# 🦁 CHIMERA: DIAGNOSING SHORTCUT LEARNING IN VISUAL-LANGUAGE UNDERSTANDING

**Ziheng Chi**[*1]  **Yifan Hou**[*1]  **Chenxi Pang**[2]  **Shaobo Cui**[3]  **Mubashara Akhtar**[1]  **Mrinmaya Sachan**[1]

{zihchi, yifan.hou, mubashara.akhtar, mrinmaya.sachan }@inf.ethz.ch

chenxipang@google.com, shaobo.cui@epfl.ch

ETH Zürich[1], Google DeepMind[2], EPFL[3]

## ABSTRACT

Visual language, exemplified by diagrams, conveys symbolic information in a visual format rather than a linear stream of words, making them especially challenging for AI models to process. While recent evaluations suggest that vision-language models (VLMs) perform well on diagram-related benchmarks, their reliance on knowledge, reasoning, or modality shortcuts raises concerns about whether they genuinely understand and reason over diagrams. To address this gap, we introduce CHIMERA, a comprehensive test suite comprising 7,500 high-quality diagrams sourced from Wikipedia; each diagram is annotated with its symbolic content represented by semantic triples along with multi-level questions designed to assess four fundamental aspects of diagram comprehension: *entity recognition*, *relation understanding*, *knowledge grounding*, and *visual reasoning*. We use CHIMERA to measure the presence of three types of shortcuts in visual question answering: (1) the *visual-memorization shortcut*, where VLMs rely on memorized visual patterns; (2) the *knowledge-recall shortcut*, where models leverage memorized factual knowledge instead of interpreting the diagram; and (3) the *Clever-Hans shortcut*, where models exploit superficial language patterns or priors without true comprehension. We evaluate 15 open-source VLMs from 7 model families on CHIMERA and find that their seemingly strong performance largely stems from shortcut behaviors – visual-memorization shortcuts have slight impact, knowledge-recall shortcuts play a moderate role, and Clever-Hans shortcuts contribute significantly. These findings expose critical limitations in current VLMs and underscore the need for more robust evaluation protocols that benchmark genuine comprehension of complex visual inputs (e.g., diagrams) rather than question-answering shortcuts.[1]

## 1 INTRODUCTION

Visual language enables communication through structured visual elements such as symbols, icons, and spatial relationships. Diagrams are a fundamental form of visual language, used in domains such as science, education, and engineering to convey complex information compactly and intuitively (Greenspan & Shanker, 2009; Anderson et al., 2011; Zdebik, 2012; Marriott & Meyer, 2012). Comprehending diagrams requires a wide range of abilities, from basic visual recognition to complex reasoning, making it a particularly challenging task for AI systems (Seo et al., 2014; Kembhavi et al., 2016; Lu et al., 2021). Understanding how vision-language models (VLMs) interpret and reason over diagrams is thus both conceptually challenging and practically important: it reveals current limitations and guides the design of future multimodal systems (Li, 2023). While recent VLMs have shown impressive results on diagram-related evaluation (Xue et al., 2024; Liu et al., 2024b; Bai et al., 2025; Meta, 2024; Google, 2025; Agrawal et al., 2024; Microsoft, 2025), these works often focus narrowly on performance and lack a structured evaluation of the step-by-step reasoning process. More importantly, they do not systematically address shortcut behaviors, such as relying on memorized

---

*   Equal contribution

[1]Our code and data are publicly available.

Figure 1: An example from CHIMERA showcasing three modalities (visual, semantic, and textual modality) and four evaluation tasks: entity recognition, relation understanding, knowledge grounding, and visual reasoning.

patterns or language priors that can inflate scores without true comprehension (Goyal et al., 2017; Bleeker et al., 2024; Hou et al., 2025). This highlights the need for a test suite that not only measures accuracy, but also disentangles how models comprehend diagrams, from basic recognition to abstract reasoning, while controlling for potential shortcuts.

Motivated by semiotics, the study of how meaning is conveyed through signs, we represent the diagram content using semantic triples, enabling consistent alignment across three modalities: the original diagram, i.e., *visual modality*; visualized triples, i.e., *semantic modality*; and sentences, i.e., *textual modality*. Building on Peirce's theory of semiosis, which models interpretation as linking signs to objects through reasoning (Peirce, 1935; Morris, 1938), we frame diagram comprehension as a four-stage process: *entity recognition*, *relation understanding*, *knowledge grounding*, and *visual reasoning*. This structured perspective reflects the key cognitive steps required for VLMs to move from surface recognition to deeper multimodal understanding.

We introduce CHIMERA, a fine-grained test suite designed to evaluate the abilities of VLMs to interpret and reason about diagrams with meticulous annotations of both diagram content and evaluation questions. To construct the dataset, we collect diagram images from Wikipedia (Burns et al., 2023), and filter out unsuitable images such as natural photographs using MetaCLIP (Xu et al., 2024a). We then employ VLMs to annotate each diagram with its domain and type, further removing low-quality samples. For semantic content, we use Gemini (Google, 2024) to describe the key information conveyed in the diagram and use it to generate semantic triples and four levels of questions. To ensure annotation reliability, we perform multi-round consistency checks under different settings to filter out ambiguous or inconsistent diagrams. In total, CHIMERA comprises 7,500 diagrams (6,000 training, 1,500 test), each enriched with semantic triples and four levels of questions—targeting entity recognition, relation understanding, knowledge grounding, and visual reasoning (see Fig. 1). Notably, although Wikipedia may overlap with VLM training data, this design choice is intentional: by leveraging commonly seen images, CHIMERA is positioned to *expose shortcut learning behaviors*. In contrast to using novel or out-of-distribution diagrams, which may simply cause models to fail, our test suite reveals *how current models succeed through superficial cues rather than genuine understanding*.

Then, we revisit the shortcut behaviors in visual question answering (VQA) under the diagram comprehension scenario, and categorize them into three distinct types. First, models could rely on image priors, memorizing visual information from training data and using it directly during inference, without genuinely understanding the diagram content (Jayaraman et al., 2024; Li et al., 2024). We refer to this as the *visual-memorization shortcut*. Second, models could exploit language priors, which we further divide into two subtypes. Given that diagrams often convey factual or domain-specific knowledge, a model could simply recognize high-level visual patterns and rely on pre-trained language knowledge to answer the question without actually understanding the diagram (Hou et al., 2025; Zang et al., 2024). We refer to this as the *knowledge-recall shortcut*. In addition to that, models can also learn to exploit superficial patterns in the language of the questions or answer options, arriving at correct answers without using the visual input at all (Goyal et al., 2017; Bleeker et al.,

2024). We call this behavior the *Clever-Hans shortcut*, drawing analogy to the phenomenon where models appear to perform well by exploiting spurious cues rather than genuine understanding.

Using CHIMERA, we evaluate 15 open-source VLMs from 7 model families to analyze their core abilities and behavioral patterns in diagram comprehension. We compare model performance on visual modality and semantic modality. Surprisingly, VLMs perform slightly better on visually complex real diagrams than on the simpler, cleaner semantic graphs. This counterintuitive result suggests that the *visual-memorization shortcut* exists. Models could exploit memorized visual patterns from pretraining, but their impact is **slight**. The *knowledge-recall shortcut* is unlikely to affect entity recognition, but it is more plausible in the remaining three tasks, which are more knowledge-intensive. However, our results show that VLMs perform obviously worse on entity recognition than on the other three tasks, despite it being the simplest and most fundamental. This performance gap supports that the knowledge-recall shortcut occurs **moderately** in the latter tasks. Given that entity recognition is relatively free from knowledge-based shortcuts, we investigate the *Clever-Hans shortcut* in this task. Specifically, we evaluate VLMs without providing the diagram, using only the question and answer options. Surprisingly, some models could even achieve comparable performance as when the diagram is present, suggesting that they rely heavily on spurious linguistic patterns in the prompt. This provides strong evidence that the Clever-Hans shortcut is **significant**.

These findings reveal that the seemingly strong diagram reasoning performance of current VLMs is largely driven by shortcut behaviors rather than genuine comprehension. Among the three types of shortcuts, the Clever-Hans shortcut is the most severe. Our analysis exposes fundamental limitations in current open-source VLMs and underscores the need for more robust evaluation frameworks. Achieving human-level visual understanding remains a long and challenging journey.

## 2    CHIMERA

In this section, we first outline the test suite design, followed by describing the construction process in detail and presenting the results of human evaluation.

### 2.1    DESIGN FOUNDATIONS: SEMIOTICS AND SEMIOSIS

We motivate our test-suite design, deriving three modalities and four semiosis-aligned tasks from semiotic theory, and show how an in-domain setup exposes and disentangles shortcut behaviors.

**Semiotic Foundation: Three Modalities for Probing Shortcut Use.**    Our test suite is grounded in semiotic theory, the study of how meaning is constructed and interpreted through signs and representations (Peirce, 1935; Morris, 1938; Cullum-Swan & Manning, 1994). According to Charles Sanders Peirce, signs are broadly categorized into three types: icons (representing meaning through visual resemblance), symbols (through learned or conventional associations), and indexes (through direct causal links, e.g., smoke signals fire) (Yakin & Totu, 2014). While diagrams may not include all sign types, many flexibly use combinations of icons, symbols, and indexes to construct meaning.

Inspired by this semiotic framework, we design three modalities in our test suite that recast the same diagram content through different representational lenses. The *visual modality* presents the original diagram image; the *semantic modality* transforms iconic signs into symbolic form by representing the diagram as a structured graph of semantic triples; and the *textual modality* further abstracts this information by expressing the triples as natural language statements, converting indexical or context-dependent cues into symbolic language. Each modality conveys equivalent content but varies in surface cues and representational abstraction (see Fig. 1).

This design enables us to probe whether models genuinely understand diagram content or rely on modality-specific shortcuts. For example, if a model performs well only on the visual modality, but not on the equivalent semantic or textual inputs, it may suggest visual memorization or pattern-matching, rather than true comprehension. In contrast, consistent performance across modalities would indicate deeper, format-invariant understanding. Grounding the test suite in semiotics thus provides not only a cognitively informed structure, but also a principled way to evaluate modality alignment and shortcut behaviors in these models.

**Semiosis Foundation: Four Interpretive Processes for Diagnosing Shortcut Use.** In addition to representational modality, our test suite design is also guided by Peirce's theory of semiosis, a dynamic, triadic process by which a sign (e.g., a diagram) represents an object (real-world referent) and produces an interpretant (meaning in the interpreter's mind) (Peirce, 1935; Peirce et al., 1992). This process unfolds in four stages: first recognizing entity objects, then interpreting relations among them, grounding them in related knowledge, and finally drawing inferences via reasoning. Each stage reflects a core cognitive function in human diagram understanding.

We mirror this process with four evaluation tasks in our test suite (Fig. 2). *Entity recognition* corresponds to identifying visual elements and mapping them to real-world objects. *Relation understanding* assesses the ability to extract and interpret structural relationships between entities. *Knowledge grounding* tests whether the model can connect the diagram content to broader conceptual or domain knowledge. *Visual reasoning* targets the highest level of abstraction, integrating grounded elements to infer or derive conclusions.



Figure 2: Diagram comprehension process inspired by semiosis.

This decomposition enables fine-grained diagnosis of reasoning shortcuts. For instance, success on grounding and reasoning tasks without corresponding recognition and relation understanding may reveal reliance on background knowledge rather than visual interpretation. By aligning test suite tasks with the semiosis process, we isolate where models succeed by reasoning versus when they default to shortcut strategies.
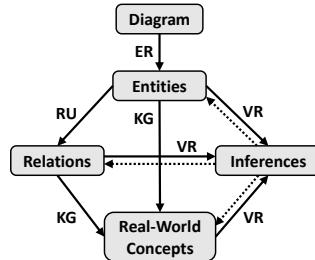
**Shortcut Exposure through In-Domain Design.** While recent studies report strong performance of VLMs on diagram-related tasks (Masry et al., 2022; Wang et al., 2024a; Lu et al., 2024), others reveal their brittleness in complex visual reasoning or generalization to new formats (Miyai et al., 2024; Sim et al., 2025; Hou et al., 2025). A key factor underlying this discrepancy is the presence or absence of shortcuts in the test suite design. To intentionally expose such shortcuts, we construct CHIMERA from Wikipedia diagrams, a source heavily represented in VLM pretraining. This choice increases the chance that models can exploit memorized content, language priors, or pattern-based biases. Far from being a flaw, this setup is critical for our analysis: if models fail even with such familiar inputs, it strongly indicates deeper reasoning limitations. If they succeed, our modality-controlled ablations and task-level consistency checks help determine whether that success is genuine or shortcut-driven.

