# ADAPTING COMMUNICATING MLLMS ON THE FLY IN REFERRING EXPRESSION TASKS

Anonymous authors

Paper under double-blind review

### ABSTRACT

Multimodal Large Language Models (MLLMs) exhibit varying comprehension levels in language and perception that complicate interacting with a diverse population of agents, similar to how miscommunication happens in humans, e.g., because intentions are not always known. In this work, we investigate whether MLLMs can adapt to the perceptual weaknesses of the communication partners in an online manner, i.e. change the way they describe their environment in a way that is understandable to their partner while communicating with them, via reinforcement learning. We experiment with two tasks: referring expression identification (REI) and referring expression segmentation (RES), where a speaker agent has to describe an object, and a listener has to identify it. To be successful, the speaker agent must discern the comprehension level of the listener and adapt accordingly, especially when the listener suffers from perceptual weaknesses such as color blindness or blurred vision. Unlike traditional offline alignment methods for LLMs, we fine-tune a Multimodal LLM (MLLM) online to adapt to other agents' conceptual understanding. Our experiments with four MLLMs on four datasets show that online adaptation is feasible in both REI and RES settings.

025 026 027

004

010 011

012

013

014

015

016

017

018

019

021

### 1 INTRODUCTION

028 029

Large Language Models (LLMs) and by extension Multimodal Large Language Models (MLLMs) have demonstrated remarkable capabilities across a variety of tasks (Bubeck et al., 2023; Piergiovanni et al., 2023; Alayrac et al., 2022; Team, 2024; Anil et al., 2023) When catering MLLMs with different architectures, e.g. vision backbones, language backbones, trained with different datasets etc, in our daily lives, we may notice the variability in their comprehension levels related to taskspecific concepts, i.e. what resonates with some MLLMs might not be clear to others. Disparities may exist both in their natural language understanding, e.g., some might understand expert terminology while another might require descriptive explanations, and in the perceptual understanding of visual information, e.g., some might have disabilities such as blurred vision or color blindness.

038 In this work, we focus on enabling MLLMs to adapt to perceptual misunderstandings of their communication partners, e.g., not perceiving colors correctly and therefore not responding to color at-040 tributes presented to them. Specifically, we fine-tune the MLLM online, i.e. on-the-fly, while it is 041 interacting with another MLLM, based on its observed behavior. We model sequential interactions 042 between pairs of agents during a vision-language referring expression tasks which is used as an en-043 vironment for both adaptation and evaluation. Given one or two images, the speaker agent needs 044 to describe the discriminating features of a target object, while the listener agent has to identify the correct object based on this description. To enhance overall task performance, the speaker has to learn which feature of the image allows the listener agent to discriminate the target object and adapt 046 its communication based on the visual concepts understood by the listeners. We consider a referring 047 expression identification (REI) task, where the listener has to identify one target image from a set of 048 two images, and a referring expression segmentation (RES) task, where the listener has to segment the target object within a single image correctly. We present both settings in Fig. 1. 050

We employ several open-source MLLMs, namely LLaVA-7B, LLaVA-13B (Liu et al., 2023b),
 Qwen (Bai et al., 2023), and PaliGemma (Beyer et al., 2024) as the speaker and listener agents
 where the difference in MLLM capabilities and pre-training datasets simulate significant diversity.
 In addition, we introduce perceptual weaknesses to some listeners by providing them with blurred

054 or grayscaled images to further increase listener variety. As the benchmark, we take inspiration 055 from Corona et al. (2019), but create a more realistic setting by modeling the interactions as free-056 form text, adding image transformations to simulate challenging adaptation scenarios, and scaling it 057 to MLLMs. We evaluate the REI task on CLEVR (Johnson et al., 2017), CUB (Wah et al., 2011), 058 and ImageNet (Deng et al., 2009), while we use the RefCOCO (Kazemzadeh et al., 2014) dataset to implement the RES task. We adapt the MLLMs on the fly using PPO (Schulman et al., 2017), KTO (Ethayarajh et al., 2024), and NLPO (Ramamurthy et al., 2023) developed originally as prefer-060 ence learning methods for LLMs when fine-tuning the LoRA adapters (Hu et al., 2022). Contrary to 061 the typical use case of these algorithms for preference optimization (Ouyang et al., 2022; Ahmadian 062 et al., 2024) where a carefully curated offline dataset of human preferences is collected, we test their 063 efficacy during online interactions which is a more realistic and noisier setting. 064

Our contributions are as follows: 1) We introduce a flexible framework for evaluating four MLLMs and adapting them on the fly using four RL algorithms on natural-language-based communication tasks on four datasets to test their efficacy in online adaptation to a diverse set of communication partners. 2) We provide insights into the decision-making process of MLLMs finding that concepts related to color and shape are most important for performing well on these tasks. 3) Through extensive experimental results on two different communication tasks, four MLLMs, and four datasets, we show that adaptation is possible both the REI and RES task.

071 072

073

074

## 2 RELATED WORK

075 A number of methods aim for parameter efficient adaption of large (language) models, which adapt a subset or an additional set of the parameters. LoRA (Hu et al., 2022) and its variants (Zhang et al., 076 2023; Lialin et al., 2023; Liu et al., 2023a; Wu et al., 2024; Sheng et al., 2023; yang Liu et al., 2024) 077 add a trainable residual low rank adaption for each matrix in the network, potentially quantizing it (Dettmers et al., 2024; Xu et al., 2024; Li et al., 2024). In contrast, sparse methods (Ben Za-079 ken et al., 2022; Ansell et al., 2021) only adapt small subsets of the parameters. Adapter based methods (Pfeiffer et al., 2020) train adapter layers and yet another approach is to train a completely 081 separate ladder side networks (Sung et al., 2022; Mercea et al., 2024). As we aim to adapt large 082 multimodal models online, we use LoRA (Hu et al., 2022) for adaptation. 083

For adapting an MLLM to obtain a desired functionality, such as the ability to adapt to a listener 084 online, different RL methods (Snell et al., 2023; Ziegler et al., 2019; Ramamurthy et al., 2023) can 085 be used. Proximal policy optimization (PPO) (Schulman et al., 2017) is an on-policy actor critic algorithm, which is extended by NLPO (Ramamurthy et al., 2023). It restricts the action space to 087 a nucleus of most likely tokens. In contrast KTO (Ethayarajh et al., 2024) directly optimizes the 088 LLM from binary preferences. On the other hand DPO (Rafailov et al., 2024) requires positive and 089 negative pairs for the same context. All of the methods apart from DPO use a single reward per 090 generation making them suitable for our task, thus, we compare their performance. Similar to our 091 work (Guo et al., 2024; Liu et al., 2024) perform (online) adaption based on model feedback in the 092 context of generic model alignment, while we focus on personalization to individual conversational 093 partners and their misunderstandings.

094 Personalizing generative language models has been studied for a long time, often viewed in the context of building an efficient conversational partner in dialogue systems (Serban et al., 2015; Song 096 et al., 2019; Zhang et al., 2019). In contrast, (Ma et al., 2023) reviews several theory of mind 097 (TOM) based approaches to personalization, such as (Takmaz et al., 2023) which proposes a plug-098 and-play TOM based on an explicit simulator, that updates a copy of the model weights on the fly. Similarly, (Raileanu et al., 2018) internally models the behavior of the listener. In contrast, we 099 only update a small amount of parameters using LoRA and do not need to simulate the listeners 100 behavior. (Wang et al., 2024a) adapts the speaker and listener differently, but studies the text-only 101 task, whereas we consider a multi-modal image reference game. We follow an online approach, 102 while (Ma et al., 2021; Zhong et al., 2022) personalizes chatbots by learning from large-scale user 103 dialogue history. 104

Image identification tasks have been studied in visual dialogue settings in (de Vries et al., 2016; Ni et al., 2021; Alaniz et al., 2021; Das et al., 2016). Our work extends this, by incorporating impairments in the communication. (Corona et al., 2019) has studied conceptual image understanding through a reference game, but we extend their attribute constrained setting to free text generation.



Figure 1: The speaker tries to identify a target object, but its pre-trained policy is not aware of misunderstandings of the listener agents, e.g., color blindness. Through interaction with the listener, the speaker learns on-the-fly to mention the shape instead of color because the listener is color-blind. The left interaction illustrates the REI task, while the right interaction shows the RES task.

