### Can a Unimodal Language Agent Provide Preferences to Tune a Multimodal Vision-Language Model?

Anonymous ACL submission

### Abstract

This paper proposes a method that enables a language agent to give feedback to a visionlanguage model (VLM) to adapt text generation to the agent's preferences. Using our proposed method, we find that the VLM can supply multimodal scene descriptions to help the LLM better understand multimodal context. Our method leads to improvements of more than 13% in absolute accuracy compared to the baseline multimodal approach. Extensive experiments provide insights on how and why the method works and its limitations.

#### 1 Introduction

006

007

011

012

014

017

024

Large Language Models (LLMs) or Language Agents have emerged as powerful tools for processing and generating textual data. They have significantly advanced natural language processing by achieving near-human-level performance across a wide range of text-centric tasks, including text classification, reasoning, and open-domain question answering, content generation, and others (Guo et al., 2025; Brown et al., 2020; Radford et al., 2018; Dam et al., 2024; Qian et al., 2025; Wang et al., 2025; Jansen et al., 2025; Huang et al., 2025). An active area of research is to make these agents multimodal (Li et al., 2023; Achiam et al., 2023; Liu et al., 2024) to support a wider range of human-AI interactive tasks. However, a key drawback persists: these models typically require training from scratch, or extensive retraining of preexisting LLMs, which is costly and prohibitive to train for many.

To address this issue, we propose a slightly different approach: to describe the multimodal world in text for a language agent to process. By enabling this, we would only need to train lightweight adaptors that describe multimodal features with text. Additionally, this would give us an interpretable 'latent' text modality to help us understand what the language agent's decision was based on – which is broadly impactful for safe AI systems. Once the multimodal information has been described in text, we can take advantage of an agent's reasoning capabilities to incorporate it into decision making. 041

042

043

044

045

047

049

052

053

055

059

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

078

081

However, there are challenges to making this work properly. First it is unclear what the agent needs from different modalities to make the right decision, such as for classification. For example we found that when converting visual modalities to text using a captioning model, and using this as the visual context for an agent to process can lead to an overall decrease in performance (See Tab. 1). It is crucial that multimodal context is incorporated because it can modify the meaning of inputs from a single modality (Baltrušaitis et al., 2018). Another challenge is that it is unclear how we can describe the multimodal in the way the agent expects. Because, such a method would need to be adaptable to many modalities at scale and account for agent preferences.

To address these issues, we propose a novel framework in which we ask the agent to provide preference feedback of VLM outputs. We then optimize the VLM output using DPO (Rafailov et al., 2023). We design several high-level prompts to enable this process. We evaluate our method on the multimodal social reasoning dataset MUStARD (MUltimodal SARcasm Detection) dataset (Castro et al., 2019) which contains properties that can benchmark contextual multimodal understanding (Liang et al., 2023).

With our proposed method, we are able to demonstrate more than 13% in absolute performance over are baseline multimodal approach. We demonstrate that it is possible for agents to give preference feedback and that it can be used to improve the training of a multimodal vision-to-text model. We conduct extensive experiments and find that naively incorporating text descriptions can hurt agent performance. We additionally find that the

- 084
- 08

880

100

101

102

103

104

105

106

107

108

110

111

112

agent's behavior can be quite sensitive to modality descriptions.

### 2 Related Work

### 2.1 Reinforcement Learning from AI Feedback (RLAIF)

Reinforcement Learning from Human Feedback (RLHF) has become a standard technique for aligning large language models (LLMs) with human preferences by training a reward model on humanannotated preference pairs (Christiano et al., 2017; Stiennon et al., 2020; Ouyang et al., 2022). More recently, Bai et al. (2022) and Lee et al. (2023) introduced Reinforcement Learning from AI Feedback (RLAIF), in which off-the-shelf LLMs replace costly human labelers to generate soft preference labels. In RLAIF, a strong LLM is prompted to rank pairs of candidate outputs. Empirical results show that RLAIF matches or exceeds RLHF on summarization, helpfulness and harmlessness benchmarks while dramatically reducing annotation cost and improving scalability (Lee et al., 2023; Gilardi et al., 2023). Unlike RLAIF, which leverages AI feedback to optimize model outputs for purely language-generation tasks, our work uses LLM-driven preferences to shape the inputs to an LLM-based classifier in a multimodal setting. We use the resulting preference pairs, not to finetune the LLM's generation policy, but to refine the vision-language model itself so that it produces better descriptions.

### 2.2 LLM-Based Multimodal Reasoning

113 LLMs have shown remarkable capabilities in understanding and generating text, prompting research 114 into their potential for multimodal reasoning. Yang 115 et al. (2022) explored LLMs for multimodal tasks 116 by converting images into textual captions using 117 models like CLIP Radford et al. (2021), then feed-118 ing these captions into LLMs for tasks such as 119 visual question answering. Similarly, Zhang et al. 120 (2024b) used VLMs to generate textual descrip-121 tions of images, which were then processed by 122 LLMs for reasoning tasks. Kamrul Hasan et al. 123 (2023) propose TextMI, an innovative framework 124 designed to convert acoustic and visual informa-126 tion into textual descriptions, allowing these cues to be effectively processed by text-based models like 127 BERT. These studies focus on static image-based 128 tasks and do not address dynamic video contexts 129 or complex social cues like sarcasm. Our approach 130

extends this line of work by applying LLM-based reasoning to video sarcasm detection, a task requiring temporal and contextual understanding. By using VLM-generated descriptions of video clips, we test whether LLMs can approximate multimodal reasoning from text alone.

