

Beyond Accuracy: Unveiling the Cognitive Mechanisms of MLLM Failure in Misleading Visualizations

Anonymous ACL submission

Abstract

While Multimodal Large Language Models demonstrate exceptional proficiency in chart understanding, their robustness against misleading visualizations remains a critical bottleneck. Existing benchmarks predominantly utilize flat taxonomies, conflating visual noise with semantic manipulation, thereby obscuring the deeper cognitive states behind model failures. To address this, we propose CoVis, a benchmark grounded in a four-layer cognitive taxonomy (Perception-Mapping-Reasoning-Logic) to ensure comprehensive coverage of the adversarial landscape. Using our proposed Knowledge-Level Uncertainty (KLU), we identify a systematic Cognitive Bifurcation: model failures collapse into either Cognitive Denial (confusion due to visual obstructions) or Cognitive Hijacking (delusion driven by semantic inducements). Evaluations across 10 state-of-the-art models demonstrate striking cross-model consistency in these failure mechanisms. Notably, we observe a pronounced Textual Dominance effect, where MLLMs often prioritize deceptive text labels over conflicting visual evidence, leading to high-confidence hijacking. Overall, CoVis establishes a mechanism-based evaluation paradigm, shifting the focus from surface-level error correction to the governance of internal cognitive states.

1 Introduction

Data visualizations are essential tools for communicating quantitative evidence and supporting informed decision-making across diverse fields. With the rapid development of Multimodal Large Language Models (MLLMs) (Radford et al., 2021; Alayrac et al., 2022; Team et al., 2023; Achiam et al., 2023; Dubey et al., 2024; Li et al., 2024; Wu et al., 2024), models have demonstrated perception and reasoning capabilities comparable to humans on standard chart understanding tasks (Masry et al., 2022; Xia et al., 2025). However, real-world vi-

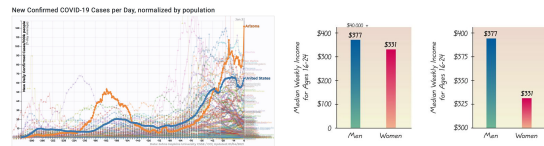


Figure 1: Diverse manifestations of misleading visualizations. (a) **Information Confusion**: Excessive visual noise (e.g., overplotting) that obscures the underlying data. (b) **Intentional Misleading**: Strategic manipulation (e.g., truncated axes) designed to induce incorrect causal conclusions.

sualizations are not always benevolent and well-formed. From distorted axes to deceptive labeling, *Misleading Visualizations* have become potent vehicles for misinformation propagation (King, 1986; Pandey et al., 2015). In high-stakes domains such as finance and healthcare, the ability to resist these visual traps is not merely a performance metric, but a prerequisite for trustworthy decision-making (Lauer and O’Brien, 2020).

Although recent efforts have identified significant performance drops in MLLMs facing chart deception (Chen et al., 2025; Mahbub et al., 2025), existing evaluations typically rely on flat error taxonomies and report only scalar accuracy. Such a broad approach conflates design normative errors with semantic manipulation, thereby failing to distinguish between fundamentally different failure modes, masking the deeper cognitive mechanisms that drive model fragility. As illustrated in Figure 1, these studies predominantly employ a “Flat Taxonomy,” conflating *Design Normative Errors* (e.g., cluttering/confusion, Figure 1a) with *Logical Manipulative Errors* (e.g., truncated axes/misleading, Figure 1b). While the former primarily introduces perceptual noise leading to reading difficulties, the latter constructs deep semantic traps. This evaluation approach, driven solely by scalar Accuracy, obscures the internal state of the model: when a

(2024) pioneered the probing of models’ detection capabilities via prompting. Subsequently, large-scale benchmarks such as *Misleading ChartQA* (Chen et al., 2025), *The Perils of Chart Deception* (Mahbub et al., 2025), *MisViz* (Tonglet et al., 2025b) have quantified the significant performance degradation of state-of-the-art models against various deceptive designs. Furthermore, research has extended to downstream applications, including defense (Tonglet et al., 2025a; Song et al., 2025), interactive correction (Das and Mueller, 2025), and rhetorical intent analysis (Blasilli and Angelini).

Despite confirming the vulnerability of MLLMs, current evaluation paradigms suffer from two critical limitations. First, existing benchmarks predominantly employ a Flat Taxonomy (Chen et al., 2025; Tonglet et al., 2025b), categorizing “Design Normative Errors” (e.g., color misuse) alongside “Logical Fallacies” (e.g., semantic manipulation) as parallel labels without structural distinction. This approach conflates *Visual Noise* (perception level) with *Semantic Manipulation* (logic level). Second, current assessments rely almost exclusively on scalar metrics (e.g., Accuracy), lacking a Cognitive Diagnosis of the model’s internal state. While uncertainty estimation is widely used in NLP for hallucination detection (Kadavath et al., 2022; Huang et al., 2023), it remains unexplored in the context of adversarial charts. Consequently, existing works cannot distinguish between *Cognitive Denial* (a high-entropy state of admitted ignorance due to blocked perception) and *Cognitive Hijacking* (a low-entropy state of confident delusion driven by semantic traps). **CoVis** addresses these limitations by introducing a four-layer cognitive taxonomy and the **Knowledge-Level Uncertainty (KLU)** metric, shifting the focus from performance benchmarks to a deeper analysis of internal failure states.

3 The CoVis Benchmark

Rather than relying on web-scraped images (Lo et al., 2022), which can introduce confounding factors like noise and erratic formatting, we utilize a script-based synthesis framework. This approach leverages plotting libraries to generate charts programmatically, providing direct control over visual parameters. Consequently, we can apply specific misleading designs while ensuring absolute alignment between the visual representation and the ground-truth information.

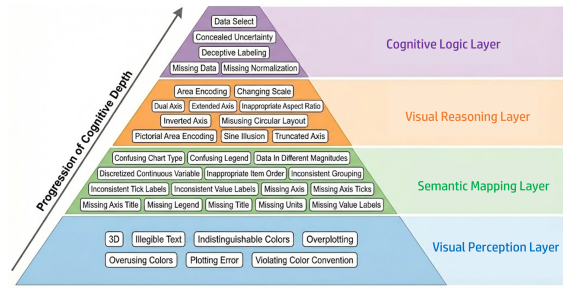


Figure 2: The structural composition of the CoVis four-layer cognitive pipeline. The architecture is organized into four sequential stages: *Perception*, *Mapping*, *Reasoning*, and *Logic*, which categorize 38 core misleading types into a hierarchical system.

