
Roboflow100-VL: A Multi-Domain Object Detection Benchmark for Vision-Language Models

Peter Robicheckaux^{1,*}, Matvei Popov^{1,*}, Anish Madan², Isaac Robinson¹,
Joseph Nelson¹, Deva Ramanan², Neehar Peri²
¹ Roboflow, ² Carnegie Mellon University

Abstract

Vision-language models (VLMs) trained on internet-scale data achieve remarkable zero-shot detection performance on common objects like car, truck, and pedestrian. However, state-of-the-art models still struggle to generalize to out-of-distribution tasks (e.g. material property estimation, defect detection, and contextual action recognition) and imaging modalities (e.g. X-rays, thermal-spectrum data, and aerial images) not typically found in their pre-training. Rather than simply re-training VLMs on more visual data (the dominant paradigm for few-shot learning), we argue that one should align VLMs to new concepts with annotation instructions containing a few visual examples *and* rich textual descriptions. To this end, we introduce Roboflow100-VL, a large-scale collection of 100 multi-modal datasets with diverse concepts not commonly found in VLM pre-training. Notably, state-of-the-art models like GroundingDINO and Qwen2.5-VL achieve less than 2% zero-shot accuracy on challenging medical imaging datasets within Roboflow100-VL, demonstrating the need for few-shot concept alignment. Our code and dataset are available on GitHub and Roboflow.

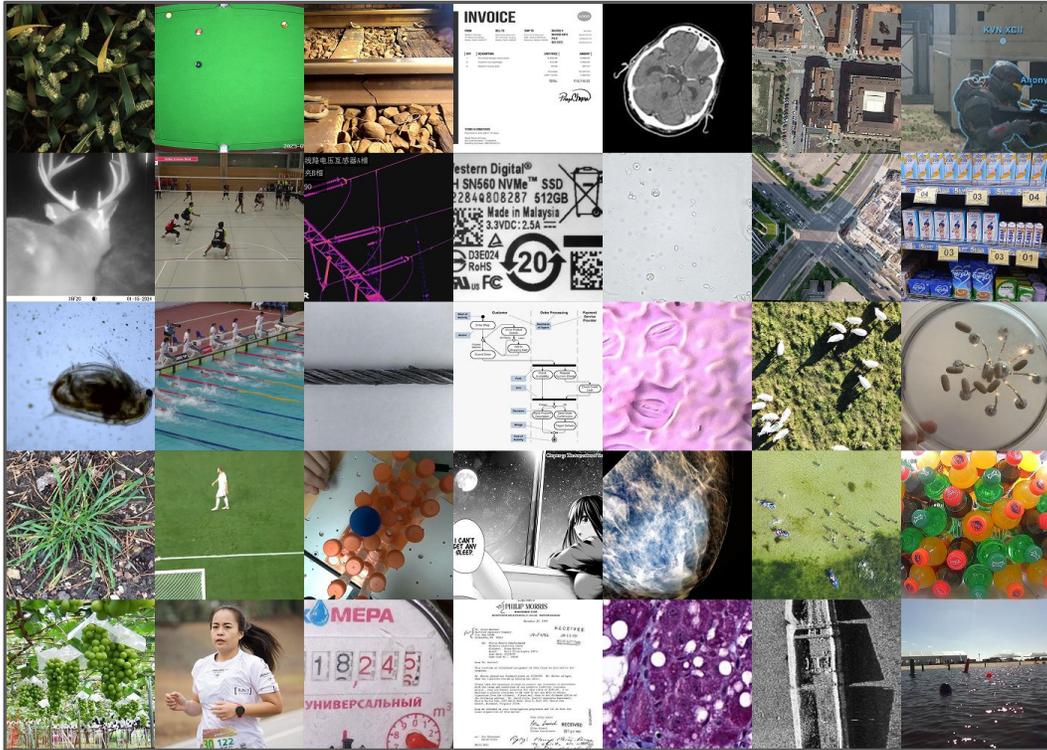
1 Introduction

Vision-language models (VLMs) trained on web-scale datasets achieve remarkable zero-shot performance on many popular academic benchmarks [54, 27, 45]. However, the performance of such foundation models varies greatly when evaluated in-the-wild, particularly on out-of-distribution tasks (e.g. material property estimation, defect detection, and contextual action recognition) and imaging modalities (e.g. X-rays, thermal spectrum data, and aerial imagery). In this paper, we introduce Roboflow100-VL, a large-scale multi-domain dataset to benchmark state-of-the-art VLMs on hundreds of diverse concepts not typically found in internet pre-training.

Status Quo. Foundation models are often trained on large-scale datasets curated from diverse sources around the web. However, despite their scale and diversity, these pre-training datasets still follow a long-tail distribution [41], causing foundation models to generalize poorly to rare concepts [36]. A common approach for improving the performance of VLMs is to scale up training data and model size [1]. However, we argue that some data will always remain out-of-distribution, whether due to being sequestered from the internet or being created after the model’s training cutoff [47], motivating the need to learn new concepts from a few examples.

Evaluating Out-of-Distribution Generalization. Existing benchmarks primarily assess VLM performance through multi-modal visual question answering (VQA) and common sense reasoning [27, 55, 45]. However, we argue that evaluating model performance on compositional reasoning benchmarks alone does not effectively measure generalization to out-of-distribution tasks. To address this limitation, we introduce Roboflow100-VL, a large-scale detection benchmark comprised of 100

*Equal Contribution



Flora & Fauna Sports Industrial Document Medical Aerial Other

Figure 1: Roboflow100-VL Dataset. We identify a set of 100 challenging datasets from Roboflow Universe that contain concepts not typically found in internet-scale pre-training. To simplify analysis, we cluster these 100 datasets using per-dataset CLIP [38] embeddings into seven categories. We visualize examples from each of these categories above. Furthermore, we also generate multi-modal instructions for each dataset with a few visual examples and rich textual descriptions per class to facilitate few-shot concept alignment.

multi-modal datasets from diverse domains (Fig. 1). Importantly, we carefully curate Roboflow100-VL such that it cannot be solved by simply prompting state-of-the-art models with class names. Specifically, we include datasets where classes are labeled using scientific names (e.g. liver fibrosis and steatosis), acronyms (e.g. DIP and MCP), context-dependent names (e.g. detecting a block vs. set in the context of volleyball), material properties (e.g. metal vs. hard plastic), and diverse imaging modalities (Fig. 2). We posit that models must leverage multi-modal contextual information (presented in the form of multi-modal annotator instructions) to align to target concepts in Roboflow100-VL.

Multi-Modal Annotator Instructions. Annotating large-scale datasets is an iterative process that often requires extensive discussions between data curators and annotators to clarify class definitions and ensure label consistency. These (often multi-modal) labeling instructions provide rich contextual information not provided by class names alone. We argue that aligning foundation models to target concepts can be principally addressed through the lens of few-shot learning by presenting vision-language models with visual examples and rich textual descriptions per class. Importantly, this approach mirrors how we align human annotators to concepts of interest with few-shot multi-modal examples [3, 30].

