

FiVL: A Framework for Improved Vision-Language Alignment

Anonymous CVPR submission

Paper ID 12778

Abstract

Large Vision Language Models (LVLMs) have achieved significant progress in integrating visual and textual inputs for multimodal reasoning. However, a recurring challenge is ensuring these models utilize visual information as effectively as linguistic content when both modalities are necessary to formulate an accurate answer. We hypothesize that hallucinations arise due to the lack of effective visual grounding in current LVLMs. This issue extends to vision-language benchmarks, where it is difficult to make the image indispensable for accurate answer generation, particularly in vision question-answering tasks. In this work, we introduce FiVL, a novel method for constructing datasets designed to train LVLMs for enhanced visual grounding and to evaluate their effectiveness in achieving it. These datasets can be utilized for both training and assessing an LVLM's ability to use image content as substantive evidence rather than relying solely on linguistic priors, providing insights into the model's reliance on visual information. To demonstrate the utility of our dataset, we introduce an innovative training task that outperforms baselines alongside a validation method and application for explainability.

this issue has been to introduce a visual grounding mechanism to the model [4, 26, 27, 30, 32, 33].

Visual grounding aims to achieve more precise alignment between visual attention and semantic concepts within the model. A common method for grounding involves using bounding boxes, represented as a sequence of numerical numbers, to specify a particular region of an image. This enables the user to query specific parts of the image and the model to reference image locations within its generated response [4, 26, 30]. Bounding boxes, however, are coarse coordinates and unable to highlight objects or abstract concepts in finer detail. Additionally, the causal relationship between the generated bounding box and the response is not interpretable, and may instead be irrelevant to how the response is actually produced. Recent work have addressed these concerns by applying pixel-level grounding through the use of segmentation masks instead [27, 32, 33].

Training models with pixel-level grounding requires datasets that provide fine-grained visual alignment between images and text. However, such datasets are scarce, and prior work have often constructed custom datasets alongside model development [23, 27, 32, 33]. To address these challenges and simplify the training of visually grounded LVLMs, we introduce a novel method, FiVL (Framework for Improved Vision-Language Alignment) for constructing datasets with visual-concept alignment. Different from prior methods, our framework leverages existing datasets by employing SoTA segmentation models while prioritizing visual information that is crucial for accuracy. The main contributions of this paper are as follows:

- We introduce a framework for augmenting existing datasets with segmentation masks that align visual regions and the corresponding text
- We propose a novel pretraining task that jointly trains text and vision tokens to generate grounding masks
- Utilizing the proposed dataset, we construct a framework based on perturbations to assess a model's dependency on visual input and the necessity of visual context for providing accurate responses
- Finally, we leverage the dataset for an application focused on explainability.

079

2. Related Work

080

LVLM and Visual Grounding. Building upon LLMs, LVLMs extend their capabilities to a multimodal context by incorporating visual perception into the generation process, with notable models such as GPT-4o [25], LLaVA [18], Qwen2-VL [30], and many others [5, 6, 37], demonstrating advanced visual reasoning ability. Additionally, to link language output with visual input, some LVLMs employ grounding mechanisms to enhance multimodal interaction by allowing the model to reference specific regions of an image. This visual grounding has been achieved through the prediction of bounding boxes coordinates, as seen in models such as Kosmos-2 [26], Shikra [4], BuboGPT [34], Ferret [31], Qwen2-VL [30], and Groma [23]. To obtain a fine-grained localization of objects and semantic concepts pixel-level grounding has subsequently proposed in models such as Llava-Grounding [32], GLaMM [27], and GROUND-HOG [33].

097

Visually Grounded Datasets. Training LVLMs require large-scale visual instruction-following data [18]. However, these datasets focus on the task of visual and language reasoning and generally do not have fine-grained image segmentation annotations. Prior work has mainly constructed custom datasets to train their respective grounded LVLM models. In [23], a custom dataset, Groma Instruct, was constructed by prompting GPT-4V to generate grounded conversations based on 30K samples with region annotations from COCO [15] and VG [12]. Llava-Grounding [32] curated the Grounded Visual Chat (GVC) dataset by matching class labels of ground truth bounding boxes from COCO to noun phrases in conversations from LLaVA-Instruct-150K [18] using GPT-4. The Grounding-anything Dataset (GranD) was specifically constructed to train GLaMM [27] and utilized an object detection model to obtain visual entities that were then used to generate grounded dense captions through an LLM. A grounded visual instruction tuning dataset, M3G2, was proposed to train the GROUNDHOG model [33]. There, the authors curated a dataset consisting of 2.5M text-image pairs for visually grounded instruction tuning derived and augmented from 27 existing datasets.

119

3. Our Framework

120

3.1. Data Collection Pipeline

121

We created grounded datasets for both training and evaluation, building upon existing vision-question-answer and instruction datasets. Each sample in the original datasets was augmented with key expressions, along with their corresponding bounding box indices and segmentation masks within the images. The data collection pipeline proceeded as follows:

Key Expression Retrieval. The initial stage of data collection focused on identifying key expressions within each question-answer pair, using GPT-4o. These expressions are specific words or phrases that would be unattainable without the visual context provided by the image, such as object names, attributes, or spatial relations. We provided only the text of the question-answer pairs, omitting the images, and prompted GPT-4o to detect essential expressions with a custom-designed prompt (Appendix A). Using only questions and answers without visual cues allows GPT-4o to rely solely on linguistic context to determine whether certain words could be evoked based on text alone. This approach can help filter language-based answers from those needing visual context, while being computationally efficient. This process yielded a robust set of expressions, capturing the unique elements in each conversation that are closely tied to the visual information.

Bounding Box and Segmentation Masks. To accurately associate key expressions with specific regions in each image, we used the GroundedSAM pipeline [28], which employs the GroundingDINO-tiny model [20] for initial expressions localization generating bounding box indices, followed by the Segment Anything vit-huge model [11] for precise segmentation mask creation. Each key expression was mapped to its relevant visual region, creating high-quality segmentation maps. If multiple segments corresponded to a single phrase, they were consolidated into a unified mask assigned to each token within the phrase, to maintain consistency across annotations. During this process, we filtered out redundant key phrases, retaining only unique tokens to enhance the dataset's precision. Additionally, if a segmentation mask overlapped by more than 95% with another mask in the same sample, only one of the masks was retained. This filtering step ensured that each segmentation map uniquely represented essential visual regions, avoiding unnecessary redundancy and improving annotation clarity. No synthetic images or simulations were used; instead, we enriched existing images with detailed annotations that strengthen the visual grounding of language.

3.2. Training Dataset

Our training dataset, FiVL-Instruct, is built upon the LLaVA-1.5-mix-665K instruction tuning dataset [16], a public vision-language instruction dataset containing 665K structured dialogues between users and GPT. Most interactions begin with a user-provided image, followed by questions related to the visual content, with GPT offering responses, each question-answer pair is referred as a turn.

