# YesBut: A High-Quality Annotated Multimodal Dataset for evaluating Satire Comprehension capability of Vision-Language Models

Abhilash Nandy, Yash Agarwal, Ashish Patwa, Millon Madhur Das, Aman Bansal, Ankit Raj, Pawan Goyal, Niloy Ganguly

IIT Kharagpur        Univ. of Massachusetts        Haldia Institute of Tech.

EMNLP 2024

## Background

- **Satire in Images:** Satirical images blend humor and irony, making them hard to interpret
- **Research Gap:** Few studies focus on full satire detection in images
- **Model Challenges:** VL models struggle with satire due to complex visual-textual cues
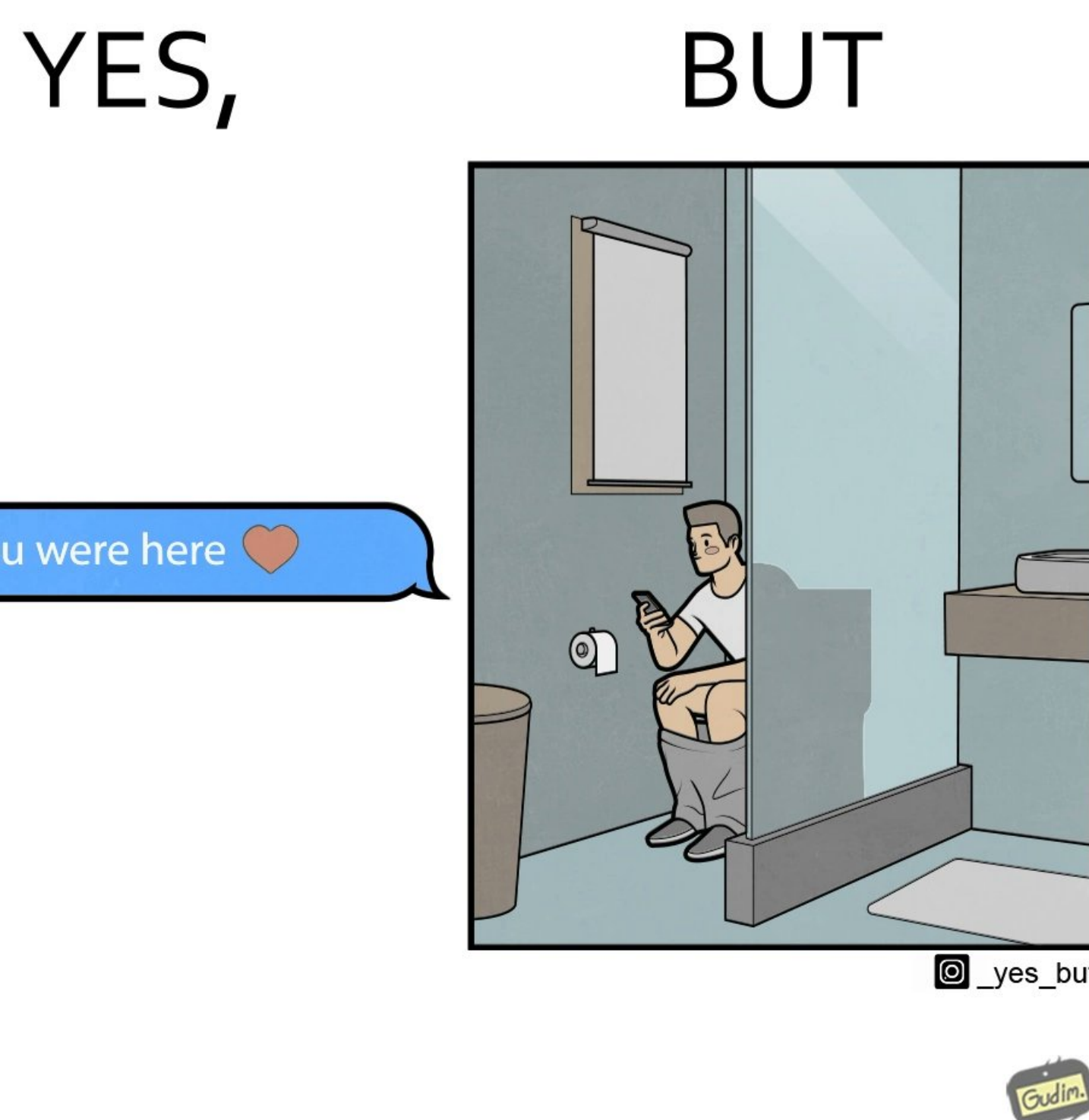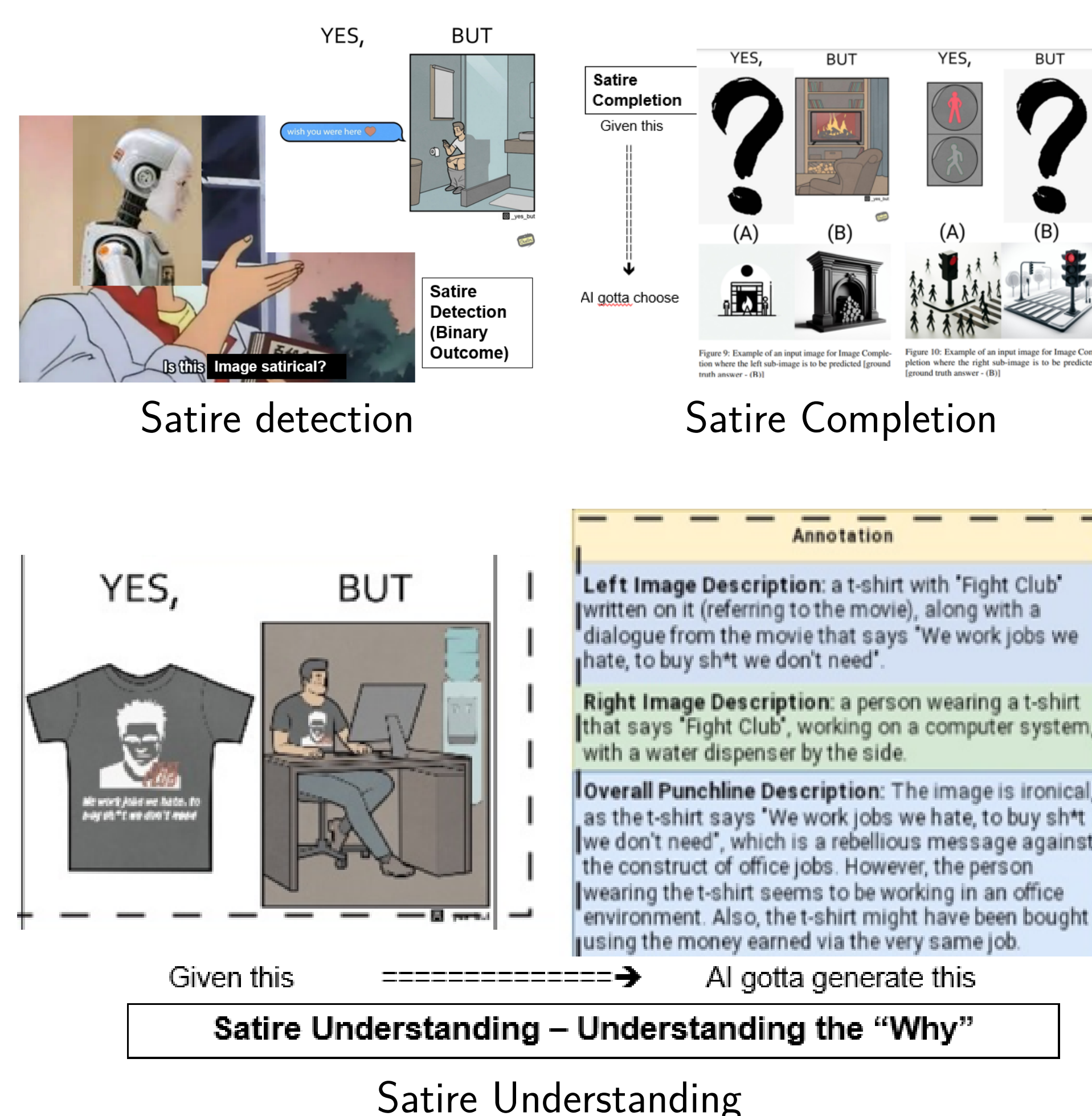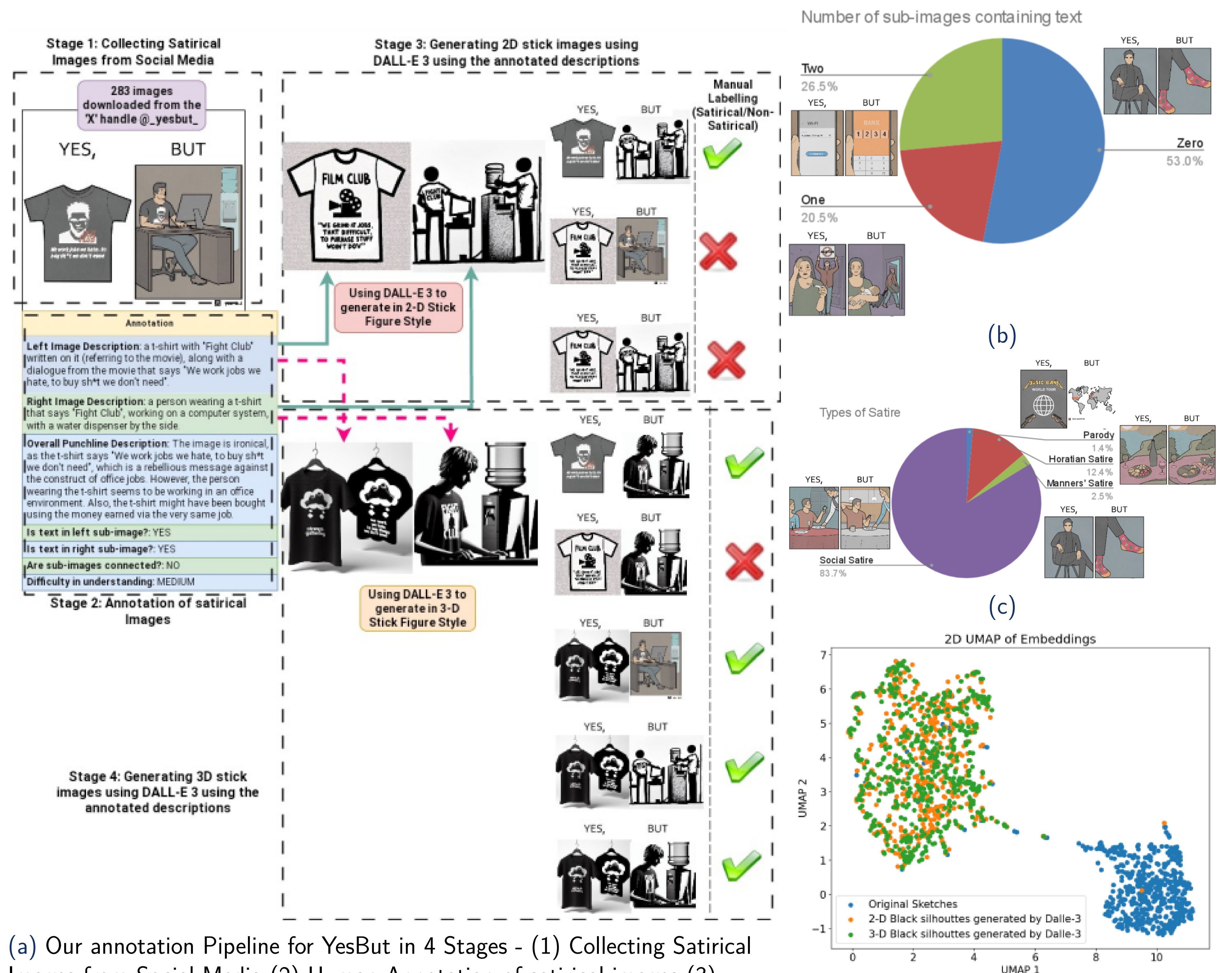