In summary, the structure of CHIMERA, grounded in semiotics and semiosis, not only reflects how humans understand diagrams, but also enables rigorous analysis of when and how VLMs fail to replicate that process. This design lays the foundation for systematically dissecting and diagnosing shortcut learning behaviors in visual language understanding.

## 2.2 TEST SUITE CONSTRUCTION

We build our test suite data in three stages: diagram cleaning, tagging, and annotation (semantic triples and question-answer pairs). An illustration of our construction pipeline is given in Fig. 3.

**Diagram Cleaning.** To build our test suite, we extract images from WikiWeb2M (Burns et al., 2023), a large-scale corpus of English Wikipedia pages. Since many images are irrelevant to diagrams, we apply a filtering process using MetaCLIP (Xu et al., 2024a), combining one positive prompt and six negative prompts. Only images consistently classified as diagrams are retained, resulting in approximately 100k candidate images. Details are provided in § B.1.

**Diagram Tagging.** Diagrams vary widely in type and domain due to their role in knowledge transfer. To structure our test suite, we use VLMs (Molmo and LLaMA) to tag each diagram by its type and subject domain (Fig. 3). After aggregating four annotations per image, we group the most common tags into 12 categories across two groups: statistical (e.g., bar chart, line graph) and scientific (e.g., biology, physics). Only diagrams with consistent tags are retained, yielding around 60k images. Full tagging prompts and category details are provided in § B.2.
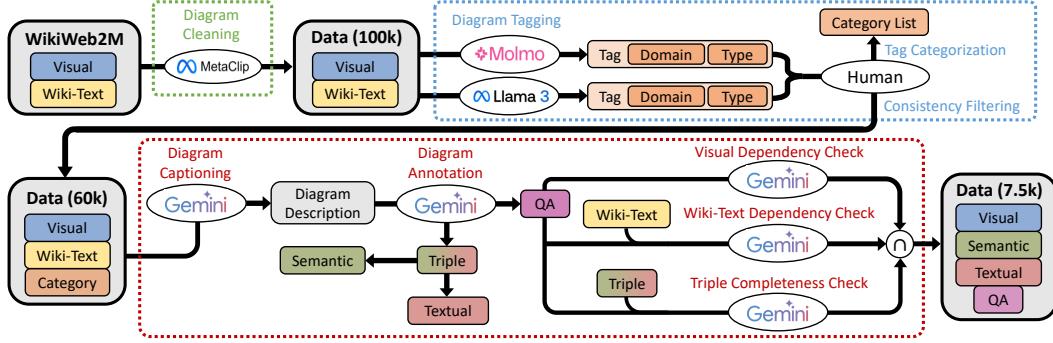
Figure 3: Overview of our test suite data construction pipeline. First, starting from the WikiWeb2M dataset, we use MetaCLIP to remove non-diagram images, resulting in 100k diagrams. Second, we apply Molmo and LLaMA for tagging, and then derive a fixed category list and filter inconsistent results, yielding 60k diagrams. Third, we prompt Gemini to caption diagrams and annotate semantic triples and QA pairs. We then apply three rounds of quality checks, producing a final dataset containing 7.5k high-quality diagrams.

**Diagram Annotation.** We posit that the information and knowledge that a diagram conveys can be naturally formalized by a knowledge graph, that is, a set of *semantic triples* (Lassila & Swick, 1999), where each triple contains a head entity, a relation, and a tail entity. In addition to using the diagram as the information carrier (i.e., *visual modality*), we can also represent the information directly by visualizing the semantic triples or transforming it to textual sentences.

Our test suite includes two core parts of annotations: semantic triples and question–answer (QA) pairs (Fig. 3). To ensure high-quality and consistent annotation, we adopt a two-step pipeline using Gemini-2.0-Flash (Google, 2024) as the annotation backbone. In the first step, we prompt the model to generate a detailed description of each input diagram. These prompts are tailored to different diagram groups and enriched with in-context examples to encourage accurate and specific descriptions. To reduce hallucinations and improve factual grounding, we also provide the associated Wikipedia text to the model as the supplementary input.

In the second step, we use the generated descriptions to extract semantic triples and generate QA pairs. To ensure that the resulting annotations are both accurate and visually grounded, we apply a three-stage consistency check: (1) we discard examples if questions can be answered without the image; (2) we verify that questions remain unanswerable when only Wikipedia text is available; and (3) we confirm that the semantic triples alone are sufficient to answer the questions. Only diagrams that pass all three checks are retained. After filtering, the final test suite comprises 6,000 diagrams for training and 1,500 for testing. All evaluations in this paper are conducted on the test set. Additional details, including prompt templates and filtering criteria, are provided in § B.3.

## 2.3 HUMAN EVALUATION

Despite implementing several statistical verification methods to ensure annotation quality, automatically generated annotations may still lack consistency and accuracy. To further assess the reliability of our test suite, we conduct a round of human evaluation following the automatic annotation process. Unlike the earlier verification, which focused on the independence of Wikipedia text, this evaluation emphasizes the correctness and reliability of the QA annotations. We evaluate each data point along three key dimensions:

- **Visual Dependency**: We assess whether each question truly requires the diagram to be answered, rather than relying on commonsense or background knowledge. An annotation is labeled as *Fully Dependent* if all questions rely on visual content, and *Partially Dependent* if at least one question can be answered without referring to the diagram.

- **QA Correctness**: We evaluate whether the questions are clearly phrased, contextually grounded, and whether the provided answers are correct. Each data point is labeled as *Perfectly Valid* or *Slightly Flawed*, depending on whether any question contains a factual error.

- **Triple Completeness**: We verify whether the annotated semantic triples accurately and sufficiently capture the key information in the diagram. Data points are labeled as *Totally Sufficient* if the triples are complete and correct, and *Marginally Insufficient* if an essential triple is missing or inaccurate.

Table 1: Human evaluation results on 300 diagrams across three dimensions: visual dependency, QA correctness, and triple completeness. Scores reflect the percentage of diagrams rated under each category by four annotators (A, B, C, D), showing overall strong annotation quality with minor variations in strictness.

| Score Ratio (%) | Visual Dependency | | QA Correctness | | Triple Completeness | |
|---|---|---|---|---|---|---|
| | Fully Dependent | Partially Dependent | Perfectly Valid | Slightly Flawed | Totally Sufficient | Marginally Insufficient |
| Annotator A | 85.3 | 14.7 | 92.0 | 8.0 | 86.0 | 14.0 |
| Annotator B | 100.0 | 0.0 | 99.3 | 0.7 | 80.7 | 19.3 |
| Annotator C | 78.7 | 21.3 | 87.3 | 12.7 | 70.7 | 29.3 |
| Annotator D | 95.3 | 4.7 | 96.0 | 4.0 | 82.7 | 17.3 |

We evenly sample 20% of the test set (300 diagrams) across categories and assign them to four expert annotators (A, B, C, and D). As shown in Tab. 1, the majority of annotations are consistently rated as *Fully Dependent*, *Perfectly Valid*, and *Totally Sufficient*. While minor differences exist among annotators in terms of strictness, the overall results confirm that the test suite annotations are of high quality and suitable for reliable evaluation.

## 3 DIAGRAM COMPREHENSION EVALUATION

In this section, we first present the overall evaluation results on our test suite. We then delve deeper into a central open question: *how do VLMs actually comprehend complex images such as diagrams?* One hypothesis posits that VLMs achieve genuine understanding, while the alternative suggests that their performance is largely driven by shortcut behaviors. To investigate this, we analyze three typical shortcut types: *visual-memorization shortcut*, *knowledge-recall shortcut*, and *Clever-Hans shortcut* using CHIMERA as a diagnostic tool.

### 3.1 OVERALL EVALUATION

**Experiment Setup.** We evaluate 15 models from 7 model families, covering both academic and industrial models across a range of parameter scales. We select the Qwen2.5-VL (simplified as Qwen) series (3B, 7B, 32B, 72B) (Bai et al., 2025), the LLaMA3.2-Vision-Instruct (simplified as LLaMA) series (11B, 90B) (Meta, 2024), the Gemma3 series (4B, 12B, 27B) (Google, 2025), the LLaVA-1.6 series (7B, 13B, 34B) (Liu et al., 2024b), as well as three standalone models: Pixtral-12B (Agrawal et al., 2024), Phi-4 5.6B (Microsoft, 2025), and BLIP-3 4B (Xue et al., 2024). More details about the model, the evaluation setting (e.g., prompts) can be found in § C.1.

**Overall Results.** We report average accuracy across 15 models in Tab. 2, with detailed results provided in § C.2. Models are evaluated across three input modalities—visual (original diagram), semantic (visualized triples), and textual (sentence-form triples)—and four tasks: entity recognition (ER), relation understanding (RU), knowledge grounding (KG), and visual reasoning (VR). Overall, VLMs perform best with textual inputs across all tasks, while accuracy drops significantly for visual and semantic modalities, revealing clear room for improvement in diagram comprehension.

Table 2: Average accuracy of 15 VLMs on CHIMERA across three input modalities and four tasks.

| Accuracy (%) | ER | RU | KG | VR |
|---|---|---|---|---|
| Visual Modality | 80.6 | 85.8 | 87.7 | 85.7 |
| Semantic Modality | 76.1 | 84.2 | 88.0 | 84.6 |
| Textual Modality | 89.5 | 91.4 | 92.9 | 90.1 |

## 3.2 VISUAL-MEMORIZATION SHORTCUT: DO VLMS ANSWER USING MEMORIZED VISUAL PATTERNS?

With the increasing model capacity, recent studies suggest that VLMs could memorize training data (e.g., diagrams) and rely on this memorized content for inference, rather than genuine comprehension (Jayaraman et al., 2024; Li et al., 2024). We refer to this behavior as the *visual-memorization shortcut*, where a model bypasses reasoning by exploiting memorized visual patterns.

**Experiment Design.** To investigate whether VLMs rely on the visual-memorization shortcut for diagram comprehension, we leverage the multimodal design of CHIMERA. Each diagram in the test suite is annotated with semantic triples, which are visualized as semantic modality inputs, i.e., structured and simplified versions of the original diagrams. Compared to real diagrams (visual modality), semantic graphs eliminate noise and layout ambiguity, offering a clearer path for reasoning.

If a model is not relying on memorized visual patterns, we would expect it to perform worse on real diagrams than on the cleaner, more structured semantic modality. In contrast, if the visual-memorization shortcut is in use, models might perform better on the visual modality, indicating reliance on memorized diagram appearances rather than actual visual reasoning. Additionally, we treat the textual modality (i.e., sentences generated from triples) as an upper-bound reference, since it presents all essential information in the most language-friendly form for VLMs.

**Evaluation Results.** Fig. 4 reports the average accuracy across all tasks and models. Detailed results are in § C.2. As expected, performance on the textual modality is the highest, confirming the language-centric nature of current VLMs. However, a surprising pattern emerges: models perform slightly better on the visual modality than on the semantic modality ($\approx$ 2%), comparing to the gap between textual modality and the visual modality. Despite being more complex and less structured, real diagrams yield better performance than their simplified semantic counterparts. This contradicts the intuition that structured, noise-free inputs should facilitate better reasoning.



Figure 4: Average performance across models and tasks on different modalities. The overall performance on the visual modality is **slightly** better than that on the semantic modality.

**Takeaways.** These results suggest that VLMs do make **slight** use of the visual-memorization shortcut when performing diagram comprehension ($\approx$ 2%). While the relative gap is not large, the fact that models outperform on real diagrams despite their complexity implies some level of visual shortcuts. The shortcut appears limited but measurable, and it could become more pronounced in settings where training and evaluation data overlap.