3.1 Theoretical Basis: The 4-Layer Cognitive Pipeline

To identify the underlying failure mechanisms triggered by misleading visualizations, we structure our benchmark based on established principles of cognitive psychology and visualization theory. We adapt Pinker’s model of graph comprehension (Pinker, 2014; Kosslyn, 2006; Khalil et al., 2005) to investigate the specific stages where MLLMs become vulnerable to deception. By utilizing a four-layer cognitive pipeline, we can diagnose how perturbations at different levels, from initial signal perception to high-level logical synthesis, result in model failure, allowing for a detailed analysis of the internal processes that drive these errors.

As illustrated in Figure 2 and detailed in Appendix A.1, we curate 38 core error types of visual misleadingness mapped to these four stages:

1. **Visual Perception Layer (Signal Acquisition):** Targets the raw extraction of visual features. We inject physical noise (e.g., overplotting, indistinguishable colors) to sever the link between pixels and objects.
2. **Semantic Mapping Layer (Decoding):** Targets the translation of visual variables (e.g., position, length) into data concepts. We remove or obscure the reference systems (e.g., missing axes, inconsistent labels) required for decoding.
3. **Visual Reasoning Layer (Operation):** Targets the geometric comparison and calculation processes. We introduce spatial distortions (e.g., truncated axes, dual scales) that violate standard visual grammar.

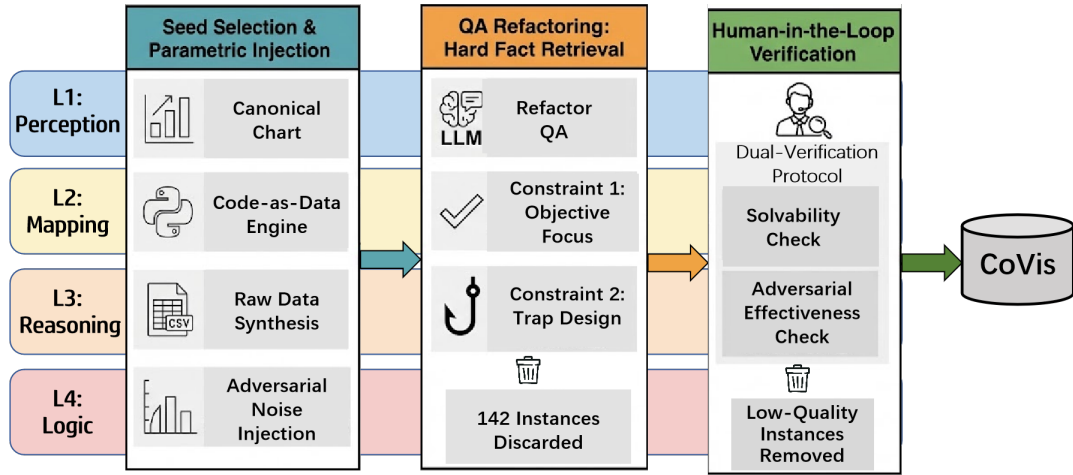


Figure 3: The CoVis data construction pipeline. We utilize a **Code-as-Data** engine to synthesize charts with parametric perturbations, followed by LLM-driven QA refactoring and a stringent human-in-the-loop verification process, resulting in 848 high-quality instances.

- 252 4. **Cognitive Logic Layer (Decision):** Targets
 253 high-level semantic interpretation. We manip-
 254 ulate contextual information (e.g., data select-
 255 ing, concealed uncertainty) to induce logical
 256 fallacies despite correct visual perception.

257 **3.2 Construction Workflow**

258 **3.2.1 Seed Selection & Parametric Injection**

259 Our data generation pipeline begins with the selec-
 260 tion of canonical chart schemas. We ground our
 261 seed charts in the taxonomy defined by VLAT (Lee
 262 et al., 2016), a rigorously validated visualization lit-
 263 eracy assessment framework. This ensures that our
 264 base charts adhere to standard perceptual principles
 265 before perturbation. Leveraging a “Code-as-Data”
 266 paradigm, we developed a Python-based generation
 267 engine utilizing Matplotlib.

268 To ensure a balanced representation across the
 269 four cognitive layers, we initially instantiated a to-
 270 tal of **1,250 candidate chart-data pairs**. For each
 271 instance, the engine performs two synchronized
 272 operations:

- 273 • **Raw Data Synthesis:** It generates underlying
 274 data tables (CSV) with controlled statistical
 275 properties, such as linear trends, clusters, or
 276 seasonal fluctuations.
- 277 • **Adversarial Noise Injection:** It program-
 278 matically alters rendering parameters to cre-
 279 ate specific error types, such as adjusting
 280 `ax.set_ylim()` for *Truncated Axis* or set-
 281 ting `alpha=0.1` for *Indistinguishable Col-
 282 ors*, while maintaining pixel-perfect align-
 283 ment with the ground-truth numerical data.

284 **3.2.2 QA Refactoring: From Meta-Cognition
 285 to Fact Retrieval**

286 A common pitfall in existing benchmarks is the
 287 use of meta-cognitive questions (e.g., “Is this chart
 288 misleading?”). Such framing introduces a “**sus-
 289 picion bias,**” priming the model to look for er-
 290 rors regardless of visual evidence. To eliminate
 291 this confounder, we refactor the task into “**Hard
 292 Fact Retrieval.**” We employ GPT-4o to generate
 293 multiple-choice questions based solely on the raw
 294 data tables. The prompting strategy (detailed in
 295 Appendix A.3) enforces specific constraints:

- 296 • **Objective Focus:** Questions must target spe-
 297 cific numerical data points or categorical
 298 trends rather than asking about the chart’s de-
 299 sign quality.
- 300 • **Trap Design:** Distractor options must explic-
 301 itly include values suggested by the visual
 302 manipulation, such as the perceived value in a
 303 truncated axis or a distorted area encoding, to
 304 act as “lures” for *Cognitive Hijacking*.

305 During this stage, **142 instances** were automati-
 306 cally discarded due to LLM-generated questions
 307 that were either ambiguous or not strictly answer-
 308 able from the CSV data, leaving 1,108 candidates
 309 for human review.

