# LVLM-COUNT: ENHANCING THE COUNTING ABILITY OF LARGE VISION-LANGUAGE MODELS

**Anonymous authors**
Paper under double-blind review

## ABSTRACT

Counting is a fundamental skill for various visual tasks in real-life applications, requiring both object recognition and robust counting capabilities. Despite their advanced visual perception, large vision-language models (LVLMs) struggle with counting tasks, especially when the number of objects exceeds those commonly encountered during training. We enhance LVLMs' counting abilities using a divide-and-conquer approach, breaking counting problems into sub-counting tasks. Unlike prior methods, which do not generalize well to counting datasets on which they have not been trained, our method performs well on new datasets without any additional training or fine-tuning. We demonstrate that our approach enhances counting capabilities across various datasets and benchmarks.

## 1 INTRODUCTION

Counting is a key cognitive task with broad applications in industry, healthcare, and environmental monitoring (De Almeida et al., 2015; Guerrero-Gómez-Olmedo et al., 2015; Paul Cohen et al., 2017; Lempitsky & Zisserman, 2010). It improves manufacturing, inventory, and quality control, ensures safety in medical settings, and helps manage resources in environmental efforts (Wang & Wang, 2011; Zen et al., 2012; Arteta et al., 2016). Recent advancements in prompt-based models enable counting of unlimited object varieties without visual exemplars. Although current trained, and text-prompt-based, counting models by Dai et al. (2024); Amini-Naieni et al. (2024) perform well on the datasets they are trained on, they face the following challenges. First, they require fine-tuning on new datasets. Second, since the concepts in the counting datasets are limited, these models do not generalize well to counting questions that involve complex reasoning. There are also training-free models by Shi et al. (2024), but overall they have weaker performance compared to trained models. On the other hand, large vision-language models (LVLMs), such as GPT-4o (Achiam et al., 2023), which also do not need dataset specific training, show very good performance in counting low numbers of objects, usually less than 20, across most datasets. However, their performance deteriorates for larger numbers of objects, regardless of the dataset.
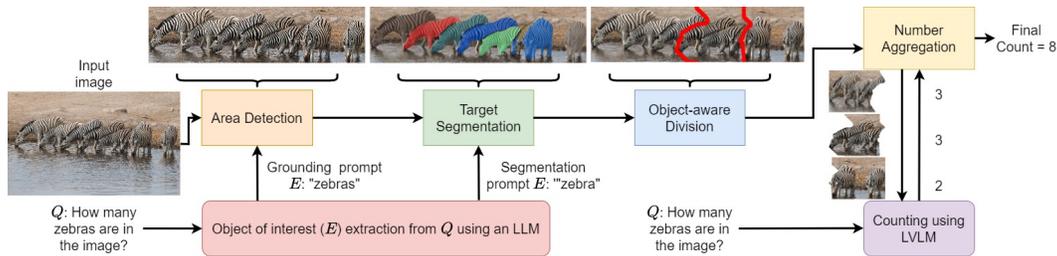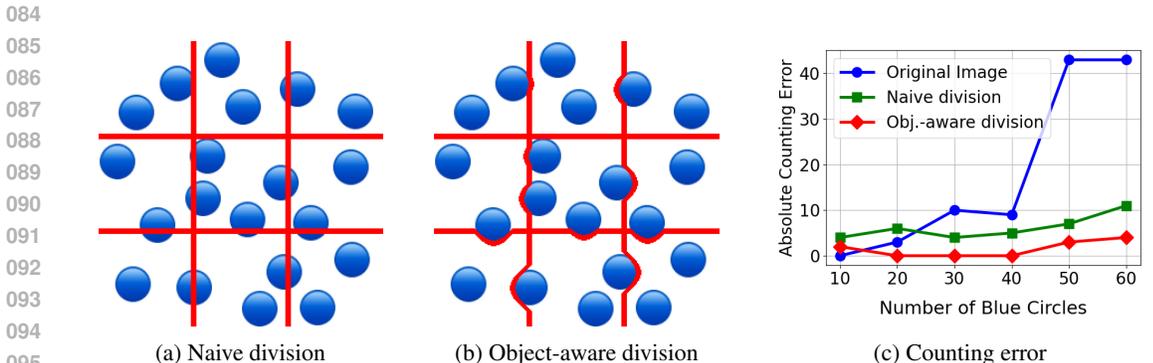


Figure 1: Illustration of our proposed pipeline. First, the area of interest, such as a set of zebras, is detected by extracting the object of interest ($E$) from a prompt question ($Q$) using a large language model (LLM) which is the same as LVLM in our work. Then, $E$ and the image are provided as input to an area detection model, such as the one by Liu et al. (2023). Second, any objects corresponding to $E$ are segmented. Third, in the object-aware division step, we use the segmentation masks to divide the areas of objects of interest without cutting through the objects. Finally, the number of objects of interest in each sub-image is computed using an LVLM, and the results are aggregated.

We enhance the accuracy of LVLMs to count objects in images by leveraging their reasoning power within a divide-and-conquer method. The LVLMs allow us to handle diverse objects and complex counting questions, while our divide-and-conquer method alleviates challenges associated with counting large numbers of objects. Inspired by prior work on the rapid and accurate estimation of small quantities by Chattopadhyay et al. (2017), we divide an image into sub-images, and prompt the LVLM to count the objects of interest in each sub-image. The counts from each sub-image are then aggregated to make the final prediction. Our proposed workflow is illustrated in Figure 1.

Initially, in our pipeline, the object of interest is extracted from the input question using an LLM. The area containing the object of interest is detected in the image by a grounding model, such as Liu et al. (2023), and then cropped. The cropping step is important since it removes irrelevant context from the image. Secondly, using an object detection model by Liu et al. (2023), and a segmentation model by Kirillov et al. (2023), the segmentation masks of the objects of interest are created. Thirdly, we use a mechanism that divides the image into multiple sub-images without cutting through the objects of interest. We call this mechanism object-aware division. The division positions can be either set manually or determined automatically using unsupervised and non-parametric methods. In the latter case, we automatically select the division positions based on object masks, then we treat the object-aware division as a path-finding problem, avoiding objects as obstacles. A black-white image is built by converting all the masks into black and the rest of the image into white pixels. The binary image is converted into a graph where only white pixels are connected as nodes. Using the $A^*$ algorithm (Russell & Norvig, 2016), a path is found from one end to the other end of the image, ensuring objects remain intact. Finally, using an LVLM as a counting tool, the objects of interest in the sub-images are counted and aggregated. Our contributions are summarized below:

1. We propose LVLM-Count that leverages the strengths of large vision-language models (LVLMs) in visual perception for counting. Our method is a prompt-based counting approach that, in addition to simple counting problems, can handle complex cases as well. Our method does not require any additional training or fine-tuning on any dataset. LVLM-Count outperforms existing state-of-the-art models across various datasets and benchmarks.

2. We propose a solution for object-aware division. Accurate division is crucial, as parts of cut objects can lead to over-counting (see Figure 2a). To the best of our knowledge, this is the first method to divide images without cutting through objects of interest.



(a) Naive division      (b) Object-aware division      (c) Counting error

Figure 2: Comparison of the naive and the object-aware division. The objects of interest are the circles. In Figure 2a, we illustrate a naive division of the input image, which is divided into equally sized sub-images with straight lines. In Figure 2b, we illustrate the object-aware division, which avoids cutting through circles. In Figure 2c, we illustrate the counting error of GPT-4o for images with randomly positioned circles. The absolute counting error is the absolute difference between the ground truth and the number predicted by GPT-4o. The results are averaged over three trials.

As a minor contribution, we create a new benchmark to address two drawbacks of existing benchmarks. Prior datasets either feature simple counting tasks, e.g., counting "strawberries", or include complex questions with small numbers of objects. To address both issues, we develop a challenging benchmark for counting emoji icons. The subtle variations within emoji classes make this benchmark uniquely difficult. Since none of the models have been exposed to it, this benchmark serves as a fair test of the performance of counting methods for various attributes and complex concepts.

## 2 RELATED WORK

Early counting models, referred to as class-specific, targeted counting problems for certain categories (Arteta et al., 2016; Babu Sam et al., 2022; Mundhenk et al., 2016; Xie et al., 2018), such as cars, people, or cells. Later, with the emergence of stronger vision models and large-scale datasets, class-agnostic methods were proposed that could count objects from a wide variety of categories. However, most existing class-agnostic, or open-world, models require visual exemplars of the target objects (Đukić et al., 2023; Gong et al., 2022; Lin et al., 2022; Liu et al., 2022; Lu et al., 2019; Nguyen et al., 2022; Ranjan et al., 2021; Shi et al., 2022; Yang et al., 2021; You et al., 2023). The concept of divide and conquer has also been used in early work(Xiong et al., 2019; Chattopadhyay et al., 2017). However, this early work requires training dedicated models to utilize this concept.

**Text-based counting-specific trained models.** With the advent of vision-language foundation models such as CLIP and GroundingDINO, text-based open-world methods have been proposed. Leveraging the rich textual and visual feature extraction capabilities of foundation models, obtained through web-scale training, the text-based counting methods by Amini-Naieni et al. (2023); Dai et al. (2024); Kang et al. (2024); Amini-Naieni et al. (2024) have started to demonstrate comparable or superior accuracy. GroundingREC by Dai et al. (2024) is an open-world model built on top of GroundingDINO (Liu et al., 2023), and introduces an additional task called referring expression counting. Concurrently, Amini-Naieni et al. (2024) proposed a method that also builds on GroundingDINO but adds an extra image-text fusion module in the input, enabling the model to accept text and/or visual exemplars to determine the target.

**Models without counting-specific training.** Shi et al. (2024) introduce TFOC, a counting model that does not require any counting-specific training. Instead, they cast the counting problem as a prompt-based segmentation task, using SAM (Kirillov et al., 2023) to obtain segmentation maps that determine the output number. Another group of models that do not need further training to count are LVLMs. State-of-the-art models such as GPT-4o (Achiam et al., 2023), Gemini 1.5 Pro (Reid et al., 2024), and Claude 3.5 Sonnet (Anthropic, 2024) show strong performance in counting small numbers of objects, although their performance degrades with larger numbers of objects. In comparison to prior work, our method (LVLM-Count) is an open-world, text-based counting approach that does not require any counting-specific training. Although we use pre-trained GroundingDINO and SAM in our method, we differ from prior work by not treating the counting task as a segmentation or detection problem. Instead, we use LVLMs as a counting tool.
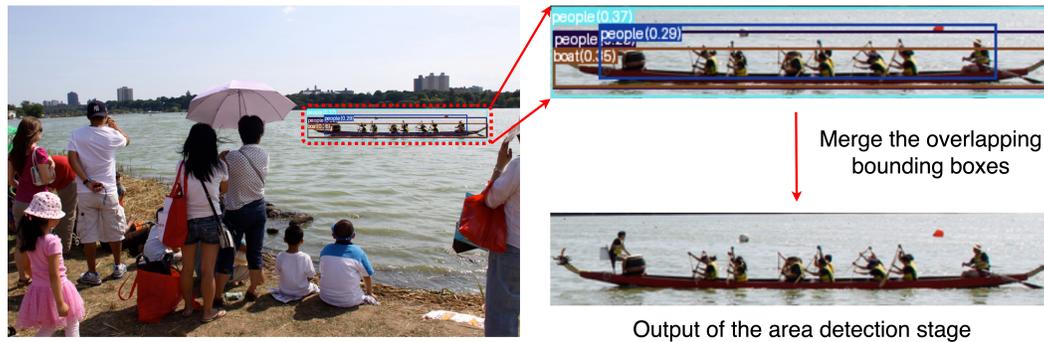
## 3 LVLM-COUNT

Our proposed method aims to answer counting questions by dividing an image into sub-images while avoiding cuts through objects of interest. LVLM-Count consists of four key stages. First, in the "Area Detection" stage, we localize areas containing relevant objects. Second, in the "Target Segmentation" stage, we identify and segment these local areas. Third, in the "Object-aware Division" stage, we divide the localized areas into sub-images without cutting through the segmented objects. Finally, the LVLM counts the target objects in each sub-image and aggregates the results. Figure 1 illustrates the workflow of our method, which we will detail in the following subsections. Note that in this section, we illustrate each stage using example images. These images are are used solely for illustration purposes. LVLM-Count may not outperform other methods on these examples.

### 3.1 AREA DETECTION

In this part of the pipeline, we assume that we are given a counting question $Q$ along with an image. The question $Q$ contains an expression $E$ that specifies a set of objects of interest. The expression $E$ distinguishes these objects from objects of other categories or the same category but with different attributes present in the image. By employing an LLM, the expression $E$ is extracted from $Q$. For example, let $Q$ be "How many people are in the boat?". $Q$ is given to an LLM, which is prompted to return the expression $E$ "people in the boat", referring to the objects we want to count. After $E$ is extracted, it is given as input to GroundingDINO along with the image. The output of GroundingDINO is a set of bounding boxes that have relevance to $E$ beyond a certain threshold. These bounding boxes often overlap and typically contain repeated objects. Thus, all the overlapping output bounding boxes are merged. After merging, a set of non-overlapping areas of interest may
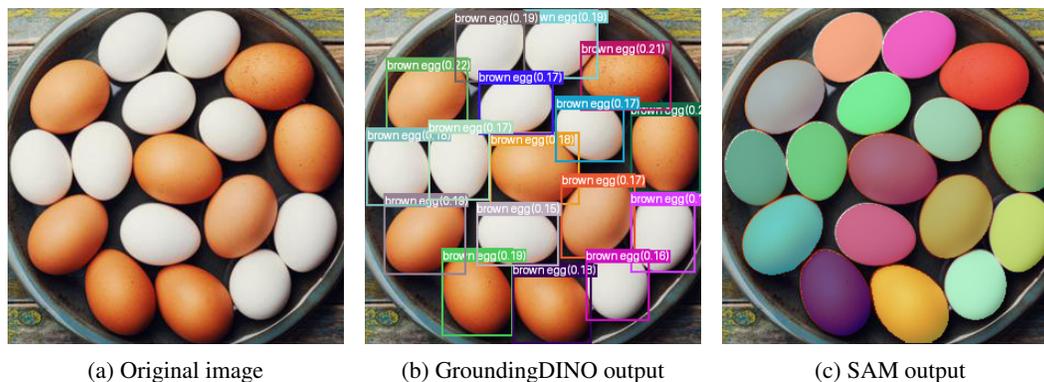
remain. We consider the non-overlapping areas as "detected areas", which are then cropped to be passed to the next stage. Note that the area detection stage is important as it extracts the area with the most relevant context for the counting question. This process is illustrated in Figure 3.



