

# CAP : Evaluation of Persuasive and Creative Image Generation

Aysan Aghazadeh, Adriana Kovashka Department of Computer Science - University of Pittsburgh





# Overview

### **Prompt** → **Ad Image**

- Prompt: Action-Reason Message (AR) I should {action} because {reason}
- Ad Image: Creative and Persuasive

### Commercial and Public Service Ads



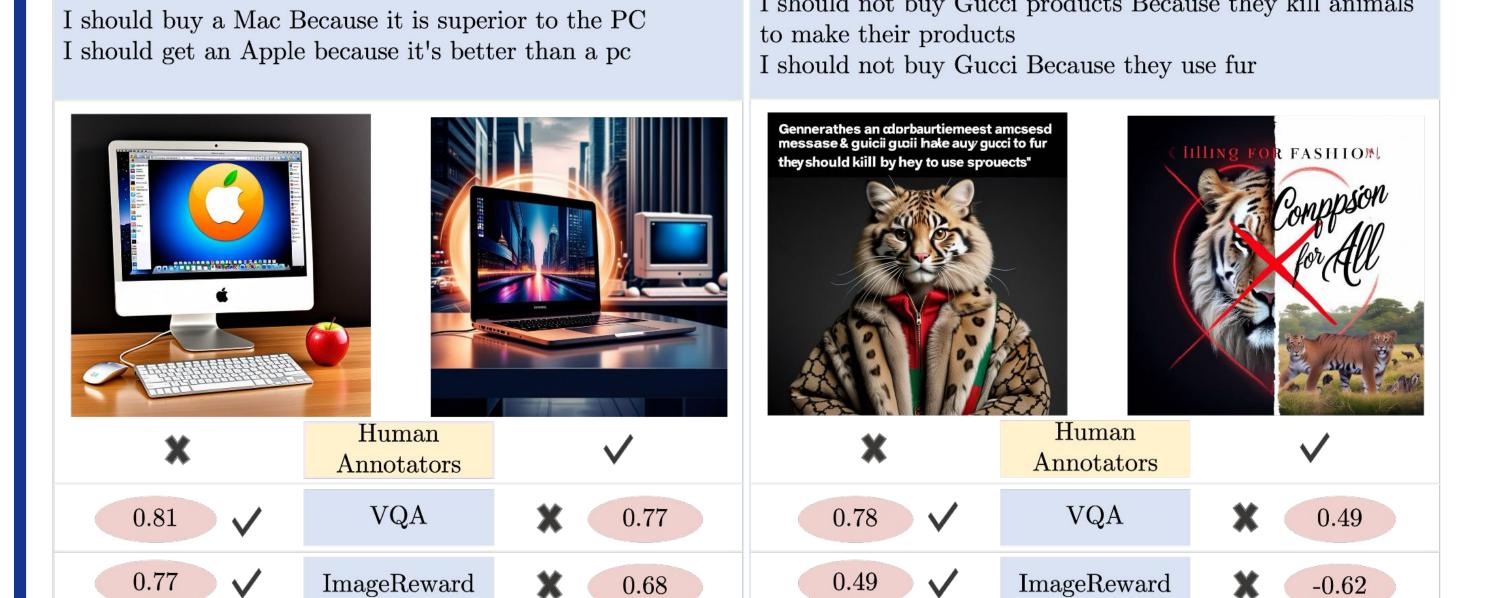
# • Prompt is implicit: Different images convey same AR message → Visual Matching is not enough