### 3 ADAPTING THE SPEAKER ON THE FLY IN REFERRING EXPRESSION TASKS

We present a framework for referring expression communication tasks (Figure 1) where a speaker agent describes images to a listener agent using visual concepts. The "speaker" is a single learner that participates in sequences of K episodes describing an image to a group of "listeners".

**Referring Expression Identification (REI) Task.** In the REI task, each episode involves the speaker  $\pi^{(s)}$  and listener  $\pi^{(l)}$  being presented with a pair of images  $[x_k^t, x_k^c]$ . The speaker is assigned one image as the target  $x_k^t$  and the other as a confounding image  $x_k^c$ . The speaker then generates a description  $m_k^{(s)}$  as a message to the listener for it to make its guess regarding the target's identity  $m_k^{(l)}$ , i.e. left or right image. The speaker will observe whether its description led to a correct or incorrect guess via a reward  $r_k \in \{+1, -1\}$  communicated for every episode.

**Referring Expression Segmentation (RES) Task.** In the RES task, the speaker  $\pi^{(s)}$  and listener  $\pi^{(l)}$  are presented with a single image  $x_k$  in each episode. The speaker additionally receives the bounding box of a target object  $o_k^t$  for which the speaker generates a description  $m_k^{(s)}$  with the intention to identify the object in the context of the image. Given the speaker's message, the listener generates a segmentation mask  $m_k^{(l)}$  as a guess regarding the target object in the image. The intersection over union (IoU) metric between the predicted and the ground-truth segmentation masks serves as a reward of the episode for the speaker.

Based on this feedback from the reward alone, the speaker's goal is to change its policy  $\pi^{(s)*}$ , i.e., adapt its image description, to maximize the success rate of the listener agent to solve the referring expression task. To further increase the difficulty of each task, any listener may suffer from a perceptual weaknesses, i.e., color blindness or blurry vision, which is unknown to the speaker.

Since the listener operates as a black box from the perspective of the speaker, pinpointing the source of errors when they exhibit unexpected behavior can be challenging. When the listener makes an incorrect guess, identifying the source of the error becomes difficult, e.g., it could be a lack of comprehension in language, or the visual concepts used to describe the image.

When the listener fails to guess the correct object, the speaker should explore different descriptions
to find a policy tailored for the listener. In this work, we examine, whether LLM adaptation methods
can successfully find policies that maximize task performance for a diverse set of listener agents
solely from the reward signal in this multimodal, i.e. vision and language-based, framework.

158

121

122

123

124 125 126

127 128

129

- 159 3.1 ONLINE MLLM ADAPTATION
- 161 To perform well on the referring expression tasks, the speaker agent needs to adapt to the listener in an online setting during ongoing interactions. After each episode the speaker can update its

177

178

179

197

205



Figure 2: Speaker is asked to describe an object in the context of the REI or RES task. The description is passed to the Listeners which need to decide which image was described. Depending on the correctness of the decision, the Speaker receives a sparse reward and updates its LoRA weights to maximize the reward. For each type of Listener, we have a distinct set of LoRA weights.

weights based on the reward provided by the listener's response. Through these rewards, the speaker
 increases the likelihood of generating descriptions that are adapted to the capabilities of the listener.

Reinforcement learning from human feedback (RLHF) (Ouyang et al., 2022; Christiano et al., 2017;
Stiennon et al., 2020; Ahmadian et al., 2024) is a popular technique to adapt LLMs to human preferences. Typically, a dataset of human preferences is collected, before a RLHF algorithm is applied
either offline or through training a reward model to update the parameters of the LLM or MLLM for
better human alignment. In this work, we explore how well RLHF algorithms extend to an online
setting which is more challenging because the reward data is not carefully annotated and can be
noisy, e.g., when the listener misunderstands the description, but still guesses correctly.

**Proximal Policy Optimization (PPO)** (Schulman et al., 2017) is an on-policy actor-critic algorithm that treats language generation as Markov Decision Process (MDP) where at each state  $s_t$  in the sequence (current context), the next action  $a_t$  is chosen (token), until at the end of the sequence T a reward r is observed. As is typical in RL, the discounted expected reward of the policy is optimized  $\mathbb{E}_{\pi}[\sum_{t=0}^{T} \gamma^t r(s_t, a_t)]$  with  $\gamma$  as the discount factor. PPO starts from the initially pre-trained MLLM  $\pi_{\theta} = \pi_0$  and updates the policy using the following loss:

$$\mathcal{L}_{\text{PPO}}(\pi_{\theta_k}, \pi_{\theta_{k-1}}) = \mathbb{E}_{a_t, s_t \sim \pi_{\theta_k}} \left[ \min \left( \phi_{\pi_{\theta_{k-1}}}^{\pi_{\theta_k}} A^{\pi_{\theta_{k-1}}}, \operatorname{clip}(\phi_{\pi_{\theta_{k-1}}}^{\pi_{\theta_k}}, 1 - \epsilon, 1 + \epsilon) A^{\pi_{\theta_{k-1}}}) \right) \right]$$
(1)

where  $\phi_{\pi\theta_{k-1}}^{\pi\theta_k} = \frac{\pi_{\theta_k}(a_t|s_t)}{\pi_{\theta_{k-1}}(a_t|s_t)}$ ,  $\epsilon$  is a hyperparameter and  $A^{\pi\theta}$  is the advantage function that estimates whether the current action is better than average.

As suggested by Wu et al. (2021), a token-level penalty  $\text{KL}(\pi_q || \pi_p) = (\log \pi_p(a_t | s_t) - \log \pi_q(a_t | s_t))$  regularizes the reward function. This avoids large deviations from the pre-trained MLLM, i.e. the initial policy  $\pi_0$ . The updated reward is computed as:

$$\hat{r}(s_t, a_t) = r(s_t, a_t) - \beta \text{KL}\left(\pi_{\theta} || \pi_0\right)$$
(2)

where the KL coefficient  $\beta$  is a hyperparameter.

**Natural Language Policy Optimization (NLPO)** (Ramamurthy et al., 2023) extends PPO by restricting the action-space with a reduced number of tokens. This is achieved by freezing a masked policy  $\pi_{\psi}$  every  $\mu$  steps and sampling sentences during training from this masked policy. NLPO employs top-*p* sampling for  $\pi_{\psi}$  which limits the sampled tokens to the smallest subset of tokens with cumulative probability greater than the probability *p*. This additional constraints restricts the sampled sentences to be closer to the masked policy, a snapshot of a previous policy, preventing large deviations and divergence.

**Kahneman-Tversky Optimization (KTO)** (Ethayarajh et al., 2024) takes inspiration from prospect theory and proposes to directly optimize the LLM from binary preferences similar to

DPO (Rafailov et al., 2024), instead of performing RLHF. In contrast to DPO, it does not require paired preference data. The loss function is defined as:

$$L^{+}_{\mathrm{KTO}}(\pi_{\theta}, \pi_{0}) = \mathbb{E}_{a_{t}, s_{t} \sim \pi_{\theta}} [\lambda^{+} (1 - \sigma(\beta(\log \phi_{\pi_{0}}^{\pi_{\theta}} - \mathbb{E}_{s' \sim \pi_{\theta}} [\mathrm{KL}(\pi_{\theta} \| \pi_{0})])))] \quad \text{if } r = +1 \qquad (3)$$

$$L^{-}_{\mathrm{KTO}}(\pi_{\theta},\pi_{0}) = \mathbb{E}_{a_{t},s_{t}\sim\pi_{\theta}}[\lambda^{-}(1-\sigma(\beta(\mathbb{E}_{s'\sim\pi_{\theta}}[\mathrm{KL}(\pi_{\theta}\|\pi_{0})] - \log\phi_{\pi_{0}}^{\pi_{\theta}})))] \quad \text{if } r = -1 \qquad (4)$$

that depends on whether a generated sentence produced a +1 or -1 reward.  $\lambda^{+/-}$  are hyperparameters for the two loss terms respectively. Since we do not have a static dataset, we sample sentences on-policy and shuffle the context, i.e. image input and prompt, within each batch for the KL term.

RL algorithms are known to be unstable (Ouyang et al., 2022; Christiano et al., 2017; Ahmadian et al., 2024) which is why KL terms have been introduced for fine-tuning LLMs. Nonetheless, a potential danger that can arise from this is that the policy of the speaker may diverge and start to generate unusual sentences which exploit the listener agent. These sentences may not describe the images correctly, or deviate from being grammatically correct, but enumerations of words instead. Careful selection of hyperparameters is generally important for success with any of these algorithms.