We augmented the original LLaVA-1.5-mix-665K dataset by integrating the key expressions and their segmentation masks according to the pipeline outlined in Section 3.1. Not every FiVL-Instruct sample includes a key expres-

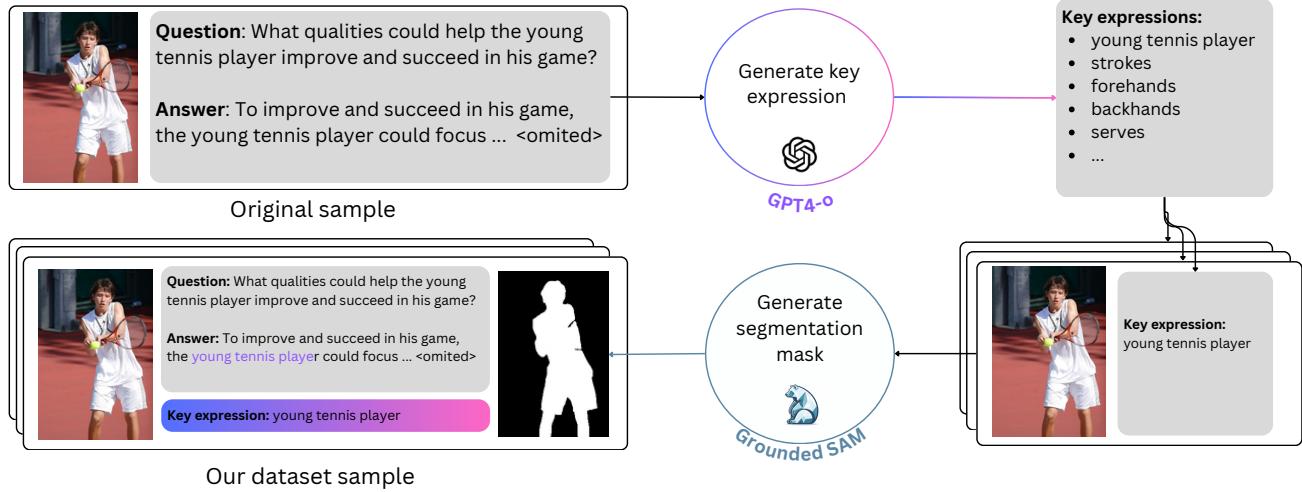


Figure 1. Dataset Collection Overview: The top figure illustrates that GPT is initially used to extract key expressions, while the bottom figure demonstrates that grounded SAM is used to extracts segmentations based on the provided keyword and image.

179 sion. For such cases, we retained the original data point
 180 unchanged to maintain the dataset size for training. In our
 181 dataset, each conversation consists of multiple turns with an
 182 average of ten turns. Images are present in 94% of the con-
 183 versations. Across the entire dataset, we collected around
 184 1.5 million unique segmentation masks, resulting in an av-
 185 erage of 2.3 segmentation masks per conversation. This cor-
 186 responds to about 2.3 million key expressions, averaging
 187 3.5 key expressions per conversation. The range spans from
 188 0 key expression for samples without images or identified
 189 keywords to a maximum 33 per conversation. On average,
 190 each key expression contains 2.4 words. The average area
 191 of the segmentation is 28% of the image.

192 3.3. Evaluation Datasets

193 To assess the visual reliance of various LVLMs, we cre-
 194 ated three benchmark datasets derived from the validation
 195 sets of POPE [14], VQAv2 [8], and GQA [10]. These aug-
 196 mented datasets, named FiVL-POPE, FiVL-VQAv2, and
 197 FiVL-GQA will be used to assess the extent to which mod-
 198 els rely on visual information for accurate responses. We
 199 selected these benchmarks because they each requires dif-
 200 ferent levels of image reliance. POPE assesses sensitivity to
 201 visual perturbations, GQA evaluates understanding of de-
 202 tailed scene relationships, and VQAv2 tests visual ground-
 203 ing for diverse question types. Together, they offer a well-
 204 rounded assessment of how much models depend on visual
 205 information to answer accurately.

206 Following the procedure outlined in Section 3.1, for each
 207 sample, we identified key expressions that can be extracted
 208 from the question but also the answer and produce the re-
 209 lated segmentation mask. We included the answers as well,
 210 as many samples feature the key element in the question,

211 and a significant portion of answers are simple responses
 212 like “Yes,” “No,” or numerical values. Section 5.2 describes
 213 our use of these augmented datasets to evaluate how ef-
 214 fectively models focus on visual cues. Similar to FiVL-
 215 Instruct dataset, GPT-4o does not consistently identify key
 216 expressions from each question-answer pair of these three
 217 datasets. Moreover, even for samples with key expressions,
 218 not all have corresponding segmentation masks generated
 219 by the GroundedSAM pipeline. This indicates that some of
 220 the questions from these datasets might not depend on the
 221 images. In our evaluation set, we filter those cases without
 222 extracted keywords or segmentation. Table 1 compares the
 223 size of our datasets with the original datasets and Table 2
 224 shows the other statistic of our datasets.

	VQA-v2	GQA	POPE
Original	9,999	12,280	9,000
FiVL	4,040	11,660	5,870

Table 1. Sizes of the evaluation datasets after filtering out samples lacking key expressions or segmentation masks.

	FiVL-VQAv2	FiVL-GQA	FiVL-POPE
Key expressions	1.27	1.5	1
Segmentation masks	3.79	4.71	3.48
% of masked pixels	24%	21%	16%

Table 2. Statistics of our evaluation datasets. First row details the average number of key expressions per sample, second row describes the average number of distinct segmentation masks per sample and last row describes the average percentage of the pixels that were masked per sample.

225

4. Method Evaluation

226 To ensure the quality of our datasets and effectiveness
 227 in training LVLMs for visual grounding, we conducted a
 228 multi-step evaluation process on the training dataset de-
 229 scribed in Section 3.2. This included both human-based
 230 evaluations and automated assessments, allowing us to val-
 231 idate the relevance and accuracy of the key tokens and their
 232 alignment with visual content. Below, we outline the key
 233 components of our evaluation strategy.

234

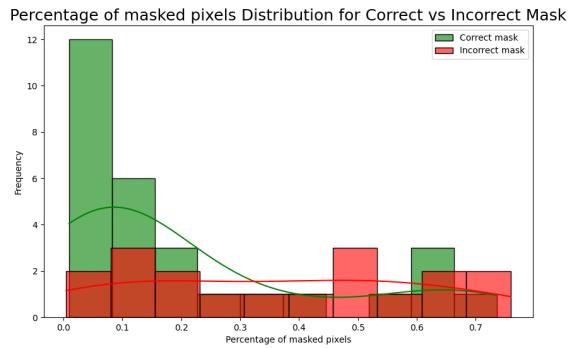
4.1. Human Evaluation

235 We conducted a manual evaluation in order to validate the
 236 coherency of the key expressions as well as the relevancy
 237 of the segmentation maps with respect to the formers. For
 238 each sample, we presented to the annotators one random
 239 key expression with its associated segmentation map. An-
 240 notators were asked three questions: *whether the key ex-*
241 pression aligns with the definition provided in Section 3.1, if
242 the segmentation map is relevant to the key expression, and
243 whether the sample is of good quality (does the text makes
 244 sense, is the answer related to the question...). In total, 557
 245 unique samples were annotated by 12us different annota-
 246 tors. The accumulated results of the three types of question
 247 is that 77% of the annotators labeled the samples as overall
 248 good data points. In detail, for the key expression evalua-
 249 tion, 75% of key expressions were marked as essential to the
 250 answer’s meaning. For the segmentation map evaluation,
 251 58% of segmentation map were annotated as relevant to the
 252 key expression. This can be explained by the fact that some
 253 key expressions might inherently be more abstract or com-
 254 plex or by the performance of the GroundedSAM pipeline.
 255 We also find that the quality of the segmentation is related to
 256 the size of the segmentation. Finally, if we compute the key
 257 expressions and segmentations score only for the samples
 258 annotated as “good data point overall”, 85% of the data are
 259 with valid key expressions and 69 % are with relevant seg-
 260 mentation masks. Specifically, Figure 2 indicates that when
 261 the segmented mask occupies less than 20% of the image,
 262 annotators were more likely to consider the segmentation
 263 relevant. To train our model (see Section 5.1), we selected
 264 the segmentation masks based on their size. We also made
 265 ablations by using FiVL-Instruct filtered to exclude big seg-
 266 mentation masks and compared it with that of training on
 267 the unfiltered version and observed no difference in per-
 268 formance. Regardless, this metric can be used as a threshold
 269 filtering method for future applications of the dataset. A
 270 snapshot of our API is provided in Appendix A.