## 3.3 LANGUAGE SHORTCUTS

In addition to relying on visual memorization, VLMs may also exploit shortcuts derived from the language prior patterns and knowledge embedded in the language modeling component rather than performing genuine multimodal reasoning. We divide such language-based shortcuts into two distinct types: (1) The *knowledge-recall shortcut*, where models retrieve factual or commonsense knowledge from pretraining to answer questions, bypassing the diagram. (2) The *Clever-Hans shortcut*, where models rely on superficial linguistic patterns in questions or answer options, independent of any grounded understanding. In this section, we analyze these two shortcuts in turn.

### 3.3.1 KNOWLEDGE-RECALL SHORTCUT: DO VLMS USE MEMORIZED KNOWLEDGE?

A common form of language-based shortcut is the knowledge shortcut, where VLMs draw on memorized background knowledge or commonsense associations from pretraining instead of interpreting the visual content (Hou et al., 2025; Zang et al., 2024).

**Experiment Design.** To assess the presence of knowledge shortcuts, we analyze VLM performance across the four tasks in CHIMERA: entity recognition (ER), relation understanding (RU), knowledge grounding (KG), and visual reasoning (VR). As the most fundamental and prerequisite step in diagram
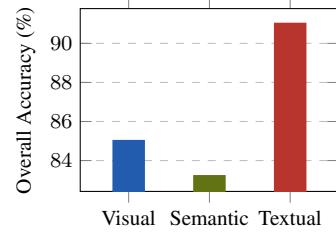
comprehension (Fig. 2), The entity recognition task is highly localized and visual, making it unlikely to benefit from knowledge-recall shortcuts. In contrast, other three tasks involve deeper reasoning and are more likely to draw on factual knowledge stored in the model. Intuitively, if a model engages in genuine visual comprehension, we would expect the highest accuracy on entity recognition, followed by decreasing performance on the more complex tasks. However, if a model performs worse on the recognition but better on other tasks, it suggests a reliance on memorized knowledge rather than true visual understanding, an indicator of knowledge-recall shortcuts.



**(a)** Average performance across all models.   **(b)** Performance of Qwen across four model sizes.
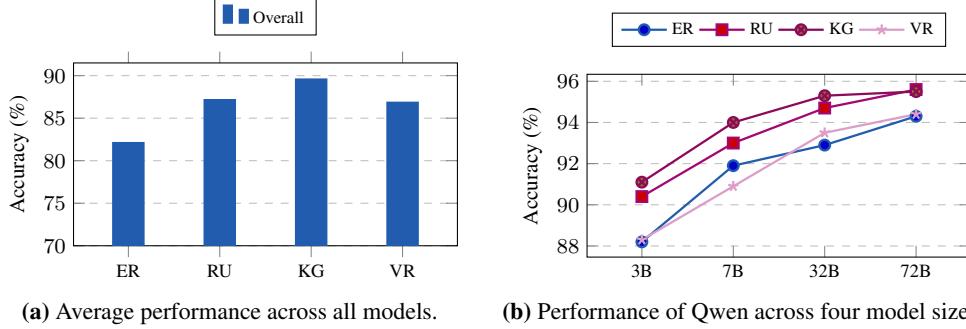
Figure 5: The overall evaluation accuracy for 15 VLMs and the accuracy of four Qwen2.5-VL models on the four tasks. VLMs perform on entity recognition much worse than that on the other three tasks. For Qwen models, larger model is more likely to have smaller gaps between entity recognition and other tasks.

**Quantitative Results.**   As shown in Fig. 5a, VLMs surprisingly perform worst on entity recognition, while achieving higher accuracy on relation understanding, knowledge grounding, and visual reasoning ($\approx 5\%$). This contradicts the intuition that simpler, recognition-level tasks should be easier. The pattern suggests that VLMs rely on memorized knowledge to handle semantically richer tasks, rather than building understanding through visual parsing. Furthermore, as shown in Fig. 5b, this trend holds consistently across the Qwen model family (from 3B to 72B), with larger models exhibiting smaller performance gaps. This indicates that larger VLMs are less likely to be susceptible to knowledge-recall shortcuts. One possible reason is that their larger language backbones contribute more to processing the visual information they perceive, rather than merely expanding the pool of stored knowledge they can draw upon.



Figure 6: Model responses for a diagram of the largest evaluated VLM (i.e., LLaMA-90B). The model fails to recognize the basic, simple elements in the diagram while providing correct answers for more complex questions.

**Qualitative Evidence.**   Fig. 6 illustrates a representative failure case from LLaMA-90B. The model incorrectly classifies a scatter plot as a line graph, i.e., failing in basic visual recognition, yet proceeds to correctly describe complex trends in the data and even offer projections and possible data sources. This behavior reinforces the hypothesis that the model bypasses perception and relies instead on memorized knowledge patterns to perform diagram comprehension.

**Takeaways.**   Both quantitative trends and qualitative examples support the conclusion that knowledge-recall shortcuts occur **moderately** in current VLMs ($\approx 5\%$). These shortcuts are observed across model sizes and tend to be more pronounced in larger models. While they help models answer knowledge-intensive questions, this often comes at the expense of genuine visual comprehension.

8

### 3.3.2 CLEVER-HANS SHORTCUT: DO VLMS RELY ON SUPERFICIAL LANGUAGE PATTERNS?

Another widely observed form of shortcut in visual question answering is the Clever-Hans shortcut, where models exploit superficial patterns in the input text (i.e., the question and answer options), rather than relying on visual input (Goyal et al., 2017; Agrawal et al., 2018; Cadène et al., 2019; Bleeker et al., 2024). This shortcut is particularly insidious because the model can appear accurate by exploiting linguistic regularities, even when the visual input is missing or irrelevant.

**Experimental Design.** To isolate the Clever-Hans shortcut from other language priors (e.g., factual knowledge), we focus on the entity recognition task in CHIMERA. Our earlier analysis shows that this task is less influenced by the knowledge-recall shortcut, making it an ideal case for probing the effects of shallow language pattern exploitation.

We compare model performance under two conditions: (1) the standard setting with access to the original diagram, and (2) a blank-image setting where no visual information is provided. Since each question in CHIMERA is multiple-choice with four options, the expected accuracy from random guessing is approximately 25%. Any significant improvement above this baseline in the absence of visual input suggests the presence of Clever-Hans behavior.
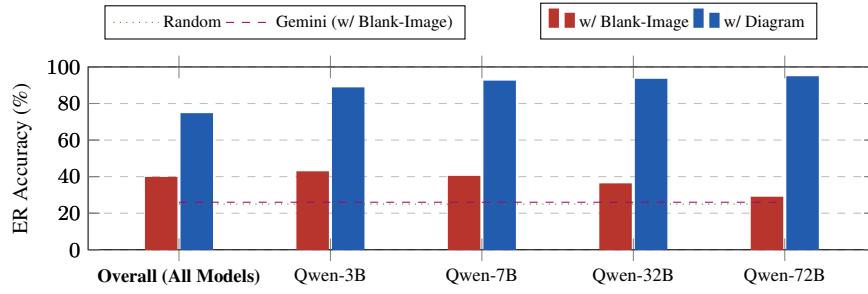


Figure 7: Entity recognition accuracy under normal VQA and blank-image settings. **Overall** is the average value for all 15 models. Results show that Qwen models have strong reliance on language-only cues. Besides, larger models exhibiting slightly less susceptibility to the Clever-Hans shortcut.

**Quantitative Results.** Fig. 7 presents entity recognition accuracy across VLMs under both settings. The performance gap between the two settings reflects the extent to which models rely on language-only cues embedded in the questions and options. Interestingly, we observe that larger models tend to rely less on the Clever-Hans shortcut. For example, Qwen-VL-72B shows a worse performance under the w/ blank-image setting compared to Qwen-VL-3B. This trend suggests that increased model capacity may improve multimodal grounding, making models more reliant on actual visual content.

**Qualitative Results.** Fig. 8 presents a representative example from LLaMA-90B (entity recognition). When the diagram input is removed, the model still generates a confident and contextually reasonable answer by relying solely on the question phrasing and the content of the answer options. Notably, the response lacks any reference to visual content or spatial cues, indicating that the model is not engaging in genuine diagram interpretation. Instead, it is leveraging superficial language patterns, a clear instance of the Clever-Hans shortcut, highlighting its dependence on linguistic biases rather than multimodal understanding.



Figure 8: Response of LLaMA-90B on the entity recognition task. Even without a valid diagram input, the model examines the question and options and makes an educated guess based on superficial language patterns.

**Takeaways.** These results provide strong evidence that Clever-Hans shortcuts are **significant** in open-source VLMs ($\approx$ 15%), particularly among smaller models. Even without valid visual input, models achieve non-trivial accuracy by exploiting linguistic biases. While larger models show some improvement in resisting this behavior, the shortcut remains a significant barrier to robust

multimodal reasoning. Addressing it will require improved training signals, more carefully designed datasets, and evaluation protocols that explicitly discourage reliance on language-only cues.

## 4 CONCLUSION

We introduce CHIMERA, a comprehensive test suite for diagram comprehension in VLMs, with carefully annotated multimodal inputs and multi-level tasks. Unlike prior work, it enables fine-grained analysis across modalities and diagram comprehension stages. Our evaluation of 15 VLMs reveals that much of their success stems from language-based shortcuts, especially Clever-Hans behaviors, rather than genuine diagram understanding. These insights highlight key limitations in current open-source models and offer guidance for building more robust, interpretable, and multimodal systems.

## ETHICS STATEMENT

This work does not involve human subjects, sensitive data, or potentially harmful applications. All datasets used are publicly available and widely adopted in the vision-language and reasoning communities. We follow best practices in data handling, model evaluation, and reproducibility, and adhere to the ICLR Code of Ethics in all aspects of our research.

## REPRODUCIBILITY STATEMENT

All diagrams, semantic triples, and multi-level QA annotations in CHIMERA are released under a permissive license, along with detailed documentation of the data collection and annotation pipeline. To facilitate replication of our experiments, we provide training and evaluation scripts, model prompts, and preprocessing utilities. The benchmark design is modular, making it straightforward to extend to new domains or models. Additionally, we describe implementation details, hyperparameters, and evaluation procedures in the appendix, enabling others to reproduce our reported results.

## REFERENCES

Aishwarya Agrawal, Dhruv Batra, Devi Parikh, and Aniruddha Kembhavi. Don't just assume; look and answer: Overcoming priors for visual question answering. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pp. 4971–4980. Computer Vision Foundation / IEEE Computer Society, 2018. doi: 10.1109/CVPR.2018.00522. URL http://openaccess.thecvf.com/content_cvpr_2018/html/Agrawal_Dont_Just_Assume_CVPR_2018_paper.html.

Pravesh Agrawal, Szymon Antoniak, Emma Bou Hanna, Baptiste Bout, Devendra Chaplot, Jessica Chudnovsky, Diogo Costa, Baudouin De Monicault, Saurabh Garg, Theophile Gervet, Soham Ghosh, Amélie Héliou, Paul Jacob, Albert Q. Jiang, Kartik Khandelwal, Timothée Lacroix, Guillaume Lample, Diego Las Casas, Thibaut Lavril, Teven Le Scao, Andy Lo, William Marshall, Louis Martin, Arthur Mensch, Pavankumar Muddireddy, Valera Nemychnikova, Marie Pellat, Patrick Von Platen, Nikhil Raghuraman, Baptiste Rozière, Alexandre Sablayrolles, Lucile Saulnier, Romain Sauvestre, Wendy Shang, Roman Soletskyi, Lawrence Stewart, Pierre Stock, Joachim Studnia, Sandeep Subramanian, Sagar Vaze, Thomas Wang, and Sophia Yang. Pixtral 12b, 2024. URL https://arxiv.org/abs/2410.07073.

Michael Anderson, Bernd Meyer, and Patrick Olivier. *Diagrammatic representation and reasoning*. Springer Science & Business Media, 2011.