Figure 3: Illustration of the area detection step of LVLM-Count. For this image, $Q$ is set to "How many people are in the boat". The LLM that is used in this step returns an $E$ which is "people in the boat". $E$ and the original image are given as input to GroundingDINO, which returns some bounding boxes (left and upper right images) that are merged to form the final detected area.

## 3.2 Target Segmentation

The cropped images from the first stage contain objects of interest, and the ultimate goal is to divide them without cutting through those objects. However, a prerequisite for implementing such a mechanism is to first detect and localize the objects of interest. Each cropped image is fed into an open-world detection model along with $E$. The output of the open-world detection model produces a bounding box for each object of interest. The bounding boxes are then given as input to a segmentation model, which returns segmentation masks for the objects within each bounding box. We illustrate an example of this process in Figure 4.



(a) Original image       (b) GroundingDINO output       (c) SAM output

Figure 4: Illustration of the target segmentation step of LVLM-Count. The goal is to produce all the instance masks for $E$ set to "brown egg". Figure 4a, together with $E$, is given as input to GroundingDINO, which produces the output shown in Figure 4b. Figure 4b is then given as input to SAM, which produces the output shown in Figure 4c.

**How to determine the bounding boxes.** To determine the bounding boxes, we use GroundingDINO and set the bounding box probability threshold to a low value to avoid missing any objects. The bounding boxes alone cannot help with the object-aware division of the cropped images due to their rigid structure, which includes redundant areas in the vicinity of the object and, in the worst case, overlaps with other bounding boxes. Our goal is to precisely locate the pixels of an object of interest.

4

**How to determine the segmentation masks.** We use a pre-trained segmentation model, i.e., SAM (Kirillov et al., 2023), for the segmentation task, which accepts bounding boxes as prompts. It produces a mask covering the most prominent object within a bounding box. We run non-maximum suppression on the masks produced by SAM to remove the masks corresponding to less certain bounding boxes that have large overlaps with other masks. The masks produced for each cropped image are then passed to the next stage. For GroundingDINO we set the bounding box threshold to a very low value to avoid missing any objects of interest, especially if the objects belong to a rare category that GroundingDINO has encountered less frequently during pre-training. As shown in Figure 4, a low threshold value can lead to false positives, making the number of masks unreliable for the counting problem. Nonetheless, the false positive masks have no negative effect on the object-aware division in the next stage. The only consequence is that, in addition to the target masks, the division paths will not cut through the false positive masks either. In other words, extra masks are of no concern as long as all the target masks are identified.

### 3.3 OBJECT-AWARE DIVISION

In this stage, the cropped image is divided into appropriate sub-images so that no object of interest is cut by the dividing paths. The core idea is that the dividing paths should not intersect the pixels covered by the masks corresponding to the objects of interest. This step consists of two sub-steps. First, we decide the starting and ending points of the paths. Second, we draw the paths. Below, we describe how we approach these two sub-steps.

**How to determine the starting and ending points of the paths.** We consider the first sub-step as a hyper-parameter, which can be set in various ways depending on the type of dataset. We utilize two approaches in our experiments. The first approach is unsupervised and non-parametric, where we use a clustering method to obtain the start and end points of the paths. We describe the clustering approach in more detail below. The second approach is to simply use a pre-determined number of equidistant points. This is helpful when we know a priori that objects are uniformly distributed across the image. In this case, we fix $k$ equidistant vertical and horizontal coordinates in the image as the points of the dividing paths for the entire dataset. In our experiments, we specify which approach we use for each dataset. Nonetheless, we will illustrate in the experimental results that even if we do not treat this as a hyperparameter and always use the unsupervised and non-parametric approach, LVLM-Count is quite effective.
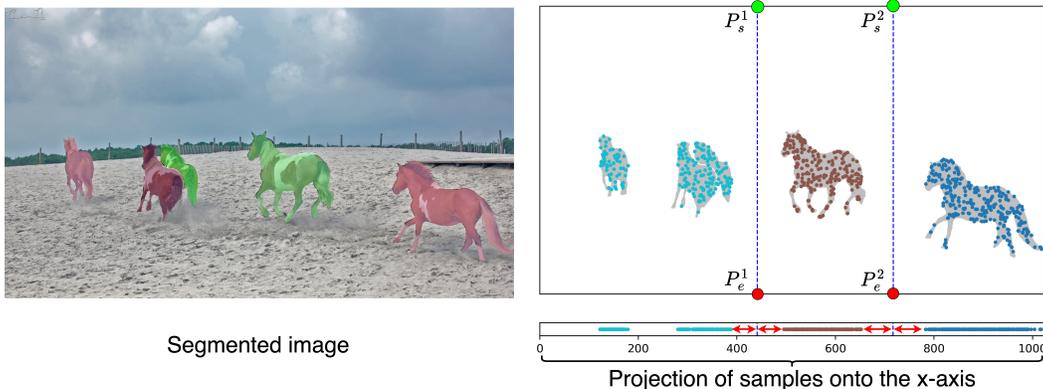


Segmented image | Projection of samples onto the x-axis

Figure 5: Illustration of the unsupervised and non-parametric method to obtain the division points $(P_s^1, P_e^1)$, and $(P_s^2, P_e^2,)$. The object of interest is "horse". Firstly, the target masks are produced. Then, a few pixels are sampled (shown as points inside the segmented objects) from the pixels composing each mask. The samples are projected onto the $x$-axis. The projected points are clustered using mean-shift clustering. The point in the middle of two consecutive clusters is considered a vertical division point. The straight vertical lines are obtained using our path-finding method.

**Description of the unsuperivsed and non-parametric approach.** A few pixels are sampled from each of the masks. To determine the points of the vertical division paths, the samples taken from the masks are projected onto the $x$-axis. The projected points are automatically clustered using a

non-parametric mean-shift algorithm. Once the clusters are identified, the point between the point with the highest $x$ value in one cluster and the point with the lowest $x$ value in the subsequent cluster is considered the $x$-coordinate of a vertical division path. Using this technique, we obtain the appropriate coordinates for the vertical paths, as well as the number of paths. For example, if there is only one cluster, no vertical division is required, and if there are two clusters, one vertical path will divide the image into two parts. We illustrate this approach in Figure 5. Note that we chose the vertical axis purely for convenience in illustration. The same process can be applied along the $y$-axis to obtain the $y$-coordinates of the horizontal division paths.

**How to draw the paths.** In effect, knowing the $x$-coordinate of a vertical path means that the coordinates of its endpoints are known. In particular, assuming height $h$ for an image crop, we consider $P_s = (x, 0)$ and $P_e = (x, h)$ as the start and end points, respectively. In an ideal case where there are no masks in the path of a straight line connecting the two points, this line will be drawn by connecting all the pixels on the straight path. However, there are potential masks that can be considered obstacles blocking the path. In other words, beginning from $P_s$, the line needs to go around these obstacles to reach $P_e$. Consequently, we treat this as a 2-dimensional path-finding problem. To solve the problem, we build a 2D binary map, $I_B$, where the pixels covered by the masks are turned into black, indicating them as obstacles, and all the other pixels are turned into white, showing they are open for passage. This binary image $I_B$ is mapped into a graph $G$, where each white pixel is a node, and it is connected to all of its white neighboring pixels. We use the $A^*(G, P_s, P_e, g)$ search algorithm to find a path that connects $P_s$ to $P_e$, where the heuristic $g$ is set to be Manhattan distance. The output of $A^*$ is a set of connected pixels that go around the obstacles and connect $P_s$ to $P_e$, creating an object-aware division path, as shown in Figure 6. The path-finding algorithm is run for all division coordinates. Finally, we draw the image contours based on these division paths and take the area surrounded by each contour as a resulting sub-image.
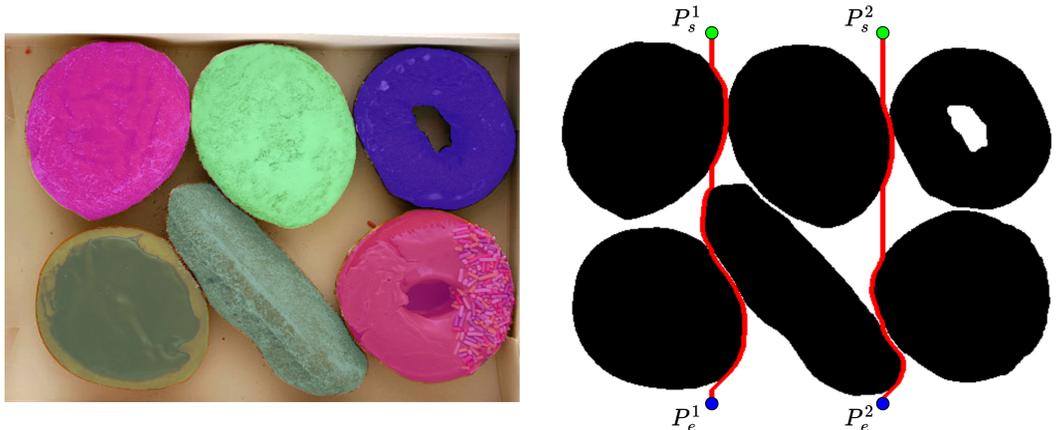


Figure 6: Illustration of object-aware division. The masks for "donut" are obtained and turned into a black-and-white image. A dividing path is found by connecting $P_s$ to $P_e$ using the $A^*$ search algorithm in a graph that corresponds to the binary image, where the only nodes in the graph are white pixels, which are fully connected. The path is mapped back to the pixel domain.

## 3.4 TARGET COUNTING

At this point, all the sub-images obtained from the cropped areas are gathered. Since these are small partitions of the original image, they might lack the desired resolution. Thus, using a super-resolution model can alleviate this problem for the sub-images. In this work, we use Real-ESRGAN (Wang et al., 2021). Then, question $Q$ and each sub-image are given as input to an LVLM. At the end of the loop, the recorded numbers for the sub-images are aggregated to form the final answer. [1].

---

[1]Note that using super-resolution before counting the objects in the sub-images is also the default behavior of our method. For a specific dataset where we observe that the resolution is preserved after division, we deactivate super-resolution, as the cost and usage of LVLM APIs is much higher. Finally, for images with a

# 4 EXPERIMENTS

In this section, we present the performance results of our method on a counting-specific dataset, an open-ended counting benchmark with two types of questions, simple and complex, and a challenging counting benchmark that we propose using emoji icons. We compare the results to state-of-the-art models which have been specifically trained on counting datasets, and we also compare to state-of-the-art models which have not been trained on counting datasets.

## 4.1 DATASETS AND BENCHMARKS

**FSC-147 (Ranjan et al., 2021).** FSC-147 is a counting dataset that contains 6135 images, spanning 147 different object categories such as kitchen utensils, office supplies, vehicles, and animals. The number of objects in each image ranges from 7 to 3731, with an average of 56 objects per image. The dataset is split into training, validation, and test sets. A total of 89 object categories are assigned to the training set, 29 to the validation set, and 29 to the test set, with different categories in each split. The training set contains 3659 images, with the validation and test sets containing 1286 and 1190 images, respectively. For each image in the test set, a single category name is given, and the expected output is the number of instances from the category.

**Open-ended Counting Benchmark.** TallyQA (Acharya et al., 2019) is an open-ended counting dataset that includes complex counting questions involving relationships between objects, attribute identification, reasoning, and more. TallyQA is quite a large dataset, with the train set having $249,318$ questions and the test set having $22,991$ simple and $22,991$ complex counting questions. Please refer to Appendix C and Figure 9 for more information on simple and complex categorization of the questions. The number of objects in each image ranges from 0 to 15. The dataset has a heavy bias towards a low number of objects (see Figure 8 in the Appendix). To alleviate this bias and create a benchmark for efficiently measuring the simple and complex open-ended capabilities of a counting model, we randomly sample 10 questions per ground truth count. This results in 155 simple and 149 complex open-ended counting questions in total. It is important to note that the bias in TallyQA is pronounced; for most ground truth values greater than 10, there are fewer than 10 samples available in the entire test set.