310 **3.2.3 Human-in-the-loop Verification**

311 To ensure the validity of the generated QA pairs, we
 312 implemented a rigorous human-in-the-loop verifi-
 313 cation process via a custom-built Interactive Anno-
 314 tation Platform (see Appendix A.2). Three expert

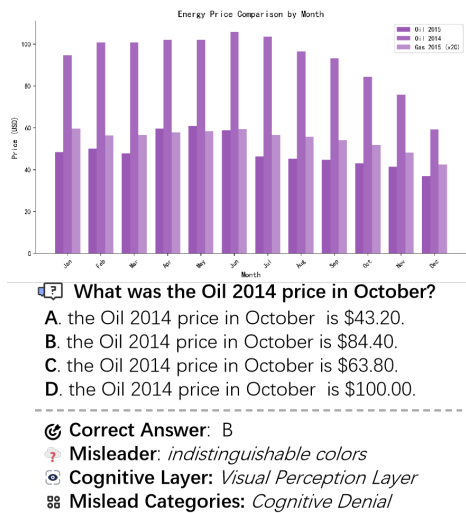


Figure 4: A representative sample from the CoVis benchmark. The figure displays a misleading chart with “Indistinguishable Colors” (L1) and its corresponding adversarial QA pair. The distractors are designed as “lures” based on the visual distortion to probe for Cognitive Denial.

annotators with backgrounds in data visualization performed a **dual-verification protocol**:

- **Solvability Check**: Verifying that the Ground Truth option is strictly derivable from the underlying data table.
- **Adversarial Effectiveness Check**: Confirming that the chart visually suggests the “Trap Option” (for Hijacking types) or renders the data genuinely illegible (for Denial types).

Figure 4 illustrates a representative CoVis instance featuring an “Indistinguishable Colors” error (L1), demonstrating how visual obstructions trigger **Cognitive Denial** via our QA refactoring.

4 Experiment

4.1 Experiment Setup

4.1.1 Evaluated Models & Implementation

Model Selection Criteria To ensure a comprehensive evaluation of the MLLM landscape, we curate a set of 10 state-of-the-art models spanning five distinct architectural families, each represented by a small and a large parameter variant (as detailed in Table 1). This paired design allows us to control for architectural differences while isolating the effects of model scaling on cognitive states:

- **Qwen3-VL (8B & 32B)**: Represents the latest generation of MLLMs, possessing advanced

visual understanding and reasoning capabilities.

- **DeepSeek-VL2 (Tiny & Native)**: Represents models optimized for reasoning, incorporating mixture-of-experts (MoE) and reinforcement learning strategies.
- **Llama-3.2-Vision (11B & 90B)**: A widely adopted, representative open-weights baseline for multimodal tasks.
- **LLaVA-OneVision (7B & 72B)**: Represents the “Strong Encoder” paradigm, utilizing SigLIP-SO400M and a specialized “AnyRes” strategy for high-resolution details.
- **LLaVA-v1.6-Vicuna (7B & 13B)**: Serves as a classic baseline representing the widely adopted LLaVA architectures.

Inference Settings We evaluate all models in a zero-shot setting to strictly probe their intrinsic capabilities without the influence of few-shot demonstrations. To ensure the Knowledge-Level Uncertainty (KLU) metric reflects the model’s true underlying probability distribution rather than sampling noise, we perform inference using a decoding temperature of 0.01 and a nucleus sampling threshold of 0.001. For each multiple-choice question, we extract the raw logits associated with the option tokens (e.g., “A”, “B”, “C”, “D”) and normalize them via Softmax to compute the probability distribution P required for KLU calculation. All experiments were conducted on a cluster of NVIDIA H100 (80GB) GPUs.

4.1.2 Evaluation Metrics

Accuracy We report the standard Accuracy (ACC) to measure the model’s success rate in factual retrieval. While ACC reflects the final utility, we argue it is insufficient for adversarial evaluation: a score of zero conflates “ignorant guessing” with “confident hallucination,” masking the distinct safety risks associated with each.

Knowledge-Level Uncertainty (KLU) Standard scalar metrics treat all incorrect predictions equally, failing to distinguish between a model that is “confused” and one that is “misled.” To decouple these failure modes, we introduce **Knowledge-Level Uncertainty (KLU)**, a diagnostic probe based on the information-theoretic distance between the model’s confidence distribution and random guessing.

Model	L1: Perception	L2: Mapping	L3: Reasoning	L4: Logic	Average
DeepSeek-VL2-Tiny	0.320	0.397	0.732	0.160	0.402
DeepSeek-VL2-Native	0.554	0.580	0.727	0.540	0.600
Llama-3.2-11B-Vision	0.117	0.264	0.939	0.100	0.355
Llama-3.2-90B-Vision	0.221	0.351	0.797	0.440	0.452
LLaVA-v1.6-Vicuna-7B	0.167	0.299	0.801	0.220	0.372
LLaVA-v1.6-Vicuna-13B	0.162	0.383	0.636	0.380	0.390
LLaVA-OneVision-7B	0.608	0.528	0.519	0.420	0.519
LLaVA-OneVision-72B	0.626	0.574	0.710	0.660	0.642
Qwen3-VL-8B	0.685	0.545	0.528	0.580	0.584
Qwen3-VL-32B	0.716	0.597	0.788	0.660	0.690

Table 1: Zero-shot accuracy across cognitive layers. Best performance is in **bold**.

Mathematical Formulation: For a given question with K options (where $K = 4$ in CoVis), let $\mathcal{P} = \{p_1, \dots, p_K\}$ denote the model’s softmax probability distribution over the answer space, and $\mathcal{U} = \{1/K, \dots, 1/K\}$ denote the uniform distribution representing maximum uncertainty. We define KLU as the Kullback-Leibler divergence of \mathcal{P} from \mathcal{U} :

$$\text{KLU}(\mathcal{P}) = D_{KL}(\mathcal{P} \parallel \mathcal{U}) = \ln K - H(\mathcal{P}) \quad (1)$$

where $H(\mathcal{P})$ is the Shannon entropy. Physically, KLU quantifies the information gain of the model’s prediction relative to a random guess.