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

### 2.3 Preference Optimization in Model Alignment

Direct Preference Optimization (DPO) has emerged as a powerful technique for aligning models with human or AI-generated preferences, bypassing the need for complex reward modeling (Rafailov et al., 2023). In the context of VLMs, Zhou et al. (2024) employed DPO to fine-tune vision-language models by generating preference data using GPT-4V, targeting hallucination reduction. Zhang et al. (2024a) applied DPO to video multimodal models, using language model rewards to optimize instruction-following capabilities. While these works leverage DPO for VLM alignment, they focus on direct multimodal input processing rather than text-based abstractions. Our study innovates by integrating DPO into an LLM feedback loop to enhance the structure of VLM-generated textual descriptions for a specific task. Unlike prior work, which uses DPO to optimize model outputs directly, we employ DPO to refine the input descriptions, enabling the LLM to better interpret multimodal contexts and improve classification accuracy. This feedback-driven approach to structuring textual inputs is a novel contribution to the application of preference optimization.

### 3 Methodology

Our methodology for preference tuning follows a similar pipeline as RLHF (Ouyang et al., 2022) and other preference tuning methods (Rafailov et al., 2023). There are two important differences: first, we do not use any demonstration (i.e. agent does not generate text of what is preferred since we assume it does not have access to the visual modality); and second, we only use a unimodal language agent for preference feedback.

To answer whether a unimodal language agent can provide feedback to a multimodal vision language model (VLM) we setup an experiment with the following methodology (visualized in Fig. 1):

1. Generating Captions with a VLM: We first<br/>ask the VLM (e.g. Llava) to generate multiple178179



Figure 1: An overview of the training pipeline. (a) First, a VLM generates diverse descriptions with a set of prompts; Second, (b) descriptions are ranked using an LLM Agent to generate preference dataset; Lastly, (c) generated preference dataset is used for DPO-based optimization to fine tune the VLM. The agent model stays fixed, and we fine tune the VLM using LoRA.

180

. . .

186

18/

188

189

190

191 192

193

194

195

196

197

198

199

201

206

210

i.e.:

 $\pi_{\theta}: x_v \in \mathbb{R}^{v_w \times v_h \times v_t} \to x_t \in \mathbb{R}^{l \times h},$ 

3.1 Generating Captions with a VLM

The vision language model (VLM) is a model

which maps a short video segment into a caption -

language descriptions of a video.

variations.

all experiments.

2. Generating Agent Preferences: We then ask

3. Preference Optimization: Using the prefer-

ences provided by the language agent, we op-

timize the VLM using preference optimization

(e.g. DPO). The agent model stays fixed for

a language agent (e.g. DeepSeek R1) to rank

generated text descriptions using a few prompt

where the input video segment is a sequence of  $v_t$ images with width and height of  $v_w$  and  $v_h$ . The output caption  $(x_t)$  is of length l with embedding dimension h. The images are sampled uniformly (8 frames per video segment frames).

To increase the diversity of the generated prompts for the video content, we design five distinct prompts. We target different aspects of the video including general content, emotional cues, sarcasm indicators, and facial expressions. For each video, we apply a conversational template combining both the video content and a set of crafted prompts. We run the model across these multiple prompts to obtain a set of five unique description outputs, which are then stored for downstream tasks. For example, we ask the VLM to "Describe what is happening in this video in detail." or "Describe the facial expressions in the video that might indicate contrasting emotions. Keep the description brief". See Appendix A in the appendix for additional prompts that were used. We do not use any sampling strategies (such as varying temperature or top-k) as we noticed that these contributed minimally to differences in generated captions. Instead we only rely on the distinct prompt formulations.

211

212

213

214

215

216

217

218

219

220

221

222

224

225

226

227

228

229

230

232

233

234

235

236

237

238

241

242

243

244

### 3.2 Generating Agent Preferences

We are interested in assessing in 1) whether agents are capable of providing preferences to tune another model, and 2) how capable they are. To do so, we propose three conditions for evaluating the preferences provided by agents. First, agents provide preferences *with* knowledge of what is the groundtruth class of the dataset. Second, agents provide preferences on VLM outputs *without* knowing what is the true class of the dataset. And lastly, we provide in-context examples of preferences without knowledge of ground truth.

# **3.2.1** Agent preferences with ground truth knowledge (With GT)

In the second approach, we incorporate the known ground-truth sarcasm label to guide the LLM's evaluation. For each video, we ask the agent to provide its rating of 5 captions:

You are tasked with evaluating five descriptions of one video to determine how helpful each is in predicting whether the given utterance is {sarcasm}. Below are the previous conversation (if any), the

3

```
original utterance, and the five
    descriptions.
Context: {context}
Utterance: {utterance}
Video Descriptions:
1. {description_1}
2. {description_2}
3. {description_3}
4. {description_4}
5. {description_5}
Please rank each description based on how
    helpful it is for predicting that the
    utterance is {sarcasm}. Assign a score
    between 1 and 10, where:
- 1: Not helpful at all
- 10: Very helpful
Always reply **exactly** in this JSON format:
{
    "score_1":<your integer score here>
    "score_2":<your integer score here>
    "score_3":<your integer score here>
    "score_4":<your integer score here>
    "score_5":<your integer score here>
}
```

247 248

249

251

252

256

259

260

261

263

265

266

268

269

270

271

274

275

279

282

285

289 290

296

297

298

302

305

306

### 3.2.2 Agent preferences without ground truth knowledge (Without GT)

We first aim to evaluate whether an LLM can independently assess the likelihood of sarcasm in a scene using only visual text description and conversational context, without providing any groundtruth labels. For each video, we supply the LLM with the generated video description paired with previous context and original utterance from the scene.

