Both Text and Images Leaked! A Systematic Analysis of Data Contamination in Multimodal LLMs

Dingjie Song^{†12} Sicheng Lai^{†1} Mingxuan Wang¹ Shunian Chen¹ Lichao Sun² Benyou Wang¹

Abstract

The rapid advancement of multimodal large language models (MLLMs) has significantly enhanced performance across benchmarks. However, data contamination-unintentional memorization of benchmark data during model training-poses critical challenges for fair evaluation. Existing detection methods for unimodal large language models (LLMs) are inadequate for MLLMs due to multimodal data complexity and multiphase training. We systematically analyze multimodal data contamination using our analytical framework, MM-DETECT, which defines two contamination categories-unimodal and crossmodal-and effectively quantifies contamination severity across multiple-choice and caption-based Visual Question Answering tasks. Evaluations on twelve MLLMs and five benchmarks reveal significant contamination, particularly in proprietary models and older benchmarks. Crucially, contamination sometimes originates during unimodal pre-training rather than solely from multimodal fine-tuning. Our insights refine contamination understanding, guiding evaluation practices and improving multimodal model reliability.

1. Introduction

The development of MLLMs has exceeded expectations (Liu et al., 2023a; Lin et al., 2023), showcasing extraordinary performance on various multimodal benchmarks (Lu et al., 2022; Liu et al., 2023b; Song et al., 2024), even surpassing human performance. However, due to the partial obscurity associated with MLLMs training (OpenAI, 2023; Reid et al., 2024), it remains challenging to definitively ascertain the impact of training data on model performance, despite some works showing the employment of the training set of certain datasets (Liu et al., 2023a; Chen et al., 2023; Bai et al., 2023b). The issue of data contamination, occurring when training or test data of benchmarks is exposed during the model training phase (Xu et al., 2024), could potentially instigate inequitable performance comparisons among models. This not only creates a dilemma for users in model selection but also poses a significant hurdle to further advancements in this domain.

Existing contamination detection methods primarily focus on LLMs (Yeom et al., 2018; Deng et al., 2024; Dong et al., 2024), showing limitations when applied to MLLMs, due to their multimodal data complexity and multi-stage training processes (Liu et al., 2023a; Li et al., 2023). Thus, systematic analytical frameworks tailored explicitly for multimodal contamination are urgently needed.

In this study, we address three key questions:

- How can we effectively quantify and detect multimodal data contamination?
- What is the degree of contamination across different MLLMs and benchmark datasets?
- When is contamination predominantly introduced—during unimodal pre-training or multimodal fine-tuning?

To comprehensively answer these questions, we first define **Multimodal Data Contamination**, as it pertains to the modality of data sources exposed to the MLLMs, into two categories: *Unimodal Contamination* and *Cross-modal Contamination*. Subsequently, we unveil a detection framework designed explicitly as an analytical tool, **MM-DETECT**, which incorporates two methods, *Option Order Sensitivity Test* and *Slot Guessing for Perturbed Caption*, designed to handle two common types of Visual Question Answering (VQA) tasks: multiple-choice and caption-based questions, respectively.

To corroborate the validity and sensitivity of our approach, we deliberately induce contamination in MLLMs, simulat-

^{*}Equal contribution ¹School of Data Science, The Chinese University of Hong Kong, Shenzhen, Guangdong, China ²Department of Computer Science and Engineering, Lehigh University, Bethlehem, Pennsylvania, USA. Correspondence to: Benyou Wang <wangbenyou@cuhk.edu.cn>, Lichao Sun <lis221@lehigh.edu>.

Published at Data in Generative Models Workshop: The Bad, the Ugly, and the Greats (DIG-BUGS) at ICML 2025, Vancouver, Canada. Copyright 2025 by the author(s).

ing realistic contamination scenarios. Experimental results demonstrate the effectiveness of MM-DETECT in identifying varying contamination degrees. Our evaluations on twelve widely-used MLLMs across five prevalent VQA datasets reveal significant contamination among both proprietary and open-source models. Critically, contamination is not only prevalent in multimodal training data but also can originates from unimodal pre-training phases, impacting older benchmarks disproportionately.

In summary, this work provides the first systematic analytical examination of multimodal data contamination, making the following explicit analytical contributions: 1) We analytically characterize multimodal contamination into clearly defined unimodal and cross-modal categories, introducing MM-DETECT as an essential analytical tool. 2) We systematically quantify how benchmark leakage inflates performance metrics, providing clear insights into dataset and model susceptibility to contamination. 3) We present novel analytical insights indicating that contamination not solely emerges during the multimodal training stage but could also from unimodal pre-training stage, critically refining current understandings of contamination dynamics.

2. Preliminaries

We formally define the multimodal data contamination and outline the unique challenges associated with its detection.

2.1. Definition of Multimodal Data Contamination

In contrast to single-modal contamination, multimodal contamination may arise from both unimodal and multimodal data sources. The training data for MLLMs generally consists of pure text pre-training data D_{pretrain} and multimodal alignment or instruction-following data D_{vision} . Consider an instance (x, i, y) from a benchmark dataset D, where xrepresents the text input, i is the image input, and y is the label. Data contamination in MLLMs can be categorized into the following two cases:

- Unimodal Contamination: The pair (x, y) or the input x appears in D_{pretrain} .
- Cross-modal Contamination: The triplet (x, i, y) or the pair (x, i) appears in D_{vision} .

In both cases, models trained on these data may gain an unfair advantage.

2.2. Challenges in Multimodal Detection

The challenges of multimodal contamination detection mainly arise from two aspects.

Challenge I: Inefficiency of Unimodal Methods. Despite the prevalence of unimodal detection methods, their application in multimodal scenarios often encounters difficulties. For example, retrieval-based methods (Brown et al., 2020; Touvron et al., 2023a) attempt to detect contamination by retrieving large-scale corpora used for model training. Yet, they struggle when retrieving multimodal information. Similarly, logits-based methods (Shi et al., 2024; Yeom et al., 2018) rely on observing the distribution of low-probability tokens in model outputs, but the disparity in token probability distributions is less pronounced in instruction-tuned MLLMs. Masking-based methods (Deng et al., 2024), which assess training contamination by evaluating a model's ability to predict specific missing or masked text, face challenges when images in multimodal samples provide clues, leading to overestimated contamination detection. Finally, comparison-based methods (Dong et al., 2024) that measure contamination by comparing model outputs with benchmark data prove to be ineffective for image caption tasks due to low output similarity. To validate these inefficiencies, we have conducted comprehensive experiments with compelling results, which are detailed in Appendix **B**.

Challenge II: Multi-stage Training in MLLMs. Another challenge in detecting contamination in MLLMs is the multistage nature of their training (Yin et al., 2023). Each stage may be subject to data contamination. 1) Initially, the pretraining corpus could contain the textual components of questions from benchmark samples. Moreover, in certain native multimodal model training (Reid et al., 2024), samples may be entirely exposed. 2) Subsequently, during multimodal fine-tuning, the model may utilize training samples of some benchmarks, leading to skewed performance improvements. 3) Furthermore, some models employ extensive mixed image-text data from the internet for modality alignment training (Lin et al., 2023; Bai et al., 2023b), potentially introducing additional contamination. Given the challenges, the development of an effective detection framework for multimodal contamination becomes an urgent need.