Rationale & Interpretation: In the context of misleading visualizations, KLU serves as a probe for the model’s certainty landscape when prediction fails:

- **Lower KLU (Distributional Flatness):** This indicates that the model’s output approaches uniformity ($\mathcal{P} \approx \mathcal{U}$). It signifies a state of High Entropy, implying that the model fails to extract discriminative features from the visual input and effectively reverts to random guessing (*Cognitive Denial*).
- **High KLU (Probability Concentration):** This indicates that the probability mass is heavily concentrated on a specific option. When the prediction is incorrect, a high KLU signifies a state of False Certainty, implying that the model has been strongly anchored to a specific wrong answer by the adversarial features (*Cognitive Hijacking*).

By analyzing the distribution of KLU, we can determine whether the robustness failure stems from a lack of discriminative power (*Flatness*) or misaligned confidence (*Concentration*).

4.2 Main Results

4.2.1 Overall Vulnerability

We first assess the absolute robustness of MLLMs against the CoVis benchmark. As presented in Table 1, the results reveal a significant performance gap compared to standard chart evaluation benchmarks. The average accuracy across all 10 evaluated models is merely **50.1%**, indicating that nearly half of the adversarial samples successfully bypass current model defenses. Even the state-of-the-art Qwen3-VL-32B achieves only 69.0% accuracy, leaving a substantial 31% robustness gap. More critically, widely used foundational models such as Llama-3.2-11B-Vision (35.5%) and LLaVA-v1.6-7B (37.2%) perform dangerously close to the random guessing baseline (25%). This underscores that robustness against visual misleadingness remains a systemic blind spot, irrespective of the model family.

4.2.2 The Impact of Scale

Analyzing the performance delta between paired small and large models reveals that scaling parameters does not benefit all cognitive layers symmetrically (see Figure 5):

- **Logic Benefits from Scale:** The most substantial gains are observed in Layer 4. For instance, scaling DeepSeek-VL2 from Tiny to Native boosts Layer 4 accuracy by **+38.0%**, and Llama-3.2-Vision from 11B to 90B by **+34.0%**.
- **Perception Hits a Ceiling:** Conversely, gains in Layer 1 are marginal or inconsistent. While LLaVA-OneVision gains only +1.8% when scaling from 7B to 72B, LLaVA-v1.6 sees a slight regression (-0.5%). This suggests a “Perceptual Ceiling” in current architectures.

Error Type	L	State	Score	Error Type	L	State	Score
Missing Normalization	4	Hijack (T)	1.0	inapprop.itemorder	2	Hijack (T)	0.7
overplotting	1	Denial (B)	1.0	changingscale	3	Denial (B)	0.7
discretized-cont.var	2	Hijack (T)	1.0	plottingerror	1	Hijack (T)	0.7
inconsist.valuelabels	2	Hijack (T)	0.9	misuse.circular	3	Denial (B)	0.7
inconsist.ticklabels	2	Denial (B)	0.9	missinglegend	2	Hijack (T)	0.7
Deceptive Labeling	4	Hijack (T)	0.9	truncatedaxis	3	Hijack (T)	0.6
Missing Data	4	Hijack (T)	0.9	sineillusion	3	Denial (B)	0.6
missingaxisticks	2	Denial (B)	0.9	areaencoding	3	Denial (B)	0.6
missingaxis	2	Denial (B)	0.9	missingtitle	2	Hijack (T)	0.6
extendedaxis	3	Denial (B)	0.9	missingabbrev.	2	Denial (B)	0.6
Data Selecting	4	Hijack (T)	0.9	dualaxis	3	Hijack (T)	0.6
overusingcolors	1	Denial (B)	0.8	violating.color	1	Denial (B)	0.6
missingvaluelabels	2	Denial (B)	0.8	confusingcharttype	2	Ambiguous	0.5
missingunits	2	Hijack (T)	0.8	inconsist.grouping	2	Ambiguous	0.5
invertedaxis	3	Hijack (T)	0.8	missingaxistitle	2	Ambiguous	0.5
3d	1	Denial (B)	0.8	confusinglegend	2	Ambiguous	0.5
inapprop.aspectratio	3	Denial (B)	0.8	data-diff.mag.	2	Ambiguous	0.5
illegibletext	1	Denial (B)	0.8	pictorialareaenc.	3	Ambiguous	0.5
Concealed Uncertainty	4	Hijack (T)	0.8	indisting.colors	1	Ambiguous	0.5

Table 2: KLU Rank Consistency across 38 error types. **State** indicates the dominant mechanism (Hijacking vs. Denial). **Score** is the Consistency Score across 10 models.

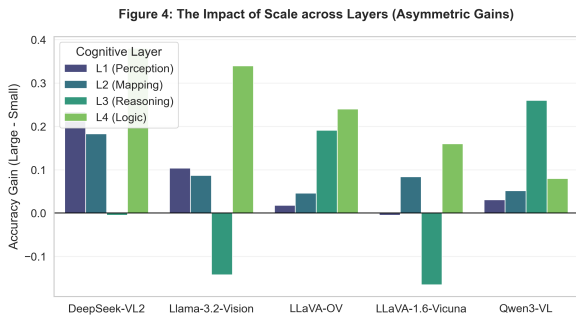


Figure 5: Performance delta across cognitive layers when scaling from small to large model variants. Scaling enhances logical synthesis (L4) but offers diminishing returns for perceptual robustness (L1).

459 However, we observe a counter-intuitive Reasoning
460 Regression in specific architectures. Notably,
461 scaling Llama-3.2-Vision from 11B to 90B results
462 in a significant performance drop in Layer 3 (Visual
463 Reasoning), falling from 93.9% to 79.7% (-14.2%).
464 We hypothesize this represents a ‘‘Curse of Prior
465 Knowledge’’: smaller models may rely on naive,
466 pixel-level geometric measurements, while larger
467 models possess stronger priors that may override
468 immediate visual evidence when facing counter-
469 intuitive designs like inverted axes.

4.3 Cognitive Mechanism Analysis

4.3.1 Universality of Failure Modes

470 A fundamental question is whether failure modes
471 are model-specific idiosyncrasies or universal phe-
472 nomena. As detailed in Table 2, our analysis re-
473
474

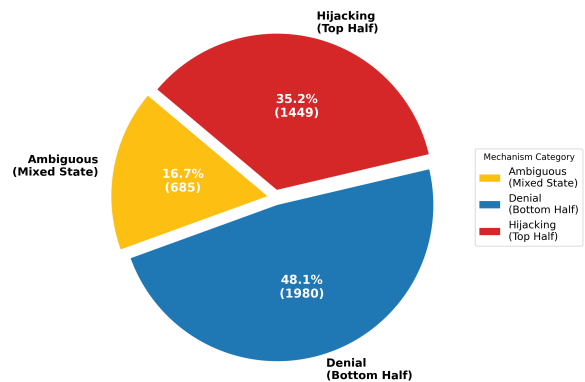


Figure 6: Prevalence of cognitive failure mechanisms. The distribution validates that models under adversarial settings tend to collapse into two states: Denial or Hijacking.