271

4.2. Automatic Evaluation

272 Manual human evaluation is time-consuming. Inspired by
 273 recent applications using GPT as an evaluation tool [35],
 274 we designed two prompting techniques to automatically as-
 275 sess the quality of extracted keywords and segmentation



276 Figure 2. Impact of the size of the segmentation mask. Compari-
 277 son of the Percentage of masked pixels Distributions for Correctly
 278 and Incorrectly annotated Masks

279 masks based on a given keyword. Both evaluations were
 280 conducted on a randomly sampled set of 1,957 keywords
 281 and their corresponding segmentations from FiVL-Instruct
 282 dataset.

4.2.1. Keyword Evaluation

283 We prompt the model to verify if the keyword is important
 284 to the question and depending on the image. Additionally,
 285 we ask the model to rate the degree of importance on a scale
 286 from 0 to 10. The full prompt is provided in Figure 4. We
 287 report three metrics. The first is the *Importance Ratio* =
 288 76% representing the percentage of extracted expressions
 289 classified as key expressions. This result is close to human
 290 evaluation, which is 75%. The second is the *Overall Im-*
291 portance Degree = 6.8, which indicates the average impor-
 292 tance score across all keywords, regardless of whether GPT-
 293 4o classifies them as important or not important. The third
 294 metric is the *Importance Degree of Important Keywords* =
 295 9.0, which calculates the average importance ratio of key-
 296 words identified as important by GPT-4o. These metrics
 297 indicate the high quality of our keywords.

4.2.2. Segmentation Evaluation

298 Given a keyword, we aim to evaluate whether our segmen-
 299 tation for this keyword is accurate. We designed two prompts
 300 to assess the quality of the segmentation: first, we check
 301 if the segmentation content adequately covers the keyword
 302 (Seg1); second, we verify that the inverse of the segmen-
 303 tation does not contain any content related to the keyword
 304 (Seg2). Both prompts are given in Figure 4. Results show
 305 that only for Seg1 = 46% of the cases GPT-4o capture the
 306 keywords in the segmentation. On one hand, this results
 307 aligns with the manual annotations and can be addressed
 308 in the same manner. On the other hand, we found that
 309 segmentations classified as good often involve specific ob-
 310 jects (e.g., tennis players, bears). In contrast, segmentations
 311 classified as bad are often abstract concepts (e.g., water

You are given a question, a word/phrase and an image. Please rate the importance degree from 0-10 scale ([OID]).
Note that
- 0 means not important at all and 10 means very important.
- Important word/phrase means that this word/phrase is closely related to the image and the question, and it could not be evoked without the use of the image (IR).
- If the question does not relate to the image, in other words, the answer does not depend on the image content, then any words are not important.

Question: {question}

A word: {word}

Only answer important or not important, and the importance degree from 0-10?

Figure 3. Keyword Verification Prompt for GPT-4o.

[Seg1] You are given a part of the image and a word/phrase, do you think this is a good segmentation that the given part of the image covers this word/phrase?
Word/phrase: {word}
Answer only "yes" or "no".

[Seg2] You are given a part of the image and a word/phrase, do you see any part of the image that is related to the word?
Word/phrase: {word}
Answer only "yes" or "no".

Figure 4. Segmentation Verification Prompt for GPT-4o.

311 pressure, mental game, splashing), descriptive words (e.g.,
312 unique, uneven ground), or complex actions (e.g., walking
313 over logs). These types of words are difficult to link to a
314 specific part of an image when the full image context is not
315 provided. This also highlights the limitations of the first
316 type of evaluation prompt. In our training setup, we will
317 take this into account and only focus on the nouns of the
318 key expressions. In *Seg2* = 72% of cases, the model de-
319 termines that the inverse of the segmentation is irrelevant to
320 the keywords, accurately recognizing that without the seg-
321 mented mask, the key expressions are not present in the im-
322 age. This measures the fact that we do not miss key objects
323 in our segmentation maps. If 2 objects appear in the im-
324 age not at the same positions, we make sure that our maps
325 contain both of them.

5. Applications of FiVL Dataset

In this section, we describe three ways to utilize our datasets. Section 5.1 describes how FiVL can be used as a training dataset and the resulting models not only achieve better performance but also has one more capability than the baseline model: generate segmentation maps. Section 5.2 introduces FiVL as a tool for evaluating the visual reliance of LVLMs. Section 5.3 shows that FiVL can assist the interpretability of models.

5.1. Training

We introduce here a novel pretraining task, Vision Modeling. To assess the effectiveness of this task, we fine-tuned an LLM, specifically, LLaVA-1.5-7b [17], on FiVL-Instruct. For training our model, we used only key expressions that appeared verbatim in the answers for each turn, focusing exclusively on noun-based key expressions. By leveraging this dataset, we ensured alignment with widely used benchmarks in vision-language research, facilitating comparability and reproducibility across studies.

Method. In order to train Large Language Models and even some LVLMs (such as LLaVA), the cross-entropy loss is used to minimize the difference between the predicted probability distribution of the text output and the true distribution of the target words. LLaVA training strategy is in two stages: the first (pretraining) trains a projector which aims to align visual and textual representations, while the second (finetuning) performs only language modeling on the textual outputs of the LM head. In this work, we propose to also guide the visual outputs of the last linear layer (the LM head). Figure 7 plots were created by taking the *argmax* of the vision logits, mapping it back to the text vocabulary, and highlighting the image patches that relate to that token. We can see from that Figure, column on the right, that the baseline is already able to capture a relevant token from the vocabulary with respect to the image patches. However, since it was not trained to do so, we obtained a lot of different tokens that realize the argmax for each patch, some of them not being relevant. These examples were extracted from samples of LLaVA-1.5-mix-665K instruction dataset. Our dataset gives us access to segmentation maps inherently provided by the vision logits argmax. We augment the Language modeling cross-entropy loss with a Vision modeling (VM) cross-entropy loss where each patch that belongs to a segmentation map is trained to predict the related keyword from the vocabulary.

We denote by x the input and y the logits with respect to each token. The logits are the outputs of the last linear layer that projects the last hidden states to the vocabulary space:

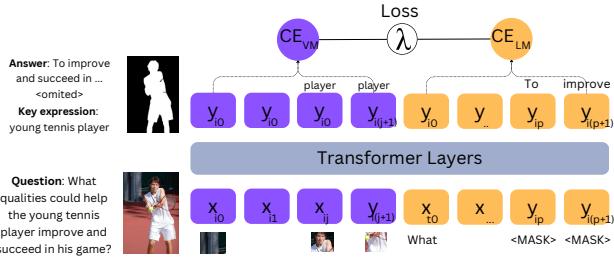


Figure 5. Overview of Vision Modeling pretraining task

$$\begin{aligned} 374 \quad x &= (x_{i_0}, x_{i_1}, \dots, x_{i_{N_i}}, x_{t_0}, x_{t_1}, \dots, x_{t_N}), \\ 375 \quad y &= (y_{i_0}, y_{i_1}, \dots, y_{i_{N_i}}, y_{t_0}, y_{t_1}, \dots, y_{t_N}), \end{aligned}$$

376 where N_i is the number of image tokens, N_t the number
377 of text tokens, N the total length. x_i are the inputs embedding
378 that relate to the image tokens and x_t to the text tokens;
379 $y_i \in \mathbb{R}^{N_i \times \text{vocabulary_size}}$ represents visual logits, while
380 $y_t \in \mathbb{R}^{N_t \times \text{vocabulary_size}}$ represents textual logits.

