# CLIP-UP: CLIP-Based Unanswerable Problem Detection
# for Visual Question Answering

## Supplementary Material

## A. Implementation details

### A.1. CLIP-UP-Emb

We begin by describing the setup used to train CLIP-UP-Emb. Unless otherwise specified, the training settings are identical across both multiple-choice and open-ended tasks, and across all evaluated VLMs.

We use the AdamW optimizer [41] with a weight decay of 0.0001, a cosine learning rate schedule [40], and a total of 3 training epochs. The learning rate starts at 0.0625 and decays to 0 over the course of training. We use a batch size of 8 for LLaVA-1.5-7B, Phi-3.5-Vision and InternVL3-1B, and a batch size of 4 for LLaVA-NeXT-13B, Ovis2-16B, and InternVL3-8B. Gradient checkpointing is applied in all settings to reduce GPU memory usage. For Phi-3.5-Vision, we additionally apply gradient clipping to a maximum norm of 0.5.

As we observed that training produces separability between embedding vectors from answerable and unanswerable inputs (see Fig. 4), we begin multiple-choice training with one warm-up epoch using supervised contrastive loss [28], aiming to separate "answerable" and "unanswerable" projected embeddings. This stage uses a batch size of 128, a temperature of 0.07, and a constant learning rate of 0.0005. Open-ended training skips the warm-up phase.

LLaVA models use float16 precision, while other models use bfloat16 precision. Phi-3.5-Vision, InternVL3-1B/8B and Ovis2-16B were set to process up to 4, 6, and 1 image crops, respectively.

We generate the correlation vectors using Structure-CLIP [25]. Its embedding dimension is 768, resulting in correlation vectors with a total dimension of 3072 for multiple-choice VQA (with four concatenated correlation vectors) and 768 for open-ended VQA (with a single correlation vector). The learned linear projection layer includes a bias term and operates in bfloat16 precision.

For multiple-choice VQA, training and inference are conducted using the base prompt setting, where the model receives the VQA inputs without additional instructions. For open-ended questions, the instruction "Answer the question using a single word or phrase" is included. All experiments are performed using greedy decoding.

Training was done on two NVIDIA GeForce RTX 3090 GPUs for LLaVA-1.5-7B, Phi-3.5-Vision, Ovis2-16B and InternVL3 models, and on a single NVIDIA RTX A6000 GPU for LLaVA-NeXT-13B.

### A.2. CLIP-UP-EmbLoRA

We describe the CLIP-UP-EmbLoRA setting. CLIP-UP-EmbLoRA combines both embedding injection and InjLoRA fine-tuning. For training the embedding injection, we use the same CLIP-UP-Emb configuration described above.

For InjLoRA, projections ($\mathcal{P}'$ in Eq. (3)) are trained using the same settings, and LoRA fine-tuning is applied to all linear layers in the LLM component of the VLM. For all the LoRA layers and the learnable residual matrices ($C$ in Eq. (3)), we use a cosine learning rate schedule that decays to 0, with a linear warmup over the first 3% of training steps. Learning rate starts at $1 \times 10^{-5}$ for LLaVA models and Phi-3.5-Vision, $4 \times 10^{-5}$ for InternVL3 models, and $1 \times 10^{-4}$ for Ovis2-16B. For multiple-choice VQA, we use a LoRA rank of 8 and a LoRA alpha of 16, while for open-ended VQA, we use a rank of 32 and an alpha of 64.

We note that since InjLoRA introduces signal-dependent modifications to the LoRA layers, the LoRA weights cannot be merged into the base model at inference time.

### A.3. Structure-CLIP

We generate the correlation vectors using Structure-CLIP [25]. As Structure-CLIP's weights are not published, we fine-tune CLIP ViT-L/14@336 [51] to replicate it.

Fine-tuning is performed on CLIP ViT-L/14@336px [51] for one epoch on a single NVIDIA A100 GPU, over the MS COCO dataset [34] with augmentations by [25]. To reduce memory usage, we freeze the first 9 transformer blocks of the image encoder and the first 21 transformer blocks of the text encoder. The Knowledge-Enhanced Encoder (KEE) component is fine-tuned following the procedure in [25]. We use a learning rate of $3 \times 10^{-6}$, a batch size of 16, a weight decay of 0.1, and a KEE Knowledge weight of 0.2. In inference, we use the fine-tuned image and text encoders of Structure-CLIP without the additional KEE.

### A.4. LoRA fine-tuning baseline

In the main paper, we compare CLIP-UP methods to LoRA fine-tuning [23]. For multiple-choice UPD on LLaVA models, we followed the setup proposed in [45]: we used their published LLaVA-NeXT-13B weights, and for LLaVA-1.5-7B, which was not evaluated in their work, we reproduced the fine-tuning process using their LLaVA-NeXT-13B settings and training data.

Phi-3.5-Vision, InternVL3 models and Ovis2-16B were

neither fine-tuned in [45], so we followed their recommended LoRA training recipe: for Phi-3.5-Vision, a learning rate of $2 \times 10^{-4}$, batch size of 64, LoRA rank of 32, and LoRA alpha of 16; for InternVL3 models, a learning rate of $4 \times 10^{-5}$, batch size of 64, LoRA rank of 16, and LoRA alpha of 32; for Ovis2-16B, a learning rate of $1 \times 10^{-4}$, batch size of 4, LoRA rank of 32, and LoRA alpha of 64.

For open-ended VQA, we applied the multiple-choice recipes with our own training data.

## B. Training datasets

### B.1. Multiple-choice training dataset

We provide details about the dataset we created for training CLIP-UP on multiple-choice VQA. The goal was to create a compact high-quality UPD training dataset. We do not use the fine-tuning dataset from [45] as it is too large (10,000 samples), lacks IASD samples, and, upon our manual inspection, found to be of insufficient quality.

The dataset is organized into multiple-choice VQA question pairs, each consisting of an answerable question and its corresponding unanswerable variant. The training set contains 263, 159, and 277 question pairs for AAD, IASD, and IVQD, respectively (a total of 526, 318, and 554 samples). The validation set contains 30 pairs for each category. We do not include a test set.

Unlike the training set, each question in the validation set is augmented with $n$ repetitions ($n$ is the number of options), each with a different circular shift of the options, enabling dual accuracy evaluation. Consequently, the validation set contains a total of 204, 232, and 226 questions for AAD, IASD, and IVQD, respectively.

Data were created with a different process for each unanswerability category, as we explain below. All data were sourced from public training sets to ensure no leakage with public benchmarks and test sets. For all categories, questions were generated with four options. Most questions were left unchanged, but some were modified to include fewer options. We also ensured that the correct option varies (*e.g.*, it is not always "A").

Note that the dataset was constructed in a straightforward manner, resulting in structurally uniform questions, as we assumed this would suffice for training CLIP-UP. This simplicity highlights CLIP-UP's robustness and suggests that a more diverse dataset could further improve performance.

#### B.1.1. AAD data

The AAD data consist of 293 pairs of questions: 143 sourced from the A-OKVQA dataset [54], and 150 generated using GPT-4o mini [48] based on MS COCO [34].

Our goal is to have standard questions with exactly one correct answer option, while all others are clearly incorrect. This ensures that AAD unanswerable questions may be generated by removing the correct answer option, leaving no valid answer in the answer options set. Note that this condition is not always met, as many multiple-choice questions are intentionally designed to be challenging, requiring the selection of the best option from several plausible ones.

We began by creating the standard questions, selecting 143 multiple-choice VQA items from the A-OKVQA training dataset [54]. We manually examined the data to include only questions with exactly one correct answer option.

We created 150 additional standard questions using the following process: we first sampled examples from MS COCO training set [34] (2017 split). Each sample consists of an image and five ground truth captions, from which we randomly selected one. Next, we used GPT-4o mini [48] to generate three incorrect captions for each sample. GPT-4o mini was given an image and its correct caption, and instructed to output a multiple-choice VQA question asking to select the correct caption, with four answer options: a correct one (the ground truth caption) and three incorrect ones (generated by GPT-4o mini). See the instruction used in Fig. 6a. To diversify the data, we alternated between two question formats: "Which caption describes the image?" and "Which one is the correct caption for this image?". As with the A-OKVQA questions, we included only standard questions with exactly one correct answer option.

After obtaining 293 standard multiple-choice VQA questions from both sources, we created the AAD counterparts by removing the correct answer option from each standard question. See Fig. 5a for an example.