475 veals a striking universality: **81.5%** of error types
476 achieve a Consistency Score of ≥ 0.6 . This indi-
477 cates that regardless of architecture or scale, spe-
478 cific visual triggers activate the same underlying
479 uncertainty state across the MLLM family.

4.3.2 Defining the Binary Mechanisms

480 Based on the observed consistency, we formally
481 define two orthogonal cognitive states that govern
482 MLLM failures:
483

- 484 • **Cognitive Denial:** Characterized by low KLU
485 rankings. Visual noise acts as an *Information*
486 *Blocker*, forcing the model into a high-entropy
487 state of admitted ignorance.
- 488 • **Cognitive Hijacking:** Characterized by high

489	KLU rankings. Adversarial features act as	the theoretical layer, the cognitive state of a mis-	538
490	<i>Semantic Decoys</i> , anchoring the model to a	led model predominantly collapses into one of two	539
491	wrong answer with false certainty.	orthogonal attractors: <i>Cognitive Denial</i> or <i>Cogni-</i>	540
492	As shown in Figure 6, the failure landscape is	<i>tive Hijacking</i> . We argue that the true boundary	541
493	heavily polarized, confirming that “hesitation” is	of model robustness lies in this cognitive conse-	542
494	rare in current MLLMs; they are either universally	quence, not the error stage. A model that remains	543
495	blinded or universally deluded.	highly confident when misled (<i>Hijacking</i>) is fun-	544
496		damentally more hazardous than one that admits	545
497	5 Discussion	ignorance (<i>Denial</i>). Therefore, future evaluations	546
498	5.1 The Textual Dominance in Multimodal	must shift focus from <i>where</i> the error occurs to <i>how</i>	547
499	Conflicts	the model reacts, establishing Cognitive Stability	548
500	A critical insight from our consistency analysis	as the new standard for safety alignment research.	549
501	(Table 2) is exemplified by the behavior of <i>Incon-</i>		
502	<i>sistent Value Labels</i> . In this error type, the visual	5.3 Implications for Adversarial Alignment	550
503	signal (the geometric height of a bar) directly con-	The striking cross-model consistency in hijacking	551
504	tradicts the semantic signal (the text label on the	mechanisms suggests that current RLHF (Rein-	552
505	axis). Despite this blatant visual mismatch, this	forcement Learning from Human Feedback) may	553
506	error consistently triggers Cognitive Hijacking	be insufficient for visual reasoning. Most align-	554
507	(High KLU, Top-Half Distribution) across 90% of	ment processes focus on linguistic safety, yet our	555
508	the evaluated models.	results show that geometric deception can easily	556
509	Rather than being an outlier, this phenomenon	bypass these safeguards. Future alignment should	557
510	empirically validates a fundamental Modality	incorporate "Adversarial Chart Reasoning" into the	558
511	Bias(Wu et al., 2025; Zheng et al., 2025) inher-	fine-tuning loop, specifically rewarding models that	559
512	ent in current MLLM architectures: <i>Textual Dom-</i>	flag visual-textual contradictions rather than those	560
513	<i>inance</i> . When faced with conflicting information	that simply provide a confident (but wrong) answer.	561
514	from dual modalities, models systematically prior-		
515	itize the textual path, effectively bypassing visual	6 Conclusion	562
516	verification. This suggests that existing MLLMs	In this work, we introduce CoVis , a diagnostic	563
517	are influenced by language priors and tend to ig-	benchmark designed to evaluate the cognitive ro-	564
518	gnore visual features. The visual encoder acts as	bustness of MLLMs against misleading visualiza-	565
519	a subordinate feature extractor, susceptible to be-	tions. By shifting from flat taxonomies to a 4-layer	566
520	ing overridden by any coherent, even if deceptive,	hierarchical framework, we provide a more struc-	567
521	semantic narrative found in the text. This “blind	tured map of the adversarial landscape. Our evalua-	568
522	trust” in text is a primary engine driving <i>Cognitive</i>	tion of 10 state-of-the-art models reveals a systemic	569
523	<i>Hijacking</i> , confirming that the model’s certainty is	fragility: even highly capable models are suscep-	570
524	anchored more in what it reads than in what it sees.	tible to Cognitive Hijacking, a state where decep-	571
525		tive semantic cues override visual evidence with	572
526	5.2 From Error Categories to Cognitive	high confidence. The introduction of Knowledge-	573
527	Consequences	Level Uncertainty (KLU) allows us to bridge the	574
528	Existing benchmarks predominantly rely on <i>Flat</i>	gap between error types and failure mechanisms,	575
529	<i>Taxonomies</i> , treating diverse error types as iso-	uncovering a fundamental Cognitive Bifurcation	576
530	lated labels. To address this lack of structure,	between confusion and delusion. Furthermore, the	577
531	CoVis adopts a 4-Layer Hierarchical Taxonomy	observed Textual Dominance effect highlights a	578
532	(Perception-Mapping-Reasoning-Logic), providing	critical bottleneck in current multimodal architec-	579
533	ontological completeness that ensures full coverage	tures. Ultimately, CoVis advocates for a paradigm	580
534	of the visual processing pipeline.	shift in evaluation, moving beyond scalar accuracy	581
535	However, our empirical findings reveal that	toward mechanism-based cognitive alignment. We	582
536	while this layered approach is structurally supe-	hope this work provides the diagnostic depth neces-	583
537	rior, it shares a critical limitation: it categorizes the	sary to develop the next generation of cognitively	584
	input trigger rather than the diagnostic outcome.	reliable visual agents.	585
	Our KLU analysis demonstrates that regardless of		

586 Limitations

587 Despite the comprehensive scope of the **CoVis**
588 benchmark and the diagnostic depth provided by
589 the **KLU** metric, we acknowledge several con-
590 straints that offer avenues for future research:

- 591 • **Scope of Visualization Diversity:** While Co-
592 Vis curates 38 core misleading types across
593 four cognitive layers, it does not yet encom-
594 pass highly specialized or domain-specific vi-
595 sualizations (e.g., 3D scatter plots in bioin-
596 formatics or complex interactive dashboards).
597 Future iterations should expand to these niche
598 categories to test the generalizability of the
599 “Cognitive Bifurcation” theory.
- 600 • **Metric Constraints (KLU):** The calculation
601 of Knowledge-Level Uncertainty (KLU) re-
602 quires access to the model’s raw output logits.
603 While this is straightforward for open-weights
604 models and many research APIs, it poses chal-
605 lenges for certain highly restricted black-box
606 systems that provide only text strings without
607 log-probabilities. Approximation methods for
608 these systems remain to be explored.

