomb sweeping and prayers. Read right to left, as a viewer unrolled it, successive scenes reveal the lifestyle of all
 488 levels of the society from rich to poor as well as economic activities in rural areas and the city, and offer glimpses of period
 489 clothing and architecture. The painting is considered to be the most renowned work among all Chinese paintings, and it has
 490 been called "China's Mona Lisa."



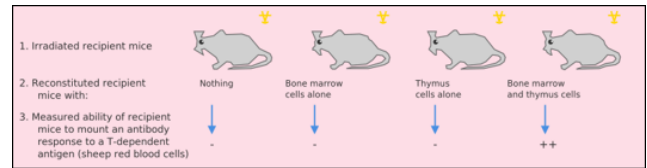
491 (a) Original image (b) VCoder (c) GPT-5 (d) GPT-4.1

492 **Question:** Is the man going to fall down? **Answer:** no

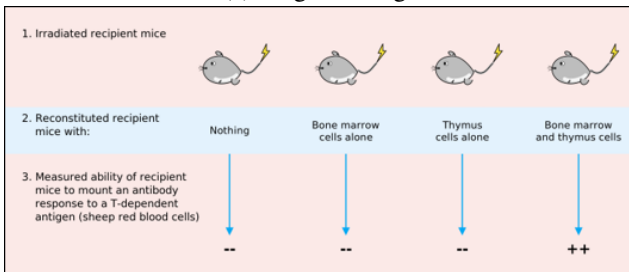
493 **C.1.2. MMMU**



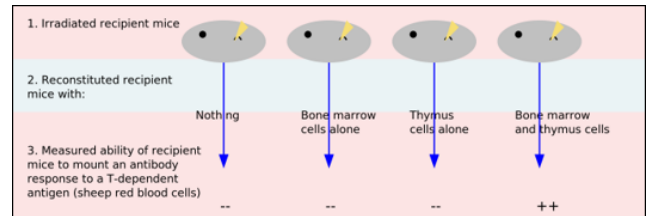
494 (a) Original image



(b) VCoder



(c) Gemini-2.5-Pro



(d) GPT-4o

495 **Question:** For your independent research, you transferred lymphocyte populations between syngeneic mice. You irradiated
 496 recipients first to ablate (get rid of) existing lymphocytes, then transferred defined cell populations from donors of same
 497 genetic background. The result is shown in . What does this experiment tell us? (A) Both B cells and T cells can produce
 498 antibodies. (B) Both B cells and T cells have memory functions. (C) Both B cells and T cells are required for an antibody
 499 response. (D) B cells are required for an antibody response in the absence of T cells. (E) B cells and T cells are co-localized
 500 and produce synergetic effects in bone marrow and thymus.
 501 Answer with the option's letter from the given choices directly.

502 **Answer:** C

I.R.	F.R.	Position of anchor point	Remarks
6.520	2.724	Outside the figure	The zero of the counting disc crossed the fixed index mark once in the clockwise direction.
1.222	7.720	Inside the figure	The zero of the counting disc crossed the fixed index mark twice in the anticlock direction.

504 (a) Original image

I.R.	F.R.	Position of anchor point	Remarks
6.520	2.724	Outside the figure	The zero of the counting disc crossed the fixed index mark once in the clockwise direction.
1.222	7.720	Inside the figure	The zero of the counting disc crossed the fixed index mark twice in the anticlockwise direction.

(b) VCoder

I.R.	F.R.	Position of anchor point	Remarks
6.520	2.724	Outside the figure	The zero of the counting disc crossed the fixed index mark once in the clockwise direction.
1.222	7.720	Inside the figure	The zero of the counting disc crossed the fixed index mark twice in the anticlock direction.

(c) Gemini-2.5-Pro

I.R.	F.R.	Position of anchor point	Remarks
6.520	2.724	Outside the figure	The zero of the counting disc crossed the fixed index mark once in the clockwise direction.
1.222	7.720	Inside the figure	The zero of the counting disc crossed the fixed index mark twice in the anticlock direction.

(d) GPT-4o

505

Question: Calculate the area of the zero circle with the following data: Assume that the tracing arm of the planimeter was so set that one revolution of the measuring wheel measures 100 cm^2 on the paper. Answer the question using a single word or phrase.