#### B.1.2. IASD data

The IASD data consist of 189 pairs of questions. In the case of IASD, there are no specific constraints on the standard questions. However, for unanswerable questions, the textual question (the question itself, *e.g.*, "What color is the dress?") and the answer options set must be incompatible.

Similar to the AAD case, we used standard questions from the A-OKVQA training dataset [54] and ones generated with GPT-4o mini [48]. To create the unanswerable counterpart for each standard question, the original answer set was replaced with one from another randomly selected standard question. We then manually examined the data to include only pairs where the textual question is genuinely incompatible with the unanswerable answer options set. See Fig. 5b for an example.

#### B.1.3. IVQD data

The IVQD data consist of 307 pairs of questions: 42 sourced from the fine-tuning data by [45], and 265 generated using GPT-4o mini [48] based on MS COCO [34] and TextCaps [55].

Our goal is to have pairs of multiple-choice VQA questions where the textual question conveys some specific information about the image. This allows generating unanswerable IVQD questions by replacing the image with an-

Figure 5. Pairs of standard and unanswerable multiple-choice VQA questions from our multiple-choice dataset for (a) AAD, (b) IASD, and (c) IVQD.

other image that is incompatible with the information in the textual question (in contrast, non-specific questions like "What emotion does this image convey?" are compatible with most images).

The 42 pairs sourced from the fine-tuning data by [45] include corresponding standard and IVQD unanswerable

You are an assistant with the task of creating multiple-choice questions about images. You will be given an image, and its correct caption. The correct caption is the correct answer to the question "Which one is the correct caption of this image?".

Your job is to create 3 distractors that are incorrect captions for the image. Note that the distractors must be incorrect. This means that if we will take off the correct option, there will be no correct distractor that might describe the image.

The output should be in the form of a python dictionary, with 6 entries: "question" containing the question, "image_id" containing an image id as integer (that will be given as input), "A" containing the correct caption, and "B", "C", "D" containing (each) the 3 distractors.

Here is an output for example: {"question": "Which one is the correct caption of this image?", "image_id": 57703, "A": "A man and two women walking their dogs and hiking in the woods.", "B": "A group of people camping near a lake with their pets.", "C": "Three hikers climbing a mountain trail with no animals in sight.", "D": "Two women and a child having a picnic in a grassy field."}

(a)

You are an assistant with the task of creating a "specific" question about an image. You will be given a caption of an image (without the image itself), and you should phrase a question that can be answered using the information in this caption. The question must be phrased so it delivers some information about the image, thus it will not be relevant for any image. In addition, the information in the caption must be necessary to answer the question. You may deliver only some information about the caption, and not all of it, use your judgment. Please try to output long answers when possible.

The output should be in the form of a python dictionary, with 3 entries: "image_id" containing an image id as integer (that will be given as input), "question" containing the question, and "answer" containing the answer.

For you to understand, here are some examples. Each example contains input and output, an additional undesired output with an explanation:

Example 1:
Input: {"image_id": 32677, "caption": "A dog and a cat sleeping next to each other."}
Output: {"image_id": 32677, "question": "What animals are sleeping in the image?", "answer": "A dog and a cat."}
Undesired output: {"image_id": 32677, "question": "What is in the image?", "answer": "A dog and a cat."}
Explanation: "What is in the image?" may be applied for any image, and thus it is an undesired question.

Example 2:
Input: {"image_id": 32678, "caption": "A yellow happy emoji."}
Output: {"image_id": 32678, "question": "What emotion does this emoji express?", "answer": "Happiness."}
Undesired output: {"image_id": 32678, "question": "What emotion does this image express?", "answer": "Happiness."}
Explanation: Mentioning a specific object, emoji, implies that there must be an emoji in the image. On the other end, "What emotion does this image express?" may be applied for any image (one may say any image conveys some emotion).

Example 3:
Input: {"image_id": 34512, "caption": "An image of the Empire State Building."}
Output: {"image_id": 34512, "question": "What is the name of the building in the image?", "answer": "The Empire State Building."}
Undesired output: {"image_id": 34512, "question": "What place is it in the image?", "answer": "The Empire State Building."}
Explanation: Mentioning a specific object, building, implies that there must be a building in the image. On the other end, "What place is it in the image?" may be applied for almost any image.

(b)

Figure 6. The instructions given to GPT-4o mini for (a) generating incorrect answer options for AAD multiple-choice questions and (b) generating image-specific questions for IVQD multiple-choice questions.

questions. We manually ensured that in all pairs, the textual question conveys image-specific information and is genuinely incompatible with the image in the unanswerable item.

For the 265 other question pairs, we generated standard questions using the following process: similar to the AAD case, we sampled examples from MS COCO training set [34] (2017 split), but also from TextCaps training set [55]. Each sample consists of an image and five ground truth captions, from which we randomly selected one. Next, we used GPT-4o mini [48] to generate an image-specific textual question from each caption. GPT-4o mini was given a caption (without the image) and instructed to output an image-specific textual question related to the caption along with the correct answer. See Fig. 6b for the instruction used. Then, we used GPT-4o mini to create three incorrect answer options for each question by providing it with the image, question and correct answer as input, and instructing it similarly to the AAD case.

To create the unanswerable counterpart for each standard question, we replaced the image with one from another randomly selected standard question. The data were manually reviewed to include only pairs where the textual question is image-specific and genuinely incompatible with unanswerable IVQD image. See Fig. 5c for an example.

### B.2. Open-ended training dataset

The open-ended VQA training set consists of 700 corresponding answerable-unanswerable question pairs sampled from the TDIUC training set [27], with 50 pairs allocated for validation. Unanswerable questions were first drawn from the "Absurd" category. Each question was then paired with an answerable question from a valid (non-absurd) category about the same image.

## C. Evaluation details

### C.1. Multiple-choice baselines

We provide details on the three prompt engineering settings from [45]: (1) Base Prompt Setting: uses only the multiple-choice VQA prompt without additional instructions. As this setting does not explicitly encourage choosing an answer, it identifies unanswerable questions better than the original setup; (2) Additional-Option Setting: adds an option depending on the unanswerability category ("None of the above" for AAD and IASD, or "The image and question are irrelevant" for IVQD), and includes the original instruction. This setting ensures that a correct answer is always present and encourages the model to select one; (3) Additional-Instruction Setting: adds an instruction to encourage withholding an answer when appropriate. The instructions vary by the unanswerability category and are similar to the extra option in the Additional-Option Setting.

Note that settings (2) and (3) assume knowledge of the input's unanswerability category, which is not the case in real-world scenarios. They are thus meant to test models' capabilities via prompt engineering rather than serve as practical solutions.

### C.2. Multiple-choice evaluation

We conducted all the multiple-choice UPD evaluations ourselves, including those of the prompt engineering methods. For all experiments that were also performed by Miyai et al. [45], our results closely align with theirs.

Multiple-choice UPD evaluation requires extracting the selected option from the model's prediction. We followed the extraction approach described in [45]: each VLM prediction is first processed using a string matching algorithm, and if this fails, GPT-3.5 (gpt-3.5-turbo-0125 [47]) is employed with a tailored prompt to extract the selected option. We introduced slight modifications to the string matching algorithm to improve efficiency and accuracy, and reduce calls to GPT-3.5. To ensure a fair comparison, all results were evaluated using our modified string matching extraction algorithm.

### C.3. Open-ended evaluation

We evaluate on two open-ended VQA test datasets. The first is sampled from the RGQA benchmark [66] and consists of 2,500 pairs of corresponding answerable-unanswerable questions. Pairs were randomly drawn from all four RGQA subsets (CLIP-easy, PT-easy, CLIP-hard, and PT-hard), with 625 pairs from each subset. Each pair shares the same image but contains a different question.

The second dataset consists of 2,500 pairs of corresponding answerable-unanswerable questions sampled from the TDIUC benchmark [27]. The data was collected similarly to the open-ended training set, with the key difference that here the samples were drawn from the TDIUC test set.

We adopt the LAVE evaluation metric [43], which leverages an LLM for open-ended VQA scoring and has been shown to align better with human judgment than alternative metrics. We use GPT-3.5 (gpt-3.5-turbo-0125 [47]) as the LLM.

