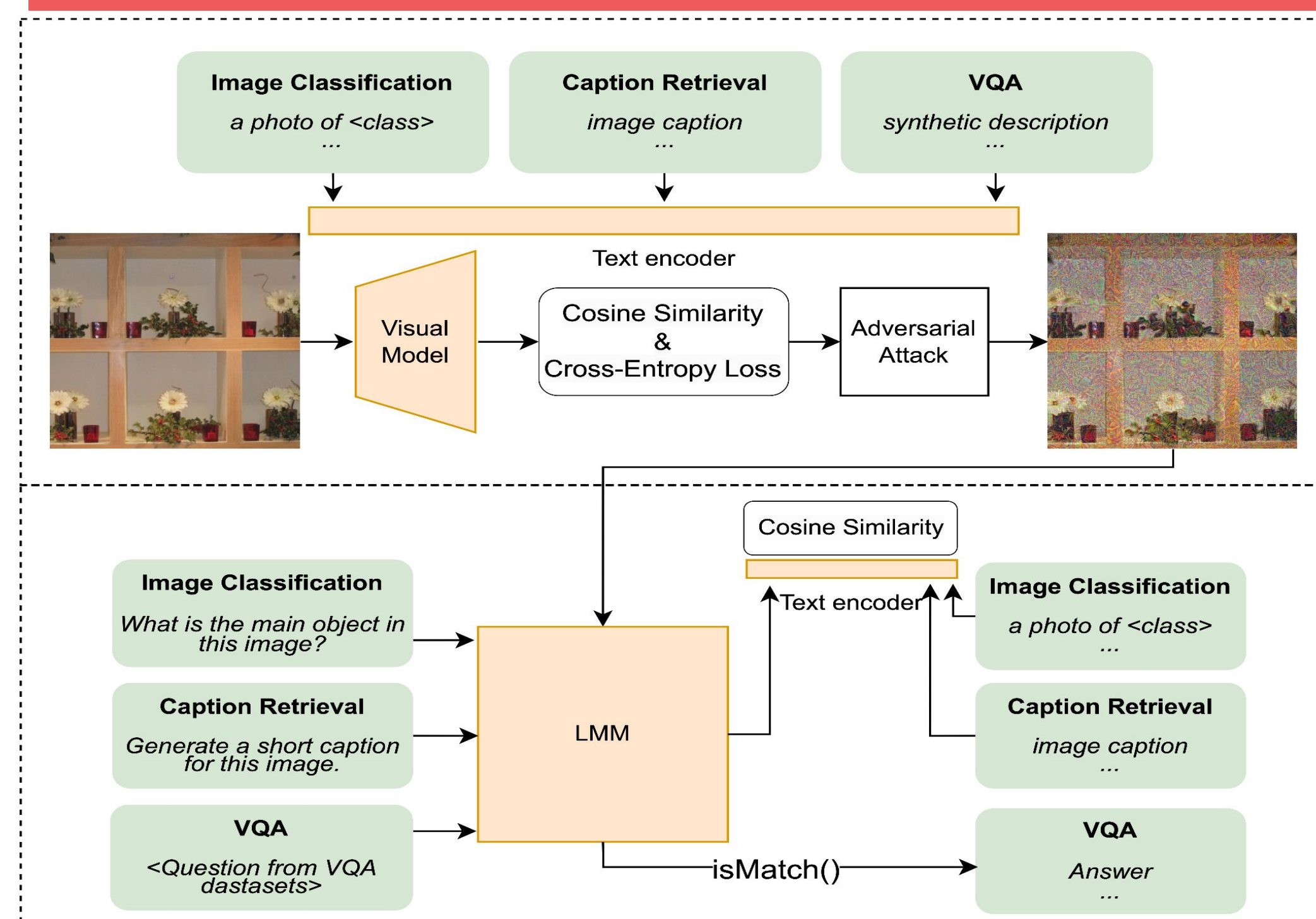## Introduction

Overview:

A comprehensive study of the performance and behavior of Large Multimodal Models (LMM)s under white-box visual adversarial attack.

Takeaways:

- LMMs are generally vulnerable to visual adversarial perturbations.
- Visual adversarial attacks are not *universal:* they are less effective when the query targets different visual contents.
- Adding additional textual context improves LMMs' robustness against visual adversarial input

## Setup

## Empirical Results

While the model completely misinterprets the object, it can still correctly answer questions that are not directly related with what's being attacked

## VQA

Table 1. Accuracy percent dropped post-attack comparing to pre-attack. We can observe LMMs' VQA accuracy drops much less than that of their visual encoders'.

| dataset | LLaVA1.5 | | BLIP2-T5 | | InstructBLIP | |
|---|---|---|---|---|---|---|
| | VQA | Vis.Enc. | VQA | Vis.Enc. | VQA | Vis.Enc. |
| MME | 29.3 | 97.3 | 32.7 | 97.0 | 29.7 | 97.0 |
| POPE | 19.0 | 99.3 | 19.3 | 98.3 | 25.7 | 98.3 |
| ScienceQA | 4.4 | 98.7 | 5.8 | 99.7 | 5.1 | 99.7 |
| TextVQA | 23.3 | 99.3 | 25.3 | 99.7 | 39.3 | 99.7 |
| VQAV2 | 27.4 | 97.0 | 29.3 | 99.0 | 35.0 | 99.0 |