Contributions. We present three major contributions. First, we introduce Roboflow100-VL, a large-scale, multi-domain benchmark designed to evaluate vision-language models (VLMs) on challenging real-world use cases. We evaluate state-of-the-art models on our benchmark in zero-shot, few-shot, semi-supervised, and fully-supervised settings. Our extensive experiments highlight the difficulty

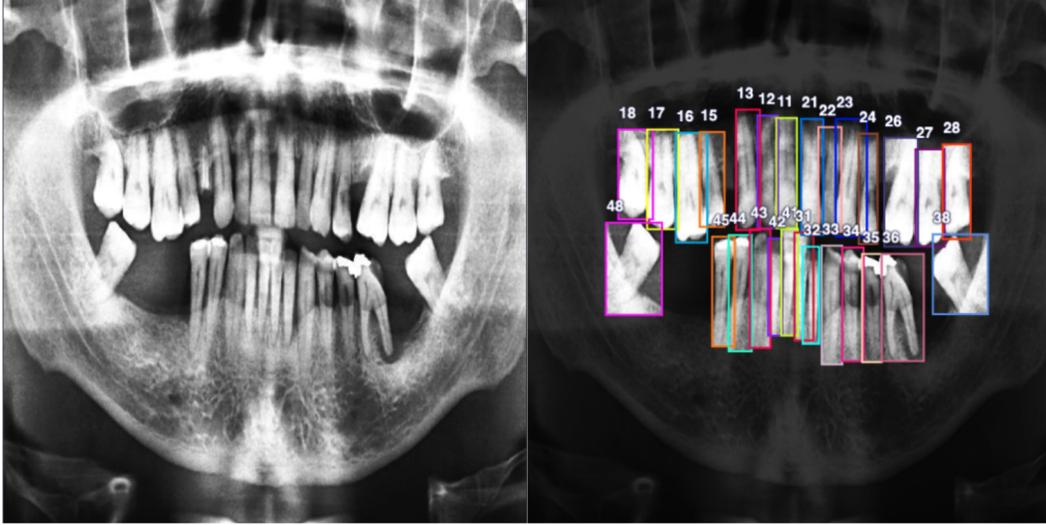


Figure 2: **Hard Examples in Roboflow100-VL.** Our dataset is particularly challenging because it is difficult to detect objects in Roboflow100-VL using class-names alone. Specifically, we select datasets where classes are labeled using scientific names, acronyms, context-dependent names, material properties, and diverse imaging modalities. For example, ufba-425-wniel-fsod-izom has numerical class-names which refer to ISO 3950 [15], a standardized dental ontology. We posit that models must leverage multi-modal contextual annotations to address such hard examples.

of adapting VLMs to out-of-distribution tasks and reveal the limitations of current state-of-the-art methods. Lastly, we host a challenge at CVPR 2025 in conjunction with the Workshop on Visual Perception via Learning in An Open World to encourage broad community involvement in addressing this challenging problem.

2 Related Works

Vision Language Models are trained using large-scale, weakly supervised image-text pairs sourced from the web. Although many vision-language models primarily focus on classification [38] or image understanding, recent methods address spatial understanding with open-vocabulary detectors. Early approaches adapted VLMs for object detection by classifying specific image regions [11, 12] or integrating detection components into frozen [20] or fine-tuned [33, 32, 9] encoders. In contrast, RegionCLIP [58] employs a multi-stage training strategy that involves generating pseudo-labels from captioning data, performing region-text contrastive pre-training, and fine-tuning on detection tasks. GLIP [23] treats detection as a phrase grounding problem by using a single text query for the entire image. Detic [59] improves long-tail detection performance by utilizing image-level supervision from ImageNet [40]. Notably, recent VLMs achieve remarkable zero-shot performance and are widely used as “black box” models in diverse downstream applications [29, 37, 19, 34]. Multi-modal large language models (MLLMs) such as Qwen2.5-VL [2] and Gemini Flash 2.0 [8] frame spatial understanding as a text generation task. Interestingly, such MLLMs perform worse at object detection than task-specific models like GroundingDINO [26].

Fine-Tuning Vision-Language Models is crucial for adapting foundation models to downstream tasks [14, 56, 10]. Traditional fine-tuning methods, such as linear probing [6, 13] and full fine-tuning [49, 50] can be computationally expensive. Instead, parameter-efficient approaches like CLIP-Adapter [10] and Tip-Adapter [57] optimize lightweight MLPs while keeping encoders frozen. Although prior few-shot learners commonly used meta-learning [52], more recent approaches show that transfer learning generalizes better to novel categories [48]. In particular, [30, 35] demonstrate that transfer learning can be effectively used to fine-tune foundation models using a few multi-modal examples. More recently, in-context learning [51] demonstrates promising results for test-time few-

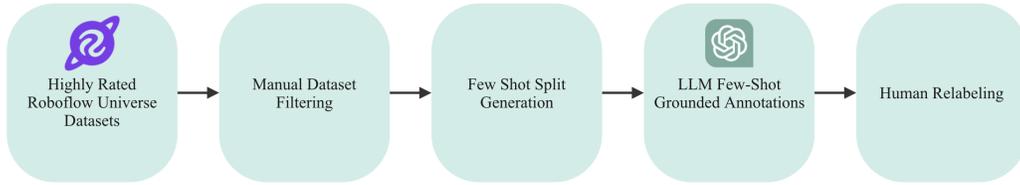


Figure 3: **Dataset Curation.** We begin by sorting all object detection datasets on Roboflow Universe by stars as a proxy for quality and usefulness to the community. Next, we manually filter all datasets with common classes, datasets where images only have a single focal object, or datasets with watermarks. We generate 10-shot splits following the protocol defined by Wang et.al. [48], where we find a subset of images with 10 total instances per class. We use these 10-shot splits to generate visually grounded “annotator instructions”, which allow VLMs to perform object detection from language and vision grounding. Finally, human labelers verify that all images within a dataset follow consistent annotation policies (e.g. bounding-box fit, semantic legibility of class names, and completeness of annotation instructions).

shot adaptation without gradient-based fine-tuning. We explore such test-time fine-tuning strategies in the context of multi-modal large language models [8, 2].