381 In Language Modeling (LM), only y_t related to the
382 answer are trained. We propose to also train y_i related to the
383 segmented piece. For instance: if given a picture of a bowl
384 of soup the question is *What qualities could help the young*
385 *tennis player improve and succeed in his game?*. Along
386 with it, our dataset provides the segmentation map related
387 to the *young tennis player*. The LM loss would only guide
388 the relevant tokens y_t to be *To improve and succeed in ...*
389 *;omited;*. In our method, we also do vision modeling by
390 training each visual logit corresponding to the segmented
391 mask to refer to the noun from the key expression: *player*
392 from the text vocabulary. See Figure 5 for a visual expla-
393 nation of the method. In order to create the vision labels
394 we proceed like such: for each sample, each image token
395 will be assigned to exactly one token in the text vocabulary.
396 The selection is based on the size of the mask (we take the
397 smallest) and the type of the keyword (we filter only nouns).
398 If an object (containing several patches) gets more than one
399 text token assigned to it, we randomly select one. That way,
400 for each image patch, there is maximum one key token that
401 describes the patch. Image patches that do not have a re-
402 lated keytoken are ignored in the loss, similar to LM. We
403 then compute the cross-entropy, CE_{VM} between the cre-
404 ated vision labels and the visual logits. We then compute
405 a weighted sum between this loss and the cross-entropy re-
406 lated to language modeling, CE_{LM} . The resulting loss is
407 computed as such:

$$408 \quad L = \lambda * CE_{VM} + (1 - \lambda) * CE_{LM}, \text{ where } \lambda \in [0, 1]$$

409 **Results.** We conducted multiple experiments to determine
410 the optimal hyperparameters. We finetuned LLaVa-v1.5-

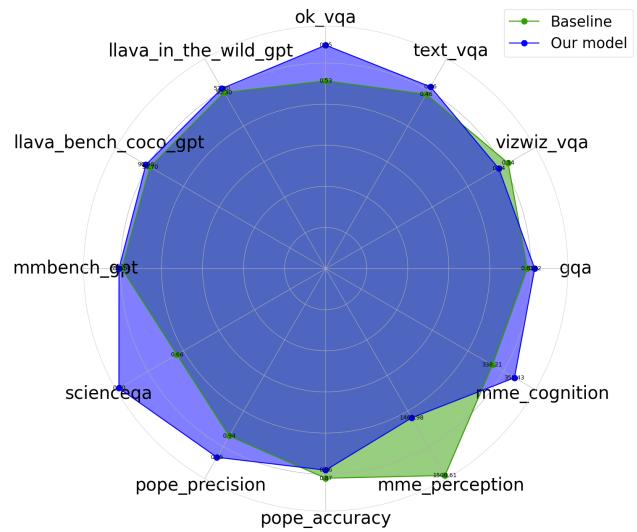


Figure 6. Our model trained on FiVL-Instruct evaluated on various benchmarks compared to the baseline.

7b from scratch using our augmented dataset. We used the
411 trained multimodal projector and started from Vicuna-v1.5-
412 7b [36] weights. We maintained the original training setup
413 (batch size, number of epochs, etc.) and primarily focused
414 on experimenting with different learning rates and λ . The
415 best results were achieved with a learning rate of 2e-5, the
416 same as in the original setup, and λ set to 0.1. (Results were
417 equivalently good for 0.2). Figure 6 shows how we outper-
418 formed the baseline in different benchmarks. We outper-
419 form the baseline for OK-VQA [24], for which our model
420 achieves 0.55 of accuracy compared to 0.53 for the base-
421 line, MME-cognition [7] for which our model achieves a
422 score of 351 compared to 338 for the baseline, for POPE
423 precision [14] for which our model achieves a precision of
424 0.95 compared to 0.94 for the baseline, ScienceQA [22],
425 for which our model gets an accuracy of 0.70 compared to
426 0.66, for MMBench [21] where we obtain 51.6 compared
427 to 50.1, LLaVA-Bench-COCO [17], where we obtain 96
428 against 93.7 and for LLaVA-in-the-wild [17] for which we
429 obtained 52.6 against 50.3. For Text-VQA [29], VizWiz-
430 VQA [9], GQA [10], POPE accuracy [14] we get compa-
431 rable accuracies of respectively 0.46, 0.54, 0.62, 0.86. We
432 do notice a small drop in performance for MME-perception
433 [7] from 1470 to 1509.

434 As an additional outcome of the training, we also ob-
435 tain “segmentation maps” by predicting text tokens using
436 an argmax operation over vision logits, which are the out-
437 puts at the positions of vision tokens. As a simple observa-
438 tion, averaging over 100 examples, the baseline predicts 74
439 different tokens overall (with lots of unrelated tokens such
440

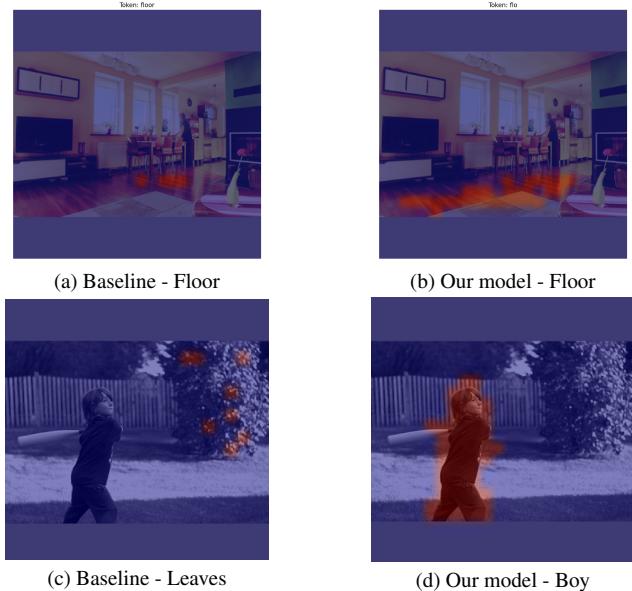


Figure 7. Example of predicted text from vision logits and their corresponding region in the image. The left column displays the vision logits from the baseline model, while the right column shows the output from our model.

as "a", "*", "is" etc.), while our model only encompasses 9 tokens. This demonstrates potential in leveraging visual logits for segmentation, as shown in Figure 7.

5.2. Visual Reliance Evaluation

FiVL datasets also allow us to measure *Visual Reliance* by performing perturbation based evaluation: first evaluating model accuracy on the original, unmasked images, then assessing performance on the masked images. Then, we design a *Visual Reliance Score* (Eq.1), which indicates the percentage of drop in accuracy from the original to the masked image version, with higher scores indicating stronger dependency on visual input. A higher score indicates a greater drop in performance from the original image to the masked image, suggesting a stronger reliance of the model on visual information for generating answers.

$$\text{Visual Reliance Score} = \frac{\text{accuracy}_{\text{original}} - \text{accuracy}_{\text{perturb}}}{\text{accuracy}_{\text{original}}} \quad (1)$$

To validate that FiVL is a good dataset to evaluate visual reliance, we generated a control dataset with random masking. In this control set, each image includes a black rectangular mask of the same size as the key expression mask but placed at a random location within the image. This approach provides a comparison to determine whether performance declines specifically due to masking critical visual areas or simply from general occlusion.