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

Maurits J. R. Bleeker, Mariya Hendriksen, Andrew Yates, and Maarten de Rijke. Demonstrating and reducing shortcuts in vision-language representation learning. *CoRR*, abs/2402.17510, 2024. doi: 10.48550/ARXIV.2402.17510. URL https://doi.org/10.48550/arXiv.2402.17510.

Andrea Burns, Krishna Srinivasan, Joshua Ainslie, Geoff Brown, Bryan A. Plummer, Kate Saenko, Jianmo Ni, and Mandy Guo. Wikiweb2m: A page-level multimodal wikipedia dataset, 2023. URL https://arxiv.org/abs/2305.05432.

Rémi Cadène, Corentin Dancette, Hédi Ben-Younes, Matthieu Cord, and Devi Parikh. Rubi: Reducing unimodal biases for visual question answering. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 839–850, 2019. URL https://proceedings.neurips.cc/paper/2019/hash/51d92be1c60d1db1d2e5e7a07da55b26-Abstract.html.

BETS Cullum-Swan and Peter Manning. Narrative, content, and semiotic analysis. *Handbook of qualitative research*, pp. 463–477, 1994.

Google. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context, 2024. URL https://arxiv.org/abs/2403.05530.

Google. Gemma 3 technical report, 2025. URL https://arxiv.org/abs/2503.19786.

Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the V in VQA matter: Elevating the role of image understanding in visual question answering. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pp. 6325–6334. IEEE Computer Society, 2017. doi: 10.1109/CVPR.2017.670. URL https://doi.org/10.1109/CVPR.2017.670.

Stanley I Greenspan and Stuart Shanker. *The first idea: How symbols, language, and intelligence evolved from our primate ancestors to modern humans*. Da Capo Press, 2009.

Yifan Hou, Buse Giledereli, Yilei Tu, and Mrinmaya Sachan. Do vision-language models really understand visual language?, 2025. URL https://arxiv.org/abs/2410.00193.

Bargav Jayaraman, Chuan Guo, and Kamalika Chaudhuri. Déjà vu memorization in vision-language models. In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang (eds.), *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024. URL http://papers.nips.cc/paper_files/paper/2024/hash/5ab6f836f464d0f4e4f6aaa523249280-Abstract-Conference.html.

Kushal Kafle, Brian Price, Scott Cohen, and Christopher Kanan. Dvqa: Understanding data visualizations via question answering, 2018. URL https://arxiv.org/abs/1801.08163.

Samira Ebrahimi Kahou, Vincent Michalski, Adam Atkinson, Akos Kadar, Adam Trischler, and Yoshua Bengio. Figureqa: An annotated figure dataset for visual reasoning, 2018. URL https://arxiv.org/abs/1710.07300.

Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali Farhadi. A diagram is worth a dozen images, 2016. URL https://arxiv.org/abs/1603.07396.

Aniruddha Kembhavi, Minjoon Seo, Dustin Schwenk, Jonghyun Choi, Ali Farhadi, and Hannaneh Hajishirzi. Are you smarter than a sixth grader? textbook question answering for multimodal machine comprehension. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 5376–5384, 2017. doi: 10.1109/CVPR.2017.571.

Jayant Krishnamurthy, Oyvind Tafjord, and Aniruddha Kembhavi. Semantic parsing to probabilistic programs for situated question answering, 2016. URL https://arxiv.org/abs/1606.07046.

Alexander Kuhnle and Ann Copestake. Shapeworld - a new test methodology for multimodal language understanding, 2017. URL `https://arxiv.org/abs/1704.04517`.

Ora Lassila and Ralph R. Swick. Resource Description Framework (RDF) Model and Syntax Specification, 1999. URL `http://www.w3.org/TR/1999/REC-rdf-syntax-19990222/`.

Sicong Leng, Hang Zhang, Guanzheng Chen, Xin Li, Shijian Lu, Chunyan Miao, and Lidong Bing. Mitigating object hallucinations in large vision-language models through visual contrastive decoding. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*, pp. 13872–13882. IEEE, 2024. doi: 10.1109/CVPR52733.2024.01316. URL `https://doi.org/10.1109/CVPR52733.2024.01316`.

Fei-Fei Li. *The Worlds I See: Curiosity, Exploration, and Discovery at the Dawn of AI*. Flatiron books: a moment of lift book, 2023.

Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 292–305, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.20. URL `https://aclanthology.org/2023.emnlp-main.20/`.

Yongqi Li, Wenjie Wang, Leigang Qu, Liqiang Nie, Wenjie Li, and Tat-Seng Chua. Generative cross-modal retrieval: Memorizing images in multimodal language models for retrieval and beyond. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pp. 11851–11861. Association for Computational Linguistics, 2024. doi: 10.18653/V1/2024.ACL-LONG.639. URL `https://doi.org/10.18653/v1/2024.acl-long.639`.

Fuxiao Liu, Xiaoyang Wang, Wenlin Yao, Jianshu Chen, Kaiqiang Song, Sangwoo Cho, Yaser Yacoob, and Dong Yu. Mmc: Advancing multimodal chart understanding with large-scale instruction tuning, 2024a. URL `https://arxiv.org/abs/2311.10774`.

Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, January 2024b. URL `https://llava-vl.github.io/blog/2024-01-30-llava-next/`.

Junpeng Liu, Yifan Song, Bill Yuchen Lin, Wai Lam, Graham Neubig, Yuanzhi Li, and Xiang Yue. Visualwebbench: How far have multimodal llms evolved in web page understanding and grounding?, 2024c. URL `https://arxiv.org/abs/2404.05955`.

Pan Lu, Liang Qiu, Jiaqi Chen, Tanglin Xia, Yizhou Zhao, Wei Zhang, Zhou Yu, Xiaodan Liang, and Song-Chun Zhu. Iconqa: A new benchmark for abstract diagram understanding and visual language reasoning. In Joaquin Vanschoren and Sai-Kit Yeung (eds.), *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual*, 2021. URL `https://datasets-benchmarks-proceedings.neurips.cc/paper/2021/hash/d3d9446802a44259755d38e6d163e820-Abstract-round2.html`.

Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for science question answering, 2022. URL `https://arxiv.org/abs/2209.09513`.

Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL `https://openreview.net/forum?id=KUNzEQMWU7`.

Kim Marriott and Bernd Meyer. *Visual language theory*. Springer Science & Business Media, 2012.

Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A benchmark for question answering about charts with visual and logical reasoning, 2022. URL https://arxiv.org/abs/2203.10244.

Minesh Mathew, Viraj Bagal, Rubèn Pérez Tito, Dimosthenis Karatzas, Ernest Valveny, and C. V Jawahar. Infographicvqa, 2021a. URL https://arxiv.org/abs/2104.12756.

Minesh Mathew, Dimosthenis Karatzas, and C. V. Jawahar. Docvqa: A dataset for vqa on document images, 2021b. URL https://arxiv.org/abs/2007.00398.

Meta. The llama 3 herd of models, 2024. URL https://arxiv.org/abs/2407.21783.

Nitesh Methani, Pritha Ganguly, Mitesh M. Khapra, and Pratyush Kumar. Plotqa: Reasoning over scientific plots, 2020. URL https://arxiv.org/abs/1909.00997.

Microsoft. Phi-4-mini technical report: Compact yet powerful multimodal language models via mixture-of-loras, 2025. URL https://arxiv.org/abs/2503.01743.

Anand Mishra, Shashank Shekhar, Ajeet Kumar Singh, and Anirban Chakraborty. Ocr-vqa: Visual question answering by reading text in images. In *2019 International Conference on Document Analysis and Recognition (ICDAR)*, pp. 947–952, 2019. doi: 10.1109/ICDAR.2019.00156.

Atsuyuki Miyai, Jingkang Yang, Jingyang Zhang, Yifei Ming, Qing Yu, Go Irie, Yixuan Li, Hai Li, Ziwei Liu, and Kiyoharu Aizawa. Unsolvable problem detection: Evaluating trustworthiness of vision language models. *CoRR*, abs/2403.20331, 2024. doi: 10.48550/ARXIV.2403.20331. URL https://doi.org/10.48550/arXiv.2403.20331.

Charles William Morris. Foundations of the theory of signs. In *International encyclopedia of unified science*, pp. 1–59. Chicago University Press, 1938.

Charles Peirce, Christian S., and Nathan House J. W. Kloesel. *The Essential Peirce: Selected Philosophical Writings Vol. 1*. Indiana University Press, Bloomington, 1992.

Charles Sanders Peirce. Logic as semiotic: The theory of signs. In Charles Sanders Peirce (ed.), *Philosophical Writings*. Dover Publications, 1935.

Shailaja Keyur Sampat, Yezhou Yang, and Chitta Baral. Visuo-linguistic question answering (vlqa) challenge, 2020. URL https://arxiv.org/abs/2005.00330.

Min Joon Seo, Hannaneh Hajishirzi, Ali Farhadi, and Oren Etzioni. Diagram understanding in geometry questions. In Carla E. Brodley and Peter Stone (eds.), *Proceedings of the Twenty-Eighth AAAI Conference on Artificial Intelligence, July 27 -31, 2014, Québec City, Québec, Canada*, pp. 2831–2838. AAAI Press, 2014. doi: 10.1609/AAAI.V28I1.9146. URL https://doi.org/10.1609/aaai.v28i1.9146.

Mong Yuan Sim, Wei Emma Zhang, Xiang Dai, and Biaoyan Fang. Can vlms actually see and read? A survey on modality collapse in vision-language models. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (eds.), *Findings of the Association for Computational Linguistics, ACL 2025, Vienna, Austria, July 27 - August 1, 2025*, pp. 24452–24470. Association for Computational Linguistics, 2025. URL https://aclanthology.org/2025.findings-acl.1256/.

Alane Suhr, Mike Lewis, James Yeh, and Yoav Artzi. A corpus of natural language for visual reasoning. In Regina Barzilay and Min-Yen Kan (eds.), *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pp. 217–223, Vancouver, Canada, July 2017. Association for Computational Linguistics. doi: 10.18653/v1/P17-2034. URL https://aclanthology.org/P17-2034/.

Ryota Tanaka, Kyosuke Nishida, and Sen Yoshida. Visualmrc: Machine reading comprehension on document images, 2021. URL https://arxiv.org/abs/2101.11272.

Xiaoxuan Wang, Ziniu Hu, Pan Lu, Yanqiao Zhu, Jieyu Zhang, Satyen Subramaniam, Arjun R. Loomba, Shichang Zhang, Yizhou Sun, and Wei Wang. Scibench: Evaluating college-level scientific problem-solving abilities of large language models. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net, 2024a. URL `https://openreview.net/forum?id=bq1JEgioLr`.

Zirui Wang, Mengzhou Xia, Luxi He, Howard Chen, Yitao Liu, Richard Zhu, Kaiqu Liang, Xindi Wu, Haotian Liu, Sadhika Malladi, Alexis Chevalier, Sanjeev Arora, and Danqi Chen. Charxiv: Charting gaps in realistic chart understanding in multimodal llms, 2024b. URL `https://arxiv.org/abs/2406.18521`.

Hu Xu, Saining Xie, Xiaoqing Ellen Tan, Po-Yao Huang, Russell Howes, Vasu Sharma, Shang-Wen Li, Gargi Ghosh, Luke Zettlemoyer, and Christoph Feichtenhofer. Demystifying clip data, 2024a. URL `https://arxiv.org/abs/2309.16671`.

Zhengzhuo Xu, Sinan Du, Yiyan Qi, Chengjin Xu, Chun Yuan, and Jian Guo. Chartbench: A benchmark for complex visual reasoning in charts, 2024b. URL `https://arxiv.org/abs/2312.15915`.

Le Xue, Manli Shu, Anas Awadalla, Jun Wang, An Yan, Senthil Purushwalkam, Honglu Zhou, Viraj Prabhu, Yutong Dai, Michael S Ryoo, Shrikant Kendre, Jieyu Zhang, Can Qin, Shu Zhang, Chia-Chih Chen, Ning Yu, Juntao Tan, Tulika Manoj Awalgaonkar, Shelby Heinecke, Huan Wang, Yejin Choi, Ludwig Schmidt, Zeyuan Chen, Silvio Savarese, Juan Carlos Niebles, Caiming Xiong, and Ran Xu. xgen-mm (blip-3): A family of open large multimodal models, 2024. URL `https://arxiv.org/abs/2408.08872`.