Answer: 1970.6

506

507

508

509

	2014	2015
Sales	\$ 9,402	\$10,091
Depreciation	1,350	1,351
Cost of goods sold	3,235	3,672
Other expenses	767	641
Interest	630	724
Cash	4,931	6,244
Accounts receivable	6,527	7,352
Short-term notes payable	953	895
Long-term debt	16,152	19,260
Net fixed assets	41,346	42,332
Accounts payable	5,179	5,022
Inventory	11,604	11,926
Dividends	1,147	1,261

(a) Original image

	2014	2015
Sales	\$ 9,402	\$10,091
Depreciation	1,350	1,351
Cost of goods sold	3,235	3,672
Other expenses	767	641
Interest	630	724
Cash	4,931	6,244
Accounts receivable	6,527	7,352
Short-term notes payable	953	895
Long-term debt	16,152	19,260
Net fixed assets	41,346	42,332
Accounts payable	5,179	5,022
Inventory	11,604	11,926
Dividends	1,147	1,261

(b) VCoder

	2014	2015
Sales	\$ 9,402	\$10,091
Depreciation	1,350	1,351
Cost of goods sold	3,235	3,672
Other expenses	767	641
Interest	630	724
Cash	4,931	6,244
Accounts receivable	6,527	7,352
Short-term notes payable	953	895
Long-term debt	16,152	19,260
Net fixed assets	41,346	42,332
Accounts payable	5,179	5,022
Inventory	11,604	11,926
Dividends	1,147	1,261

(c) Gemini-2.5-Pro

	2014	2015
Sales	\$ 9,402	\$10,091
Depreciation	1,350	1,351
Cost of goods sold	3,235	3,672
Other expenses	767	641
Interest	630	724
Cash	4,931	6,244
Accounts receivable	6,527	7,352
Short-term notes payable	953	895
Long-term debt	16,152	19,260

(d) GPT-4o

510

Question: For 2015, calculate the cash flow from assets(1) -----, cash flow to creditors(2) -----, and cash flow to stockholders(3) ----- (A) 1): -\$493.02 (2):-\$2,384 (3):\$1,890.98 (B) 1): \$1843.98 (2):-\$2,384 (3):\$493.02 (C) 1): -\$493.02 (2):-\$2,384 (3):-\$1,890.98

Answer with the option's letter from the given choices directly.

Answer: C

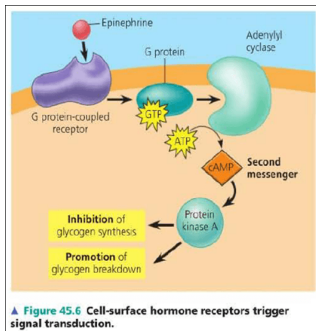
511

512

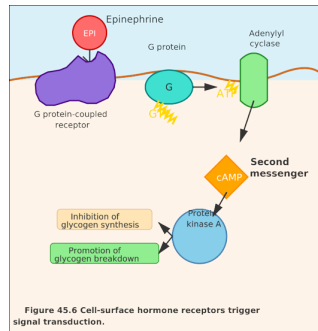
513

514

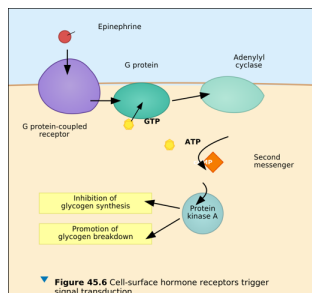
515



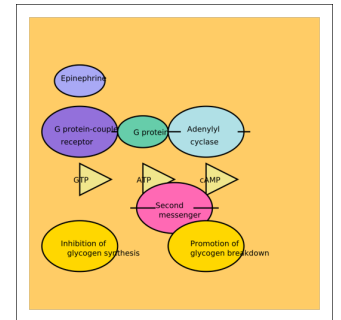
(a) Original image



(b) VCoder



(c) Gemini-2.5-Pro



(d) GPT-4o

516

Question: Which of the following correctly describes the reception stage of this signal transduction pathway? (A) epinephrine binds to a g-protein coupled receptor protein present in the cell membrane (B) the g protein changes shape, is activated, activates adenylyl cyclase, which activates cAMP, which activates protein kinases (C) protein kinases phosphorylate molecules (D) glycogen synthesis is inhibited and glycogen breakdown is promoted

Answer with the option's letter from the given choices directly.

Answer: A

517

518

519

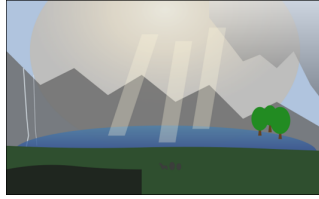
520

521

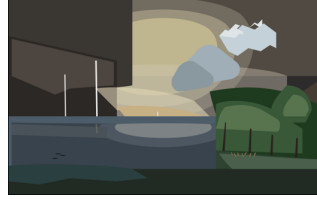
522



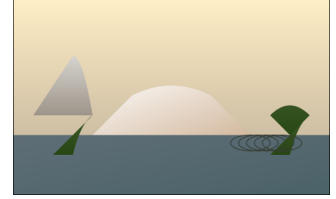
(a) Original image



(b) VCoder



(c) Gemini-2.5-Pro



(d) GPT-4o

523

524
525
526
527

Question: The painting on the right focuses on the (A) contribution of Native Americans to landscape preservation (B) implementation of the Homestead Act (C) impact of the gold rush on landscape development (D) idea of Manifest Destiny
Answer with the option's letter from the given choices directly.