230 231 232

219 220 221

222

223

### 3.2 EFFICIENT ADAPTATION OF THE SPEAKER AGENT

233 Online adaptation of an MLLM does not only require a suitable optimization algorithm, but it should 234 also be feasible in terms of update speed and flexibility as a common use-case may involve a speaker 235 agent interacting with several listeners in parallel. As full-fine-tuning MLLMs is computationally expensive, we adapt these methods by using a parameter-efficient fine tuning method. Given the ver-236 satility of LoRA (Hu et al., 2022) for both the visual domain and the text domain, and its simplicity, 237 we employ it in our architecture. We add LoRA adapters on each linear layer in the LLM-module 238 of the network. As a result, the total number of tuneable parameters are orders of magnitude smaller 239 than the total number of parameters in the MLLM. One can initialize one set of LoRA adapters for 240 each listeners and effortlessly swap out LoRA parameters when interacting with multiple listeners. 241

We employ LLaVA-7B as the speaker model for all experiments because it fits into the memory of a
 single GPU while training with LoRA adapters. Since the listener runs in inference mode, we also
 evaluate on LLaVA-13B, Qwen , and PaliGemma to increase listener diversity.

### 246 4 EXPERIMENTS

247

252

245

We first introduce our experimental setting, i.e. our datasets, the agents, the training, and evaluation
protocol. Then we present the weaknesses and strengths of current MLLMs when dealing with the
visual-language referring expression tasks. Finally, we provide extensive experiments into adapting
a speaker model to different listeners on four different datasets using three algorithms.

253 4.1 EXPERIMENTAL SETTING

254 Datasets. We propose a framework for referring expression tasks on four datasets: CLEVR(Johnson 255 et al., 2017), CUB (Wah et al., 2011), ImageNet(Deng et al., 2009) for REI, and RefCOCO 256 (Kazemzadeh et al., 2014) for RES. CLEVR contains images with objects of varying attributes (size, 257 color, material), requiring fine-grained reasoning to distinguish between different CLEVR scenes. 258 CUB and ImageNet feature natural images with more conversationally relevant concepts. For REI 259 on these datasets, we sample two images, randomly select one as the target, and ask the speaker to 260 describe it in contrast to the other image. We shuffle their order when presenting the images to the listener to avoid trivial solutions, such as "the left image is the target image". Further, we ensure 261 the images come from different classes for CUB and ImageNet. For RES, we employ RefCOCO 262 which extends COCO (Lin et al., 2014) with human-annotated referring expressions and bounding 263 box/segmentation mask annotations. This task requires contrasting a specific detail within an im-264 age's context, posing a different challenge from REI. To visually prompt the speaker on the target 265 object, following Shtedritski et al. (2023) we use a red circle as big as the ground truth bounding 266 box. 267

Agents. Our experiments consider pairs of agents: a speaker and a listener. Specifically, we use
 LLaVA-1.5-7B (Liu et al., 2023b) as the speaker across all adaptation experiments, providing a good balance between its pre-trained capabilities to bootstrap from and a model size that allows us



Figure 3: Performance for various agents on ground-truth descriptions with all attributes and with sets of three attributes for CLEVR.





All The image contains a small purple rubber cylinder, a attributes small cyan metal cylinder, a small blue rubber cylinder and a large brown rubber cylinder.

w/o The image contains a small rubber cylinder and a color small metal cylinder.

- o The image contains a purple rubber cylinder, a
- size cyan metal cylinder and a brown rubber cylinder.w/o The image contains a small purple cylinder, a small
- material cyan cylinder, a small blue cylinder and a large brown cylinder.

Figure 4: Example of ground-truth descriptions (right) on CLEVR for the target image (left) with all attributes, and with sets of three attributes.

to fine-tune LoRA adapters on a single A100 40GB GPU. As listener agents, we employ LLaVA-1.5-7B, LLaVA-1.5-13B, Qwen (7B)(Bai et al., 2023) for REI, and PaliGemma (3B) (Beyer et al., 2024) for RES, which is the only open model of reasonable size capable of producing segmentation masks as output. Each listener model has distinct capabilities when it comes to image and language recognition, with Qwen being the weakest one. This diversity in listener agents simulates a population of agents, testing the speaker's ability to adapt its language effectively.

To introduce an additional challenge, we induce perceptual weaknesses in the listener agents: "color blindness" (grayscaled images) and "blurred vision" (Gaussian blur). These weaknesses require the speaker, which receives unaltered images, to adapt its language to account for concepts that are not recognizable by the listener agent.

296 **Training and evaluation.** We train the speaker (LLaVA-7B) with LoRA adapters on all linear layers 297 of the LLM, keeping the vision module fixed. During online adaptation, we play three episodes 298 before updating the parameters using PPO, NLPO, or KTO algorithms, resulting in a batch size of 3 299 which maximizes our memory usage. The speaker is trained for 1800 interactions (600 update steps) and evaluated on a held-out test set of 300 episodes per dataset. We use the average success rate 300 as evaluation metric for REI and mean IoU for the RES task. Each experiment combines a specific 301 speaker-listener pair either with or without perceptual weaknesses. We provide additional details 302 about the MLLM prompts in Supp. B. 303

304

281

283 284

286

287

288

289

290

291

305 306

### 4.2 EVALUATING LISTENERS WITH GROUND-TRUTH DESCRIPTIONS ON CLEVR

CLEVR's detailed scene descriptions allow us to construct a ground-truth (GT) speaker agent for
 the REI task that produces image descriptions with perfect perception and reasoning abilities. This
 enables us to evaluate listeners given an ideal speaker. The produced descriptions mention all at tributes that appear at least once in the target image, but do not exist in the confounding image. We
 also ablate the GT speaker by omitting one attribute type, measuring the importance of each attribute
 for REI. Examples of these image descriptions are shown in Fig. 4.

313 We evaluate our listener agents alongside GPT-4V, to obtain a reference for a state-of-the-art MLLM, 314 and present the results in Fig. 3. We observe that when all attributes are present, GPT-4V performs best (0.99), followed by LLaVA-7B and LLaVA-13B (0.83 and 0.73), with QWEN being the weak-315 est model (0.63). Removing size and material attributes has little impact on performance, except 316 for a slight increase in LLaVA-13B and QWEN's scores, indicating that size information is more 317 confusing than helpful for these models, possibly because of perspective. In contrast, omitting shape 318 information significantly affects GPT-4V's performance (from 0.99 to 0.84), while the other listeners 319 are less affected, showing that GPT-4V is more sensitive to shape than other models. 320

Most notably, removing color information results in significant performance drops across all listeners, highlighting its importance for solving the REI task on CLEVR. These findings demonstrate that different MLLMs prioritize different attributes and have varying capabilities, as shown in Fig. Even GPT-4V struggles to solve the task without color or shape information.

w/o The image contains a small purple rubber object, a shape small cyan metal object, a small blue rubber object and a large brown rubber object.



Figure 5: Comparing NLPO, PPO, KTO, GT on CLEVR. ZSL: no training was involved. Normal: no perceptual impairment, Blur: Blurry vision, Color blind: Vision with no color. P-value of statistical significance test w.r.t. ZSL: (< 0.1), \* (< 0.05), \*\* (< 0.01), \*\*\* (< 0.001)



Figure 6: Results on the CUB (Top) and ImageNet (Bottom) datasets (REI task). ZSL means that no training was involved. Perceptual weakness refers to the visual impairment applied to the listener. P-value of statistical significance test w.r.t. ZSL: (< 0.1), \* (< 0.05), \*\* (< 0.01), \*\*\* (< 0.001)

### 4.3 COMPARING LISTENERS AND ADAPTATION METHODS ON REI TASK

**REI on CLEVR.** As shown in Fig. 5, when we do not adapt the speaker in the zero-shot learning (ZSL) setting, listener models achieve modest performance. The LLaVA-13B listener achieves the highest performance with an accuracy of 0.58. Introducing color blindness decreases performance for both LLaVA models, while blurred vision has little impact. Qwen performs weakest both with and without perceptual weaknesses, i.e., it struggles to understand the descriptions of LLaVA-7B.

