

## A Appendix

The Appendix provides a Model Card [27] for OWLv2 as well as additional methodological details, hyperparameters, and results. At the end of the Appendix, we provide qualitative examples of the self-training data and model predictions.

The Appendix is structured as follows:

A.1	Model Card	14
A.2	Human-Curated Label Space	15
A.3	Machine-Generated Label Space	15
A.4	Combined Label Space	16
A.5	Augmentations for Self-Training	16
A.6	Token Dropping	17
A.7	Further Efficiency Improvements	18
A.8	Model Hyperparameters	18
A.9	Additional Results	19
A.9.1	Fixed Average Precision	19
A.9.2	Per-Dataset ODinW Results	19
A.9.3	Fine-Tuning Robustness Trade-Off for OWLv2 L/14	19
A.10	Qualitative Examples	19

## A.1 Model Card

<b>Model Summary</b>	
Model Architecture	OWL v2 is an open-vocabulary object detector based on OWL-ViT [26]. It consists of an image encoder with a Vision Transformer [17] architecture, a text encoder with a similar Transformer architecture, and heads that predict bounding boxes and label scores from provided images and text queries.
Input(s)	An image and a list of free-text object descriptions (queries).
Output(s)	A list of bounding boxes and a score for each box/query pair.
<b>Usage</b>	
Application	The model is intended for open-vocabulary object detection.
Known Caveats	(1) Confidence scores of predictions are not intended to be compared across text queries. While the training loss encourages cross-query calibration for <i>seen</i> queries, scores for <i>unseen</i> queries are not calibrated. Further, the mean Average Precision (mAP) metric does not measure cross-query calibration, so higher mAP does not imply better cross-query calibration. Also see Section 5. (2) Fine-tuning the model creates a trade-off between the performance on fine-tuned texts and unseen texts. See Section 4.6 for details.
<b>System Type</b>	
System Description	This is a standalone model.
Upstr. Dependencies	None.
Downstr. Dependencies	None.
<b>Implementation Frameworks</b>	
Hardware & Software	Hardware: TPU [13] v2 or v3 (for B- and L-sized models) or v4 (for G-sized models). Software: JAX [3], Flax [11], Scenic [7].
Compute Requirements	Reported in Section 4.5.
<b>Model Characteristics</b>	
Model Initialization	The model is initialized from pre-trained language CLIP [30] or SigLIP [43] checkpoints.
Model Status	This is a static model trained on an offline dataset.
Model Stats	The largest OWLv2 model has 2.3B parameters, of which 2B are used for the image encoder and 300M for the text encoder (the heads have a negligible number of parameters). We also trained models with 430M and 150M parameters.
<b>Data Overview</b>	
Training dataset	The model is self-trained on bounding boxes predicted by the original OWL-ViT L/14 model [26] on the WebLI dataset [4]. Details on the annotation procedure are provided in Section 3.1.
Evaluation & Fine-tuning Dataset	Open-vocabulary object detection performance is evaluated using the LVIS [10] and ODinW13 [21] datasets. As indicated in Table 1, some models are fine-tuned on the “base” annotations of LVIS, i.e. only annotations for “frequent” and “common” object categories as defined in the official annotations [10]. None of our models have seen any human annotations for LVIS “rare” categories, such that LVIS mAP <sub>rare</sub> measures zero-shot performance.
<b>Evaluation Results</b>	
Evaluation Results	Reported in Table 1.
<b>Model Usage &amp; Limitations</b>	
Sensitive Use	The model detects objects matching free-text descriptions. This capability should not be used for unethical use cases such as surveillance.
Known Limitations	Reported in Section 5.
Ethical Considerations	Reported in Section 5.

## A.2 Human-Curated Label Space

The human-curated label space was obtained by merging common dataset class lists with the Python code below.

```
1 # Dataset class names, as available e.g. from TensorFlow Datasets.
2 # For Visual Genome, we used the 1600 most common label strings.
3 LVIS_CLASS_NAMES = [...]
4 OBJECTS365_CLASS_NAMES = [...]
5 OPEN_IMAGES_V4_BOXABLE_CLASS_NAMES = [...]
6 VISUAL_GENOME_CLASS_NAMES = [...]
7
8 queries = (
9     LVIS_CLASS_NAMES
10    + OBJECTS365_CLASS_NAMES
11    + OPEN_IMAGES_V4_BOXABLE_CLASS_NAMES
12 )
13
14 # Remove duplicates:
15 queries = set([q.lower() for q in queries])
16
17 # Remove plural forms:
18 remove = set()
19 for singular in queries:
20     plurals = [singular + 's', singular + 'es']
21     for plural in plurals:
22         if plural in queries:
23             remove.add(plural)
24
25 # Same queries for all images:
26 queries = list(queries.difference(remove))
```