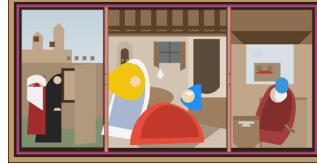
Answer: D



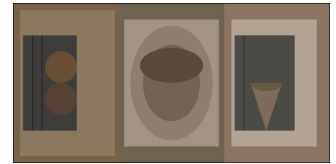
(a) Original image



(b) VCoder



(c) Gemini-2.5-Pro



(d) GPT-4o

528

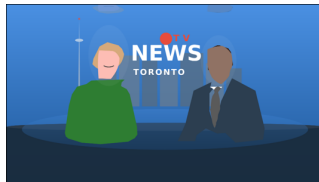
529
530
531

Question: Both works come from which art-historical period? (A) Baroque (B) Renaissance (C) Rococo (D) Classical
Answer with the option's letter from the given choices directly.

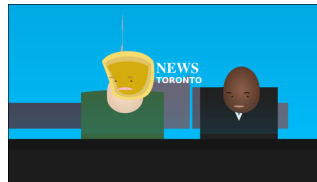
Answer: B



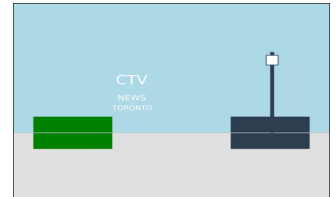
(a) Original image



(b) VCoder



(c) Gemini-2.5-Pro



(d) GPT-4o

532

533
534
535
536
537

Question: Refer to the figure, which term best describes the practice where students take on the role of television or newspaper reporters and interview characters from the book to retell an event from a range of perspectives? (A) News Program (B) Readers Theatre (C) Hot Seat (D) News
Answer with the option's letter from the given choices directly.

Answer: A



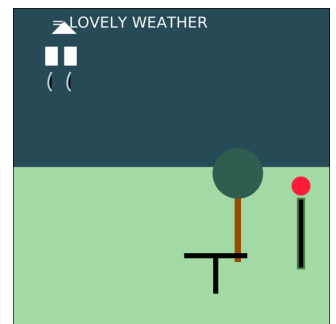
(a) Original image



(b) VCoder



(c) Gemini-2.5-Pro



(d) GPT-4o

538

539
540
541
542
543

Question: Refer to the description, which type of irony is depicted when a person says or writes one thing and means another, or uses words to convey a meaning opposite to the literal meaning? (A) verbal irony (B) situational irony (C) foreshadowing (D) dramatic irony
Answer with the option's letter from the given choices directly.

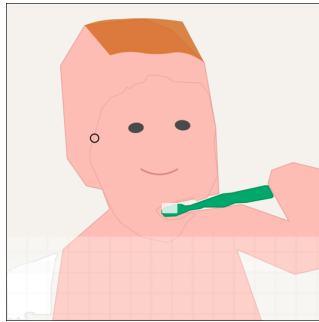
Answer: A

C.1.3. CV-Bench

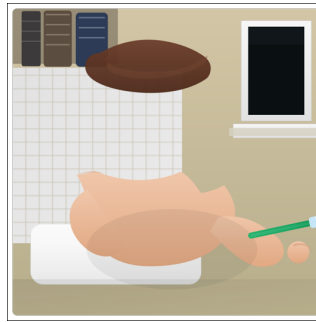
544



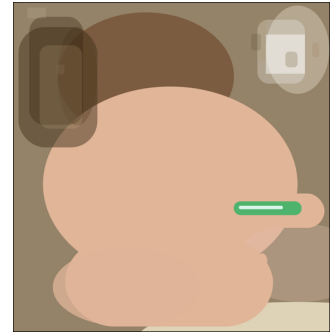
(a) Original image



(b) VCoder



(c) GPT-5



(d) GPT-4.1

545

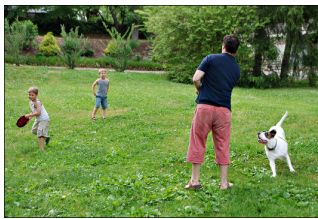
Question: How many persons are in the image? Select from the following choices. (A) 2 (B) 3 (C) 0 (D) 1 Answer with the option's letter from the given choices directly.