KTO-based adaptation significantly improves performance for LLaVA-7B and LLaVA-13B (peaking at 0.69 and 0.67). Qwen also sees smaller improvements to 0.57. PPO-based adaptation yields
smaller gains, while NLPO shows little improvement over zero-shot learning, except when Qwen
is the listener. Testing these algorithms with perceptual weaknesses reveals reduced performance
increases due to the harder task for the speaker. Blurred vision is generally easier to handle than
color blindness, with KTO performing the best overall.

Compared to using GT descriptions for evaluating the listeners (0.67/0.82/0.85), there is a significant
 gap to the best adaptation results with KTO (0.57/0.67/0.69) even with normal vision. This suggests
 that the REI task is challenging enough for further research in online adaptation of MLLMs.



Figure 7: Qualitative results on CUB and CLEVR when the speaker interacts with a colorblind listener. We present the descriptions generated by the untrained agents (ZSL) and the descriptions obtained after training (Adapted). After adaption, the speaker avoids color attributes.

REI on natural images. Fig. 6 presents the adaptation results on CUB and ImageNet using natural images. We observe that all listeners perform well in ZSL, with LLaVA-13B achieving an accuracy of 0.86 (CUB) and 0.87 (ImageNet). The MLLMs are likely more familiar with such natural images making it easier for the speaker to pick out differences and the listener to recognize them. However, there is still a large gap to Qwen with 0.63/0.73 for CUB/ImageNet.

In general, adaptation methods provide a boost in performance for all listeners. While KTO-based adaptation excels on ImageNet, all three algorithms perform similarly well on CUB. Perceptual weaknesses have a larger impact on CUB, with removing color having the highest effect on performance. On ImageNet both weaknesses only slightly decrease the performance. This is consistent across listeners and algorithms.

In conclusion, online adaptation is possible for every tested agent and algorithm on the REI task. 410 However, listener capabilities influence improvements, and different algorithms perform better on 411 different datasets and listeners. Overall, KTO seems to work best when considering all experiments. 412 At the same time, none of the existing algorithms are able to find a policy that achieves results close 413 to the of the GT agent leaving room for improvement. Moreover, we find that adaptation on blurred 414 or grayscale images can reach or surpass zero-shot learning performance on normal images, which is 415 a desirable outcome in scenarios where we want to avoid a disadvantage for agents with perceptual 416 weaknesses. This applies to a lesser degree on ImageNet, and was not generally true on CUB, where 417 achieving this target could be an promising direction within the REI task framework.

418 419

396

397 398 399

### 4.4 ADAPTING TO PALIGEMMA ON THE RES TASK

420 421

On the referring expression segmentation task, we adapt the LLaVA-7B speaker to PaliGemma as
listener on the RefCOCO dataset. In Fig. 9, we report the mean intersection over union (mIoU) for
ZSL, PPO, NLPO, and KTO together with probing the PaliGemma listener with the ground truth
(GT) referring expressions created by humans that come with the dataset.

We find that the RES task poses a particular challenge to some adaptation algorithms, because neither
PPO or NLPO can significantly improve over the zero-shot descriptions in normal, blurred, and
grayscaled images. Only KTO manages to obtain an improvement from 0.34 to 0.44 for normal
images, from 0.28 to 0.41 in blurry images, and from 0.28 to 0.40 in grayscale images. At the same
time, the GT descriptions still outperform the KTO adapted speaker reaching 0.63/0.56/0.61 mIoU
in the three settings respectively. Thus, we conclude that there is still room for improvement for
online adaptation to reach closer to human performance.



Figure 8: Qualitative results of the RES task on RefCOCO with a LLaVA-7B speaker and coloblind PaliGemma listener.

460 When inducing perceptual weaknesses on the 461 PaliGemma listener, ZSL performance degrades, but 462 to a lesser degree than for the REI task. This is ex-463 pected because objects that are contrasted in RES are 464 often easier to identify by their relation in the scene, 465 e.g., where it is located spatially rather than by color 466 or shapes. As a result, PaliGemma can deal with blur 467 and grayscale relatively well in this context. Apart 468 from KTO being the best adaptation algorithm for 469 RES, we also find that KTO can adapt to perceptu-470 ally weakened listeners to improve over ZSL perfor-471 mance of the normal listener.

472 473

475

458 459

- 4.5 QUALITATIVE
- 474 ANALYSIS ON COLORBLIND LISTENER
- LLaVA-7B speaker and PaliGemma listener. (< 0.01), \*\*\* (< 0.001)

0.6

0.2

0.0

mIoU 0.4

476 In Fig. 7, we show qualitative results on CUB, and

P-value w.r.t. ZSL: (< 0.1), \* (< 0.05), \*\*

Normal

ZSL

NLPO

PaliGemma

Blur

PPO

KTO 

Figure 9: mIoU on RefCOCO for RES with

\*\*>

Color-blind

GT

477 CLEVR, by contrasting generated descriptions before and after adaptation on the REI task when interacting with a colorblind listener. 478

479 We observe that color attribute is mentioned predominantly before adaptation, and, apart from refer-480 ring to "black" and "white", completely avoided after adaptation. On CUB for instance, the speaker 481 mentions the "yellow chest" and "yellow beak" to discriminate the birds in the zero-shot setting, 482 and learns to focus the description more on the surrounding scene and action performed by the bird 483 to discriminate the two images after adaptation. On CLEVR, descriptions similarly contains many references to the color attributes in the initial descriptions, but they do not mention colors after adap-484 tation. In contrast, the adapted descriptions focus on the overall count of the objects and are more 485 concise than the original ones. Moreover, zero-shot descriptions sometimes mix objects from both images, e.g., description in the third row mentions "red cube" and "blue cube" from the left image,
and "green sphere" and "yellow sphere" from the right image. After adaptation this behaviour is
suppressed and the speaker focuses more on the target image.

In Fig. 8, we show examples of the adaptation on RefCOCO for the RES task, again when the listener 490 is colorblind. The first two rows exemplify how mentioning color can confuse the listener, e.g., in 491 the second row, where the listener segments the incorrect baseball player because it cannot attribute 492 the "blue" uniform to the correct one. After adaptation, not mentioning the colors and focusing on 493 other aspects, such as the "suit and tie" in the first example, allows the listener to more accurately 494 segment the target. Interestingly, there are a few examples where the visual prompting through the 495 red circle (Shtedritski et al., 2023) can cause incorrect descriptions mentioning the circle which is 496 not visible to the listener. However, online adaptation can also correct for this failure case as seen in the third row, where the speaker correctly refers to the "zebra on left". 497

In conclusion, from these qualitative examples, we observe that the speaker learns to correctly iden tify the perceptual weakness of the listener, and adapts its description accordingly to be more effective in its communication.

501 502 503

504

### 5 LIMITATIONS

505 As it is widely known in the literature (Ouyang et al., 506 2022; Christiano et al., 2017; Ahmadian et al., 2024), RL 507 algorithms tend to be unstable when the reward signal is 508 noisy, or the actions space immense. During this study, 509 we have observed that there is a divergence effect during online adaptation. Fig. 10 exemplifies this divergence 510 effect on CLEVR dataset for LLaVA-13B which is rep-511 resentative of the observations on other datasets and with 512 other listeners. For all our experiments, we report the per-513 formance after 1800 episodes. However, Fig. 10 shows 514 the peak performance is sometimes achieved at different 515 times during training due to the variance in online adapta-516 tion. One potential reason for this is the online nature of 517 gathering training samples. The constantly changing pol-518 icy during training affects the generated data, which in 519 turn influences the future policy and exploration of possi-520 ble descriptions. With the large actions space of MLLMs, it is challenging to keep these effects in check. 521



Figure 10: Divergence effect on CLEVR for LLaVA-13B. The performance fluctuates instead of monotonically improving.

### 6 CONCLUSION

524 525

522 523

In this work, we introduce a framework for two referring expression tasks (REI/RES) involving com-526 municating MLLM agents. On these tasks, we study how MLLM agents can adapt to one another 527 on-the-fly. Our online adaptation setting is significantly more challenging than aligning MLLMs on 528 carefully collected offline datasets, while opening up new applications that require individual per-529 sonalization. Every communication partner understands language and concepts required to solve the 530 tasks at different levels and we introduce perceptual weaknesses to further control for agent variety. 531 The referring expression tasks pose a challenge to currently available MLLMs, especially for images 532 with fine-grained differences, and when precise segmentation is required. All the adaptation algo-533 rithms we have tested could improve task performance on REI with KTO working the best overall 534 and being the only one achieving improvements on the RES task. These results show that, 1) it is possible to improve over the initial pre-trained policy by learning about the listener capabilities, and 536 2) we can perform this learning in an online setting. However, we also observe that current methods 537 do not monotonically improve during the training process, and cannot find an "optimal" policy, since we have demonstrated that better ones exist with our GT agent experiments. With our task setting, 538 we want to encourage further research on how to make online adaptation of MLLM effective and practically viable to extend to real-world scenarios for MLLM personalization.