Benchmarking Vision-Language Models is of significant interest to the community. State-of-the-art VLMs are typically evaluated using benchmarks such as MMStar [4], MMMU [55], MME [24], ScienceQA [28], MMBench [27], MM-Vet [54], Seed-Bench [21], and MMVP [46]. These benchmarks evaluate a broad set of vision-language tasks, including fine-grained perception, reasoning, common sense knowledge, and problem solving in various domains. However, existing evaluations primarily focus on multi-modal understanding in the context of Visual Question Answering (VQA). In contrast, Roboflow100-VL evaluates VLM detection accuracy given a few visual examples and rich textual descriptions. Prior VLM grounding benchmarks like RefCOCO [53] often focus on referential grounding of common object categories. Recent efforts like ODinW [22] consider more challenging scenarios by sourcing real-world data from Roboflow [7]. However, we find that state-of-the-art methods achieve high zero-shot accuracy on ODinW [2], suggesting that these datasets may not be well suited for evaluating foundational few-shot object detection [30].

3 Roboflow100-VL Benchmark

As shown in Fig. 1, Roboflow100-VL consists of diverse datasets not typically found in internet-scale pre-training. We highlight our data curation procedure, and present several baselines to evaluate state-of-the-art models in zero-shot, few-shot, semi-supervised, and fully-supervised settings.

3.1 Creating Roboflow100-VL

We source our datasets from Roboflow Universe, a community-driven platform that hosts diverse open-source datasets created to solve real-world computer vision tasks. With more than 500,000 public datasets spanning medical imaging, agriculture, robotics, and manufacturing, we focus on selecting high-quality datasets not commonly found in internet-scale pre-training (e.g. COCO [25], Objects365 [42], GoldG [17], CC4M [43]) to better assess VLM generalization to rare concepts. When selecting candidates for Roboflow100-VL, we prioritized datasets where images contained multiple objects, ensuring more realistic evaluation beyond classification. In addition, we sought out datasets with semantically ambiguous class names (e.g. “button” can refer to both clothing and electronics) to encourage algorithms to leverage multi-modal annotator instructions rather than simply relying on labels. We manually validate the labeling quality of each dataset to ensure exhaustive annotations. In cases without exhaustive annotations, we manually re-annotate the dataset to the best of our ability (Fig 3).

Multi-Modal Annotation Generation. Annotator instructions offer precise class definitions and visual examples that help clarify annotation policies (e.g. by highlighting typical cases, corner cases, and negative examples) and improve labeling accuracy. Despite providing significant value during

Dataset Type	# Classes	# Images	# Anno.
Aerial	29	11,627	186,789
Document	88	21,418	127,129
Flora & Fauna	70	46,718	441,677
Industrial	122	29,758	205,627
Medical	77	16,369	125,433
Sports	36	8,443	58,508
Other	142	29,816	210,328
All	564	164,149	1,355,491

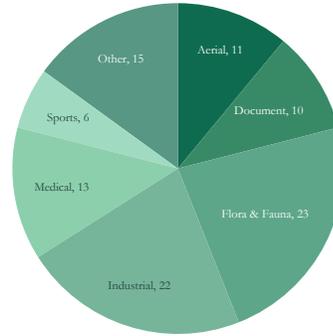


Figure 4: **Dataset Statistics.** The table on the left provides details on the number of classes, images, and annotations across different dataset types. The figure on the right illustrates the distribution of dataset types by count.

the labeling process, few datasets publicly release these annotator instructions. Recognizing the importance of these instructions in aligning humans with target concepts of interest, we generate multi-modal annotator instructions for all 100 datasets within Roboflow100-VL.

We prompt GPT-4o [1] to generate an initial set of annotator instructions, providing in-context examples based on the nulImages annotator guidelines. Our prompt includes a structured output template, along with dataset metadata, class names, and few-shot visual examples per class. In practice, we find that GPT-4o often overlooks the few-shot images and instead relies heavily on class names to generate class descriptions. Notably, GPT-4o struggles when class names are uninformative and sometimes produces overly vague instructions that, while correct, lack useful detail. To address this, we manually verify all generated annotator instructions to mitigate hallucinations and incorporate additional informative details missed by the model.

Dataset Statistics. Figure 4 (right) presents an overview of the different dataset types within Roboflow100-VL, detailing the number of classes, images, and annotations per cluster. Roboflow100-VL contains a total of 564 classes and 164,149 images, with over 1.3 million annotations. The “Other” category has the highest number of classes (142), followed by “Industrial” (122) and “Flora & Fauna” (70). Despite having fewer classes, the “Flora & Fauna” category has the highest number of images (46,718) and annotations (441,677), indicating a higher density of labeled data. Figure 4 (left) provides a visual representation of class distribution, reinforcing the dominance of the “Other”, “Industrial”, and “Flora & Fauna” categories. In contrast, “Sports” has the fewest classes (36) and the least representation in Roboflow100-VL. Despite consisting of 100 datasets, Roboflow100-VL has about half the number of images as COCO [25], making this an approachable benchmark for the academic community.

3.2 State-of-the-Art Baselines

We train and evaluate all models on each dataset within Roboflow100-VL independently. Importantly, we do not tune any parameters or modify zero-shot prompts per-dataset.

Zero-Shot Baselines prompt models with expressive descriptions or class names [31] to guide foundation models toward target concepts. However, the effectiveness of zero-shot prompting depends on the pre-training data: If the target class name is semantically meaningful and aligns well with the model’s foundational pre-training, performance is strong; otherwise, the model fails catastrophically. We benchmark the zero-shot performance of Detic [59], GroundingDINO [26], MQ-GLIP [51], QwenV2.5-VL [2] and Gemini Flash 2.0 [8].

Few-Shot Baselines. We evaluate two types of few-shot baselines: visual prompting and multi-modal prompting. Visual prompting uses images of target concepts that are difficult to describe through text as prompts to help models learn novel concepts in-context. For example, while “hard plastic” is a broad and ambiguous category that is hard to define textually, providing image examples improves concept alignment. Typically, visual prompts are tokenized and fed as inputs to a frozen VLM. Here, we apply MQ-GLIP [51] with image prompting. Multi-modal prompting combines language and

visual prompts to leverage multi-modal features. Intuitively, using both text and images yields better alignment than using either modality alone. In the case of “soft plastic”, ambiguous concepts can be clarified with textual descriptions (e.g., “thin plastic film” and “plastic bag”) alongside visual examples. Both visual and language prompts are tokenized and separately fed into a frozen VLM. We evaluate MQ-GLIP [51], and Gemini Flash 2.0 [8] by prompting models with class names, few-shot images, and annotator instructions.

Semi-Supervised Baselines. We evaluate variants of YOLO [16, 18] and YOLO with STAC [44] trained on 10% of each dataset in Roboflow100-VL. STAC generates high-confidence pseudo-labels for localized objects in unlabeled images and updates the model by enforcing consistency through strong augmentations.

