

Multimodal Inconsistency Reasoning (MMIR): A New Benchmark for Multimodal Reasoning Models

Anonymous ACL submission

Abstract

Existing Multimodal Large Language Models (MLLMs) are predominantly trained and tested on consistent visual-textual inputs, leaving open the question of whether they can handle inconsistencies in real-world, layout-rich content. To bridge this gap, we propose the Multimodal Inconsistency Reasoning (MMIR) benchmark to assess MLLMs’ ability to detect and reason about semantic mismatches in artifacts such as webpages, presentation slides, and posters. MMIR comprises 534 challenging samples, each containing synthetically injected errors across five reasoning-heavy categories: Factual Contradiction, Identity Misattribution, Contextual Mismatch, Quantitative Discrepancy, and Temporal/Spatial Incoherence. We evaluate six state-of-the-art MLLMs, showing that models with dedicated multimodal reasoning capabilities, such as o1, substantially outperform their counterparts while open-source models remain particularly vulnerable to inconsistency errors. Detailed error analyses further show that models excel in detecting inconsistencies confined to a single modality, particularly in text, but struggle with cross-modal conflicts and complex layouts. Probing experiments reveal that single-modality prompting, including Chain-of-Thought (CoT) and Set-of-Mark (SoM) methods, yields marginal gains, revealing a key bottleneck in cross-modal reasoning. Our findings highlight the need for advanced multimodal reasoning and point to future research on multimodal inconsistency.

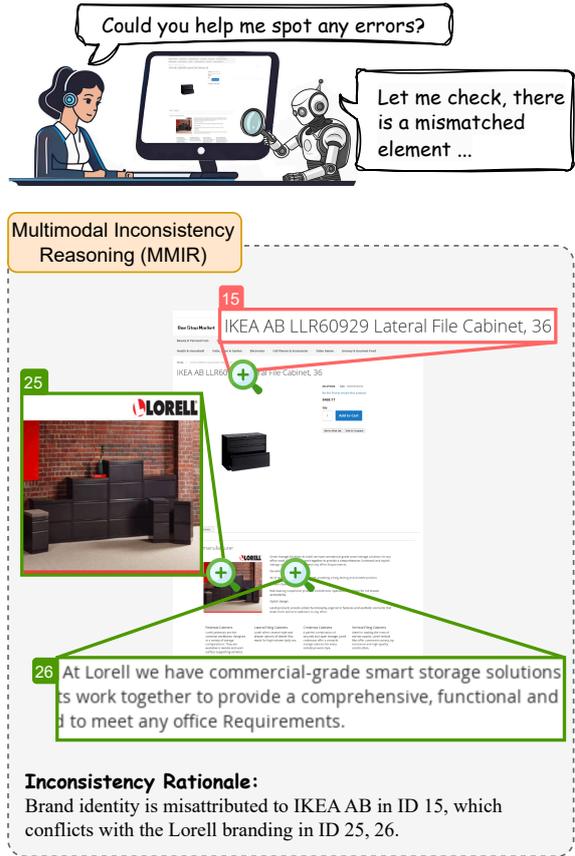


Figure 1: **An illustration of multimodal inconsistency reasoning on a webpage.** An agent examines a webpage where the brand “IKEA AB” is mentioned, but other elements clearly refer to “Lorell.” Detecting this brand identity misattribution requires the ability to compare text fields across different sections of the page and reconcile them with accompanying images or context—an inherently multimodal reasoning task.

1 Introduction

Recent advances in Large Language Models (LLMs) have demonstrated impressive reasoning abilities across a variety of tasks (OpenAI, 2024b; Guo et al., 2025; Kojima et al., 2022; Wei et al., 2022). Building on pre-trained LLMs, Multimodal Large Language Models (MLLMs) are fast evolving. However, they usually face greater challenges as they need to reason across different modalities,

especially when inconsistencies (i.e., mismatched or contradictory contents) exist. We find that, being primarily trained and evaluated on consistent visual-textual inputs, existing MLLMs are largely untested in scenarios where the input contains misaligned or contradictory information—a situation that is common in real-world scenarios. For example, in Figure 1, a user presents a web page

044
045
046
047
048
049
050
051

052	containing conflicting visual and textual elements,	104
053	asking the model to identify errors.	105
054	To comprehensively evaluate the ability of	106
055	MLLMs in reasoning over multimodal inconsis-	107
056	tency, we introduce the Multimodal Inconsistency	108
057	Reasoning Benchmark (MMIR) . MMIR is the	109
058	first framework dedicated to evaluating how ef-	110
059	fectively MLLMs can reason about and identify	111
060	semantic mismatches within complex, layout-rich	112
061	content with interleaved image and text compo-	113
062	nents. Our benchmark is built on a diverse collec-	114
063	tion of real-world artifacts (e.g. websites, slides,	115
064	posters) which have been augmented with syn-	116
065	thetic inconsistencies —realistic inconsistency er-	117
066	rors injected into their original structures. These	118
067	inconsistency errors span a range of reasoning-	
068	heavy categories: <i>Factual Contradiction</i> , <i>Iden-</i>	119
069	<i>tity Misattribution</i> , <i>Contextual Mismatch</i> , <i>Quan-</i>	120
070	<i>titative Discrepancy</i> , and <i>Temporal/Spatial Inco-</i>	121
071	<i>herence</i> —posing a next-level reasoning challenge	122
072	for models. For example, resolving a <i>Identity</i>	
073	<i>Misattribution</i> involves verifying entity alignment	123
074	across modalities, while <i>Quantitative Discrepancy</i>	124
075	requires cross-referencing chart data with textual	125
076	claims. By challenging models to detect such in-	126
077	consistencies, MMIR forces them to perform intri-	127
078	cate reasoning that goes well beyond simple pattern	128
079	recognition. This benchmark not only exposes the	129
080	limitations of current MLLMs in handling real-	
081	world challenges of reasoning over multimodal	130
082	content with inconsistency, but also provides a plat-	131
083	form for developing more robust multimodal rea-	132
084	soning systems.	133
085	In our experiments, we evaluated the ad-	134
086	vanced multimodal reasoning model o1 (OpenAI,	135
087	2024b) and five other state-of-the-art MLLMs:	
088	GPT-4o (OpenAI, 2024a), Qwen2.5-VL (Team,	136
089	2025), LLaVA-NeXT (Liu et al., 2024b), In-	
090	ternVL2.5 (Chen et al., 2024) and Phi-3.5-	137
091	Vision (Abdin et al., 2024) using MMIR’s 534 test	138
092	samples. The results overall underscore that cur-	139
093	rent MLLM models struggle with multimodal in-	140
094	consistency reasoning. Specifically, there is a stark	141
095	contrast between proprietary and open-source mod-	142
096	els. The open-source models evaluated only reach	143
097	less than 25% accuracy. o1 with strong reasoning	144
098	capability achieves the overall best performance	145
099	with over 50% accuracy.	146
100	To further understand the benchmarking results,	147
101	we conduct analysis based on the inconsistency	148
102	category, modality, and layout complexity of the	149
103	artifact. We find the proprietary models excel in	150
	identifying factual contradiction and identity mis-	151
	attribute types of inconsistency and inconsistency	152
	within a single modality, either image or text. Last	
	but not least, we investigate some approaches to en-	
	hance the model performance in our probing exper-	
	iment. The results indicate that text-based Chain-	
	of-Thought prompting and visual-based prompt-	
	ing (Set-of-Mark annotations) offer minimal and	
	sometimes adverse effects, whereas an iterative	
	multimodal interleaved reasoning strategy shows	
	promising gains. Overall, these results highlight a	
	critical bottleneck in the ability of MLLMs to per-	
	form robust, integrated reasoning—a key challenge	
	for future research.	
	Our contributions are threefold:	
	• We introduce MMIR, a novel benchmark that	
	targets the critical yet underexplored task of	
	multimodal inconsistency reasoning in layout-	
	rich content.	
	• We perform a comprehensive evaluation of	
	one leading multimodal reasoning model and	
	five state-of-the-art MLLMs, revealing sig-	
	nificant gaps in their ability to detect in-	
	consistency errors with detailed error analyses	
	across multiple error types, modalities, and	
	layout complexities.	
	• We provide detailed probing analyses that ex-	
	pose key challenges—from perceptual short-	
	comings to reasoning bottlenecks—and pro-	
	pose a framework that iteratively refines pre-	
	dictions by jointly leveraging visual and tex-	
	tual modalities.	
	2 Related Work	
	Multimodal Understanding and Reasoning	
	Multimodal Large Language Models (MLLMs)	
	process multimodal inputs by first processing vi-	
	sual inputs with pre-trained vision encoders such	
	as CLIP (Radford et al., 2021) to extract features,	
	and then projecting them into the textual repre-	
	sentation space with adapters (Liu et al., 2024a; Li	
	et al., 2023a). Significant efforts have been made to	
	bridge the gap between vision and text modalities	
	via integrating more cross-modality data such as in-	
	terleaved image-text sequences and visual ground-	
	ing data (Alayrac et al., 2022; Chen et al., 2023;	
	Peng et al., 2023). Also, some recent works de-	
	velop MLLMs with improved nuanced multimodal	
	abilities, such as Optical Character Recognition	
	(OCR) (Bai et al., 2023; Liu et al., 2024b), layout	

153 understanding (Feng et al., 2024; Fan et al., 2024a),
154 Graphic User Interface (GUI) interpretation (Liu
155 et al., 2024c; Team, 2025).