## D. Parsing and classification of question prompts

This section describes the rule-based algorithm mentioned in the main paper. The algorithm serves two purposes: first, to classify whether a textual input is a multiple-choice question, an open-ended question, or neither. This classification determines whether, and which, correlation vector should be generated and integrated into the VLM. It can be applied to InjLoRA, and even to standard LoRA fine-tuned models, to decide whether LoRA weights should be used (as long as LoRA weights are not merged into the base model).

| Method | AAD | | | IASD | | | IVQD | | | Average | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Stand. | UPD | Dual | Stand. | UPD | Dual | Stand. | UPD | Dual | Stand. | UPD | Dual |
| Base Setting | 69.02 | 1.71 | 1.59 | 66.49 | 19.70 | 12.73 | 63.48 | 0.28 | 0.28 | 66.33 | 7.23 | 4.87 |
| CLIP-UP-Emb ($\alpha = 0.0$) | 69.02 | 1.46 | 1.34 | 66.27 | 19.59 | 12.84 | 63.76 | 0.56 | 0.56 | 66.35 | 7.20 | 4.91 |
| CLIP-UP-Emb ($\alpha = 0.1$) | 68.90 | 1.59 | 1.46 | 66.70 | 19.48 | 12.95 | 63.48 | 0.28 | 0.28 | 66.36 | 7.12 | 4.90 |
| CLIP-UP-Emb ($\alpha = 0.2$) | 69.15 | 1.46 | 1.46 | 66.59 | 20.13 | 13.28 | 62.92 | 0.28 | 0.28 | 66.22 | 7.29 | 5.01 |
| CLIP-UP-Emb ($\alpha = 0.3$) | 69.15 | 1.59 | 1.46 | 66.92 | 20.35 | 13.49 | 63.20 | 0.56 | 0.28 | 66.42 | 7.50 | 5.08 |
| CLIP-UP-Emb ($\alpha = 0.4$) | 69.27 | 1.22 | 1.22 | 66.81 | 19.15 | 12.51 | 64.04 | 0.56 | 0.28 | 66.71 | 6.98 | 4.67 |
| CLIP-UP-Emb ($\alpha = 0.5$) | 69.39 | 0.73 | 0.73 | 67.03 | 12.19 | 7.94 | 64.04 | 0.56 | 0.28 | 66.82 | 4.49 | 2.98 |
| CLIP-UP-Emb ($\alpha = 0.6$) | 67.68 | 36.95 | 33.17 | 66.16 | 45.05 | 30.90 | 62.64 | 32.58 | 22.19 | 65.49 | 38.19 | 28.75 |
| CLIP-UP-Emb ($\alpha = 0.7$) | 63.78 | 57.93 | 45.00 | 62.46 | 80.41 | 51.80 | 60.67 | 68.82 | 45.51 | 62.30 | 69.05 | 47.44 |
| CLIP-UP-Emb ($\alpha = 0.8$) | 61.71 | 65.00 | 47.07 | 59.74 | 88.25 | 53.75 | 58.99 | 78.65 | 49.44 | 60.15 | 77.30 | 50.09 |
| CLIP-UP-Emb ($\alpha = 0.9$) | 61.46 | 67.07 | 47.68 | 59.09 | 90.75 | 54.19 | 58.15 | 83.15 | 51.40 | 59.57 | 80.32 | 51.09 |
| Original CLIP-UP-Emb ($\alpha = 1.0$) | 61.22 | 67.68 | 47.80 | 59.52 | 90.64 | 54.73 | 58.43 | 82.58 | 51.12 | 59.72 | 80.30 | 51.22 |

Table 4. Standard-dual accuracy trade-off control results on LLaVA-1.5-7B.

Second, if the input is a multiple-choice question, the algorithm parses it to separate the textual question and answer options, a step necessary for generating the correlation vectors.

The algorithm relies on simple string matching and assumes a specific structure of multiple-choice question prompts: a question followed by answer options, each preceded by a letter (*e.g.*, "A"). For example, "What animal is by the flowers? A. Dog B. Rabbit C. Cat."

In the first step, the algorithm checks whether the input is a multiple-choice question by detecting for the presence of "A." and "B." (since a question must have at least two options). If these are present in the input, the algorithm proceeds to the next step, where it parses the input: the question is the text before "A.", the first answer option is the text between "A." and "B.", and so on for the remaining answer options. If the input is not classified as multiple-choice, the algorithm simply checks for the presence of a question mark to determine whether it is an open-ended question (or not a question at all).

Since the algorithm relies on string matching, it can be easily adjusted to support different multiple-choice input formats (*e.g.*, options denoted with numbers instead of letters). Moreover, the algorithm could easily be replaced with a more sophisticated approach, such as leveraging the LLM component of the VLM for more robust detection. We however found it unnecessary given the simplicity of the parsing task on our test data.

## E. Standard-dual accuracy trade-off control

Although CLIP-UP enhances VLMs' UPD capabilities, reflected in improved dual accuracy, this comes at the cost of reduced standard accuracy, introducing a trade-off between the two. In some applications, particularly those requiring high reliability (*e.g.*, medical VQA systems), a steep drop in standard accuracy may be unacceptable. This highlights the need for a mechanism to control this trade-off.

We introduce a simple inference-time method, requiring no retraining, for controlling the trade-off. This is done by interpolating CLIP-UP-Emb's embedding vector ($\mathbf{e}$ from Eq. (2)) with random noise:

$$\mathbf{e} = \alpha\mathbf{e} + (1 - \alpha)\mathbf{e}_{noise}, \tag{4}$$

where $\alpha \in [0, 1]$ controls the interpolation strength and $\mathbf{e}_{noise}$ is sampled from the standard normal distribution.

Tab. 4 shows the MM-UPD results on LLaVA-1.5-7B. First, using only noise ($\alpha = 0.0$) performs similarly to the base setting, thus not altering model behavior, motivating the idea of noise interpolation. As $\alpha$ increases, dual accuracy improves while standard accuracy decreases, reflecting a controllable trade-off. For some $\alpha$ values, adjusting CLIP-UP for higher standard accuracy still preserves strong dual performance. For example, $\alpha = 0.7$ yields a $2.58\%$ gain in standard accuracy with only a $3.78\%$ drop in dual accuracy.

## F. Additional results

### F.1. Full MM-UPD results

Tabs. 5 and 6 present the complete multiple-choice MM-UPD results, including LLaVA-NeXT-13B and InternVL3-1B, as well as results for each unanswerability category (AAD, IASD, and IVQD). Figs. 7 and 8 show VLM responses to answerable and unanswerable questions, with and without CLIP-UP.

### F.2. Additional baselines

We include three additional baselines in Tab. 5. First, we present results of LLaVA-1.5-7B LoRA fine-tuning following [45], using CLIP-UP's multiple-choice training data.

| Method | AAD | | | IASD | | | IVQD | | | Average | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Stand. | UPD | Dual | Stand. | UPD | Dual | Stand. | UPD | Dual | Stand. | UPD | Dual |
| Original Model | 70.24 | 0.00 | 0.00 | 67.79 | 0.33 | 0.33 | 65.17 | 0.00 | 0.00 | 67.73 | 0.11 | 0.11 |
| Base Setting | 69.02 | 1.71 | 1.59 | 66.49 | 19.70 | 12.73 | 63.48 | 0.28 | 0.28 | 66.33 | 7.23 | 4.87 |
| Additional-Option | 68.41 | 48.54 | 40.85 | 65.72 | 78.24 | 51.58 | 63.76 | 26.97 | 21.91 | 65.96 | 51.25 | 38.11 |
| Additional-Instruction | 68.54 | 33.90 | 27.80 | 65.61 | 65.18 | 42.76 | 63.48 | 31.74 | 23.60 | 65.88 | 43.61 | 31.39 |
| Correlation Vectors Classifier | 30.24 | 99.63 | 30.00 | 29.16 | 63.87 | 17.30 | 46.08 | 90.59 | 28.65 | 35.16 | 84.70 | 25.32 |
| LoRA Fine-Tuning | 64.63 | 56.34 | 43.78 | 61.92 | 87.81 | 54.73 | 61.24 | 85.96 | 52.25 | 62.60 | 76.70 | 50.25 |
| LoRA Fine-Tuning (CLIP-UP data) | 66.71 | 51.10 | 40.49 | 63.33 | 88.25 | 55.71 | 61.80 | 67.70 | 44.66 | 63.95 | 69.02 | 46.95 |
| CLIP-UP-Emb (ours) | 61.22 | 67.68 | 47.80 | 59.52 | 90.64 | 54.73 | 58.43 | 82.58 | 51.12 | 59.72 | 80.30 | 51.22 |
| CLIP-UP-EmbLoRA (ours) | 62.44 | 71.71 | **49.27** | 60.17 | 93.47 | **56.58** | 59.55 | 84.83 | **53.37** | 60.72 | 83.34 | **53.07** |

