# SCRAMBLe : Enhancing Multimodal LLM Compositionality with Synthetic Preference Data

Samarth Mishra*, Kate Saenko, Venkatesh Saligrama

BOSTON UNIVERSITY

ICCV OCT 19-23, 2025 HONOLULU HAWAII

## Can Multimodal LLMs distinguish between



a dog chasing a cat

a cat chasing a dog

On Winoground, a benchmark of such questions,
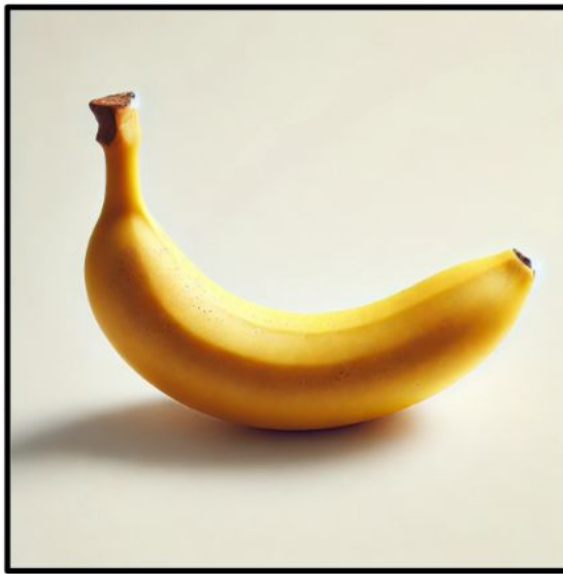MLLM (GPT-4V) performance = 33%
Human performance = 85%

## Compositionality:

Identifying parts/concepts and how they are composed to give rise to a given scene
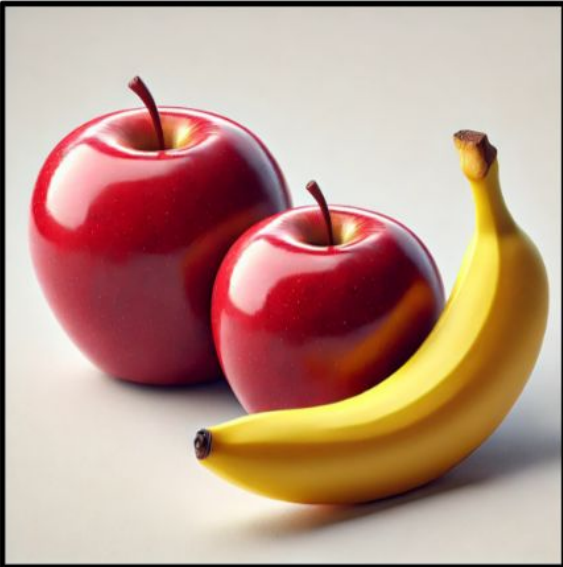
**Atomic Concepts**



