

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

Anonymous ACL submission

## 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 multi-phase training. We systematically analyze multimodal data contamination using our analytical framework, MM-DETECT, which defines two contamination categories—unimodal and cross-modal—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),

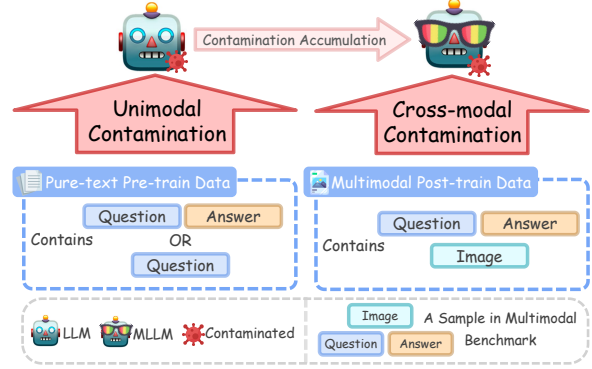


Figure 1: An analytical breakdown illustrating different forms and origins of multimodal data contamination across distinct training stages of MLLMs.

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*, as illustrated in Figure 1. 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, simulating 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 originate 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:

- We analytically characterize multimodal contamination into clearly defined unimodal and cross-modal categories, introducing MM-DETECT as an essential analytical tool.
- We systematically quantify how benchmark leakage inflates performance metrics, providing clear insights into dataset and model susceptibility to contamination.
- 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, as depicted in Figure 1. The training data for MLLMs generally consists of pure text pre-training data  $D_{\text{pretrain}}$  and multimodal alignment or instruction-following data  $D_{\text{vision}}$ . Consider an instance  $(x, i, y)$  from a benchmark dataset  $D$ , where  $x$  represents 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  $D_{\text{pretrain}}$ .
- **Cross-modal Contamination:** The triplet  $(x, i, y)$  appears in  $D_{\text{vision}}$ .

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 A.

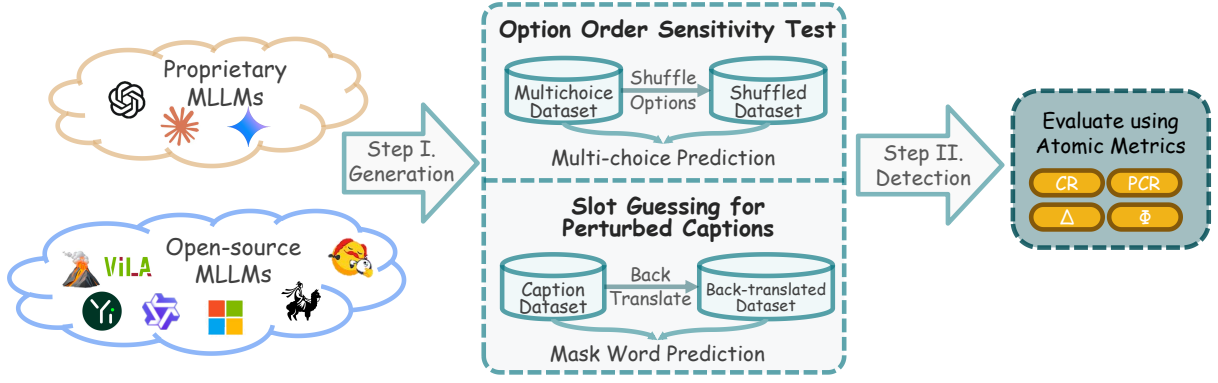


Figure 2: The overview of proposed MM-DETECT framework.

### Challenge II: Multi-stage Training in MLLMs.

Another challenge in detecting contamination in MLLMs is the multi-stage 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 semantic-irrelevant perturbations. As depicted in Figure 2, this framework operates in two primary steps:

- The first step is to generate perturbed datasets using two methods: *Option Order Sensitivity Test* (§3.1) and *Slot Guessing for Perturbed*

*Captions* (§3.2), 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, as shown in Figure 3, it indicates potential contamination, leading the model to memorize a certain canonical order of the options.