### 540 REFERENCES 541

574

575

576 577

- Arash Ahmadian, Chris Cremer, Matthias Gallé, Marzieh Fadaee, Julia Kreutzer, Olivier Pietquin, 542 Ahmet Ustün, and Sara Hooker. Back to basics: Revisiting REINFORCE style optimization for 543 learning from human feedback in llms. In ACL, 2024. 544
- Stephan Alaniz, Diego Marcos, and Zeynep Akata. Learning decision trees recurrently through 546 communication. In CVPR, 2021. 547
- 548 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, 549 Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob L. Menick, 550 Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, 551 Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karén Simonyan. Flamingo: a visual 552 language model for few-shot learning. In NeurIPS, 2022. 553
- 554 Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Millican, David Silver, Slav Petrov, Melvin 555 Johnson, Ioannis Antonoglou, Julian Schrittwieser, Amelia Glaese, Jilin Chen, Emily Pitler, Tim-556 othy P. Lillicrap, Angeliki Lazaridou, Orhan Firat, James Molloy, Michael Isard, Paul Ronald Barham, Tom Hennigan, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong Xu, Ryan 558 Doherty, Eli Collins, Clemens Meyer, Eliza Rutherford, Erica Moreira, Kareem Ayoub, Megha 559 Goel, George Tucker, Enrique Piqueras, Maxim Krikun, Iain Barr, Nikolay Savinov, Ivo Danihelka, Becca Roelofs, Anaïs White, Anders Andreassen, Tamara von Glehn, Lakshman Yagati, 561 Mehran Kazemi, Lucas Gonzalez, Misha Khalman, Jakub Sygnowski, and et al. Gemini: A family of highly capable multimodal models. In arXiv preprint arXiv:2312.11805, 2023. 563
- 564 Alan Ansell, E. Ponti, Anna Korhonen, and Ivan Vulic. Composable sparse fine-tuning for crosslingual transfer. In ACL, 2021. 565
- 566 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, 567 Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, 568 Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi 569 Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng 570 Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, 571 Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, 572 Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. Qwen technical report. In arXiv preprint 573 arXiv:2309.16609, 2023.
  - Elad Ben Zaken, Yoav Goldberg, and Shauli Ravfogel. BitFit: Simple parameter-efficient finetuning for transformer-based masked language-models. In ACL, 2022.
- Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel Salz, 578 Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschannen, Emanuele Bugliarello, et al. Paligemma: A versatile 3b vlm for transfer. In arXiv preprint arXiv:2407.07726, 2024.
- 580 Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Ka-581 mar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott M. Lundberg, Harsha Nori, Hamid Palangi, 582 Marco Túlio Ribeiro, and Yi Zhang. Sparks of artificial general intelligence: Early experiments 583 with GPT-4. In arXiv preprint arXiv:2303.12712, 2023. 584
- 585 Paul F. Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep 586 reinforcement learning from human preferences. In NeurIPS, 2017.
- Rodolfo Corona, Stephan Alaniz, and Zeynep Akata. Modeling conceptual understanding in image 588 reference games. In NeurIPS, 2019.
- Abhishek Das, Satwik Kottur, Khushi Gupta, Avi Singh, Deshraj Yadav, José M. F. Moura, Devi Parikh, and Dhruv Batra. Visual dialog. In CVPR, 2016. 592
- Harm de Vries, Florian Strub, A. P. Sarath Chandar, Olivier Pietquin, H. Larochelle, and Aaron C. Courville. Guesswhat?! visual object discovery through multi-modal dialogue. In CVPR, 2016.

- 594 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale 595 hierarchical image database. In CVPR, 2009. 596 Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning 597 of quantized llms. In NeurIPS, 2024. 598 Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Kto: Model 600 alignment as prospect theoretic optimization. In arXiv preprint arXiv:2402.01306, 2024. 601 602 Shangmin Guo, Biao Zhang, Tianlin Liu, Tianqi Liu, Misha Khalman, Felipe Llinares-López, Alexandre Ramé, Thomas Mesnard, Yao Zhao, Bilal Piot, Johan Ferret, and Mathieu Blondel. 603 Direct language model alignment from online ai feedback. In arXiv preprint arXiv:2402.04792, 604 2024. 605 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, 607 and Weizhu Chen. Lora: Low-rank adaptation of large language models. In ICLR, 2022. 608 Justin Johnson, Bharath Hariharan, Laurens Van Der Maaten, Li Fei-Fei, C Lawrence Zitnick, and 609 Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual 610 reasoning. In CVPR, 2017. 611 612 Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara Berg. Referitgame: Referring to 613 objects in photographs of natural scenes. In EMNLP, 2014. 614 Yixiao Li, Yifan Yu, Chen Liang, Pengcheng He, Nikos Karampatziakis, Weizhu Chen, and Tuo 615 Zhao. Loftq: Lora-fine-tuning-aware quantization for large language models. In ICLR, 2024. 616 617 Vladislav Lialin, Namrata Shivagunde, Sherin Muckatira, and Anna Rumshisky. Relora: High-rank 618 training through low-rank updates. In NeurIPS workshops, 2023. 619 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr 620 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In ECCV, 2014. 621 622 Liu, Jing Liu, Toshiaki Koike-Akino, Pu Wang, Matthew Brand, Ye Wang, and Kieran Parsons. 623 Loda: Low-dimensional adaptation of large language models. In *NeurIPS workshops*, 2023a. 624 Aiwei Liu, Haoping Bai, Zhiyun Lu, Xiang Kong, Simon Wang, Jiulong Shan, Mengsi Cao, and 625 Lijie Wen. Direct large language model alignment through self-rewarding contrastive prompt 626 distillation. In ACL, 2024. 627 628 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*, 629 2023b. 630 Zhengyi Ma, Zhicheng Dou, Yutao Zhu, Hanxun Zhong, and Ji rong Wen. One chatbot per person: 631 Creating personalized chatbots based on implicit user profiles. In SIGIR, 2021. 632 633 Ziqiao Ma, Jacob Sansom, Run Peng, and Joyce Chai. Towards a holistic landscape of situated 634 theory of mind in large language models. In EMNLP, 2023. 635 Otniel-Bogdan Mercea, Alexey Gritsenko, Cordelia Schmid, and Anurag Arnab. Time-memory-and 636 parameter-efficient visual adaptation. In CVPR, 2024. 637 638 Jinjie Ni, Tom Young, Vlad Pandelea, Fuzhao Xue, Vinay Vishnumurthy Adiga, and E. Cambria. 639 Recent advances in deep learning based dialogue systems: a systematic survey. In Artificial Intelligence Review, 2021. 640 641 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, 642 Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser 643 Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan 644 Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. 645 In NeurIPS, 2022. 646
- 647 Jonas Pfeiffer, Ivan Vulic, Iryna Gurevych, and Sebastian Ruder. Mad-x: An adapter-based framework for multi-task cross-lingual transfer. In *EMNLP*, 2020.