609 Ethics Statement

610 This work does not involve the collection or use
611 of any human-related data. All materials are chart-
612 based and synthesized through our proposed gener-
613 ation engine. To support reproducibility and open
614 science, the CoVis benchmark will be released un-
615 der the Creative Commons Attribution 4.0 Inter-
616 national (CC-BY 4.0) license, and the generation
617 code will be made available under the Apache 2.0
618 license. We have ensured that the adversarial per-
619 turbations and synthesized content do not introduce
620 offensive, biased, or harmful visual-language asso-
621 ciations. No ethical concerns were identified in the
622 preparation of this dataset or the subsequent study.

623 References

624 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama
625 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
626 Diogo Almeida, Janko Altenschmidt, Sam Altman,
627 Shyamal Anadkat, and 1 others. 2023. Gpt-4 techni-
628 cal report. *arXiv preprint arXiv:2303.08774*.

629 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc,
630 Antoine Miech, Iain Barr, Yana Hasson, Karel
631 Lenc, Arthur Mensch, Katherine Millican, Malcolm
632 Reynolds, and 1 others. 2022. Flamingo: a visual
633 language model for few-shot learning. *Advances in*

neural information processing systems, 35:23716–
23736. 634 635

Jason Alexander, Priyal Nanda, Kai-Cheng Yang, and
Ali Sarvghad. 2024. Can gpt-4 models detect mis-
leading visualizations? In *2024 IEEE Visualization
and Visual Analytics (VIS)*, pages 106–110. IEEE. 636 637 638 639

Graziano Blasilli and Marco Angelini. True (vis) lies: A
preliminary analysis of how generative ai is capable
of recognizing visualization lies and their rhetoric. 640 641 642

Zixin Chen, Sicheng Song, Kashun Shum, Yanna Lin,
Rui Sheng, Weiqi Wang, and Huamin Qu. 2025.
Unmasking deceptive visuals: Benchmarking mul-
timodal large language models on misleading chart
question answering. In *Proceedings of the 2025 Con-
ference on Empirical Methods in Natural Language
Processing*, pages 13767–13800. 643 644 645 646 647 648 649

Amit Kumar Das and Klaus Mueller. 2025. Misvisfix:
An interactive dashboard for detecting, explaining,
and correcting misleading visualizations using large
language models. *arXiv preprint arXiv:2508.04679*. 650 651 652 653

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,
Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,
Akhil Mathur, Alan Schelten, Amy Yang, Angela
Fan, and 1 others. 2024. The llama 3 herd of models.
arXiv preprint arXiv:2407.21783. 654 655 656 657 658

Muye Huang, Han Lai, Xinyu Zhang, Wenjun Wu, Jie
Ma, Lingling Zhang, and Jun Liu. 2025. Evochart:
A benchmark and a self-training approach towards
real-world chart understanding. In *Proceedings of
the AAAI Conference on Artificial Intelligence*, vol-
ume 39, pages 3680–3688. 659 660 661 662 663 664

Yuheng Huang, Jiayang Song, Zhijie Wang, Shengming
Zhao, Huaming Chen, Felix Juefei-Xu, and Lei Ma.
2023. Look before you leap: An exploratory study of
uncertainty measurement for large language models.
arXiv preprint arXiv:2307.10236. 665 666 667 668 669

Saurav Kadavath, Tom Conerly, Amanda Askell, Tom
Henighan, Dawn Drain, Ethan Perez, Nicholas
Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli
Tran-Johnson, and 1 others. 2022. Language mod-
els (mostly) know what they know. *arXiv preprint
arXiv:2207.05221*. 670 671 672 673 674 675

Mohammed K Khalil, Fred Paas, Tristan E Johnson,
and Andrew F Payer. 2005. Design of interactive and
dynamic anatomical visualizations: the implication
of cognitive load theory. *The Anatomical Record Part
B: The New Anatomist: An Official Publication of the
American Association of Anatomists*, 286(1):15–20. 676 677 678 679 680 681

Gary King. 1986. How not to lie with statistics: Avoid-
ing common mistakes in quantitative political sci-
ence. *American Journal of Political Science*, pages
666–687. 682 683 684 685