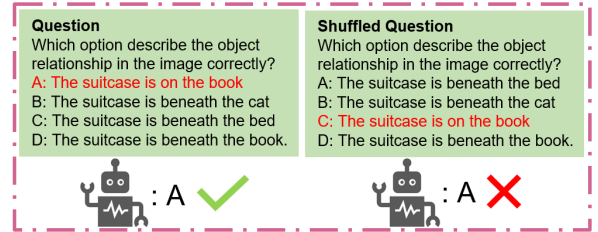


Figure 3: An example of **Option Order Sensitivity Test** applied to a contaminated model.

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

### 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.



Figure 4: A simple example shows the procedure.

As shown in Figure 4, 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  $d_i$  ( $i \in \{1, \dots, n\}$ ) 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 function, where we use the Google Translate API for Python to implement back-translation, to  $S_i$ :<sup>1</sup>

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

resulting in a paraphrased version  $S'_i$ . Next, we perform keyword extraction<sup>2</sup> on both  $S_i$  and  $S'_i$ :

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

where  $K_i$  and  $K'_i$  denote the extracted keywords from  $S_i$  and  $S'_i$ , respectively. We then apply a

<sup>1</sup>A quantitative analysis of the semantic and lexical similarity between the original and back-translated captions is provided in Appendix B.1.

<sup>2</sup>We employ the Stanford POS Tagger (Toutanova and Manning, 2000), targeting nouns, adjectives, and verbs, as they encapsulate the core meaning of the sentences.

masking function  $f_{\text{mask}}$  to replace the extracted keywords with a placeholder token [MASK]:

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

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

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

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

### 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:

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

**Step 2: Contamination Degree Analysis.** This step quantifies the model’s contamination degree based on the performance 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  $\Delta$ :

$$\Delta = PCR - CR$$

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  $\Delta$  suggests potential extensive leakage in the benchmark dataset, leading to highly perturbation-sensitive model performance.

**Instance Level Metric.** Despite a non-significant or positive  $\Delta$ , 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**  $\Phi$  is then obtained by dividing  $X$  by the dataset size:

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

where a larger  $\Phi$  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 subquestions:

**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 ( $D_{\text{tuning}}$ ) 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.

Table 1 shows that incorporating contaminated data during training increases both the model’s per-

Models	ScienceQA Test Set			NoCaps Val. Set		
	CR	PCR	$\Delta$	CR	PCR	$\Delta$
Baseline	61.4	61.5	0.01	33.0	32.1	-0.9
Early Cont.	71.5	68.1	<b>-3.4</b>	37.5	32.0	<b>-5.5</b>
Mid Cont.	69.4	67.3	<b>-2.1</b>	38.5	35.1	<b>-3.4</b>
Late Cont.	70.2	66.9	<b>-3.3</b>	38.7	32.6	<b>-6.1</b>

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

formance 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 NoCaps-contaminated models show increases of 5.2% and 1.1%. Moreover, all contaminated models demonstrate a marked decrease in  $\Delta$ , 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.

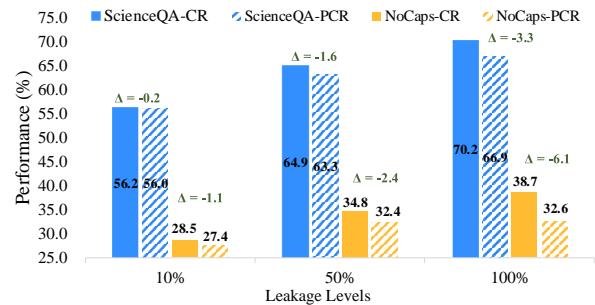


Figure 5: Performance and atomic metrics evaluated under varying leakage levels on the ScienceQA test set and NoCaps validation set.