683

- 648 A. J. Piergiovanni, Isaac Noble, Dahun Kim, Michael S. Ryoo, Victor Gomes, and Anelia Angelova. 649 Mirasol3b: A multimodal autoregressive model for time-aligned and contextual modalities. In 650 arXiv preprint arXiv:2311.05698, 2023. 651
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea 652 Finn. Direct preference optimization: Your language model is secretly a reward model. In 653 NeurIPS, 2024. 654
- 655 Roberta Raileanu, Emily L. Denton, Arthur Szlam, and Rob Fergus. Modeling others using oneself 656 in multi-agent reinforcement learning. In ICML, 2018.
- Rajkumar Ramamurthy, Prithviraj Ammanabrolu, Kianté Brantley, Jack Hessel, Rafet Sifa, Chris-658 tian Bauckhage, Hannaneh Hajishirzi, and Yejin Choi. Is reinforcement learning (not) for natural 659 language processing: Benchmarks, baselines, and building blocks for natural language policy 660 optimization. In ICLR, 2023. 661
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy 662 optimization algorithms. In *arXiv preprint arXiv:1707.06347*, 2017. 663
- Iulian Serban, Alessandro Sordoni, Yoshua Bengio, Aaron C. Courville, and Joelle Pineau. Building 665 end-to-end dialogue systems using generative hierarchical neural network models. In AAAI, 2015. 666
- 667 Ying Sheng, Shiyi Cao, Dacheng Li, Coleman Hooper, Nicholas Lee, Shuo Yang, Christopher Chou, Banghua Zhu, Lianmin Zheng, Kurt Keutzer, Joseph E. Gonzalez, and Ion Stoica. S-lora: Serving 668 thousands of concurrent lora adapters. In arXiv preprint arXiv:2311.03285, 2023. 669
- 670 Aleksandar Shtedritski, Christian Rupprecht, and Andrea Vedaldi. What does clip know about a red 671 circle? visual prompt engineering for vlms. In ICCV, 2023. 672
- Charles Burton Snell, Ilya Kostrikov, Yi Su, Mengjiao Yang, and Sergey Levine. Offline rl for 673 natural language generation with implicit language q learning. In ICLR, 2023. 674
- 675 Haoyu Song, Weinan Zhang, Jingwen Hu, and Ting Liu. Generating persona consistent dialogues 676 by exploiting natural language inference. In AAAI, 2019. 677
- Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, 678 Dario Amodei, and Paul F. Christiano. Learning to summarize from human feedback. In NeurIPS, 679 2020. 680
- 681 Yi-Lin Sung, Jaemin Cho, and Mohit Bansal. Lst: Ladder side-tuning for parameter and memory 682 efficient transfer learning. In NeurIPS, 2022.
- Ece Takmaz, Nicolo' Brandizzi, Mario Giulianelli, Sandro Pezzelle, and Raquel Fern'andez. Speak-684 ing the language of your listener: Audience-aware adaptation via plug-and-play theory of mind. 685 In ACL, 2023. 686
- 687 Chameleon Team. Chameleon: Mixed-modal early-fusion foundation models. In arXiv preprint 688 arXiv:2405.09818, 2024.
- 689 Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd 690 birds-200-2011 dataset. 2011.
- 692 Jian Wang, Chak Tou Leong, Jiashuo Wang, Dongding Lin, Wenjie Li, and Xiao-Yong Wei. Instruct 693 once, chat consistently in multiple rounds: An efficient tuning framework for dialogue. In ACL, 2024a. 694
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, 696 Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng 697 Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. arXiv preprint arXiv:2409.12191, 2024b. 699
- Jeff Wu, Long Ouyang, Daniel M. Ziegler, Nisan Stiennon, Ryan Lowe, Jan Leike, and 700 Paul F. Christiano. Recursively summarizing books with human feedback. In arXiv preprint arXiv:2109.10862, 2021.

702	Yichao Wu Yafei Xiang Shuning Huo Yulu Gong and Penghao Liang Lora-sp: Streamlined
703	internet in a section from the constant of the tender of least loss between the section of the tender of tender
	partial parameter adaptation for resource-efficient line-tuning of large language models. In <i>arXiv</i>
704	preprint arXiv:2403.08822, 2024.
705	

- Yuhui Xu, Lingxi Xie, Xiaotao Gu, Xin Chen, Heng Chang, Hengheng Zhang, Zhensu Chen, Xiaopeng Zhang, and Qi Tian. Qa-lora: Quantization-aware low-rank adaptation of large language models. In *ICLR*, 2024.
- Shih yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, KwangTing Cheng, and Min-Hung Chen. Dora: Weight-decomposed low-rank adaptation. In *ICML*, 2024.
- Longteng Zhang, Lin Zhang, Shaohuai Shi, Xiaowen Chu, and Bo Li. Lora-fa: Memory-efficient low-rank adaptation for large language models fine-tuning. In *arXiv preprint arXiv:2308.03303*, 2023.
- Yizhe Zhang, Xiang Gao, Sungjin Lee, Chris Brockett, Michel Galley, Jianfeng Gao, and William B.
   Dolan. Consistent dialogue generation with self-supervised feature learning. In *arXiv preprint* arXiv:1903.05759, 2019.
- Hanxun Zhong, Zhicheng Dou, Yutao Zhu, Hongjin Qian, and Ji rong Wen. Less is more: Learning to refine dialogue history for personalized dialogue generation. In *NAACL-HLT*, 2022.
- Daniel M. Ziegler, Nisan Stiennon, Jeff Wu, Tom B. Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. In *arXiv* preprint arXiv:1909.08593, 2019.

# 756 SUPPLEMENTARY MATERIAL

## 758 A BROADER IMPACT

759 760

In this work we study the capabilities of a speaker to adapt to a listener. We considered MLLMs adapting to other MLLMs, but one could apply these methods also for adapting MLLMs to humans. If such techniques were used to adapt MLLMs to humans, people with malicious intend could purposefully teach the MLLMs to produce harmful or otherwise undesirable content. Online adaptation could effectively overwrite previously learned safety measures of the alignment phase. A possible solution could involve intertwining or following online adaptation with alignment training. Additional research is require to measure both opportunities and risks in this scenario.

In the setting where we adapt an MLLM agent to another MLLM agent, malicious actors could try to exploit systems employing MLLMs by programmatically learning to maximize a desired action of the target MLLM. These "hacks" or "jailbreaks" are a security concern for everyone deploying MLLM, especially if they are deployed adapting to the users. As a result, research on defense mechanisms is just as important as developing more advanced ways to enable personalization.

On the other hand, we believe that allowing MLLMs to adapt to the specific needs of a users can enable new use cases and improve inclusion across diverse population groups. More effective communication towards users with disabilities could lower the barrier of entry and learning curve to bring MLLM technology and their advancement to a broad audience.

### 777 778

779

781 782

783

784

785

788

789

790 791

792

793

794

796

797

798

800

801 802

803 804

805

## **B** MLLM PROMPTING DETAILS

The referring expression identification (REI) task starts with the speaker generating a description for the target image. The prompt given to the speaker is:

Write a description for the left/right image, such that it can be differentiated from the right/left image, but do not talk about the right/left image. Do not name which image you are describing.

Subsequently, with the help of the speaker's response, the listener generates a sentence containingits guess. For LLaVA listener agents, we use the query template:

```
Does this sentence: {}' \{m^{(s)}\}' describe the left image or the right image? Do not explain your reasoning.
```

where  $\{m^{(s)}\}\$  is replaced with the description written by the speaker. On the other hand, Qwen gets the prompt:

Which image does the sentence  ${}'\{m^{(s)}\}'$  describe? A. Picture 1 B. Picture 2.

After receiving the listener's answer, the reward is computed by looking for keywords, i.e. "left, A, 1" and "right, B, 2", and comparing it with the ground truth label.

For the referring expression segmentation (RES) task, the prompt given to the speaker is:

Write a short description for the highlighted object.

The PaliGemma listener is then prompted with:

segment  $'\{m^{(s)}\}'$ 

where "segment" is a PaliGemma specific keyword to induce its segmentation capabilites. The
model proceeds to output tokens that can be translated to a segmentation mask. We calculate the
intersection over union (IoU) between the predicted segmentation mask and the ground truth segmentation mask as a reward for the speaker. Since KTO requires a binary reward, we binarize the
IoU values with a threshold of 0.5.



Figure 11: Performance for ground-truth descriptions with blurred vision (left) and colorblindness (right).

Since LLaVA models can only take a single picture as input, we concatenate the images horizontally and add a white bar between them before feeding them to LLaVA-7B and LLaVA-13B. As a result of this step, LLaVA refers to the images as left or right. No such processing is necessary with Qwen, as it can handle multiple images in a single query. Qwen automatically labels them as picture 1 and picture 2.

### C GROUND-TRUTH DESCRIPTIONS WITH PERCEPTUALLY WEAKENED LISTENERS

We present the evaluation of the GT speaker against listeners with perceptual weakness in Fig. 11.
We observed that both blurry and grayscale images cause a significant drop in performance, with the latter having the greatest impact.

When all attributes are mentioned, blurring decreases the scores of LLaVA-7B and Qwen from 837 0.83 and 0.63 to 0.74 and 0.53. LLaVA-13B maintains its accuracy of 0.73. When the speaker 838 additionally does not mention any color attributes in the description, the accuracy of all listeners 839 drop to near-random performance (i.e., 0.5), with LLaVA-13B performing best at 0.56 accuracy. 840 This result indicates that colors are vital for agents with blurry vision. Removing shapes from 841 the descriptions increases the scores by a small margin in all cases, which suggests this information 842 could be confusing in the presence of blur. Additionally, LLaVA models gain a few percent accuracy 843 when materials are not mentioned in the description. Finally, we would like to highlight that Qwen 844 achieves at most 0.55 score in this setup, which is very close to random guessing.