Figure: Satire conveyed through a social media image

## Okay, Yes, But what is YesBut?

- **YesBut Dataset**: A one-of-a-kind multimodal dataset with 2,547 images, containing both satirical and non-satirical images, enriched with diverse artistic styles!
- **Challenges**: We introduced unique tasks like Satirical Image Detection, Satirical Image Understanding, and Satirical Image Completion—each pushing the boundaries of current Vision-Language (VL) models.
- **Benchmarking Results**: Even cutting-edge VL models struggle with our tasks, showing the complexity of understanding irony, humor, and societal satire!

## Tasks Evaluated



Satire detection

Satire Completion

Satire Understanding

Satire Understanding – Understanding the "Why"

## How was YesBut curated?



(a) Our annotation Pipeline for YesBut in 4 Stages - (1) Collecting Satirical Images from Social Media (2) Human Annotation of satirical images (3) Generating 2D stick images using DALL-E 3 and annotated descriptions (4) Generating 3D stick images using DALL-E 3 and annotated descriptions

(d) 2D UMAP Representations of CLIP Image representations of YesBut sub-images

Figure: Left: Our annotation pipeline. Right: Distribution of satirical images based on content and annotated descriptions, and UMAP representations.

## Experiments and Results

|  | TEST ACC. | F1 SCORE |
|---|---|---|
| LLaVA (0-shot) | 53.67 | 48.64 |
| LLaVA (0-shot, CoT) | 52.22 | 46.87 |
| Kosmos-2 (0-shot) | 42.56 | 59.71 |
| Kosmos-2 (0-shot, CoT) | 56.97 | 20.35 |
| MiniGPT4 (0-shot) | 48.29 | 49.33 |
| MiniGPT4 (0-shot, CoT) | 48.88 | 50.61 |
| GPT4 (0-shot) | 55.44 | 55.13 |
| GPT4 (0-shot, CoT) | 48.29 | 42.32 |
| Gemini (0-shot) | 50.82 | 48.29 |
| Gemini (0-shot, CoT) | 46.36 | 38.93 |

Table 3: Evaluation of different VL models on the Satirical Image Detection task

|  | TEST ACC. |
|---|---|
| LLaVA (0-shot) | 51.33 |
| LLaVA (0-shot, CoT) | 56.55 |
| Kosmos-2 (0-shot) | 54.67 |
| Kosmos-2 (0-shot, CoT) | 53.33 |
| MiniGPT4 (0-shot) | 40 |
| MiniGPT4 (0-shot, CoT) | 60.67 |
| GPT4 (0-shot) | 58.67 |
| GPT4 (0-shot, CoT) | 57.33 |
| Gemini (0-shot) | 61.11 |
| Gemini (0-shot, CoT) | 61.81 |

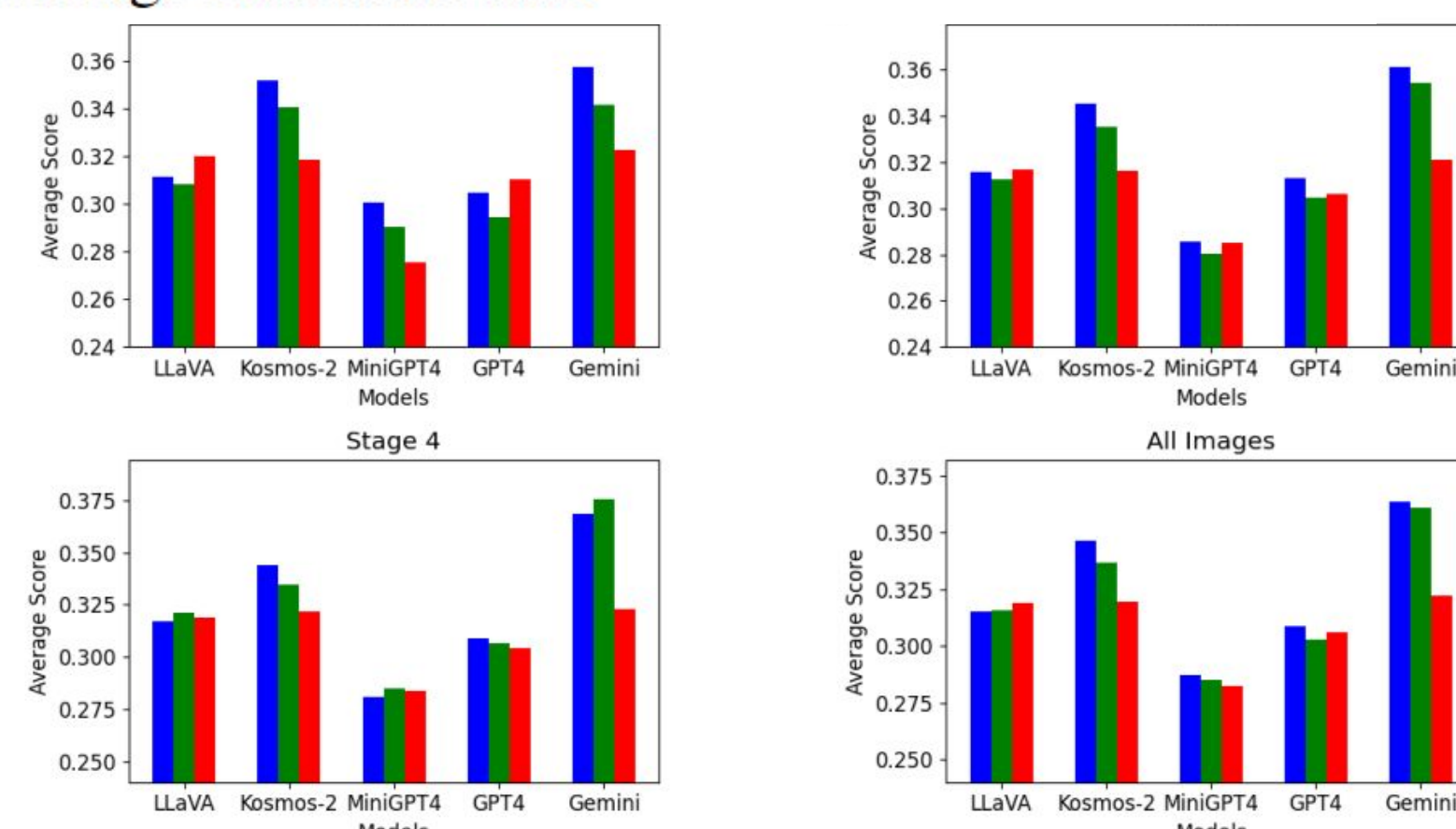Table 4: Evaluation of different VL models on the Satirical Image Completion task