(a) LLaVA-1.5-7B

| Method | AAD | | | IASD | | | IVQD | | | Average | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Stand. | UPD | Dual | Stand. | UPD | Dual | Stand. | UPD | Dual | Stand. | UPD | Dual |
| Original Model | 76.71 | 0.00 | 0.00 | 73.23 | 0.11 | 0.00 | 71.35 | 0.00 | 0.00 | 73.76 | 0.04 | 0.00 |
| Base Setting | 72.32 | 23.78 | 17.80 | 69.75 | 49.62 | 31.66 | 68.82 | 44.66 | 33.15 | 70.30 | 39.35 | 27.54 |
| Additional-Option | 75.85 | 18.41 | 18.05 | 72.47 | 39.28 | 29.92 | 70.79 | 46.35 | 38.20 | 73.04 | 34.68 | 28.72 |
| Additional-Instruction | 67.07 | 48.66 | 38.29 | 63.87 | 87.81 | 57.02 | 68.82 | 71.91 | 54.49 | 66.59 | 69.46 | 49.93 |
| Chain-of-Thought | 60.00 | 60.50 | 42.80 | 56.40 | 70.80 | 43.90 | 59.00 | 75.30 | 47.50 | 58.47 | 68.87 | 44.73 |
| Self-Reflection | 66.20 | 50.00 | 37.80 | 62.60 | 55.80 | 36.70 | 59.80 | 61.50 | 39.00 | 62.87 | 55.77 | 37.83 |
| LoRA Fine-Tuning | 69.15 | 58.54 | 47.56 | 65.51 | 91.19 | 59.85 | 67.42 | 86.24 | **59.55** | 67.36 | 78.66 | 55.65 |
| CLIP-UP-Emb (ours) | 62.07 | 83.90 | 54.02 | 58.54 | 95.65 | 55.71 | 57.87 | 92.70 | 55.06 | 59.49 | 90.75 | 54.93 |
| CLIP-UP-EmbLoRA (ours) | 67.07 | 75.73 | **54.27** | 63.22 | 95.76 | **60.50** | 63.48 | 84.83 | 53.65 | 64.59 | 85.44 | **56.14** |

(b) LLaVA-NeXT-13B

| Method | AAD | | | IASD | | | IVQD | | | Average | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Stand. | UPD | Dual | Stand. | UPD | Dual | Stand. | UPD | Dual | Stand. | UPD | Dual |
| Original Model | 80.73 | 1.22 | 1.22 | 77.48 | 0.00 | 0.00 | 77.53 | 0.00 | 0.00 | 78.58 | 0.41 | 0.41 |
| Base Setting | 79.51 | 1.95 | 1.95 | 76.28 | 0.54 | 0.44 | 75.56 | 0.28 | 0.28 | 77.12 | 0.92 | 0.89 |
| Additional-Option | 80.24 | 23.41 | 21.95 | 77.04 | 31.56 | 24.27 | 75.28 | 64.33 | 51.97 | 77.52 | 39.77 | 32.73 |
| Additional-Instruction | 77.93 | 31.95 | 27.93 | 74.76 | 46.25 | 32.86 | 74.72 | 72.75 | 56.18 | 75.80 | 50.32 | 38.99 |
| LoRA Fine-Tuning | 61.83 | 71.95 | 47.93 | 59.52 | 95.21 | 56.69 | 58.71 | 93.82 | 56.46 | 60.02 | 86.99 | 53.69 |
| CLIP-UP-Emb (ours) | 62.68 | 81.95 | 52.80 | 59.52 | 88.25 | 52.67 | 60.67 | 93.26 | 56.74 | 60.96 | 87.82 | 54.07 |
| CLIP-UP-EmbLoRA (ours) | 72.32 | 77.44 | **60.61** | 69.10 | 91.40 | **63.00** | 69.10 | 91.01 | **63.20** | 70.17 | 86.62 | **62.27** |

(c) Phi-3.5-Vision

Table 5. Full results (%) on MM-UPD [45] multiple-choice VQA for (a) LLaVA-1.5-7B, (b) LLaVA-NeXT-13B, and (c) Phi-3.5-Vision. Metrics include circular standard, UPD, and dual accuracies.

Fine-tuning was conducted under the same settings as the original setup, but with 3 epochs instead of one, for a fair comparison with CLIP-UP methods. This setting achieves reasonable performance but is inferior to the CLIP-UP methods, and to the original LoRA fine-tuning setup that uses more data. This suggests that CLIP-UP methods are more data-efficient.

Second, we tested a simple classifier-based baseline: a logistic regression classifier trained on our correlation vectors determines whether the input is answerable or unan-

swerable. If classified as unanswerable, the output is "I cannot answer"; otherwise, the VLM generates the response. The classifier itself labels only 46.27% of circular standard questions as answerable, which already sets a dual accuracy upper bound that indicates that the baseline underperforms CLIP-UP methods on all VLMs. For completeness, we still evaluate the full baseline, with results for LLaVA-1.5-7B shown in Tab. 5a (line 5). The baseline substantially underperforms CLIP-UP variants, suggesting that injection is necessary and that the correlation vectors alone do not cap-

| Method | AAD | | | IASD | | | IVQD | | | Average | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Stand. | UPD | Dual | Stand. | UPD | Dual | Stand. | UPD | Dual | Stand. | UPD | Dual |
| Original Model | 76.59 | 0.00 | 0.00 | 74.10 | 0.00 | 0.00 | 74.44 | 0.00 | 0.00 | 75.04 | 0.00 | 0.00 |
| Base Setting | 71.22 | 5.85 | 5.37 | 67.03 | 6.53 | 2.50 | 67.98 | 2.53 | 1.69 | 68.74 | 4.97 | 3.19 |
| Additional-Option | 75.49 | 36.83 | 35.73 | 73.01 | 41.35 | 30.79 | 74.44 | 26.40 | 20.79 | 74.31 | 34.86 | 29.10 |
| Additional-Instruction | 74.02 | 0.12 | 0.12 | 71.93 | 1.96 | 1.09 | 72.19 | 0.00 | 0.00 | 72.71 | 0.69 | 0.40 |
| LoRA Fine-Tuning | 69.27 | 62.20 | 51.10 | 67.25 | 85.96 | 58.43 | 66.29 | 80.62 | 53.09 | 67.60 | 76.26 | 54.21 |
| CLIP-UP-Emb (ours) | 64.63 | 64.02 | 46.22 | 60.39 | 76.71 | 47.12 | 61.80 | 89.89 | 56.18 | 62.27 | 76.87 | 49.84 |
| CLIP-UP-EmbLoRA (ours) | 71.95 | 71.95 | **56.10** | 68.99 | 86.40 | **58.54** | 70.22 | 88.76 | **61.52** | 70.39 | 82.37 | **58.72** |

(a) InternVL3-1B

| Method | AAD | | | IASD | | | IVQD | | | Average | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Stand. | UPD | Dual | Stand. | UPD | Dual | Stand. | UPD | Dual | Stand. | UPD | Dual |
| Original Model | 91.10 | 0.00 | 0.00 | 87.27 | 0.00 | 0.00 | 88.48 | 0.00 | 0.00 | 88.95 | 0.00 | 0.00 |
| Base Setting | 88.17 | 37.07 | 35.73 | 84.00 | 52.77 | 45.81 | 86.80 | 43.54 | 36.80 | 86.32 | 44.46 | 39.45 |
| Additional-Option | 90.73 | 51.34 | 50.49 | 86.94 | 73.45 | 63.11 | 88.76 | 74.72 | 66.57 | 88.81 | 66.50 | 60.06 |
| Additional-Instruction | 87.93 | 59.88 | 56.95 | 84.87 | 92.38 | 78.24 | 86.52 | 88.20 | 76.12 | 86.44 | 80.15 | 70.44 |
| LoRA Fine-Tuning | 88.05 | 66.34 | 62.68 | 84.00 | 94.56 | 79.33 | 86.80 | 85.39 | 72.75 | 86.28 | 82.10 | 71.59 |
| CLIP-UP-Emb (ours) | 86.22 | 80.00 | 73.17 | 82.05 | 96.30 | 78.78 | 83.99 | 96.63 | 80.90 | 84.09 | 90.98 | 77.62 |
| CLIP-UP-EmbLoRA (ours) | 87.07 | 80.49 | **74.15** | 83.35 | 98.04 | **81.50** | 87.08 | 94.10 | **82.02** | 85.83 | 90.88 | **79.22** |