Based on the discussion above, we have designed a detection method specifically tailored for multimodal contamination, with a particular focus on VQA tasks. Additionally, we have developed a heuristic method to trace the introduction of contamination across different training phases.

3. Detection Framework: MM-DETECT

We introduce the multimodal contamination detection framework, **MM-DETECT**, designed explicitly to support our systematic analysis of contamination phenomena. The core philosophy of MM-DETECT is to detect the unusual discrepancies in model performance before and after semanticirrelevant perturbations. The framework operates in two primary steps:

- The first step is to generate perturbed datasets using two methods: *Option Order Sensitivity Test* (§C.1) and *Slot Guessing for Perturbed Captions* (§C.2), tailored for multiple-choice and image captioning tasks, respectively.
- The second step involves the application of predefined metrics to detect contamination (§3.3), conducting thorough analyses at both the dataset and instance levels.

3.1. Option Order Sensitivity Test

This method is based on a reasonable and intuitive premise that if the model's performance is highly sensitive to the order of the options, it indicates potential contamination, leading the model to memorize a certain canonical order of the options.

3.2. Slot Guessing for Perturbed Caption

This method is based on the intuition that if a model can predict a missing and important part of a sentence but fails with the back-translated version (from English to Chinese, then back to English), it likely indicates that the model has encountered the original sentence during training.

As shown in Figure 5, the keywords identified are "woods" and "bike". Since the image contains "woods", a correct guess by the model may stem from its multimodal capabilities rather than data contamination. However, if the model fails to predict "bike", which is also present in the image, this may indicate potential leakage of this instance.

For details of the two methods, see Appendix C and D.

3.3. Detection Metrics

Detection Metrics serve as the core analytical instruments within MM-DETECT. Having introduced two detection methods, we now delineate the atomic metrics for the detection pipeline, which consists of two primary steps.

Step 1: Correct Rate Calculation. This step assesses the model's performance on benchmark D before and after perturbation. We denote the correct rate (CR) and perturbed correct rate (PCR) uniformly for both Option Order Sensitivity Test (using Accuracy) and Slot Guessing (using Exact Match). Here, N and N' are the counts of correct answers before and after perturbation, respectively. They are calculated as: N N'

$$CR = \frac{N}{|D|}, \quad PCR = \frac{N'}{|D|}$$

Step 2: Contamination Degree Analysis. This step quantifies the model's contamination degree based on the per-

formance variation pre- and post-perturbation. Specifically, we introduce two metrics to evaluate contamination at both dataset and instance levels.

Dataset Level Metric. We evaluate the reduction in atomic metrics, denoted as Δ :

$$\Delta = PCR - CR$$

This reduction indicates the model's familiarity or memory of the original benchmark relative to the perturbed set, thereby offering insights into potential contamination at the **dataset level**. A significant negative Δ suggests potential extensive leakage in the benchmark dataset, leading to highly perturbation-sensitive model performance.

Instance Level Metric. Despite a non-significant or positive Δ , contamination may still occur at the instance level, as some instances may still have been unintentionally included during training. To identify such instances, we compute X, the count of cases where the model provided correct answers before perturbation but incorrect answers after. The **instance leakage metric** Φ is then obtained by dividing X by the dataset size:

$$\Phi = \frac{X}{|D|},$$

where a larger Φ indicates a higher likelihood of instance leakage.

Compared to methods relying solely on accuracy or perplexity, MM-DETECT explicitly highlights performance drop after perturbations, preventing exaggeration or underestimation of contamination. Moreover, it offers advantages of lower computational overhead, higher sensitivity, and effective black-box applicability, thus serving as an essential analytical toolkit in our study.

4. Evaluating MM-DETECT with Intentional Contamination

This section tackles our first overarching research question — How can we effectively quantify and detect multimodal data contamination? To operationalise this goal, we break RQ1 into three sub-questions:

SQ1 (Effectiveness) Is MM-DETECT able to detect contamination regardless of where it is injected?

SQ2 (Sensitivity) How finely can MM-DETECT measure different leakage levels?

SQ3 (Bias Diagnostic) When training-set data leak, can MM-DETECT reveal the evaluation bias?

We answer these sub-questions by adopting the LLaVA framework and training a suite of 7B-parameter models with intentionally contaminated data during the visual-instruction tuning phase. The contamination protocol and data split follow §5.1.

4.1. MM-DETECT is An Effective Detector

We reproduced the LLaVA-1.5-7B experiment to obtain a baseline model without contamination. Recognizing that contamination can occur anywhere in the training data, we inserted contaminated samples into the visual instruction tuning dataset (D_{tuning}) at three positions, early, mid, and late, creating two groups of contaminated training sets using 1340 ScienceQA test samples or 1000 NoCaps validation samples. Corresponding models, termed Early Cont., Mid Cont., and Late Cont., were then trained for comparison with the baseline.

Models	Scien	ScienceQA Test Set			NoCaps Val. Set			
Widdels	CR	PCR	Δ	CR	PCR	Δ		
Baseline	61.4	61.5	0.01	33.0	32.1	-0.9		
Early Cont.	71.5	68.1	-3.4	37.5	32.0	-5.5		
Mid Cont.	69.4	67.3	-2.1	38.5	35.1	-3.4		
Late Cont.	70.2	66.9	-3.3	38.7	32.6	-6.1		

Table 1. Detection results on contamination using the ScienceQA test set and NoCaps validation set.

Table 1 shows that incorporating contaminated data during training increases both the model's performance and its sensitivity to perturbations. Compared with the baseline, ScienceQA-contaminated models exhibit average increases in CR and PCR of 9.0% and 5.9%, while NoCapscontaminated models show increases of 5.2% and 1.1%. Moreover, all contaminated models demonstrate a marked decrease in Δ , confirming that MM-DETECT effectively identifies data contamination.

4.2. MM-DETECT is Sensitive and Fine-grained

We evaluated MM-DETECT's sensitivity by varying leakage levels in the training set. Using the fully contaminated model as our baseline, we trained additional models with moderate and minimal contamination, by inserting reduced amounts (10% and 50%) of contaminated data at the late position of the training set, to assess leakage impact.



Figure 1. Metrics evaluated under varying leakage levels on the ScienceQA test set and NoCaps validation set.

As illustrated in Figure 1, increasing contamination from 10% to 50% to 100% results in corresponding increases in CR and PCR, alongside progressively larger Δ values. The findings confirm that our framework can accurately

differentiate between varying leakage levels in datasets.

4.3. MM-DETECT Diagnoses Evaluation Bias from Training-set Leakage

We investigated whether MM-DETECT can detect training set leakage by comparing models trained with and without benchmark data contamination. For the ScienceQA experiment, we appended 2000 ScienceQA training samples to the training dataset, creating a contaminated model. For the COCO experiment, we removed the COCO-Caption2017 training data from the original training dataset, resulting in a model without leakage.

Model	Dataset	CR	PCR	Δ
Clean	ScienceQA	61.4	61.5	0.01
Leaked	ScienceQA	64.3	63.8	- 0.5
Clean	COCO-Caption2017	32.5	31.9	-0.6
Leaked	COCO-Caption2017	38.1	34.9	-3.2

Table 2. Performance of models trained without (Clean) and with (Leaked) training set contamination.

Table 2 compares the models' performance. On the ScienceQA test set, the contaminated model outperforms the clean model by 2.9% in CR and 2.3% in PCR, with a Δ of -0.5. On the COCO-Caption2017 validation set, the model trained with COCO data shows a Δ of -3.2. The results indicate that training set leakage inflates performance and that MM-DETECT effectively detects it.