Persuasion and Creativity are subjective

# Alignment of Image and Message (AIM)

### Motivation

Visual Matching vs Semantic Matching

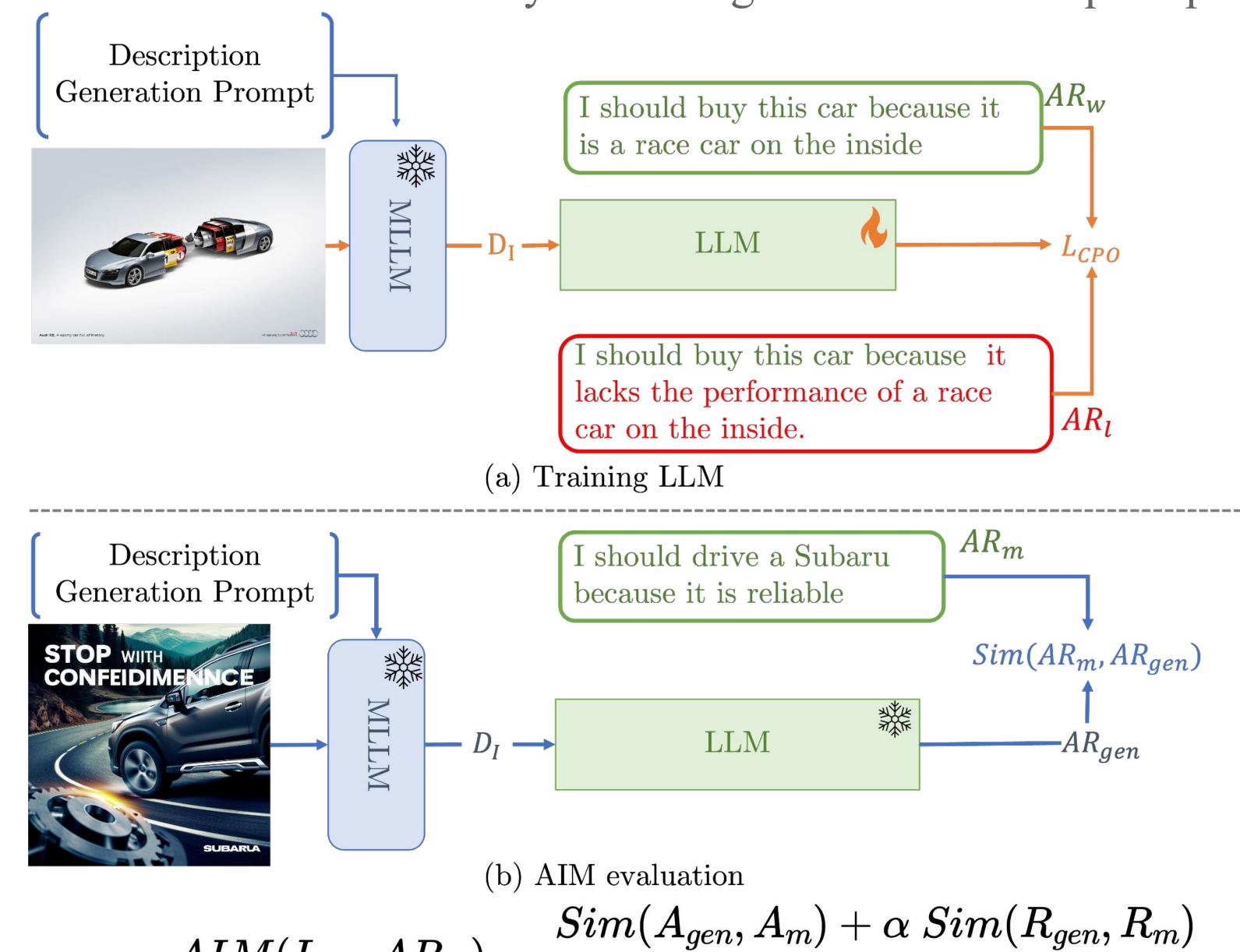


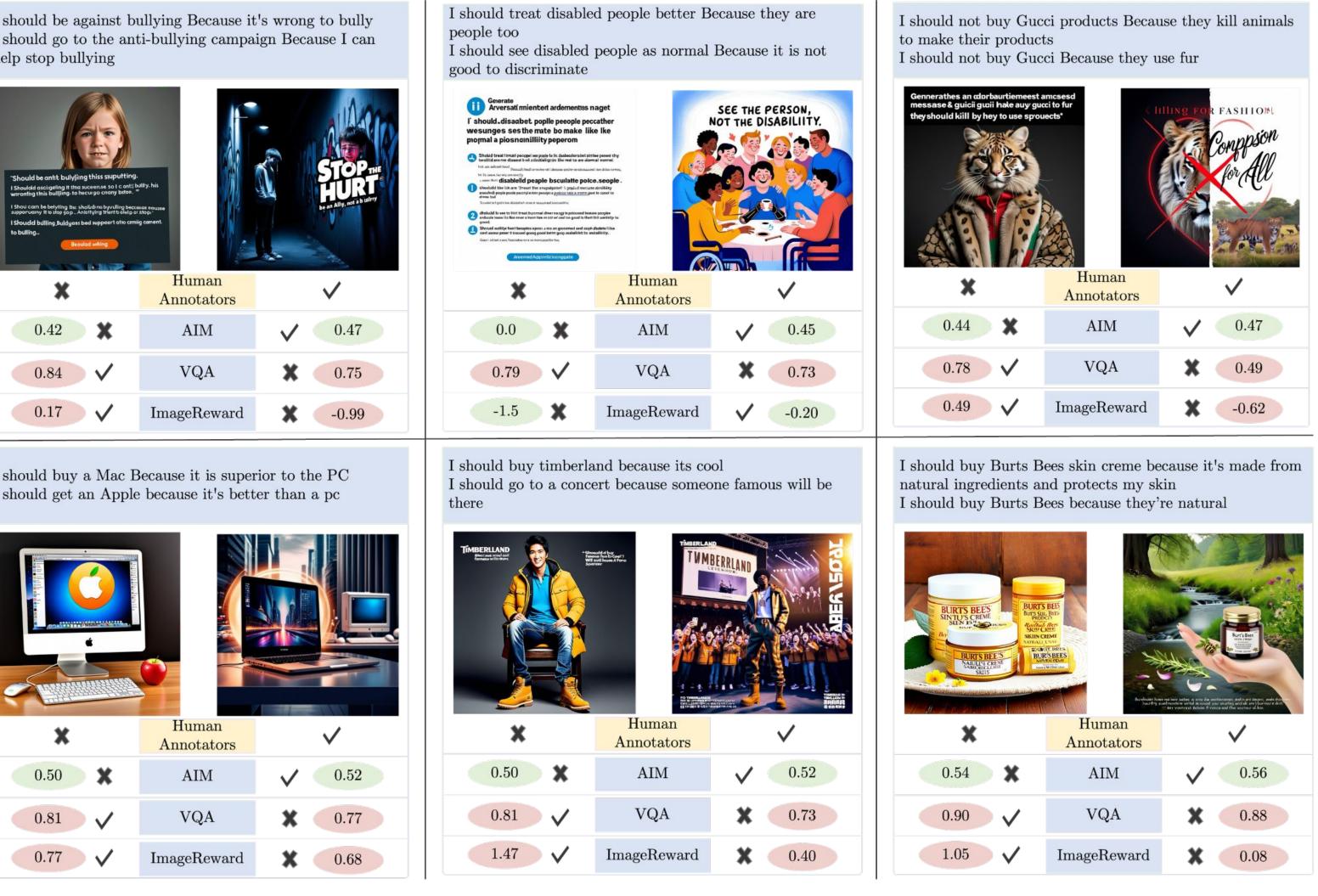
### Related works

- Image-Image Similarity FID (Heusel, et al.), IS (Salimans, et al.)
- Image-Text Similarity:
  - CLIPScore (Hessel, et al.): sim(CLIP embeddings)
  - Training LLMs/VLMs VQA-score (Lin, et al.): P(Yes Prompt, I) from LLM ImageReward (Xu, et al.): Reward Model
  - O Zero-shot LLMs: DSG (Cho, et al.)

## **Evaluation Method**

- We fine-tune LLM to interpret the image semantically accurate
- AR generation: Image description generation→AR generation
- Score: Weighted average of action text similarity and reason text-text similarity between generated AR and prompt