As illustrated in Figure 5, increasing contamination from 10% to 50% to 100% results in corresponding increases in CR and PCR, alongside progressively larger  $\Delta$  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	$\Delta$
Clean	ScienceQA	61.4	61.5	0.01
Leaked	ScienceQA	64.3	63.8	<b>-0.5</b>
Clean	COCO-Caption2017	32.5	31.9	<b>-0.6</b>
Leaked	COCO-Caption2017	38.1	34.9	<b>-3.2</b>

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  $\Delta$  of -0.5. On the COCO-Caption2017 validation set, the model trained with COCO data shows a  $\Delta$  of -3.2. The results indicate that training set leakage inflates performance and that MM-DETECT effectively detects it.

#### Takeaways

*Both training and test set leakage can result in unfairness, and the degree of contamination can be detected through MM-DETECT effectively.*

## 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-8b<sup>3</sup>, 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-2024-06-20<sup>4</sup>.

**Benchmark Datasets.** Our analysis leverages two multi-choice datasets: ScienceQA (Lu et al., 2022) and MMStar (Chen et al., 2024a), along with

<sup>3</sup><https://www.adept.ai/blog/fuyu-8b>

<sup>4</sup><https://www.anthropic.com/news/claude-3-5-sonnet>

three caption datasets: COCO-Caption2017 (Lin et al., 2015), NoCaps (Agrawal et al., 2019), and Vintage<sup>5</sup>. 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  $\Delta$  with higher  $\Phi$  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  $\Delta$  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  $\Delta$  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-4o 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

<sup>5</sup><https://huggingface.co/datasets/SilentAntagonist/vintage-artworks-60k-captioned>

<sup>6</sup>Based on intentional contamination experiments in §4.1, the degrees on multi-choice datasets are defined as follows:  $\Delta \in (-1.6, -0.2]$  for minor leakage,  $\Delta \in (-2.9, -1.6]$  for partial leakage, and  $\Delta \leq -2.9$  for severe leakage.

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

Model	ScienceQA Training Set				ScienceQA Test Set				MMStar Validation Set			
Metric	CR	PCR	$\Delta$	$\Phi$	CR	PCR	$\Delta$	$\Phi$	CR	PCR	$\Delta$	$\Phi$
<i>Open-source MLLMs</i>												
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	<b>14.9</b>	28.2	27.0	-1.2	<b>17.7</b>
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
<i>Proprietary MLLMs</i>												
GPT-4o	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	<b>15.3</b>	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  $\Delta$  or  $\Phi$ ; color codes denote contamination degree: **green** for minor leakage, **yellow** for partial leakage, and **red** for severe leakage.<sup>6</sup>

Model	COCO Validation Set				NoCaps Validation Set				Vintage Training Set			
Metric	CR	PCR	$\Delta$	$\Phi$	CR	PCR	$\Delta$	$\Phi$	CR	PCR	$\Delta$	$\Phi$
<i>Open-source MLLMs</i>												
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	<b>23.3</b>	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
<i>Proprietary MLLMs</i>												
GPT-4o	58.1	54.4	-3.7	<b>23.1</b>	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	<b>28.3</b>
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  $\Delta$  or  $\Phi$ ; color codes denote contamination degree: **green** for minor leakage, **yellow** for partial leakage, and **red** for severe leakage.<sup>7</sup>

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.

#### Takeaways

Multimodal data contamination, at both dataset and instance levels, is prevalent in open-source and proprietary MLLMs across multi-choice and image caption datasets.

## 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 during the multimodal training phase, 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

We employed a heuristic approach based on the intuition that if an LLM 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 idfics2), Phi-3-small-128k-instruct (Abdin et al., 2024) (used by Phi-3-vision), Yi-6B (AI et al., 2024) (used by Yi-VL), and Internlm2-7B (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 B.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 Fuyu-8B and proprietary models since their unimodal LLM components and training data remain undisclosed.