5. Assessing the Extent of Contamination in MLLMs

In this section, we systematically quantify the extent of contamination across various MLLMs and benchmarks, addressing our second research question — What is the degree of contamination?

5.1. Setup

Models. We conducted extensive evaluations on nine open-source MLLMs, including LLaVA-1.5-7B (Liu et al., 2023a), VILA1.5-3B (Lin et al., 2023), Qwen-VL-Chat (Bai et al., 2023b), fuyu-8b⁵, idefics2-8b (Laurençon et al., 2024), Phi-3-vision-128k-instruct (Abdin et al., 2024), Yi-VL-6B (AI et al., 2024), InternVL2-8B (Chen et al., 2023; 2024b), DeepSeek-VL2-Tiny (Wu et al., 2024), as well as three proprietary MLLMs: GPT-4o-2024-08-06 (OpenAI, 2023), Gemini-1.5-Pro-002 (Reid et al., 2024), and Claude-3.5-Sonnet-2024-06-20⁶.

Benchmark Datasets. Our analysis leverages two multichoice datasets: ScienceQA (Lu et al., 2022) and MM-Star (Chen et al., 2024a), along with three caption datasets: COCO-Caption2017 (Lin et al., 2015), NoCaps (Agrawal

⁵https://www.adept.ai/blog/fuyu-8b

⁶https://www.anthropic.com/news/claude-3-5-sonnet

Both Text and Images Leaked! A Systematic Analysis of Data Contamination in Multimodal LLM

Model	Scie	nceQA '	Training	g Set	Sc	cienceQ	A Test S	Set		IStar Va	lidatior	n Set
Metric	CR	PCR	Δ	Φ	CR	PCR	$ \Delta$	Φ	CR	PCR	$ \Delta$	Φ
			0	pen-sou	rce MLI	LMs						
LLaVA-1.5-7B	59.7	58.6	-1.1	12.7	60.3	61.6	1.3	10.5	38.9	41.7	2.8	11.0
VILA1.5-3B	57.7	58.3	0.6	14.5	60.3	59.8	-0.5	14.8	38.6	37.6	-1.0	13.9
Qwen-VL-Chat	58.4	60.8	2.5	13.3	60.3	60.4	0.1	13.7	40.9	44.2	3.3	13.2
fuyu-8b	36.5	37.5	1.0	13.4	37.4	36.9	-0.5	14.9	28.2	27.0	-1.2	17.7
idefics2-8b	85.1	84.0	-1.2	3.7	84.0	84.3	0.3	2.8	48.2	49.3	1.1	7.9
Phi-3-vision-128k-instruct	90.5	90.4	-0.1	4.6	88.4	89.1	0.7	3.9	48.7	51.9	3.2	7.2
Yi-VL-6B	60.5	61.8	1.3	10.0	59.5	61.3	1.8	9.6	38.8	44.0	5.2	9.3
InternVL2-8B	94.1	93.9	-0.3	2.0	92.3	93.1	0.8	1.7	56.9	60.1	3.2	5.1
DeepSeek-VL2-Tiny	86.4	86.5	0.1	5.3	87.1	86.9	-0.2	5.3	51.1	52.1	1.0	10.7
Proprietary MLLMs												
GPT-40	69.9	70.0	0.1	2.7	69.1	69.7	0.6	2.8	48.6	50.5	1.9	9.4
Gemini-1.5-Pro	68.5	67.9	-0.6	6.6	66.5	66.2	-0.3	7.1	45.7	45.5	-0.2	9.9
Claude-3.5-Sonnet	70.3	65.0	-5.3	15.3	67.3	64.9	-2.4	12.4	36.3	36.4	0.1	15.9

Table 3. Comparison of MLLMs on multi-choice datasets. Bold values represent the most significant Δ or Φ ; color codes denote contamination degree: green for minor leakage, yellow for partial leakage, and red for severe leakage.²

Model	CO	OCO Val	idation	Set	No	Caps Va	lidation	Set	Vi	ntage Tr	aining	Set
Metric	CR	PCR	Δ	Φ	CR	PCR	Δ	Φ	CR	PCR	Δ	Φ
	Open-source MLLMs											
LLaVA-1.5-7B	34.6	34.0	-0.6	19.0	30.9	28.5	-2.4	17.9	10.8	10.1	-0.7	9.0
VILA1.5-3B	19.1	20.5	1.4	13.0	19.1	20.5	1.4	13.0	1.5	2.2	0.7	1.5
Qwen-VL-Chat	32.2	30.3	-1.9	19.2	28.7	27.3	-1.4	16.7	15.1	15.4	0.3	12.4
fuyu-8b	9.6	10.6	1.0	7.8	10.0	9.8	-0.2	8.3	2.4	3.3	0.9	2.3
idefics2-8b	43.5	42.3	-1.2	21.2	42.6	37.5	-5.1	23.3	18.5	17.0	-1.5	14.5
Phi-3-vision-128k-instruct	38.8	39.3	0.5	19.4	36.9	33.3	-3.6	19.7	17.4	11.7	-5.7	14.3
Yi-VL-6B	43.9	43.3	-0.6	19.4	37.2	36.1	-1.1	17.5	3.3	4.2	0.9	2.8
InternVL2-8B	53.3	51.9	-1.4	20.4	48.0	46.2	-1.8	20.9	28.0	28.7	0.7	18.8
DeepSeek-VL2-Tiny	23.8	21.4	-2.4	13.5	19.3	18.1	-1.2	12.2	7.5	6.9	-0.6	6.3
Proprietary MLLMs												
GPT-40	58.1	54.4	-3.7	23.1	54.2	55.1	0.9	19.4	36.3	38.4	2.1	20.1
Gemini-1.5-Pro	57.5	55.3	-2.2	21.6	51.2	52.0	0.8	18.7	46.3	41.0	-5.3	28.3
Claude-3.5-Sonnet	53.7	51.0	-2.7	21.8	50.8	51.5	0.7	20.0	35.2	33.0	-2.2	21.3

Table 4. Comparison of MLLMs on caption datasets. Bold values represent the most significant Δ or Φ ; color codes denote contamination degree: green for minor leakage, yellow for partial leakage, and red for severe leakage.⁴

et al., 2019), and Vintage⁷. MMStar and Vintage, owing to their recent inception, serve to contrast leakage levels with other datasets. We randomly selected 2000 and 1340 samples from ScienceQA's training and test sets, respectively, with 1000 samples from the other datasets. Given the unavailability of public test labels for COCO-Caption2017 and NoCaps, we used their validation sets.

5.2. Main Results

Multi-choice Datasets. Table 3 yields several conclusions: (1) Both open-source and proprietary models exhibit contamination. For example, on the ScienceQA training set, both open-source models like LLaVA-1.5-7B and idefics2-8b and proprietary model Gemini-1.5-Pro show minor contamination degree. (2) Proprietary models are more contaminated. Claude-3.5-Sonnet, for instance, registers a severe Δ with higher Φ values on both ScienceQA training and test sets, indicating extensive leakage. (3) **Training set leakage is more pronounced than test set leakage.** On the ScienceQA dataset, models generally exhibit larger Δ values in the training set, for instance, Claude-3.5-Sonnet shows $\Delta = -5.3$ on training versus $\Delta = -2.4$ on the test set, while most models have near-zero Δ on the test set. (4) **Older benchmarks are more prone to leak.** The older ScienceQA test set shows more severe leakage compared to the newer MMStar validation set.