546

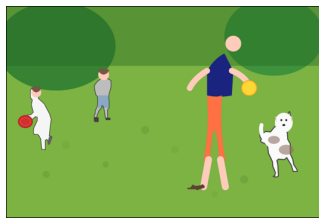
547

Answer: D

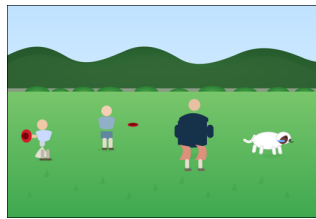
548



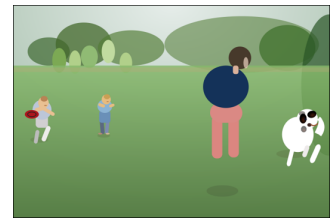
(a) Original image



(b) VCoder



(c) GPT-5



(d) GPT-4.1

549

Question: How many dogs are in the image? Select from the following choices. (A) 1 (B) 3 (C) 2 (D) 0 Answer with the option's letter from the given choices directly.

550

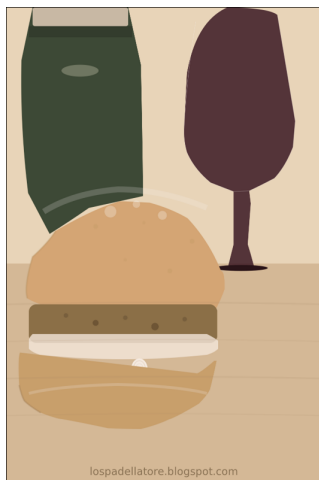
551

Answer: A

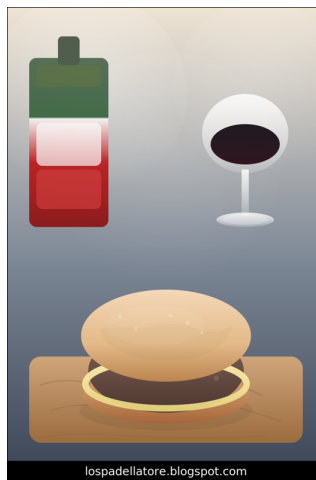
552



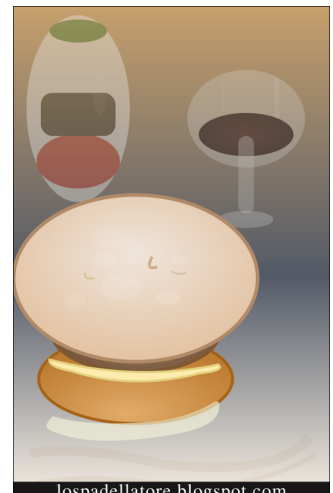
(a) Original image



(b) VCoder



(c) GPT-5



(d) GPT-4.1

553

Question: Considering the relative positions of the bottle and the wine glass in the image provided, where is the bottle located with respect to the wine glass? Select from the following choices. (A) left (B) right Answer with the option's letter from the given choices directly.

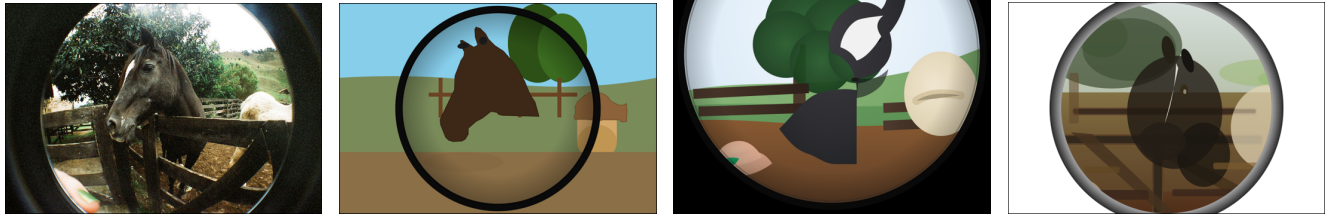
554

555

556

Answer: A

557



(a) Original image

(b) VCoder

(c) GPT-5

(d) GPT-4.1

558

Question: Considering the relative positions of the sheep and the horse in the image provided, where is the sheep located with respect to the horse? Select from (A) left (B) right
Answer with the option's letter from the given choices directly.