We use the following template to rank captions without ground truth:

```
prompt = f"""You are given a video scene
description, previous conversation
context, and an utterance from that scene.
Your job is to rate how likely the scene
contains SARCASM, on a scale of 1-10,
using this rubric:
1-5: NOT sarcastic
```

```
- 1: Strongly believe this is NOT sarcastic
- 5: Not sure, but leaning towards NOT sarcastic
6-10: Sarcastic
- 6: Not sure, but leaning towards sarcastic
- 10: Strongly believe this IS sarcastic
Video description: {video_description}
Context: {context}
Target Utterance: {utterance}
Please provide your response in this exact format:
FINAL SCORE: [number between 1-10]"""
```

In this method, the LLM is asked to predict the likelihood of sarcasm on a scale from 1 to 10. Formally, given input q = (d, u, c), where d is the

1:	Sort <i>scores</i> list in ascending order by score
2:	if ground_truth_sarcasm is True then
3:	preferred $\leftarrow$ description with highest score (last in <i>scores</i> list)
4:	dispreferred $\leftarrow$ description with lowest score (first in <i>scores</i> list)
5:	else
6:	preferred $\leftarrow$ description with lowest score (first in <i>scores</i> list)
7:	dispreferred $\leftarrow$ description with highest score (last in <i>scores</i> list)
8:	end if

description, u is the utterance, and c is the previous context of the given video. The LLM outputs a scalar score  $s \in [1, 10]$  estimating the sarcastic nature of the scene, where 1 means that the LLM strongly believes u is not sarcastic and 10 means strongly believes u is sarcastic.

307

308

309

310

311

312

313

314

315

316

317

319

321

322

323

324

325

327

328

329

331

332

333

334

335

336

337

338

339

341

343

344

### 3.2.3 Agent preferences with in-context examples (With ICL)

In this approach, we adopt an in-context learning (ICL) approach by providing the language model with two illustrative examples to demonstrate the ranking process for video descriptions. The first example describes a scene containing sarcasm, which contains a score of 10 to indicate strong sarcastic content. The second example describes a scene without any sarcasm, with a score of 1 to represent no sarcasm. Then, we prompt the LLM to predict the likelihood of sarcasm as described in Section 3.2.2.

For our in-context examples, we selected two of the model's raw outputs from Section 3.2.2. One is of score 10 for a sarcastic scene and another is scored 1 for a non-sarcastic scene. The full prompt can be found in Appendix B.4

### 3.2.4 Sorting preferences

We provide a simple rule to determine the order of the preferences (found in Algorithm 1). These rules are applied to the Without GT and With ICL conditions. The reason we do this is to evaluate how effective agents are at incorporating the ground truth. The ranking provides a (weak) set of rules that restricts the preferences to align with the ground truth of the dataset manually.

### 3.3 Preference Optimization

Direct Preference Optimization (DPO) (Rafailov et al., 2023) is used to align VLM outputs to the preferences of the language agent. To summarize, DPO leverages a reparameterization of the reward



Figure 2: Overview of the evaluation process. In the unimodal evaluation pipeline shown in (a), only the conversation context and the final statement is given to the Agent to determine whether the final statement is sarcastic or not. (b) demonstrates how multimodal information is incorporated during the evaluation process. The VLM acts like a feature extractor by converting visual features to text descriptions. This is fed to the LLM. (c) shows how the preference tuned VLM (Optimized Vision Language Model) is incorporated into the pipeline.

function under the Bradley-Terry model, showing that the optimal policy satisfying human preferences can be obtained by maximizing a simple binary cross-entropy objective:

345

347

352

357

372

$$\mathcal{L}_{\text{DPO}}(\pi_{\theta}; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \Big[ \log \sigma \Big( \beta \log \frac{\pi_{\theta}(y_w | x) \cdot \pi_{\text{ref}}(y_l | x)}{\pi_{\text{ref}}(y_w | x) \cdot \pi_{\theta}(y_l | x)} \Big) \Big],$$
(1)

where  $y_w$  and  $y_l$  denote the preferred and dispreferred responses, respectively,  $\pi_{ref}$  is the reference policy (often the supervised fine-tuned model), and  $\beta$  controls the KL-divergence regularization strength. This formulation allows DPO to directly adjust the model's log-probabilities to satisfy human preferences without maintaining a learned reward critic.

In our case,  $\pi_{\theta}$  is the VLM we want to fine-tune and  $\pi_{ref}$  is a reference model, or an unmodified version of  $\pi_{\theta}$ . D is the dataset of preferences generated in the section 3.2. x is a sample prompt which is "Describe the video is details".  $y_w$  and  $y_l$ are the LLM's preferred and disprefered response to the prompt x. Using  $\beta$  we control the amount of divergence from the reference model  $\pi_{ref}$ .