**Emoji-Count.** Although TallyQA addresses the scarcity of complex counting questions in prior datasets to a certain degree, the range of target objects is limited, spanning from 0 to 15. To our knowledge, no counting benchmark exists for large numbers involving complex reasoning. To this end, we propose a challenging counting benchmark using emoji icons. From the 1816 standard emoji icons, we remove those that directly overlap with concepts demonstrated by other icons. We then group the remaining 1197 icons into 82 classes. In each class, there are icons from the same or similar object categories, but with subtle differences that require complex reasoning to distinguish. For each of the 82 classes, an empty $1024 \times 1024$ image is first created. This image is filled with six categories chosen randomly from the class, with each category having a random count between 30 and 50 in the image. We illustrate an example of this dataset in Figure 7.
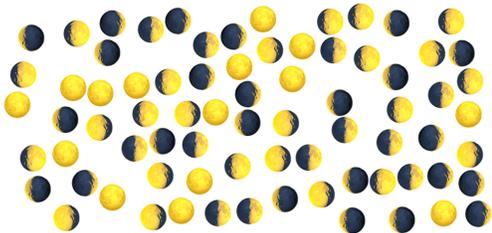
## 4.2 RESULTS

The following discusses the numerical results of our experiments with LVLM-Count on each benchmark described in Section 4.1. For visual examples of LVLM-Count's performance on each benchmark, see Appendix M. Additionally, for the ablation study and experiments on the PASCAL VOC dataset (Everingham et al., 2015), see Appendix A and Appendix K, respectively.
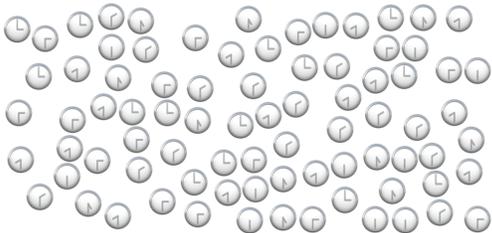
**FSC-147.** We compare the performance of LVLM-Count to an extensive list of state-of-the-art counting models on the test set of the FSC-147 dataset. We also include two baselines: i) taking the number of target segmentation masks as the final answer, and ii) giving the output of the target segmentation stage to base GPT-4o and asking it to count the masks. We run different experiments using GPT-4o, Gemini 1.5 Pro, and an open-source model Qwen2 VL 72B AWQ (Yang et al., 2024) as LVLMs. Since objects in the FSC-147 dataset are uniformly distributed across scenes, we deactivate the cluster-based method for finding division points and instead use two vertical and two

---

very large number of objects, sometimes LVLMs refuse to count, citing the large number. In those cases, an "estimate" prompt is given instead of a "count" prompt.

(a) Q: How many waning gibbous moons are there in the image? Answer: 13.

(b) Q: How many clocks at time "two-thirty" are there in the image? Answer: 15.

Figure 7: Illustration of two challenging cases from Emoji-Count. Note that, for convenience in visualization, a smaller version of the real images is depicted here. The real images are $1024 \times 1024$ and contain many more icons. In Figure 7a, the class name is "Moon Phase". To predict the correct answers for the question, the model needs accurate information about the different moon phases in addition to strong counting capability. In Figure 7b, the class name is "Clock Time". In this case, the model needs to be able to read the time shown on the clocks in addition to strong counting capability.

horizontal object-aware division paths with equidistant start and end points. We also deactivate the super-resolution. We further run an experiment with GPT-4o as the LVLM that uses the cluster-based approach along both horizontal and vertical axes to automatically determine the start and end points of the division paths. This is done to measure the performance of our method when prior knowledge is not used to set the start and end points. The rest of the process follows the procedure described in Section 3. The expression $E$ used in different stages of our method is the category name provided in the test set. A simple $Q$ in the form of "How many $E$ are there? If you don't see any, say zero." is built and given as a text prompt to the LVLM during the counting stage. The results are shown in Table 1. For this dataset, the mean absolute error (MAE) and root mean square error (RMSE) are reported. We observe that our method enhances the performance of all three LVLMs in terms of MAE. With our method, all three LVLMs outperform TFOC, which which is a training-free method. Interestingly, although the base Qwen2 VL 72B AWQ is not as powerful as its commercial counterpart, GPT-4o, it performs almost on par with GPT-4o when both use our pipeline. The results also show that even without leveraging prior information to set the start and end points of the division paths, our pipeline remains effective.

Table 1: Evaluation of state-of-the-art models on the test set of the FSC-147 dataset. The column "Trained Model" indicates if a model has been trained on FSC-147. In all tables, the results for base LVLMs and LVLM-Count are reported over three trials. Columns marked with $\Delta$ show the improvement that LVLM-Count brings over the base LVLM that it uses. For accuracy metrics, refer to Table 9. Additionally, MAE analysis across different intervals of ground truth values is provided in Appendix H.

| Method | Trained Model | MAE ↓ | Δ | RMSE ↓ | Δ |
|---|---|---|---|---|---|
| TFOC (Shi et al., 2024) | ✗ | 24.79 | - | 137.15 | - |
| VLCounter (Kang et al., 2024) | ✓ | 17.05 | - | 106.16 | - |
| CounTX(Amini-Naieni et al., 2023) | ✓ | 15.88 | - | 106.29 | - |
| DAVE$_{prm}$ (Pelhan et al., 2024) | ✓ | 14.90 | - | 103.42 | - |
| CountGD (Amini-Naieni et al., 2024) | ✓ | 14.76 | - | 120.42 | - |
| GroundingREC (Dai et al., 2024) | ✓ | **10.12** | - | 107.19 | - |
| Number of target segmentaion masks | ✗ | 44.14 | - | 154.39 | - |
| Base GPT-4o counting the target segmentaion masks | ✗ | 38.45 | - | 38.45 | - |
| GPT-4o | ✗ | 23.75 | - | 137.39 | - |
| LVLM-Count (GPT-4o as LVLM, Cluster-based on both axes) | ✗ | 17.67 | ↓ 6.08 | 90.61 | ↓ 46.78 |
| LVLM-Count (GPT-4o as LVLM) | ✗ | 16.23 | ↓ 7.52 | **76.09** | ↓ 61.3 |
| Gemini 1.5 Pro | ✗ | 25.20 | - | 108.76 | - |
| LVLM-Count (Gemini 1.5 Pro as LVLM) | ✗ | 14.85 | ↓ 10.35 | 85.60 | ↓ 23.16 |
| Qwen2 VL 72B AWQ | ✗ | 33.45 | - | 145.90 | - |
| LVLM-Count (Qwen2 VL 72B AWQ as LVLM) | ✗ | 17.77 | ↓ 15.68 | 113.06 | ↓ 32.84 |

Table 1 shows that, although we outperform models that have not been trained on FSC-147, the best-performing models are those that have been trained on this dataset. However, in subsequent experiments, we demonstrate that, while the best-performing models on FSC-147 show acceptable generalization on benchmarks they have not been trained on, they lack the generalization ability of our method on new datasets, and they are even outperformed by the base LVLM models.

**Open-ended Counting Benchmark.** We evaluate the performance of LVLM-Count on the simple and complex questions in the open-ended counting benchmark mentioned above. For LVLM-Count, we compute the start and end points of the object-aware division paths using the unsupervised and parameter-free clustering approach described in Section 3.3 on $x$-axis. This approach is necessary because the images in this benchmark are diverse, and objects are not uniformly distributed across the image. Thus, using a fixed number of object-aware division paths with equidistant start and end points does not perform well. We also conduct an experiment using GPT-4o as the LVLM, where we use the clustering method along both axes to find the start and end points.

Since these are special cases of VQA tasks, we report exact accuracy (EA), which is the standard performance metric for VQA, in addition to MAE and RMSE. We compare our method against the base GPT-4o, Gemini 1.5 Pro, and Qwen2 VL 72B models, as well as three of the best-performing trained models from Table 1, namely GroundingREC, CountGD, and DAVE$_{prm}$[2], and the only prior training-free model, TFOC. Note that due to extremely poor performance, we do not provide the original question to these models. Instead, we assist them with the $E$ extracted in our pipeline. The results are shown in Table 2. Our method improves both MAE and EA over the base GPT-4o model on both complex and simple benchmarks.

In general, however, there is barely any improvement over base LVLMs on simple questions. This is because the questions are straightforward, and the ground truths are within [0, 15]. Based on our observations from Figure 2c and Figure 18, we expect the base LVLMs to perform better in such cases. However, for complex questions, there is consistent improvement in EA across all three LVLMs. Interestingly, despite the considerable gap between the EA of the base Qwen2 VL 72B AWQ and base GPT-4o on complex questions, using LVLM-Count enables this open-source model to achieve a higher EA than the base GPT-4o. Moreover, observe that GroundingREC, CountGD, and DAVE$_{prm}$ are outperformed by the LVLM-based models on both simple and complex categories. This is because they have not been specifically trained on TallyQA. The performance gap is more pronounced on complex questions.

Table 2: Evaluation of models on the TallyQA benchmark. Numbers in parentheses under the $\Delta$ signs show the performance difference between LVLM-Count and the base LVLM it uses. Green indicates improvement, while red represents degradation. For more details about the simple and complex types of questions, please refer to Appendix C and Figure 9. For additional accuracy measures, refer to Table 10 and 11.

| Method | Simple Questions | | | Complex Questions | | |
|---|---|---|---|---|---|---|
| | EA (%)↑(Δ) | MAE ↓(Δ) | RMSE ↓(Δ) | EA (%)↑(Δ) | MAE ↓(Δ) | RMSE ↓(Δ) |
| TFOC (Shi et al., 2024) | 6.45 (-) | 8.28 (-) | 14.79 (-) | 1.34 (-) | 12.41 (-) | 22.80 (-) |
| DAVE$_{prm}$ (Pelhan et al., 2024) | 5.81 (-) | 16.06 (-) | 35.26 (-) | 3.36 (-) | 28.36 (-) | 57.56 (-) |
| GroundingREC (Dai et al., 2024) | 23.87 (-) | 2.90 (-) | 4.71 (-) | 16.78 (-) | 5.83 (-) | 10.13 (-) |
| CountGD (Amini-Naieni et al., 2024) | 41.94 (-) | 2.37 (-) | 4.56 (-) | 5.37 (-) | 9.78 (-) | 17.21 (-) |
| Number of the target segmentation masks | 27.31 (-) | 2.60 (-) | 4.10 (-) | 15.44 (-) | 4.25 (-) | 6.78 (-) |
| Base GPT-4o counting the target segmentation masks | 25.38 (-) | 2.82 (-) | 2.82 (-) | 11.41 (-) | 4.99 (-) | 4.99 (-) |
| GPT-4o | 44.30 (-) | 1.38 (-) | 2.35 (-) | 29.08 (-) | 2.60 (-) | 4.74 (-) |
| LVLM-Count (GPT-4o as LVLM, Cluster-based on both axes) | 41.29 (↓ 3.01) | 1.82 (↑ 0.44) | 3.33 (↑ 0.98) | 28.41 (↓ 0.67) | 3.18 (↑ 0.58) | 8.91 (↑ 4.17) |
| LVLM-Count (GPT-4o as LVLM, Cluster-based on $x$-axis) | 44.73 (↑ 0.43) | 1.18 (↓ 0.20) | 2.06 (↓ 0.29) | **34.68** (↑ 5.6) | 2.28 (↓ 0.32) | 4.18 (↓ 0.56) |
| Gemini 1.5 Pro | 47.10 (-) | **1.08** (-) | **1.87** (-) | 25.73 (-) | **2.13** (-) | **3.36** (-) |
| LVLM-Count (Gemini 1.5 Pro as LVLM, Cluster-based on $x$-axis) | 45.16 (↓ 1.94) | 1.43 (↑ 0.35) | 4.72 (↑ 2.85) | 26.62 (↑ 0.89) | 2.79 (↑ 0.66) | 4.70 (↑ 1.34) |
| Qwen2 VL 72B AWQ | **49.03** (-) | 1.44 (-) | 2.74 (-) | 24.61 (-) | 3.21 (-) | 5.35 (-) |
| LVLM-Count (Qwen2 VL 72B AWQ as LVLM, Cluster-based on $x$-axis) | 41.72 (↓ 7.31) | 1.66 (↑ 0.22) | 3.41 (↑ 0.67) | 30.65 (↑ 6.04) | 2.47 (↓ 0.74) | 4.35 (↓ 1) |

**Emoji-Count.** We evaluate the performance of LVLM-Count on the Emoji-Count benchmark. The results are shown in Table 3. For LVLM-Count, we deactivate the cluster-based division points and, similar to the FSC-147 experiment, use two vertical and horizontal object-aware division paths with

---

[2]GroundingREC in Table 2 does not use the same model weights as those trained for FSC-147. It was trained on the REC-8K dataset (Dai et al., 2024), which is specifically designed for referring expression counting, making it stronger for complex counting questions. CountGD and DAVE$_{prm}$ use the same model weights as those used for FSC-147, as the authors have not provided model weights trained on a dataset with referring expressions.

equidistant start and end points. We also run an experiment with GPT-4o where LVLM-Count uses the clustering approach along both axes to determine the start and end points of the division paths. We report MAE and RMSE and compare against the base LVLMs, TFOC, DAVE_{prm}, CountGD, and GroundingREC. Note that none of the models have been exposed to this benchmark. Furthermore, it is a challenging benchmark as it requires understanding complex concepts. We observe that prior counting models perform poorly because, for any object of interest in the image, these models tend to count all the objects and cannot distinguish between different icons. Nonetheless, the base LVLMs show reasonable performance, but since the number of objects of interest is large in this dataset, all three base LVLMs are outperformed by LVLM-Count. The results for the clustering approach show that LVLM-Count is quite effective even when no prior knowledge is used to set the start and end points of the division paths.