Caption Datasets. Table 4 yields several conclusions: (1) **Both open-source and proprietary models exhibit contamination on caption datasets.** For example, in the COCO Validation Set, open-source models such as DeepSeek-VL2-Tiny and proprietary models like GPT-40

⁷https://huggingface.co/datasets/

SilentAntagonist/vintage-artworks-60k-captioned

⁷Based on §4.1, the degrees on multi-choice datasets are defined as: $\Delta \in (-1.6, -0.2]$ for minor leakage, $\Delta \in (-2.9, -1.6]$ for partial leakage, and $\Delta \leq -2.9$ for severe leakage.

⁷Based on §4.1, the degrees on caption datasets are defined as: $\Delta \in (-2.4, -1.1]$ for minor leakage, $\Delta \in (-5.0, -2.4]$ for partial leakage, and $\Delta \leq -5.0$ for severe leakage.

record a significant contamination degree. (2) **Leakage levels vary significantly by benchmark.** For example, on the NoCaps Validation Set, open-source models exhibit more pronounced contamination degree than proprietary models, whereas the trend reverses on the COCO Validation Set. These findings confirm that caption datasets are vulnerable to leakage, with proprietary models generally exhibiting more pronounced contamination effects.

6. Identifying the Origin of Contamination in MLLMs

In this section, we address our third research question — **When is contamination predominantly introduced?** Although the training data for some MLLMs is openly documented, an important question remains: if contamination does not arise from the visual training, could it stem from the unimodal (pre-training) phase, as defined in §2.1? To address this possibility, we examined the underlying LLMs of the evaluated MLLMs and conducted a series of experiments (§6.1). We also explored the origins of cross-modal contamination arising during visual instruction tuning (§6.2).

6.1. Heuristic Detection of Unimodal Pre-training Contamination

Model	Accuracy	Φ_M
LLaMA2-7b (LLaVA-1.5 & VILA)	25.6	11.0
Qwen-7B (Qwen-VL)	13.2	13.2
Internlm2-7B (InternVL2)	11.0	5.1
Mistral-7B-v0.1 (idefics2)	10.7	7.9
Phi-3-small-128k-instruct (Phi-3-vision)	6.1	7.2
Yi-6B (Yi-VL)	3.4	9.3

Table 5. Contamination rates of LLMs used by MLLMs. Φ_M denotes the Φ of the respective MLLMs.

We employ a heuristic approach based on the intuition that if LLMs can correctly answer an **image-required** question **without the image** when **random guessing is effectively inhibited**, it may indicate the leakage of that instance.

Experiment Setup. We used MMStar as the benchmark, where every question relies on visual input for correct answers. The tested models include LLaMA2-7B (Touvron et al., 2023b) (used by LLaVA-1.5 and VILA), Qwen-7B (Bai et al., 2023a) (used by Qwen-VL), Mistral-7B-v0.1 (Jiang et al., 2023) (used by idefics2), Phi-3-small-128kinstruct (Abdin et al., 2024) (used by Phi-3-vision), Yi-6B (AI et al., 2024) (used by Yi-VL), and InternIm2-7B (Cai et al., 2024) (used by InternVL2). To inhibit random guessing, we appended the prompt "If you do not know the answer, output I don't know" to the instructions. A sanity check in Appendix G.2 confirms that this uncertainty clause effectively suppresses lucky guesses, validating its inclusion in our main protocol. Accuracy - the frequency with which models correctly answer questions without image input - is reported as the primary metric. Note that we did not evaluate

Model	ScienceQA	COCO Caption	Nocaps
Phi-3-Vision	0.7	0.5	-3.6
VILA	-0.5	1.4	1.4
Idefics2	0.3	-1.2	-5.1
LLaVA-1.5	1.3	-0.6	-2.4
Yi-VL	1.8	-0.6	-1.1
DeepSeek-VL2	-0.2	-2.4	-1.2
Qwen-VL-Chat	0.1	-1.9	-1.4
InternVL2	0.8	-1.4	-1.8

Table 6. Overlap between the training data of MLLMs and the benchmarks, as well as the contamination degree Δ of MLLMs on benchmarks. Green signifies no overlap, yellow suggests potential overlap, and Red indicates partial or entire overlap.

Fuyu-8B and proprietary models since their unimodal LLM components and training data remain undisclosed.

Main Results. Table 5 yields several conclusions: (1) **Contamination occurs in LLM.** All models exhibit varied contamination rates, indicating that their pre-training data likely included text from multimodal benchmarks. (2) **Elevated LLM contamination correlates with increased MLLM leakage.** For instance, VILA1.5-3B and Qwen-VL-Chat exhibit significant Φ values that mirror their underlying LLM contamination levels. These findings suggest that contamination in these MLLMs may originate partly from the LLMs' pre-training phase, rather than solely from multimodal training.

6.2. Analyzing Cross-modal Contamination in Multimodal Fine-tuning

To investigate the origins of cross-modal contamination, we scrutinize the visual instruction tuning data of MLLMs. We delve into the construction process of three benchmark datasets: ScienceQA, COCO Caption, and Nocaps, comparing them with the training data and its sources of various open-source MLLMs to analyze the degree of overlap.

As Table 6 illustrates, MLLMs marked in red and yellow typically exhibit a significant contamination degree. Yet, even MLLMs labeled in green aren't exempt from the risk of cross-modal contamination. This is because some models have been trained on large-scale interleaved imagetext datasets (e.g., OBELICS (Laurençon et al., 2023)), datasets derived from online sources (e.g., Conceptual Caption (Sharma et al., 2018)), or in-house data. Furthermore, some models haven't fully disclosed their training data, which may lead to overlooked potential leaks. A more detailed analysis is shown in the appendix F.

7. Conclusion and Future Work

We have presented MM-DETECT, a lightweight, black-box framework that systematically detects and quantifies multimodal data contamination in MLLMs. Our experiments reveal that leakage—originating both in large-scale text pretraining and in multimodal fine-tuning—can significantly inflate performance and undermine fair evaluation.