We compared the performance of two models, LLaVA-v1.5-13b and Qwen2-VL-7B-Instruct [30], across the three evaluation datasets we created.

Compare FiVL and Random Perturbation. Table 3 shows that across all benchmarks and models, the perturbation based on FiVL masks causes a significantly larger performance drop compared to random perturbation, indicating that our bounding boxes capture meaningful visual content relevant to the questions and FiVL represent good testbeds for visual reliance. One thing to note is that for LLaVA-13B on VQA-v2, we see that masking with random bounding box even improves over the original images. We assume that sometimes random masking may hide non-essential parts of the image, allowing the model to better focus on the remaining regions, which likely contain the intended area of focus.

Comparison LLaVA and Qwen Models. For the VQA-v2 and POPE, we observe that LLaVA-13B shows a greater drop in accuracy compared to Qwen2-VL-7B-Instruct. For GQA, we see the opposite behavior. Since LLaVA-13B exhibits a larger drop on more datasets and the average drop across all three datasets is a bit higher compared to Qwen2-VL (0.513 versus 0.497, respectively), this suggests that the LLaVA-13B model has a higher visual reliance on visual information than the LLaVA model. However, we acknowledge that three datasets may not be sufficient to draw a definitive conclusion, and expanding the evaluation set is a goal for future work.

Compare Three Evaluation Sets. Across models, we observe a larger performance drop on VQA-v2 compared to POPE, and a larger drop on POPE compared to GQA, indicating that VQA-v2 relies most heavily on the image for answering questions, followed by POPE and then GQA.

5.3. Explainability

Here, we also show that FiVL can assist the interpretability of black box models. We produce a summary plot that displays a vision-alignment metric calculated across all heads and layers, like introduced in [2]. Figure 9a offers a detailed view of how the model grounds vision and language through

	VQA-v2		GQA		POPE	
	FiVL	Random	FiVL	Random	FiVL	Random
LLaVA-13B	0.72	-0.05	0.33	0.03	0.49	0.02
Qwen2-VL-7B	0.64	0.07	0.38	0.03	0.47	0.02

Table 3. Percentage of drop in performance using the perturbation based method. Using FiVL bounding boxes to introduce perturbations compared to random perturbations across different benchmarks and models.

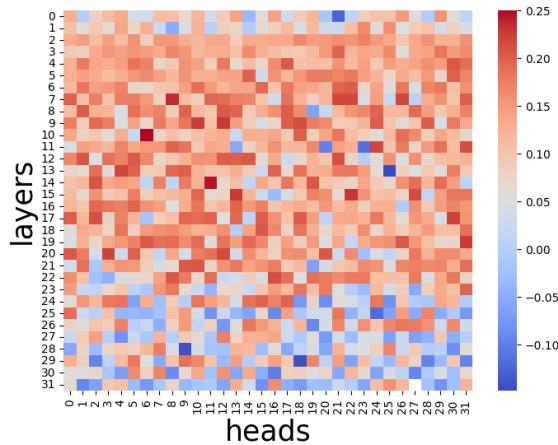


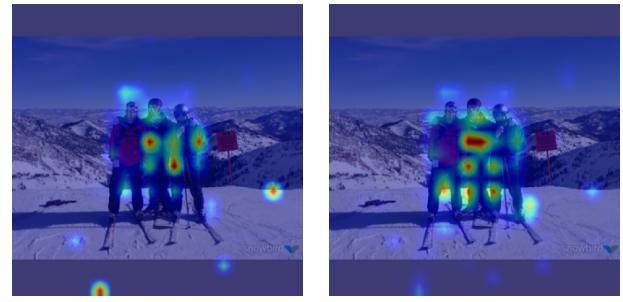
Figure 8. Head summary for VL alignment computed using the Spearman correlation between the segmentation of a token and its related vision attention

the attentions of each head across all layers. For an attention matrix of size $(N_{\text{layers}}, N_{\text{heads}}, N_i + N_t, N_i + N_t)$, The head summary calculates the statistical mean over the last two dimensions, producing a plot with dimensions of $(N_{\text{layers}}, N_{\text{heads}})$ averaged for 500 samples. For a given question, image, key expression and related segmentation mask from the FiVL-Instruct dataset, we generate the answer using LVA-VA-v1.5-7b. We then identify if the key expression is in the answer or in the question. If so, we probe each head by computing the Spearman correlation between the segmentation mask $(\sqrt{N_i}, \sqrt{N_t})$. and the attention to the corresponding key expression tokens in the Vision-to-Language attention component $(1, 1, N_i, 1)$ (first dimension selects the layer, second the head and the last dimension corresponds to the key token) for each head. This is performed on the language model component but not on the vision component of LLaVA. In this way, we identify the attention heads that ground the most the two modalities by performing a function similar to object segmentation. Figure 9 shows the head summary and the corresponding language-vision attention weights related to the key expression tokens displayed as a heatmap over the image. The head summary shows that the heads achieving the strongest vision-language alignment are in the early layers. This might be due to the fact that the input to this transformer is the output of multimodal projector of LLaVA, which is designed specifically to align these two modalities. The head summary indicates that heads (10,6) and (14,11) are effective at aligning vision with language. Figures 9a and 9b show from which patches of the image the token *girl* gets the most attention, clearly focusing on the girl depicted in the image.



(a) Attention Head (10,6) of the token "girl". Q - Who are the two people playing Frisbee in the image? A - The two people playing Frisbee in the image are a man and a little girl

(b) Attention Head (14,11) of the token "girl". Q - Who are the two people playing Frisbee in the image? A - The two people playing Frisbee in the image are a man and a little girl



(c) Attention Head (10,6) of the token *three*. Q - A - There are three people in the image

(d) Attention Head (14,11) of the token *three*. Q - A - There are three people in the image

Figure 9. Attention heatmaps overlaid on the original images for attention heads (10,6) and (14,11), which have a high Spearman Correlation, to probe vision-language alignment.

6. Conclusion

In this paper, we introduced FiVL, a framework designed to enhance vision-language alignment and visual focus in large vision-language models. We applied our approach across key stages of an LVLM training workflow: training, evaluation, and explainability. By training a LLaVA model using our FiVL dataset, we saw improvement in a majority of test benchmarks. Our evaluation datasets facilitated comparisons between models regarding their reliance on images to answer questions and provided insight into the degree of image dependency needed across benchmarks. Finally, our explainability application enables users to identify attention heads that excel in vision-language alignment, allowing for a deeper understanding of the internal workings of LVLMs and potential model refinement. Our training method produced a built-in feature that segments the image. Future work could include the evaluation of these segmentations and the enhancement of grounded LVLMs using FiVL.