Fully-Supervised Baselines. We benchmark YOLOv8 [16], YOLOv11 [18], and LW-DETR [5] on all datasets within Roboflow100-VL. YOLOv8, developed by Ultralytics, builds on the YOLOv5 architecture with improvements in model scaling and architectural refinements. YOLOv11 adds more architecture improvements, validated on COCO. LW-DETR is a lightweight detection transformer that outperforms YOLO models for real-time object detection. Its architecture consists of a ViT encoder, a projector, and a shallow DETR decoder. This baseline serves as an upper bound on performance, though in rare cases, few-shot foundation models may surpass it when the target dataset only has a few examples. For all models, we follow the standard established in [7] and train for 100 epochs with batch size 16.

4 Experiments

We conduct extensive experiments to evaluate the performance of state-of-the-art models on Roboflow100-VL. We present our zero-shot, few-shot, semi-supervised, and fully supervised results below. See Appendix A for additional implementation details.

Datasets and Metrics. Each dataset is independently evaluated using AP. We report the average accuracy per super-category to simplify analysis. Roboflow100-VL includes datasets that are out-of-distribution from typical internet-scale pre-training data, making it particularly challenging (even for VLMs). To construct the few-shot split, we follow the K -shot dataset creation process established by [48]. To construct the semi-supervised split, we randomly sample 10% of the training set. Importantly, all methods are evaluated on the same test set.

4.1 Empirical Analysis of Results.

We benchmark state-of-the-art methods and present our results from Table 1 below.

State-of-the-Art Zero-Shot and Few-Shot Models Struggle on Roboflow100-VL. Roboflow100-VL is a much harder dataset than prior open-vocabulary object detection benchmarks. Specifically, GroundingDINO achieves 49.2 mAP on ODinW35, but only reaches 15 mAP on Roboflow100-VL. Similar trends can be seen with Qwen and OWL-ViT2. Furthermore, both zero-shot and few-shot models perform significantly worse on Roboflow100-VL than on COCO, suggesting that our dataset curation policies highlight a data bias in VLM pre-training towards common categories.

Open-Vocabulary Object Detectors Outperform MLLMs. We find that open-vocabulary object detectors like Detic, GroundingDINO, OWL-ViT2, and MQ-GLIP consistently outperform multi-modal LLMs (MLLMs) like Qwen 2.5 VL, Gemini Flash 2.0, despite MLLMs pre-training on orders of magnitude more data. This highlights the advantage of task-specific architectures over generalist models.

Multi-Modal Annotator Instructions Provide Limited Benefit. Somewhat surprisingly, state-of-the-art MLLMs struggle to benefit from multi-modal annotator instructions. In fact, prompting with instructions provides inconsistent benefit compared to prompting with class names. Intuitively, we expect annotator instructions to improve object detection performance by resolving semantic ambiguity in class names and providing rich contextual information. However, we posit that this performance decline can be attributed to the fact that MLLMs are instruction-tuned for open vocabulary detection with rigid prompt structures, making it difficult to effectively leverage additional contextual information.

Table 1: **Roboflow100-VL Benchmarks.** We evaluate the zero-shot, few-shot, semi-supervised, and fully-supervised performance of state-of-the-art methods on the Roboflow100-VL benchmark. We find that Roboflow100-VL is particularly challenging for zero-shot and few-shot approaches, with most methods struggling to achieve 10% mAP averaged over all 100 datasets. In contrast, we find that semi-supervised learners are able to reach nearly 80% of the performance of fully supervised models using 10% labeled data.

Method	Aerial	Document	Flora & Fauna	Industrial	Medical	Sports	Other	All
Zero-Shot								
Detic [59]	12.2	4.5	17.9	6.0	0.8	7.6	11.2	9.5
GroundingDINO [26]	21.8	7.9	28.2	10.3	2.1	13.0	18.1	15.7
OWL-ViT2 [32] (Class Names Only)	15.7	8.0	18.4	6.8	1.2	10.2	10.5	10.6
MQ-GLIP-Text [51](Class-Names Only)	11.9	9.7	22.6	7.7	1.4	9.2	14.1	12.0
Qwen 2.5 VL (72B) [2] (Class Names Only)	4.4	5.8	9.2	4.5	2.21	8.8	5.7	5.8
Qwen 2.5 VL (72B) [2] (Instructions Only)	4.7	6.2	13.3	5.4	1.2	9.8	8.1	7.3
Gemini Flash 2.0 [8] (Class Names Only)	6.0	2.9	18.1	3.9	1.1	5.0	9.5	7.8
Gemini Flash 2.0 [8] (Instructions Only)	3.0	1.6	9.4	1.9	0.3	2.8	5.5	4.1
Few-Shot (10 shots)								
MQ-GLIP-Image [51] (Images Only)	4.4	3.0	13.0	3.8	1.4	7.4	6.8	6.2
MQ-GLIP [51] (Class Names + Images)	11.9	9.2	22.6	7.7	1.4	9.3	14.1	12.0
Gemini Flash 2.0 [8] (Instructions + Images)	2.3	1.15	4.9	4.0	1.10	0.09	2.3	2.9
Semi-Supervised (10% Labels)								
YOLOv8n [16]	35.0	35.7	42.0	51.7	29.5	32.0	38.8	40.0
YOLOv8n [16] w/ STAC [44]	39.0	39.8	45.0	53.5	33.2	36.2	44.0	43.5
YOLOv8s [16]	39.4	40.5	42.5	53.5	34.4	40.9	44.0	43.5
YOLOv8s [16] w/ STAC [44]	41.1	42.5	45.5	55.8	36.4	43.5	46.7	45.8
YOLOv8m [16]	39.7	42.7	44.1	54.1	33.7	45.3	46.7	44.8
YOLOv8m [16] w/ STAC [44]	41.6	45.7	46.1	55.6	35.8	47.1	49.2	46.8
Fully-Supervised								
YOLOv8n [16]	52.8	57.6	55.5	66.4	51.2	52.3	57.5	57.4
YOLOv1n [18]	52.1	57.4	55.2	66.5	51.8	52.7	57.6	57.3
YOLOv8s [16]	55.4	60.0	56.9	67.5	52.8	55.0	60.0	59.2
YOLOv1s [18]	54.3	60.3	56.8	67.6	51.8	56.0	60.1	59.0
LW-DETRs [5]	52.0	59.4	55.1	68.2	51.4	54.8	57.7	58.0
YOLOv8m [16]	56.5	62.2	57.3	67.4	52.1	57.2	60.8	59.8
YOLOv1m [18]	55.1	61.8	57.1	68.4	51.9	56.4	60.8	59.7
LW-DETRm [5]	51.9	57.8	56.0	66.5	51.2	53.6	57.8	57.5

Semi-Supervised Learners Are Data Efficient. We find that leveraging simple semi-supervised learning algorithms like STAC [44] significantly improve model performance when learning with limited labels. In a majority (8 out of 14) of combinations of model size and data domain, using a semi-supervised method yielded at least as much improvement in mAP as stepping up a model size. For example, training a YOLOv8n on 10% labeled data with STAC achieves the same performance as a YOLOv8s trained on 10% labeled data.