156 As MLLMs typically leverage pre-trained large
157 language models (LLMs) as the backbone, they
158 inherent strong textural reasoning abilities from
159 the advanced LLMs (Floridi and Chiriatti, 2020;
160 Touvron et al., 2023; Bai et al., 2023; Taori et al.,
161 2023; Chowdhery et al., 2023; OpenAI, 2024a;
162 Team, 2024). To further enhance the reasoning
163 ability of MLLMs, increasing efforts have focused
164 on improving MLLMs in multimodal reasoning.
165 The proprietary model, o1 (OpenAI, 2024b) first
166 realize strong multimodal reasoning with reason-
167 ing process similar to the Chain-of-Thought (Wei
168 et al., 2022) and other following works have also
169 explored the multimodal reasoning either through
170 training (Wu and Xie, 2024; Qi et al., 2024; Shao
171 et al., 2024) or prompting (Zhang et al., 2023,
172 2024b; Zheng et al., 2023).

173 **Multimodal Reasoning Benchmarks** To evalu-
174 ate the reasoning capabilities of MLLMs, numer-
175 ous benchmarks have been developed with vari-
176 ous focuses. Broad-coverage benchmarks such as
177 MM-Bench (Liu et al., 2024d), MMMU (Yue et al.,
178 2024) and MM-Vet (Yu et al., 2024) cover compre-
179 hensive reasoning challenges in real life scenarios,
180 offering holistic insights into model performance.
181 Others are developed with focuses on specific per-
182 spectives, such as TextVQA (Singh et al., 2019),
183 POPE (Li et al., 2023b) and MATHVERSE (Zhang
184 et al., 2024a) respectively challenge models with
185 tasks in domains of reasoning about text, objects,
186 mathematics in multimodal contexts. Recently, ad-
187 ditional benchmarks have emerged targeting arti-
188 ficially created multipanel images—such as posters
189 and screenshots—that combine several subfigures
190 in structured layouts (Fan et al., 2024b; Hsiao et al.,
191 2025), which require models to analyze spatial re-
192 lationships and hierarchical structures in complex vi-
193 sual contexts. However, current multi-modal bench-
194 marks assume visual-text alignment, overlooking
195 detecting critical errors of vision-language incon-
196 sistency in the input - a key challenge in real-world
197 scenarios. Instead, we evaluate MLLMs’ ability
198 to detect and localize such inconsistency via the
199 proposed MMIR benchmark.

200 **Inconsistency Checking** Existing works on tasks
201 related to checking or verifying inconsistency in
202 the input are primarily in the language domain. For
203 example, fact-checking (Thorne et al., 2018) re-

204 quires a model to first retrieve evidence and then
205 decide if a claim is supported, where the model
206 must reason if contradictory information existed in
207 the retrieved corpus. One step further, summary in-
208 consistency detection (Laban et al., 2022) focuses
209 on flagging any errors in summaries that create
210 contradictions regardless of correctness, including
211 incorrect use or hallucination of entities. As mod-
212 ern language models prosper, inconsistencies are
213 found existing within their outputs (Ravichander
214 et al., 2020) and across different outputs of para-
215 phrased queries (Elazar et al., 2021), and efforts
216 have been made towards the evaluation of those
217 inconsistencies (Fabbri et al., 2021; Wang et al.,
218 2020; Lattimer et al., 2023). In our research, we
219 lead efforts in detecting inconsistencies in the field
220 of vision and language.

221 3 MMIR

222 The MMIR benchmark is designed to assess how
223 effectively MLLMs can detect and localize seman-
224 tic mismatches within complex, layout-rich arti-
225 facts. Unlike conventional benchmarks that assume
226 coherent visual–textual inputs, MMIR challenges
227 models with realistic errors that require deep, cross-
228 modal reasoning. In MMIR, errors are defined and
229 categorized along five semantic dimensions:

230 *A. Factual Contradiction:* Direct conflict be-
231 tween two elements (text–text, text–image, or im-
232 age–image) within the modified content.

233 *B. Identity Misattribution:* Mislabeling of enti-
234 ties (objects, locations, brands, people) that conflict
235 with other elements.

236 *C. Contextual Mismatch:* Tonal, thematic, or
237 situational incompatibility between elements.

238 *D. Quantitative Discrepancy:* Numerical or sta-
239 tistical inconsistencies between elements.

240 *E. Temporal/Spatial Incoherence:* Implied time-
241 lines, dates, or spatial relationships that are impos-
242 sible or conflicting.

243 Figure 2 provides one example from each error
244 type across web, office, and poster artifacts, illus-
245 trating the diverse challenges MMIR poses.

246 3.1 Data Curation

247 MMIR’s data is curated through a four-stage
248 pipeline (Figure 3), ensuring high-quality, diverse,
249 and challenging test cases.

250 **Artifact Collection and Parsing** We begin by
251 manually selecting a total of 521 original arti-
252 facts from two domains: 349 webpages (sub-

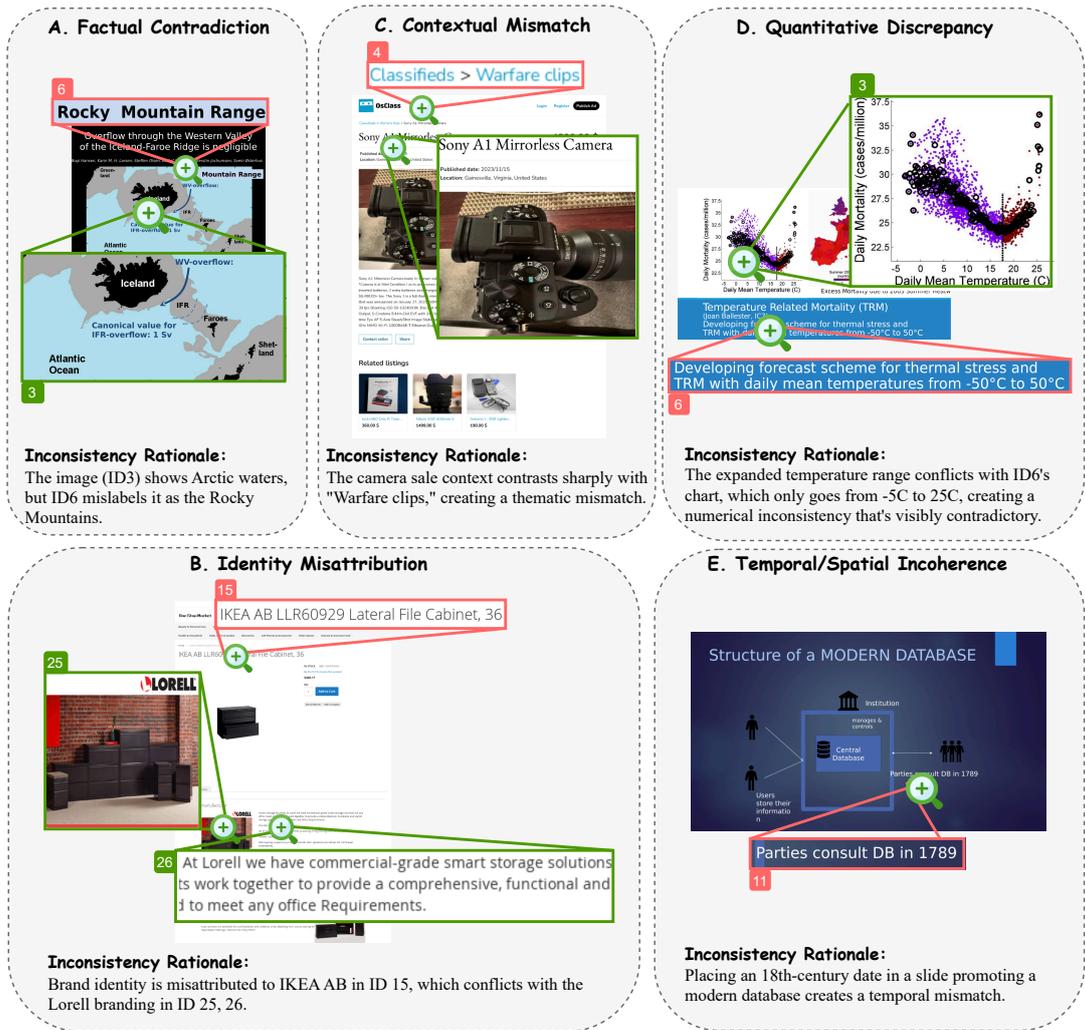


Figure 2: There are five inconsistency categories in the MMIR benchmark, posing diverse challenges.