535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552

553 **References**

- [1] Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadal-lah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Ben-haim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Qin Cai, Vishrav Chaudhary, Dong Chen, Dongdong Chen, Weizhu Chen, Yen-Chun Chen, Yi-Ling Chen, Hao Cheng, Parul Chopra, Xiyang Dai, Matthew Dixon, Ronen Eldan, Victor Fragoso, Jianfeng Gao, Mei Gao, Min Gao, Amit Garg, Allie Del Giorno, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Wenxiang Hu, Jamie Huynh, Dan Iter, Sam Ade Jacobs, Mojan Javaheripi, Xin Jin, Nikos Karam-patiakis, Piero Kauffmann, Mahoud Khademi, Dongwoo Kim, Young Jin Kim, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Yunsheng Li, Chen Liang, Lars Liden, Xihui Lin, Zeqi Lin, Ce Liu, Liyuan Liu, Mengchen Liu, Weishung Liu, Xiaodong Liu, Chong Luo, Piyush Madan, Ali Mahmoudzadeh, David Majercak, Matt Mazzola, Caio César Teodoro Mendes, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Liliang Ren, Gustavo de Rosa, Corby Rosset, Sam-budha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacroce, Shital Shah, Ning Shang, Hiteshi Sharma, Yelong Shen, Swadheen Shukla, Xia Song, Masahiro Tanaka, Andrea Tupini, Praneetha Vaddamanu, Chunyu Wang, Guanhua Wang, Lijuan Wang, Shuohang Wang, Xin Wang, Yu Wang, Rachel Ward, Wen Wen, Philipp Witte, Haiping Wu, Xiaoxia Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu, Weijian Xu, Jilong Xue, Sonali Yadav, Fan Yang, Jianwei Yang, Yifan Yang, Ziyi Yang, Dong-han Yu, Lu Yuan, Chenrudong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyra Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and Xiren Zhou. Phi-3 technical report: A highly capable language model locally on your phone, 2024. 12
- [2] Estelle Aflalo, Meng Du, Shao-Yen Tseng, Yongfei Liu, Chenfei Wu, Nan Duan, and Vasudev Lal. Vi-interpret: An interactive visualization tool for interpreting vision-language transformers, 2022. 7
- [3] Pravesh Agrawal, Szymon Antoniak, Emma Bou Hanna, Baptiste Bout, Devendra Chaplot, Jessica Chudnovsky, Diogo Costa, Baudouin De Monicault, Saurabh Garg, Theophile Gervet, Soham Ghosh, Amélie Héliou, Paul Jacob, Albert Q. Jiang, Kartik Khandelwal, Timothée Lacroix, Guillaume Lample, Diego Las Casas, Thibaut Lavril, Teven Le Scao, Andy Lo, William Marshall, Louis Martin, Arthur Mensch, Pavankumar Muddireddy, Valera Nemychnikova, Marie Pellat, Patrick Von Platen, Nikhil Raghuraman, Baptiste Rozière, Alexandre Sablayrolles, Lucile Saulnier, Romain Sauvestre, Wendy Shang, Roman Soletskyi, Lawrence Stewart, Pierre Stock, Joachim Studnia, Sandeep Subramanian, Sagar Vaze, Thomas Wang, and Sophia Yang. Pixtral 12b, 2024. 12
- [4] Keqin Chen, Zhao Zhang, Weili Zeng, Richong Zhang, Feng Zhu, and Rui Zhao. Shikra: Unleashing multi-modal llm's referential dialogue magic. *arXiv preprint arXiv:2306.15195*, 2023. 1, 2
- [5] Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhang-wei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*, 2024. 2
- [6] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. InstructBLIP: Towards general-purpose vision-language models with instruction tuning. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. 2
- [7] Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, Yunsheng Wu, and Rongrong Ji. Mme: A comprehensive evaluation benchmark for multimodal large language models, 2024. 6
- [8] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering, 2017. 3
- [9] Danna Gurari, Qing Li, Abigale J. Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P. Bigham. Vizwiz grand challenge: Answering visual questions from blind people, 2018. 6
- [10] Drew A. Hudson and Christopher D. Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering, 2019. 3, 6
- [11] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. Segment anything. *arXiv:2304.02643*, 2023. 2
- [12] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision*, 123:32–73, 2017. 2
- [13] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models, 2023. 12
- [14] Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models, 2023. 3, 6
- [15] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In David Fleet, Tomas Pajdla, Bernt Schiele, and Tinne Tuytelaars, editors, *Computer Vision – ECCV 2014*, pages 740–755, Cham, 2014. Springer International Publishing. 2
- [16] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning, 2024. 2
- [17] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023. 5, 6, 12
- [18] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024. 1, 2
- [19] Hanchao Liu, Wenyuan Xue, Yifei Chen, Dapeng Chen, Xiutian Zhao, Ke Wang, Liping Hou, Rongjun Li, and Wei Peng. 9

- 670 A survey on hallucination in large vision-language models.
arXiv preprint arXiv:2402.00253, 2024. 1
- 671 [20] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao
672 Zhang, Jie Yang, Qing Jiang, Chunyuan Li, Jianwei Yang,
673 Hang Su, Jun Zhu, and Lei Zhang. Grounding dino: Marry-
674 ing dino with grounded pre-training for open-set object de-
675 tection, 2024. 2
- 676 [21] Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang
677 Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He,
678 Ziwei Liu, Kai Chen, and Dahua Lin. Mmbench: Is your
679 multi-modal model an all-around player?, 2024. 6
- 680 [22] Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei
681 Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and
682 Ashwin Kalyan. Learn to explain: Multimodal reasoning via
683 thought chains for science question answering, 2022. 6
- 684 [23] Chuofan Ma, Yi Jiang, Jianne Wu, Zehuan Yuan, and Xiao-
685 juan Qi. Groma: Localized visual tokenization for grounding
686 multimodal large language models. In *European Conference
687 on Computer Vision*, pages 417–435. Springer, 2024. 1, 2
- 688 [24] Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and
689 Roozbeh Mottaghi. Ok-vqa: A visual question answering
690 benchmark requiring external knowledge, 2019. 6
- 691 [25] OpenAI. Gpt-4o system card, 2024. 1, 2, 12
- 692 [26] Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan
693 Huang, Shuming Ma, and Furu Wei. Kosmos-2: Grounding
694 multimodal large language models to the world. arXiv
695 preprint arXiv:2306.14824, 2023. 1, 2
- 696 [27] Hanoona Rasheed, Muhammad Maaz, Sahal Shaji, Abdel-
697 rahman Shaker, Salman Khan, Hisham Cholakkal, Rao M.
698 Anwer, Eric Xing, Ming-Hsuan Yang, and Fahad S. Khan.
699 GLaMM: Pixel grounding large multimodal model. In *Pro-
700 ceedings of the IEEE/CVF Conference on Computer Vision
701 and Pattern Recognition (CVPR)*, pages 13009–13018, June
702 2024. 1, 2
- 703 [28] Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kun-
704 chang Li, He Cao, Jiayu Chen, Xinyu Huang, Yukang Chen,
705 Feng Yan, Zhaoyang Zeng, Hao Zhang, Feng Li, Jie Yang,
706 Hongyang Li, Qing Jiang, and Lei Zhang. Grounded sam:
707 Assembling open-world models for diverse visual tasks,
708 2024. 2
- 709 [29] Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang,
710 Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus
711 Rohrbach. Towards vqa models that can read, 2019. 6
- 712 [30] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan,
713 Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin
714 Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui
715 Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Jun-
716 yang Lin. Qwen2-vl: Enhancing vision-language model’s
717 perception of the world at any resolution, 2024. 1, 2, 7
- 718 [31] Haoxuan You, Haotian Zhang, Zhe Gan, Xianzhi Du, Bowen
719 Zhang, Zirui Wang, Liangliang Cao, Shih-Fu Chang, and
720 Yinfei Yang. Ferret: Refer and ground anything anywhere
721 at any granularity. In *The Twelfth International Conference
722 on Learning Representations*, 2024. 2
- 723 [32] Hao Zhang, Hongyang Li, Feng Li, Tianhe Ren, Xueyan
724 Zou, Shilong Liu, Shijia Huang, Jianfeng Gao, Chunyuan
725 Li, Jainwei Yang, et al. Llava-Grounding: Grounded visual
726 chat with large multimodal models. In *European Conference
727 on Computer Vision*, pages 19–35. Springer, 2024. 1, 2
- 728 [33] Yichi Zhang, Ziqiao Ma, Xiaofeng Gao, Suhaila Shakiah,
729 Qiaozi Gao, and Joyce Chai. Groundhog: Grounding large
730 language models to holistic segmentation. In *Proceedings of
731 the IEEE/CVF Conference on Computer Vision and Pattern
732 Recognition (CVPR)*, pages 14227–14238, June 2024. 1, 2
- 733 [34] Yang Zhao, Zhiping Lin, Daquan Zhou, Zilong Huang, Jiashi
734 Feng, and Bingyi Kang. Bubogpt: Enabling visual ground-
735 ing in multi-modal llms. arXiv preprint arXiv:2307.08581,
736 2023. 2
- 737 [35] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan
738 Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan
739 Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with
740 mt-bench and chatbot arena. *Advances in Neural Information
741 Processing Systems*, 36:46595–46623, 2023. 4
- 742 [36] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan
743 Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan
744 Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonza-
745 lez, and Ion Stoica. Judging llm-as-a-judge with mt-bench
746 and chatbot arena, 2023. 6
- 747 [37] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mo-
748 hamed Elhoseiny. MiniGPT-4: Enhancing vision-language
749 understanding with advanced large language models. In *The
750 Twelfth International Conference on Learning Representa-
751 tions*, 2024. 2