Table 3: Evaluation of state-of-the-art models on the Emoji-Count benchmark. Columns marked with $\Delta$ show the improvement LVLM-Count brings over the base LVLM it uses. Please refer to Table 12 for more accuracy metrics.

| Method | MAE $\downarrow$ | $\Delta$ | RMSE $\downarrow$ | $\Delta$ |
|---|---|---|---|---|
| TFOC (Shi et al., 2024) | 64.64 | - | 87.45 | - |
| DAVE_{prm} (Pelhan et al., 2024) | 198.99 | - | 208.08 | - |
| CountGD (Amini-Naieni et al., 2024) | 137.93 | - | 156.80 | - |
| GroundingREC (Dai et al., 2024) | 36.16 | - | 51.88 | - |
| Number of the target segmentation masks | 82.47 | - | 107.98 | - |
| Base GPT-4o counting the target segmentation masks | 107.72 | - | 162.12 | - |
| GPT-4o | 22.51 | - | 35.94 | - |
| LVLM-Count (GPT-4o as LVLM, Cluster-based on both axes) | 11.10 | $\downarrow$ 11.41 | 24.23 | $\downarrow$ 11.71 |
| LVLM-Count (GPT-4o as LVLM) | **9.82** | $\downarrow$ 12.69 | **19.46** | $\downarrow$ 16.48 |
| Gemini 1.5 Pro | 18.17 | - | 27.83 | - |
| LVLM-Count (Gemini 1.5 Pro as LVLM) | 14.35 | $\downarrow$ 3.82 | 23.42 | $\downarrow$ 4.41 |
| Qwen2 VL 72B AWQ | 82.41 | - | 186.32 | - |
| LVLM-Count (Qwen2 VL 72B AWQ as LVLM) | 20.67 | $\downarrow$ 61.73 | 34.48 | $\downarrow$ 151.84 |

## 5 LIMITATIONS AND FUTURE WORK

A limitation of LVLM-Count is the area detection stage. If the cropped areas do not provide enough context, performance may suffer. This is especially true for complex questions that require understanding the relationships between all the objects and the background in an image. This opens up opportunities for future work on improved area detection methods or the use of a context provider to complement the area detection stage. Another limitation arises with images containing thousands of objects. After one iteration of division, a significant number of objects may still remain in the sub-images. A potential solution is to apply additional iterations of division on the sub-images; however, the low resolution of these sub-images may make this infeasible. Developing a solution to maintain resolution is another direction for future work. Finally, in some cases, sub-images do not contain any objects of interest. The LVLM occasionally predicts a non-zero value in such instances. This is a weakness of LVLMs that requires special consideration to improve their performance.

## REFERENCES

Manoj Acharya, Kushal Kafle, and Christopher Kanan. TallyQA: Answering complex counting questions. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 8076–8084, 2019.

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

Niki Amini-Naieni, Kiana Amini-Naieni, Tengda Han, and Andrew Zisserman. Open-world text-specifed object counting. In *34th British Machine Vision Conference 2023, BMVC 2023, Aberdeen, UK, November 20-24, 2023*. BMVA, 2023. URL https://papers.bmvc2023.org/0510.pdf.

Niki Amini-Naieni, Tengda Han, and Andrew Zisserman. CountGD: Multi-modal open-world counting. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2024.

Anthropic. The Claude 3 model family: Opus, Sonnet, Haiku. `https://www-cdn.anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/Model_Card_Claude_3.pdf`, 2024. Accessed: 2024-09-30.

Carlos Arteta, Victor Lempitsky, and Andrew Zisserman. Counting in the wild. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part VII 14*, pp. 483–498. Springer, 2016.

Deepak Babu Sam, Abhinav Agarwalla, Jimmy Joseph, Vishwanath A Sindagi, R Venkatesh Babu, and Vishal M Patel. Completely self-supervised crowd counting via distribution matching. In *European Conference on Computer Vision*, pp. 186–204. Springer, 2022.

Prithvijit Chattopadhyay, Ramakrishna Vedantam, Ramprasaath R Selvaraju, Dhruv Batra, and Devi Parikh. Counting everyday objects in everyday scenes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1135–1144, 2017.

Siyang Dai, Jun Liu, and Ngai-Man Cheung. Referring expression counting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16985–16995, 2024.

Paulo RL De Almeida, Luiz S Oliveira, Alceu S Britto Jr, Eunelson J Silva Jr, and Alessandro L Koerich. PKLot–A robust dataset for parking lot classification. *Expert Systems with Applications*, 42(11):4937–4949, 2015.

Mark Everingham, SM Ali Eslami, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The PASCAL visual object classes challenge: A retrospective. *International journal of computer vision*, 111:98–136, 2015.

Shenjian Gong, Shanshan Zhang, Jian Yang, Dengxin Dai, and Bernt Schiele. Class-agnostic object counting robust to intraclass diversity. In *European Conference on Computer Vision*, pp. 388–403. Springer, 2022.

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. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 6904–6913, 2017.

Ricardo Guerrero-Gómez-Olmedo, Beatriz Torre-Jiménez, Roberto López-Sastre, Saturnino Maldonado-Bascón, and Daniel Onoro-Rubio. Extremely overlapping vehicle counting. In *Pattern Recognition and Image Analysis: 7th Iberian Conference, IbPRIA 2015, Santiago de Compostela, Spain, June 17-19, 2015, Proceedings 7*, pp. 423–431. Springer, 2015.

Philipp Kainz, Martin Urschler, Samuel Schulter, Paul Wohlhart, and Vincent Lepetit. You should use regression to detect cells. In *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18*, pp. 276–283. Springer, 2015.

Seunggu Kang, WonJun Moon, Euiyeon Kim, and Jae-Pil Heo. VLCounter: Text-aware visual representation for zero-shot object counting. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 2714–2722, 2024.

Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4015–4026, 2023.

Antti Lehmussola, Pekka Ruusuvuori, Jyrki Selinummi, Heikki Huttunen, and Olli Yli-Harja. Computational framework for simulating fluorescence microscope images with cell populations. *IEEE transactions on medical imaging*, 26(7):1010–1016, 2007.

Victor Lempitsky and Andrew Zisserman. Learning to count objects in images. *Advances in Neural Information Processing Systems*, 23, 2010.

Wei Lin, Kunlin Yang, Xinzhu Ma, Junyu Gao, Lingbo Liu, Shinan Liu, Jun Hou, Shuai Yi, and Antoni B Chan. Scale-prior deformable convolution for exemplar-guided class-agnostic counting. In *The British Machine Vision Conference (BMVC)*, pp. 313, 2022.

Chang Liu, Yujie Zhong, Andrew Zisserman, and Weidi Xie. CounTR: Transformer-based generalised visual counting. In *BMVA Press*, 2022.

Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding DINO: Marrying DINO with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023.

Erika Lu, Weidi Xie, and Andrew Zisserman. Class-agnostic counting. In *Computer Vision–ACCV 2018: 14th Asian Conference on Computer Vision, Perth, Australia, December 2–6, 2018, Revised Selected Papers, Part III 14*, pp. 669–684. Springer, 2019.

T Nathan Mundhenk, Goran Konjevod, Wesam A Sakla, and Kofi Boakye. A large contextual dataset for classification, detection and counting of cars with deep learning. In *Computer Vision– ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part III 14*, pp. 785–800. Springer, 2016.

Thanh Nguyen, Chau Pham, Khoi Nguyen, and Minh Hoai. Few-shot object counting and detection. In *European Conference on Computer Vision*, pp. 348–365. Springer, 2022.

Joseph Paul Cohen, Genevieve Boucher, Craig A Glastonbury, Henry Z Lo, and Yoshua Bengio. Count-ception: Counting by fully convolutional redundant counting. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pp. 18–26, 2017.

Jer Pelhan, Vitjan Zavrtanik, Matej Kristan, et al. DAVE-A: Detect-and-verify paradigm for low-shot counting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 23293–23302, 2024.

Penguin Research. Penguin research webpage, 2016. URL https://www.robots.ox.ac.uk/~vgg/data/penguins/. Accessed: 2024-11-23.

Viresh Ranjan, Udbhav Sharma, Thu Nguyen, and Minh Hoai. Learning to count everything. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3394–3403, 2021.

Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024.

Stuart J Russell and Peter Norvig. *Artificial intelligence: a modern approach*. Pearson, 2016.

Min Shi, Hao Lu, Chen Feng, Chengxin Liu, and Zhiguo Cao. Represent, compare, and learn: A similarity-aware framework for class-agnostic counting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9529–9538, 2022.

Zenglin Shi, Ying Sun, and Mengmi Zhang. Training-free object counting with prompts. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 323–331, 2024.

Nikola Đukić, Alan Lukežič, Vitjan Zavrtanik, and Matej Kristan. A low-shot object counting network with iterative prototype adaptation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 18872–18881, 2023.

Meng Wang and Xiaogang Wang. Automatic adaptation of a generic pedestrian detector to a specific traffic scene. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR) 2011*, pp. 3401–3408. IEEE, 2011.

Xintao Wang, Liangbin Xie, Chao Dong, and Ying Shan. Real-ESRGAN: Training real-world blind super-resolution with pure synthetic data. In *International Conference on Computer Vision Workshops (ICCVW)*, 2021.

Weidi Xie, J Alison Noble, and Andrew Zisserman. Microscopy cell counting and detection with fully convolutional regression networks. *Computer methods in biomechanics and biomedical engineering: Imaging & Visualization*, 6(3):283–292, 2018.

Haipeng Xiong, Hao Lu, Chengxin Liu, Liang Liu, Zhiguo Cao, and Chunhua Shen. From open set to closed set: Counting objects by spatial divide-and-conquer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 8362–8371, 2019.

An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2024.

Shuo-Diao Yang, Hung-Ting Su, Winston H Hsu, and Wen-Chin Chen. Class-agnostic few-shot object counting. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 870–878, 2021.

Zhiyuan You, Kai Yang, Wenhan Luo, Xin Lu, Lei Cui, and Xinyi Le. Few-shot object counting with similarity-aware feature enhancement. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 6315–6324, 2023.

Gloria Zen, Negar Rostamzadeh, Jacopo Staiano, Elisa Ricci, and Nicu Sebe. Enhanced semantic descriptors for functional scene categorization. In *Proceedings of the 21st International Conference on Pattern Recognition*, pp. 1985–1988. IEEE, 2012.

# A  ABLATION STUDY

We examine the effect of each of the following stages in our method: i) area detection and ii) object-aware division (note that the object-aware division necessitates the inclusion of target segmentation stage). The experiments are designed to investigate the effect of each stage individually, as well as when the stages are combined in our pipeline. Additionally, we run an experiment for a case where both stages of area-detection and target segmentation are excluded. In this case, the object-aware division can not be performed. Thus, images are divided by equidistant straight lines into subimages. We give the name of naive division to such an approach. Moreover, we run another experiment where area detection is in place but the target segmentation is excluded and naive division is applied on the detected areas We run the ablation scenarios for two LVLMs: GPT-4o, and Gemini 1.5 Pro. Furthermore, the effect of super resolution at the final stage of LVLM-Count is also investigated for the case that GPT-4o is the LVLM. For the experiments, we randomly sample 4 images from each category in the FSC-147 test set (29 categories) and report the performance on the resulting subset containing 116 samples. In Table 4 we show the results of the ablation experiments.

Table 4: Ablation study for LVLM-Count on 116 samples from the FSC-147 test dataset. Columns marked with $\Delta$ show the performance difference between an entry and the base LVLM it uses. Green indicates improvement, while red represents degradation.

| Method | MAE ↓ | $\Delta$ | RMSE ↓ | $\Delta$ |
|---|---|---|---|---|
| GPT-4o | 18.58 | - | 56.84 | - |
| GPT-4o + Naive division | 33.95 | ↑ 15.37 | 88.75 | ↑ 31.91 |
| GPT-4o + Area Detection + Naive division | 34.82 | ↑ 16.24 | 73.19 | ↑ 16.35 |
| GPT-4o + Area detection | 19.27 | ↑ 0.69 | 70.64 | ↑ 13.8 |
| GPT-4o + Object-aware division | 14.92 | ↓ 3.66 | 31.56 | ↓ 25.28 |
| GPT-4o + Area detection + Object-aware divsion (equiv. to LVLM-Count) | **11.03** | ↓ 7.55 | **21.02** | ↓ 35.82 |
| GPT-4o + Area detection + Object-aware divsion + Super resolution | 12.30 | ↓ 6.28 | 28.98 | ↓ 27.86 |
| Gemini 1.5 Pro | 22.07 | - | 64.34 | - |
| Gemini 1.5 Pro + Naive division | 41.31 | ↑ 19.24 | 76.56 | ↑ 12.22 |
| Gemini 1.5 Pro + Area Detection + Naive division | 43.71 | ↑ 21.64 | 75.47 | ↑ 11.13 |
| Gemini 1.5 Pro + Area detection | 23.04 | ↑ 0.97 | 61.94 | ↓ 2.4 |
| Gemini 1.5 Pro + Object-aware division | 15.39 | ↓ 6.68 | 38.15 | ↓ 26.19 |
| Gemini 1.5 Pro + area detection + Object-aware division (equiv. to LVLM-Count) | **13.64** | ↓ 8.43 | **36.56** | ↓ 27.78 |

## B  BIAS TO SMALL NUMBERS IN DATASETS AND PERFORMANCE OF LVLMs

In our experiments, LVLMs (GPT-4o) are able to make correct predictions when the number of items to be counted is small, but errors increase as the ground truth number grows. Although, we cannot be certain, one likely reason for this behavior in LVLMs is that, during training, the counting questions these models encounter are heavily biased toward small numbers. As an example, in Figure 8 we show the distribution of 'How many' questions in some well-known VQA datasets.



(a) TallyQA Dataset (Acharya et al., 2019)    (b) VQAv2 Dataset (Goyal et al., 2017)

Figure 8: Distribution of answers in two VQA datasets.

## C  DEFINITION AND EXAMPLE OF SIMPLE AND COMPLEX COUNTING TASKS

Acharya et al. (2019) were among the first to formally categorize counting questions into simple and complex types. They applied a linguistic rule: first, they removed any substrings such as "...in the photo?" or "...in the image?". Then, they used SpaCy to perform part-of-speech tagging on the remaining substring. They classified a question as simple if it contained only one noun, no adverbs, and no adjectives; otherwise, they deemed it complex. This rule classifies questions such as "How many dogs?" as simple and "How many brown dogs?" as complex. Following this rule, they built two splits for the TallyQA dataset: a simple split and a complex split. Since we sample our open-ended counting benchmark from the TallyQA simple and complex splits, we adopt the same classification criteria for this benchmark.