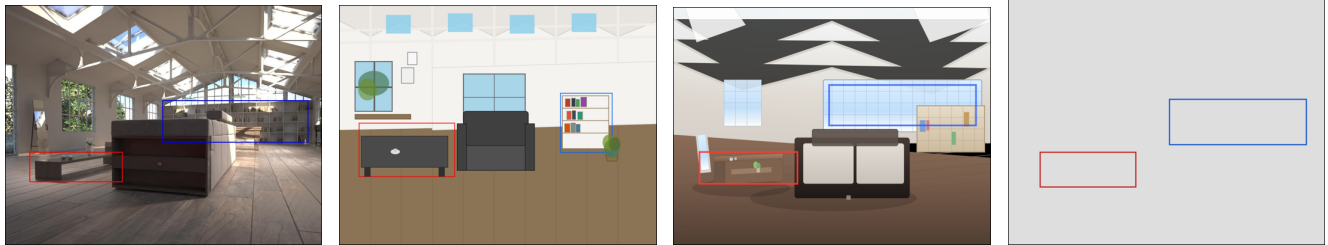
Answer: B

559

560

561

562



(a) Original image

(b) VCoder

(c) GPT-5

(d) GPT-4.1

563

Question: Which object is closer to the camera taking this photo, the table (highlighted by a red box) or the bookcase (highlighted by a blue box)? (A) table (B) bookcase
Answer with the option's letter from the given choices directly.

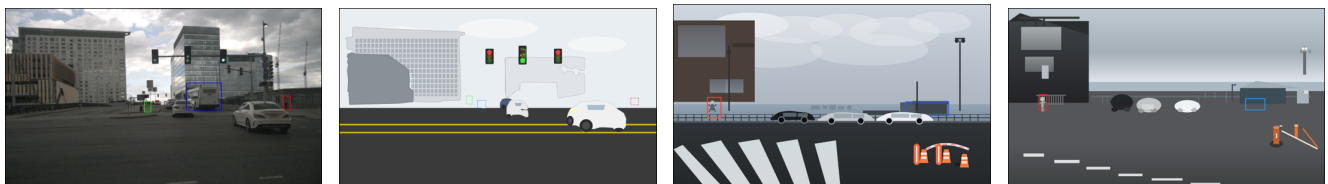
Answer: A

564

565

566

567



(a) Original image

(b) VCoder

(c) GPT-5

(d) GPT-4.1

568

Question: Estimate the real-world distances between objects in this image. Which object is closer to the traffic cone (highlighted by a red box), the trailer (highlighted by a blue box) or the bus (highlighted by a green box)? (A) trailer (B) bus
Answer with the option's letter from the given choices directly.

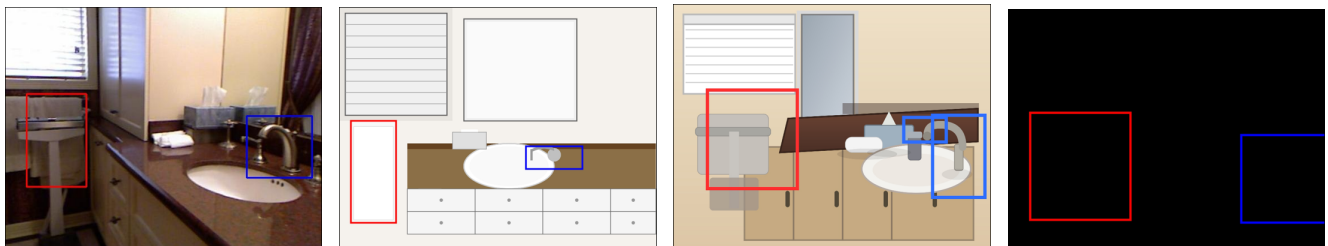
Answer: A

569

570

571

572



(a) Original image

(b) VCoder

(c) GPT-5

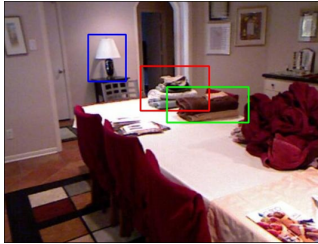
(d) GPT-4.1

573

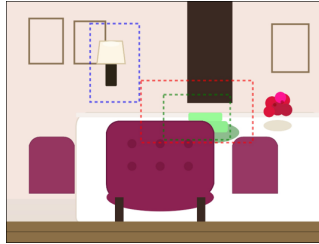
Question: Which object is closer to the camera taking this photo, the towel (highlighted by a red box) or the faucet (highlighted by a blue box)? (A) towel (B) faucet
Answer with the option's letter from the given choices directly.