Apple

Banana

**Infinite Compositions - Nothing New**



Two apples and a banana

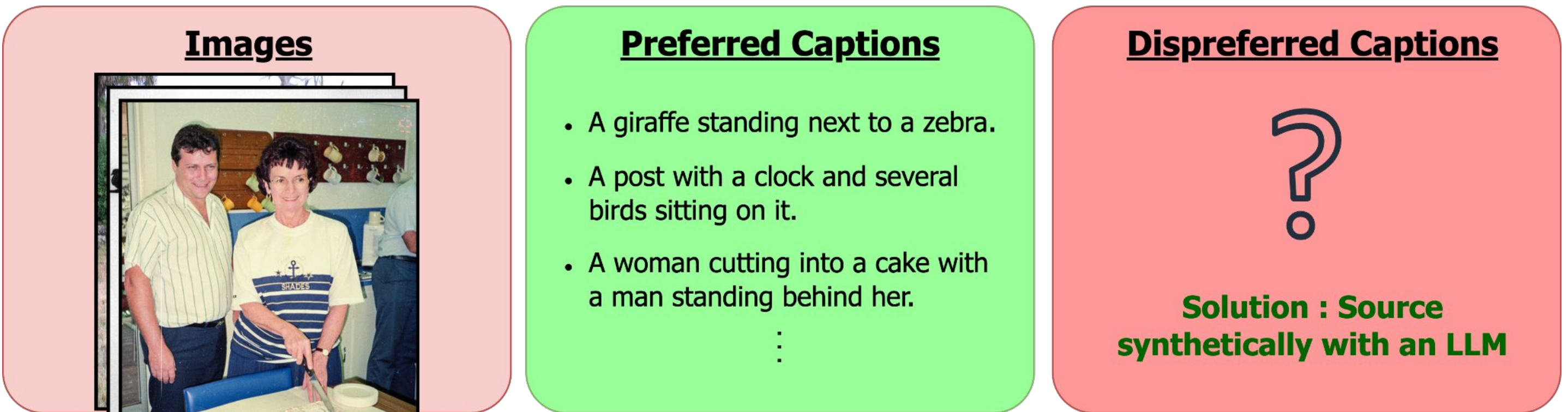An apple and two bananas

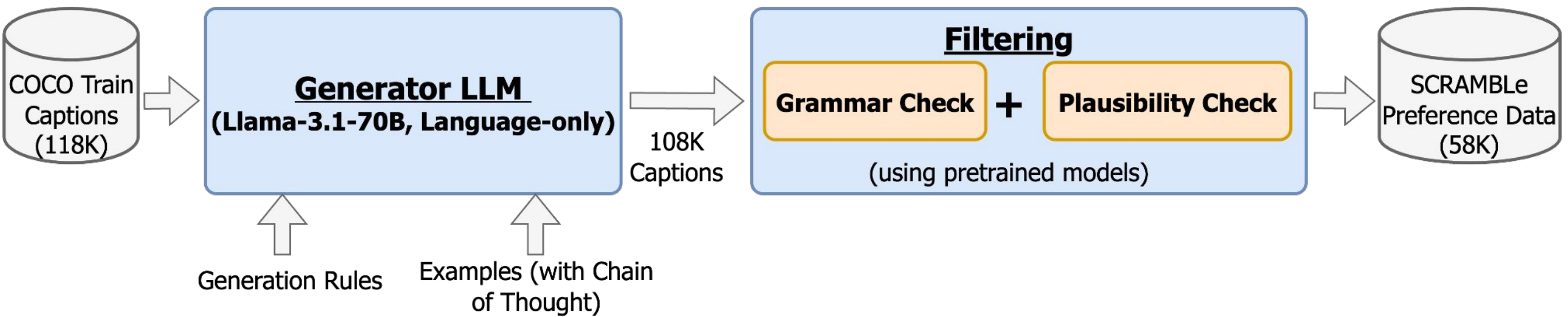A banana on an apple

An apple on a banana

## Improving Compositionality with Preference Tuning



Caption this

**Preferred**
A cat is chasing a dog through a park

**Dispreferred**
A dog is chasing a cat through a park

## Data Source : Existing Image Caption Dataset (COCO)

**Images**



**Preferred Captions**
- A giraffe standing next to a zebra.
- A post with a clock and several birds sitting on it.
- A woman cutting into a cake with a man standing behind her.
⋮

**Dispreferred Captions**
?
**Solution : Source synthetically with an LLM**

## Synthetic Data Generation



COCO Train Captions (118K) → **Generator LLM** (Llama-3.1-70B, Language-only) → 108K Captions → **Filtering** [Grammar Check + Plausibility Check] (using pretrained models) → SCRAMBLe Preference Data (58K)

Generation Rules    Examples (with Chain of Thought)

## Winoground Accuracy



*GPT-4V eval based on verbose responses and LLM as a judge [1]. In contrast LLaVA and Molmo use VQAScore [2] based on model output logprobs.