#### Agent Evaluation Process 3.4

We wish to see if the agent is properly accounting for the visual context to make a decision about a multimodal property. Therefore, both the textual description of the visual context as well as the original text-based caption should be available to

the agent for it to make its decision. We design a prompt to do so and to provide a decision as following-

373

374

375

376 377

384

385

391

393

394

396

397

400

401

402

403

404

405

406

407

408

409

410

```
prompt = f"""
Based on the following previous conversation in
    the video, decide whether the final
    utterance is sarcastic or not.
Context: {context}
Utterance:
   {speaker}: "{utterance}"
Your job is to rate how likely the scene
    contains SARCASM, on a scale of 1-10, using
                                                           386
    this rubric:
1-5: NOT sarcastic
                                                           390
- 1: Strongly believe this is NOT sarcastic
- 5: Not sure, but leaning towards NOT sarcastic
                                                           392
6-10: Sarcastic
- 6: Not sure, but leaning towards sarcastic
- 10: Strongly believe this IS sarcastic
                                                           395
Always reply **exactly** in this JSON format:
{
    "score":<your integer score here>
}
"""
                                                           399
```

Depending on the evaluation condition, the video\_description comes from the output of the VLM. The utterance is provided using the text modality. Additionally, any conversation context is also provided. In this way, we can evaluate how well the agent can integrate the visual descriptive context with the utterance. We present the conditions of evaluation in Fig. 2. Once the LLM provides its output, we assess our model's performance using standard classification metrics (pre-

494

495

496

497

498

499

500

501

502

503

504

505

506

458

#### 411 412

413

414

436

437

438

439

440

441

449

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

cision, recall, accuracy) following the tasks of the dataset.

### 4 Experimental Setup

### 4.1 Dataset

For our experiments, we use the MUltimodal SAR-415 casm Detection (MUStARD) dataset, a multimodal 416 sarcasm detection (binary classification) corpus in-417 troduced by Castro et al. (2019). MUStARD con-418 tains 690 short video clips drawn from popular 419 TV shows such as Friends, The Big Bang The-420 ory, The Golden Girls, and Sarcasmaholics Anony-421 mous. The dataset is carefully balanced, with an 499 equal number of sarcastic and non-sarcastic sam-423 ples. Each video is annotated with high-quality 424 sarcasm labels derived through a two-step man-425 ual annotation process, achieving improved inter-426 annotator agreement. Each utterance includes three 427 aligned modalities: text, which is manually or semi-428 automatically transcribed dialogue; audio which is 429 the raw speech signals containing prosodic cues; 430 and video which contains frame-level visual fea-431 tures representing facial expressions and gestures. 432 For our experiments, we split the data into a train-433 ing dataset of 551 videos and a testing dataset of 434 139 videos. 435

### 4.2 VLM model

In our experiment, we fine tune the LLaVA-Video model, a large-scale video-language model developed to advance instruction-following tasks in multimodal settings (Zhang et al., 2024c). The model architecture uses a specialized video representation technique called SlowFast, allowing the system to process up to three times more frames than standard methods within GPU memory constraints. To generate video descriptions during evaluation we use the following prompt for all methods:

Describe the speaker's nonverbal cues, the context, and any mismatches between them.

### 4.3 Agent Model

We used the 7B parameter Deepseek-R1 model as the language agent. The agent model stays fixed through all parts of the training and evaluation.

### 4.4 Training Details:

We conducted our experiments on a system equipped with six NVIDIA GeForce RTX 4090 GPUs (24GB VRAM each) using CUDA 12.6. We used CUDA for GPU acceleration with fallback to CPU. To fine-tune our model, hyperparameters that impacted computational constraints were chosen to balance algorithmic performance and computational feasibility.

The learning rate was set to  $1 \times 10^{-5}$  to ensure stable convergence. Due to memory constraints, a batch size of 1 was used during training. The model was trained for a total of 5 epochs to achieve sufficient adaptation to the dataset while avoiding overfitting. We utilized Direct Preference Optimization (DPO) with a temperature parameter, ( $\beta$ ) of 0.1 to control the divergence from the reference model. For training the VLM, we used the generic prompt "Describe the video in detail". We fine-tuned the VLM using Low-Rank Adaptation (LoRA) with a rank of r = 4, a scaling factor of  $\alpha = 16$ , and a dropout probability of 0.1. LoRA was applied to the query, key, and value projection layers (q\_proj, k\_proj, and v\_proj).

To ensure consistency and mitigate randomness in the evaluation process, we carefully tuned the inference parameters. The temperature was set to 0.00, ensuring more deterministic outputs. This configuration was chosen to maintain reliability in evaluation and ensure meaningful comparisons across different captions.

### 5 Result Analysis

### 5.1 Main Results

We evaluated our approach using three distinct methods. First, we feed the utterance to the LLM without any multimodal information. Then we combined the description generated by the base LLaVA-NeXT-Video model with the utterance from the dataset, and then sent the combined descriptionutterance pair to the DeepSeek evaluator for sarcasm analysis. Finally, we performed the evaluation procedure using descriptions generated by our DPO fine-tuned LLaVA-NeXT-Video model (Figure 2). Tab. 1 presents the sarcasm detection accuracy results for our DPO-fine-tuned model compared to the base LLaVA-NeXT-Video model on the MUStARD dataset.

# **5.2 Impact of the Temperature Parameter (**β) on DPO Fine-Tuning

Here we highlight the effect of different  $\beta$  values for our DPO fine tuning approach. Higher values of  $\beta$  heavily penalize any deviation from the reference policy, whereas lower values encourage the model to diverge more to satisfy the reward signal.

Model	Accuracy	Precision	Recall	
Baseline				
Utterance Only	61.2	62.0	64.0	
Visual Only	48.9	52.0	19.0	
Multimodal	52.9	58.0	40.0	
Preference Tuned				
Visual Only				
With GT	47.48	49.0	40.0	
Without GT	50.36	71.0	7.0	
With ICL	51.80	65.0	15.0	
Multimodal				
With GT	58.7	58.0	75.0	
Without GT	66.9	70.0	65.0	
With ICL	65.1	76.0	50.0	

Table 1: Performance comparison of baseline (pretrained VLM) compared to a few preference tuned conditions. While the visual-only modality did not outperform an utterance-only prediction, it improves over the descriptions obtained from an unoptimized VLM. Additionally, we see that there is better contextual information that the agent can make use of which led to the improvements in multimodal performance. **Note:** The preference-tuned models have to do with how the training dataset was obtained (see methods). GT information was not available at test time, nor during prompting or training of the VLM.