574
575
576
577

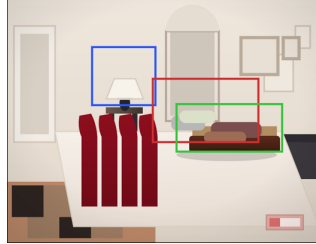
Answer: B



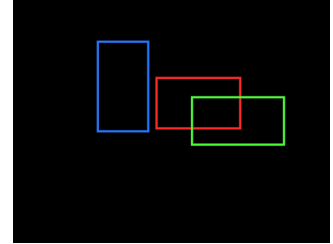
(a) Original image



(b) VCoder



(c) GPT-5



(d) GPT-4.1

578

Question: Estimate the real-world distances between objects in this image. Which object is closer to the clothes (highlighted by a red box), the lamp (highlighted by a blue box) or the towel (highlighted by a green box)? (A) lamp (B) towel
Answer with the option's letter from the given choices directly.

579
580
581
582
583

Answer: B

C.2. VCoder Individual Components

584

In this section, we present ablation studies visualizing the contribution of individual components in VCoder. For each example, we show: (a) the original reference image, (b) the initial rendered output without any refinement, (c) the output after applying visual tools, and (d) the final output after the revision module. These progressive visualizations demonstrate how each component incrementally improves the quality and accuracy of the generated images.

585
586
587
588

<p>Solve the following equations:</p> <p>1) $8x + 11 = 4x + 14$</p> <p>2) $7d - 4 = 11d - 9$</p>	<p>Solve the following equation:</p> <p>1) $8\Box + 11 = 4\Box + 14$</p> <p>2) $7\Box - 4 = 11\Box - 9$</p>	<p>Solve the following equation:</p> <p>1) $8x + 11 = 4x + 14$</p> <p>2) $7d - 4 = 11d - 9$</p>	<p>Solve the following equation:</p> <p>1) $8x + 11 = 4x + 14$</p> <p>2) $7d - 4 = 11d - 9$</p>
(a) Original image	(b) Initial rendered	(c) w. visual tools	(d) w. revision

589

Question: What is d in the last equation? **Answer:** $1.25 / \frac{5}{4}$.