In Figure 9, we provide examples of a simple and a complex VQA question from TallyQA. The VQA question in Figure 9a is selected from the simple split of TallyQA. It is a straightforward counting question, merely asking for all instances of the animal. Figure 9b, on the other hand, is selected from the complex split of TallyQA. In addition to counting, it requires detecting context and distinguishing between consonants and vowels.

## D  INCORRECT EXAMPLES IN THE FSC-147 DATASET

We observed some incorrect instances in the FSC-147 dataset. For example, see Figure 10. In these cases, the category names which are provided are incorrect. For these instances, the extensive knowledge embedded within the LVLMs employed in our approach proves to be a disadvantage. These models detect inconsistencies and provide a count of zero as the output, whereas other methods are misled by superficial similarities and mistakenly count the objects. A more thorough study is required to detect all the incorrect examples in FSC-147.

(a) How many cows are visible?

(b) How many consonants are there on the green sign?

Figure 9: A simple (left image) and a complex (right image) question from TallyQA.



(a) Category name in FSC-147: apples, Answer in FSC-147: 182, Ground truth: 0.

(b) Category name in FSC-147: sunglasses, Answer in FSC-147: 81, Ground truth: 0.

Figure 10: Two erroneous examples from FSC-147. Figure 10a shows a number of crabapples while the category name in the FSC-147 is apples. Figure 10b features a number of glasses while the category name in FSC-147 is sunglasses.

## E    ALLEVIATING THE PROBLEM OF INACCURATE DETECTIONS AT THE AREA DETECTION AND TARGET SEGMENTATION STAGES OF LVLM-COUNT

One of the limitations of LVLM-Count is that if the area detection fails to detect a relevant area, the objects of interest in that area will not be counted. Another limitation arises when the target segmentation fails to segment all the instances of the target category of objects. In such a scenario, non-segmented objects of interest might be cut through by division lines, causing the predicted total number to increase due to being counted multiple times. One simple, yet quite effective, approach that we use to overcome this limitation to a significant degree is setting the detection threshold of GroundingDINO, which is used in both of the mentioned stages, to an extremely low value. This significantly reduces the chance of missing an area or object of interest, although it might lead to some false positives. However, note that false positives do not harm the performance of LVLM-Count, as their only effect is that the object-aware division lines avoid cutting through them.

To demonstrate the effectiveness of this approach, we evaluate it on the Penguin dataset (Penguin Research, 2016). The goal in this dataset is to count penguins in images. The challenging Penguin dataset consistently exhibits heavy occlusion and complex background patterns that can easily be

15

mistaken for penguins Arteta et al. (2016). This dataset consists of two splits: the mixed-site split, in which images from the same camera can appear in both the training and testing sets, and the separate-site split, in which images in each set strictly belong to different cameras. Images in this dataset are annotated by multiple annotators, where each annotator might identify a different number of penguins due to the challenges in locating them within the images. Since annotators usually undercount the penguins, similar to Arteta et al. (2016), we take the maximum number of penguins among the annotations as the ground truth and calculate MAE and RMSE with respect to this value. For more details about the dataset and the calculation of metrics, we refer the reader to Arteta et al. (2016).

Since both splits are very large, we randomly select 100 samples from each split. To preserve balance, the probability of selecting an image with a certain ground-truth annotation number is inversely proportional to the frequency of that annotation number in the entire split, excluding images with zero penguins. We set the area detection and target segmentation thresholds to 0.5, and then we run them with both thresholds set to 0.1. The results are presented in Table 5 and Table 6. We observe that for both splits, when the detection threshold is high, the MAE generally increases to a value higher than the MAE of the base LVLMs. However, setting the thresholds to a low value improves the MAE for all three LVLMs.
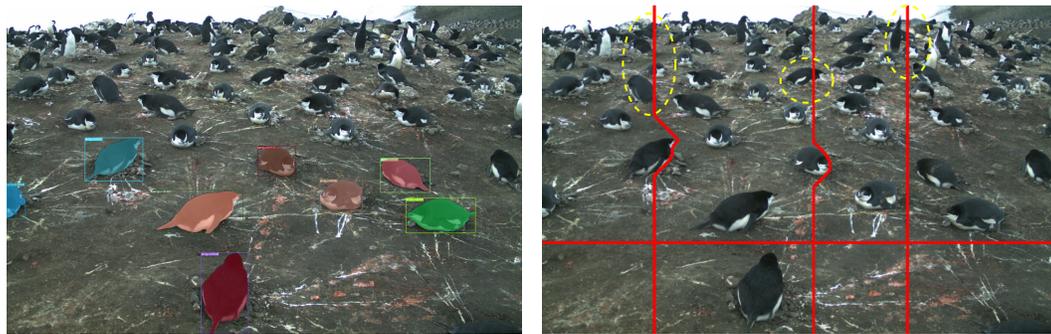
Table 5: Penguin dataset - Mixed sites. Columns marked with $\Delta$ show the performance difference between LVLM-Count and the base LVLM it uses. Green indicates improvement, while red represents degradation.

| Method | MAE (Max)↓ | $\Delta$ | RMSE (Max)↓ | $\Delta$ |
|---|---|---|---|---|
| GPT-4o | 26.45 | - | 37.53 | - |
| LVLM-Count(GPT-4o - Thredshold=0.5) | 36.06 | ↑ 9.61 | 49.23 | ↑ 11.7 |
| LVLM-Count(GPT-4o - Thredshold=0.1) | 23.54 | ↓ 2.91 | 42.24 | ↑ 4.71 |
| Gemini 1.5 Pro | 37.23 | - | 49.15 | - |
| LVLM-Count(Gemini 1.5 Pro - Threshold=0.5) | 42.42 | ↑ 5.19 | 50.39 | ↑ 1.24 |
| LVLM-Count(Gemini 1.5 Pro - Threshold=0.1) | 20.09 | ↓ 17.14 | 29.95 | ↓ 19.2 |
| Qwen2 VL 72B AWQ | 34.25 | - | 55.70 | - |
| LVLM-Count(Qwen2 VL 72B AWQ - Threshold=0.5) | 37.13 | ↑ 2.87 | 53.71 | ↓ 1.99 |
| LVLM-Count(Qwen2 VL 72B AWQ - Threshold=0.1) | 18.22 | ↓ 16.03 | 28.76 | ↓ 26.94 |

Table 6: Penguin dataset - Separate sites. Columns marked with $\Delta$ show the performance difference between LVLM-Count and the base LVLM it uses. Green indicates improvement, while red represents degradation.

| Method | MAE (Max)↓ | $\Delta$ | RMSE (Max)↓ | $\Delta$ |
|---|---|---|---|---|
| GPT4o | 36.01 | - | 46.64 | - |
| LVLM-Count(GPT-4o - Thredshold=0.5) | 43.25 | ↑ 7.24 | 52.17 | ↑ 5.53 |
| LVLM-Count(GPT-4o - Thredshold=0.1) | 27.11 | ↓ 8.9 | 36.72 | ↓ 9.92 |
| Gemini 1.5 Pro | 46.67 | - | 57.26 | - |
| LVLM-Count(Gemini 1.5 Pro - Threshold=0.5) | 33.59 | ↑ 13.08 | 46.23 | ↑ 11.03 |
| LVLM-Count(Gemini 1.5 Pro - Threshold=0.1) | 28.50 | ↓ 18.17 | 41.39 | ↓ 15.87 |
| Qwen2 VL 72B AWQ | 44.96 | - | 66.78 | - |
| LVLM-Count(Qwen2 VL 72B AWQ - Threshold=0.5) | 51.17 | ↑ 6.21 | 75.49 | ↑ 8.71 |
| LVLM-Count(Qwen2 VL 72B AWQ - Threshold=0.1) | 24.83 | ↓ 20.13 | 38.92 | ↓ 27.86 |

In order to also provide a visual insight into this limitation, we chose an image from the dataset where distinguishing penguins from the background features is challenging for GroundingDINO. Figure 11 shows the performance of our method on this image, with the area detection and target segmentation thresholds both set to 0.5. We can observe that the failure of the target segmentation stage in segmenting all instances of penguins in Figure 11a has resulted in division paths cutting through several penguins inside the areas surrounded by yellow shapes in Figure 11b. Consequently, these penguins are counted twice, causing a larger error.

On the other hand, Figure 12 shows the performance of LVLM-Count when the area detection and target segmentation thresholds are set to 0.1. Figure 12a illustrates that, while a low threshold has helped to segment all instances of the penguins, it has also caused some false positive areas to be
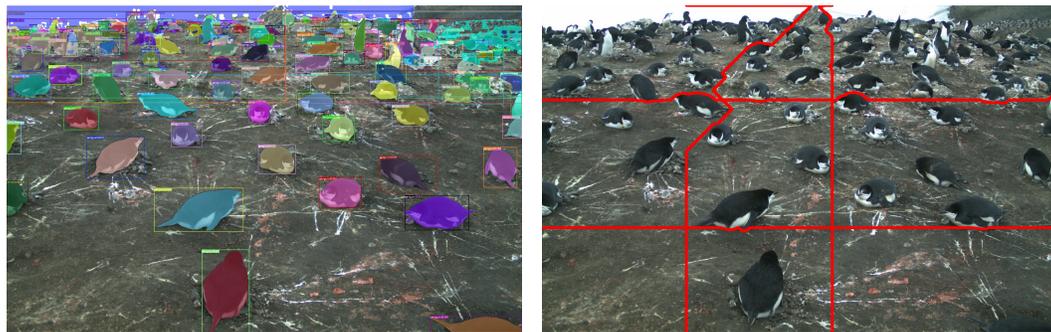
segmented. Nonetheless, it can be observed in Figure 12b that the only effect of false positive segmentation masks is that the division paths avoid cutting through those regions as well. Note that for both images, the cluster-based approach was used to automatically find the start and end points of the division paths based on the arrangement of the segmentation masks in the scene. Comparing Figure 11 and Figure 12, we can clearly see that reducing the detection threshold is a quite successful technique to overcome limitations in the initial stages such as target segmentation.



(a) With a threshold of 0.5, target segmentation stage has failed to segment all penguins

(b) Areas surrounded by yellow shapes contain several penguins bisected by the dividing lines.

Figure 11: In this figure, the detection threshold in the target segmentation stage is set to 0.5. For this example, the cluster-based approach has been used to find the start and end points of the division paths. Ground truth (Max) = 106, LVLM-Count prediction = 131.



(a) With a low threshold such as 0.1, some false positive areas have also been segmented.

(b) The division paths avoid cutting through the false positive areas as well.

Figure 12: In this figure, the detection threshold in the target segmentation stage is set to 0.1. For this example, a cluster-based approach was used to find the start and end points of the division paths. Ground truth (Max) = 106, LVLM-Count prediction = 110

Additionally, Figure 13 shows how using a low detection threshold can help alleviate the limitations in the area detection stage for a challenging image taken in unfavorable weather conditions. Figure 13a shows the output of the area detection stage when the threshold is set to 0.5. The output contains two detected areas that, in combination, fail to cover all the penguins, resulting in error. However, Figure 13b shows the output of the area detection stage when the threshold is set to 0.1, successfully giving the area containing all the penguins.

# F LVLM-COUNT'S POWER IN HANDLING MULTIPLE OBJECT CATEGORIES IN THE SAME IMAGE

LVLM-Count is a highly effective method for handling counting tasks that involve multiple objects in the same image. Its strength in such scenarios stems from the capabilities of LVLMs to answer numerous visual questions about an image and its objects. Depending on the given text prompt, it can

(a) The area detection threshold set to 0.5          (b) The area detection threshold set to 0.1

Figure 13: This figure shows that setting the detection threshold in the area detection stage to a very low value (Figure 13b) significantly alleviates the area detection limitations at the area detection stage.

count instances of a single object category among others or instances of multiple object categories simultaneously. In this section, we demonstrate how LVLM-Count performs in counting different objects of interest, determined simply by a prompt, using an image with multiple object categories.

The image in Figure 14 contains three object categories: person, cow, and horse. In the top row, the object of interest is "cow." We prompt LVLM-Count to count the cows. First, the masks are produced through the initial stages of our pipeline, and then the cluster-based approach is used to automatically determine the start and end points of the division paths. It can be observed that horses have also been masked as cows. Nonetheless, this does not negatively impact the final answer; it merely causes the division lines to avoid cutting through the horses as well. The counting in LVLM-Count is performed by an LVLM (GPT-4o in this figure) and does not rely on the masks. We observe that GPT-4o successfully counts the number of cows in the resulting subimages, leading to the correct final answer.

In the middle row of Figure 14, the object of interest is "person." LVLM-Count again successfully counts the number of people accurately. A more interesting case is the bottom row of Figure 14, where both cows and persons are objects of interest. We prompt LVLM-Count to count the number of "cows and persons." Similar to the first row of the figure, there are false positive masks here as well. However, LVLM-Count successfully counts the number of instances from both categories combined since the counting is ultimately performed by the LVLM. Note that the number of objects in this image is limited, and GPT-4o might answer these questions correctly without the need for the LVLM-Count pipeline. This image has been chosen to illustrate LVLM-Count's power in handling multiple objects in a counting task rather than for comparison with the baseline LVLM.