( $\beta$ ) Value	Accuracy %	Precision	Recall
0.1	60.45	63.0	55.0
0.5	57.14	59.0	55.0
0.8	53.33	54.0	47.0

Table 2: Effect of  $(\beta)$  on DPO Fine-Tuning. At  $\beta = 0.1$ , we achieve best results. As we increase the value of  $\beta$  we observe that accuracy precision, and recall becomes lower.

To measure the impact of the DPO fine-tuning parameter, ( $\beta$ ) with the preference dataset- "Without GT", we fixed the training prompt and inference prompt to "Describe the video in detail.". Then we ran the three experiments with  $\beta$  set to 0.2, 0.5 and 0.8. As  $\beta$  increased, the model's accuracy steadily declined from 60.45% to 53.33% (Tab. 2). This indicates that stronger regularization toward the base policy yields more conservative but less adaptive behavior.

### 5.3 Impact of Prompt Variations on performance

507

508

509

510

511

512

513

514

515

516

517

518

519

520

522

523

To assess how different prompting strategies impact DPO fine-tuning for a specific task, we experiment with three of training prompt variants: 1) *Generic prompt*; 2)*Task-specific prompt*; 3) *No prompt*. We use  $\beta$  of 0.1 for these experiments.

### **5.3.1 Effect of Prompt in DPO Fine-tuning:**

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

First, we assess how different training prompts affect a "Without GT" preference-tuned model by comparing three configurations: 1) Training with a generic prompt ("Describe the video in detail."). 2) Training with a task-specific prompt ("Describe the speaker's nonverbal cues, the context, and any mismatches between them."). Lastly, 3) Training with no prompt.

We used the same prompt for inference as used in the training process. We observe that training with the generic prompt yields the highest accuracy (60.45%) (Tab. 4).

#### 5.3.2 Effect of Prompt During Inference:

Here we evaluate inference-time prompt effects on different preference tuned models while keeping the training prompt the same (*"Describe the video in detail"*). Notably, when the model is trained with the generic prompt, switching to the task-specific inference prompt significantly boosts performance(66.9%) significantly. This demonstrate that a general training instruction combined with a focused, task-specific inference prompt can yield superior results (Tab. 5).

### 5.4 Impact of Different Modalities

In this section we highlight how different modalities affect task performance. We compare the agent's performance across utterance-only, visualonly and multimodal inputs. In multimodal input we include both the utterance and description of the scene. The utterance-only model reaches 61.2% accuracy, while we see performance degradation with visual-only information (Tab. 1). Then we investigate how descriptions generated by three different preference-tuned models ("Without GT," "With GT," and "With ICL") help to alter the performance of the agent. With our training process using "Without GT" preference tuned model we get the highest accuracy (66.9%)in the multimodal setting.

#### 5.5 What changed in the VLM outputs?

In this section, we explore how the VLM outputs changed as a result of the preference optimization which led to the improvements in agent prediction performance. Since the dataset being examined is a multimodal social reasoning dataset, we examine psychometric properties (i.e. emotional properties) extracted by the VLM. We present some summary statistics in Tab. 3. We found that for the without

	Overall			Agent pos pred			Agent neg pred		
VLM Model	Acc %	Length	Tone	Acc %	Length	Tone	Acc %	Length	Tone
Baseline	52.9	162.1	47.1	58.0	148.0	41.4	50.0	170.3	50.3
Preference Tuned	577	147.0	44.5	57.0	140.7	<i>41.6</i>	575	145 7	50.0
With GT Without GT With ICL	66.9 65.1	147.8 170.8 120.2	44.5 54.4 38.0	57.8 70.1 75.6	148.7 171.2 132.2	41.6 52.8 43.0	63.8 59.5	145.7 170.4 113.8	56.0 35.2

Table 3: Summary statistics comparing text output from the VLM model for the MUStARD. Three conditions are compared: 1) overall statistics, 2) statistics of VLM outputs which led the agents to predict a positive label, and 3) summary of what caused the agent to predict a negative label. Tone scores are obtained from LIWC-22. After the VLM was preference tuned, the tone of the VLM also changed significantly.

Training Prompt	Acc %	Prec.	Recall
"Describe the speaker's nonverbal cues, the context, and any mismatches between them."	58.33	60.0	63.0
"Describe the video in detail	60.45	63.0	55.0
No Prompt	47.79	50.0	49.0

Table 4: Effect of Different Prompts on DPO Fine-Tuning. We choose the "With GT" preference tuned model to test the effect of different prompts in the DPO fine-tuning process. We notice that, with a generic prompt in the training process we can acheive higher accuracy.

GT and with ICL conditions, there was a large difference in tone for the tuned VLM outputs. We also found significant variability in the lengths of the generated text with ICL being the shortest. The VLM trained on preferences of the LLM when given the ground truth tended to have a shorter token length. In general the agent was more accurate with positive instances of sarcasm than negative ones.