Model	Accuracy	$\Phi_M$
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 (idfics2)	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.  $\Phi_M$  denotes the  $\Phi$  of the respective MLLMs.

**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  $\Phi$  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 image-text

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: Depiction of the overlap between the training data of MLLMs and the benchmarks, as well as the contamination degree  $\Delta$  of MLLMs on benchmarks. **Green** signifies no overlap, **yellow** suggests potential overlap, and **Red** indicates partial or entire overlap.

datasets (e.g., OBELICS (Laurenson 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 in benchmark datasets.

### Takeaways

*The contamination in MLLMs may not only stem from cross-modal contamination but also from unimodal contamination, both of which can significantly impact the overall performance.*

## 7 Conclusion and Future Work

In this study, we systematically analyzed multimodal data contamination in MLLMs through our proposed detection framework, MM-DETECT. We demonstrated that MM-DETECT effectively quantifies and detects varying contamination degrees, revealing significant performance biases induced by benchmark leakage. Importantly, we identified that contamination originates notably from both unimodal pre-training and multimodal fine-tuning phases, impacting the reliability and fairness of multimodal evaluations.

Future work will focus on two key areas:

- Firstly, standardizing the use of multimodal datasets and reporting potential contamination impacts to minimize contamination, thereby enhancing data consistency and quality.
- Secondly, creating a continuously updated benchmarking system for the ongoing evaluation of multimodal model performance.

This will support advancements and broader applications in this field.

## Limitations

We acknowledge several limitations in our work. First, this work is limited to discussions around visual modalities, and does not yet cover other modalities such as audio or video. Second, we only selected widely used and representative multimodal datasets for detection, including multiple-choice datasets and caption datasets, without testing additional datasets, such as open-ended generation and cloze questions. However, we speculate that the method *Slot Guessing for Perturbed Caption* may also apply to other types of image-feature-analyzing benchmarks. Third, the effectiveness of *Option Order Sensitivity Test* can be undermined by option shuffling, which, while potentially improving model performance, is computationally expensive and may increase the training cost. Fourth, as a perturbation-based black-box detector, MM-DETECT might underestimate contamination if a model generalizes sufficiently to answer perturbed questions correctly. Although dataset-level evaluations reduce this risk, completely eliminating such false-negative cases remains an open challenge.

## References

- Marah I Abidin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat S. Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Caio César Teodoro Mendes, Weizhu Chen, Vishrav Chaudhary, Parul Chopra, Allie Del Giorno, Gustavo de Rosa, Matthew Dixon, Ronen Eldan, Dan Iter, Amit Garg, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Jamie Huynh, Mojan Javaheripi, Xin Jin, Piero Kauffmann, Nikos Karampatziakis, Dongwoo Kim, Mahoud Khademi, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Chen Liang, Weishung Liu, Eric Lin, Zeqi Lin, Piyush Madan, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacroce, Shital Shah, Ning Shang, Hiteshi Sharma, Xia Song, Masahiro Tanaka, Xin Wang, Rachel Ward, Guanhua Wang, Philipp Witte, Michael Wyatt, Can Xu, Jiahang Xu, Sonali Yadav, Fan Yang, Ziyi Yang, Donghan Yu, Chengruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and Xiren Zhou. 2024. [Phi-3 technical report: A highly capable language model locally on your phone](#). *CoRR*, abs/2404.14219.
- Harsh Agrawal, Karan Desai, Yufei Wang, Xinlei Chen, Rishabh Jain, Mark Johnson, Dhruv Batra, Devi Parikh, Stefan Lee, and Peter Anderson. 2019. [no-caps: novel object captioning at scale](#). In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*. IEEE.
01. AI, :, Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, Kaidong Yu, Peng Liu, Qiang Liu, Shawn Yue, Senbin Yang, Shiming Yang, Tao Yu, Wen Xie, Wenhao Huang, Xiaohui Hu, Xiaoyi Ren, Xinyao Niu, Pengcheng Nie, Yuchi Xu, Yudong Liu, Yue Wang, Yuxuan Cai, Zhenyu Gu, Zhiyuan Liu, and Zonghong Dai. 2024. [Yi: Open foundation models by 01.ai](#). *Preprint*, arXiv:2403.04652.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. 2023a. [Qwen technical report](#). *CoRR*, abs/2309.16609.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023b. [Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond](#). *Preprint*, arXiv:2308.12966.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). *CoRR*, abs/2005.14165.
- Zheng Cai, Maosong Cao, Haojiong Chen, Kai Chen, Keyu Chen, Xin Chen, Xun Chen, Zehui Chen, Zhi Chen, Pei Chu, Xiaoyi Dong, Haodong Duan, Qi Fan, Zhao Fei, Yang Gao, Jiaye Ge, Chenya Gu, Yuzhe Gu, Tao Gui, Aijia Guo, Qipeng Guo, Conghui He, Yingfan Hu, Ting Huang, Tao Jiang, Penglong Jiao, Zhenjiang Jin, Zhikai Lei, Jiaxing Li, Jingwen Li, Linyang Li, Shuaibin Li, Wei Li, Yining Li, Hongwei Liu, Jiangning Liu, Jiawei Hong, Kaiwen Liu, Kuikun Liu, Xiaoran Liu, Chengqi Lv, Haijun Lv, Kai Lv, Li Ma, Runyuan Ma, Zerun Ma, Wenchang Ning, Linke Ouyang, Jiantao Qiu, Yuan Qu, Fukai Shang, Yunfan Shao, Demin Song, Zifan Song, Zhihao Sui, Peng Sun, Yu Sun, Huanze Tang, Bin Wang,