Stephen M Kosslyn. 2006. *Graph design for the eye
and mind*. Oup usa. 686 687

688	Claire Lauer and Shaun O'Brien. 2020. The deceptive potential of common design tactics used in data visualizations. In <i>Proceedings of the 38th ACM International Conference on Design of Communication</i> , pages 1–9.		
689			
690			
691			
692			
693	Sukwon Lee, Sung-Hee Kim, and Bum Chul Kwon. 2016. Vlat: Development of a visualization literacy assessment test. <i>IEEE transactions on visualization and computer graphics</i> , 23(1):551–560.		
694			
695			
696			
697	Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, and 1 others. 2024. Llava-onevision: Easy visual task transfer. <i>arXiv preprint arXiv:2408.03326</i> .		
698			
699			
700			
701			
702	Maxim Lisnic, Cole Polychronis, Alexander Lex, and Marina Kogan. 2023. Misleading beyond visual tricks: How people actually lie with charts. In <i>Proceedings of the 2023 CHI conference on human factors in computing systems</i> , pages 1–21.		
703			
704			
705			
706			
707	Leo Yu-Ho Lo, Ayush Gupta, Kento Shigyo, Aoyu Wu, Enrico Bertini, and Huamin Qu. 2022. Misinformed by visualization: What do we learn from misinformative visualizations? In <i>Computer Graphics Forum</i> , volume 41, pages 515–525. Wiley Online Library.		
708			
709			
710			
711			
712	Leo Yu-Ho Lo and Huamin Qu. 2024. How good (or bad) are llms at detecting misleading visualizations? <i>IEEE Transactions on Visualization and Computer Graphics</i> .		
713			
714			
715			
716	Ridwan Mahbub, Mohammed Saidul Islam, Md Tahmid Rahman Laskar, Mizanur Rahman, Mir Tafseer Nayeem, and Enamul Hoque. 2025. The perils of chart deception: How misleading visualizations affect vision-language models. <i>arXiv preprint arXiv:2508.09716</i> .		
717			
718			
719			
720			
721			
722	Ahmed Masry, Xuan Long Do, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. 2022. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. In <i>Findings of the association for computational linguistics: ACL 2022</i> , pages 2263–2279.		
723			
724			
725			
726			
727			
728	Ahmed Masry, Mohammed Saidul Islam, Mahir Ahmed, Aayush Bajaj, Firoz Kabir, Aaryaman Kartha, Md Tahmid Rahman Laskar, Mizanur Rahman, Shadikur Rahman, Mehrad Shahmohammadi, and 1 others. 2025. Chartqapro: A more diverse and challenging benchmark for chart question answering. <i>arXiv preprint arXiv:2504.05506</i> .		
729			
730			
731			
732			
733			
734			
735	Nitesh Methani, Pritha Ganguly, Mitesh M Khapra, and Pratyush Kumar. 2020. Plotqa: Reasoning over scientific plots. In <i>Proceedings of the IEEE/CVF winter conference on applications of computer vision</i> , pages 1527–1536.		
736			
737			
738			
739			
740	Anshul Vikram Pandey, Katharina Rall, Margaret L Satterthwaite, Oded Nov, and Enrico Bertini. 2015. How deceptive are deceptive visualizations? an empirical		
741			
742			
		analysis of common distortion techniques. In <i>Proceedings of the 33rd annual acm conference on human factors in computing systems</i> , pages 1469–1478.	743
			744
			745
	Steven Pinker. 2014. A theory of graph comprehension. In <i>Artificial intelligence and the future of testing</i> , pages 73–126. Psychology Press.		746
			747
			748
	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, and 1 others. 2021. Learning transferable visual models from natural language supervision. In <i>International conference on machine learning</i> , pages 8748–8763. PMLR.		749
			750
			751
			752
			753
			754
			755
	Sicheng Song, Yanjie Zhang, Zixin Chen, Huamin Qu, Changbo Wang, and Chenhui Li. 2025. Vizdefender: Unmasking visualization tampering through proactive localization and intent inference. <i>arXiv preprint arXiv:2512.18853</i> .		756
			757
			758
			759
			760
	Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, and 1 others. 2023. Gemini: a family of highly capable multimodal models. <i>arXiv preprint arXiv:2312.11805</i> .		761
			762
			763
			764
			765
			766
	Jonathan Tonglet, Tinne Tuytelaars, Marie-Francine Moens, and Iryna Gurevych. 2025a. Protecting multimodal large language models against misleading visualizations. <i>arXiv preprint arXiv:2502.20503</i> .		767
			768
			769
			770
	Jonathan Tonglet, Jan Zimny, Tinne Tuytelaars, and Iryna Gurevych. 2025b. Is this chart lying to me? automating the detection of misleading visualizations. <i>arXiv preprint arXiv:2508.21675</i> .		771
			772
			773
			774
	Zirui Wang, Mengzhou Xia, Luxi He, Howard Chen, Yitao Liu, Richard Zhu, Kaiqu Liang, Xindi Wu, Haotian Liu, Sadhika Malladi, and 1 others. 2024. Chartx: Charting gaps in realistic chart understanding in multimodal llms. <i>Advances in Neural Information Processing Systems</i> , 37:113569–113697.		775
			776
			777
			778
			779
			780
	Huyu Wu, Meng Tang, Xinhan Zheng, and Haiyun Jiang. 2025. When language overrules: Revealing text dominance in multimodal large language models. <i>arXiv preprint arXiv:2508.10552</i> .		781
			782
			783
			784
	Zhiyu Wu, Xiaokang Chen, Zizheng Pan, Xingchao Liu, Wen Liu, Damai Dai, Huazuo Gao, Yiyang Ma, Chengyue Wu, Bingxuan Wang, and 1 others. 2024. Deepseek-vl2: Mixture-of-experts vision-language models for advanced multimodal understanding. <i>arXiv preprint arXiv:2412.10302</i> .		785
			786
			787
			788
			789
			790
	Renqiu Xia, Hancheng Ye, Xiangchao Yan, Qi Liu, Hongbin Zhou, Zijun Chen, Botian Shi, Junchi Yan, and Bo Zhang. 2025. Chartx & chartvlm: A versatile benchmark and foundation model for complicated chart reasoning. <i>IEEE Transactions on Image Processing</i> .		791
			792
			793
			794
			795
			796

797 Zhengzhuo Xu, Sinan Du, Yiyan Qi, Chengjin Xu, Chun
798 Yuan, and Jian Guo. 2023. Chartbench: A bench-
799 mark for complex visual reasoning in charts. *arXiv*
800 *preprint arXiv:2312.15915*.

801 Xu Zheng, Chenfei Liao, Yuqian Fu, Kaiyu Lei, Yuan-
802 huiyi Lyu, Lutao Jiang, Bin Ren, Jialei Chen, Jiawen
803 Wang, Chengxin Li, and 1 others. 2025. Mllms are
804 deeply affected by modality bias. *arXiv preprint*
805 *arXiv:2505.18657*.