753 **Appendix**754 **A. System prompts for key expressions re-
755 trieval**

We use GPT-4o via the Azure OpenAI API to extract the key expressions of the datasets we considered. In this section, we share the prompts used for this step of the data collection. We had to use slightly different prompts for the training datasets compared to the evaluation datasets. In the training datasets, where instructions are open-ended question-answer pairs, the key expressions are often found in the answer. However, in the evaluation datasets, we encountered questions that required specific types of responses (yes/no questions, counting etc...). In these cases, the key expressions are typically found in the question instead. For references, we have provided the prompt used for training dataset in Figure 11 and prompts used for evaluation datasets VQA-V2, GQA, and POPE in Figure 12. For each benchmark we use different examples that suit the best to the types of questions. See Figure 13 for FiVL-VQAv2 and Figure 14 for FiVL-GQA and FiVL-POPE.

773 **B. Training details**

To train our model, we used Nvidia RTX A6000 GPUs using the hyperparameters from Table 4

Batch Size	4
Number of GPUs	8
Gradient Accumulation	4
Number of epoch	1
LLaVA Image Size	576
Optimizer	SGD
Learning Rate	$2e - 5$
λ_{VM}	0.1
BF16	True
LR scheduler	cosine
Vision Tower	openai/clip-vit-large-patch14-336
Language Model	lmsys/vicuna-7b-v1.5

Table 4. Hyperparameters to train our model

776 **C. Performance of the segmentation maps in-
777 herently provided by our model**

To evaluate the segmentation ability of our FiVL model, we evaluated Intersection-Over-Union (IoU) on a subset of 10,000 images from the GQA-val dataset. For each sample, we perform an inference using the baseline LLaVA-7b and our model. From the outputs, we retrieve the visual logits for each visual token, we assigned a text token from the vocabulary corresponding to the maximum logit probability, referred to as the max-v token. By aggregating all

image tokens associated with each max-v token, we effectively generated a segmentation mask for each represented text token, like describe in Section 5.1. Additionally, as ground truth to compare against, we employed Grounded-SAM to produce segmentation maps given each max-v token. Grounded-SAM was implemented using the IDEA-Research/grounning-Dino-Tiny model with thresholds set at 0.2, 0.4, and 0.6, followed by facebook/sam-vit-huge with a threshold of 0.0. The Intersection over Union (IoU) score was computed between the FiVL-generated segmentation masks and the corresponding Grounded-SAM masks to quantitatively assess alignment. To provide a comparative analysis, we also computed IoU scores for the segmentation masks produced by the baseline model. As detailed in Table 5, across all thresholds, FiVL generated approximately 7 times fewer max-v tokens per image compared to the baseline model (column "#tokens/sample"), indicating more concise and semantically meaningful segmentation. FiVL also showed significant improvement in average IoU scores (column IoU), increasing approximately three times: from 0.05 to 0.18 at a threshold of 0.2, from 0.06 to 0.21 at 0.4, and from 0.09 to 0.24 at 0.6, showcasing its superior ability to generate precise and coherent segmentation masks. In general, across all thresholds, the baseline generates significantly more max-v tokens per image, resulting in a higher number of samples with segmentation maps found by Grounded-SAM (column samples). Finally, the percentage of tokens processed by Grounded-SAM is substantially higher for our model compared to the baseline (column processed), indicating that the max-v tokens retrieved by our model were more meaningful than those from the baseline. Figure 10 shows the segmentation maps we obtained for the max-v token describing each image. For example for the example 10a, we computed the argmax of the tokens highlighted in red, and it corresponded to the token "bear" in the vocabulary

	Thresh	IoU	# tokens/sample	# samples	#processed
Baseline	0.2	0.05	73.3	10,000	0.89
Our Model		0.18	10.3	10,000	0.96
Baseline	0.4	0.06	73.3	10,000	0.40
Our Model		0.21	10.3	9,983	0.65
Baseline	0.6	0.09	73.4	9,326	0.08
Our Model		0.24	10.6	8,604	0.26

Table 5. Performance of the segmentation maps inherently provided by our model

822 **D. Additional Visual Reliance evaluations**

To expand our benchmarks and gain a broader understanding of model/benchmark performance, we evaluated five additional models on FiVL-VQAv2, FiVL-POPE, and FiVL-GQA. This helps to assess the generalizability of our approach across more models. Table 6 shows our results for



(a) Bear



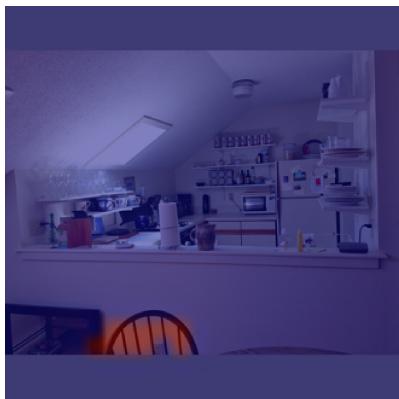
(b) Bird



(c) Birds



(d) Bott



(e) Chair



(f) Dog



(g) People



(h) Train



(i) Water

Figure 10. Segmentations produced inherently by our model. Each figure corresponds to the max-v token specified in caption. Max-v token being the token realizing the maximum for each highlighted patch

828 the 2 models presented in Section 5.2 as well as for LLaVA-
829 v1.5-7b [17], GPT4o [25], BLIP-2 [13], Pixtral-12B [3]
830 and Phi3-Vision[1], which are state-of-the arts multimodal

models. In bold, are the highest visual reliance scores per model and across all benchmarks. Underlined are the highest visual reliance scores across models, given a benchmark.