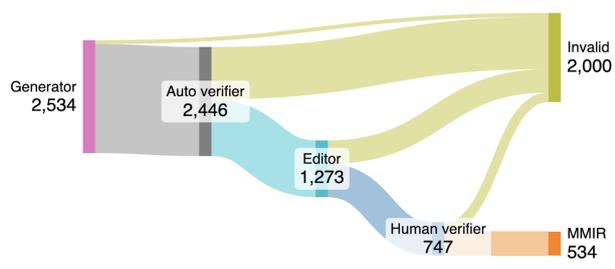


Figure 3: MMIR Data filtering process.

253
254
255
256
257
258
259
260
261
262

categories: shopping, classifieds, wiki) from VisualWebArena (Koh et al., 2024) and 172 presentations from Zenodo (European Organization For Nuclear Research and OpenAIRE, 2013), categorized into Office (sub-categories: slides, charts, diagrams) and Posters. Each artifact A_i is parsed using either using Document Object Model (DOM) or the python-pptx library to extract a set of elements $E_i = \{e_j\}_{j=1}^{n_i}$, where each element e_j is assigned a unique ID id_j and labeled with its type, content,

and a bounding box showing location information. Additionally, each artifact is paired with a Set-of-Marks (SoM) annotation A_i^{SoM} derived from E_i . This structured metadata forms the basis for subsequent error injection and question-answer curation.

Synthetic Inconsistency Generation To simulate real-world errors, we prompt an MLLM, o1-1217 (OpenAI, 2024b), as a generator with the annotated artifact and its element set $\{A_i^{SoM}, E_i\}$. The generator produces 2,534 proposals, each comprising a formatted edit instruction, the ground-truth element or element pair introducing the inconsistency:

$$GT \in \{id_j\} \cup \{(id_j, id_k) | j \neq k\},$$

the inconsistency error type, and the accompanying rationale. Following a self-evaluation loop (details in Appendix A.2), 2,446 valid proposals are retained.

Table 1: **MMIR Statistics.** Breakdown of the dataset by artifact category and error type.

Category	#Questions	Ave. #Elements
Artifact Categories		
Web	240	38.8
- Shopping	108	46.1
- Wiki	28	44.9
- Classifieds	104	29.5
Office	223	9.1
- Slides	102	9.4
- Tables/Charts	61	4.1
- Diagrams	60	13.9
Poster	71	27.6
Total	543	24.9
Error Categories		
Factual Contradiction	138	-
Identity Misattribution	84	-
Contextual Mismatch	141	-
Quantitative Discrepancy	76	-
Temporal/Spatial Incoherence	95	-
Total	543	-

Automated Editing and Human Verification

An auto-verification process then filters these proposals based on format and backend constraints (e.g., ensuring the target elements are editable), reducing the candidate set to 1,273, and saves low-level edit details, such as the path of the new image for an image edit, as inputs to the editor.

An automated editor-implemented using the Chrome DevTools Protocol (CDP) for web pages and python-pptx for presentations-executes the approved edits, generating for each successful operation a modified pair: $\{A'_i, E'_i\}$ where A'_i represents the modified artifact and E'_i contains the updated element metadata after the edit. For each pair, a descriptive caption set C_i is generated, where each caption within C_j details the element ID, location, and content summary of e'_j . These captions serve as references for later evaluation. More details on the verifier and editor are provided in Appendix A.3.

Finally, human experts review 747 edited samples, resulting in a final dataset of 534 validated quintuples: $D_{MMIR} = \{S'_i, E'_i, GT_i, category_i, rationale_i\}_{i=1}^{534}$, ensuring that only realistic and challenging samples remain. Table 1 provides a detailed breakdown by artifact type, subcategory, and error type. For example, webpages are further divided into shopping, wiki, and classifieds, each with its average number of elements, while errors are distributed across the five defined categories. Notably, the average word count in multiple-choice questions is 382.6, whereas open-ended responses are fixed at 59 words.

3.2 Evaluation

MMIR assesses a model’s ability to *detect inconsistency*, i.e., identifying and localizing semantic mismatches where elements deviate from their expected roles within an artifact. To assess the model’s performance comprehensively, each of the 534 test samples is provided to models under two distinct settings:

Open-Ended Setting Models receive the artifact A'_i with a fixed prompt Q_{open_ended} and generate a free-form response that identifies the semantic mismatch. This formulation evaluates the model’s ability to detect inconsistencies without relying on predefined answer options, thereby testing its unsupervised perception and reasoning.

Multiple-Choice Setting Models receive the artifact A'_i , but now with a combined prompt $Q_{MCQ} = (Q_{open_ended}, C_i)$. Each candidate in C_i is a textual description of an element. The model must select, from these options, the element(s) corresponding to the introduced inconsistency.

Evaluation Setup For the MCQ setting, we utilize regular expressions to compare the MLLM’s predicted answers against the ground truth, using accuracy as our metric. For the open-ended setting, o1-mini (0912) is employed as an LLM judge (Hsu et al., 2023; Hackl et al., 2023; Liu et al., 2023) to map the model’s free-form response back to the most likely ground-truth element IDs. The predicted IDs are then compared against GT_i to calculate accuracy.

4 Experiments and Analysis

We first evaluate the advanced multimodal reasoning model o1 (OpenAI, 2024b) and five other state-of-the-art MLLMs: GPT-4o (OpenAI, 2024a), Qwen2.5-VL (Team, 2025), LLaVA-NeXT (Liu et al., 2024b), InternVL2.5 (Chen et al., 2024) and Phi-3.5-Vision (Abdin et al., 2024) on the MMIR benchmark. We implement open-source models using their default settings and select the 1217 version of o1 and the 1120 version of GPT-4o for evaluation. Model implementation details are provided in Appendix B. We then examine error patterns across different inconsistency types and layout complexities and finally explore how prompting strategies affect multimodal reasoning under the open-ended setting.

Table 2: **The accuracy of six MLLMs under the two evaluation settings.** Proprietary models demonstrate higher performance as well as larger performance gain in the MCQ setting.

Models	Open-ended				Multiple-choice			
	Web	Office	Poster	Overall	Web	Office	Poster	Overall
<i>Proprietary Models</i>								
o1 (1217)	47.91	59.19	38.73	51.40	47.91	58.52	46.47	52.15
GPT-4o (1120)	25.00	42.60	30.98	33.14	37.29	58.96	47.88	47.75
<i>Open-sourced Models</i>								
Qwen2.5-VL-7B	8.54	29.14	11.97	17.60	14.37	33.18	16.90	22.56
LLaVA-NeXT-7B	10.20	21.97	7.04	14.70	11.45	25.33	5.63	16.47
InternVL2.5-8B	7.70	24.21	4.92	14.23	9.37	23.54	11.97	15.63
Phi-3.5-Vision-4B	6.87	24.43	7.04	14.23	1.66	8.52	0.00	4.30

4.1 Main Results

As shown in Table 2, proprietary models (o1 and GPT-4o) significantly outperform open-source alternatives, though all models exhibit substantial room for improvement. Appendix A.4 shows a qualitative example with question-answer and model response.

Performance Gap Between Reasoning, Proprietary and Open-Source Models. In both open-ended and MCQ settings, the reasoning o1 model substantially outperforms the rest, surpassing all open-source models by over 30%. The other proprietary model GPT-4o, although missing the explicit reasoning ability of o1, outperforms open-source alternatives, reflecting stronger multimodal alignment and reasoning capabilities.

Impact of Semantic Cues. GPT-4o sees a 14.61% accuracy boost in the MCQ setting with additional element descriptions as options, narrowing its gap with o1 from 18.26% to just 4.4%. This indicates that GPT-4o relies heavily on semantic context when available.

Inconsistent Gains for Open-Source Models. Most open-source models gain moderate or little accuracy when provided with MCQ-style prompts. Phi-3.5-Vision-4B experiences a 9.93% drop, suggesting weaker reasoning capacity and less effective use of textual cues. The gap between proprietary and open-source models widens further in MCQ (from 27.08% to 35.21%), highlighting the persistent challenge of integrating perceptual grounding with logical inference.

4.2 Error Analysis

4.2.1 Results Across Inconsistency Categories and Modalities

To investigate how different types of inconsistencies affect model performance, we show the results

across the category and modality of inconsistency in Figure 4.

Inconsistency Categories Figure 4(a) breaks down accuracy by the five inconsistency error categories. Proprietary models (o1, GPT-4o) outperform open-source models across the board, but the gap is particularly pronounced for *Factual Contradictions* and *Identity Misattribution*, implying that high-capacity models may have stronger factual grounding and entity recognition. Interestingly, *Temporal/Spatial Incoherence* also poses a substantial challenge for all models, highlighting a limitation in reasoning about time and space coherence.