Supervised Object Detectors Overfit Training and Architecture Decisions to COCO. Real-time object detectors are often optimized for COCO, assuming better performance on COCO translates to real-world improvements. However, real-world datasets (such as those in Roboflow100-VL) are often much smaller and more diverse than COCO, challenging this assumption. Specifically, although Roboflow100-VL has half as many images as COCO, it has more than seven times as many classes (Fig. 4). Interestingly, we find that models that achieved higher performance on COCO did not necessarily improve real-world performance on Roboflow100-VL – both within and across model families. For example, YOLOv11 outperforms YOLOv8 on COCO but underperforms across all three tested sizes (nano, small, medium) on Roboflow100-VL. This suggests that newer YOLO models may be overfitting to COCO. We find similar trends with LW-DETR. Lastly, we find that increasing model size leads to smaller performance improvements on Roboflow100-VL compared to COCO. The performance difference between the smallest and largest models within a model family is within 2.5 mAP, suggesting that simply increasing model capacity may not lead to significant performance gains on Roboflow100-VL.

4.2 Limitations and Future Work

Reliance on Crowdsourced Annotations. All our datasets are sourced from Roboflow Universe, a community platform where anyone can upload dataset annotations. Although this allows us to source diverse datasets, it introduces uncertainty regarding overall annotation quality. While we manually inspect all datasets to ensure quality to the best of our ability, verifying annotations in specialized domains like medical imaging remains a significant challenge.

Generated Annotator Instructions May Not Reflect Real Instructions. Our annotator instructions are automatically generated by GPT-4o and are manually verified for correctness. However, they may not fully reflect the nuances of real-world instructions typically developed alongside dataset collection. We encourage the community to release real annotator instructions generated through iterative discussions between annotators and stakeholders. Furthermore, although our annotator instructions provide high-level class descriptions, they often do not directly incorporate image evidence to identify typical cases, edge cases, and negative examples.

5 Conclusion

In this paper, we introduce Roboflow100-VL, a large-scale benchmark to evaluate state-of-the-art VLMs on concepts not typically found in internet-scale pre-training. Roboflow100-VL is curated to evaluate detection performance on out-of-distribution tasks (e.g. material property estimation, defect detection, and contextual action recognition) and imaging modalities (e.g. X-rays, thermal spectrum data, and aerial imagery) using a few visual examples and rich textual descriptions. We find that state-of-the-art models struggle on this challenging benchmark, demonstrating the limitations of existing methods, and highlighting opportunities for future work.

6 Acknowledgments

This work was supported in part by compute provided by NVIDIA, and the NSF GRFP (Grant No. DGE2140739).

References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altmenschmidt, Sam Altman, Shyamal Anadkat, et al. “Gpt-4 technical report”. In: *arXiv preprint arXiv:2303.08774* (2023).
- [2] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. “Qwen2. 5-vl technical report”. In: *arXiv preprint arXiv:2502.13923* (2025).
- [3] Nadine Chang, Francesco Ferroni, Michael J Tarr, Martial Hebert, and Deva Ramanan. “Thinking Like an Annotator: Generation of Dataset Labeling Instructions”. In: *arXiv preprint arXiv:2306.14035* (2023).
- [4] Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao, Dahua Lin, et al. “Are we on the right way for evaluating large vision-language models?” In: *arXiv preprint arXiv:2403.20330* (2024).
- [5] Qiang Chen, Xiangbo Su, Xinyu Zhang, Jian Wang, Jiahui Chen, Yunpeng Shen, Chuchu Han, Ziliang Chen, Weixiang Xu, Fanrong Li, et al. “LW-DETR: a transformer replacement to yolo for real-time detection”. In: *arXiv preprint arXiv:2406.03459* (2024).
- [6] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. “A simple framework for contrastive learning of visual representations”. In: *International conference on machine learning*. PMLR. 2020, pp. 1597–1607.
- [7] Floriana Ciaglia, Francesco Saverio Zuppichini, Paul Guerrie, Mark McQuade, and Jacob Solawetz. “Roboflow 100: A rich, multi-domain object detection benchmark”. In: *arXiv preprint arXiv:2211.13523* (2022).
- [8] Google DeepMind. *Introducing Gemini 2.0: our new AI model for the agentic era*. Dec. 2024. URL: <https://blog.google/technology/google-deepmind/google-gemini-ai-update-december-2024/>.
- [9] Yu Du, Fangyun Wei, Ziheng Zhang, Miaoqing Shi, Yue Gao, and Guoqi Li. “Learning to prompt for open-vocabulary object detection with vision-language model”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 14084–14093.
- [10] Peng Gao, Shijie Geng, Renrui Zhang, Teli Ma, Rongyao Fang, Yongfeng Zhang, Hongsheng Li, and Yu Qiao. “Clip-adapter: Better vision-language models with feature adapters”. In: *International Journal of Computer Vision* 132.2 (2024), pp. 581–595.
- [11] Xiuye Gu, Tsung-Yi Lin, Weicheng Kuo, and Yin Cui. “Open-vocabulary object detection via vision and language knowledge distillation”. In: *arXiv preprint arXiv:2104.13921* (2021).
- [12] Xiuye Gu, Tsung-Yi Lin, Weicheng Kuo, and Yin Cui. “Open-vocabulary object detection via vision and language knowledge distillation”. In: *arXiv preprint arXiv:2104.13921* (2021).
- [13] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. “Masked autoencoders are scalable vision learners”. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2022, pp. 16000–16009.
- [14] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. “Lora: Low-rank adaptation of large language models”. In: *arXiv preprint arXiv:2106.09685* (2021).
- [15] International Organization for Standardization. *ISO 56007:2024 Dentistry designation system for teeth and areas of the oral cavity*. 2024. URL: <https://www.iso.org/standard/68292.html>.
- [16] Glenn Jocher, Jing Qiu, and Ayush Chaurasia. *Ultralytics YOLO*. Version 8.0.0. Jan. 2023. URL: <https://github.com/ultralytics/ultralytics>.
- [17] Gaoussou Youssouf Kebe, Pdraig Higgins, Patrick Jenkins, Kasra Darvish, Rishabh Sachdeva, Ryan Barron, John Winder, Donald Engel, Edward Raff, Francis Ferraro, and Cynthia Matuszek. “A Spoken Language Dataset of Descriptions for Speech-Based Grounded Language Learning”. In: *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*. 2021.
- [18] Rahima Khanam and Muhammad Hussain. “Yolov11: An overview of the key architectural enhancements”. In: *arXiv preprint arXiv:2410.17725* (2024).
- [19] Mehar Khurana, Neehar Peri, Deva Ramanan, and James Hays. “Shelf-Supervised Multi-Modal Pre-Training for 3D Object Detection”. In: *arXiv preprint arXiv:2406.10115* (2024).