573

574

577

582

583

584

585

586

587

589

593

594

597

We also saw that sometimes the VLM would hallucinate. It is not surprising since the language agent has no way to judge whether the explanation is faithful to the original video or not – which also aligns with data collection processes of RLHF. For example, in Tab. 6 found in the Appendix, we see that the VLM is describing a voiceover that does not exist. As a rough measure of hallucination, we count the occurrence of the word "voice" and found that the baseline model generates captions with 0 occurrence of "voice". However, the without GT, with GT, and with ICL models generate 5, 34, and 1 occurrences of this word. Despite the hallucinations, it is interesting that this resulted in a large performance improvement even when the visual descriptions were factually inaccurate. We will ex-

Model	Accuracy	Precision	Recall
Inference Prompt: P1			
Base	52.9	53.0	53.0
With GT	57.7	59.0	65.0
Without GT	60.5	63.0	55.0
With ICL	60.0	64.0	51.0
Inference Prompt:P2			
Base	51.9	58.0	40.0
With GT	57.7	58.0	75.0
Without GT	66.9	70.0	65.0
With ICL	65.1	76.0	50.0

Table 5: Effect of different inference prompts in VLM. Here, prompt, P1="Describe the video in detail"; P2="Describe the speaker's nonverbal cues, the context, and any mismatches between them." We observe that training with a generic prompt, we can achieve better performance with a task specific prompt during the inference time

plore how to mitigate these sorts of hallucinations in the future.

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

### 6 Conclusion

This work investigated the potential of off-the-shelf unimodal LLMs for multimodal inference by transforming non-textual inputs into rich, descriptive text. Using sarcasm detection as a case study, we found that directly combining these textual descriptions with utterances initially caused  $\sim 8\%$  drop in accuracy compared to the utterance-only baseline—highlighting the challenge of approximating multimodal reasoning through text alone.

However, optimizing the description generator to align with the LLM's preference, we not only fully recovered the initial loss but also surpassed the original utterance-only baseline by an additional  $\sim 5\%$  accuracy. This demonstrates that LLMs can effectively generate the preference signals needed to produce better textual descriptions for realizing multimodal inference tasks.

641

647

654

664

### 7 Limitation

This study demonstrates the potential of Large Lan-619 guage Models (LLMs) to interpret multimodal con-621 texts using text-only descriptions for sarcasm detection, achieving promising results in sarcasm detection. Although our experiments demonstrate that offline reinforcement-learning (RL) fine-tuning can 624 significantly improve policy performance on static 625 datasets, this design lacks of real-time adaptation. In the current framework, DPO processes preference data in batch mode, optimizing the model based on a static dataset. This offline approach restricts the model's ability to adapt dynamically to new or evolving data patterns, which is critical for real-world applications. To address this limitation, a primary focus can be on transitioning from 633 offline to online RL fine-tuning. Online RL would enable the model to learn and adapt in real-time as new data becomes available, enhancing its responsiveness and accuracy in dynamic contexts. Another aspect is the hallucinations we noticed in 638 the VLM outputs. We will also explore why this is and develop ways to mitigate this.

### 8 Ethics Statement

This study uses publicly available datasets (MUStARD) and pre-trained models, without collecting or processing personal or sensitive data. The proxies for video data are generated from dataset content following appropriate academic use guidelines. No human subjects were involved, and all models comply with their respective licenses. This research is intended for academic exploration and is not suitable for direct deployment in sensitive or decision-critical applications without further evaluation.

### References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, and 1 others. 2022. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*.
- Tadas Baltrušaitis, Chaitanya Ahuja, and Louis-Philippe Morency. 2018. Multimodal machine learning: A

survey and taxonomy. *IEEE transactions on pattern analysis and machine intelligence*, 41(2):423–443.

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

711

712

713

714

715

716

717

718

719

- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, and 1 others. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Santiago Castro, Devamanyu Hazarika, Verónica Pérez-Rosas, Roger Zimmermann, Rada Mihalcea, and Soujanya Poria. 2019. Towards multimodal sarcasm detection (an \_obviously\_ perfect paper). *arXiv preprint arXiv:1906.01815*.
- Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30.
- Sumit Kumar Dam, Choong Seon Hong, Yu Qiao, and Chaoning Zhang. 2024. A complete survey on llmbased ai chatbots. *arXiv preprint arXiv:2406.16937*.
- Fabrizio Gilardi, Meysam Alizadeh, and Maël Kubli. 2023. Chatgpt outperforms crowd workers for text-annotation tasks. *Proceedings of the National Academy of Sciences*, 120(30):e2305016120.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shi-rong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in Ilms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Zeyi Huang, Yuyang Ji, Xiaofang Wang, Nikhil Mehta, Tong Xiao, Donghyun Lee, Sigmund Vanvalkenburgh, Shengxin Zha, Bolin Lai, Licheng Yu, and 1 others. 2025. Building a mind palace: Structuring environment-grounded semantic graphs for effective long video analysis with llms. *arXiv preprint arXiv:2501.04336*.
- Jacqueline A Jansen, Artür Manukyan, Nour Al Khoury, and Altuna Akalin. 2025. Leveraging large language models for data analysis automation. *PloS one*, 20(2):e0317084.
- Md Kamrul Hasan, Md Saiful Islam, Sangwu Lee, Wasifur Rahman, Iftekhar Naim, Mohammed Ibrahim Khan, and Ehsan Hoque. 2023. Textmi: Textualize multimodal information for integrating non-verbal cues in pre-trained language models. *arXiv e-prints*, pages arXiv–2303.
- Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas Mesnard, Johan Ferret, Kellie Lu, Colton Bishop, Ethan Hall, Victor Carbune, Abhinav Rastogi, and 1 others. 2023. Rlaif vs. rlhf: Scaling reinforcement learning from human feedback with ai feedback. *arXiv preprint arXiv:2309.00267*.

Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. Blip-2: Bootstrapping language-image pretraining with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR.