Halina Sendera Mohd Yakin and Andreas Totu. The semiotic perspectives of peirce and saussure: A brief comparative study. *Procedia-Social and Behavioral Sciences*, 155:4–8, 2014.

Yuan Zang, Tian Yun, Hao Tan, Trung Bui, and Chen Sun. Pre-trained vision-language models learn discoverable visual concepts. *CoRR*, abs/2404.12652, 2024. doi: 10.48550/ARXIV.2404.12652. URL `https://doi.org/10.48550/arXiv.2404.12652`.

Jakub Zdebik. Deleuze and the diagram. *Deleuze and the Diagram*, pp. 1–256, 2012.

Peng Zhang, Yash Goyal, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Yin and yang: Balancing and answering binary visual questions, 2016. URL `https://arxiv.org/abs/1511.05099`.

## A   RELATED WORKS

**Diagram Question Answering (DQA).**   Diagram Question Answering (DQA) is a specialized subfield of Visual Question Answering (VQA), where the input image is a schematic, symbolic, or abstract diagram rather than a natural scene (Hou et al., 2025). These diagrams commonly convey structured, domain-specific knowledge—such as scientific processes, mathematical relations, or logical systems—making DQA a valuable testbed for evaluating a model's ability to perform symbolic interpretation and structured visual reasoning.

**Benchmarks on Statistical and Analytical Diagrams.**   One major category of DQA benchmarks focuses on statistical or analytical charts, such as bar graphs, line plots, and scatter plots. These tasks require models to extract numerical values, recognize trends, and reason over structured visual features. Notable datasets in this area include FigureQA (Kahou et al., 2018), DVQA (Kafle et al., 2018), PlotQA (Methani et al., 2020), ChartQA (Masry et al., 2022), MMC (Liu et al., 2024a), ChartBench (Xu et al., 2024b), and CharXiv (Wang et al., 2024b).

**Benchmarks on Visually Structured Content.**   Another category evaluates visually structured content, particularly infographics and document-like formats. These include images such as posters, book covers, webpages, and scientific figures, where layout-aware reasoning is critical. Datasets like OCR-VQA (Mishra et al., 2019), DocVQA (Mathew et al., 2021b), InfographicVQA (Mathew et al., 2021a), VisualMRC (Tanaka et al., 2021), and VisualWebBench (Liu et al., 2024c) target the integration of visual structure and textual information.

**Benchmarks from Educational and Instructional Diagrams.**   Several DQA benchmarks are derived from science education and domain-specific instructional content, often sourced from textbooks or learning platforms. These diagrams are rich and require external knowledge integration. Key datasets in this space include AI2D (Kembhavi et al., 2016), FoodWebs (Krishnamurthy et al., 2016), TQA (Kembhavi et al., 2017), VLQA (Sampat et al., 2020), and ScienceQA (Lu et al., 2022).

**Benchmarks on Synthetic and Abstract Diagrams.**   A final class of benchmarks uses synthetic or abstract diagrams to isolate core reasoning skills. These datasets typically involve geometric primitives or symbolic representations that are free from real-world biases. NLVR (Suhr et al., 2017) and ShapeWorld (Kuhnle & Copestake, 2017) focus on compositional and spatial reasoning, while Zhang et al. (2016) and IconQA (Lu et al., 2021) test high-level relational and symbolic inference through minimalistic, abstract scenes.

## B   DETAILS OF TEST SUITE CONSTRUCTION

### B.1   DIAGRAM CLEANING

To construct a comprehensive diagram test suite, we source images from one of the largest open-source knowledge bases: Wikipedia. Specifically, we use WikiWeb2M (Burns et al., 2023), a large-scale dataset containing over 2 million English Wikipedia webpages with diverse images, rich textual content, and structured metadata.

However, WikiWeb2M includes many non-diagram images such as human portraits, logos, and natural scenes. To isolate true diagrammatic content, we design a binary classification pipeline based on MetaCLIP (Xu et al., 2024a). We construct one descriptive prompt to identify diagrams and six complementary prompts to exclude non-diagram content. Each image is evaluated across these prompts, and only those classified as diagrams in all negative prompt settings are retained. This conservative strategy ensures high precision in diagram selection. The full list of prompts used in this filtering process is provided in Fig. 9. After filtering, we retain approximately 100,000 diagram candidates for further processing.

### B.2   DIAGRAM TAGGING

Since diagrams serve as versatile tools for knowledge transfer, they span a wide variety of types and subject domains. To better organize our test suite and support structured annotation, we use two

vision-language models (Molmo-7B and LLaMA-3.2-7B) to tag each diagram with both its type and associated knowledge domain (Fig. 9). The full prompt templates used for tagging are available in Figs. 10 to 12.

We repeat the tagging process twice with both models, resulting in four independent annotations per image. We then manually analyze the distribution of tags and consolidate the most frequent ones into 12 categories. These are divided into two groups:

- **Statistical Group**: Includes four types of statistical diagrams — Bar Chart, Line Graph, Pie Chart, and Map.
- **Scientific Group**: Includes eight types of non-statistical diagrams categorized by academic disciplines — Biology, Chemistry, Computer Science, Mathematics, Physics, Astronomy, History, and Music.

To ensure label consistency and reliability, we retain only diagrams with consistent tags across all four annotations. This filtering results in a curated set of approximately 60,000 diagrams.

### B.3 Diagram Annotation

Our test suite contains two core forms of annotation: semantic triples and question–answer (QA) pairs, which together capture both the content of the diagram and the levels of comprehension required.

To ensure annotation quality, we use Gemini-2.0-Flash (Google, 2024) as the primary annotation model in a structured two-step process.

**Step 1: Diagram Description.**   To simplify the downstream annotation and improve quality, we first prompt Gemini to generate a detailed description of each diagram. This intermediate step provides a structured foundation from which semantic triples and QA pairs are derived. Since triple extraction and QA generation emphasize different semantic aspects of a diagram, the description prompts are carefully designed to highlight relevant content.

To reduce hallucination—an inherent issue in large models (Li et al., 2023; Leng et al., 2024)—we supplement each image with its corresponding Wikipedia text to provide factual grounding. Moreover, we design tailored prompts for different diagram groups (e.g., statistical vs. scientific) and include in-context examples to guide the model away from vague or generic outputs. Full prompt details are in Figs. 13 to 16.

**Step 2: Semantic Triples and QA Pairs.**   Using the diagram description, we prompt Gemini again to extract semantic triples and generate multiple-choice QA pairs. Detailed prompt designs are available in Figs. 17 to 20.

To ensure the quality of the QA annotations, we implement a three-stage consistency check:

- **Visual Dependency Check (No Image)**: The model attempts to answer questions without seeing the diagram. If it succeeds, the question likely does not depend on the visual content.
- **Wiki-Text Independency Check (No Image + Wiki-Text)**: The model is shown the Wikipedia context but not the image. The question should remain unanswerable.
- **Triple Completeness Check (No Image + Triples)**: The model is given only textual sentences derived from the semantic triples. The question should be answerable in this setting.

Each setting is evaluated twice with shuffled answer choices to minimize bias. We consider a diagram as "succeeded" if the model selects the correct answer in both runs, and as "failed" if it make mistakes in either run.

We discard diagrams:

- That succeed in the entity recognition task in the first two checks, indicating that the QA annotation is not image-dependent.
- That fail in any of the four tasks (ER, RU, KG, VR) in the third check, indicating that triples are incomplete.

After applying these filters, we retain a total of 7,500 diagrams, though the category distribution remains imbalanced. From this pool, we curate a balanced test set of 1,500 diagrams and a training set of 6,000 diagrams. Comprehensive category-wise statistics are presented in Tab. 3. [2]

Table 3: Number of diagrams per category in the test dataset and training dataset.

| Category | Test Set | Training Set |
|---|---|---|
| Bar Chart | 150 | 900 |
| Line Graph | 150 | 350 |
| Pie Chart | 150 | 0 |
| Map | 150 | 2000 |
| Biology | 150 | 900 |
| Chemistry | 150 | 1600 |
| Computer Science | 150 | 0 |
| Mathematics | 150 | 150 |
| Physics | 150 | 100 |
| Others | 150 | 0 |

## C  SUPPLEMENTARY RESULTS

### C.1  EXPERIMENT SETUP DETAILS

#### C.1.1  MODEL LIST

We evaluate a diverse set of vision-language Models (VLMs) on our test suite. Our selection encompasses both industry-developed models from leading AI companies such as Google, Meta, Alibaba, and Microsoft, as well as representative open-source models from the academic community. For certain model families, we include multiple variants with different parameter scales to facilitate comparative analysis. The following models are evaluated in our test suite.

**Qwen-2.5-VL**  (Bai et al., 2025) is a multimodal model series developed by Alibaba, featuring a native dynamic-resolution Vision Transformer with window attention, enabling efficient processing of high-resolution images and long-form videos. It supports precise object grounding with absolute coordinates and demonstrates strong capabilities in document parsing, chart interpretation, and temporal event localization. In our experiments, we evaluate four variants of Qwen2.5-VL with 3B, 7B, 32B, and 72B parameters.

**LLaMA-3.2**  (Meta, 2024) is a large-scale foundation model family developed by Meta. It introduces multimodal capabilities, integrating image, video, and speech understanding via modular adapters. For vision, it employs a pretrained image encoder, connected to the language model through a cross-attention-based vision adapter. This compositional setup allows the system to process image-text pairs without modifying the core language model. In our experiments, we evaluate two variants of LLaMA-3 with 11B, and 90B parameters.

**Gemma-3**  (Google, 2025) is a multimodal model series developed by Google DeepMind, supporting vision, long-context reasoning, and multilingual understanding. It adopts a decoder-only architecture with grouped-query attention and introduces a local-to-global attention mechanism to reduce KV-cache memory overhead during long-context inference. For vision processing, it can handle flexible image resolutions. In our experiments, we evaluate three variants of Gemma-3 with 4B, 12B, and 27B parameters.

**Pixtral**  (Agrawal et al., 2024) is a multimodal language model developed by Mistral. It features a custom vision encoder trained from scratch, capable of ingesting images at their native resolution and aspect ratio, and supports flexible tokenization strategies. The model employs RoPE-2D position encoding in the vision encoder and uses a decoder-only architecture based on Mistral NeMo. In our experiments, we evaluate the 12B variant.

---

[2]Our data license is CC-BY-4.0.

**Phi-4** (Microsoft, 2025) is a multimodal model developed by Microsoft, extending the Phi-4 series to support text, vision, and speech/audio modalities. It employs a novel Mixture-of-LoRAs architecture that integrates modality-specific adapters without modifying the frozen language backbone, thus preserving its strong language capabilities. In our experiments, we evaluate the 5.6B variant.

**BLIP-3 (xGen-MM)** (Xue et al., 2024) is a multimodal model series developed by Salesforce, designed to unify training objectives and scale vision-language understanding through a simplified architecture. The framework replaces the Q-Former in previous models with a scalable perceiver resampler, enabling efficient any-resolution vision token sampling and supporting interleaved multimodal inputs. In our experiments, we evaluate the 4B variant.

**LLaVA-1.6** (Liu et al., 2024b) is a multimodal model series that enhances visual reasoning, OCR, and world knowledge while maintaining a lightweight architecture. It introduces higher input resolutions and refined visual instruction tuning, enabling better understanding of complex visual scenes. In our experiments, we evaluate three variants of LLaVA-1.6 with 7B, 13B, and 34B parameters.

### C.1.2 PROMPT PIPELINE

For question answering, we design a three-step, systematic, rule-based evaluation pipeline. In the first step, the model is presented with the input multimodal data and a corresponding question, and is prompted to analyze and answer the question in a step-by-step manner. In the second step, given the full preceding context, the model is instructed to produce a final, conclusive answer in the form of a multiple-choice selection (i.e., A, B, C, or D). To address potential limitations in instruction-following abilities (especially in smaller models), we introduce a third step that automatically extracts the final answer from the model's generated response in Step 2. This is achieved using a set of robust regular expressions and response-processing workflows that identify key phrases, such as numeric values and conclusion markers, to ensure accurate answer extraction and matching. An example of the three-step pipeline is shown in Fig. 21.