## A.3 Machine-Generated Label Space

The machine-generated label space was obtained from the image-associated text, for each image separately, using the Python code below. Figure A3 shows example pseudo-annotations using the N-gram label space.

```
1 from typing import Iterable, List
2 import nltk
3
4 # Stopwords from nltk.corpus.stopwords.words('english'):
5 STOPWORDS_EN = frozenset({
6     'a', 'about', 'above', 'after', 'again', 'against', 'all', 'am', 'an',
7     'and', 'any', 'are', 'as', 'at', 'be', 'because', 'been', 'before', 'being',
8     'below', 'between', 'both', 'but', 'by', 'can', 'did', 'do', 'does',
9     'doing', 'don', 'down', 'during', 'each', 'few', 'for', 'from', 'further',
10    'had', 'has', 'have', 'having', 'he', 'her', 'here', 'hers', 'herself',
11    'him', 'himself', 'his', 'how', 'i', 'if', 'in', 'into', 'is', 'it', 'its',
12    'itself', 'just', 'me', 'more', 'most', 'my', 'myself', 'no', 'nor', 'not',
13    'now', 'of', 'off', 'on', 'once', 'only', 'or', 'other', 'our', 'ours',
14    'ourselves', 'out', 'over', 'own', 's', 'same', 'she', 'should', 'so',
15    'some', 'such', 't', 'than', 'that', 'the', 'their', 'theirs', 'them',
16    'themselves', 'then', 'there', 'these', 'they', 'this', 'those', 'through',
17    'to', 'too', 'under', 'until', 'up', 'very', 'was', 'we', 'were', 'what',
18    'when', 'where', 'which', 'while', 'who', 'whom', 'why', 'will', 'with',
19    'you', 'your', 'yours', 'yourself', 'yourselves'
20 })
21
22 # These words were found by manually going through the most common 1000 words
23 # in a sample of alt-texts and selecting generic words without specific meaning:
24 COMMON_GENERIC_WORDS = frozenset({
25     'alibaba', 'aliexpress', 'amazon', 'available', 'background', 'blog', 'buy',
```

```

26     'co', 'com', 'description', 'diy', 'download', 'facebook', 'free', 'gif',
27     'hd', 'ideas', 'illustration', 'illustrations', 'image', 'images', 'img',
28     'instagram', 'jpg', 'online', 'org', 'original', 'page', 'pdf', 'photo',
29     'photography', 'photos', 'picclick', 'picture', 'pictures', 'png', 'porn',
30     'premium', 'resolution', 'royalty', 'sale', 'sex', 'shutterstock', 'stock',
31     'svg', 'thumbnail', 'tumblr', 'tumblr', 'tumblr', 'twitter', 'uk', 'uploaded', 'vector',
32     'vectors', 'video', 'videos', 'wallpaper', 'wallpapers', 'wholesale', 'www',
33     'xxx', 'youtube'
34 })
35
36 def _is_all_stopwords(ngram: Iterable[str]) -> bool:
37     return set(ngram).issubset(STOPWORDS_EN)
38
39
40 def _get_ngrams(
41     caption: str, max_num_queries: int, max_ngram_len: int
42 ) -> List[str]:
43     """Returns image caption ngrams as queries."""
44
45     # Make lower-case:
46     caption = caption.lower()
47
48     # Remove common generic words:
49     words = [w for w in caption.split() if w not in COMMON_GENERIC_WORDS]
50
51     queries = []
52     for ngram in nltk.everygrams(words, max_len=max_ngram_len):
53         # Don't use ngram if it only consists of stop words:
54         if _is_all_stopwords(ngram):
55             continue
56         queries.append(' '.join(ngram))
57         if len(queries) == max_num_queries:
58             break
59     return queries
60
61 # Example command to get queries for one image:
62 queries = _get_ngrams(caption, max_num_queries=300, max_ngram_len=10)

```

#### A.4 Combined Label Space

When merging pseudo-annotations obtained with human-curated and machine-generated queries, it is important to consider that human-curated queries tend to be closer to the training distribution of the annotator and therefore tend to have higher scores than pseudo-annotations based on machine-generated queries. Simply merging annotations from the two label spaces and filtering them with the same confidence threshold would therefore retain primarily annotations based on human-curated queries. To achieve a more even balance when using the combined label space (“N-grm+curated” in Table 1), we therefore re-scaled scores of pseudo-annotations obtained with the human-curated queries by a factor of 0.3 before applying the same confidence threshold to all (human-curated and machine-generated) annotations.