- [20] Weicheng Kuo, Yin Cui, Xiuye Gu, AJ Piergiovanni, and Anelia Angelova. “F-vlm: Open-vocabulary object detection upon frozen vision and language models”. In: *arXiv preprint arXiv:2209.15639* (2022).
- [21] Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. “Seed-bench: Benchmarking multimodal llms with generative comprehension”. In: *arXiv preprint arXiv:2307.16125* (2023).
- [22] Chunyuan Li, Haotian Liu, Liunian Harold Li, Pengchuan Zhang, Jyoti Aneja, Jianwei Yang, Ping Jin, Houdong Hu, Zicheng Liu, Yong Jae Lee, and Jianfeng Gao. “ELEVATER: A Benchmark and Toolkit for Evaluating Language-Augmented Visual Models”. In: *Neural Information Processing Systems* (2022).
- [23] Liunian Harold Li, Pengchuan Zhang, Haotian Zhang, Jianwei Yang, Chunyuan Li, Yiwu Zhong, Lijuan Wang, Lu Yuan, Lei Zhang, Jenq-Neng Hwang, et al. “Grounded language-image pre-training”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 10965–10975.
- [24] Zijiang Liang, Yanjie Xu, Yifan Hong, Penghui Shang, Qi Wang, Qiang Fu, and Ke Liu. “A Survey of Multimodal Large Language Models”. In: *Proceedings of the 3rd International Conference on Computer, Artificial Intelligence and Control Engineering*. 2024, pp. 405–409.
- [25] 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: *Computer vision—ECCV 2014: 13th European conference, zurich, Switzerland, September 6-12, 2014, proceedings, part v 13*. Springer. 2014, pp. 740–755.
- [26] 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”. In: *arXiv preprint arXiv:2303.05499* (2023).
- [27] Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. “Mmbench: Is your multi-modal model an all-around player?” In: *European conference on computer vision*. Springer. 2024, pp. 216–233.
- [28] Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. “Learn to explain: Multimodal reasoning via thought chains for science question answering”. In: *Advances in Neural Information Processing Systems* 35 (2022), pp. 2507–2521.
- [29] Yechi Ma, Neehar Peri, Shuoquan Wei, Wei Hua, Deva Ramanan, Yanan Li, and Shu Kong. “Long-Tailed 3D Detection via 2D Late Fusion”. In: *arXiv preprint arXiv:2312.10986* (2023).
- [30] Anish Madan, Neehar Peri, Shu Kong, and Deva Ramanan. “Revisiting few-shot object detection with vision-language models”. In: *Advances in Neural Information Processing Systems* 37 (2024), pp. 19547–19560.
- [31] Sachit Menon and Carl Vondrick. “Visual Classification via Description from Large Language Models”. In: *The Eleventh International Conference on Learning Representations (ICLR)*. 2023.
- [32] Matthias Minderer, Alexey Gritsenko, and Neil Houlsby. “Scaling Open-Vocabulary Object Detection”. In: *arXiv preprint arXiv:2306.09683* (2023).
- [33] Matthias Minderer, Alexey Gritsenko, Austin Stone, Maxim Neumann, Dirk Weissenborn, Alexey Dosovitskiy, Aravindh Mahendran, Anurag Arnab, Mostafa Dehghani, Zhuoran Shen, et al. “Simple open-vocabulary object detection”. In: *European Conference on Computer Vision*. Springer. 2022, pp. 728–755.
- [34] Aljosa Osep, Tim Meinhardt, Francesco Ferroni, Neehar Peri, Deva Ramanan, and Laura Leal-Taixe. “Better Call SAL: Towards Learning to Segment Anything in Lidar”. In: *ECCV*. 2024.
- [35] Hongpeng Pan, Shifeng Yi, Shouwei Yang, Lei Qi, Bing Hu, Yi Xu, and Yang Yang. “The Solution for CVPR2024 Foundational Few-Shot Object Detection Challenge”. In: *arXiv preprint arXiv:2406.12225* (2024).
- [36] Shubham Parashar, Zhiqiu Lin, Tian Liu, Xiangjue Dong, Yanan Li, Deva Ramanan, James Caverlee, and Shu Kong. “The Neglected Tails in Vision-Language Models”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2024, pp. 12988–12997.
- [37] Neehar Peri, Achal Dave, Deva Ramanan, and Shu Kong. “Towards Long-Tailed 3D Detection”. In: 2023.

- [38] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. “Learning transferable visual models from natural language supervision”. In: *International conference on machine learning*. PmLR. 2021, pp. 8748–8763.
- [39] Roboflow. *Roboflow Inference*. URL: <https://github.com/roboflow/inference>.
- [40] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. “Imagenet large scale visual recognition challenge”. In: *International journal of computer vision* 115 (2015), pp. 211–252.
- [41] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. “Laion-5b: An open large-scale dataset for training next generation image-text models”. In: *Advances in neural information processing systems* 35 (2022), pp. 25278–25294.
- [42] Shuai Shao, Zeming Li, Tianyuan Zhang, Chao Peng, Gang Yu, Xiangyu Zhang, Jing Li, and Jian Sun. “Objects365: A Large-Scale, High-Quality Dataset for Object Detection”. In: *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*. 2019, pp. 8429–8438. DOI: 10.1109/ICCV.2019.00852.
- [43] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. “Conceptual Captions: A Cleaned, Hypernymed, Image Alt-text Dataset For Automatic Image Captioning”. In: *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Melbourne, Australia: Association for Computational Linguistics, 2018, pp. 2556–2565.
- [44] Kihyuk Sohn, Zizhao Zhang, Chun-Liang Li, Han Zhang, Chen-Yu Lee, and Tomas Pfister. “A simple semi-supervised learning framework for object detection”. In: *arXiv preprint arXiv:2005.04757* (2020).
- [45] Tristan Thrush, Ryan Jiang, Max Bartolo, Amanpreet Singh, Adina Williams, Douwe Kiela, and Candace Ross. “Winoground: Probing vision and language models for visio-linguistic compositionality”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 5238–5248.
- [46] Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. “Eyes wide shut? exploring the visual shortcomings of multimodal llms”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2024, pp. 9568–9578.
- [47] Liyuan Wang, Xingxing Zhang, Hang Su, and Jun Zhu. “A comprehensive survey of continual learning: Theory, method and application”. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2024).
- [48] Xin Wang, Thomas E. Huang, Trevor Darrell, Joseph E Gonzalez, and Fisher Yu. “Frustratingly Simple Few-Shot Object Detection”. In: *International Conference on Machine Learning (ICML)*. 2020.
- [49] Yu-Xiong Wang, Deva Ramanan, and Martial Hebert. “Growing a brain: Fine-tuning by increasing model capacity”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2017, pp. 2471–2480.
- [50] Wenhao Wu, Zhun Sun, and Wanli Ouyang. “Revisiting classifier: Transferring vision-language models for video recognition”. In: *Proceedings of the AAAI conference on artificial intelligence*. Vol. 37. 3. 2023, pp. 2847–2855.
- [51] Yifan Xu, Mengdan Zhang, Chaoyou Fu, Peixian Chen, Xiaoshan Yang, Ke Li, and Changsheng Xu. “Multi-modal queried object detection in the wild”. In: *Advances in Neural Information Processing Systems* 36 (2024).
- [52] Xiaopeng Yan, Ziliang Chen, Anni Xu, Xiaoxi Wang, Xiaodan Liang, and Liang Lin. “Meta r-cnn: Towards general solver for instance-level low-shot learning”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2019, pp. 9577–9586.
- [53] Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. “Modeling context in referring expressions”. In: *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part II 14*. Springer. 2016, pp. 69–85.
- [54] Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. “Mm-vet: Evaluating large multimodal models for integrated capabilities”. In: *arXiv preprint arXiv:2308.02490* (2023).