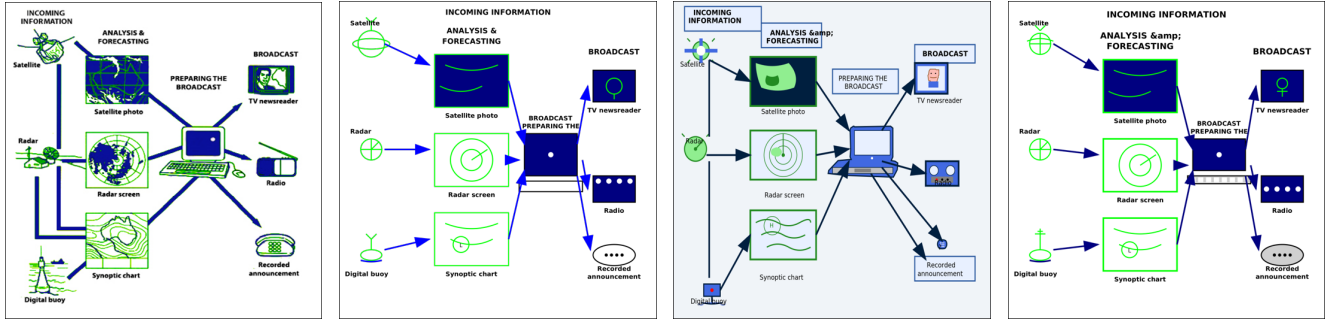
590

<table border="1"> <tr><td>$4 + 7 =$</td><td>$7 + 2 =$</td></tr> <tr><td>$2 + 2 =$</td><td>$6 + 1 =$</td></tr> <tr><td>$9 + 3 =$</td><td>$3 + 8 =$</td></tr> </table>	$4 + 7 =$	$7 + 2 =$	$2 + 2 =$	$6 + 1 =$	$9 + 3 =$	$3 + 8 =$	<table border="1"> <tr><td>$4 + 7 =$</td><td>$7 + 2 =$</td></tr> <tr><td>$2 + 2 =$</td><td>$6 + 1 =$</td></tr> <tr><td>$9 + 3 =$</td><td>$3 + 8 =$</td></tr> </table>	$4 + 7 =$	$7 + 2 =$	$2 + 2 =$	$6 + 1 =$	$9 + 3 =$	$3 + 8 =$	<table border="1"> <tr><td>$4 + 7 =$</td><td>$7 + 2 =$</td></tr> <tr><td>$2 + 2 =$</td><td>$6 + 1 =$</td></tr> <tr><td>$9 + 3 =$</td><td>$3 + 8 =$</td></tr> </table>	$4 + 7 =$	$7 + 2 =$	$2 + 2 =$	$6 + 1 =$	$9 + 3 =$	$3 + 8 =$	<table border="1"> <tr><td>$4 + 7 = 11$</td><td>$7 + 2 = 9$</td></tr> <tr><td>$2 + 2 = 4$</td><td>$6 + 1 = 7$</td></tr> <tr><td>$9 + 3 = 12$</td><td>$3 + 8 = 11$</td></tr> </table>	$4 + 7 = 11$	$7 + 2 = 9$	$2 + 2 = 4$	$6 + 1 = 7$	$9 + 3 = 12$	$3 + 8 = 11$
$4 + 7 =$	$7 + 2 =$																										
$2 + 2 =$	$6 + 1 =$																										
$9 + 3 =$	$3 + 8 =$																										
$4 + 7 =$	$7 + 2 =$																										
$2 + 2 =$	$6 + 1 =$																										
$9 + 3 =$	$3 + 8 =$																										
$4 + 7 =$	$7 + 2 =$																										
$2 + 2 =$	$6 + 1 =$																										
$9 + 3 =$	$3 + 8 =$																										
$4 + 7 = 11$	$7 + 2 = 9$																										
$2 + 2 = 4$	$6 + 1 = 7$																										
$9 + 3 = 12$	$3 + 8 = 11$																										
(a) Original image	(b) Initial rendered	(c) w. visual tools	(d) w. revision																								

591

Question: What is the answer to the second equation on the right? **Answer:** 12

592



593

(a) Original image

(b) Initial rendered

(c) w. visual tools

(d) w. revision

594

Question: The diagram below shows how the Australian Bureau of Meteorology collects up-to-the-minute information on the weather in order to produce reliable forecasts. Write a report for a university lecturer describing the information shown below. Write at least 150 words.

595

596

597

Answer: The figure illustrates the process used by the Australian Bureau of Meteorology to forecast the weather. There are four stages in the process, beginning with the collection of information about the weather. This information is then analysed, prepared for presentation, and finally broadcast to the public. Looking at the first and second stages of the process, there are three ways of collecting weather data and three ways of analysing it. Firstly, incoming information can be received by satellite and presented for analysis as a satellite photo. The same data can also be passed to a radar station and presented on a radar screen or synoptic chart. Secondly, incoming information may be collected directly by radar and analysed on a radar screen or synoptic chart. Finally, drifting buoys also receive data which can be shown on a synoptic chart. At the third stage of the process, the weather broadcast is prepared on computers. Finally, it is delivered to the public on television, on the radio, or as a recorded telephone announcement.

598

599

600

601

602

603

604

605



606

(a) Original image

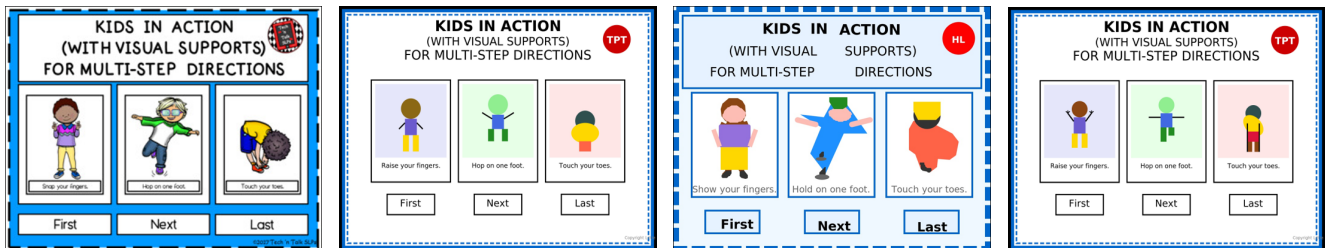
(b) Initial rendered

(c) w. visual tools

(d) w. revision

607

Question: What should I do before cutting herbs, sausage, and mushrooms? **Answer:** milk



608

(a) Original image

(b) Initial rendered

(c) w. visual tools

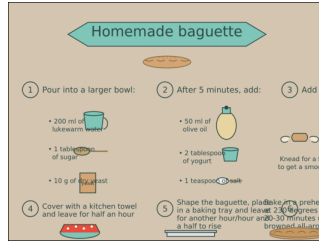
(d) w. revision

609

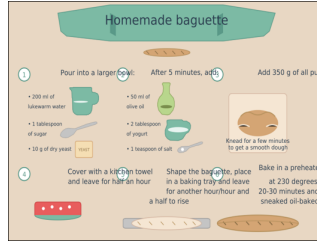
Question: What should kids do after snap fingers? **Answer:** hop on one foot