References

- Abdin, M. I., Jacobs, S. A., Awan, A. A., Aneja, J., Awadallah, A., Awadalla, H., Bach, N., Bahree, A., Bakhtiari, A., Behl, H. S., Benhaim, A., Bilenko, M., Bjorck, J., Bubeck, S., Cai, M., Mendes, C. C. T., Chen, W., Chaudhary, V., Chopra, P., Giorno, A. D., de Rosa, G., Dixon, M., Eldan, R., Iter, D., Garg, A., Goswami, A., Gunasekar, S., Haider, E., Hao, J., Hewett, R. J., Huynh, J., Javaheripi, M., Jin, X., Kauffmann, P., Karampatziakis, N., Kim, D., Khademi, M., Kurilenko, L., Lee, J. R., Lee, Y. T., Li, Y., Liang, C., Liu, W., Lin, E., Lin, Z., Madan, P., Mitra, A., Modi, H., Nguyen, A., Norick, B., Patra, B., Perez-Becker, D., Portet, T., Pryzant, R., Qin, H., Radmilac, M., Rosset, C., Roy, S., Ruwase, O., Saarikivi, O., Saied, A., Salim, A., Santacroce, M., Shah, S., Shang, N., Sharma, H., Song, X., Tanaka, M., Wang, X., Ward, R., Wang, G., Witte, P., Wyatt, M., Xu, C., Xu, J., Yadav, S., Yang, F., Yang, Z., Yu, D., Zhang, C., Zhang, C., Zhang, J., Zhang, L. L., Zhang, Y., Zhang, Y., Zhang, Y., and Zhou, X. Phi-3 technical report: A highly capable language model locally on your phone. CoRR, abs/2404.14219, 2024. doi: 10.48550/ARXIV.2404.14219. URL https: //doi.org/10.48550/arXiv.2404.14219.
- Agrawal, H., Desai, K., Wang, Y., Chen, X., Jain, R., Johnson, M., Batra, D., Parikh, D., Lee, S., and Anderson, P. nocaps: novel object captioning at scale. In 2019 IEEE/CVF International Conference on Computer Vision (ICCV). IEEE, October 2019. doi: 10.1109/iccv. 2019.00904. URL http://dx.doi.org/10.1109/ICCV. 2019.00904.
- AI, ., :, Young, A., Chen, B., Li, C., Huang, C., Zhang, G.,
 Zhang, G., Li, H., Zhu, J., Chen, J., Chang, J., Yu, K.,
 Liu, P., Liu, Q., Yue, S., Yang, S., Yang, S., Yu, T., Xie,
 W., Huang, W., Hu, X., Ren, X., Niu, X., Nie, P., Xu, Y.,
 Liu, Y., Wang, Y., Cai, Y., Gu, Z., Liu, Z., and Dai, Z. Yi:
 Open foundation models by 01.ai, 2024.
- Bai, J., Bai, S., Chu, Y., Cui, Z., Dang, K., Deng, X., Fan, Y., Ge, W., Han, Y., Huang, F., Hui, B., Ji, L., Li, M., Lin, J., Lin, R., Liu, D., Liu, G., Lu, C., Lu, K., Ma, J., Men, R., Ren, X., Ren, X., Tan, C., Tan, S., Tu, J., Wang, P., Wang, S., Wang, W., Wu, S., Xu, B., Xu, J., Yang, A., Yang, H., Yang, J., Yang, S., Yao, Y., Yu, B., Yuan, H., Yuan, Z., Zhang, J., Zhang, X., Zhang, Y., Zhang, Z., Zhou, C., Zhou, J., Zhou, X., and Zhu, T. Qwen technical report. *CoRR*, abs/2309.16609, 2023a. doi: 10.48550/ARXIV.2309.16609. URL https://doi.org/10.48550/arXiv.2309.16609.
- Bai, J., Bai, S., Yang, S., Wang, S., Tan, S., Wang, P., Lin, J., Zhou, C., and Zhou, J. Qwen-vl: A versatile vision-language model for understanding, localization,

text reading, and beyond, 2023b. URL https://arxiv. org/abs/2308.12966.

- Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei, D. Language models are few-shot learners. *CoRR*, abs/2005.14165, 2020. URL https://arxiv.org/abs/2005.14165.
- Cai, Z., Cao, M., Chen, H., Chen, K., Chen, K., Chen, X., Chen, X., Chen, Z., Chen, Z., Chu, P., Dong, X., Duan, H., Fan, Q., Fei, Z., Gao, Y., Ge, J., Gu, C., Gu, Y., Gui, T., Guo, A., Guo, Q., He, C., Hu, Y., Huang, T., Jiang, T., Jiao, P., Jin, Z., Lei, Z., Li, J., Li, J., Li, L., Li, S., Li, W., Li, Y., Liu, H., Liu, J., Hong, J., Liu, K., Liu, K., Liu, X., Lv, C., Lv, H., Lv, K., Ma, L., Ma, R., Ma, Z., Ning, W., Ouyang, L., Qiu, J., Qu, Y., Shang, F., Shao, Y., Song, D., Song, Z., Sui, Z., Sun, P., Sun, Y., Tang, H., Wang, B., Wang, G., Wang, J., Wang, J., Wang, R., Wang, Y., Wang, Z., Wei, X., Weng, Q., Wu, F., Xiong, Y., and et al. Internlm2 technical report. *CoRR*, abs/2403.17297, 2024. doi: 10.48550/ARXIV.2403.17297. URL https: //doi.org/10.48550/arXiv.2403.17297.
- Chen, L., Li, J., Dong, X., Zhang, P., Zang, Y., Chen, Z., Duan, H., Wang, J., Qiao, Y., Lin, D., et al. Are we on the right way for evaluating large vision-language models? *arXiv preprint arXiv:2403.20330*, 2024a.
- Chen, Z., Wu, J., Wang, W., Su, W., Chen, G., Xing, S., Zhong, M., Zhang, Q., Zhu, X., Lu, L., Li, B., Luo, P., Lu, T., Qiao, Y., and Dai, J. Internvl: Scaling up vision foundation models and aligning for generic visuallinguistic tasks. arXiv preprint arXiv:2312.14238, 2023.
- Chen, Z., Wang, W., Tian, H., Ye, S., Gao, Z., Cui, E., Tong, W., Hu, K., Luo, J., Ma, Z., et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*, 2024b.
- Deng, C., Zhao, Y., Tang, X., Gerstein, M., and Cohan, A. Investigating data contamination in modern benchmarks for large language models. In Duh, K., Gómez-Adorno, H., and Bethard, S. (eds.), *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), NAACL 2024, Mexico City, Mexico, June 16-21, 2024*, pp. 8706–8719. Association for Computational Linguistics, 2024. doi: 10.18653/V1/2024.NAACL-LONG.482. URL https: //doi.org/10.18653/v1/2024.naacl-long.482.

- Dong, Y., Jiang, X., Liu, H., Jin, Z., Gu, B., Yang, M., and Li, G. Generalization or memorization: Data contamination and trustworthy evaluation for large language models. In Ku, L., Martins, A., and Srikumar, V. (eds.), *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, pp. 12039–12050. Association for Computational Linguistics, 2024. URL https: //aclanthology.org/2024.findings-acl.716.
- Jiang, A. Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D. S., de las Casas, D., Bressand, F., Lengyel, G., Lample, G., Saulnier, L., Lavaud, L. R., Lachaux, M.-A., Stock, P., Scao, T. L., Lavril, T., Wang, T., Lacroix, T., and Sayed, W. E. Mistral 7b, 2023. URL https: //arxiv.org/abs/2310.06825.
- Laurençon, H., Tronchon, L., Cord, M., and Sanh, V. What matters when building vision-language models? *CoRR*, abs/2405.02246, 2024. doi: 10.48550/ARXIV. 2405.02246. URL https://doi.org/10.48550/arXiv. 2405.02246.
- Laurençon, H., Saulnier, L., Tronchon, L., Bekman, S., Singh, A., Lozhkov, A., Wang, T., Karamcheti, S., Rush, A. M., Kiela, D., Cord, M., and Sanh, V. Obelics: An open web-scale filtered dataset of interleaved image-text documents, 2023. URL https://arxiv.org/abs/2306. 16527.
- Li, J., Li, D., Savarese, S., and Hoi, S. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference* on machine learning, pp. 19730–19742. PMLR, 2023.
- Lin, J., Yin, H., Ping, W., Lu, Y., Molchanov, P., Tao, A., Mao, H., Kautz, J., Shoeybi, M., and Han, S. VILA: on pre-training for visual language models. *CoRR*, abs/2312.07533, 2023. doi: 10.48550/ARXIV. 2312.07533. URL https://doi.org/10.48550/arXiv. 2312.07533.
- Lin, T.-Y., Maire, M., Belongie, S., Bourdev, L., Girshick, R., Hays, J., Perona, P., Ramanan, D., Zitnick, C. L., and Dollár, P. Microsoft coco: Common objects in context, 2015.
- Liu, H., Li, C., Li, Y., and Lee, Y. J. Improved baselines with visual instruction tuning. *CoRR*, abs/2310.03744, 2023a. doi: 10.48550/ARXIV.2310.03744. URL https: //doi.org/10.48550/arXiv.2310.03744.
- Liu, Y., Duan, H., Zhang, Y., Li, B., Zhang, S., Zhao, W., Yuan, Y., Wang, J., He, C., Liu, Z., Chen, K., and Lin, D. Mmbench: Is your multi-modal model an allaround player? *CoRR*, abs/2307.06281, 2023b. doi: 10.48550/ARXIV.2307.06281. URL https://doi.org/ 10.48550/arXiv.2307.06281.