- [55] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. “Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2024, pp. 9556–9567.
- [56] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. “Adding conditional control to text-to-image diffusion models”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2023, pp. 3836–3847.
- [57] Renrui Zhang, Rongyao Fang, Wei Zhang, Peng Gao, Kunchang Li, Jifeng Dai, Yu Qiao, and Hongsheng Li. “Tip-adapter: Training-free clip-adapter for better vision-language modeling”. In: *arXiv preprint arXiv:2111.03930* (2021).
- [58] Yiwu Zhong, Jianwei Yang, Pengchuan Zhang, Chunyuan Li, Noel Codella, Liunian Harold Li, Luwei Zhou, Xiyang Dai, Lu Yuan, Yin Li, et al. “Regionclip: Region-based language-image pretraining”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 16793–16803.
- [59] Xingyi Zhou, Rohit Girdhar, Armand Joulin, Philipp Krähenbühl, and Ishan Misra. “Detecting twenty-thousand classes using image-level supervision”. In: *European Conference on Computer Vision*. Springer. 2022, pp. 350–368.

A Implementation Details

Detic. We use Detic [59] with a SWIN-L backbone for all zero-shot experiments. Additionally, we use the model checkpoint trained on LVIS, COCO and ImageNet-21K. We use class names provided as text prompts for Detic’s CLIP classifier.

GroundingDINO. We use GroundingDINO [26] with pretrained weights from mmdetection MM-GroundingDINO-L*. We prompt the model with all the class names combined into a single prompt.

MQ-GLIP. MQ-Det [51] proposes a learnable module that enables multi-modal prompting. We choose GLIP with a SWIN-L backbone as the underlying detection model for our experiments. We use the model checkpoint trained on Objects365, FourODs, GoldG, Cap24M. Lastly, we use class names as the text prompts and few-shot visual examples as visual prompts.

OWL-ViT2. We use OWL-ViT2 [32] as implemented in the Roboflow inference package [39]. We prompt the model with all the class names combined into a single prompt.

Qwen-2.5VL. We conduct all experiments using Qwen2.5-VL’s 72B model via API. We prompt the model based on guidelines from Qwen’s official documentation:

```
"Outline the position of each of the following objects: (list of categories), and output all the coordinates in JSON format."
```

We also evaluated Qwen with dataset-specific annotator instructions for improved detection performance. We modified the above prompt to include:

```
"Use the annotator instructions for the dataset that this image belongs to, to aid better detection: (dataset’s annotator instructions) Outline the position of each of the following objects: (list of categories), and output all the coordinates in JSON format."
```

We implemented a robust parser to handle minor JSON formatting errors while preserving correct predictions. To speed up inference, we limited Qwen2.5-VL to only generate a maximum of 6144 tokens.

Gemini Flash 2.0. We conducted all experiments using the Gemini Flash 2.0 API. We prompt the model based on guidelines from Gemini’s official documentation:

```
"Return bounding boxes as a JSON array with labels. Never return masks or code fencing. Detect all instances of all objects requested by prompt."  
"Detect the 2d bounding boxes of the following objects: (list of categories)"
```

Similar to Qwen2.5-VL, we prompted the model with annotator instructions and few-shot visual examples (10-examples per class). We use the following system prompt for all three modes.

```
"Return bounding boxes as a JSON array with labels. Never return masks or code fencing. Use the attached (dataset name) dataset annotator instructions and the few-shot examples as a reference for better predictions."
```

We structured all prompts to place annotator instructions and few-shot examples first, while keeping the original simple class-based detection prompt at the end to maintain a familiar format for the model immediately before generation.

We limited all prediction generations to 6144 tokens to speed up inference. We implement a robust parser to handle minor JSON formatting errors. In some cases with many few-shot image examples, the API failed to return a valid response for requests of excessive size. We simply assign a score of 0 AP in such cases.

YOLOv8 and YOLOv11. We train our YOLOv8 [16] and YOLOv11 [18] family of models using the Ultralytics package with default parameters. We use a batch size of 16 and train for the default of 100 epochs.

STAC. We follow the training protocol defined by Sohn et. al. [44]. First, we train a teacher model on the labeled subset of the data. Then, we use the teacher model to pseudo-label the remaining unlabeled subset of the data. We keep all detections above a confidence C , where the confidence is tuned to maximize the F1 score of the teacher model on a validation set. Finally, we combine the subset of data with true ground truth labels and the subset with pseudo-labels to form a training set for a student model of the same architecture. We train this student model until convergence with heavy augmentations. We use the same hyperparameters as our supervised YOLOv8 and YOLOv11 implementation.

B Sample Annotation Instructions

We present sample annotator instructions below. We use dataset metadata, class names and few-shot visual examples and prompt GPT-4o [1] to generate annotator instructions. We then manually verify that the instructions accurately describe the few-shot examples. These annotator instructions are from recode-waste-czvmg-fsod-yxsw.

```
# Overview
- [Introduction] (#introduction)
- [Object Classes] (#object-classes)
  - [Aggregate] (#aggregate)
  - [Cardboard] (#cardboard)
  - [Hard Plastic] (#hard-plastic)
  - [Metal] (#metal)
  - [Soft Plastic] (#soft-plastic)
  - [Timber] (#timber)

# Introduction
This dataset is designed for waste classification within different material classes. The goal is to accurately identify and annotate different types of waste materials for sorting and recycling purposes. The classes represented are: Aggregate, Cardboard, Hard Plastic, Metal, Soft Plastic, and Timber.

# Object Classes

## Aggregate
### Description
Aggregate refers to small, granular materials, often irregular in shape with rough surfaces. They generally appear as pieces of stone or concrete.

### Instructions
Annotate all visible portions of aggregate items. Ensure to include entire objects even if occluded by other materials, estimating boundaries if necessary. Exclude dust or very fine particles that do not form distinct objects.

## Cardboard
### Description
Cardboard objects are typically flat and have a layered texture. They may appear as boxes or sheets.

### Instructions
Annotate only distinguishable pieces of cardboard, focusing on their flat surfaces and any visible layering. Do not annotate cardboard that is part of another object or soiled beyond recognition.

## Hard Plastic
```

Description

Hard plastics are rigid and maintain their shape. They can be cylindrical, tubular, or robust objects often found in industrial contexts.