- Guoteng Wang, Jiaqi Wang, Jiayu Wang, Rui Wang, Yudong Wang, Ziyi Wang, Xingjian Wei, Qizhen Weng, Fan Wu, Yingdong Xiong, and et al. 2024. [Internlm2 technical report](#). *CoRR*, abs/2403.17297.
- Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao, Dahua Lin, et al. 2024a. Are we on the right way for evaluating large vision-language models? *arXiv preprint arXiv:2403.20330*.
- Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. 2024b. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*.
- Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo, Tong Lu, Yu Qiao, and Jifeng Dai. 2023. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. *arXiv preprint arXiv:2312.14238*.
- Chunyu Deng, Yilun Zhao, Xiangru Tang, Mark Gestein, and Arman Cohan. 2024. [Investigating data contamination in modern benchmarks for large language models](#). In *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, pages 8706–8719. Association for Computational Linguistics.
- Yihong Dong, Xue Jiang, Huanyu Liu, Zhi Jin, Bin Gu, Mengfei Yang, and Ge Li. 2024. [Generalization or memorization: Data contamination and trustworthy evaluation for large language models](#). In *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, pages 12039–12050. Association for Computational Linguistics.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Léo R. Lavau, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. [Mistral 7b](#). *Preprint*, arXiv:2310.06825.
- Hugo Laurençon, Léo Tronchon, Matthieu Cord, and Victor Sanh. 2024. [What matters when building vision-language models?](#) *CoRR*, abs/2405.02246.
- Hugo Laurençon, Lucile Saulnier, Léo Tronchon, Stas Bekman, Amanpreet Singh, Anton Lozhkov, Thomas Wang, Siddharth Karamcheti, Alexander M. Rush, Douwe Kiela, Matthieu Cord, and Victor Sanh. 2023. [Obelics: An open web-scale filtered dataset of interleaved image-text documents](#). *Preprint*, arXiv:2306.16527.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. [Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models](#). In *International conference on machine learning*, pages 19730–19742. PMLR.
- Ji Lin, Hongxu Yin, Wei Ping, Yao Lu, Pavlo Molchanov, Andrew Tao, Huizi Mao, Jan Kautz, Mohammad Shoyebi, and Song Han. 2023. [VILA: on pre-training for visual language models](#). *CoRR*, abs/2312.07533.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. 2015. [Microsoft coco: Common objects in context](#). *Preprint*, arXiv:1405.0312.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023a. [Improved baselines with visual instruction tuning](#). *CoRR*, abs/2310.03744.
- Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, Kai Chen, and Dahua Lin. 2023b. [Mmbench: Is your multi-modal model an all-around player?](#) *CoRR*, abs/2307.06281.
- Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Taffjord, Peter Clark, and Ashwin Kalyan. 2022. Learn to explain: Multimodal reasoning via thought chains for science question answering. In *The 36th Conference on Neural Information Processing Systems (NeurIPS)*.
- OpenAI. 2023. [GPT-4 technical report](#). *CoRR*, abs/2303.08774.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.
- Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*.
- Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. 2018. Conceptual captions: A cleaned, hypernamed, image alt-text dataset for automatic image captioning. In *Proceedings of ACL*.
- Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi Chen, and Luke Zettlemoyer. 2024. [Detecting pretraining data from large language models](#). In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.