(b) InternVL3-8B

| Method | AAD | | | IASD | | | IVQD | | | Average | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Stand. | UPD | Dual | Stand. | UPD | Dual | Stand. | UPD | Dual | Stand. | UPD | Dual |
| Original Model | 89.51 | 16.95 | 16.95 | 86.29 | 12.40 | 11.43 | 89.33 | 0.56 | 0.56 | 88.38 | 9.97 | 9.65 |
| Base Setting | 87.93 | 44.27 | 43.54 | 84.11 | 48.75 | 42.22 | 89.04 | 46.35 | 41.57 | 87.03 | 46.46 | 42.44 |
| Additional-Option | 89.27 | 53.78 | 52.20 | 85.64 | 72.36 | 62.68 | 87.64 | 86.52 | 75.84 | 87.52 | 70.89 | 63.57 |
| Additional-Instruction | 88.05 | 64.76 | 62.56 | 84.77 | 83.57 | 71.06 | 88.20 | 92.13 | 81.46 | 87.01 | 80.15 | 71.69 |
| LoRA Fine-Tuning | 86.71 | 71.59 | 66.83 | 83.57 | 95.76 | 79.98 | 87.36 | 95.51 | **83.15** | 85.88 | 87.62 | 76.65 |
| CLIP-UP-Emb (ours) | 85.12 | 79.63 | 72.56 | 81.39 | 96.95 | 78.89 | 83.99 | 98.31 | 82.30 | 83.50 | 91.63 | 77.92 |
| CLIP-UP-EmbLoRA (ours) | 88.54 | 78.54 | **75.12** | 84.87 | 98.80 | **83.79** | 89.04 | 93.26 | 82.87 | 87.48 | 90.20 | **80.59** |

(c) Ovis2-16B

Table 6. Full results (%) on MM-UPD [45] multiple-choice VQA for (a) InternVL3-1B, (b) InternVL3-8B, and (c) Ovis2-16B. Metrics include circular standard, UPD, and dual accuracies.

ture all relevant information.

Finally, we compare CLIP-UP to two prompt engineering methods proposed in [45]. The first employs zero-shot Chain-of-Thought [29] reasoning by appending the phrase "Let's think step by step" to the multiple-choice VQA prompt, encouraging the model to reason more carefully. The second uses self-reflection [26] by prompting the model to evaluate its own response. Although both techniques improve performance, CLIP-UP methods outperform them by a significant margin. The results for these methods are taken directly from [45] and reported only for LLaVA-NeXT-13B.

### F.3. Additional open-ended results

Tab. 7 presents the RGQA results for LLaVA-NeXT-13B, InternVL3-1B, and Ovis2-16B, not shown in the main pa-per. Fig. 9 shows model responses to RGQA answerable and unanswerable questions, with and without CLIP-UP.

Tab. 8 presents results on open-ended questions from the TDIUC [27] test set. CLIP-UP-EmbLoRA outperforms the baselines on three VLMs, while on most others the performance gap is small (under 0.6%). Notably, accuracies are high for all training-involved methods, as training was performed on samples from the TDIUC training set.

### F.4. Ruling out potential RGQA bias

We rule out the possibility of a CLIP-induced bias in the RGQA evaluation. RGQA [66] comprises four subsets: CLIP-easy, CLIP-hard, PT-easy, and PT-hard, which differ in how unanswerable questions are generated. In CLIP-easy, unanswerable questions are created by pairing texts

| Method | LLaVA-NeXT-13B | | | InternVL3-1B | | |
|---|---|---|---|---|---|---|
| | Stand. | UPD | Dual | Stand. | UPD | Dual |
| Original Model | 73.02 | 12.90 | 9.84 | 64.40 | 11.10 | 7.46 |
| Prompt Engineering | 68.32 | 33.76 | 22.64 | 66.20 | 15.72 | 11.02 |
| LoRA Fine-Tuning | 65.16 | 60.62 | 38.14 | 56.68 | 60.04 | **31.10** |
| CLIP-UP-Emb (ours) | 67.84 | 62.06 | **40.20** | 41.40 | 73.24 | 27.66 |
| CLIP-UP-EmbLoRA (ours) | 55.26 | 75.94 | 39.42 | 57.08 | 56.62 | 29.30 |

Table 7. Results (%) on RGQA [66] open-ended VQA for LLaVA-NeXT-13B and InternVL3-1B. Metrics include standard, UPD, and dual accuracies.

| Method | LLaVA-1.5-7B | | | LLaVA-NeXT-13B | | | Phi-3.5-Vision | | |
|---|---|---|---|---|---|---|---|---|---|
| | Stand. | UPD | Dual | Stand. | UPD | Dual | Stand. | UPD | Dual |
| Original Model | 84.84 | 10.50 | 9.10 | 87.70 | 11.66 | 10.28 | 85.16 | 5.64 | 4.70 |
| Prompt Engineering | 85.04 | 13.46 | 11.34 | 86.18 | 73.46 | 63.44 | 84.44 | 46.86 | 39.88 |
| LoRA Fine-Tuning | 81.20 | 99.82 | 81.06 | 84.62 | 99.96 | 84.58 | 83.66 | 88.54 | 74.04 |
| CLIP-UP-Emb (ours) | 81.42 | 98.46 | 80.24 | 86.68 | 99.76 | 86.48 | 82.38 | 98.98 | 81.64 |
| CLIP-UP-EmbLoRA (ours) | 85.68 | 98.86 | **84.58** | 87.00 | 99.58 | **86.58** | 85.70 | 99.54 | **85.36** |

| Method | InternVL3-1B | | | InternVL3-8B | | | Ovis2-16B | | |
|---|---|---|---|---|---|---|---|---|---|
| | Stand. | UPD | Dual | Stand. | UPD | Dual | Stand. | UPD | Dual |
| Original Model | 86.98 | 6.16 | 5.22 | 89.06 | 14.04 | 12.42 | 88.64 | 7.14 | 6.56 |
| Prompt Engineering | 87.32 | 16.10 | 13.66 | 84.70 | 95.28 | 80.80 | 46.66 | 99.92 | 46.58 |
| LoRA Fine-Tuning | 87.78 | 99.74 | **87.52** | 89.54 | 99.82 | **89.36** | 89.84 | 99.88 | **89.72** |
| CLIP-UP-Emb (ours) | 83.58 | 99.24 | 82.94 | 89.28 | 98.68 | 88.22 | 88.52 | 99.74 | 88.26 |
| CLIP-UP-EmbLoRA (ours) | 87.76 | 99.08 | 86.92 | 89.80 | 99.42 | 89.30 | 87.90 | 99.84 | 87.74 |

Table 8. Results (%) on TDIUC [27] open-ended VQA. Metrics include standard, UPD, and dual accuracies.

and images with low CLIP similarity, while in CLIP-hard unanswerable questions are created from pairs with high CLIP similarity. In contrast, PT-easy and PT-hard are constructed by modifying standard questions through random (PT-easy) or adversarial (PT-hard) word replacements.

Thus, VLMs equipped with CLIP-UP may benefit from questions in the CLIP-easy category. Although CLIP-hard is designed to have the opposite effect, the advantage on CLIP-easy may be more significant, potentially biasing our RGQA results in favor of CLIP-UP.

To examine this potential bias, we report in Tab. 9 the RGQA results on the PT-easy and PT-hard subsets only, which are not influenced by CLIP similarity. We observe that evaluation trends are consistent across the full RGQA test set and the PT-only subsets (compare with Tabs. 2 and 7): CLIP-UP-Emb or CLIP-UP-EmbLoRA outperform LoRA fine-tuning on all models except InternVL3-1B. This rules out the CLIP-easy bias concern.

## F.5. Additional standard VQA evaluation

We evaluated CLIP-UP on standard multiple-choice questions from SEED-Bench [31] to further assess standard ac-

curacy drops. We randomly sampled 1,000 examples from SEED-Bench's image categories and augmented each with circular duplicates, yielding 4,000 samples. Tab. 10 shows that CLIP-UP-EmbLoRA standard accuracy drops are minimal: up to 6.1%, and typically 1-3%.

## F.6. Training and inference times

We report training and inference times on LLaVA-1.5-7B. For CLIP-UP-Emb, training takes 23.4 minutes on multiple-choice data and 19.7 minutes on open-ended data. For CLIP-UP-EmbLoRA, training takes 27.2 minutes for multiple-choice and 22.4 minutes for open-ended data.