## Hmm, could this be overfitting?

Seemingly not. Models retain or improve general QA performance (on Control Benchmarks)

| Model Name | Compositionality Benchmarks | | | Control Benchmarks | |
|---|---|---|---|---|---|
| | Winoground | COLA | ConMe | SEED-Bench | MM-Vet |
| LLaVA-1.5-13B | 36.5 | 49.5 | 62.3 | **68.23** | 36.2 $\pm$ 0.3 |
| + SCRAMBLe | **39.3** | **55.7** | **64.5** | 68.19 | **38.6 $\pm$ 0.1** |
| MoLMo-7B | 49.5 | 57.1 | 72.2 | 74.04 | 59.3 $\pm$ 0.2 |
| + SCRAMBLe | **54.8** | **60.5** | **74.6** | **74.61** | **60.9 $\pm$ 0.4** |

## Chats with SCRAMBLe-Molmo



What is the man holding in the image?

**Molmo**
The man in the image is **holding a hose with a spray nozzle. He's using it to spray water** onto the side of a black pickup truck, likely as part of a car washing process.

**SCRAMBLe-Molmo**
The man in the image is **holding a white mitt. This mitt is likely being used as part of the car washing process, to apply soap or water to the vehicle's surface as he cleans.**



Can you describe the image?

**Molmo**
The image shows a young woman sitting on a bed in a bedroom..... The woman is barefoot and **has a silver plate on her lap.** Her feet are positioned differently.....

**SCRAMBLe-Molmo**
The image shows a young woman sitting on the edge of a bed in a bedroom..... **There's a small silver tray on the bed next to her.** The room has white walls and a light gray wooden floor.....

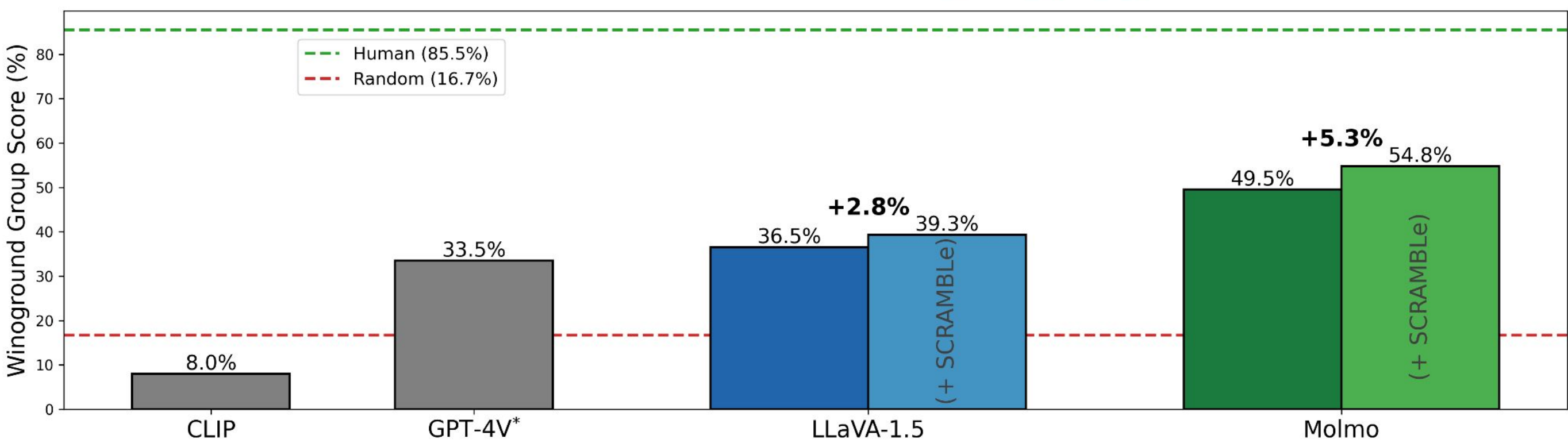## References
[1] Mitra, Chancharik, et al. "Compositional chain-of-thought prompting for large multimodal models." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2024.
[2] Lin, Zhiqiu, et al. "Evaluating text-to-visual generation with image-to-text generation." European Conference on Computer Vision. Cham: Springer Nature Switzerland, 2024.

* Samarth is now affiliated with Waymo. Work done while he was at Boston University