Inconsistency Modalities In Figure 4(b), we examine how accuracy varies by the modality of the inconsistency. Overall, single-modality errors (those involving only one text or image field) yield the highest performance, with text-text inconsistencies proving especially tractable—likely because these language-centric models excel at purely textual reasoning. Next in difficulty are inter-modality errors (image-text), which require partial cross-modal integration but can still leverage textual anchors. Finally, image-image inconsistencies pose the greatest challenge, as they demand more advanced visual understanding and the ability to reconcile two distinct visual elements without the benefit of textual cues. These findings highlight that while language-focused models cope relatively well with purely textual conflicts, their capacity for deep visual or cross-modal reasoning remains underdeveloped.

4.2.2 Impact of Layout Complexity

We further examine the relationship between model accuracy and the number of elements in an artifact. To ensure statistical significance, we only include data points where at least 10 samples share

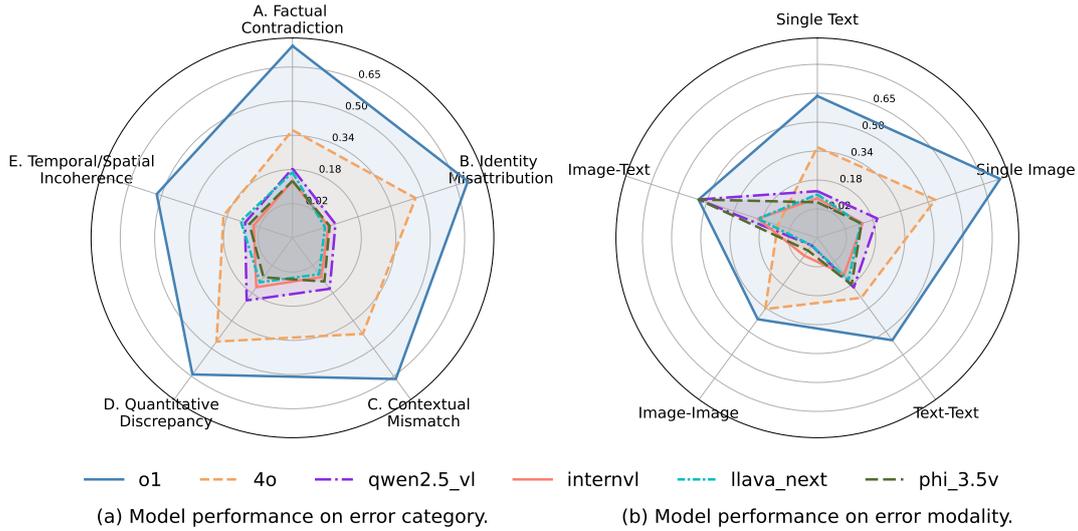


Figure 4: Fine-grained analysis of model performance.

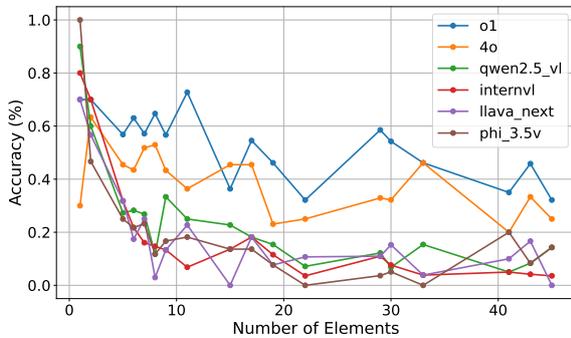


Figure 5: Model performance on layout complexity.

Table 3: Probing results of different prompting methods. Performance of each prompting method is directly compared with the vanilla setting. Gains are in blue and drops are in red.

Models	Vanilla	+ CoT	+ SoM	+ Both	MM-CoT
<i>Proprietary Models</i>					
o1 (1217)	51.40	–	-0.66	–	+0.09
GPT-4o (1120)	33.14	–	+5.34	–	+4.40
<i>Open-sourced Models</i>					
Qwen2.5-VL-7B	17.60	+0.28	+0.09	+0.28	+4.59
LLaVA-NeXT-7B	14.70	-1.78	-2.53	-0.47	+3.65
InternVL2.5-8B	14.23	+2.24	-0.66	-1.41	-0.85
Phi-3.5-Vision-4B	14.23	-0.38	+0.47	+0.84	+0.65

the same element count. As shown in Figure 7, the overall trend suggests that handling visually dense, information-rich artifacts remains a major challenge for current MLLMs. (1) Performance declines sharply as the number of elements increases, highlighting the difficulty in parsing cluttered layouts. (2) Proprietary models maintain higher accuracy in simpler layouts but degrade similarly in highly dense artifacts, indicating limitations in spatial reasoning. Open-source models struggle even in low-complexity settings, reinforcing the gap in perception and layout-aware inference.

4.3 Probing on Prompting Methods

We further investigate whether textual or visual prompts can alleviate the reasoning bottleneck. Table 3 compares *Chain-of-Thought (CoT)* prompting (Wei et al., 2022) and *Set-of-Mark (SoM)* visual augmentation (Yang et al., 2023), as well as their combination. We also explored an interleaved multimodal reasoning strategy, which we term *Multimodal Interleaved CoT (MM-CoT)* to further in-

tegrate and refine reasoning across both visual and textual modalities.

4.3.1 Chain-of-Thought (CoT) Prompting

To assess whether explicit reasoning instructions can enhance performance, we apply CoT prompting (Wei et al., 2022) to the four open-sourced models (benchmarked proprietary models have API guides to not include additional CoT prompting).

As shown in Table 3, CoT prompting yields negligible or even negative effects on accuracy. This suggests that simply injecting explicit reasoning steps is insufficient when the underlying model lacks strong cross-modal alignment or robust logical inference mechanisms.

4.3.2 Set-of-Mark (SoM) Prompting

We next examine the effect of SoM visual prompting (Yang et al., 2023). By overlaying bounding boxes onto the artifact screenshots (example in Figure 6), we aim to enhance the models' ability to perceive and localize elements.

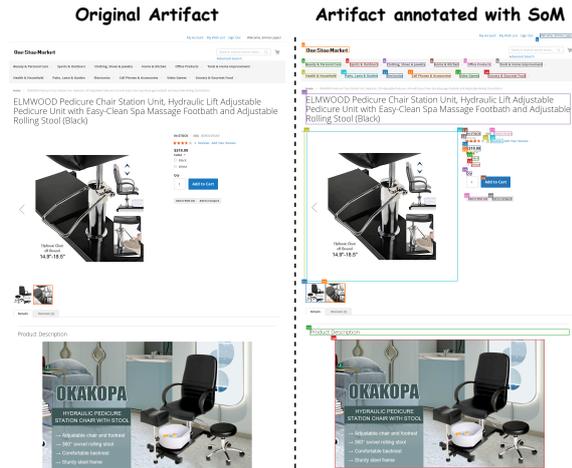


Figure 6: Example of original artifact in MMIR (left) and artifact annotated with Set-of-Mark in the probing analysis (right).

The result shows that these additional visual cues yield moderate improvements for GPT-4o (5.34%) yet confuse the rest of the models, leading to little or even slightly degraded performance, likely because the additional visual cues interfere with the model’s initial perception.

When combined with CoT prompting, SoM provides little gains for some open-source models but remains largely inconsistent or even detrimental for others. This indicates that simply stacking CoT and SoM techniques does not guarantee improved performance, underscoring the need for more sophisticated strategies to unify visual cues with explicit reasoning steps.

4.3.3 Multimodal Interleaved CoT (MM-CoT)

Our previous analyses indicate that single-modality prompts (CoT or SoM) often yield minimal or even detrimental gains in the open-ended setting when models receive no textual hints about which elements might be inconsistent. We hypothesize that MMIR tasks demand *iterative* reasoning that tightly integrates both visual and textual modalities. To address this, we propose *Multimodal Interleaved CoT (MM-CoT)*, a two-stage approach explicitly designed to weave visual cues into a step-by-step reasoning process:

Stage 1: Initial Candidate Generation The model receives the same input in Stage 1 as in the open-ended setting, generating its top five predictions (along with associated reasoning). Using o1-mini (0912) to interpret these responses, we map each prediction back to one or a pair of ele-

ment IDs from the artifact’s metadata C_i . We then highlight the bounding boxes of those elements on the artifact image, producing a SoM-annotated version to be used in the next stage.

Stage 2: Multimodal Refinement The model is subsequently given the SoM-annotated artifact from Stage 1, alongside the textual reasoning it generated previously. This additional visual context helps the model refine its earlier predictions, integrating both the visual bounding-box annotations and the initial textual reasoning to arrive at a final answer.

Results As shown in Table 3, MM-CoT outperforms all other prompting methods. GPT-4o, for example, improves by 4.40% over its vanilla baseline, while open-source models gain an average of around 2% improvements. These findings underscore the importance of iterative cross-modal reasoning: once textual inferences guide which visual elements to focus on, SoM annotations become more informative, and the overall reasoning process becomes more accurate. Although the bounding boxes used for SoM are derived from ground-truth references, this probing experiment demonstrates that *interleaved* multimodal interaction is a promising direction for closing the reasoning gap in challenging, inconsistency-heavy scenarios.