- Lu, P., Mishra, S., Xia, T., Qiu, L., Chang, K.-W., Zhu, S.-C., Tafjord, O., Clark, P., and Kalyan, A. Learn to explain: Multimodal reasoning via thought chains for science question answering. In *The 36th Conference on Neural Information Processing Systems (NeurIPS)*, 2022.
- OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023. doi: 10.48550/ARXIV.2303.08774. URL https://doi.org/10.48550/arXiv.2303.08774.
- Papineni, K., Roukos, S., Ward, T., and Zhu, W.-J. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pp. 311–318, 2002.
- Reid, M., Savinov, N., Teplyashin, D., Lepikhin, D., Lillicrap, T., Alayrac, J.-b., Soricut, R., Lazaridou, A., Firat, O., Schrittwieser, J., et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024.
- Reimers, N. and Gurevych, I. Sentence-bert: Sentence embeddings using siamese bert-networks. arXiv preprint arXiv:1908.10084, 2019.
- Sharma, P., Ding, N., Goodman, S., and Soricut, R. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings* of ACL, 2018.
- Shi, W., Ajith, A., Xia, M., Huang, Y., Liu, D., Blevins, T., Chen, D., and Zettlemoyer, L. Detecting pretraining data from large language models. In *The Twelfth International Conference on Learning Representations*, *ICLR 2024, Vienna, Austria, May 7-11, 2024*. Open-Review.net, 2024. URL https://openreview.net/ forum?id=zWqr3MQuNs.
- Song, D., Chen, S., Chen, G. H., Yu, F., Wan, X., and Wang,
 B. Milebench: Benchmarking mllms in long context. arXiv preprint arXiv:2404.18532, 2024.
- Toutanvoa, K. and Manning, C. D. Enriching the knowledge sources used in a maximum entropy part-of-speech tagger. In Schütze, H. and Su, K. (eds.), Joint SIGDAT Conference on Empirical Methods in Natural Language Processing and Very Large Corpora, EMNLP 2000, Hong Kong, October 7-8, 2000, pp. 63–70. Association for Computational Linguistics, 2000. doi: 10.3115/1117794.1117802. URL https://aclanthology.org/W00-1308/.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., Bikel, D., Blecher, L., Canton-Ferrer, C., Chen, M., Cucurull, G., Esiobu, D., Fernandes, J., Fu, J., Fu, W., Fuller, B., Gao, C., Goswami, V., Goyal, N.,

Hartshorn, A., Hosseini, S., Hou, R., Inan, H., Kardas, M., Kerkez, V., Khabsa, M., Kloumann, I., Korenev, A., Koura, P. S., Lachaux, M., Lavril, T., Lee, J., Liskovich, D., Lu, Y., Mao, Y., Martinet, X., Mihaylov, T., Mishra, P., Molybog, I., Nie, Y., Poulton, A., Reizenstein, J., Rungta, R., Saladi, K., Schelten, A., Silva, R., Smith, E. M., Subramanian, R., Tan, X. E., Tang, B., Taylor, R., Williams, A., Kuan, J. X., Xu, P., Yan, Z., Zarov, I., Zhang, Y., Fan, A., Kambadur, M., Narang, S., Rodriguez, A., Stojnic, R., Edunov, S., and Scialom, T. Llama 2: Open foundation and fine-tuned chat models. *CoRR*, abs/2307.09288, 2023a. doi: 10.48550/ARXIV.2307.09288. URL https://doi.org/10.48550/arXiv.2307.09288.

- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., Bikel, D., Blecher, L., Canton-Ferrer, C., Chen, M., Cucurull, G., Esiobu, D., Fernandes, J., Fu, J., Fu, W., Fuller, B., Gao, C., Goswami, V., Goyal, N., Hartshorn, A., Hosseini, S., Hou, R., Inan, H., Kardas, M., Kerkez, V., Khabsa, M., Kloumann, I., Korenev, A., Koura, P. S., Lachaux, M., Lavril, T., Lee, J., Liskovich, D., Lu, Y., Mao, Y., Martinet, X., Mihaylov, T., Mishra, P., Molybog, I., Nie, Y., Poulton, A., Reizenstein, J., Rungta, R., Saladi, K., Schelten, A., Silva, R., Smith, E. M., Subramanian, R., Tan, X. E., Tang, B., Taylor, R., Williams, A., Kuan, J. X., Xu, P., Yan, Z., Zarov, I., Zhang, Y., Fan, A., Kambadur, M., Narang, S., Rodriguez, A., Stojnic, R., Edunov, S., and Scialom, T. Llama 2: Open foundation and fine-tuned chat models. CoRR, abs/2307.09288, 2023b. doi: 10.48550/ARXIV.2307.09288. URL https: //doi.org/10.48550/arXiv.2307.09288.
- Wu, Z., Chen, X., Pan, Z., Liu, X., Liu, W., Dai, D., Gao, H., Ma, Y., Wu, C., Wang, B., Xie, Z., Wu, Y., Hu, K., Wang, J., Sun, Y., Li, Y., Piao, Y., Guan, K., Liu, A., Xie, X., You, Y., Dong, K., Yu, X., Zhang, H., Zhao, L., Wang, Y., and Ruan, C. Deepseek-vl2: Mixture-of-experts visionlanguage models for advanced multimodal understanding, 2024. URL https://arxiv.org/abs/2412.10302.
- Xu, R., Wang, Z., Fan, R.-Z., and Liu, P. Benchmarking benchmark leakage in large language models. *arXiv preprint arXiv:2404.18824*, 2024.
- Yeom, S., Giacomelli, I., Fredrikson, M., and Jha, S. Privacy risk in machine learning: Analyzing the connection to overfitting. In 31st IEEE Computer Security Foundations Symposium, CSF 2018, Oxford, United Kingdom, July 9-12, 2018, pp. 268–282. IEEE Computer Society, 2018. doi: 10.1109/CSF.2018.00027. URL https://doi.org/ 10.1109/CSF.2018.00027.
- Yin, S., Fu, C., Zhao, S., Li, K., Sun, X., Xu, T., and Chen, E. A survey on multimodal large language models. *arXiv* preprint arXiv:2306.13549, 2023.

A. Illustrative Overview and Framework Visualization



Figure 2. An analytical breakdown illustrating different forms and origins of multimodal data contamination across distinct training stages of MLLMs.