CLIP-UP-Emb's impact on inference time is minimal, as the new embedding vector is generated once per input and cached for reuse. On the IVQD sub-benchmark from MM-UPD, CLIP-UP-Emb reduces total inference time (15.1 vs. 17.8 minutes) as it generates shorter responses, despite a 16% increase in per-token time (0.0519 seconds per token vs. 0.0451 seconds). For CLIP-UP-EmbLoRA, inference time increases due to the LoRA weights not being merged into the base model. In this case, inference takes 21.4 minutes with 0.0748 seconds per-token.

| Method | LLaVA-1.5-7B | | | LLaVA-NeXT-13B | | | Phi-3.5-Vision | | |
|---|---|---|---|---|---|---|---|---|---|
| | Stand. | UPD | Dual | Stand. | UPD | Dual | Stand. | UPD | Dual |
| Original Model | 66.64 | 11.40 | 6.92 | 71.56 | 12.68 | 9.28 | 67.96 | 9.36 | 6.48 |
| Prompt Engineering | 66.52 | 15.32 | 9.32 | 66.88 | 35.60 | 22.76 | 68.04 | 28.40 | 19.20 |
| LoRA Fine-Tuning | 54.96 | 63.20 | 26.84 | 64.24 | 59.72 | 34.84 | 66.36 | 39.76 | 24.24 |
| CLIP-UP-Emb (ours) | 36.60 | 76.44 | 23.52 | 66.04 | 59.76 | **35.72** | 58.24 | 57.92 | 27.04 |
| CLIP-UP-EmbLoRA (ours) | 56.36 | 58.76 | **27.96** | 53.16 | 74.04 | 34.72 | 65.32 | 54.72 | **32.72** |

| Method | InternVL3-1B | | | InternVL3-8B | | | Ovis2-16B | | |
|---|---|---|---|---|---|---|---|---|---|
| | Stand. | UPD | Dual | Stand. | UPD | Dual | Stand. | UPD | Dual |
| Original Model | 63.08 | 10.32 | 7.16 | 69.72 | 13.52 | 10.12 | 68.96 | 13.08 | 8.76 |
| Prompt Engineering | 65.20 | 15.08 | 10.28 | 67.00 | 47.12 | 29.40 | 28.84 | 93.92 | 24.40 |
| LoRA Fine-Tuning | 55.16 | 57.32 | **26.36** | 65.56 | 56.80 | 33.44 | 64.64 | 63.16 | 36.52 |
| CLIP-UP-Emb (ours) | 39.32 | 70.64 | 22.32 | 66.40 | 49.40 | 31.24 | 61.24 | 66.00 | 36.44 |
| CLIP-UP-EmbLoRA (ours) | 55.32 | 54.88 | 25.48 | 65.20 | 60.20 | **36.68** | 62.12 | 69.08 | **38.72** |

Table 9. Results (%) on PT-easy and PT-hard categories from RGQA [66] open-ended VQA. Metrics include standard, UPD, and dual accuracies.

| Method | LLaVA-1.5 | LLaVA-Ne | Phi-V-3.5 | InVL3-1B | InVL3-8B | Ovis2-16B |
|---|---|---|---|---|---|---|
| Original Model | 57.30 | 64.60 | 61.70 | 62.50 | 72.50 | 73.50 |
| CLIP-UP-Emb (ours) | 51.10 | 58.50 | 55.00 | 50.00 | 69.00 | 69.20 |
| CLIP-UP-EmbLoRA (ours) | 51.20 | 60.50 | 59.20 | 60.00 | 71.30 | 70.90 |

Table 10. Results (%) on multiple-choice answerable VQA from SEED-Bench [66], reported with circular standard accuracy. Evaluated models are LLaVA-1.5-7B, LLaVA-NeXT-13B, Phi-3.5-Vision, InternVL3-1B, InternVL3-8B, and Ovis2-16B.

## G. Additional ablation studies

Tab. 11 presents the full ablation results, including performance for each unanswerability category and additional experiments. We examine the effect of the training data by evaluating CLIP-UP-Emb when trained only on data from a single challenge (CLIP-UP-Emb-AAD/IASD/IVQD), but tested on all challenges (lines 1–3 in Tab. 11). As expected, each model performs best on the challenge it was trained on, and outperforms CLIP-UP-Emb (trained on all challenges). Each such specific model also shows gains, although limited, on the other challenges. For example, training on AAD data yields reasonable performance for IASD, a point also observed in [45]. We postulate this is because challenges are interrelated. For instance, IASD may be seen as an extreme case of AAD, where the answer options are not only incorrect but also irrelevant to the question.

We also test the effect of using correlation vectors from Structure-CLIP on CLIP-UP-Emb trained with standard LoRA, and on CLIP-UP-EmbLoRA (lines 8–9 and 10–11). Consistent with the ablations presented in the main paper, removing the correlation vectors leads to a performance drop.

## H. Limitations

Although shown to significantly enhance VLMs' UPD performance, CLIP-UP has several limitations to be acknowledged. It depends on the quality of the CLIP signal, which introduces potential shortcomings since CLIP is known to struggle with issues such as attribute binding and spatial reasoning [64]. While we mitigate some of these challenges by using Structure-CLIP, others remain.

To analyze these limitations, we report dual accuracy, true positive rate (TPR), and true negative rate (TNR) on different MM-UPD question categories in Tab. 12, for three VLMs enhanced with CLIP-UP-EmbLoRA. "Positive" refers to the VLM choosing to refrain from answering. Thus, TPR measures how often a VLM withholds answers to unanswerable questions (*i.e.*, responds with "I cannot answer"), while TNR measures how often it chooses to answer answerable questions.

Some categories are more difficult, such as Relation Reasoning and Logical Reasoning, yielding lower dual accuracies. These are also categories where Structure-CLIP is expected to struggle, leading to lower TNRs (*i.e.*, over-abstention). However, stronger models (InternVL3-8B and Ovis2-16B) show more stable TNRs on these categories,

| Method | AAD | | | IASD | | | IVQD | | | Average | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Stand. | UPD | Dual | Stand. | UPD | Dual | Stand. | UPD | Dual | Stand. | UPD | Dual |
| CLIP-UP-Emb-AAD | 63.78 | 87.44 | **59.15** | 61.48 | 71.16 | 44.07 | 55.06 | 53.93 | 32.58 | 60.11 | 70.84 | 45.27 |
| CLIP-UP-Emb-IASD | 63.41 | 48.66 | 34.76 | 60.72 | 92.82 | **56.15** | 60.11 | 62.92 | 41.29 | 61.41 | 68.13 | 44.07 |
| CLIP-UP-Emb-IVQD | 56.10 | 23.05 | 9.27 | 53.21 | 51.69 | 25.90 | 58.15 | 88.20 | **51.97** | 55.82 | 54.31 | 29.05 |
| CLIP-UP-Emb w/ const. signal | 59.27 | 59.39 | 43.05 | 56.69 | 88.03 | 49.51 | 54.49 | 78.09 | 43.26 | 56.82 | 75.17 | 45.27 |
| CLIP-UP-Emb w/ similarities | 62.32 | 55.61 | 42.44 | 59.85 | 90.64 | 54.62 | 58.99 | 75.28 | 46.63 | 60.39 | 73.84 | 47.90 |
| CLIP-UP-Emb w/ CLIP ViT-L/14 | 50.73 | 65.12 | 35.00 | 47.01 | 89.88 | 42.66 | 46.07 | 87.92 | 39.33 | 47.94 | 80.97 | 39.00 |
| CLIP-UP-Emb (ours) | 61.22 | 67.68 | 47.80 | 59.52 | 90.64 | 54.73 | 58.43 | 82.58 | 51.12 | 59.72 | 80.30 | **51.22** |
| CLIP-UP-Emb (const. signal) + LoRA | 61.46 | 58.66 | 37.56 | 58.11 | 93.47 | 53.54 | 56.74 | 86.24 | 47.75 | 58.77 | 79.46 | 46.28 |
| CLIP-UP-Emb + LoRA | 61.22 | 72.56 | **49.27** | 59.30 | 92.38 | **55.17** | 58.99 | 85.11 | **52.81** | 59.84 | 83.35 | **52.42** |
| CLIP-UP-EmbLoRA (const. signals) | 62.93 | 59.88 | 41.10 | 59.30 | 93.80 | 55.39 | 58.15 | 88.48 | 50.56 | 60.13 | 80.72 | 49.02 |
| CLIP-UP-EmbLoRA (ours) | 62.44 | 71.71 | **49.27** | 60.17 | 93.47 | **56.58** | 59.55 | 84.83 | **53.37** | 60.72 | 83.34 | **53.07** |

Table 11. Full ablation results (%) on LLaVA-1.5-7B. Best dual accuracies for each CLIP-UP setting are bolded.