5 Discussion and Conclusion

In this work, we introduce the *Multimodal Inconsistency Reasoning Benchmark (MMIR)* to evaluate how well MLLMs detect and localize semantic mismatches in complex real-world artifacts. MMIR challenges models across five error categories and two reasoning settings for a detailed assessment of multimodal reasoning. Our experiments show that even advanced proprietary models struggle with open-ended inconsistency detection. Although providing natural-language descriptions in a multiple-choice format offers modest gains, standard prompting techniques (e.g., Chain-of-Thought and Set-of-Mark) yield inconsistent or negative effects, while a proposed Multimodal Interleaved CoT (MM-CoT) method that iteratively refines reasoning by integrating visual and textual modalities, yielding greater performance improvements. Despite these advances, significant challenges remain, motivating further research on robust multimodal reasoning for real-world inconsistency detection.

549 Limitations

550 While MMIR provides a rigorous framework for
551 evaluating multimodal inconsistency reasoning, it
552 is not without its limitations. Annotating and ver-
553 ifying inconsistencies in layout-rich artifacts re-
554 mains a labor-intensive process. Although MMIR’s
555 pipeline integrates automated editing and verifica-
556 tion, the overall scale is still limited by the need for
557 careful human review. Although these domains cap-
558 ture a range of layouts and content types, they do
559 not encompass the full variety of real-world multi-
560 modal artifacts (e.g., multi-page documents, social
561 media feeds, or mobile application interfaces). On
562 the other hand, synthetic error generation—while
563 effective for systematically introducing controlled
564 inconsistencies—may not perfectly mirror the nu-
565 anced mistakes that occur in human-generated con-
566 tent. This could lead to discrepancies between
567 model performance on MMIR and in truly open-
568 ended, real-world scenarios. Scaling up the dataset
569 to cover broader domains, more intricate layouts,
570 and diverse error types would strengthen its ability
571 to serve as a comprehensive benchmark for real-
572 world multimodal inconsistency detection.

573 References

574 Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed
575 Awadallah, Ammar Ahmad Awan, Nguyen Bach,
576 Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat
577 Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck,
578 Sébastien Bubeck, Martin Cai, Qin Cai, Vishrav
579 Chaudhary, Dong Chen, Dongdong Chen, Weizhu
580 Chen, Yen-Chun Chen, Yi-Ling Chen, Hao Cheng,
581 Parul Chopra, Xiyang Dai, Matthew Dixon, Ronen
582 Eldan, Victor Fragoso, Jianfeng Gao, Mei Gao,
583 Min Gao, Amit Garg, Allie Del Giorno, Abhishek
584 Goswami, Suriya Gunasekar, Emman Haider, Jun-
585 heng Hao, Russell J. Hewett, Wenxiang Hu, Jamie
586 Huynh, Dan Iter, Sam Ade Jacobs, Mojan Javaheripi,
587 Xin Jin, Nikos Karampatziakis, Piero Kauffmann,
588 Mahoud Khademi, Dongwoo Kim, Young Jin Kim,
589 Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi
590 Li, Yunsheng Li, Chen Liang, Lars Liden, Xihui
591 Lin, Zeqi Lin, Ce Liu, Liyuan Liu, Mengchen Liu,
592 Weishung Liu, Xiaodong Liu, Chong Luo, Piyush
593 Madan, Ali Mahmoudzadeh, David Majercak, Matt
594 Mazzola, Caio César Teodoro Mendes, Arindam Mi-
595 tra, Hardik Modi, Anh Nguyen, Brandon Norrick,
596 Barun Patra, Daniel Perez-Becker, Thomas Portet,
597 Reid Pryzant, Heyang Qin, Marko Radmilac, Liliang
598 Ren, Gustavo de Rosa, Corby Rosset, Sambudha Roy,
599 Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil
600 Salim, Michael Santacrose, Shital Shah, Ning Shang,
601 Hiteshi Sharma, Yelong Shen, Swadheen Shukla, Xia
602 Song, Masahiro Tanaka, Andrea Tupini, Praneetha
603 Vaddamanu, Chunyu Wang, Guanhua Wang, Lijuan

Wang, Shuohang Wang, Xin Wang, Yu Wang, Rachel
604 Ward, Wen Wen, Philipp Witte, Haiping Wu, Xiaoxia
605 Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu,
606 Weijian Xu, Jilong Xue, Sonali Yadav, Fan Yang,
607 Jianwei Yang, Yifan Yang, Ziyi Yang, Donghan Yu,
608 Lu Yuan, Chenruidong Zhang, Cyril Zhang, Jianwen
609 Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan
610 Zhang, and Xiren Zhou. 2024. [Phi-3 technical report: A highly capable language model locally on your phone](#). *Preprint*, arXiv:2404.14219. 613

Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc,
614 Antoine Miech, Iain Barr, Yana Hasson, Karel
615 Lenc, Arthur Mensch, Katherine Millican, Malcolm
616 Reynolds, et al. 2022. Flamingo: a visual language
617 model for few-shot learning. *Advances in neural
618 information processing systems*, 35:23716–23736. 619

Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang,
620 Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou,
621 and Jingren Zhou. 2023. Qwen-vl: A versatile
622 vision-language model for understanding, localiza-
623 tion, text reading, and beyond. *arXiv preprint
624 arXiv:2308.12966*, 1(2):3. 625

Keqin Chen, Zhao Zhang, Weili Zeng, Richong Zhang,
626 Feng Zhu, and Rui Zhao. 2023. Shikra: Unleashing
627 multimodal llm’s referential dialogue magic. *arXiv
628 preprint arXiv:2306.15195*. 629

Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu,
630 Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye,
631 Hao Tian, Zhaoyang Liu, Lixin Gu, Xuehui Wang,
632 Qingyun Li, Yiming Ren, Zixuan Chen, Jiapeng Luo,
633 Jiahao Wang, Tan Jiang, Bo Wang, Conghui He, Bo-
634 tian Shi, Xingcheng Zhang, Han Lv, Yi Wang, Wenqi
635 Shao, Pei Chu, Zhongying Tu, Tong He, Zhiyong
636 Wu, Hui Deng, Jiaye Ge, Kaiming Chen, Min Dou,
637 Lewei Lu, Xizhou Zhu, Tong Lu, Dahu Lin, Yunfeng
638 Qiao, Jifeng Dai, and Wenhai Wang. 2024. [Expanding performance boundaries of open-source multi-modal models with model, data, and test-time scaling](#). *ArXiv*, abs/2412.05271. 642

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin,
643 Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul
644 Barham, Hyung Won Chung, Charles Sutton, Sebas-
645 tian Gehrmann, et al. 2023. Palm: Scaling language
646 modeling with pathways. *Journal of Machine Learn-
647 ing Research*, 24(240):1–113. 648

Yanai Elazar, Nora Kassner, Shauli Ravfogel, Abhi-
649 lasha Ravichander, Eduard Hovy, Hinrich Schütze,
650 and Yoav Goldberg. 2021. Measuring and improving
651 consistency in pretrained language models. *Transac-
652 tions of the Association for Computational Linguis-
653 tics*, 9:1012–1031. 654

European Organization For Nuclear Research and Ope-
655 nAIRE. 2013. [Zenodo](#). 656

Alexander R Fabbri, Chien-Sheng Wu, Wenhao Liu,
657 and Caiming Xiong. 2021. Qafacteval: Improved
658 qa-based factual consistency evaluation for summa-
659 rization. *arXiv preprint arXiv:2112.08542*. 660