### C.1.3 HUMAN EVALUATION GUIDELINES

The guideline for the human evaluation of the data annotation quality assessment is given below.

- **Visual Dependency.** Evaluate whether answering the questions requires visual reference to the diagram. *Fully Dependent* means all questions rely on visual information (e.g., labels, layout, spatial structure). *Partially Dependent* indicates that **at least one question** could be answered without seeing the diagram, using commonsense or background knowledge.

- **QA Correctness.** Assess the overall quality of the four QA pairs. *Perfectly Valid* means all QA pairs are accurate, clear, and grounded in the diagram. *Slightly Flawed* means **at least one QA pair** contains minor issues such as ambiguity, hallucination, or poor phrasing.

- **Triple Completeness.** Examine how well the knowledge triples represent the information in the diagram. *Totally Sufficient* indicates that the triple set is comprehensive, factually correct, and well-structured. *Marginally Insufficient* means that **at least one triple** misses important details, includes minor errors, or lacks clarity.

### C.1.4 PROJECT COST

In our test suite, most experiments are conducted on NVIDIA GPUs, including RTX 3090 and A100, with the specific hardware selected based on model size. For Llama-3.2-90B only, we leverage the Together AI inference API to perform evaluation. Additionally, since we only perform inference on VLMs, we use `torch.bfloat16` precision for all tasks for reducing GPU memory usage.

We report the computation resources to clean and annotate our test suite. Besides, we report the computing cost for our evaluation. We measure the computation cost by GPU Hours and the financial cost for API models in Tab. 4.

### C.2 DETAILED RESULTS

Table 4: The cost of building our test suite and evaluation on our test suite.

| Task | Model | Data | Type | Cost |
|---|---|---|---|---|
| Diagram Cleaning | MetaCLIP | 2M | H100 | 200 GPU hours |
| Diagram Tagging | Molmo & LLaMA3.2 | 100k | RTX3090 | 400 GPU hours |
| Diagram Annotation | Gemini | 60k | Google API | 8,000 USD |
| Consistency Checking | Gemini | 60k | Google API | 12,000 USD |
| Evaluation | 14 VLMs | 1.5k | RTX3090/A100 | 100 GPU hours |
| | LLaMA-90B | | TogetherAI API | 400 USD |

Table 5: Comparative evaluation of multiple vision-language models across real, synthetic, and textual modalities on four tasks. The best-performing result is highlighted in **bold**, and the second-best is underlined. Note that ER, RU, KG, and VR denote *entity recognition*, *relation understanding*, *knowledge grounding*, and *visual reasoning*.

| Model | Visual Modality | | | | Semantic Modality | | | | Textual Modality | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | ER | RU | KG | VR | ER | RU | KG | VR | ER | RU | KG | VR |
| Qwen2.5-VL-3B (Bai et al., 2025) | 88.2 | 90.4 | 91.1 | 88.3 | 87.5 | 91.0 | 94.2 | 90.3 | 89.8 | 91.9 | 92.0 | 89.4 |
| Qwen2.5-VL-7B (Bai et al., 2025) | 91.9 | 93.0 | 94.0 | 90.9 | 88.3 | 92.8 | 93.9 | 89.3 | 92.9 | 93.7 | 93.1 | 91.1 |
| Qwen2.5-VL-32B (Bai et al., 2025) | 92.9 | 94.7 | 95.3 | 93.5 | **93.8** | **95.3** | **97.4** | **95.9** | 95.7 | 96.3 | **98.2** | 95.6 |
| Qwen2.5-VL-72B (Bai et al., 2025) | **94.3** | **95.6** | 95.5 | **94.4** | 92.3 | 95.0 | 97.1 | 94.6 | 95.5 | **97.1** | 97.9 | 95.7 |
| LLaMA3.2-11B (Meta, 2024) | 71.6 | 74.6 | 78.1 | 75.9 | 67.1 | 74.2 | 78.8 | 72.9 | 84.5 | 89.5 | 90.7 | 88.6 |
| LLaMA3.2-90B (Meta, 2024) | 89.9 | 91.8 | 94.5 | 92.5 | 81.3 | 90.1 | 93.2 | 88.6 | 95.7 | 96.3 | 96.9 | 94.6 |
| Gemma3-4B (Google, 2025) | 83.7 | 85.7 | 88.1 | 84.9 | 77.9 | 82.3 | 87.5 | 84.1 | 88.0 | 87.1 | 89.3 | 86.9 |
| Gemma3-12B (Google, 2025) | 90.1 | 92.9 | 93.9 | 92.4 | 87.0 | 90.4 | 94.9 | 91.5 | 93.0 | 94.1 | 95.1 | 92.9 |
| Gemma3-27B (Google, 2025) | 91.9 | 95.0 | 95.1 | 93.5 | 90.3 | 93.3 | 96.4 | 94.0 | **96.2** | 96.5 | 96.1 | **95.8** |
| LLaVA1.6-7B (Liu et al., 2024b) | 50.6 | 60.3 | 65.1 | 62.0 | 44.8 | 57.1 | 62.8 | 57.8 | 71.7 | 78.1 | 82.4 | 77.8 |
| LLaVA1.6-13B (Liu et al., 2024b) | 63.3 | 75.5 | 81.1 | 78.9 | 56.4 | 71.3 | 79.7 | 75.0 | 82.8 | 86.9 | 89.8 | 85.6 |
| LLaVA1.6-34B (Liu et al., 2024b) | 81.0 | 84.9 | 88.7 | 86.2 | 71.1 | 83.2 | 89.3 | 86.5 | 91.2 | 92.7 | 94.5 | 91.6 |
| Pixtral-12B (Agrawal et al., 2024) | 89.0 | 88.9 | 89.5 | 88.7 | 77.3 | 85.3 | 90.1 | 88.3 | 92.5 | 94.3 | 95.7 | 91.9 |
| Phi4-5.6B (Microsoft, 2025) | 88.2 | 90.6 | 90.1 | 89.0 | 83.2 | 88.9 | 90.0 | 85.9 | 88.9 | 90.3 | 93.7 | 91.3 |
| BLIP3-4B (Xue et al., 2024) | 42.9 | 72.7 | 75.0 | 75.0 | 42.9 | 72.7 | 75.0 | 75.0 | 83.9 | 86.4 | 88.6 | 82.0 |

## C.3 PROMPT EXAMPLES

---

### Prompt for Diagram Cleaning
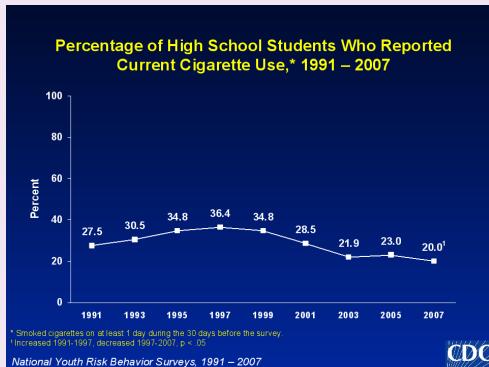
**Positive Prompt:**

- A visual representation of information or data, explicitly intended for educational or scientific purposes. This includes flowcharts, circuit diagrams, architectural blueprints, and graphs, characterized by clear labeling and structured layout for easy understanding of complex concepts.

**Negative Prompts:**

- An **image of a company or brand logo**, designed to be a simple yet distinctive symbol that represents a company or product. Logos often consist of stylized letterforms, abstract geometric shapes, or a combination of both, and are designed to be easily recognizable even at small sizes. They usually feature a limited color palette and lack detailed textual information.

- An **image depicting natural landscapes**, including forests, mountains, rivers, or beaches, characterized by vivid natural colors and organic forms without any superimposed text or symbols.

- A **photograph** of one or several human beings, focusing on the face or figure, often capturing expression, personality, and mood, without any overlay of graphical information or text.

- **Images of old books, pages, or manuscripts**, primarily showing textual content in a historical or literary context, often with visible textures of paper and traditional fonts.

- A **screenshot** from a computer or mobile device, typically showing a user interface with icons, menus, and open applications, which may include web pages, software programs, or mobile apps.

- An **image with minimal visual content**, often appearing as a solid color background with sparse elements like one or two letters or one or two simple shapes. These images lack detail and complexity, presenting very basic or stark visual information with no significant features or recognizable patterns.

---

Figure 9: We perform six rounds of binary classification. In each round, an image is classified as a diagram or not by comparing its embedding with the embeddings of the two text prompts using MetaCLIP. Only images consistently classified as positive examples—that is, diagrams—across all rounds are retained.

## Prompt for Tagging (Step 1: Captioning)



**Percentage of High School Students Who Reported Current Cigarette Use,* 1991 – 2007**

Data points: 27.5 (1991), 30.5 (1993), 34.8 (1995), 36.4 (1997), 34.8 (1999), 28.5 (2001), 21.9 (2003), 23.0 (2005), 20.0 (2007)

\* Smoked cigarettes on at least 1 day during the 30 days before the survey.
† Increased 1991-1997, decreased 1997-2007, p < .05

*National Youth Risk Behavior Surveys, 1991 – 2007*

**System:** You are a diagram description assistant. Your task is to provide a detailed and structured description of the given diagram. Focus on aspects that might help to tag its domain (e.g., Biology, Chemistry, History) and type (e.g., Bar Chart, Flow Chart, Map).

**Context:** The diagram is sourced from Wikipedia, and here is some background information. Use the Wikipedia information above only if the diagram alone does not provide enough clarity or context. Always give priority to the information directly visible in the diagram for your analysis.

- **Page Title:** Prevalence of tobacco use.

- **Page Description:** Prevalence of tobacco use is reported by the World Health Organization, which focuses on cigarette smoking due to reported data limitations. Smoking has therefore been studied more extensively than any other form of consumption.Smoking is generally five times more prevalent among men than women; however, the gender gap differs across countries and is smaller in younger age groups. (text truncated due to space)

- **Diagram Description:** None.

**Instruction:** The description must be organized into the following three sections:

- **Content:** Describe key visual elements, labels, and any prominent features in the diagram.

- **Layout:** Explain how the elements are arranged (e.g., hierarchical, circular, linear) and the overall structure.

- **Function:** Indicate the likely purpose of the diagram (e.g., explaining a process, showing relationships, presenting data).

Figure 10: Before predicting tags for the diagrams, we conduct a captioning step. We instruct the VLM to act as a diagram description assistant and provide it with contextual information from Wikipedia, including the page title, page description, and diagram description (if available). The model is then prompted to focus on describing the content, layout, and function of the diagram.

## Prompt for Tagging (Step 2: Open-Ended Prediction)

**System:** You are a diagram tagging assistant. Your task is to analyze a diagram and identify its domain and type.

**Context:** The description of the diagram is provided for your reference:

- **Content:** The diagram appears to be a line graph depicting trends over time. It shows data points connected by lines, representing changes in a specific measure from 1991 to 2007. The graph includes numerical values on the y-axis and years on the x-axis. There are likely labels for the y-axis and x-axis, as well as a title at the top of the graph.

- **Layout:** The layout of the diagram is typical of a line graph. The vertical axis (y-axis) represents percentages, while the horizontal axis (x-axis) represents years. The data points are plotted along the x-axis and connected by lines to show the trend over time. The title is likely positioned at the top of the graph, providing context for the data being presented.

- **Function:** The function of this diagram is to visually represent and illustrate trends in a specific measure over a 16-year period. It allows viewers to quickly understand how the measured value has changed from 1991 to 2007. The use of a line graph makes it easy to see patterns, trends, and changes in the data over time, which is particularly useful for analyzing long-term data sets and identifying any significant shifts or fluctuations in the measured variable.

**Instruction:** Now analyze the diagram and provide its domain and type:

- **Domain:** The domain should be a specific field or area of knowledge. Its examples include Biology, Chemistry, Physics, Astronomy, History, etc.

- **Type:** The type should describe the nature of the diagram. Its examples include Bar Chart, Flow Chart, Table, Map, Logo, etc.