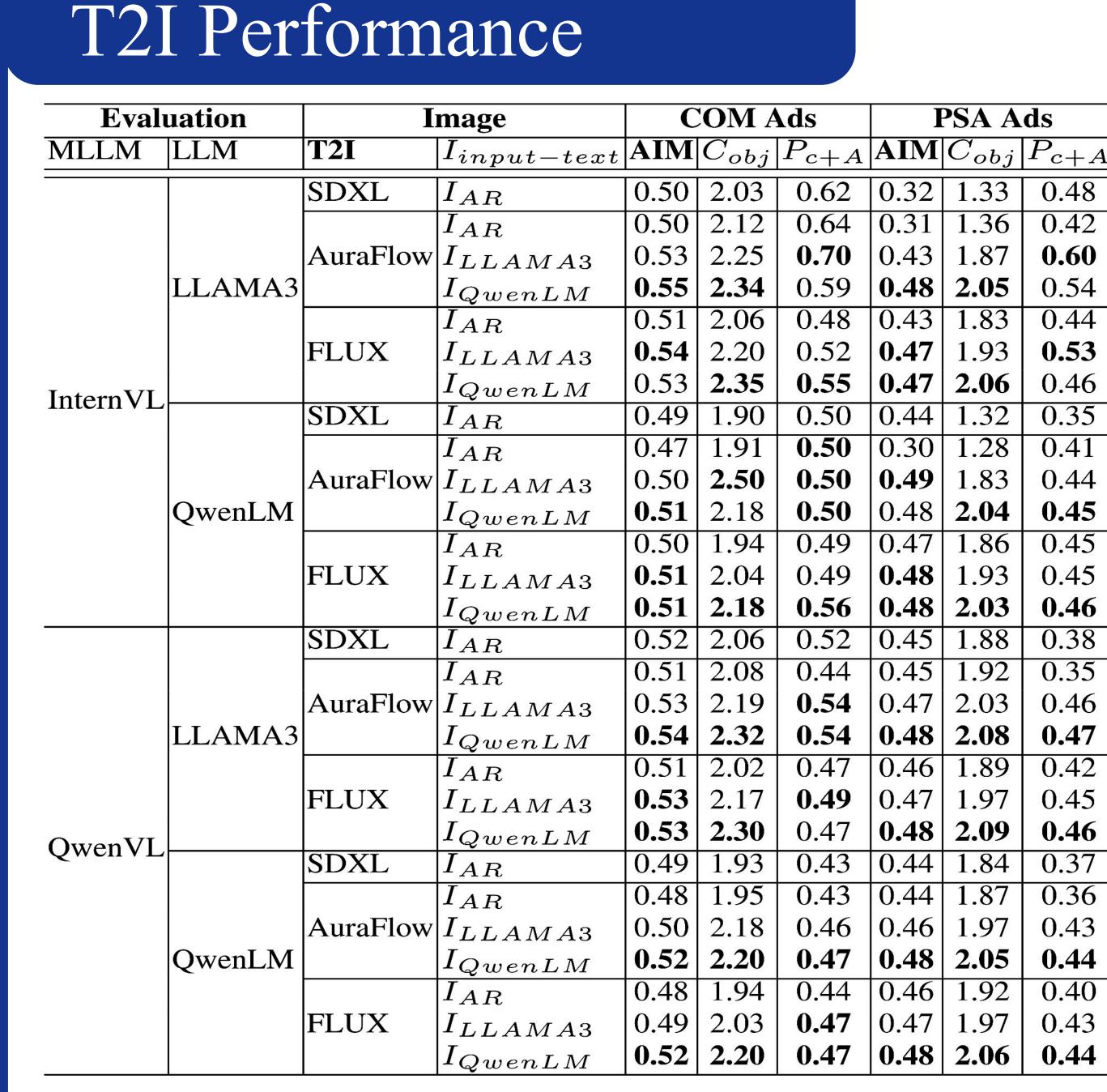


## Results

- Baselines do not capture semantic mismatch
- +030/2 higher agreement than hagalin

+93% higher agreement	than baselines
COM: 57%	Annotators
PSA: 140%	H, ImageReward H, VQAScore
Low 0-shot agreement	H, CLIPScore H, AIM (InternVL,
(Krippendorff's Alpha)	H, AIM (InternVL, H, AIM (InternVL,

COM: 57%	Annotators	COM PSA  All
PSA: 140%	H, ImageReward H, VQAScore H, CLIPScore	0.12     0.06     0.11       0.38     0.24     0.31       0.04     0.34     0.17
Low 0-shot agreement (Krippendorff's Alpha)  • Importance of training	H, AIM (InternVL, LLAMA3) (0-shot) H, AIM (InternVL, LLAMA3) H, AIM (InternVL, QwenLM)	
	H1, H2	0.86   0.85   0.86



- Low AIM, C<sub>obi</sub>, and P<sub>C+A</sub> on implicit prompt
- Lower AIM, C<sub>obi</sub>, and P<sub>C+A</sub> on PSA ads (more implicit prompts)
- Higher scores when using explicit descriptions generated by LLMs

# Creativity of Images

# Creativity Criteria

## **Evaluation Method**

 Uniqueness → Distant from visual elements in AR (Sim with elements ↓)