661	Yue Fan, Lei Ding, Ching-Chen Kuo, Shan Jiang, Yang Zhao, Xinze Guan, Jie Yang, Yi Zhang, and Xin Eric Wang. 2024a. Read anywhere pointed: Layout-aware gui screen reading with tree-of-lens grounding. <i>arXiv preprint arXiv:2406.19263</i> .	717
662		718
663		719
664		720
665		
666	Yue Fan, Jing Gu, Kaiwen Zhou, Qianqi Yan, Shan Jiang, Ching-Chen Kuo, Yang Zhao, Xinze Guan, and Xin Wang. 2024b. Muffin or chihuahua? challenging multimodal large language models with multipanel vqa. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 6845–6863.	721
667		722
668		723
669		724
670		725
671		
672		
673	Weixi Feng, Wanrong Zhu, Tsu-jui Fu, Varun Jampani, Arjun Akula, Xuehai He, Sugato Basu, Xin Eric Wang, and William Yang Wang. 2024. Layoutgpt: Compositional visual planning and generation with large language models. <i>Advances in Neural Information Processing Systems</i> , 36.	726
674		727
675		728
676		729
677		
678		
679	Luciano Floridi and Massimo Chiriatti. 2020. Gpt-3: Its nature, scope, limits, and consequences. <i>Minds and Machines</i> , 30:681–694.	730
680		731
681		732
682	Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. <i>arXiv preprint arXiv:2501.12948</i> .	733
683		734
684		735
685		736
686		737
687	Veronika Hackl, Alexandra Elena Müller, Michael Granitzer, and Maximilian Sailer. 2023. Is gpt-4 a reliable rater? evaluating consistency in gpt-4’s text ratings. In <i>Frontiers in Education</i> , volume 8, page 1272229. Frontiers Media SA.	738
688		739
689		740
690		741
691		742
692	Yu-Chung Hsiao, Fedir Zubach, Gilles Baechler, Srinivas Sunkara, Victor Carbune, Jason Lin, Maria Wang, Yun Zhu, and Jindong Chen. 2025. Screenqa: Large-scale question-answer pairs over mobile app screenshots . <i>Preprint</i> , arXiv:2209.08199.	743
693		744
694		745
695		746
696		
697	Ting-Yao Hsu, Chieh-Yang Huang, Ryan Rossi, Sungchul Kim, C Lee Giles, and Ting-Hao K Huang. 2023. Gpt-4 as an effective zero-shot evaluator for scientific figure captions. <i>arXiv preprint arXiv:2310.15405</i> .	747
698		748
699		749
700		750
701		751
702	Jing Yu Koh, Robert Lo, Lawrence Jang, Vikram Duvvur, Ming Chong Lim, Po-Yu Huang, Graham Neubig, Shuyan Zhou, Ruslan Salakhutdinov, and Daniel Fried. 2024. Visualwebarena: Evaluating multimodal agents on realistic visual web tasks .	752
703		753
704		754
705		755
706		756
707	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. <i>Advances in neural information processing systems</i> , 35:22199–22213.	757
708		758
709		759
710		760
711		761
712	Philippe Laban, Tobias Schnabel, Paul N Bennett, and Marti A Hearst. 2022. Summac: Re-visiting nli-based models for inconsistency detection in summarization. <i>Transactions of the Association for Computational Linguistics</i> , 10:163–177.	762
713		763
714		764
715		765
716		766
	Barrett Martin Lattimer, Patrick Chen, Xinyuan Zhang, and Yi Yang. 2023. Fast and accurate factual inconsistency detection over long documents. <i>arXiv preprint arXiv:2310.13189</i> .	767
		768
		769
		770
	Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023a. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In <i>International conference on machine learning</i> , pages 19730–19742. PMLR.	767
		768
		769
		770
	Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. 2023b. Evaluating object hallucination in large vision-language models . <i>Preprint</i> , arXiv:2305.10355.	767
		768
		769
		770
	Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024a. Improved baselines with visual instruction tuning. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 26296–26306.	767
		768
		769
		770
	Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024b. Llava-next: Improved reasoning, ocr, and world knowledge .	767
		768
		769
		770
	Junpeng Liu, Tianyue Ou, Yifan Song, Yuxiao Qu, Wai Lam, Chenyan Xiong, Wenhui Chen, Graham Neubig, and Xiang Yue. 2024c. Harnessing webpage uis for text-rich visual understanding . <i>Preprint</i> , arXiv:2410.13824.	767
		768
		769
		770
	Yang Liu, Dan Iter, Yichong Xu, Shuhang Wang, Ruochen Xu, and Chenguang Zhu. 2023. G-eval: Nlg evaluation using gpt-4 with better human alignment. <i>arXiv preprint arXiv:2303.16634</i> .	767
		768
		769
		770
	Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. 2024d. Mmbench: Is your multi-modal model an all-around player? In <i>European conference on computer vision</i> , pages 216–233. Springer.	767
		768
		769
		770
	OpenAI. 2024a. Gpt-4o system card . <i>Preprint</i> , arXiv:2410.21276.	767
		768
		769
		770
	OpenAI. 2024b. Openai o1 system card . <i>Preprint</i> , arXiv:2412.16720.	767
		768
		769
		770
	Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. 2023. Kosmos-2: Grounding multimodal large language models to the world. <i>arXiv preprint arXiv:2306.14824</i> .	767
		768
		769
		770
	Ji Qi, Ming Ding, Weihang Wang, Yushi Bai, Qingsong Lv, Wenyi Hong, Bin Xu, Lei Hou, Juanzi Li, Yuxiao Dong, et al. 2024. Cogcom: Train large vision-language models diving into details through chain of manipulations. <i>arXiv preprint arXiv:2402.04236</i> .	767
		768
		769
		770
	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from	767
		768
		769
		770

771	natural language supervision. In <i>International conference on machine learning</i> , pages 8748–8763. PMLR.	Jianwei Yang, Hao Zhang, Feng Li, Xueyan Zou, Chun yue Li, and Jianfeng Gao. 2023. Set-of-mark prompting unleashes extraordinary visual grounding in gpt-4v . <i>ArXiv</i> , abs/2310.11441.	824
772			825
773	Abhilasha Ravichander, Eduard Hovy, Kaheer Suleman, Adam Trischler, and Jackie Chi Kit Cheung. 2020. On the systematicity of probing contextualized word representations: The case of hypernymy in bert. In <i>Proceedings of the Ninth Joint Conference on Lexical and Computational Semantics</i> , pages 88–102.	Weihaio Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. 2024. Mm-vet: Evaluating large multimodal models for integrated capabilities . <i>Preprint</i> , arXiv:2308.02490.	826
774			827
775			828
776			829
777			830
778			831
779	Hao Shao, Shengju Qian, Han Xiao, Guanglu Song, Zhuofan Zong, Letian Wang, Yu Liu, and Hongsheng Li. 2024. Visual cot: Advancing multi-modal language models with a comprehensive dataset and benchmark for chain-of-thought reasoning. In <i>The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track</i> .	Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhui Chen. 2024. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In <i>Proceedings of CVPR</i> .	833
780			834
781			835
782			836
783			837
784			838
785			839
786	Amanpreet Singh, Vivek Natarjan, Meet Shah, Yu Jiang, Xinlei Chen, Devi Parikh, and Marcus Rohrbach. 2019. Towards vqa models that can read. In <i>Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition</i> , pages 8317–8326.	Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun Lin, Ziyu Guo, Pengshuo Qiu, Aojun Zhou, Pan Lu, Kai-Wei Chang, Peng Gao, and Hongsheng Li. 2024a. Mathverse: Does your multi-modal llm truly see the diagrams in visual math problems? <i>Preprint</i> , arXiv:2403.14624.	840
787			841
788			842
789			843
790			844
791	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. 2023. Stanford alpaca: An instruction-following llama model.	Yuanhan Zhang, Kaiyang Zhou, and Ziwei Liu. 2023. What makes good examples for visual in-context learning? <i>Advances in Neural Information Processing Systems</i> , 36:17773–17794.	845
792			846
793			847
794			848
795	Gemini Team. 2024. Gemini: A family of highly capable multimodal models . <i>Preprint</i> , arXiv:2312.11805.		849
796			850
797	Qwen Team. 2025. Qwen2.5-vl .	Yuechen Zhang, Shengju Qian, Bohao Peng, Shu Liu, and Jiaya Jia. 2024b. Prompt highlighter: Interactive control for multi-modal llms. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 13215–13224.	851
798	James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. Fever: a large-scale dataset for fact extraction and verification. <i>arXiv preprint arXiv:1803.05355</i> .	Ge Zheng, Bin Yang, Jiajin Tang, Hong-Yu Zhou, and Sibeil Yang. 2023. Ddcot: Duty-distinct chain-of-thought prompting for multimodal reasoning in language models. <i>Advances in Neural Information Processing Systems</i> , 36:5168–5191.	852
799			853
800			854
801			855
802	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> .		856
803			857
804			858
805			859
806			860
807			861
808	Alex Wang, Kyunghyun Cho, and Mike Lewis. 2020. Asking and answering questions to evaluate the factual consistency of summaries. <i>arXiv preprint arXiv:2004.04228</i> .		
809			
810			
811			
812	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In <i>Proceedings of the 36th International Conference on Neural Information Processing Systems, NIPS '22</i> , Red Hook, NY, USA. Curran Associates Inc.		
813			
814			
815			
816			
817			
818			
819	Penghao Wu and Saining Xie. 2024. V?: Guided visual search as a core mechanism in multimodal llms. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 13084–13094.		
820			
821			
822			
823			

862 A Benchmark Details

863 This appendix provides a comprehensive overview of the MMIR benchmark. It details the dataset curation
864 process, including error category definitions, the synthetic inconsistency generation mechanism, the
865 auto-verification and human validation processes, and the task prompts for evaluation. These details are
866 intended to facilitate reproducibility and provide clarity on the inner workings of MMIR.

867 A.1 Inconsistency Error Category Definitions

868 The MMIR benchmark employs five pre-defined error categories. These categories are designed based on
869 semantic guidelines so that the generator model can propose diverse and generalizable inconsistencies
870 without being tied to any specific artifact type.