Figure 3. The overview of proposed MM-DETECT framework.

B. Inefficiency of Unimodal Methods

We demonstrate the results of traditional unimodal contamination detection methods applied to MLLMs.

B.1. Logits-base

These methods determine contamination by observing the distribution of low-probability tokens in model outputs. However, MLLMs typically undergo instruction fine-tuning, which enhances their instruction-following capabilities, leading to less significant differences in token probability distributions. As shown in Table 7, LLaVA-1.5-13b exhibits extremely low perplexity on multimodal benchmark datasets.

Dataset	Perplexity	Split
ScienceQA	1.4498	Training Set
MMStar	1.4359	Validation Set
COCO-Caption2017	1.7530	Validation Set
NoCaps	1.8155	Validation Set

Table 7. Perplexity of LLaVA-1.5-13b on various multimodal benchmarks (100 samples randomly selected from each dataset).

B.2. Masking-base

These methods involve masking phrases or sentences and providing data from the benchmark to guide the model in filling in the missing parts. However, multimodal datasets often contain images that include the masked portions of sentences, effectively providing answers to the model. This results in significantly higher success rates for MLLMs in predicting missing parts compared to unimodal language models, leading to exaggerated contamination detection. As shown in Table 8, LLaVA-1.5-13b has a high probability of Exact Match for predicting the masked word.

Dataset	Exact Match	ROUGE-L F1	Split
COCO-Caption2017	0.24	0.36	Validation Set
NoCaps	0.22	0.29	Validation Set

Table 8. Contamination detection of LLaVA-1.5-13b using TS-Guessing (question-based) on various multimodal benchmarks (100 samples randomly selected from each dataset).

B.3. Comparison-base

These methods identify contamination by comparing the similarity between models' outputs and benchmark data. However, MLLMs often undergo data augmentation, causing their outputs to diverge significantly from the labels in benchmark data, making effective contamination detection challenging. From Table 9, we can see that CDD (Contamination Detection via Output Distribution) indicates a contamination metric of 0% across all multimodal benchmark datasets.

Dataset	Contamination Metric	Split
COCO-Caption2017	0.0000%	Validation Set
NoCaps	0.0000%	Validation Set

Table 9. Contamination detection of LLaVA-1.5-13b using CDD (Contamination Detection via Output Distribution) on various multimodal benchmarks (100 samples randomly selected from each dataset).

C. Details of the Methods



Figure 4. An example of Option Order Sensitivity Test applied to a contaminated model.

C.1. Option Order Sensitivity Test

This method is based on a reasonable and intuitive premise that if the model's performance is highly sensitive to the order of the options, as shown in Figure 4, it indicates potential contamination, leading the model to memorize a certain canonical order of the options.



Figure 5. A simple example shows the procedure of **caption pertubation**.

Method Formulation. Let *D* be a dataset consisting of *n* datapoints. Each datapoint d_i ($i \in \{1, ..., n\}$) comprises a question Q_i , an associated image I_i , and a set of answer choices $A_i = \{a_i^1, a_i^2, ..., a_i^m\}$, where *m* is the number of choices and the correct answer is denoted by a_i^c .

To introduce positional variation, the set A_i is randomly shuffled to obtain a new set A'_i , ensuring that the index of the correct answer a_i^c in A'_i differs from its original position in A_i . The final prompts, before and after shuffling, are constructed by concatenating the image, question and choices:

$$P = \text{Concat}(I_i, Q_i, A_i),$$
$$P' = \text{Concat}(I_i, Q_i, A'_i),$$

where P and P' are the inputs to the model, and Q_i and I_i remain unchanged throughout this process.

C.2. Slot Guessing for Perturbed Caption

This method is based on the intuition that if a model can predict a missing and important part of a sentence but fails with the back-translated version (from English to Chinese, then back to English), it likely indicates that the model has encountered the original sentence during training.

As shown in Figure 5, the keywords identified are "woods" and "bike". Since the image contains "woods", a correct guess by the model may stem from its multimodal capabilities rather than data contamination. However, if the model fails to predict "bike", which is also present in the image, this may indicate potential leakage of this instance.

Method Formulation. Let D be a dataset containing n datapoints. Each datapoint d_i $(i \in \{1, ..., n\})$ consists of an image-caption pair, where the caption S_i describes the visual features of the corresponding image I_i . We first apply a back-translation function⁸, where we use the Google Translate API for Python to implement back-translation, to S_i :

$$S'_i = f_{\text{back-translate}}(S_i).$$

resulting in a paraphrased version S'_i . Next, we perform keyword extraction⁹ on both S_i and S'_i :

$$K_i = f_{\text{keyword}}(S_i), \quad K'_i = f_{\text{keyword}}(S'_i),$$

⁸A quantitative analysis of the semantic and lexical similarity between the original and back-translated captions is provided in Appendix G.1.

⁹We employ the Stanford POS Tagger (Toutanvoa & Manning, 2000), targeting nouns, adjectives, and verbs, as they encapsulate the core meaning of the sentences.

where K_i and K'_i denote the extracted keywords from S_i and S'_i , respectively. We then apply a masking function f_{mask} to replace the extracted keywords with a placeholder token [MASK]:

$$S_{i,\text{mask}} = f_{\text{mask}}(S_i, K_i), \ S'_{i,\text{mask}} = f_{\text{mask}}(S'_i, K'_i).$$

The final prompt guiding the model to complete the masked-word prediction can be represented as:

$$P_i = \text{Concat}(I_i, Q_i, S_{i,\text{mask}}),$$

$$P'_i = \text{Concat}(I_i, Q_i, S'_{i,\text{mask}}).$$

D. Detailed Slot Guessing Pipeline

D.1. Back-Translation

The back-translation function applies a two-step translation process to generate a paraphrased caption S'_i from the original caption S_i . In this method, we use the Google Translate API to translate the caption into **Chinese** and then back into the original language to generate the paraphrase.

Algorithm 1 Back-Translation

- 1: **Input:** Original caption S_i
- 2: Translate S_i to an intermediate language L
- 3: Translate the resulting caption back from language L to the original language
- 4: **Output:** Paraphrased caption S'_i

D.2. Keyword Extraction

We extract keywords from both the original caption S_i and the paraphrased caption S'_i using the Stanford POS Tagger. Keywords are identified as nouns (NN), adjectives (JJ), and verbs (VB), which are considered to encapsulate the core meaning of the sentence. We apply this process to both captions.

Algorithm 2 Keyword Extraction

- 1: Input: Caption S
- 2: Apply POS tagging to S to obtain tags for each word
- 3: Extract words whose POS tags are in {NN, JJ, VB}
- 4: **Output:** List of extracted keywords K

D.3. Keyword Masking

We apply a masking function to randomly select one keyword from the extracted keywords and replace it with a placeholder token [MASK]. This is done by identifying the position of the selected keyword in the sentence and substituting it with the placeholder.