845 With grayscale images, LLaVA-7B achieves the lowest score of 0.55 when shape information is 846 lacking in the descriptions, and has the highest accuracy of 0.62 with colors removed. The worst 847 and best cases for LLaVA-13B are again without shape (0.51) and without color (0.65), which have 848 a larger difference compared to the smaller version of LLaVA. Those results show color information 849 starts to confuse the models as it is useless, and mentioning shape is more important in this case. 850 Similar to blurry images, Qwen has a very low performance, with a maximum score of 0.56. These 851 observations support our previous findings that shape and color are the most important attributes for performing well on the REI task with CLEVR images. 852

853 854

855

821

822 823 824

825

826

827

828 829 830

831

832

### C.1 ADDITIONAL QUALITATIVE RESULTS ON REI

In Fig. 12 we show qualitative results for the REI task on CLEVR, CUB and ImageNet by contrasting generated descriptions before and after adaptation. In CLEVR the original description is much
longer and even if the speaker is able to mention all the objects in the image, the associated shapes
and color are oftentimes incorrect. On the other hand, after adaptation, the descriptions are much
shorter, mentioning a subset but distriminative part of the scene. The adapted policy frequently mentioning shapes ("blocks", "balls") and colors ("yellow and silver") provides additional evidence that
these attributes are important and easier to recognize for MLLMs in this context.

The ZSL descriptions generated for CUB images are generic and long, often applying to both images. The speaker tends to confuse the confounding image into the description, for instance, when



Figure 12: Qualitative results for CLEVR, CUB and ImageNet datasets. We present the descriptions generated by the untrained agents (ZSL) and the descriptions obtained after training (Adapted).

talking about the bird "facing the camera" and the "black and yellow" bird mixing the colors of both birds. In contrast, the trained agent just mentions the essential distinguishable aspects of the target images ("brown and white feathers" and "black and white bird with a red head"). Lastly, on ImageNet, one failure case of the untrained speaker is that it describes both images without clearly identifying the target. After training, it learns to focus on describing the content of the target image by itself. In conclusion, from these qualitative examples, we observe that the model learns to be more concise, focusing on the correct image and primarily mentions the relevant attributes, which more frequently include color and shape.

### D COMPUTATIONAL RESOURCES

For every experiment, we use 2x A100 40GB GPUs, where one GPU is used for the listener and the other for the speaker. Since the speaker is trained, it requires more computational resources than the listener. It is possible to fit a 13B parameter model into the memory of a single GPU in inference mode for the listener. However, training MLLM only allows models up to 7B parameters on a single GPU, even when using a parameter-efficient fine-tuning method such as LoRA. The training time depends on the lengths of sentences LLaVA generates as the speaker. Longer token sequences take more time to produce as well as to backpropagate through the model. While the length of generations usually diminishes as the speaker adapts to the listener, we also observe the the generated descriptions vary in lengths for the different dataset. Overall, a single experiment of playing 1800 REI episodes and performing 600 update steps (batch size 3) takes around 5-6 hours training time. 

### E HYPERPARAMETERS

For all experiments, we perform a grid search over a subset of hyperparameters and report the results of the best set of hyperparameters. Generally, there was no single set of hyperparameters that performed well across all experiment. The hyperparameters that we considered for grid search are: the learning rate lr, the rank r of the LoRA and the  $\alpha$  parameters in LoRA. Depending on the algorithms, datasets and models, the lr was searched in the interval [1e-7, 1e-8, 1-9], the r was searched in the interval [32, 64, 128] and the  $\alpha$  was searched in the interval [64, 128, 256, 512,



Figure 13: Number of unique words produced by the LLaVA-7B speaker before (ZSL) and after (NLPO, PPO, KTO) adaptation to different listeners.



Figure 14: Average sentence length produced by the LLaVA-7B speaker before (ZSL) and after (NLPO, PPO, KTO) adaptation to different listeners.

1024, 2048]. The remaining hyperparameters were kept fixed without performing a grid search. Specifically, for  $\beta$  in KTO we used 0.1, and for PPO and NLPO we used 0.2.  $\epsilon$  in PPO was set to 1, top-p sampling in NLPO was set to 0.9.  $\lambda^-$  and  $\lambda^+$  were set to 1.0.

### F ANALYSIS OF ADAPTED LANGUAGE

We analyze the language of the LLaVA-7B speaker before and after adaptation. Figure 13 shows the number of unique words, i.e., the speaker's vocabulary size when interacting with the different listeners. We find that PPO consistently reduces the number of unique words the speaker uses. However, when interacting with the LLaVA listeners, NLPO does not change the vocabulary size of the speaker. Similarly, KTO also retains the number of unique words when interacting with LLaVA-13B. 

Figure 14 shows the average sentence length of the LLaVA-7B speaker. The statistics follow a similar trend to the unique words. Interestingly, NLPO and KTO can even increase the average sentence length, especially when interacting with LLaVA-13B. 

In this study, our goal is to adapt to a given listener which can include a change in language char-acteristics, such as avoiding color words for the color-blind listener. This typically leads to a more concise and effective communication between MLLMs. If language drift should be avoided as much as possible, the hyperparameter for the KL term of each adapatation algorithm can be increased.

#### **QWEN2-VL EXPERIMENTS AND TASK DIFFICULTY** G

We extend our analysis to include Qwen2-VL-7B Wang et al. (2024b) as a recent open MLLM which is generally stronger than the models evaluated in the main paper. In Table 1, we report the ZSL

972	Image Pairing	Normal	B&W	Blur	Occlusion
973	Random	0.96	0.69	0.92	0.78
974	Equal #obj. & overlap	0.95	0.56	0.89	0.68
975	Min. 8 objects	0.89	0.57	0.85	0.62

Table 1: ZSL performance of Qwen2-VL on the REI task as both speaker and listener on CLEVR
for all impairments. Different image pairing strategies alter the difficulty of the task. Normal: no
perceptual impairment, Blur: Blurry vision, B&W: Vision with no color, Occlusion: Part of image not visible.

	LLaVA-7B							
Qwen2-VL	Normal	Blur	B&W					
ZSL	0.71	0.66	0.54					
KTO	0.72	0.66	0.56					
PPO	0.74	0.66	0.56					

Table 2: Results of the REI task on the CLEVR dataset. Qwen2-VL-7B is the speaker and LLaVA-7B the listener. ZSL means that no training was involved. Normal: no perceptual impairment, Blur:
Blurry vision, B&W: Vision with no color.

990 991

976

performance of Qwen2-VL as both speaker and listener on the REI task. The first row (random) is 992 the standard evaluation setting where we randomly sample two images from the CLEVR dataset. We 993 observe that it performs significantly better than any other MLLM reaching close to a perfect score 994 both without impairment and even with the blurry impairment. Additionally, we include experiments 995 on the occlusion impairment as described in Section H. To increase the difficulty of our proposed 996 task, we can alter the sampling of the image pairs. For example, in the second row we only sample 997 images with an equal number of objects and, for every episode, pick one our of 1000 image pairs 998 for which there is the most overlap in identical objects in the scene. This increases difficulty such that the results for the colorblind listener drops from 0.69 to 0.56. Another option is to always 999 sample images with at least 8 objects which is equally challenging for the colorblind listener and 1000 also increases difficulty for all other settings, e.g., listener with occlusion impairment drops from 1001 0.78 to 0.62. Overall, while strong MLLMs can often achieve a high zero-shot learning performance 1002 on the REI task, we can increase its difficulty by sampling hard image pairs. 1003

In Table 2, we adapt a Qwen2-VL speaker to a LLaVA-7B listener the REI task on CLEVR. We observe that it is generally more challenging to adapt a strong MLLM, such as Qwen2-VL. There are small improvements when adapting Qwen2-VL on a listener without impairment (+3%) or a colorblind listener (+2%), but no improvement on a listener with blurry vision.