- 871 • **A. Factual Contradiction**

872 Direct conflict between two or more elements (text–text, text–image, or image–image) within the
873 modified content.

874 *Example (Text–Text): The product title says “Caffeinated,” while the description states “Caffeine-*
875 *free.”*

876 *Example (Text–Image): The image shows a green tea bag, but the accompanying text describes a*
877 *“fruit infusion.”*

- 878 • **B. Identity Misattribution**

879 Mislabeling of entities (objects, locations, brands, people) that conflict with other elements.

880 *Example: A product lists “Country of Origin: China” while the manufacturer is described as*
881 *“Elmwood Inn (USA).”*

- 882 • **C. Contextual Mismatch**

883 Tonal, thematic, or situational incompatibility between elements.

884 *Example: A celebratory image of diplomats shaking hands is paired with an article about violent*
885 *clashes.*

- 886 • **D. Quantitative Discrepancy**

887 Numerical or statistical inconsistencies between elements.

888 *Example: A graph labeled “50% growth” shows flat bars.*

- 889 • **E. Temporal/Spatial Incoherence**

890 Implied timelines, dates, or spatial relationships that are impossible or conflicting.

891 *Example: A map labeled “North America” depicts landmarks from Europe.*

892 These definitions serve as guidelines during the synthetic inconsistency generation process, ensuring
893 that the proposed errors are semantically meaningful and cover a broad spectrum of potential real-world
894 mistakes.

895 A.2 Generator Model and Self-Evaluation Loop

896 A.2.1 Generator Model Prompt

897 To create adversarial examples, the generator model (o1, 1217) is provided with rich context consisting of
898 the annotated artifact A_i^{SOM} and its set of elements E_i . The task prompt includes detailed instructions
899 regarding the types of modifications to propose, along with the following guidelines:

- 900 • **Modification Format:** Each modification must be expressed as:

901 ““Modify [id] [original_content] [new_content]””

For image fields, the original content includes the full details (e.g., URL), and the new content is a caption starting with "Image, description: ". For text fields, the new content should be of similar length to the original.

- **Error Categories:** The generator must propose one modification per error category. If it cannot propose an inconsistency for a given category, it may skip that category.

The generator output is structured as:

$$P_m = \{ \text{edit}_m, \text{GT}_m, \text{category}_m, \text{rationale}_m \}$$

where the ground-truth GT_m is defined as:

$$\text{GT}_m \in \{ \text{id}_j \} \cup \{ (\text{id}_j, \text{id}_k) \mid j \neq k \}$$

indicating either a single-element ID (for single-element inconsistencies) or a pair of distinct element IDs (for relational inconsistencies).

A.2.2 Self-Evaluation Loop

We follow a generator-evaluator loop that refines proposals through iterative self-assessment. A simplified Python snippet of the loop function is provided below:

```
1 def loop(client, image_dir, frame_id, task: str, evaluator_prompt: str,
2 generator_prompt: str) -> tuple[str, list[dict]]:
3     """Keep generating and evaluating until requirements are met."""
4     memory = []
5     chain_of_thought = []
6
7     thoughts, result = generate(client, image_dir, frame_id, generator_prompt, task)
8     memory.append(result)
9     chain_of_thought.append({"thoughts": thoughts, "result": result})
10
11     loop_count = 1
12     while True:
13         all_pass = True
14         evaluation, feedback = evaluate(client, image_dir, frame_id,
15 evaluator_prompt, result, task)
16         for eval_line in evaluation.split("\n"):
17             if eval_line.strip() != "PASS":
18                 all_pass = False
19                 break
20
21         if all_pass or loop_count == 2:
22             return result, evaluation
23
24         context = "\n".join([
25             "Previous attempts:",
26             *[f"- {m}" for m in memory],
27             f"\nFeedback: {feedback}"
28         ])
29         thoughts, result = generate(client, image_dir, frame_id, generator_prompt,
30 task, context)
31         memory.append(result)
32         chain_of_thought.append({"thoughts": thoughts, "result": result})
33         loop_count += 1
```

In this loop, the generator produces proposals which are then evaluated against the following criteria (as specified in the evaluator prompt):

- **Category Compliance:** The edit must match the intended error category.
- **Atomic Modification:** Exactly one inconsistency should be introduced.
- **Visual Consistency:** The modified screenshot must visibly reflect the error without relying on external context.

- **Element Validity:** The referenced element IDs must exist in the artifact.

Only proposals receiving a "PASS" in the evaluation are retained. The loop iterates until either all criteria are met or a maximum of two iterations is reached.

A.2.3 Prompt details for generator-evaluator proposal generation framework

This is the task prompt as input to the o1 generator model.

```
1 task_prompt = f"""
2 <user input>
3 Your task is to modify a {category_str} to create inconsistency. For each given
4 category of inconsistency, you will propose a modification action that
5 introduces the inconsistency in the modified {category_str}.
6 Here's the information you'll have:
7 Screenshot of the current {category_str}: This is a screenshot of the {category_str},
8 with each editable element assigned a unique numerical id. Each bounding box
9 and its respective id share the same color.
10 The Observation, which lists the IDs of all editable elements on the current {
11 category_str} with their content, in the format [id] [tagType] [content],
12 separated by "\n". Each id is mapped with the id in the screenshot. tagType is
13 the type of the element, such as button, link, or textbox. For example, "[21] [
14 SPAN] [Add to Wish List]" means that there is a span with id 21 and text content
15 'Add to Wish List' on the current {category_str}. "[23] [IMG] [Image,
16 description: a beige powder on a white background, url: http://localhost:7770/
17 media/catalog/product/cache/829a59e57f886f8cf0598ffca4f8a940/B/0/B074DBMG66.0.
18 jpg]" means that there is an image on the current screen with id 23, with a
19 description of the image and its url specified.
20 Here are the categories of errors you can introduce:
21 A. Factual Contradiction - Direct conflict between two or more elements (text-text,
22 text-image, or image-image). For example, The product title says "Caffeinated,"
23 while the description states "Caffeine-free." Another example: The image shows
24 a green tea bag, but the text describes a "fruit infusion."
25 B. Identity Misattribution - Mislabeling of entities (objects, locations, brands,
26 people) that conflict with other elements. Example: Product "Country of Origin:
27 China" contradicts manufacturer info "Elmwood Inn (USA)."
28 C. Contextual Mismatch - Tonal, thematic, or situational incompatibility between
29 elements. Example: A celebratory image of diplomats shaking hands paired with an
30 article about violent clashes.
31 D. Quantitative Discrepancy - Numerical or statistical inconsistencies between
32 elements. Example: A graph labeled "50% growth" shows flat bars.
33 E. Temporal/Spatial Incoherence - Implied timelines, dates, or spatial relationships
34 that are impossible or conflicting. Example: A map labeled "North America"
35 depicts landmarks from Europe
36 Here are the rules for the modification action:
37 The modification action you can propose to introduce inconsistency must be in the
38 format of "Modify [id] [original_content] [new_content]": This action proposes
39 to edit the original field assigned with the id to the new content to introduce
40 inconsistency. If you propose to modify an image field, the [original_content]
41 field should include the full content from observation including the url; the [
42 new_content] field should be a caption describing the updated image, starting
43 with "Image, description: ", no url needed. If you propose to modify a text
44 field, the new content string should be about the same length as the original
45 text field. For each inconsistency category, you should try to propose a
46 modification action that introduces an inconsistency in that category. If you
47 can't find a way to introduce an inconsistency in a category, you can skip it.
48 Prioritize proposing edits on text fields over image fields.
49 Generate the response in the correct format. For each inconsistency, the format
50 should be:
51 <proposal>
52 <cat>[A-E]</cat> <-- Category letter
53 <ele>[ID1, ID2]</ele> <-- Conflicting element IDs
54 <mod>Modify [ID] [Original Content] [New Content]</mod> <-- Modification plan
55 <rationale>Visible conflict explanation</rationale> <-- Visual verification
56 </proposal>
57 </user input>
```

27	"""	1022
	These are prompts for the generator and evaluator model.	1023
1	evaluator_prompt = """	1024
2	Evaluate the following proposals one by one for:	1025
3	1. Category Compliance: Introduced inconsistency matches the category definition (A-E)	1026
4	2. Atomic Modification: Introduce EXACTLY ONE inconsistency without side effects	1027
5	3. Visual Consistency: Conflict visible in the modified screenshot (with NO reliance on original page knowledge or external context)	1028
6	4. Element Validity: Conflict IDs exist in observations	1029
7		1030
8	You should be evaluating only and not attempting to solve the task.	1031
9	For each proposal, only output "PASS" if all criteria are met and you have no further suggestions for improvements.	1032
10	Output your evaluation concisely in the following format.	1033
11		1034
12	<evaluation>	1035
13	PASS, NEEDS_IMPROVEMENT, or FAIL <-- For each proposal	1036
14	</evaluation>	1037
15	<feedback>	1038
16	What needs improvement and why. <-- For proposals that need improvement	1039
17	</feedback>	1040
18	"""	1041
19		1042
20	generator_prompt = """	1043
21	Your goal is to complete the task based on <user input>. If there are feedback from your previous generations, you should reflect on them to improve proposals that NEEDS_IMPROVEMENT or FAIL. Leave the PASS proposals as they are.	1044
22		1045
23		1046
24	Output your answer concisely in the following format:	1047
25		1048
26	<thoughts>	1049
27	[Your understanding of the task and feedback and how you plan to improve]	1050
28	</thoughts>	1051
29		1052
30	<response>	1053
31	[Your response here]	1054
32	</response>	1055
33	"""	1056