| Category | LLaVA-1.5-7B | | | InternVL3-8B | | | Ovis2-16B | | |
|---|---|---|---|---|---|---|---|---|---|
| | Dual | TPR | TNR | Dual | TPR | TNR | Dual | TPR | TNR |
| Coarse Perception | 75.37 | 93.87 | 94.08 | 89.08 | 96.51 | 97.50 | 86.08 | 95.59 | 98.70 |
| Attribute Reasoning | 52.68 | 93.87 | 78.39 | 79.87 | 95.51 | 94.51 | 81.88 | 96.53 | 95.60 |
| Fine-grained Perception (Instance-Level) | 55.07 | 90.49 | 88.57 | 80.18 | 92.68 | 97.48 | 80.48 | 93.52 | 98.70 |
| Fine-grained Perception (Cross-Instance) | 47.75 | 80.56 | 92.54 | 74.32 | 91.52 | 96.01 | 82.43 | 93.42 | 98.46 |
| Relation Reasoning | 38.67 | 83.50 | 87.98 | 76.67 | 94.87 | 92.19 | 79.33 | 93.20 | 99.31 |
| Logic Reasoning | 12.93 | 91.28 | 60.48 | 47.62 | 88.94 | 87.01 | 55.78 | 83.40 | 94.25 |
| All Categories | 53.17 | 89.78 | 86.87 | 78.71 | 93.90 | 95.47 | 80.24 | 93.68 | 98.05 |

Table 12. Dual accuracy, true positive rate, and true negative rate (all in %) for LLaVA-1.5-7B, InternVL3-8B, and Ovis2-16B enhanced with CLIP-UP-EmbLoRA, evaluated across different MM-UPD [45] categories.

suggesting they are less sensitive to such limitations. Fig. 10 shows an example from the Logical Reasoning category where Structure-CLIP is less indicative and CLIP-UP is limited.

**Finer-grained unanswerability** To investigate this point further, we evaluate CLIP-UP on three subsets from the multiple-choice TUBench benchmark [22]: UVQA, UCR, and UTabMWP. Unanswerable questions in TUBench were created by applying fine-grained edits to the text, keeping the questions grounded in the image but without a correct answer, with some unanswerable due to missing or indeterminate information. Specifically, UVQA includes 250 answerable-unanswerable question pairs about natural images that require nuanced reasoning for answerability assessment.

Tab. 13 presents the results on the UVQA subset. The original models and prompt engineering baselines exhibit low abstention rates, indicating that detecting unanswerable questions is particularly difficult on this benchmark. LoRA fine-tuning mostly has minimal effect on model behavior, while CLIP-UP variants generally improve dual accuracy.

However, the gains are smaller, less consistent, and often come at the cost of a steep drop in standard accuracy.

This is somewhat expected, given the fine-grained nature of UVQA answerability and the questions being in a yes/no format (converted to multiple-choice). Still, CLIP-UP produces gains, for example, standard drops in CLIP-UP on InternVL3-8B and Ovis2-16B are reasonable.

We additionally test on two harder subsets from TUBench: 480 yes/no question pairs from the UCR subset focused on code snippets, and 108 multiple-choice question pairs from the UTabMWP subset focused on tabular data. Tab. 14 shows the results on LLaVA-1.5-7B. Although CLIP-UP generally outperforms other methods, the gains are smaller.

**Further discussion** However present, these limitations are not inherent to CLIP-UP itself, as Structure-CLIP can be replaced with another CLIP variant or even a non-CLIP model that produces shared vision-language embeddings. Thus, given an alignment model that provides a suitable signal, CLIP-UP could potentially support better handling of more fine-grained types of unanswerability.

| Method | LLaVA-1.5-7B | | | LLaVA-NeXT-13B | | | Phi-3.5-Vision | | |
|---|---|---|---|---|---|---|---|---|---|
| | Stand. | UPD | Dual | Stand. | UPD | Dual | Stand. | UPD | Dual |
| Original Model | 74.40 | 0.00 | 0.00 | 67.60 | 0.00 | 0.00 | 61.20 | 0.00 | 0.00 |
| Base Setting | 73.20 | 0.00 | 0.00 | 69.60 | 1.20 | 1.20 | 61.60 | 0.00 | 0.00 |
| Additional-Option | 53.60 | 0.00 | 0.00 | 68.40 | 0.00 | 0.00 | 62.00 | 0.00 | 0.00 |
| Additional-Instruction | 70.40 | 2.80 | 1.60 | 73.20 | 1.20 | 1.20 | 62.00 | 8.00 | 6.00 |
| LoRA Fine-Tuning | 75.60 | 2.40 | 2.40 | 65.20 | 0.00 | 0.00 | 65.20 | 32.40 | 21.20 |
| CLIP-UP-Emb (ours) | 44.80 | 74.80 | **31.60** | 44.00 | 64.40 | **25.60** | 38.00 | 74.80 | **26.00** |
| CLIP-UP-EmbLoRA (ours) | 46.00 | 70.40 | **31.60** | 59.20 | 7.20 | 3.60 | 49.60 | 54.80 | 23.20 |

| Method | InternVL3-1B | | | InternVL3-8B | | | Ovis2-16B | | |
|---|---|---|---|---|---|---|---|---|---|
| | Stand. | UPD | Dual | Stand. | UPD | Dual | Stand. | UPD | Dual |
| Original Model | 70.00 | 0.00 | 0.00 | 80.80 | 0.00 | 0.00 | 76.40 | 0.00 | 0.00 |
| Base Setting | 71.60 | 0.00 | 0.00 | 78.80 | 0.00 | 0.00 | 76.40 | 5.60 | 3.20 |
| Additional-Option | 72.00 | 0.00 | 0.00 | 82.00 | 0.00 | 0.00 | 76.80 | 0.00 | 0.00 |
| Additional-Instruction | 72.80 | 0.00 | 0.00 | 79.20 | 0.00 | 0.00 | 76.00 | 16.00 | 11.60 |
| LoRA Fine-Tuning | 70.00 | 2.00 | 1.60 | 79.20 | 1.20 | 1.20 | 78.80 | 3.60 | 2.80 |
| CLIP-UP-Emb (ours) | 54.80 | 30.00 | 17.60 | 79.20 | 27.20 | **20.00** | 72.40 | 52.80 | **38.40** |
| CLIP-UP-EmbLoRA (ours) | 59.60 | 46.40 | **26.80** | 76.40 | 5.20 | 4.40 | 74.80 | 8.40 | 5.20 |

Table 13. Results (%) on UVQA [22] multiple-choice VQA subset. Metrics include standard, UPD, and dual accuracies.

| Method | UCR | | | UTabMWP | | |
|---|---|---|---|---|---|---|
| | Stand. | UPD | Dual | Stand. | UPD | Dual |
| Original Model | 49.07 | 0.00 | 0.00 | 49.50 | 0.00 | 0.00 |
| Base Setting | 50.93 | 0.00 | 0.00 | 44.50 | 1.00 | 0.50 |
| Additional-Option | 52.78 | 0.00 | 0.00 | 44.00 | 0.00 | 0.00 |
| Additional-Instruction | 47.22 | 0.00 | 0.00 | 46.00 | 6.50 | 0.50 |
| LoRA Fine-Tuning | 50.93 | 0.00 | 0.00 | 40.00 | 22.00 | 1.00 |
| CLIP-UP-Emb (ours) | 35.19 | 20.37 | 9.26 | 17.50 | 43.00 | 0.00 |
| CLIP-UP-EmbLoRA (ours) | 22.22 | 40.74 | 7.41 | 26.00 | 35.50 | 0.50 |

Table 14. Results (%) on the UCR and UTabMWP [22] multiple-choice VQA subsets for LLaVA-1.5-7B. Metrics include standard, UPD, and dual accuracies.