Dingjie Song, Shunian Chen, Guiming Hardy Chen, Fei Yu, Xiang Wan, and Benyou Wang. 2024. Milebench: Benchmarking mllms in long context. *arXiv preprint arXiv:2404.18532*.

Kristina Toutanova and Christopher D. Manning. 2000. Enriching the knowledge sources used in a maximum entropy part-of-speech tagger. In *Joint SIGDAT Conference on Empirical Methods in Natural Language Processing and Very Large Corpora, EMNLP 2000, Hong Kong, October 7-8, 2000*, pages 63–70. Association for Computational Linguistics.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023a. *Llama 2: Open foundation and fine-tuned chat models*. *CoRR*, abs/2307.09288.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023b. *Llama 2: Open foundation and fine-tuned chat models*. *CoRR*, abs/2307.09288.

Zhiyu Wu, Xiaokang Chen, Zizheng Pan, Xingchao Liu, Wen Liu, Damai Dai, Huazuo Gao, Yiyang Ma, Chengyue Wu, Bingxuan Wang, Zhenda Xie, Yu Wu, Kai Hu, Jiawei Wang, Yaofeng Sun, Yukun Li, Yishi

Piao, Kang Guan, Aixin Liu, Xin Xie, Yuxiang You, Kai Dong, Xingkai Yu, Haowei Zhang, Liang Zhao, Yisong Wang, and Chong Ruan. 2024. *Deepseek-vl2: Mixture-of-experts vision-language models for advanced multimodal understanding*. *Preprint*, arXiv:2412.10302.

Ruijie Xu, Zengzhi Wang, Run-Ze Fan, and Pengfei Liu. 2024. Benchmarking benchmark leakage in large language models. *arXiv preprint arXiv:2404.18824*.

Samuel Yeom, Irene Giacomelli, Matt Fredrikson, and Somesh Jha. 2018. *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*, pages 268–282. IEEE Computer Society.

Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. 2023. A survey on multimodal large language models. *arXiv preprint arXiv:2306.13549*.

## A Inefficiency of Unimodal Methods

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

### A.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).

### A.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).

### A.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).

## B Other Experiments

### B.1 Semantic & Lexical Similarity After Back-Translation

**Setup.** To quantify how much meaning and wording change during our *caption perturbation* step (§3.2), we applied an **English**→**Chinese**→**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 and Gurevych, 2019) as a sentence-level *semantic* score, and

- **BLEU-4** (Papineni et al., 2002) as a token-overlap *lexical* score.

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

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

Table 10: 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.

### B.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 11 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.

Setting	Accuracy (%)	# “I don’t know”
Deter (+ instruction)	25.6	238
NonDeter (- instruction)	44.8	0

Table 11: Effect of the uncertainty instruction on LLaMA2-7B.

“I don’t know” will therefore be treated as an explicit abstention in the main study, ensuring reported accuracies reflect genuine visionlanguage capabilities rather than random chance.