**Output Format:** Your output must be in the following JSON-like format. Do not provide any explanations or additional context. Only output the JSON object.
{
    "Domain": "string (must be 1 or 2 words)",
    "Type": "string (must be 1 or 2 words)"
}

Figure 11: After generating a caption for the diagram, we prompt the VLM again using the annotated content, layout, and function descriptions, and ask it to predict both a domain tag and a type tag. In this step, we adopt an open-ended setting, allowing the model to freely generate tags without any predefined options.

## Prompt for Tagging (Step 2: Multiple-Choice Prediction)

**System:** The same as Figure 11.

**Context:** The same as Figure 11.

**Instruction:** Now analyze the diagram and provide its domain and type:

- **Domain:** The domain should be a specific field or area of knowledge. Choose only one option from the following list:
    - Agriculture
    - Astronomy
    - Biology
    - Chemistry
    - Computer Science
    - Data Science
    - Environmental Science
    - Finance
    - Geography and Geology
    - Health Science
    - History
    - Mathematics
    - Music
    - Network Science
    - Operations Research
    - Physics
    - Political Science
    - Psychology
    - Sports
    - Transportation
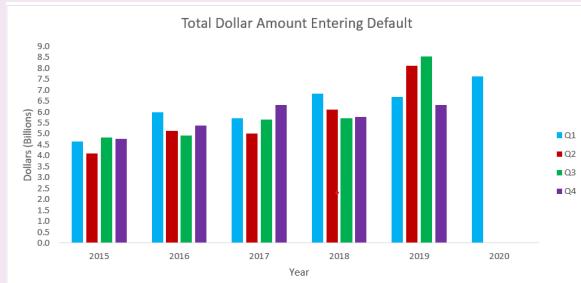    - Urban Planning

- **Type:** The type should describe the nature of the diagram. Choose only one option from the following list:
    - Bar Chart
    - Chemical Visual
    - Concept Diagram
    - Floor Plan
    - Flow Chart
    - Line Graph
    - Logo
    - Map
    - Network Chart
    - Pie Chart
    - Scatter Plot
    - Table
    - Technical Diagram
    - Timeline
    - Tree

**Output Format:** The same as Figure 11.

Figure 12: After generating open-ended tags, we apply clustering methods to analyze the tag distribution and identify a set of high-frequency tags, which are then used as options for the multiple-choice tagging setting. In this setting, we keep the instructions and context unchanged, but instead of allowing free predictions, the VLM is asked to select tags from the option list.

## Prompt for Statistical Annotation (Step 1: Captioning)

Total Dollar Amount Entering Default

**System:** You are a scene graph construction assistant. Your task is to generate a detailed language-based description of a scene graph for a provided diagram.

**Context:** The diagram is sourced from Wikipedia, and here is some background information. Use the Wikipedia information above only if the diagram alone does not provide enough clarity or context. Always give priority to the information directly visible in the diagram for your analysis.

- **Page Title:** Federal Direct Student Loan Program.
- **Page Description:** The William D. Ford Federal Direct Loan Program provides low-interest loans for students and parents to help ... (text truncated due to space)
- **Diagram Description:** Total number of dollars (in billions) entering default, 2009-2018, data source: CRS.

**Instruction:**
- Identify key elements such as axes, labels, legends, colors, and numerical values.
- Describe trends, patterns, or outliers in the data, including peaks, or correlations.
- Explain relationships between different variables if applicable.
- Describe geographical features such as colored regions and arrows if applicable.
- Use clear and structured language.

**Examples:**
- The bar representing Q3 in 2019 is the tallest among all quarters.
- The blue line in the graph shows a steady increase from 2010 to 2018.
- The dark green segment in the pie chart represents 45.9 TWh of diesel consumption.
- The shaded region in the map highlights areas with the highest population density.
- The thick arrow marks the strongest southeastern wind current towards the country.

Figure 13: Similar to the tagging stage, we conduct a captioning step before generating semantic triples in order to reduce hallucinations. We also provide the model with contextual information from Wikipedia. For statistical diagrams, we instruct the model to focus on specific features such as numerical values and data trends. To enhance the quality of output, we manually design five descriptive sentences that serve as in-context examples during prompting.

## Prompt for Statistical Annotation (Step 2: Annotation)

**System:** You are an expert information extraction assistant specializing in scene graph construction. Your task is to analyze a given diagram description and extract meaningful, structured relationships between key elements.

**Context:** The description of the diagram is provided for your reference.

**1. Key Objects: X-axis:** Represents the years from 2009 to 2018. Each year is labeled along the axis. **Y-axis:** Represents the total dollars in billions entering default. The axis is labeled "Dollars in Billions". Numerical markers are present along the axis, though precise values are not clearly visible in the image. **Bars:** Vertical bars represent the amount of dollars entering default for each year. The height of each bar corresponds to the dollar amount. **Data Labels:** Numerical values are displayed above each bar, indicating the precise amount for each year.

**2. Attributes: X-axis:** Horizontal, evenly spaced tick marks representing years. **Y-Axes:** Vertical, with numerical markers indicating billions of dollars. The scale appears to range from approximately 0 to 80 billion. **Bars:** Vertical rectangular bars, colored blue. The width of each bar is uniform. **Data Labels:** Black text, positioned above each bar.

**3. Relationships:** Each bar is associated with a year on the x-axis and a value on the y-axis. The height of the bar corresponds directly to the value indicated by the data label and represents the amount in billions of dollars entering default in that year.

**4. Structural or Hierarchical Information:** The chart is a simple bar chart.

**5. Data Trends:** The chart shows a general trend of increasing dollars entering default from 2009 to a peak, followed by a decrease and then another increase toward the end of the period (2018). Precise yearly fluctuations are observable but require more detailed numerical data. There is no clear outlier year that significantly deviates from the general pattern.
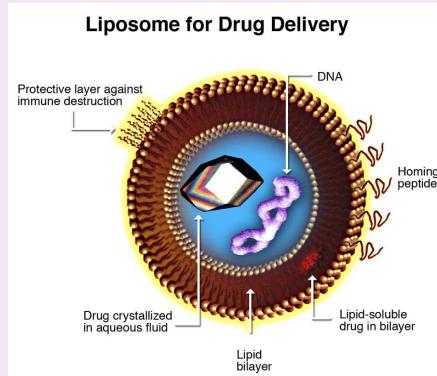
**Instruction:**

- Identify important relationships between key elements from the description.
- Structure these relationships in the form of triples with three components:
    - **Source**: The primary element (subject) in the relationship.
    - **Relationship**: The type of connection between the source and target.
    - **Target**: The secondary element (object) in the relationship.
- Ensure that:
    - Each triple represents a meaningful connection between elements.
    - The relationships are concise yet descriptive.
    - There are no duplicate, redundant, or meaningless triples.

**Output Format:** The final output must strictly follow the JSON format below:
```
{
    "1": {"Source": "Triple 1", "Relationship": "Triple 1", "Target": "Triple 1"},
    ...
    "N": {"Source": "Triple N", "Relationship": "Triple N", "Target": "Triple N"}
}
```

Figure 14: After extracting relevant information from the diagram, we prompt the model to generate a list of triples, where each triple consists of a source (head entity), a relationship (relation), and a target (tail entity). To facilitate downstream processing, we instruct the model to produce the output in JSON format.

## Prompt for Scientific Annotation (Step 1: Captioning)



**Liposome for Drug Delivery**

**System:** You are a scene graph construction assistant. Your task is to generate a detailed language-based description of a scene graph for a provided diagram.

**Context:** The diagram is sourced from Wikipedia, and here is some background information. Use the Wikipedia information above only if the diagram alone does not provide enough clarity or context. Always give priority to the information directly visible in the diagram for your analysis.

- **Page Title:** Nanomedicine.
- **Page Description:** Nanomedicine is the medical application of nanotechnology. Nanomedicine ranges from the medical applications of nanomaterials and biological devices, to nanoelectronic biosensors ... (text truncated due to space)
- **Diagram Description:** Liposomes are composite structures made of phospholipids and may contain small amounts of other molecules. Though liposomes can vary in size from low micrometer range to ... (text truncated due to space)

**Instruction:**
- Identify key objects, such as text, arrows, nodes, or data points.
- Identify attributes, such as size, color, shape, position, and numerical values.
- Explain how objects interact or relate to one another.
- Describe its overall hierarchy, structure or flow clearly if applicable.
- Use clear and structured language.

**Examples:**
- The newly discovered moon is connected to its elliptical orbit around Neptune.
- The blue alpha-helices are connected to beta-sheets through loop regions.
- The amine group ($-NH_2$) is added to the benzene ring at a new position.
- Each yellow triangular face is attached to three metallic rods at its edges.
- The E-flat note is positioned directly below the B-flat note on the staff.

Figure 15: The basic prompt framework for annotating scientific diagrams follows the same structure as that used for statistical diagrams. However, due to the inherent difference between scientific and statistical diagrams, we provide tailored instructions that emphasize features like objects, attributes, and structural hierarchy. We also include in-context examples specific to scientific content.

## Prompt for Scientific Annotation (Step 2: Annotation)

**System:** You are an expert information extraction assistant specializing in scene graph construction. Your task is to analyze a given diagram description and extract meaningful, structured relationships between key elements.

**Context:** The description of the diagram is provided for your reference.

The diagram depicts a liposome used for drug delivery. The central element is a large, circular liposome, predominantly brown-orange, representing a lipid bilayer. Inside the liposome, a light blue aqueous core contains a crystalline structure labeled "Drug crystallized in aqueous fluid" (white and iridescent) and a purple, coiled structure labeled "DNA". Several arrows connect labels to parts of the liposome:

- An arrow points from the text "Protective layer against immune destruction" to the outer edge of the liposome's lipid bilayer, indicating a protective function.

- Arrows point from the text "Lipid-soluble drug in bilayer" to the lipid bilayer itself, indicating the location of lipid-soluble drugs within the bilayer.

- Arrows point from the text "Drug crystallized in aqueous fluid" to the crystalline structure in the aqueous core.

- Arrows point from the text "Lipid bilayer" to the brown-orange lipid bilayer.

Attached to the outer edge of the liposome are several purple, wavy structures labeled "Homing peptide," suggesting a targeting mechanism. The text "Liposome for Drug Delivery" is positioned above the liposome, serving as a title. The overall structure is hierarchical, with the liposome as the central node, and various labels and arrows acting as connected nodes, describing its components and functions.

**Instruction:**
- Identify important relationships between key elements from the description.
- Structure these relationships in the form of triples with three components:
  - **Source**: The primary element (subject) in the relationship.
  - **Relationship**: The type of connection between the source and target.
  - **Target**: The secondary element (object) in the relationship.
- Ensure that:
  - Each triple represents a meaningful connection between elements.
  - The relationships are concise yet descriptive.
  - There are no duplicate, redundant, or meaningless triples.

**Output Format:** The final output must strictly follow the JSON format below:
{
    "1": {"Source": "Triple 1", "Relationship": "Triple 1", "Target": "Triple 1"},
    ...
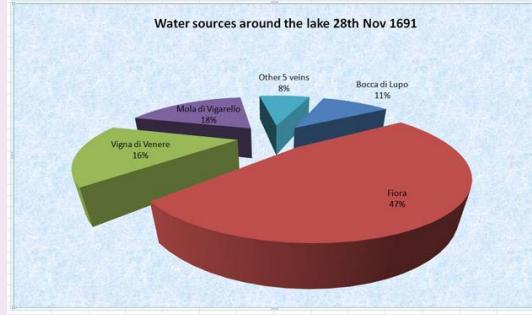    "N": {"Source": "Triple N", "Relationship": "Triple N", "Target": "Triple N"}
}

Figure 16: Similar to statistical diagrams, we provide the model with previously extracted information and ask it to generate a list of triples in JSON format.

**Prompt for QA Annotation (Step 1: Captioning)**

**System:** You are a diagram description assistant.

**Context:** The diagram is sourced from Wikipedia, and here is some background information. Use the Wikipedia information above only if the diagram alone does not provide enough clarity or context. Always give priority to the information directly visible in the diagram for your analysis.