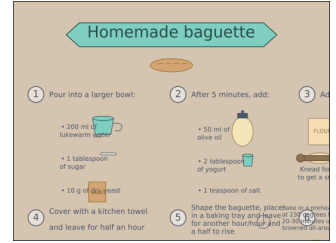
(a) Original image



(b) Initial rendered



(c) w. visual tools



(d) w. revision

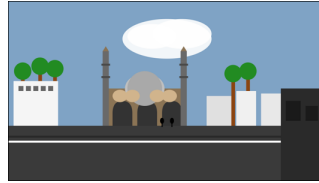
610

Question: What is the index of the step when we need to add all purpose flour? **Answer:** third / 3

611



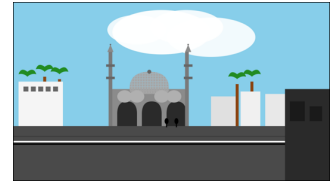
(a) Original image



(b) Initial rendered



(c) w. visual tools



(d) w. revision

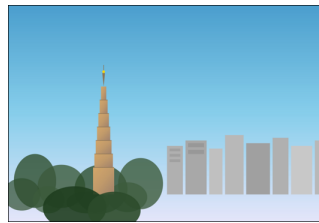
612

Question: What is the name of this landmark? **Answer:** Anbariya Mosque

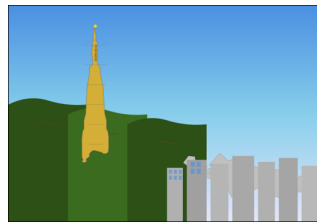
613



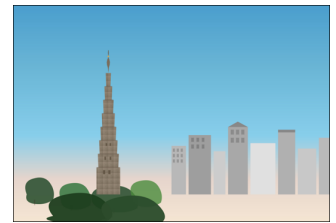
(a) Original image



(b) Initial rendered



(c) w. visual tools



(d) w. revision

614

Question: What is the name of this landmark? **Answer:** baochu pagoda

615



(a) Original image



(b) Initial rendered



(c) w. visual tools



(d) w. revision

616

Question: Can you give a short introduction to this painting?

617

Answer: Girl With A Pearl Earring (Dutch: Meisje met de parel) is an oil painting by Dutch Golden Age painter Johannes Vermeer, dated c. 1665. Going by various names over the centuries, it became known by its present title towards the end of the 20th century after the earring worn by the girl portrayed there. The work has been in the collection of the Mauritshuis in The Hague since 1902 and has been the subject of various literary and cinematic treatments..

618

619

620

621



(a) Original image



(b) Initial rendered



(c) w. visual tools



(d) w. revision

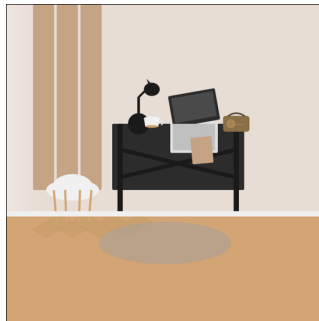
622

623

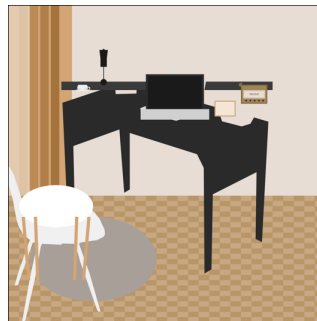
Question: What is located to the right of the shampoo? **Answer:** conditioner



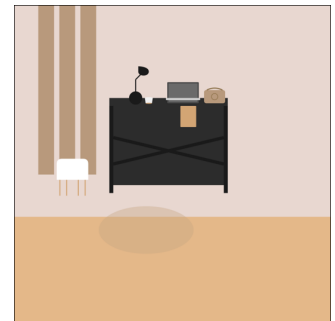
(a) Original image



(b) Initial rendered



(c) w. visual tools



(d) w. revision

624

625

Question: Is the curtain on the right side or on the left of the picture? **Answer:** left



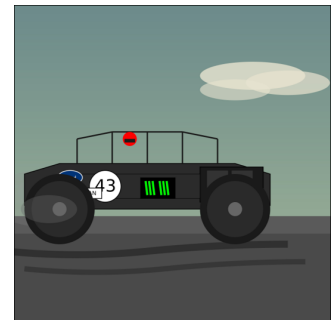
(a) Original image



(b) Initial rendered



(c) w. visual tools



(d) w. revision

626

627

Question: what is the green logo on the car? **Answer:** monster.