Algorithm 3 Keyword Masking

- 1: Input: Caption S, Keywords K
- 2: If *K* is empty then return "failed"
- 3: Randomly select a keyword k from K
- 4: Find the first occurrence of k in S
- 5: Replace k with the placeholder [MASK]
- 6: **Output:** Masked caption S_{mask}

E. Contamination Degree Analysis

Based on §4.1, the degrees on multi-choice datasets are defined as: $\Delta \in (-1.6, -0.2]$ for minor leakage, $\Delta \in (-2.9, -1.6]$ for partial leakage, and $\Delta \leq -2.9$ for severe leakage. Based on §4.1, the degrees on caption datasets are defined as:

 $\Delta \in (-2.4, -1.1]$ for minor leakage, $\Delta \in (-5.0, -2.4]$ for partial leakage, and $\Delta \leq -5.0$ for severe leakage. Details are shown in the algorithm 4.

Algorithm 4 Contamination Degree Analysis **Require:** Benchmark dataset D, Model M 1: Define contamination degree C_{Minor} , $C_{Partial}$, C_{Severe} 2: **if** *D* is multiple-choice **then** 3: Generate perturbed set D_{pert} via §C.1 4: else Generate perturbed set D_{pert} via §C.2 5: 6: end if 7: Compute CR, PCR, Δ , Φ using §3.3 8: **if** multiple-choice **then** $\begin{cases} \mathcal{C}_{\text{Minor}}, & \Delta \in (-1.6, -0.2] \\ \mathcal{C}_{\text{Partial}}, & \Delta \in (-2.9, -1.6] \\ \mathcal{C}_{\text{Severe}}, & \Delta \leq -2.9 \end{cases}$ 9: $\mathcal{C} \leftarrow$ 10: else $\begin{cases} \mathcal{C}_{\text{Minor}}, & \Delta \in (-2.4, -1.1] \\ \mathcal{C}_{\text{Partial}}, & \Delta \in (-5.0, -2.4] \end{cases}$ $\mathcal{C} \leftarrow$ 11: $C_{\text{Severe}}, \quad \Delta \leq -5.0$ 12: end if **Ensure:** CR, PCR, Δ , Φ , C

F. Detailed Overlap Analysis

It is **impractical** to quantify overlapping samples: 1) Many models do not release their complete training datasets publicly; instead, they only mention the data sources in their technical reports. 2) Even if we had access to complete training datasets, identifying specific overlapping samples using matching algorithms (such as exact match) remains challenging. This is because the original benchmarks might have undergone data augmentation before being used for model training, and multimodal benchmarks include images, both of which complicate the practical utility of matching algorithms. The feasible approach is **manually reviewing** the technical reports of these models to verify whether their training data overlaps with benchmarks, as shown in the table 10.

MLLMs	Multimodal Alignment/Pretraining Data	Supervised Fine-Tuning Data
Phi-3-Vision	Alignment Data includes FLD-5B.	Not yet released
	Open Images is one source of FLD-5B.	
	Open Images is also a source of Nocaps.	
	Therefore, there is potential overlap in Nocaps.	
VILA	No overlap	Includes RefCOCO, VQAv2, GQA
		MS COCO is a source of RefCOCO, VQAv2.
		GQA's source is Visual Genome Scene Graph, which also includes MS COCO.
		COCO Caption's source is MS COCO, and NoCaps' source includes COCO Caption.
		Therefore, there is potential overlap in COCO Caption and NoCaps.
Idefics2	Alignment Data includes SBU Captions	SFT Data includes SBU Captions: potential overlap in COCO Caption and NoCaps.
	SBU Captions' source includes Flickr	
	COCO Caption's source includes MS COCO, and MS COCO's source includes Flickr	
	NoCaps' source includes COCO Caption	
	Therefore, there is potential overlap in COCO Caption and NoCaps.	
LLaVA-1.5	Alignment Data includes SBU Captions: COCO Caption and NoCaps with potential overlap.	SFT Data includes RefCOCO, VQAv2, GQA: COCO Caption and NoCaps with potential overlap.
Yi-VL	Alignment Data includes Flickr, VQAv2, RefCOCO:	SFT Data includes GQA: COCO Caption and NoCaps with potential overlap.
	COCO Caption and NoCaps with potential overlap.	
DeepSeek-VL2	No overlap	SFT Data includes Flickr, GQA: COCO Caption and NoCaps with potential overlap.
Qwen-VL-Chat	Directly uses COCO Caption in the pretraining stage,	Not yet released
	therefore there is partial or entire overlap in COCO Caption and NoCaps.	
InternVL2	Alignment Data includes COCO Caption: partial or entire overlap in COCO Caption and NoCaps.	SFT Data includes ScienceQA, therefore there is partial or entire overlap in ScienceQA.

Table 10. Comparison of MLLMs and Their Data Sources

G. Other Experiments

G.1. Semantic & Lexical Similarity After Back-Translation

Setup. To quantify how much meaning and wording change during our *caption perturbation* step (§C.2), we applied an **English** \rightarrow **Chinese** \rightarrow **English** back-translation to every caption in three validation splits – COCO-Caption, NoCaps, and our Vintage dataset. For each original (*c*) and back-translated caption (\tilde{c}) we computed

- SBERT cosine similarity (Reimers & Gurevych, 2019) as a sentence-level semantic score, and
- BLEU-4 (Papineni et al., 2002) as a token-overlap lexical score.

Dataset	Avg. SBERT \uparrow	Avg. BLEU \uparrow	Correlation <i>r</i>
COCO Caption	0.894	0.236	0.386
NoCaps	0.887	0.264	0.410
Vintage	0.914	0.441	0.423

We additionally report the Pearson correlation between the two metrics across captions within each dataset.

Table 11. Average semantic (SBERT) and lexical (BLEU-4) similarity between original and back-translated captions, together with their Pearson correlation (r).

Key Observations.

- **High semantic preservation.** All three datasets record SBERT scores close to 0.9, indicating that back-translation keeps the *meaning* of captions largely intact; the VINTAGE split achieves the strongest preservation (0.914).
- **Substantial lexical variation.** BLEU-4 values are comparatively low, showing that wording and surface forms differ considerably—consistent with the presence of synonym substitutions and syntactic reshuffling introduced by back-translation.
- Weak yet positive coupling. Pearson correlations between the two metrics lie in the 0.38-0.42 band, suggesting only a mild positive relationship: captions that keep more tokens also tend to retain semantics, but plenty of cases preserve meaning even with low lexical overlap.

These results justify using back-translation as a *semantics-preserving yet lexically diversifying* perturbation in our contamination-detection pipeline.

G.2. Sanity Check for the "I don't know" Instruction

Setup. To verify that appending the uncertainty clause "*If you do not know the answer, output "I don't know"*." effectively suppresses random guessing, we performed a pilot experiment on 1 000 randomly sampled questions from MMSTAR. All images were removed, so a truly vision-grounded model should either fail or explicitly abstain. We evaluated the unimodal LLaMA2-7B language model under two settings:

- **Deter**: deterministic decoding with the uncertainty instruction;
- Non-Deter: deterministic decoding without the instruction.

Results. Table 12 shows that the instruction causes the model to respond "I don't know" 238 times and reduces apparent accuracy from 44.8% to 25.6% (a drop of 19.2%). This confirms that nearly half of the seemingly correct answers in the uninstructed setting are likely due to lucky guesses rather than genuine reasoning, justifying our decision to include the clause in all main experiments.

"I don't know" will therefore be treated as an explicit abstention in the main study, ensuring reported accuracies reflect genuine vision-language capabilities rather than random chance.

Setting	Accuracy (%)	# "I don't know"
Deter (+ instruction)	25.6	238
Non-Deter (- instruction)	44.8	0

Table 12. Effect of the uncertainty instruction on LLaMA2-7B.