- **Page Title:** Aqua Traiana.
- **Page Description:** The Aqua Traiana was a 1st-century Roman aqueduct built by Emperor Trajan and inaugurated on 24 June 109 AD. It channelled water from sources around Lake Bracciano, 40 kilometers north-west of Rome, to Rome in ancient Roman times but had fallen into disuse by the 17th century. (text truncated due to space)
- **Diagram Description:** None.

**Instruction:** Your task is to provide a detailed description of the diagram, addressing the following four aspects:

- **Recognition:** Identify and describe the key visual elements present in the diagram.
- **Understanding:** Explain the relationships and interactions between these elements.
- **Grounding:** Relate the diagram elements to real-world concepts or entities.
- **Reasoning:** Interpret the diagram to draw conclusions or infer information beyond what is explicitly shown.

**Output Format:** You must output your result in the following JSON-like format:
{
    "Recognition": "string or NA",
    "Understanding": "string or NA",
    "Grounding": "string or NA",
    "Reasoning": "string or NA"
}

Figure 17: Before annotating QA pairs, we prompt the model to caption the diagram. Here we provide relevant Wikipedia text and the definition of the four tasks to instruct the model to generate descriptions specific for QA annotation.

## Prompt for QA Annotation (Step 2: Annotation)

**System:** You are a question-answering annotation assistant. Your task is to analyze a diagram and annotate question-answering pairs.

**Context:** The description of the diagram is provided for your reference.

**1. Recognition:** The diagram is a 3D pie chart showing the distribution of water sources around a lake on November 28th, 1691. The chart is segmented into six sections, each representing a different water source and its percentage contribution. The sections are color-coded for easy identification. The labels for each section indicate the name of the water source and its corresponding percentage. The largest section is labeled "Fiora" and occupies 47% of the chart. The other sections are "Vigna di Venere" (16%), "Mola di Vigarello" (18%), "Bocca di Lupo" (11%), "Other 5 veins" (8%).

**2. Understanding:** The pie chart illustrates the relative proportions of water sourced from different locations around the lake on the specified date. Each slice represents a specific water source, and its size is proportional to its contribution to the total water supply. The chart visually compares the contributions of each source, highlighting the dominance of Fiora with 47% of the total water supply.

**3. Grounding:** The diagram represents the real-world distribution of water sources around a lake (likely Lake Bracciano based on the Wikipedia context) at a specific point in time. The named sources ("Fiora", "Vigna di Venere", etc.) are likely geographical locations or specific water channels feeding into the lake. The percentages represent the proportion of the total water volume coming from each source.

**4. Reasoning:** The chart shows that Fiora is the most significant water source, contributing almost half (47%) of the total water. The remaining sources contribute smaller, but still substantial, portions. This suggests that the water supply relied heavily on Fiora, while other sources played supporting roles. The data could be used to understand the hydrological characteristics of the area and the historical reliance on different water sources. The date (28th Nov 1691) implies this is a snapshot of a specific time and conditions, and the distribution might vary across seasons or years.

Figure 18: After captioning, we provide the model with the extracted information corresponding to the four tasks.

## Prompt for QA Annotation (Step 2: Annotation) (Continued)

**Instruction:** Create four multiple-choice questions based on the diagram, each focusing on one of the four aspects:

- **Recognition**: Test the identification of elements in the diagram.
- **Understanding**: Assess comprehension of the relationships or processes depicted.
- **Grounding**: Evaluate the ability to connect elements to real-world knowledge.
- **Reasoning**: Challenge inference or prediction based on the diagram.

For each question:

- Provide a clear question statement.
- Offer exactly four options labeled A, B, C, and D.
- Indicate the correct answer, which must be only one among A, B, C, or D.

**Output Format:** You must output your result in the following JSON-like format:

{
   "Recognition": {
     "Question": "string",
     "Options": { "A": "string", "B": "string", "C": "string", "D": "string" },
     "Answer": "A/B/C/D"
   },
   "Understanding": {
     "Question": "string",
     "Options": { "A": "string", "B": "string", "C": "string", "D": "string" },
     "Answer": "A/B/C/D"
   },
   "Grounding": {
     "Question": "string",
     "Options": { "A": "string", "B": "string", "C": "string", "D": "string" },
     "Answer": "A/B/C/D"
   },
   "Reasoning": {
     "Question": "string",
     "Options": { "A": "string", "B": "string", "C": "string", "D": "string" },
     "Answer": "A/B/C/D"
   }
}

Figure 19: Using the descriptive information, we instruct the model to generate one multiple-choice question for each of the four tasks. Each question is designed to include exactly four answer options with a single correct answer.

## QA Annotation Example

**Recognition**: What percentage of water sources around the lake on November 28th, 1691, came from Bocca di Lupo?
A) 11%
B) 8%
C) 16%
D) 18%
**Answer:** A

**Understanding**: What does the size of each segment in the pie chart represent?
A) The geographical location of the water source.
B) The depth of the water source.
C) The volume of water contributed by each source relative to the total.
D) The water temperature of each source.
**Answer:** C

**Grounding**: Considering the context of the chart, which of the following is MOST likely a real-world interpretation of Fiora, Vigna di Venere, and Mola di Vigarello?
A) Types of water purification methods used in 1691.
B) Names of individuals who managed the water supply.
C) Geographical locations or water channels supplying the lake.
D) Measurements of water flow rates.
**Answer:** C

**Reasoning**: Based on the data, which of the following conclusions is MOST supported?
A) The water supply was evenly distributed among all sources.
B) Fiora was the primary source of water for the lake.
C) The "Other 5 veins" contributed the least amount of water, rendering them insignificant.
D) Vigna di Venere was the most important water source besides Fiora.
**Answer:** B

Figure 20: Here we illustrate an example of the annotated results, including questions, options, and correct answers.

## Pipeline for Benchmark Evaluation



First generation immigrant population without a
Norwegian background in Norway in 2009

Legend:
- 15,000 +
- 10,000 - 14,999
- 5,000 - 9,999
- 1,000 - 4,999

**Step 1: Analysis**
**Prompt:**
You are a Visual Question Answering expert. Based on the image, analyze and answer the following question:
What color represents countries with 10,000-14,999 first-generation immigrants in 2009?

A) Dark Red.
B) Red.
C) Light Orange.
D) Light Peach.

**Output:**
Okay, I've analyzed the image. Based on the legend provided, countries with 10,000-14,999 first-generation immigrants in Norway in 2009 are represented by **B. Red**. Let me know if you have another question about the image!

**Step 2: Conclusion**
**Prompt:**
Based on the context above, give your final choice. Do not provide any explanations.
**Output:**
The final choice is B.

**Step 3: Choice**
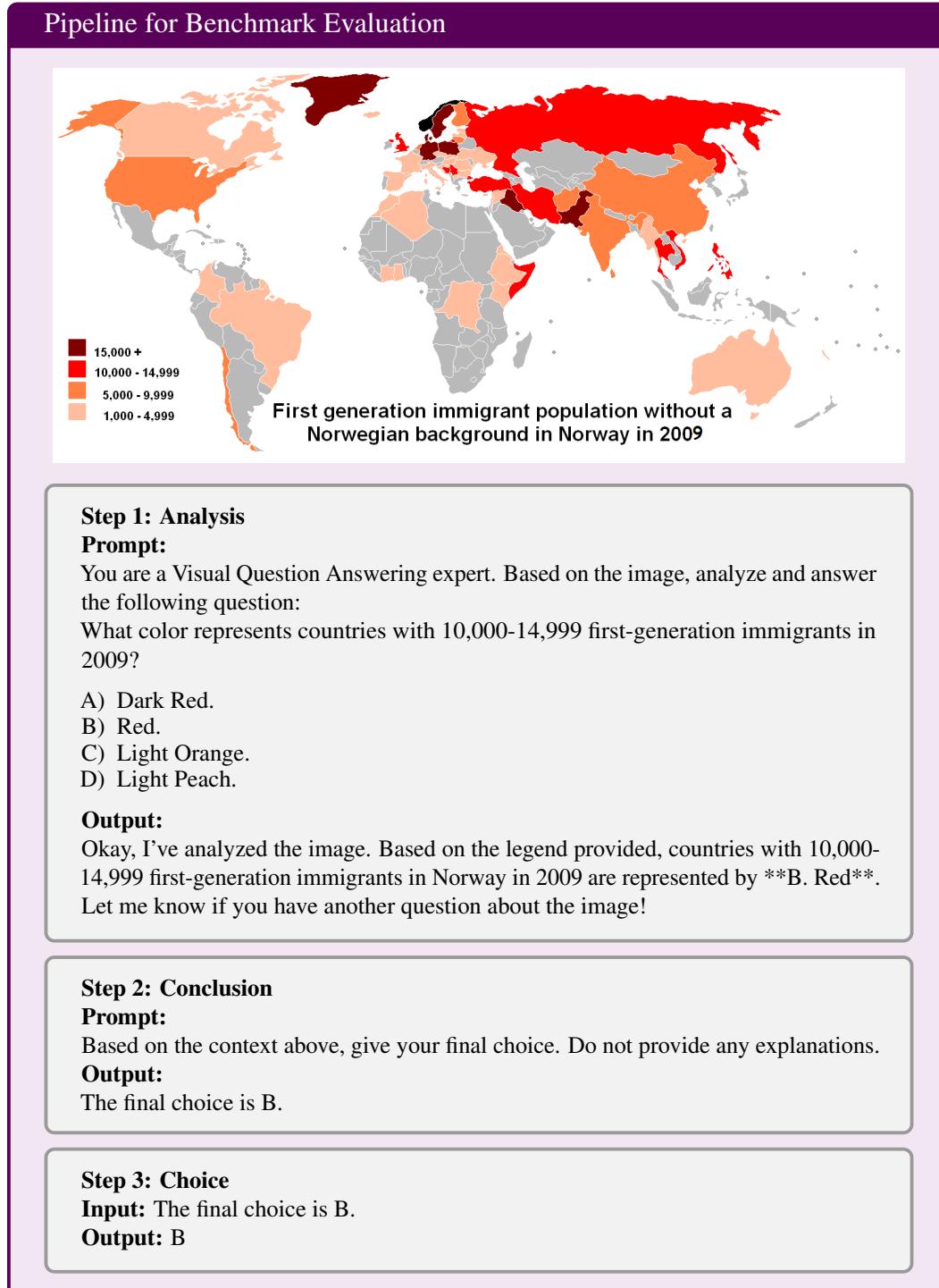**Input:** The final choice is B.
**Output:** B

Figure 21: This figure illustrates the pipeline used for benchmark evaluation. First, the model is prompted to analyze the input data and the question. Next, it is asked to produce a final conclusion based on its analysis. Finally, a response-processing step is applied to extract the answer from the conclusion text.

## BROADER IMPACT

Structured diagram data holds broad potential for advancing multimodal intelligence across both research and applied domains. The semantic annotations in our test suite, particularly the structured triples and multilevel reasoning tasks, can support a variety of downstream applications beyond evaluation. For instance, they can enable better text-to-diagram generation, where structured content such as sentences or knowledge graphs can be translated into meaningful visualizations for education, publishing, or user interfaces. Moreover, the design of our test suite, particularly its explicit separation of reasoning stages and alignment with semiotic principles, can inspire new training paradigms, such as the use of synthetic reasoning trajectories or modality-controlled supervision to improve multimodal model robustness and interpretability. We anticipate that these ideas will generalize to other structured domains, such as scientific visualization, instructional materials, and interactive agents grounded in visual knowledge.

## LIMITATIONS

While we offer a comprehensive test suite for diagram comprehension, several limitations remain. First, our dataset is constructed from Wikipedia diagrams, which, while diverse and high-quality, may not fully represent diagrams used in other domains such as medicine, engineering, or early education. This could limit generalization to domain-specific use cases. Second, although we implement rigorous consistency checks and conduct human evaluation on a subset of the data, automatic annotations, especially for complex reasoning questions, may still contain subtle noise or bias. Finally, while we identify and analyze shortcut behaviors, our diagnostic framework is correlational and does not isolate causal mechanisms behind model behavior. Future work could extend this analysis with counterfactual interventions, synthetic control diagrams, or fine-grained behavioral probing.