720

721

725

726

727

737

740

741

742

743

745

746

747

748

750

751

752

754

755 756

758

759

761

764

765

767

770

771

773

774

775

- Paul Pu Liang, Yun Cheng, Xiang Fan, Chun Kai Ling, Suzanne Nie, Richard Chen, Zihao Deng, Nicholas Allen, Randy Auerbach, Faisal Mahmood, and 1 others. 2023. Quantifying & modeling multimodal interactions: An information decomposition framework. Advances in Neural Information Processing Systems, 36:27351–27393.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26296–26306.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, and 1 others. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.
- Dongjun Qian, Kai Su, Yiming Tan, Qishuai Diao, Xian Wu, Chang Liu, Bingyue Peng, and Zehuan Yuan. 2025. Vc-llm: Automated advertisement video creation from raw footage using multi-modal llms. *arXiv preprint arXiv:2504.05673*.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, and 1 others. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PmLR.
- Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, and 1 others. 2018. Improving language understanding by generative pre-training.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36:53728– 53741.
- Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. 2020. Learning to summarize with human feedback. *Advances in neural information processing systems*, 33:3008– 3021.
- Leijie Wang, Kathryn Yurechko, Pranati Dani, Quan Ze Chen, and Amy X Zhang. 2025. End user authoring of personalized content classifiers: Comparing example labeling, rule writing, and llm prompting. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*, pages 1–21.

Zhengyuan Yang, Zhe Gan, Jianfeng Wang, Xiaowei Hu, Yumao Lu, Zicheng Liu, and Lijuan Wang. 2022. An empirical study of gpt-3 for few-shot knowledgebased vqa. In *Proceedings of the AAAI conference on artificial intelligence*, volume 36, pages 3081–3089.

776

781

782

783

784

785

786

787

788

789

790

792

793

794

795

796

797

798

799

800

801

802

803

804

805

807

808

809 810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

- Ruohong Zhang, Liangke Gui, Zhiqing Sun, Yihao Feng, Keyang Xu, Yuanhan Zhang, Di Fu, Chunyuan Li, Alexander Hauptmann, Yonatan Bisk, and 1 others. 2024a. Direct preference optimization of video large multimodal models from language model reward. *arXiv preprint arXiv:2404.01258*.
- Wenqiao Zhang, Tianwei Lin, Jiang Liu, Fangxun Shu, Haoyuan Li, Lei Zhang, He Wanggui, Hao Zhou, Zheqi Lv, Hao Jiang, and 1 others. 2024b. Hyperllava: Dynamic visual and language expert tuning for multimodal large language models. *arXiv preprint arXiv:2403.13447*.
- Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. 2024c. Video instruction tuning with synthetic data. *arXiv preprint arXiv:2410.02713*.
- Yiyang Zhou, Chenhang Cui, Rafael Rafailov, Chelsea Finn, and Huaxiu Yao. 2024. Aligning modalities in vision large language models via preference finetuning. *arXiv preprint arXiv:2402.11411*.

### A Prompts to VLM for generating video descriptions

To ensure diverse description generation using the LLaVA-NeXT-Video model, we define five distinct prompts, each crafted to elicit different aspects of the video content. These prompts are:

# Define 5 diverse prompts diverse\_prompts = [ "Describe what is happening in this video in detail." "Describe the video in such a way that it will be helpful for sarcasm detection. Try to keep the description brief.", "Describe the facial expressions and the speech tone of the speakers that is useful to understand the subtle meaning of the conversation. Try to keep the description brief.", "Describe the facial expressions in the video that might indicate contrasting emotions. Keep the description brief" "Provide a brief description of this video." ]

For each video, the model is run once per prompt, generating five separate outputs. No sampling strategies (such as varying temperature or top-k) are used; instead, diversity arises entirely from the distinct prompt formulations.

887

### B Prompts for the Language Agent (Deepseek-R1)

### **B.1** Prompt to Rank the Captions without Ground Truth

prompt = f"""Consider that the visual description of a video scene is: {video\_description}.{raw\_utterance\_text}

Please rate the likelihood that the scene contains SARCASM on a scale of 1-10.

1-5: NOT sarcastic

- 1: Strongly believe this is NOT sarcastic
- 5: Not sure, but leaning towards NOT sarcastic

6-10: Sarcastic6: Not sure, but leaning towards sarcastic

- 10: Strongly believe this IS sarcastic

Please provide your response in this exact format:

FINAL SCORE: [number between 1-10]"""

### **B.2** Prompt to Rank the Captions With Ground Truth

prompt = f"""You are tasked with evaluating five video descriptions to determine how helpful each is in predicting whether an utterance is {sarcasm\_status}. Below are the five video descriptions and the original utterance.

{raw\_utterance\_text}Video Descriptions:

- 1. {descriptions[0]}
- 2. {descriptions[1]}
- 3. {descriptions[2]}
- 4. {descriptions[3]}
- 5. {descriptions[4]}
- Please rank each description based on how helpful it is for predicting that the utterance is {sarcasm\_status}. Assign a score between 1 and 10, where: - 1: Not helpful at all
- 10: Very helpful

Provide your response in this exact format: Description 1: [score] Description 2: [score] Description 3: [score] Description 4: [score] Description 5: [score]