## G    REAL-WORLD APPLICATION OF LVLM-COUNT

As stated in Section 1, counting has numerous real-world applications, including but not limited to biology, health, industry, warehousing, and environmental monitoring. Below, we demonstrate the performance of LVLM-Count on examples from the following areas: i) biology/health, ii) industry/warehousing, and iii) environmental monitoring. We also compare its results with those of the base LVLM (GPT-4o for the figures in this section). Note that in all examples, the cluster-based approach automatically determines the start and end points of the division paths.

In Figure 15, images of two laboratory samples are analyzed using LVLM-Count. The first row shows an image from a dataset introduced by Lempitsky & Zisserman (2010), which contains simulated bacterial cells from fluorescence-light microscopy, created by Lehmussola et al. (2007). The second row shows an image from the BM dataset introduced by Kainz et al. (2015), which contains bone marrow samples from eight healthy individuals. The standard staining procedure highlights the nuclei of various cell types in blue, while other cellular components appear in shades of pink and red (Paul Cohen et al., 2017). As observed, LVLM-Count achieves much higher accuracy in counting bacterial cells and bone marrow nuclei in the top and bottom rows of Figure 15, respectively, compared to the base LVLM, particularly for the bone marrow nuclei.

Figure 14: Illustration of the ability of LVLM-Count in counting an object of interest determined my a prompt when multiple object categories exist in a single image. Top row: Object of interest is "person". Middle row: Object of interest is "cow". Bottom row: Object of interest is "person and cow"

In Figure 16, two images from industrial scenes are analyzed using LVLM-Count. The top row shows a sectional image of a stockpile of tree logs, and the bottom row shows an image from an industrial area containing barrels of various colors. For the top image, the objects of interest are the tree logs, while for the bottom image, LVLM-Count is tasked with counting the *blue* barrels. In both cases, LVLM-Count's predictions are significantly closer to the ground truth values than those of the base LVLM, particularly for the tree logs, where the ground truth number is too large for the base LVLM to estimate accurately.

Figure 17 shows an image sourced from a dataset (Penguin Research, 2016) created as part of an ongoing initiative to monitor the penguin population in Antarctica. This dataset comprises images captured hourly by a network of fixed cameras installed at over 40 locations. Over several years, this effort has accumulated over 500,000 images. Zoologists use these images to identify trends in penguin population sizes at each site, facilitating studies on potential correlations with factors such as climate change. Thus, determining the number of penguins in each image is crucial. Given the challenges of engaging human annotators to process such a vast dataset, automating the counting task is highly desirable (Arteta et al., 2016). LVLM-Count is prompted to count the number of

penguins in the image, and as observed, its predictions are significantly closer to the ground truth than those of the base LVLM.
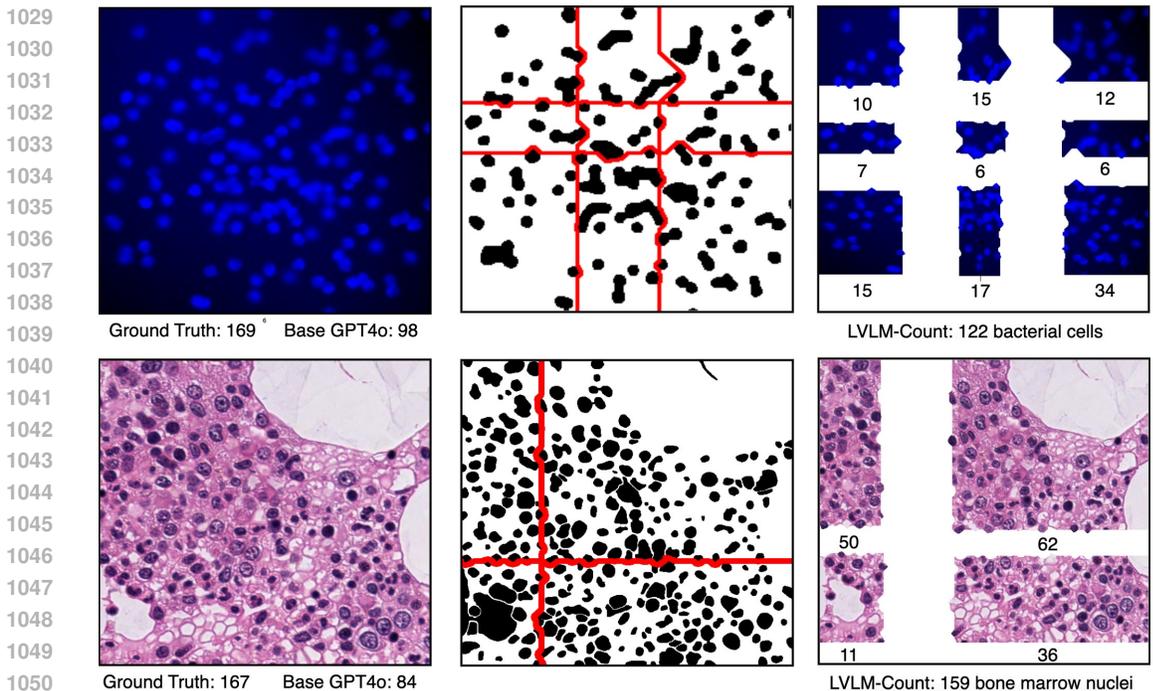


Figure 15: Performance of LVLM-Count on real-world applications in biology/health. The top row shows an image of simulated bacterial cells from fluorescence-light microscopy (Lempitsky & Zisserman, 2010), with the objects of interest being "bacterial cells." The bottom row shows an image of bone marrow, with the nuclei of various cell types highlighted in blue (Kainz et al., 2015), and the objects of interest being "bone marrow nuclei."

## H    PERFORMANCE ANALYSIS OF LVLM-COUNT FOR DIFFERENT GROUND TRUTH RANGES ON FSC-147 DATASET

To further investigate the performance of our pipeline, we divide the ground truth values in the FSC-147 test set into intervals and plot the MAE for the base GPT-4o and Gemini 1.5 Pro models, alongside the results from LVLM-Count using each model, as shown in Figure 18. The first interval contains relatively small ground truth values, a range where LVLMs already perform well. As the ground truth values increase, the base models exhibit increasingly larger errors compared to LVLM-Count, with the margin growing rapidly. This behavior is consistent with our observations of counting errors on the blue circles in Figure 2.

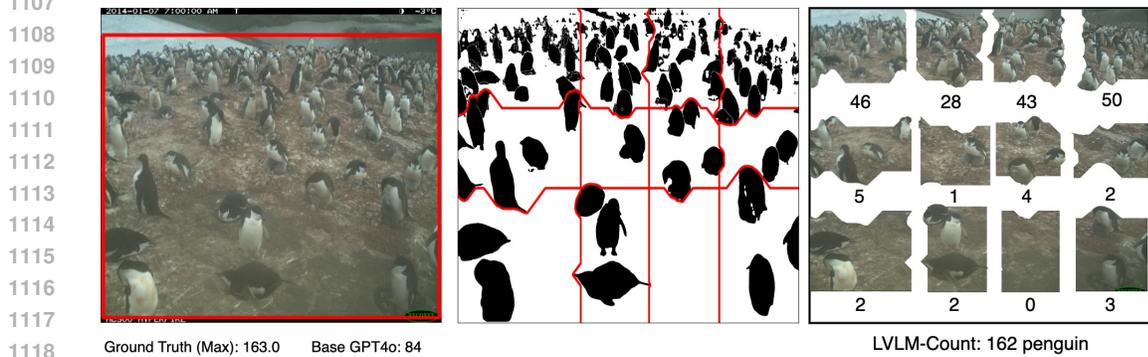## I    ILLUSTRATION OF THE WORKFLOW FOR THE ZEBRA IMAGE IN FIGURE 1

In this paper, we choose the most appropriate images for simple and understandable illustrations of the execution of each stage in our pipeline. Inevitably, this led to the selection of different images for each stage. However, in this section, for the sake of consistency, we demonstrate the same concepts illustrated in Figures 3, 4, 5, and 6 for the zebra image used in Figure 1.

The zebra image is passed to the pipeline along with the question $Q =$ "how many zebras are in the image?". First, $E =$ "zebra" is extracted using the LLM. Then the zebra image is passed to the area detection stage, where the prompt given to GroundingDINO is "zebras". The output bounding boxes are merged, and the resulting area is cropped, as illustrated in Figure 19. The cropped area is then passed to the target segmentation stage. At this stage, GroundingDINO detects the objects

Ground Truth: 205    Base GPT4o: 101          LVLM-Count: 207 tree log

Ground Truth: 27    Base GPT4o: 41          LVLM-Count: 29 blue barrel

Figure 16: Performance of LVLM-Count on real-world applications in industry/warehousing. The top row shows an image of a stockpile of tree logs, with the objects of interest being "tree logs." The bottom row shows an aerial image of an industrial area containing barrels of various colors, with the objects of interest being "*blue* barrels."
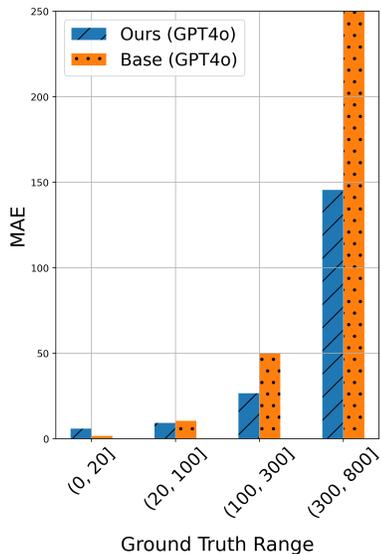


Ground Truth (Max): 163.0    Base GPT4o: 84          LVLM-Count: 162 penguin

Figure 17: Performance of LVLM-Count on real-world applications in environmental monitoring. The image is sourced from (Penguin Research, 2016), an initiative to monitor the penguin population in Antarctica, with the objects of interest being "penguins."

of interest defined by $E$ as the input prompt. SAM then uses the output bounding boxes to produce segmentation masks for the zebras, as shown in Figure 20.

After the target segmentation stage, the masks are passed to the object-aware division stage. First, the masks are used in the cluster-based approach to find the location of the start and end points of the division paths, i.e., $(P_s^1, P_e^1)$ and $(P_s^2, P_e^2)$. Then these masks are turned into black-and-white images, which, in turn, are mapped to a graph. The division paths are then found by connecting each start point to its corresponding end point by running the $A^*$ search algorithm on the graph. The found paths are mapped back into the image domain and drawn in red, as depicted in Figure 21. The image contours are determined based on the drawn red paths, and each contour's interior is masked out independently to obtain the subimages. Finally, the subimages are given to the LVLM to count the number of zebras in each.

(a) LVLM-Count using GPT-4o

(b) LVLM-Count using Gemini 1.5 Pro

Figure 18: Performance analysis of our method, LVLM-Count, on the FSC-147 test set using GPT-4o (Figure 18a) and Gemini 1.5 Pro (Figure 18b). In the first interval, both base LVLMs exhibit a lower MAE. However, in intervals with higher ground truth values, LVLM-Count achieves a lower MAE compared to the base LVLMs, and this difference increases rapidly.
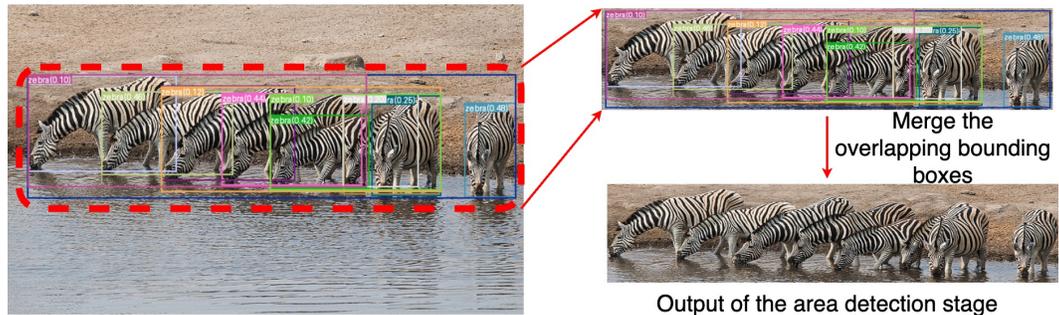


Figure 19: Illustration of the area detection step of LVLM-Count for the zebra image. For this image, $Q$ is set to "How many zebras are in the image?". The LLM used in this step returns an $E$, which is "zebra". The plural form of $E$, "zebras", and the original image are given as input to GroundingDINO, which returns some bounding boxes (left and upper right images) that are merged to form the final detected area.

## J    FALSE POSITIVE MASKS AT THE TARGET SEGMENTATION STAGE

One of the reasons we task an LVLM to count the objects in the subimages instead of using the number of generated masks at the target segmentation stage as the final count of the objects of interest is the existence of false positive masks. The GroundingDINO model is responsible for detecting the objects of interest, determined by expression $E$, and passing the output bounding boxes to SAM for producing segmentation masks. Nonetheless, GroundingDINO is not as strong as an LVLM in understanding expressions extracted from complex questions. Thus, it often returns bounding boxes for all instances of the object category mentioned in the expression, even if those instances do not satisfy other conditions in the expression.