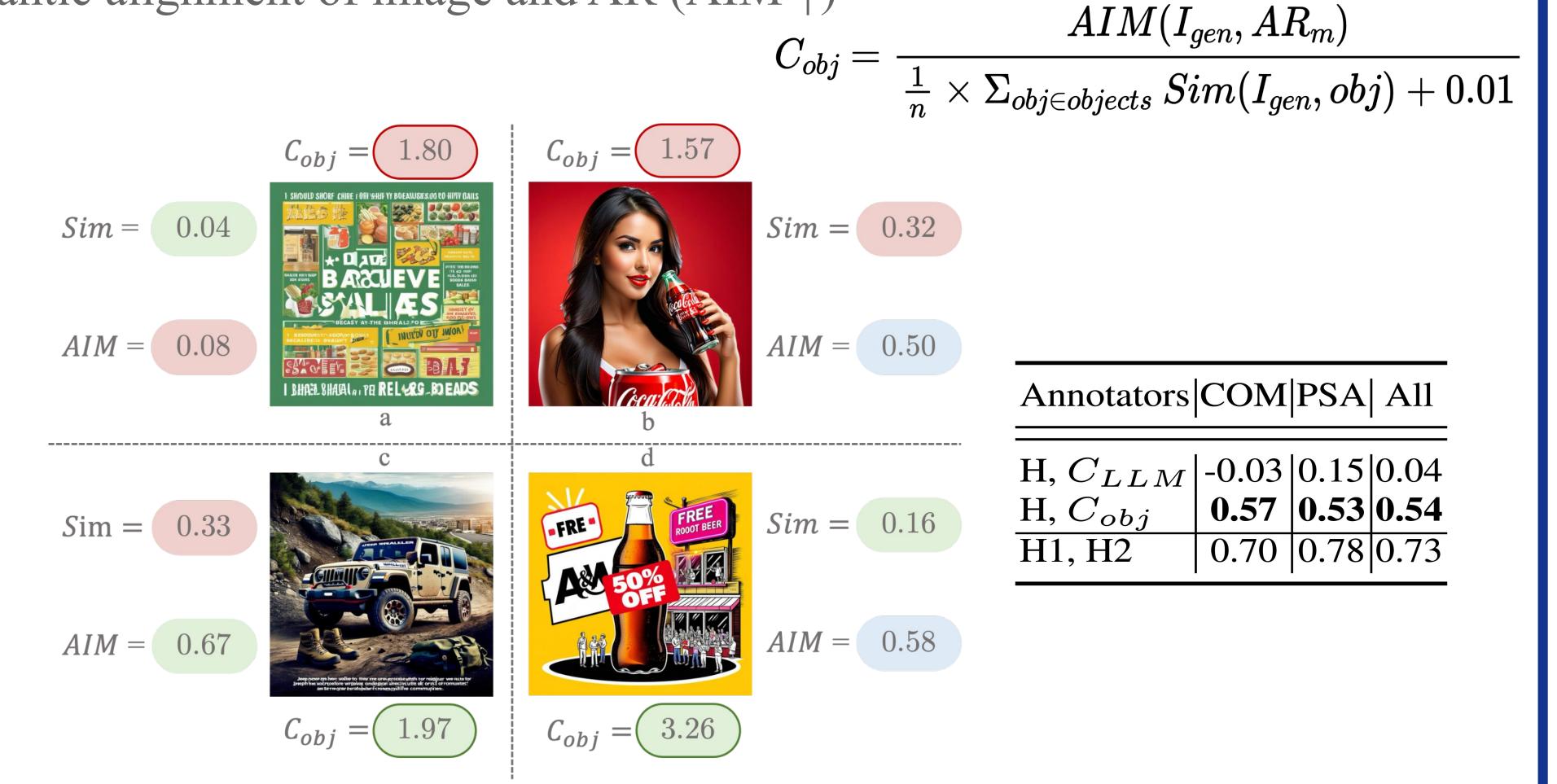
. I should not allow children to

• Conveying the message → Higher semantic alignment of image and AR (AIM ↑)

# Results

Challenge

- LLMs do not evaluate creativity accurately (No prior creativity metric)
- Agreement (Krippendorff's Alpha) 0.5 (out of 1) improvement compared to LLMs
- Balancing criteria
  - Uniqueness of Image (Sim)
  - Conveying the AR message (AIM)



# Contributions

 $AIM(I_{qen},AR_m)=rac{\sum m_{i}}{m}$ 

# Proposed CAP evaluation:

- Creativity: Balancing uniqueness and alignment
- Alignment: Evaluating what image convey instead of what it shows
- Persuasion: Incorporating marketing factors and AIM