### B.3 Prompt for Ranking Captions with Few-Shot Examples

prompt = f"""You are given a video scene description, previous conversation context, and an utterance from that scene. Your job is to rate how likely the scene contains SARCASM, on a scale of 1 to 10, using this rubric: 1 to 5: NOT sarcastic - 1: Strongly believe this is NOT sarcastic - 5: Not sure, but leaning towards NOT sarcastic 6 to 10: Sarcastic - 6: Not sure, but leaning towards sarcastic - 10: Strongly believe this IS sarcastic Always reply \*\*exactly\*\* in this format: FINAL SCORE: [number between 1 to 10] \*\*Example 1\*\* Video description: In the video, we see a man and a woman in a hospital setting. The man, dressed in a white coat and tie, is holding a clipboard and appears to be a doctor. He is speaking to the woman, who is wearing a pink shirt, and seems to be explaining something to her. The woman is smiling and nodding along, indicating that she is engaged in the conversation. However, the man's tone is sarcastic and dismissive, as he talks about the woman's medical condition with a lack of concern and even makes a joke about it. The woman seems to be taking the situation lightly, laughing along with the man's jokes, which suggests that she is either in on the joke or is not bothered by his tone. The setting is a typical hospital environment, with medical equipment visible in the background, and the overall atmosphere is light-hearted and humorous.

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911 912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

Previous conversation:

```
RACHEL: "All right, I'm outta here!"
MONICA: "I'm kidding! I'm kidding!"
```

```
RACHEL: "So were done then!"
```

Utterance: PERSON: "Almost! But first, we gotta start." FINAL SCORE: 10

\*\*Example 2\*\*

The video features two men sitting in a room, one of whom is holding a coffee cup and speaking to the other. The man holding the coffee cup is wearing a red and blue striped shirt, while the other man is wearing a white shirt. They are engaged in a conversation, and the man in the red and blue shirt is holding the coffee cup in his right hand. The room appears to be an office setting, with a desk and a chair visible in the background.

Previous conversation: HOWARD: "Do you really think you should be drinking right now?"

- RAJ: "How else am I supposed to talk to the Human Resources lady?"
- HOWARD: "I don't know. Seek professional help?"
  Utterance: RAJ: "I did. The guy at the liquor
   store said this stuff tastes great in
   coffee."
  FINAL SCORE: 1
  ---

\*\*Now you try\*\*
Video description: {desc}

{context\_text}{raw\_utterance\_text}

960 961	FINAL SCORE: [number between 1-10]
962	<b>B.4</b> Prompt to Calculate Accuracy
963 964 965 966 967 968	<pre>prompt = f""" You are given a video scene description, previous conversation context, and an utterance from that scene. Your job is to rate how likely the scene contains SARCASM, on a scale of 1-10, using this rubric:</pre>
970 971 972 973	<ul><li>1-5: NOT sarcastic</li><li>- 1: Strongly believe this is NOT sarcastic</li><li>- 5: Not sure, but leaning towards NOT sarcastic</li></ul>
974 975 976 977	6-10: Sarcastic - 6: Not sure, but leaning towards sarcastic - 10: Strongly believe this IS sarcastic
978 979 980 981	Video description: {description} Context: {context} Target Utterance: {utterance} """

### 982 B.5 Generated Output Examples

### Agent provided preferences with knowledge of GT. After training, sometimes the VLM will hallucinate. For example, we see that the VLM is talking about a voiceover that doesn't exist.

In the video, a man and a woman are seen lying on the ground, surrounded by fallen leaves, which suggests it's autumn. The man is wearing a blue sweater and the woman is in a blue top, both with their arms around each other. The man is speaking directly to the camera, and the woman is looking at it with a smile on her face. The man's tone is serious, while the woman's is more lighthearted. The man's voiceover explains that they're happy and in love, and the woman's voiceover is more playful, contrasting with his. The man's voiceover is in a serious tone, while the woman's is more lighthearted. The scene is set in a park with autumn leaves, and the man and woman are the only two people in it. The man's voiceover explains that they're happy and in love, while the woman's voiceover explains that they're happy and in love, while the woman's voiceover explains that they're happy and in love, while the woman's voiceover explains that they're happy and in love, while the woman's voiceover explains that they're happy and in love, while the woman's voiceover explains that they're happy and in love, while the woman's voiceover explains that they're happy and in love, while the woman's voiceover contrasts with his, adding a sense of humor to the scene.

#### In the following three examples, we see a different sentiment being expressed.

**Example 1, baseline model:** In the video, we see a man sitting on a couch, holding a phone to his ear. He appears to be engaged in a conversation, possibly a phone call or a video call, as he is looking at the camera with a focused expression. The room around him is casually furnished, with a couch, a chair, and some posters on the wall. The man's attire is a sweater with a pattern, and he seems to be in a relaxed setting, possibly at home. The lighting is warm, and the overall atmosphere is calm. The video captures a moment of everyday life, with the man's interaction with the phone being the main focus.

**Example 2, preference given with knowledge of GT:** In the video, we see a man sitting on a couch and talking on a phone. He's wearing a sweater with a distinctive checkered pattern and appears to be in a casual setting, possibly his living room. His expression is serious, and he's engaged in a heated conversation. The camera zooms in on his face, emphasizing his intensity. The contrast between his calm exterior and the inner turmoil is evident.

**Example 3, preference given with in-context examples:** In the video, we see a young man sitting comfortably on a couch in a living room. He is dressed in a cozy sweater and is engaged in a phone conversation. His face is lit up with a smile, suggesting a pleasant conversation. The room around him is warmly decorated, with a poster of a man on the wall and a bookshelf filled with various items. The man on the phone is holding a phone to his ear, indicating an ongoing call. The living room is well-lit, with a lamp casting a soft glow on the scene. The man's relaxed posture and the casual setting suggest a comfortable and familiar environment. The video captures a moment of everyday life, showcasing the man's engagement in a phone call while enjoying his surroundings.

Table 6: Examples of VLM generated sequences.