For example, in the top row of Figure 22, $E$ = "brown egg". However, all the eggs have been segmented regardless of their color. Thus, counting the masks results in a significant error. Interestingly, as we can see, the false positive masks do not negatively affect LVLM-Count's final answer, as

(a) GroundingDINO output  (b) SAM output

Figure 20: Illustration of the target segmentation step of LVLM-Count for the zebra image. The goal is to produce all the instance masks for $E$ set to "zebra". The cropped image from Figure 19, together with $E$, is given as input to GroundingDINO, which produces the output shown in Figure 20a. Figure 20a is then given as input to SAM, which produces the output shown in Figure 20b.
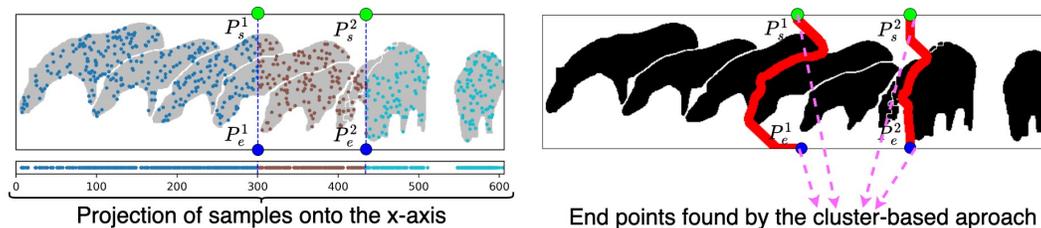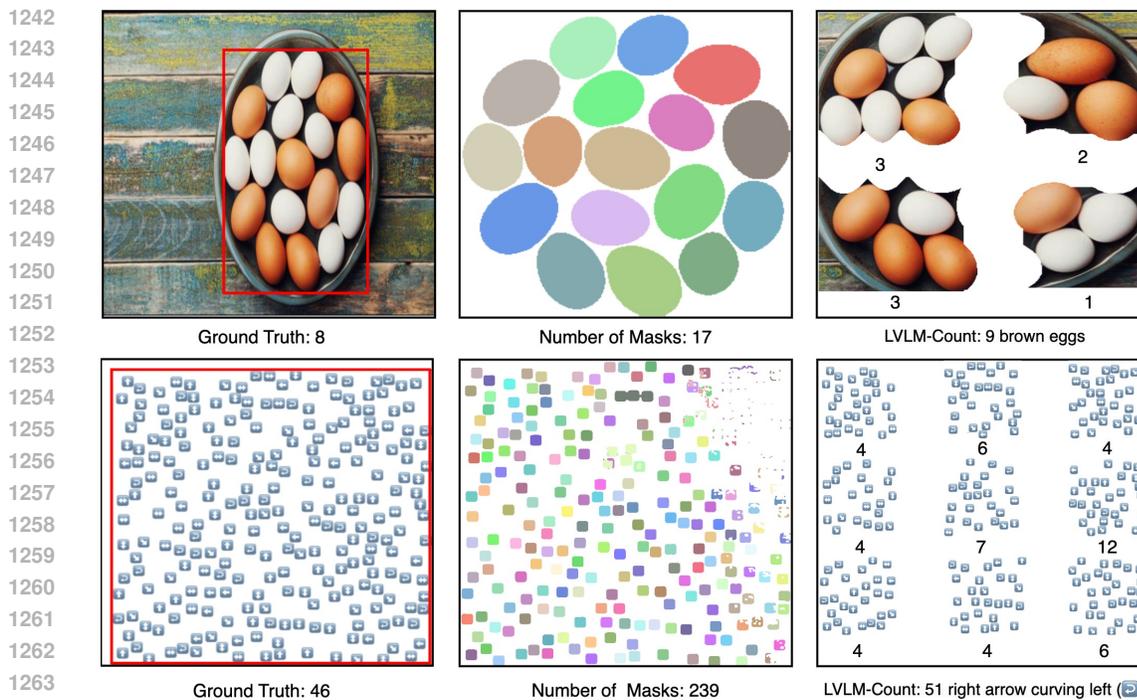


Projection of samples onto the x-axis  End points found by the cluster-based aproach

Figure 21: Left: Illustration of the unsupervised and non-parametric method to obtain the division points $(P_s^1, P_e^1)$ and $(P_s^2, P_e^2)$. First, a few pixels are sampled (shown as points inside the segmented objects) from the pixels composing each mask. The samples are projected onto the $x$-axis. The projected points are clustered using mean-shift clustering. The point in the middle of two consecutive clusters is considered a vertical division point. The straight vertical lines are drawn just for better visualization of the division points. Right: Illustration of object-aware division. The masks from Figure 20b are turned into a black-and-white image. A dividing path is found by connecting $P_s$ to $P_e$ using the $A^*$ search algorithm in a graph that corresponds to the binary image, where the only nodes in the graph are white pixels, which are fully connected. The path is mapped back to the pixel domain.

the counting is done by an LVLM at the final stage, which is much stronger than GroundingDINO at understanding referring expressions. The only effect is that the white eggs have not been cut through by the division lines either. In the bottom row, we have chosen an image from the challenging Emoji-Count benchmark. The image contains icons, all of which have an arrow but point in different directions. However, the objects of interest are only "right arrows curving left." Similar to the eggs example, taking the masks used for object-aware division results in a significant error.

## K    LVLM-COUNT'S PERFORMANCE ON A BENCHMARK SAMPLED FROM THE PASCAL VOC DATASET

The PASCAL VOC dataset (Everingham et al., 2015) is a well-known dataset in the field of computer vision, depicting everyday objects in everyday scenes. It is primarily used for tasks like object detection, classification, and segmentation. To provide more experimental results for the performance of LVLM-Count, we adapt a small subset of this dataset into the form of a counting benchmark. Similar to Chattopadhyay et al. (2017), we choose PASCAL VOC 2007 among other variants. This variant contains a training set of 2501 images, a validation set of 2510 images, and a test set of 4952 images, with 20 object categories that remain consistent across the splits.

Each image includes annotations for instances of the 20 object categories in the dataset. The most frequent count per object category (as one would expect in everyday scenes) is 0. Moreover, there is a clear bias toward lower count values. To address this bias and obtain a balanced benchmark, we first create 20 simple counting questions asking for the number of objects from each of the 20 categories for every image in the test set. Then, we randomly sample five questions for each ground truth count. Note that the ground truth counts for this sampling are based on the original annotations of the dataset. The bias toward lower numbers is so significant that for ground truth counts larger than 16, there are fewer than five questions for the entire test set. In such cases, we simply use the

Figure 22: Top row: The object of interest is "*brown* egg." However, all the eggs have been segmented because of the limitation of the GroundingDINO model in understanding complex referring expressions. Regardless, LVLM-Count provides a significantly more accurate number compared to the number of masks. Bottom row: The object of interest is "right arrows curving left." Similar to the image of the eggs, counting the number of masks results in a very large error, while LVLM-Count provides a much more accurate number.

available samples. This process resulted in 102 questions in total. Finally, we manually checked the ground truth counts and corrected them if required.

We call the resulting counting questions the PASCAL VOC counting benchmark. We evaluate the performance of LVLM-Count on this benchmark using three different LVLMs: GPT-4o, Gemini 1.5 Pro, and Qwen2 VL 72B AWQ. Note that we use a clustering approach to automatically find the start and end points of the division paths along the $x$-axis. Table 7 shows the performance of LVLM-Count in comparison with state-of-the-art counting models and the base LVLMs. We observe that LVLM-Count improves upon the base LVLMs it uses and outperforms all prior counting methods, except CountGD. Upon further inspection, we noticed that the 20 object categories in PASCAL VOC have a high overlap with the object categories of the dataset used to train CountGD.

## L    REPORT OF VARIOUS ACCURACY METRICS FOR THE PERFORMANCE OF LVLM-COUNT ON THE FSC-147 DATASET, TALLYQA BENCHMARK, AND EMOJI-COUNT BENCHMARK

This section presents various accuracy measures for the experiments reported in Tables Table 1, 2, and 3. The accuracy metrics are defined in Table 8. The observations for each table can be summarized as follows:

i) **FSC-147:** For this dataset, similar to the results in Table 1, LVLM-Count achieves higher accuracy metrics compared to the prior training-free method and each of the base LVLMs it uses. However, models specifically trained on this dataset generally achieve higher accuracy.

ii) **TallyQA Simple Benchmark:** For this benchmark, LVLM-Count achieves higher accuracy compared to all prior counting models. However, it does not surpass the accuracy of the base LVLMs

Table 7: Evaluation of state-of-the-art models on the PASCAL VOC counting benchmark. The results for base LVLMs and LVLM-Count are reported over three trials. Columns marked with $\Delta$ show the improvement brought by our method over the base LVLM that it uses.

| Method | EA (%)↑ | $\Delta$ | MAE↓ | $\Delta$ | RMSE↓ | $\Delta$ |
|---|---|---|---|---|---|---|
| TrainingFree (Shi et al., 2024) | 2.94 | - | 12.03 | - | 18.18 | - |
| GroundingRec (Dai et al., 2024) | 19.61 | - | 5.05 | - | 8.44 | - |
| DAVE$_{prm}$ (Pelhan et al., 2024) | 5.88 | - | 12.39 | - | 22.81 | - |
| CountGD (Amini-Naieni et al., 2024) | 26.47 | - | 2.81 | - | 7.01 | - |
| GPT4o | 30.39 | - | 4.46 | - | 8.35 | - |
| LVLM-Count (GPT4o as LVLM) | 31.37 | ↑ 0.98 | 3.55 | ↓ 0.91 | 7.18 | ↓ 1.17 |
| Gemini 1.5 Pro | 33.01 | - | 3.24 | - | 6.62 | - |
| LVLM-Count (Gemini 1.5 Pro as LVLM) | 39.22 | ↑ 6.21 | 3.00 | ↓ 0.24 | 6.30 | ↓ 0.32 |
| Qwen2 VL 72B AWQ | 24.18 | - | 4.83 | - | 8.84 | - |
| LVLM-Count (Qwen2 VL 72B AWQ as LVLM) | 29.08 | ↑ 4.90 | 4.12 | ↓ 0.71 | 7.59 | ↓ 1.25 |

used. This is because the questions are straightforward, and the ground truth values are limited to numbers between 0 and 15—a range where the base LVLMs excel. This observation aligns with those in Figure 2 and Figure 18.

iii) **TallyQA Complex Benchmark:** For this benchmark, LVLM-Count demonstrates significantly higher accuracies compared to prior counting models and, more importantly, consistent accuracy improvements over the base LVLMs used.

iv) **Emoji-Count:** This is a challenging benchmark due to high object counts. LVLM-Count achieves substantially higher accuracies than both the base LVLMs and prior counting models, particularly for metrics like Acc$\pm k$ where $k \geq 1$.

Table 8: Definitions of Various Accuracy Metrics. GT denotes the ground truth number.

| Metric | Definition |
|---|---|
| Acc | Percentage of answers such that $answer = GT$ |
| Acc$\pm 1$ | Percentage of answers such that $|answer - GT| \leq 1$ |
| Acc$\pm 3$ | Percentage of answers such that $|answer - GT| \leq 3$ |
| Acc$\pm 5$ | Percentage of answers such that $|answer - GT| \leq 5$ |
| Acc$\pm 10$ | Percentage of answers such that $|answer - GT| \leq 10$ |

Table 9: FSC-147 Dataset. A (↑) next to the measured accuracies for LVLM-Count indicates improvement over the base LVLM it uses, while a (↓) indicates degradation compared to the corresponding base LVLM.

| Method | Acc (%) | Acc$\pm 1$ (%) | Acc$\pm 3$ (%) | Acc$\pm 5$ (%) | Acc$\pm 10$ (%) |
|---|---|---|---|---|---|
| TFOC (Shi et al., 2024) | 9.33 | 20.17 | 34.20 | 44.37 | 61.60 |
| GroundingRec (Dai et al., 2024) | 34.03 | 51.68 | 67.65 | 75.04 | 85.13 |
| CountGD (Amini-Naieni et al., 2024) | 31.85 | 47.90 | 63.78 | 73.03 | 82.61 |
| GPT4o | 12.24 | 26.22 | 42.10 | 52.41 | 66.58 |
| LVLM-Count (GPT4o as LVLM) | 14.26 (↑) | 28.10 (↑) | 47.42 (↑) | 58.01 (↑) | 72.91 (↑) |
| Gemini 1.5 Pro | 12.97 | 26.58 | 41.57 | 51.71 | 63.78 |
| LVLM-Count (Gemini 1.5 Pro as LVLM) | 13.92 (↑) | 27.96 (↑) | 47.65 (↑) | 58.85 (↑) | 75.07 (↑) |
| Qwen2 VL 72B AWQ | 9.19 | 20.81 | 36.83 | 46.95 | 62.07 |
| LVLM-Count (Qwen2 VL 72B AWQ as LVLM) | 9.80 (↑) | 23.14 (↑) | 40.95 (↑) | 51.04 (↑) | 68.54 (↑) |

Table 10: TallyQA Simple Benchmark. A (↑) next to the measured accuracies for LVLM-Count indicates improvement over the base LVLM it uses, while a (↓) indicates degradation compared to the corresponding base LVLM.

| Method | Acc (%) | Acc±1 (%) | Acc±3 (%) | Acc±5 (%) | Acc±10 (%) |
|---|---|---|---|---|---|
| TFOC (Shi et al., 2024) | 6.45 | 17.42 | 34.84 | 59.35 | 79.35 |
| GroundingRec (Dai et al., 2024) | 23.87 | 49.03 | 73.55 | 86.45 | 93.55 |
| CountGD (Amini-Naieni et al., 2024) | 41.94 | 63.23 | 78.06 | 85.16 | 94.84 |
| GPT4o | 44.30 | 69.89 | 87.96 | 94.41 | 100.00 |
| LVLM-Count (GPT4o as LVLM) | 44.73 (↑) | 73.33 (↑) | 91.83 (↑) | 96.99 (↑) | 99.57 (↓) |
| Gemini 1.5 Pro | 47.10 | 75.27 | 92.90 | 97.63 | 100.00 |
| LVLM-Count (Gemini 1.5 Pro as LVLM) | 45.16 (↓) | 76.13 (↑) | 91.61 (↓) | 95.91 (↓) | 99.35 (↓) |
| Qwen2 VL 72B AWQ | 49.03 | 70.75 | 87.74 | 93.33 | 98.92 |
| LVLM-Count (Qwen2 VL 72B AWQ as LVLM) | 41.72 (↓) | 67.10 (↓) | 85.81 (↓) | 93.55 (↑) | 98.49 (↓) |

Table 11: TallyQA Complex Benchmark. A (↑) next to the measured accuracies for LVLM-Count indicates improvement over the base LVLM it uses, while a (↓) indicates degradation compared to the corresponding base LVLM.