Figure 5: Evaluation of Satirical Image Understanding Capability using multiple VL models at different stages (Stages 2, 3, 4) of annotation of *YesBut*, as well as, for all *YesBut* images
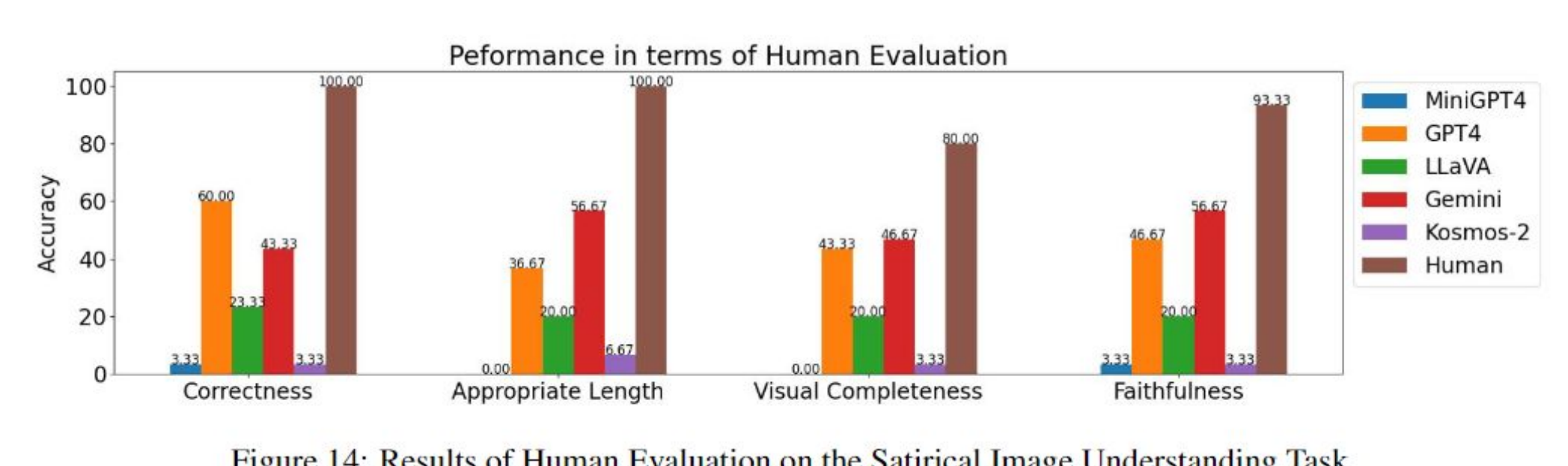
Figure 14: Results of Human Evaluation on the Satirical Image Understanding Task

**Results on Satire Detection, Understanding, and Comprehension on the YesBut Dataset**

Paper Link        YesBut Dataset        Code Base Link

## For Further Information

Arxiv: https://arxiv.org/abs/2409.13592