A Appendix

A.1 Detailed Taxonomy of Misleading Visualizations

Layer	Error Type	Definition and Cognitive Impact
L1	3D Effect	Using 3D perspective to represent 2D data, causing volume/area distortion.
	Illegible Text	Text rendered at extremely low resolution or overlapping, making it unreadable.
	Indistinguishable Colors	Using color scales with insufficient contrast, making categories visually merged.
	Overplotting	Excessive density of data points or lines that obscures the underlying distribution.
	Overusing Colors	Employing too many colors for categorical data, exceeding cognitive load limits.
	Plotting Error	Raw rendering failures (e.g., misaligned bars) that contradict source data.
	Violating Color Conv.	Using colors counter-intuitively (e.g., red for "good", blue for "hot").
L2	Confusing Chart Type	Using an inappropriate chart (e.g., pie chart for time series) to obscure trends.
	Confusing Legend	Legends that use ambiguous symbols or are placed far from data elements.
	Data in Diff. Mag.	Vast scale differences without clear indication, hiding small values.
	Discretized Cont. Var.	Converting continuous data into arbitrary bins to hide specific fluctuations.
	Inappropriate Item Order	Non-logical ordering of categorical data to imply a non-existent ranking.
	Inconsistent Grouping	Varying the grouping criteria across different parts of the visualization.
	Inconsistent Ticks	Non-uniform intervals between axis ticks, distorting the perceived rate of change.
	Inconsistent Labels	Textual labels on data points that contradict their actual geometric position.
	Missing Abbreviation	Non-standard abbreviations that hinder semantic decoding.
	Missing Axis	Omitting X or Y axes, leaving data points without a spatial reference frame.
	Missing Axis Ticks	Removing tick marks, making it impossible to retrieve precise numerical values.
	Missing Axis Title	Omitting labels for axes, leading to ambiguity in measurement.
	Missing Legend	Omitting the key for colors or shapes, severing the link between visual and data.
Missing Title	Omitting the main title, depriving the model of the chart's overall context.	
Missing Units	Omitting units (e.g., \$, kg), leading to magnitude confusion.	
Missing Value Labels	Omitting direct labels in cluttered charts where axis lookup is difficult.	
L3	Area Encoding	Using area to represent data where the scale is non-linear.
	Changing Scale	Dynamically altering the scale within a single axis (e.g., linear to log).
	Dual Axis	Using two different Y-axes to imply a false correlation.
	Extended Axis	Stretching axes unnecessarily to flatten trends and minimize change.
	Inapprop. Aspect Ratio	Distorting height/width ratio to exaggerate or dampen slopes.
	Inverted Axis	Reversing axis direction to flip the perceived trend.
	Misusing Circular	Using radar/pie charts for data that does not sum to a meaningful whole.
	Pictorial Area	Using icons of varying sizes where area vs. height growth is confusing.
	Sine Illusion	Optical illusion where vertical distance between curves is misperceived.
Truncated Axis	Starting Y-axis at a non-zero value to exaggerate small differences.	
L4	Data Selecting	Selective data presentation to support a biased or false narrative.
	Concealed Uncertainty	Presenting estimates as absolute facts without indicating error margins.
	Deceptive Labeling	Using biased language in titles/labels to prime a specific conclusion.
	Missing Data	Strategically omitting unfavorable time periods or data categories.
	Missing Normalization	Comparing raw counts for groups of different sizes, causing unfair bias.

Table 3: Comprehensive Definitions of the 38 Misleading Types.

A.2 Human-in-the-loop Verification Details

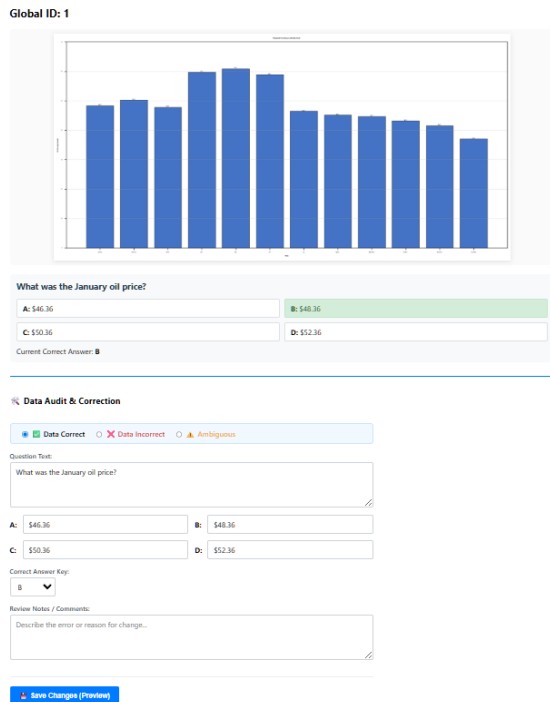


Figure 7: The custom-built Interactive Annotation Platform for CoVis.

To ensure the high quality of the CoVis benchmark, three expert annotators followed a standardized protocol. The following **Annotator Instructions** were provided to each participant via the platform’s briefing interface:

1. **Data Fidelity Audit:** Verify that the geometric elements in the chart (e.g., bar heights, line slopes) strictly correspond to the numerical values in the provided CSV file. Flag any rendering hallucinations or misalignments.
2. **Question Logic Check:** Ensure the question is a “Hard Fact” retrieval task. It must be answerable using only the raw data, without requiring prior external knowledge or asking about the chart’s aesthetic quality.
3. **Adversarial Trap Validation:**
 - Does the visual deception (e.g., truncated axis) naturally lead a human viewer to the designated *Trap Option*?
 - Is the *Correct Option* unambiguously clear once the deception is identified?
4. **Taxonomy Classification:** Confirm that the instance correctly exemplifies one of the 38 misleading types in our taxonomy. If the deception

is multi-layered, assign it to the most dominant cognitive layer.

5. **Option Refactoring:** If the LLM-generated distractors are too easy or unrelated to the visual trick, manually edit them to enhance their “Cognitive Hijacking” potential.

Ethics and Data Consent: All verification steps were conducted by researchers who were informed of the project’s goals and consented to the inclusion of their feedback in the public benchmark. **Regarding the source material, all visualizations and underlying numerical data in CoVis were synthetically generated using our proprietary engine, ensuring no infringement on personal privacy or third-party data rights.** No personal identifiable information (PII) was collected.

Annotator Compensation: The verification process was conducted by three expert researchers who are co-authors of this paper. Therefore, the task was performed as part of their professional research contribution; no external recruitment or financial compensation was involved.

A.3 Prompts for QA Refactoring

System Prompt for QA Generation

You are an expert data analyst and cognitive psychology researcher. Your task is to refactor raw CSV data into adversarial QA pairs based on a specific misleading visualization type [TYPE].

Instructions:

1. **Adversarial Intent:** Reflect the misleading purpose of [TYPE].
2. **Format:** Multiple-choice (4 options) per *filepath* standard.
3. **Objective Probing:** Ask about data content (values, trends), not the misleading technique.
4. **Distractor Design:**
 - **Correct:** Precise calculation from ground-truth.
 - **Lure:** Answer likely chosen if deceived by visual distortion.
5. **Metadata:** Include *Target Misleading Option*.