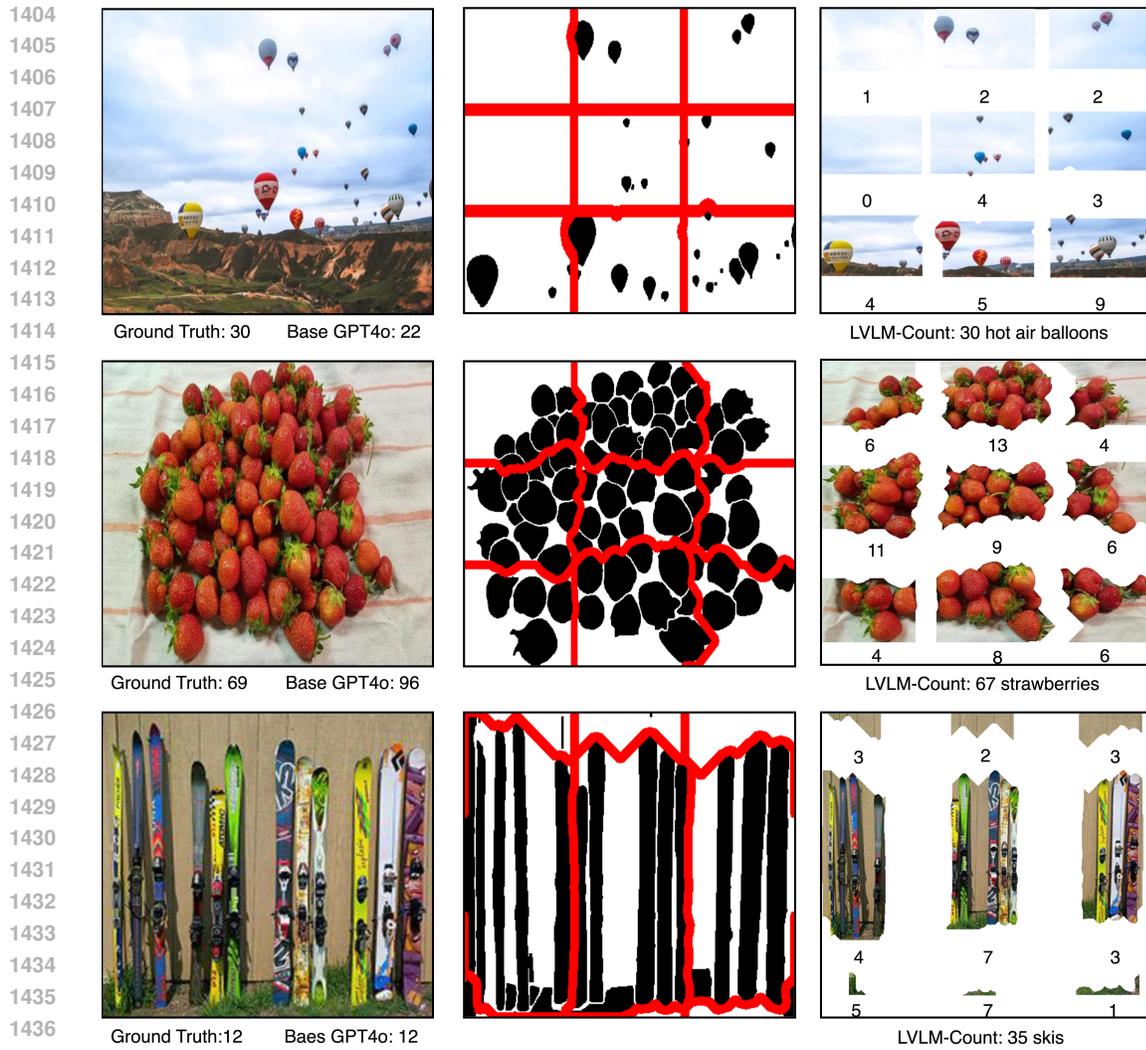| Method | Acc (%) | Acc±1 (%) | Acc±3 (%) | Acc±5 (%) | Acc±10 (%) |
|---|---|---|---|---|---|
| TFOC (Shi et al., 2024) | 1.34 | 15.44 | 32.89 | 46.31 | 66.44 |
| GroundingRec (Dai et al., 2024) | 16.78 | 28.19 | 51.68 | 64.43 | 85.23 |
| CountGD (Amini-Naieni et al., 2024) | 5.37 | 12.75 | 31.54 | 48.99 | 73.15 |
| GPT4o | 29.08 | 50.34 | 74.94 | 86.13 | 98.21 |
| LVLM-Count (GPT4o as LVLM) | 34.68 (↑) | 55.03 (↑) | 78.97 (↑) | 89.71 (↑) | 96.42 (↑) |
| Gemini 1.5 Pro | 25.73 | 25.73 | 81.66 | 91.72 | 98.88 |
| LVLM-Count (Gemini 1.5 Pro as LVLM) | 26.62 (↑) | 51.68 (↑) | 76.06 (↓) | 85.23 (↓) | 94.41 (↓) |
| Qwen2 VL 72B AWQ | 24.61 | 46.31 | 68.01 | 78.75 | 94.41 |
| LVLM-Count (Qwen2 VL 72B AWQ as LVLM) | 30.65 (↑) | 55.26 (↑) | 78.30 (↑) | 86.80 (↑) | 96.64 (↑) |

Table 12: Emoji Benchmark. A (↑) next to the measured accuracies for LVLM-Count indicates improvement over the base LVLM it uses, while a (↓) indicates degradation compared to the corresponding base LVLM.

| Method | Acc (%) | Acc±1 (%) | Acc±3 (%) | Acc±5 (%) | Acc±10 (%) |
|---|---|---|---|---|---|
| TFOC (Shi et al., 2024) | 0.24 | 0.72 | 2.89 | 4.58 | 8.43 |
| GroundingRec (Dai et al., 2024) | 4.58 | 8.43 | 12.29 | 17.11 | 24.10 |
| CountGD (Amini-Naieni et al., 2024) | 0.48 | 0.72 | 0.96 | 0.96 | 1.20 |
| GPT4o | 1.85 | 5.54 | 13.73 | 21.12 | 43.94 |
| LVLM-Count (GPT4o as LVLM) | 4.98 (↑) | 17.03 (↑) | 37.67 (↑) | 55.90 (↑) | 76.71 (↑) |
| Gemini 1.5 Pro | 2.65 | 7.23 | 14.14 | 23.37 | 42.97 |
| LVLM-Count (Gemini 1.5 Pro as LVLM) | 2.25 (↓) | 5.46 (↓) | 13.57 (↓) | 23.61 (↑) | 57.19 (↑) |
| Qwen2 VL 72B AWQ | 0.88 | 2.89 | 7.55 | 11.41 | 19.76 |
| LVLM-Count (Qwen2 VL 72B AWQ as LVLM) | 3.37 (↑) | 8.76 (↑) | 19.04 (↑) | 29.00 (↑) | 48.51 (↑) |

# M  VISUAL EXAMPLES OF LVLM-COUNT'S PERFORMANCE ON THE FSC-147 DATASET, TALLYQA BENCHMARK, AND EMOJI-COUNT BENCHMARK

This section presents several visual examples showcasing the performance of LVLM-Count on the FSC-147 dataset, the TallyQA benchmark, and the Emoji-Count benchmark. The LVLM used in the pipeline to generate these visual examples is GPT-4o. Figure 23 illustrates both successful and unsuccessful examples of LVLM-Count's performance on FSC-147. The top two rows in this figure demonstrate strong performance, where the sub-images preserve sufficient context, enabling the LVLM to predict the correct count. The last row depicts a failure case. Here, the task is to count the number of skis. However, LVLM-Count generates sub-images containing only irrelevant objects. Consequently, GPT-4o is misled by these objects and predicts an incorrect non-zero count.

Additionally, Figure 24 displays an example from the TallyQA Simple benchmark, while Figure 25 illustrates an example from the TallyQA Complex benchmark. Additionally, we include five visual

Figure 23: Illustration of good and bad examples for LVLM-Count. These examples are taken from FSC-147. The first column shows the input image. The second column shows the paths found in the object-aware division stage. The third column illustrates the obtained sub-images and the predicted count by LVLM-Count. The top two rows correspond to good examples where the base GPT-4o has a large error, but LVLM-Count provides an answer very close to the ground truth. The bottom row shows a bad example where the base GPT-4o performs better, and LVLM-Count has a large error.

examples from the Emoji-Count benchmark, split into two figures for better clarity and presentation. Figure 26 contains the first three examples, and Figure 27 includes the final two examples. In these visual examples, LVLM-Count consistently achieves more accurate results compared to the base GPT-4o.

How many people are in the image?



Ground Truth: 13    Base GPT4o: 5

LVLM-Count: 12

Figure 24: A visual example of the performance of LVLM-Count on the TallyQA Simple benchmark. The input question is: "How many people are in the image?".
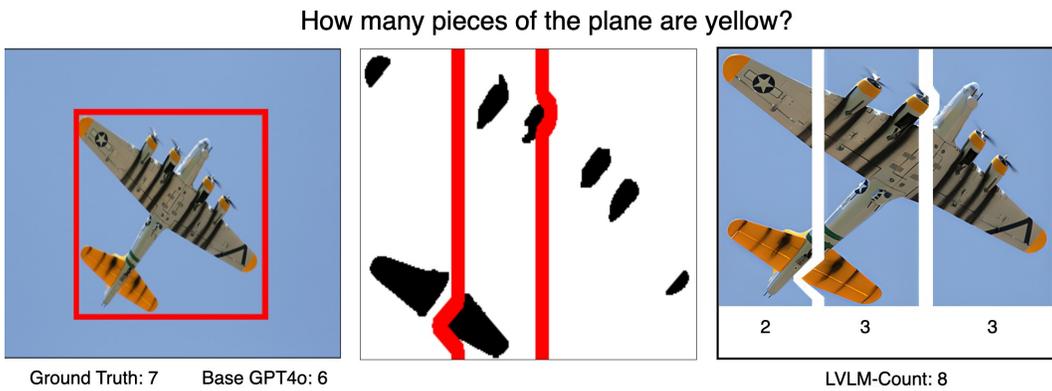
How many pieces of the plane are yellow?



Ground Truth: 7    Base GPT4o: 6

LVLM-Count: 8

Figure 25: A visual example of the performance of LVLM-Count on the TallyQA Complex benchmark. The input question is: "How many pieces of the plane are yellow?".

Figure 26: Three examples of the performance of LVLM-Count on the Emoji-Count benchmark. Top row: The object of interest is "twelve o'clock". Middle row: The object of interest is "kissing cat". Bottom row: The object of interest is "woman with red hair".
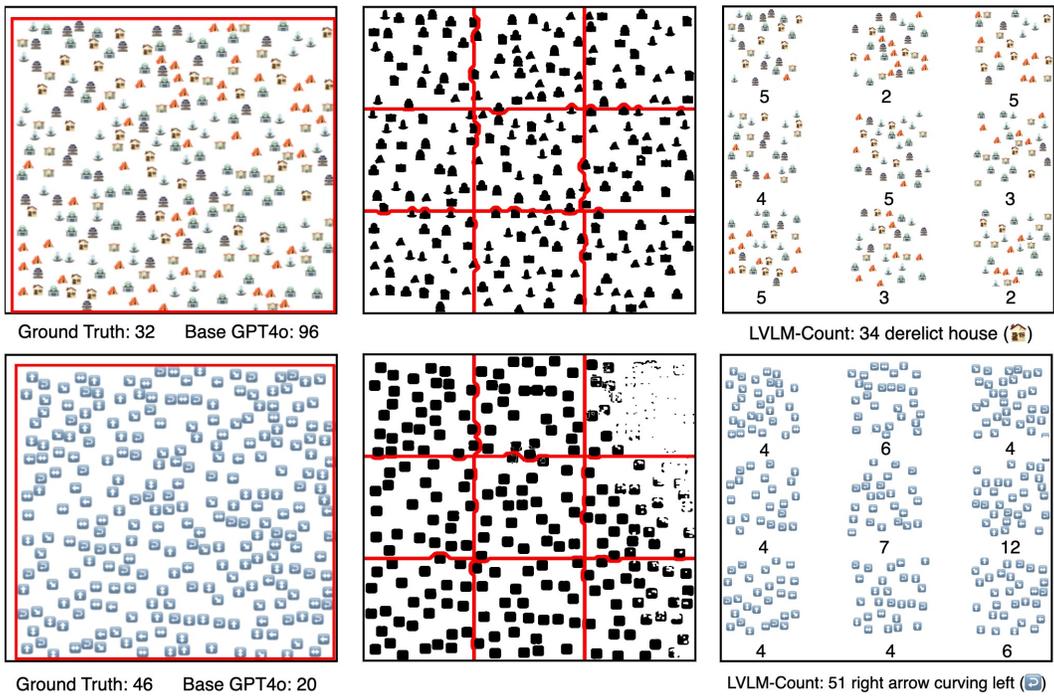
Figure 27: Two additional examples of the performance of LVLM-Count on the Emoji-Count benchmark. Top row: The object of interest is "derelict house". Bottom row: The object of interest is "right arrow curving left".