A.3 Auto-Verification and Editing Process 1063

Following proposal generation, an auto-verification step filters the proposals based on format and backend constraints. Specifically: 1064
1065

- **Edit Format Verification:** The system uses a regular expression to ensure that each proposed edit adheres to the required format: "Modify [id] [old_content] [new_content]". 1066
1067
- **Element Matching:** For web-sourced artifacts, the proposal's element ID is used to locate the corresponding element and its bounding box in the metadata. The system checks that both the content and bounding box match an editable element in the HTML/PPTX structure. For image edits, the new content (a caption) is cross-referenced against an MSCOCO image database to verify its appropriateness. 1068
1069
1070
1071
1072

Proposals that pass these checks are automatically saved for further processing. 1073

For web pages, we use the CDP to perform edit: 1074

1	# text edit	1075
2	client.send(1076
3	"Runtime.callFunctionOn",	1077
4	{	1078
5	"objectId": object_id,	1079
6	"functionDeclaration": f"function() {{ this.nodeValue = '{new_content}'; }}"	1080
	,	1081
		1082

```

1083     "arguments": [],
1084     "returnByValue": True
1085 }
1086 )
1087 # image edit
1088 with open(new_content, "rb") as image_file:
1089     img = Image.open(image_file)
1090     new_image_width, new_image_height = img.size # get original width and height
1091     for resizing
1092     aspect_ratio = new_image_width / new_image_height
1093     if w / h > aspect_ratio:
1094         w, h = w, int(w / aspect_ratio)
1095     else:
1096         w, h = int(h * aspect_ratio), h
1097     img = img.resize((w, h), Image.Resampling.LANCZOS)
1098     buffer = BytesIO()
1099     img.save(buffer, format="JPEG")
1100     buffer.seek(0)
1101     base64_image = base64.b64encode(buffer.read()).decode("utf-8")
1102     new_image = f"data:image/jpeg;base64,{base64_image}"
1103 client.send(
1104     "Runtime.callFunctionOn",
1105     {
1106         "objectId": object_id,
1107         "functionDeclaration": f"""
1108             function() {{
1109                 this.src = '{new_image}';
1110             }}
1111             """
1112         "arguments": [],
1113         "returnByValue": True
1114     }
1115 )

```

For Zenodo presentation, we use the python-pptx library:

```

1117
1118
1119 # text edit
1120 if target_shape.has_text_frame: # text edit
1121     text_frame = target_shape.text_frame
1122     for paragraph in text_frame.paragraphs:
1123         for run in paragraph.runs:
1124             if edit_info["old_content"] in run.text:
1125                 try:
1126                     run.text = run.text.replace(edit_info["old_content"], edit_info[
1127                         "new_content"])
1128                     success = True
1129                     break
1130                 except:
1131                     success = False
1132 # image edit
1133 left, top, orig_width, orig_height = target_shape.left, target_shape.top,
1134     target_shape.width, target_shape.height
1135 pic = target_shape._element
1136 pic.getparent().remove(pic)
1137 new_image_path = f"{coco_image_dir}/{edit_info['new_img_path']}"
1138 with Image.open(new_image_path) as img:
1139     new_width, new_height = img.size
1140 new_aspect = new_width / new_height
1141 orig_aspect = orig_width / orig_height
1142 if new_aspect > orig_aspect:
1143     scaled_width = orig_width
1144     scaled_height = int(scaled_width / new_aspect)
1145 else:
1146     scaled_height = orig_height
1147     scaled_width = int(scaled_height * new_aspect)
1148 new_left = left + (orig_width - scaled_width) // 2
1149 new_top = top + (orig_height - scaled_height) // 2
1150 try:
1151     slide.shapes.add_picture( # Add the new image in the same location and size
1152         new_image_path, new_left, new_top, scaled_width, scaled_height

```

```

33 )
34 success = True
35 except:
36 success = False

```

1153
1154
1155
1156

A.4 A Qualitative Example

1158

Input image to both settings



SoM-annotated image for referecne



Ground-truth answer to both settings: [6]

Question in the open-ended setting

You are given a slide that contains inconsistency/error. Identify which element(s) pose semantic inconsistency/error.

Your answer must follow the following format:

```

<ans>
"element 1\" or \"[element 1],[element 2]\" <-- detailed description of one/a pair of problematic element(s), including type (image/text), content and location
</ans>

<reason>
A brief rationale explaining your choice.
</reason>

```

Model answer in <ans> tags -> IDs mapped by LLM as judge

o1: "Rocky Mountain Range" (text, top-right label on the map) -> [6] ✓
4o: "Rocky Mountain Range" <-- text, label located on the top-right portion of the map image -> [6] ✓
Qwen2.5-VL: "Rocky Mountain Range" or "Canonical value for IFR-overflow: 1 Sv" -> [4, 6] ✗
LLaVA-NeXT: "Rocky Mountain Range" -> [6] ✓
InternVL2.5: "[element 1]" Is at the location of the text "Canonical value for IFR-overflow: 1 Sv." The issue here is semantic, as it does not provide specific meaning regarding what "1 Sv" stands for. Sv, or steradian per second, is a unit of luminous intensity, making this error in context since it is unrelated to the flow values. -> [4] ✗
Phi-3.5-Vision: "[element 1]" -> [1] ✗

Question in the multiple-choice setting

You are given a slide that contains inconsistency/error. Identify which element(s) from the given options pose semantic inconsistency/error. In the given slide, each element in the given options is assigned a unique numerical id and highlighted with bounding box sharing the same color.

Options:

- [1] Wide text span at top-center: 'Overflow through the Western Valley of the Iceland-Faroe Ridge is negligible'
- [2] Wide text span at top-center: 'Bogi Hansen, Karin M. H. Larsen, Steffen Olsen, Detlef Quadfasel, Kerstin Jochumsen, Svein Østerhus'
- [3] Image, description: the location of arctic mediterranean waters at center-center
- [4] 451.79px-wide text span at bottom-center: 'Canonical value for IFR-overflow: 1 Sv'
- [5] Compact text span at center-center: 'IFR'
- [6] 651.41px-wide text span at center-right: 'Rocky Mountain Range'
- [7] 343.31px-wide text span at center-center: 'WV-overflow:'

Your answer must follow the following format:

```

<ans>
"[ID1]" or \"[ID1],[ID2]\" <-- one problematic element ID, or two problematic element IDs separated by comma
</ans>

<reason>
A brief rationale explaining your choice.
</reason>

```

Model answer in <ans> tags

o1: [6] ✓
4o: [6] ✓
Qwen2.5-VL: [3] ✗
LLaVA-NeXT: [1] ✗
InternVL2.5: [6] ✓
Phi-3.5-Vision: [1,3] ✗

Figure 7: A test sample with model responses under the two main settings in MMIR: open-ended and multiple-choice.

B Model Application Details

1159

Here are the generation methods for the open-sourced models.

1160

For **o1** and **GPT-4o**, we utilized the API following API guidelines available at <https://platform.openai.com/docs/models#gpt-4o>.

1161

1162

For **Qwen2.5-VL**, we implemented the 7B version following the official repository: <https://github.com/QwenLM/Qwen2.5-VL>.

1163

1164

1165 For **LLaVA-NeXT**, we followed the implementation from [https://github.com/LLaVA-VL/](https://github.com/LLaVA-VL/LLaVA-NeXT)
1166 [LLaVA-NeXT](https://github.com/LLaVA-VL/LLaVA-NeXT).

1167 For **InternVL2.5** we implemented the 8B version at <https://github.com/OpenGVLab/InternVL>.

1168 For **Phi-3.5-Vision** we implemented the 4B version at [https://github.com/instill-ai/models/](https://github.com/instill-ai/models/tree/main/phi-3-5-vision)
1169 [tree/main/phi-3-5-vision](https://github.com/instill-ai/models/tree/main/phi-3-5-vision).

1170 **C Data Release**

1171 We will publicly release a comprehensive dataset that includes the artifacts and question-answer pairs
1172 in both the open-ended and multiple-choice settings. The licensing terms for the artifacts will follow
1173 those set by the respective dataset creators, as referenced in this work, while the curated artifacts will be
1174 provided under the MIT License. Additionally, our release will include standardized evaluation protocols,
1175 and evaluation scripts to facilitate rigorous assessment. The entire project will be open-sourced, ensuring
1176 free access for research and academic purposes.