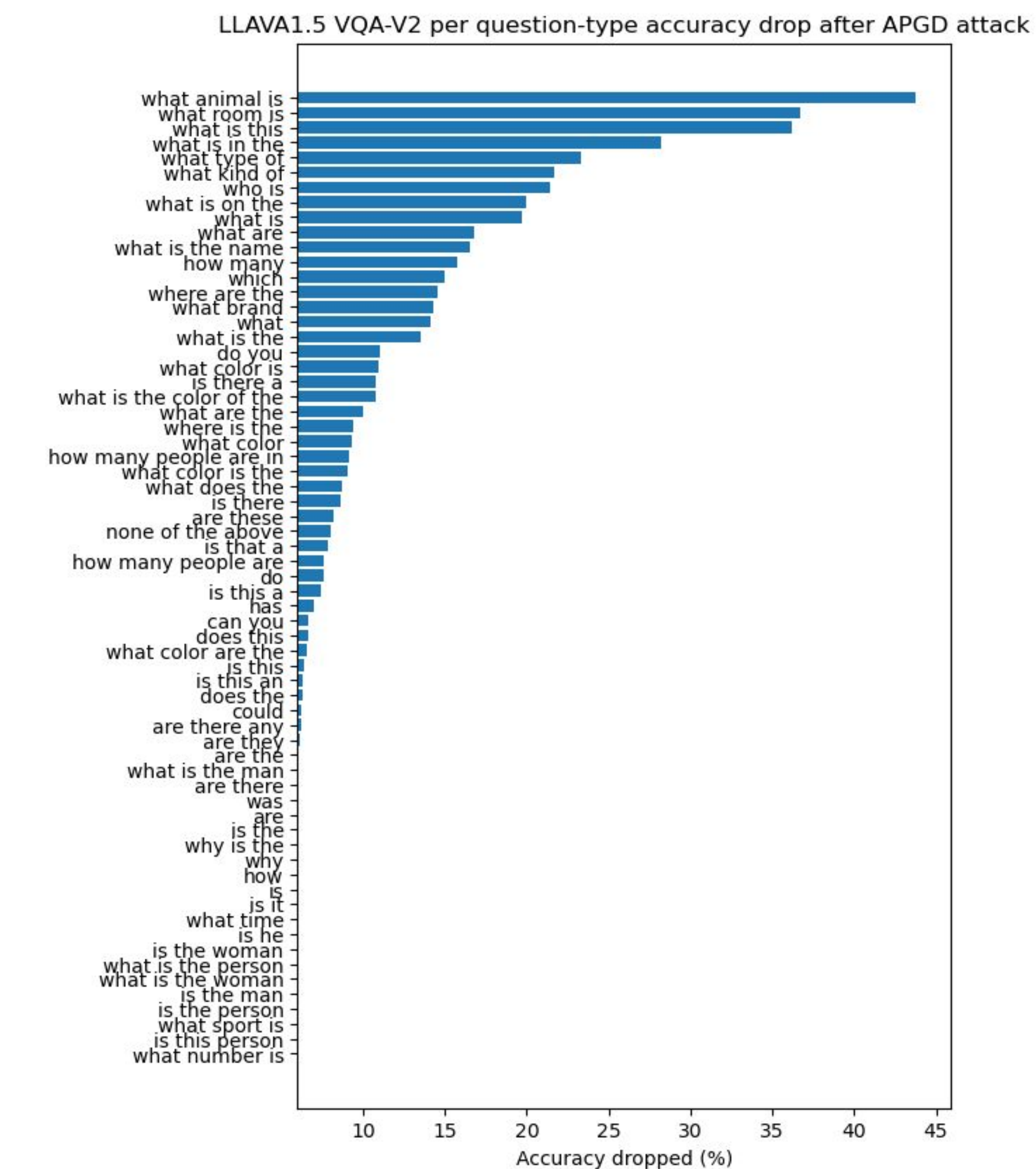VQA accuracy drops the most when questions are asking "what".

## Image Captioning/Class.

Table 2. Accuracy percent dropped post-attack comparing to pre-attack. Comparing to VQA, LMM accuracy drops are much more obvious on Image classification and caption retrieval.

| Dataset | LLaVA | BLIP2 | InstructBLIP | CLIP | BLIP |
|---|---|---|---|---|---|
| | Caption Retrieval Acc. Drop (%) | | | | |
| COCO2014 | 62 | 99 | 78 | 98 | 78 |
| | Image Classification Acc. Drop (%) | | | | |
| COCO2014 | 31 | 75 | 86 | 98 | 100 |
| Food101 | 79 | 83 | 75 | 96 | 100 |
| StanfordCars | 75 | 99 | 70 | 100 | 99 |

## Context aids LMM Robustness against adversaries

Adding context of object description helps improve LMMs' robustness against adversarial visual inputs.

We can thus 'decompose' the query into a series of existential questions, each with the corresponding object context descriptio, and select the answer based on highest token probability.