1008

# 1009 H OCCLUSION IMPAIRMENT

To extend the number of impairments, we explore occlusion as another option. For this impairment, we remove part of both images for the REI task on the listener side. Specifically, we black out the left side of the image up to a given ratio. In Table 3, we report ZSL experiments with Qwen2-VL and LLaVA-7B as both the speaker and listener. We observe that for Qwen2-VL as the speaker, occluding half of the image already reduces performance while for LLaVA-7B this only happens starting from 60% occlusion. As such, occlusion could be used as another type of impairment, for which we leave adaptation experiments to future work.

1018 1019

### I QUALITATIVE EXAMPLES OF FAILURE CASES

1020

Figure 15 shows examples of adaptation on REI tasks for the CUB and CLEVR datasets, where the trained model fails to produce descriptions that help the colorblind listener make correct guesses even after adaptation. Similar to Figure 7, we show descriptions before and after adaptation.

In the zero-shot setting, the captions include color information, which the listener cannot perceive.
 After adaptation, the models often removes color references, but sometimes fails to make any other adjustment to include supplementary information to make the images distinguishable. For example,

1026			(	Occlusi	on Rati	)	
1027	Speaker / Listener	0	0.5	0.6	0.7	0.8	0.9
1028	Qwen2-VL / Qwen2-VL	0.97	0.83	0.78	0.69	0.55	0.54
1029	Qwen2-VL / LLaVA-7B	0.70	0.58	0.55	0.51	0.49	0.52
1030	LLaVA-7B / Qwen2-VL	0.60	0.61	0.53	0.54	0.53	0.52
1031	LLaVA-7B / LLaVA-7B	0.55	0.55	0.54	0.51	0.51	0.51

Table 3: ZSL performance of different speaker-listener pairs on the CLEVR dataset when occluding the left half of each image in the REI task.



Figure 15: Qualitative results on CUB and CLEVR when the speaker interacts with a colorblind listener and the decision of the listener was wrong after adaptation.

in the first and third rows, both images match the adapted captions, making it hard for the listener to choose the correct one. This suggests that while adaptation helps in some ways by removing color information, the model cannot always introduce other relevant information.

### J RESULT TABLES

	all	w/o shape	w/o material	w/o size	w/o color
Qwen	0.63	0.63	0.67	0.65	0.56
LLaVA-13b	0.73	0.74	0.82	0.81	0.65
LLaVA-7b	0.83	0.83	0.85	0.81	0.65
GPT-4V	0.99	0.84	0.96	0.95	0.78

Table 4: Performance for various agents on ground-truth descriptions with all attributes and with sets of three attributes for CLEVR.

In Tables 4, 5, 6, 7, and 8, we report the results from Figures 3, 5, 6, and 9, respectively.

LLaVA-7B			Ι	LaVA-	Qwen			
Normal	Blur	B&W	Normal	Blur	B&W	Normal	Blur	B&
0.55	0.57	0.51	0.58	0.57	0.52	0.51	0.50	0.5
0.56	0.59	0.55	0.60	0.57	0.54	0.58	0.52	0.5
0.64	0.56	0.55	0.63	0.58	0.56	0.57	0.52	0.:
0.69	0.58	0.56	0.67	0.61	0.56	0.57	0.54	0.
0.85	0.80	0.62	0.82	0.79	0.65	0.67	0.55	0.
	Normal 0.55 0.56 0.64 0.69 0.85	LLaVA-           Normal         Blur           0.55         0.57           0.56         0.59           0.64         0.56           0.69         0.58           0.85         0.80	Normal         Blur         B&W           0.55         0.57         0.51           0.56         0.59         0.55           0.64         0.56         0.55           0.69         0.58         0.56           0.85         0.80         0.62	LLaVA-7B         I           Normal         Blur         B&W         Normal           0.55         0.57         0.51         0.58           0.56         0.59         0.55         0.60           0.64         0.56         0.55         0.63           0.69         0.58         0.56         0.67           0.85         0.80         0.62         0.82	LLaVA-7B         LLaVA-           Normal         Blur         B&W         Normal         Blur           0.55         0.57         0.51         0.58         0.57           0.56         0.59         0.55         0.60         0.57           0.64         0.56         0.55         0.63         0.58           0.69         0.58         0.56         0.67         0.61           0.85         0.80         0.62         0.82         0.79	LLaVA-7B         LLaVA-13B           Normal         Blur         B&W         Normal         Blur         B&W           0.55         0.57         0.51         0.58         0.57         0.52           0.56         0.59         0.55         0.60         0.57         0.54           0.64         0.56         0.55         0.63         0.58         0.56           0.69         0.58         0.56         0.67         0.61         0.56           0.85         0.80         0.62         0.82         0.79         0.65	Normal         Blur         B&W         Normal         Blur         B&W         Normal           0.55         0.57         0.51         0.58         0.57         0.52         0.51           0.56         0.59         0.55         0.60         0.57         0.54         0.58           0.64         0.56         0.55         0.63         0.58         0.56         0.57           0.69         0.58         0.56         0.67         0.61         0.56         0.57           0.85         0.80         0.62         0.82         0.79         0.65         0.67	Normal         Blur         B&W         Normal         Blur         B         Normal         Blur         B         W         Normal         Blur         B         W         Normal         B         B         Normal         S         Normal         B         Normal         S         Normal         B         Normal         B         Normal         S         Normal         S         Normal         S         0.50         0.50         0.50         0.50         0.50         0.50         0.52         0.50         0.51         0.52         0.51         0.52         0.51         0.50         0.51 </td

Table 5: Results of the REI task on the CLEVR dataset. LLaVA-7B is the speaker. ZSL means that no training was involved. Normal: no perceptual impairment, Blur: Blurry vision, B&W: Vision with no color.

LLaVA-7B			LLaVA-13B				Qwen		
LLaVA-7B	Normal	Blur	B&W	Normal	Blur	B&W	Normal	Blur	B&W
ZSL	0.83	0.75	0.72	0.86	0.80	0.78	0.63	0.59	0.54
NLPO	0.87	0.81	0.73	0.89	0.83	0.80	0.70	0.64	0.60
PPO	0.87	0.79	0.76	0.90	0.83	0.79	0.70	0.66	0.61
KTO	0.87	0.79	0.75	0.88	0.80	0.80	0.69	0.64	0.61
NLPO PPO KTO	0.87 0.87 0.87	0.81 0.79 0.79	0.73 0.76 0.75	0.89 0.90 0.88	0.83 0.83 0.80	0.80 0.79 0.80	0.70 0.70 0.69	0.64 0.66 0.64	0.60 0.61 0.61

Table 6: Results of the REI task on the CUB dataset. LLaVA-7B is the speaker. ZSL means that no training was involved. Normal: no perceptual impairment, Blur: Blurry vision, B&W: Vision with no color.

	L	LaVA-	7B	Ι	LLaVA-	13B		Qwen	
LLaVA-7B	Normal	Blur	B&W	Normal	Blur	B&W	Normal	Blur	B&W
ZSL	0.85	0.81	0.82	0.87	0.84	0.84	0.73	0.72	0.73
NLPO	0.86	0.82	0.83	0.87	0.84	0.83	0.81	0.76	0.77
PPO	0.85	0.82	0.83	0.87	0.84	0.85	0.81	0.75	0.80
KTO	0.89	0.84	0.88	0.92	0.94	0.96	0.81	0.75	0.76

Table 7: Results of the REI task on the ImageNet dataset. LLaVA-7B is the speaker. ZSL means that no training was involved. Normal: no perceptual impairment, Blur: Blurry vision, B&W: Vision with no color. 

	Pal	iGemm	a
LLaVA-7B	Normal	Blur	B&W
GT	0.63	0.56	0.61
ZSL	0.34	0.28	0.28
NLPO	0.34	0.29	0.29
PPO	0.34	0.28	0.29
KTO	0.44	0.41	0.40
	LLaVA-7B GT ZSL NLPO PPO KTO	Pal           LLaVA-7B         Normal           GT         0.63           ZSL         0.34           NLPO         0.34           PPO         0.34           KTO         0.44	PaliGemm           LLaVA-7B         Normal         Blur           GT         0.63         0.56           ZSL         0.34         0.28           NLPO         0.34         0.29           PPO         0.34         0.28           KTO         0.44         0.41

Table 8: Results of the RES task on the RefCOCO dataset as IoU with target segmentations. GT refers to providing human annotated referring expressions. For other experiments LLaVA-7B is the speaker. ZSL means that no training was involved. Normal: no perceptual impairment, Blur: Blurry vision, B&W: Vision with no color.