Highlighted the struggle of T2I models with implicit prompt

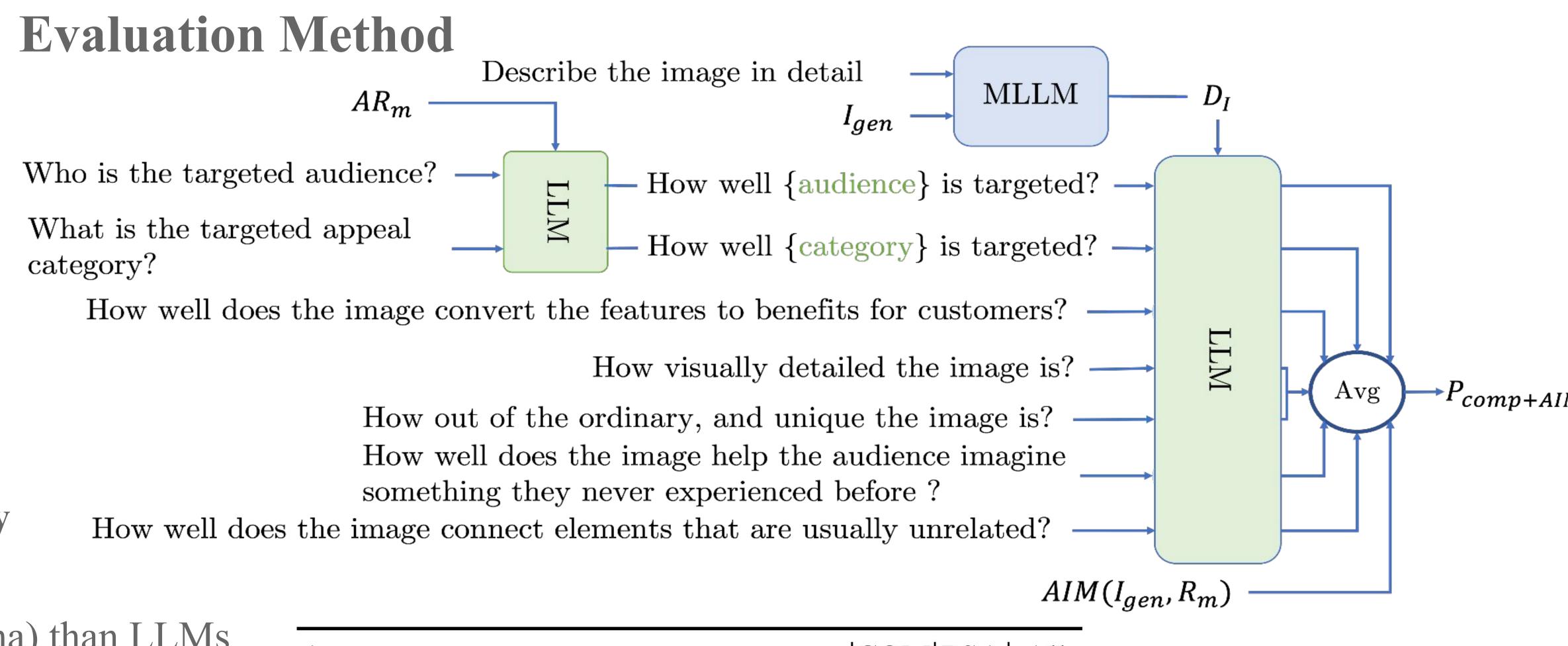
# Persuasion of Images

# Evaluation Overview

- We use effective persuasion factors analyzed in marketing research in computational way
- Reason in AR captures how to persuade the audience (Included in metric)

# Results

- LLMs do not evaluate persuasion accurately (No prior persuasion metric)
- +130% higher agreement (Krippendorff's Alpha) than LLMs COM: 220% **PSA:** 180%
- Higher agreement when adding alignment of image and reason Agreement increased by 23%
- Higher increase in COM advertisements



|COM|PSA| All Annotators 0.27 | 0.26 | 0.27  $H, P_{LLM}$  $H, P_{comp}$  (InternVL, LLAMA3) 0.83 | 0.54 | 0.65 H,  $P_{comp+AIM}$  (InternVL, LLAMA3) | 0.85 | 0.75 | 0.80 H,  $P_{comp+AIM}$  (QwenVL, LLAMA3) | 0.73 | 0.63 | 0.68 H,  $P_{comp+AIM}$  (InternVL, QwenLM) | **0.89** | 0.30 | 0.63 H,  $P_{comp+AIM}$  (QwenVL, QwenLM) | **0.89** | 0.74 | 0.74