831
832
833

834 The results unanimously indicate that, among all models,
 835 VQA-v2 appears to be the benchmark that relies most heavily
 836 on the image to answer the questions. Looking at the average
 837 performance per model across benchmarks (last column), we observe that GPT4o relies most heavily on the
 838 image as a reference for answering, followed by Pixtral-
 839 12B. We can also observe that BLIP-2 is minimally affected
 840 by the adversarial attack indicating that it does not heavily
 841 rely on the image to generate answers. This lack of reliance
 842 could explain its low accuracy on the original images for
 843 this benchmark (30 % compared to approximatively 80%
 844 for the other models). As a result, perturbing the images
 845 had little impact on its performance (29 %).
 846

	VQA-v2	GQA	POPE	Average
Qwen2-VL-7B	0.64	0.39	0.47	0.50
LLaVA-13B	0.72	0.33	0.49	0.51
LLaVA-7B	0.56	0.31	0.47	0.45
GPT4o	0.74	<u>0.63</u>	0.49	<u>0.62</u>
BLIP-2	0.52	0.23	0.03	0.26
Pixtral-12B	0.75	0.58	0.42	0.58
Phi3-V	0.60	0.33	<u>0.54</u>	0.49
Average	0.65	0.40	0.42	-

Table 6. Visual reliance scores: % of drop in performance using FiVL bounding boxes for perturbations. In **bold**: highest scores across benchmarks. Underlined: highest scores across models

FiVL-Instruct system prompt

A multimodal instruction-following dataset used for visual instruction tuning and it contains an image and a conversation. The conversation is constructed from a few turns of questions and answers regarding the image.

Given only a question and answer pair: identify short expressions from the answer which could not be generated without the image.

The expression

- expresses a visual content from the image.
- should be as short as possible.
- should not be longer than 4 words
- should not include punctuation
- should no include reference to the image

Unrelated expressions should be separated by the following string:
 ":::"

Don't add any additional information to the prompt.

For example:

Q: What are the giraffes doing in the image? <image>

A: The baby giraffe is walking next to the mother giraffe, both moving through the open area of their enclosure

The output should be as following:
 baby giraffe:::mother giraffe :::open area of their enclosure

Identify the tokens for the following:

Q: {question}

A: {answer}

Figure 11. FiVL-Instruct system prompt used for training datasets

FiVL system prompts for the evaluation datasets

A multimodal instruction-following dataset used for visual instruction tuning and it contains an image and a conversation. The conversation is constructed from a few turns of questions and answers regarding the image.

Given only a question and answer pair: identify short expressions from the answer or the question which could not be generated without the image.

The expression

- should hypothetically express an immediate visual content from image. Thus, yes/no is NOT an expected expression, and some pronouns like "this", "that", "there", and "those" are not expected expressions.
- should be as short as possible.
- should not be longer than 4 words.
- should not include punctuations.
- should not include reference to the image, like: "the man wearing a blue suit in the image" it should be: "man wearing blue suit".

Unrelated expressions should be separated by the following string:
"::::".

If there is no expressions found, your response must be "N/A".

Do not add any additional information to the prompt.

<EXAMPLES>

Identify the expressions for the following:

Q: "{question}"
A: "{answer}"

FiVL-VQAv2 - Examples**Example 1:**

Q: "Considering the limited space and packed fixtures, what suggestions can be provided for organizing the bathroom to optimize its use?"

A: "To optimize the use of the small bathroom with tightly packed fixtures, there are some organization solutions that can be recommended: Install wall-mounted storage: By using vertical space on the walls and areas above the toilet, you can create additional storage for toiletries and other bathroom essentials without taking up floor space. Utilize narrow shelves or over-the-door organizers: Placing narrow shelves or over-the-door organizers on available narrow spaces can help with the organization of toiletries and other small items. Use multi-functional items: Opt for items that serve more than one purpose, such as a toilet paper holder with a shelf or a mirror with built-in storage. Opt for compact accessories: Choose smaller or more compact versions of bathroom essentials, like toothbrush holders or soap dishes, to maximize space on countertops and around the sink. Declutter regularly: Frequently assess your bathroom supplies and remove any items that are not in use or have expired to keep the bathroom tidy and make the most of the limited space. By following these suggestions, the bathroom can appear less cluttered, and occupants can make better use of the available space."

Key Expressions: small bathroom:::tightly packed fixtures:::vertical space on the walls:::areas above the toilet:::over-the-door organizers:::narrow spaces:::space on countertops and around the sink

Example 2:

Q: "How many shadows are on the ground?"

A: "3"

Key Expressions: shadows:::3

Figure 12. FiVL-Instruct system prompt used for the evaluation datasets

Figure 13. Examples for FiVL-VQAv2

FiVL-GQA and FiVL-POPE - Examples

Example 1:

Q: "Considering the limited space and packed fixtures, what suggestions can be provided for organizing the bathroom to optimize its use?"
A: "To optimize the use of the small bathroom with tightly packed fixtures, there are some organization solutions that can be recommended: Install wall-mounted storage: By using vertical space on the walls and areas above the toilet, you can create additional storage for toiletries and other bathroom essentials without taking up floor space. Utilize narrow shelves or over-the-door organizers: Placing narrow shelves or over-the-door organizers on available narrow spaces can help with the organization of toiletries and other small items. Use multi-functional items: Opt for items that serve more than one purpose, such as a toilet paper holder with a shelf or a mirror with built-in storage. Opt for compact accessories: Choose smaller or more compact versions of bathroom essentials, like toothbrush holders or soap dishes, to maximize space on countertops and around the sink. Declutter regularly: Frequently assess your bathroom supplies and remove any items that are not in use or have expired to keep the bathroom tidy and make the most of the limited space. By following these suggestions, the bathroom can appear less cluttered, and occupants can make better use of the available space."

Key Expressions: small bathroom:::tightly packed fixtures:::vertical space on the walls:::areas above the toilet:::over-the-door organizers:::narrow spaces:::space on countertops and around the sink

Example 2:

Q: "Is there a snowboard in the image?"

A: "no"

Key Expressions: snowboard

Figure 14. Examples for GQA and POPE prompts

847 E. API of the manual evaluation

848 Figure 15 shows the API used for the manual evaluation
849 done on FiVL-Instruct. Given a question, an answer (on
850 the left) and an image with a segmentation mask (on the
851 right), the annotator had to answer the 3 following yes/no
852 questions: is { key expression } correctly represented in the
853 mask? Is {key expression} a significant word in the answer?
854 Is this example generally good to be included in the dataset?

The screenshot displays a web-based interface for dataset evaluation. On the left, a 'Conversation' window shows a human query: 'Describe the scene where the man is surfing.' and a GPT-generated response: 'The scene shows a man in a red and black wetsuit surfing a blue wave near a rocky shoreline. He is riding a surfboard in the ocean and skillfully navigating the challenging area near the rocks.' In the center, there are two images: a 'Padded Original Image' showing a surfer on a wave, and a 'Segmented Image' showing the same scene with a white silhouette of the surfer on a black background. Below these images is a list of three questions for annotation:

1. Is 'riding a surfboard' correctly represented in the mask?
2. Is 'riding a surfboard' a significant word in the answer?
3. Is this example generally good to be included in dataset?

For each question, there are four possible annotations:

Correct	Incorrect
Key Token	NOT Key Token
Yes, good data	No, bad data

At the bottom right is a 'Next Data' button, and at the bottom center is a 'Preview Labeled Data' button. To the left of the questions, there is explanatory text about significant words and examples.

Significant word: a word that relates to the question and couldn't be elicited without the image
Good example: the answer is relevant and makes sense
Bad example: gibberish, unrelated answer

Preview Labeled Data

Figure 15. Web user interface for our dataset evaluation