Finally, a more fundamental limitation relates to the reliance on image-text alignment for open-ended questions. For questions with very general prompts (*e.g.*, "What do you see in the image?"), the alignment signal may not distinguish between answerable and unanswerable cases, making CLIP-UP less effective in such scenarios.

| | Standard | Unanswerable |
|---|---|---|
| **Question** | Based on the description, how are the people in the image engaging with the game?<br>A. The group of people is physically engaging with the game by using Nintendo Wii controllers.<br>B. The group of people is physically engaging with the game by using traditional gaming controllers.<br>C. The group of people is engaging with the game by watching a screen passively.<br>D. The group of people is engaging with the game by playing a board game. | |
| LLaVA-1.5-7B | A. The group of people is physically engaging with the game by using Nintendo Wii controllers. ✓ | C. The group of people is engaging with the game by watching a screen passively. ✗ |
| LLaVA-1.5-7B + CLIP-UP-EmbLoRA | A. The group of people is physically engaging with the game by using Nintendo Wii controllers. ✓ | I cannot answer. ✓ |
| LLaVA-NeXT-13B | A. The group of people is physically engaging with the game by using Nintendo Wii controllers. ✓ | A. The group of people is physically engaging with the game by using Nintendo Wii controllers. ✗ |
| LLaVA-NeXT-13B + CLIP-UP-EmbLoRA | A. The group of people is physically engaging with the game by using Nintendo Wii controllers ✓ | I cannot answer. ✓ |
| Phi-3.5-Vision | A. ✓ | A. ✗ |
| Phi-3.5-Vision + CLIP-UP-EmbLoRA | A. ✓ | I cannot answer. ✓ |
| InternVL3-1B | A. The group of people is physically engaging with the game by using Nintendo Wii controllers. In the image, the individuals are actively participating in a video game using Nintendo Wii controllers. The person on the right is holding a controller, and the person in the middle is also holding a controller, indicating that they are actively engaged in the game. The other individuals are standing around, watching the game, but they are not using controllers; instead, they are observing the screen. ✓ | C. The group of people is engaging with the game by watching a screen passively. The image shows a group of people watching a screen passively, which suggests they are not actively participating in the game but rather observing it. This is a common scenario in gaming settings where the focus is on the content displayed on the screen rather than the gameplay itself. ✗ |
| InternVL3-1B + CLIP-UP-EmbLoRA | A. The group of people is physically engaging with the game by using Nintendo Wii controllers. ✓ | I cannot answer. ✓ |
| InternVL3-8B | A. The group of people is physically engaging with the game by using Nintendo Wii controllers. ✓ | A. The group of people is physically engaging with the game by using Nintendo Wii controllers. ✗ |
| InternVL3-8B + CLIP-UP-EmbLoRA | A. The group of people is physically engaging with the game by using Nintendo Wii controllers. ✓ | I cannot answer. ✓ |
| Ovis2-16B | The description mentions that the people are playing a game on a Nintendo Wii, which involves physical interaction with the game using motion-sensitive controllers. Therefore, the correct answer is A. ✓ | The description mentions that the group of people is physically engaging with the game, using Nintendo Wii controllers, which suggests an active and interactive form of play. Therefore, the correct answer is A. ✗ |
| Ovis2-16B + CLIP-UP-EmbLoRA | A. The group of people is physically engaging with the game by using Nintendo Wii controllers. ✓ | I cannot answer. ✓ |

Figure 7. Models predictions on a pair of standard and IVQD unanswerable multiple-choice VQA questions from MM-UPD [45], without and with CLIP-UP-EmbLoRA.

| | Standard | Unanswerable |
|---|---|---|
| **Question** |  How many cats are visible in this picture? <br> A. Three <br> B. Four <br> C. Two <br> D. One | How many cats are visible in this picture? <br> A. Three <br> B. Four <br> C. Two |
| LLaVA-1.5-7B | D. ✓ | C. two. ✗ |
| LLaVA-1.5-7B + CLIP-UP-EmbLoRA | D. one. ✓ | I cannot answer. ✓ |
| LLaVA-NeXT-13B | D. ✓ | C. two. ✗ |
| LLaVA-NeXT-13B + CLIP-UP-EmbLoRA | D. one. ✓ | I cannot answer. ✓ |
| Phi-3.5-Vision | D. ✓ | C. ✗ |
| Phi-3.5-Vision + CLIP-UP-EmbLoRA | D. ✓ | I cannot answer. ✓ |
| InternVL3-1B | D. one. ✓ | C. two. ✗ |
| InternVL3-1B + CLIP-UP-EmbLoRA | D. one. ✓ | C. two. ✗ |
| InternVL3-8B | D. one. ✓ | C. two. ✗ |
| InternVL3-8B + CLIP-UP-EmbLoRA | D. one. ✓ | I cannot answer. ✓ |
| Ovis2-16B | D. one. ✓ | In the image, there is only one cat visible, sitting in a bathtub. Therefore, the correct answer is C. two, as it is the closest option to the actual number of cats present. ✗ |
| Ovis2-16B + CLIP-UP-EmbLoRA | D. one. ✓ | I cannot answer. ✓ |

Figure 8. Models predictions on a pair of standard and AAD unanswerable multiple-choice VQA questions from MM-UPD [45], without and with CLIP-UP-EmbLoRA.

| | Standard | Unanswerable |
|---|---|---|
| |  | |
| **Question** | Where are the clothes? | Where are the pink clothes? |
| LLaVA-1.5-7B | Suitcase. ✓ | Suitcase. ✗ |
| LLaVA-1.5-7B + CLIP-UP-EmbLoRA | Suitcase. ✓ | I cannot answer. ✓ |
| LLaVA-NeXT-13B | Suitcase. ✓ | Suitcase. ✗ |
| LLaVA-NeXT-13B + CLIP-UP-EmbLoRA | Suitcase. ✓ | I cannot answer. ✓ |
| Phi-3.5-Vision | Suitcase. ✓ | Suitcase. ✗ |
| Phi-3.5-Vision + CLIP-UP-EmbLoRA | Suitcase. ✓ | I cannot answer. ✓ |
| InternVL3-1B | Suitcase. ✓ | Bag. ✗ |
| InternVL3-1B + CLIP-UP-EmbLoRA | Bag. ✓ | I cannot answer. ✓ |
| InternVL3-8B | Suitcase. ✓ | Suitcase. ✗ |
| InternVL3-8B + CLIP-UP-EmbLoRA | Suitcase. ✓ | I cannot answer. ✓ |
| Ovis2-16B | Suitcase. ✓ | Suitcase. ✗ |
| Ovis2-16B + CLIP-UP-EmbLoRA | Suitcase. ✓ | I cannot answer. ✓ |

Figure 9. Models predictions on a pair of standard and unanswerable open-ended VQA questions from RGQA [66], without and with CLIP-UP-EmbLoRA.

| | Standard | Unanswerable |
|---|---|---|
| **Question** |  Which letter marks the territory controlled by the ancient Maya civilization? A. B B. C C. A D. D | Which letter marks the territory controlled by the ancient Maya civilization? A. B B. A C. D |
| LLaVA-1.5-7B | A. ✗ | A. ✗ |
| LLaVA-1.5-7B + CLIP-UP-EmbLoRA | A. B. ✗ | A. B. ✗ |
| LLaVA-NeXT-13B | B. ✓ | C. ✗ |
| LLaVA-NeXT-13B + CLIP-UP-EmbLoRA | D.D. ✗ | C.D. ✗ |
| Phi-3.5-Vision | B. ✓ | B. ✗ |
| Phi-3.5-Vision + CLIP-UP-EmbLoRA | B. ✓ | B. ✗ |
| InternVL3-1B | To determine which letter marks the territory controlled by the ancient Maya civilization... Given this information, the correct answer is: B. ✓ | To determine which letter marks the territory controlled by the ancient Maya civilization... Given this information, the correct answer is: A. ✗ |
| InternVL3-1B + CLIP-UP-EmbLoRA | B. C. ✓ | I cannot answer. ✓ |
| InternVL3-8B | To determine which letter marks the territory controlled by the ancient Maya civilization... Therefore, the letter that marks the territory controlled by the ancient Maya civilization is D. ✗ | To determine which letter marks the territory controlled by the ancient Maya civilization... Therefore, the letter that marks the territory controlled by the ancient Maya civilization is D. ✗ |
| InternVL3-8B + CLIP-UP-EmbLoRA | B. C. ✓ | I cannot answer. ✓ |
| Ovis2-16B | The ancient Maya civilization was located in the southeastern part of Mesoamerica, which corresponds to the region marked by letter B on the map. This area includes parts of present-day Mexico, Guatemala, Belize, Honduras, and El Salvador. The Maya civilization was known for its advanced writing system, architecture, and astronomical knowledge. Therefore, the correct answer is A. B. ✗ | The letter that marks the territory controlled by the ancient Maya civilization is B. The Maya civilization was centered in the region that is now southeastern Mexico, including the Yucatán Peninsula, and parts of Guatemala, Belize, Honduras, and El Salvador. This area corresponds to the region marked by letter B on the map. ✗ |
| Ovis2-16B + CLIP-UP-EmbLoRA | D. D. ✗ | C. D. ✗ |

Figure 10. Model predictions on standard and unanswerable multiple-choice VQA questions from MM-UPD [45], illustrating the limitations of CLIP-UP-EmbLoRA.