Instructions

Annotate the entire visible area of hard plastic objects, ensuring to capture their solid structure. Avoid labeling small, indistinct pieces or any plastic that appears flexible.

Metal

Description

Metal objects are robust, often shiny or reflective. They can appear as rods, sheets, or other distinct shapes.

Instructions

Label all distinct metal parts, taking care to capture their complete form. Avoid labeling rust marks or indistinct metallic fragments lacking shape.

Soft Plastic

Description

Soft plastics are flexible and often transparent or translucent. They may appear in the form of bags or wrappers.

Instructions

Focus on full pieces of soft plastic material, ensuring to include areas with visible creases or folds indicating flexibility. Do not label pieces smaller than a recognizable package or those mixed with other materials.

Timber

Description

Timber objects are wooden, either rough or smooth, often elongated or rectangular.

Instructions

Annotate the entire visible portion of timber, focusing on the grain or wood texture. Do not label splinters or fragments that do not exhibit a clear wooden structure.

C Roboflow100-VL Datasets

We present a table with links to all datasets within Roboflow100-VL (fully-supervised and FSOD datasets) below.

Flora & Fauna	Link
aquarium-combined	FSOD, Fully Supervised
bees	FSOD, Fully Supervised
deepfruits	FSOD, Fully Supervised
exploratorium-daphnia	FSOD, Fully Supervised
grapes-5	FSOD, Fully Supervised
grass-weeds	FSOD, Fully Supervised
gwhd2021	FSOD, Fully Supervised
into-the-vale	FSOD, Fully Supervised
jellyfish	FSOD, Fully Supervised
marine-sharks	FSOD, Fully Supervised
orgharvest	FSOD, Fully Supervised
peixos-fish	FSOD, Fully Supervised
penguin-finder-seg	FSOD, Fully Supervised
pig-detection	FSOD, Fully Supervised
roboflow-trained-dataset	FSOD, Fully Supervised
sea-cucumbers-new-tiles	FSOD, Fully Supervised
thermal-cheetah	FSOD, Fully Supervised
tomatoes-2	FSOD, Fully Supervised
trail-camera	FSOD, Fully Supervised
underwater-objects	FSOD, Fully Supervised
varroa-mites-detection-test-set	FSOD, Fully Supervised
wb-prova	FSOD, Fully Supervised
weeds4	FSOD, Fully Supervised

Industrial	Link
-grccs	FSOD, Fully Supervised
13-lkc01	FSOD, Fully Supervised
2024-frc	FSOD, Fully Supervised
aircraft-turnaround-dataset	FSOD, Fully Supervised
asphalt-distress-detection	FSOD, Fully Supervised
cable-damage	FSOD, Fully Supervised
conveyor-t-shirts	FSOD, Fully Supervised
dataconvert	FSOD, Fully Supervised
deeppcb	FSOD, Fully Supervised
defect-detection	FSOD, Fully Supervised
fruitjes	FSOD, Fully Supervised
infrared-image-of-power-equipment	FSOD, Fully Supervised
ism-band-packet-detection	FSOD, Fully Supervised
110ul502	FSOD, Fully Supervised
needle-base-tip-min-max	FSOD, Fully Supervised
recode-waste	FSOD, Fully Supervised
screw-detect-classification	FSOD, Fully Supervised
smd-components	FSOD, Fully Supervised
truck-movement	FSOD, Fully Supervised
tube	FSOD, Fully Supervised
water-meter	FSOD, Fully Supervised
wheel-defect-detection	FSOD, Fully Supervised

Document	Link
activity-diagrams	FSOD, Fully Supervised
all-elements	FSOD, Fully Supervised
circuit-voltages	FSOD, Fully Supervised
invoice-processing	FSOD, Fully Supervised
label-printing-defect-version-2	FSOD, Fully Supervised
macro-segmentation	FSOD, Fully Supervised
paper-parts	FSOD, Fully Supervised
signatures	FSOD, Fully Supervised
speech-bubbles-detection	FSOD, Fully Supervised
wine-labels	FSOD, Fully Supervised

Medical	Link
canalstenosis	FSOD, Fully Supervised
crystal-clean-brain-tumors-mri-dataset	FSOD, Fully Supervised
dentalai	FSOD, Fully Supervised
inbreast	FSOD, Fully Supervised
liver-disease	FSOD, Fully Supervised
nih-xray	FSOD, Fully Supervised
spinefrxnormalvindr	FSOD, Fully Supervised
stomata-cells	FSOD, Fully Supervised
train	FSOD, Fully Supervised
ufba-425	FSOD, Fully Supervised
urine-analysis1	FSOD, Fully Supervised
x-ray-id	FSOD, Fully Supervised
xray	FSOD, Fully Supervised

Aerial	Link
aerial-airport	FSOD, Fully Supervised
aerial-cows	FSOD, Fully Supervised
aerial-sheep	FSOD, Fully Supervised
apoce-aerial-photographs-for-object-detection-of-construction-equipment	FSOD, Fully Supervised
electric-pylon-detection-in-rsi	FSOD, Fully Supervised
floating-waste	FSOD, Fully Supervised
human-detection-in-floods	FSOD, Fully Supervised
sssod	FSOD, Fully Supervised
uavdet-small	FSOD, Fully Supervised
wildfire-smoke	FSOD, Fully Supervised
zebrasatasturias	FSOD, Fully Supervised

Sports	Link
actions	FSOD, Fully Supervised
aerial-pool	FSOD, Fully Supervised
ball	FSOD, Fully Supervised
bibdetection	FSOD, Fully Supervised
football-player-detection	FSOD, Fully Supervised
lacrosse-object-detection	FSOD, Fully Supervised

Other	Link
buoy-onboarding	FSOD, Fully Supervised
car-logo-detection	FSOD, Fully Supervised
clashroyalechardetector	FSOD, Fully Supervised
cod-mw-warzone	FSOD, Fully Supervised
countingpills	FSOD, Fully Supervised
everdaynew	FSOD, Fully Supervised
fliir-camera-objects	FSOD, Fully Supervised
halo-infinite-angel-videogame	FSOD, Fully Supervised
mahjong	FSOD, Fully Supervised
new-defects-in-wood	FSOD, Fully Supervised
orionproducts	FSOD, Fully Supervised
pill	FSOD, Fully Supervised
soda-bottles	FSOD, Fully Supervised
taco-trash-annotations-in-context	FSOD, Fully Supervised
the-dreidel-project	FSOD, Fully Supervised