#### A.5 Augmentations for Self-Training

Since Web-scale image-text data differs in important aspects from human-curated detection datasets, we depart from the augmentation strategy of [26] in several ways. As described in Section 3.2, since Web images tend to be smaller and show fewer objects than e.g. LVIS images, we use stronger image mosaics with up to  $6 \times 6$  tiles (Figure A1). For the same reason, we additionally randomly resize each raw image such that its width is between  $0.5 \times$  and  $1.0 \times$  the width of the full mosaic tile, padding on the bottom and right to preserve the aspect ratio (Figure A4).

On the other hand, given the large size of our dataset, some other augmentations can be avoided: We do not use left/right flipping or random cropping during self-training. We also do not add random

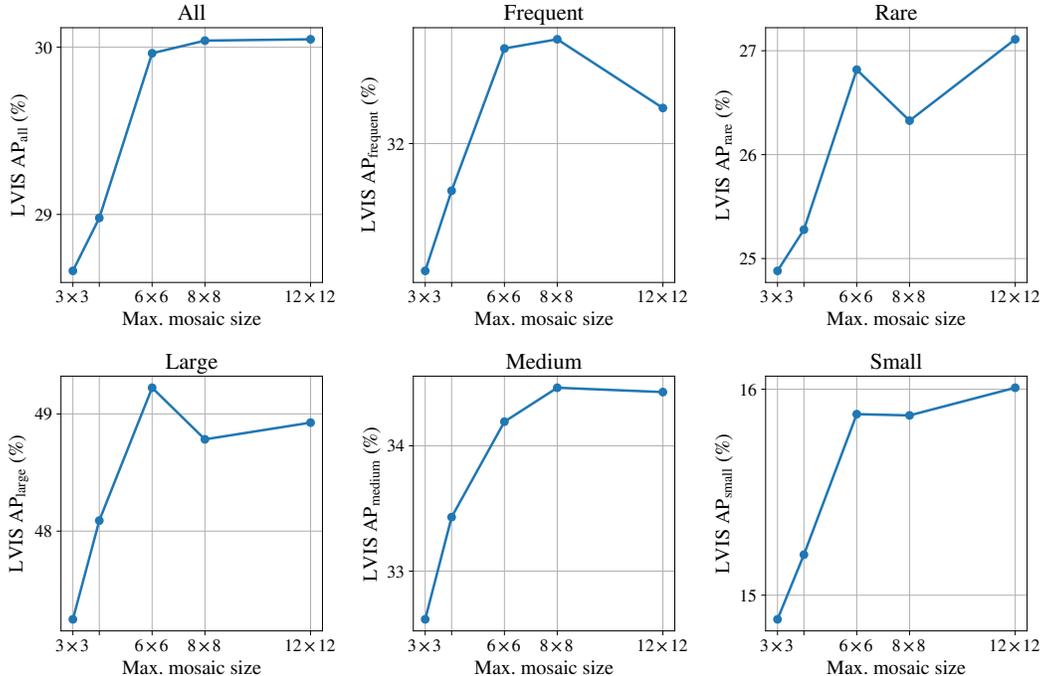


Figure A1: Sweep over mosaic sizes. OWL-ViT B/16 models were trained on pseudo-box annotations (“ngrams” label space) for 100’000 steps with different mosaic sizes. At a given “max. mosaic size”, the model is trained on equal proportions of mosaics up to that size. For example, for max. size =  $12 \times 12$ , the model receives images with  $1, 2^2, 3^2, 4^2, 6^2, 8^2$ , or  $12^2$  tiles, respectively (only sizes with prime factors 1, 2, and 3 are supported). For this figure, the model input resolution was  $768 \times 768$ . Mosaic sizes up to  $12 \times 12$  improve overall performance (mAP<sub>all</sub>) and especially “rare” and “small” object performance. The benefit may be due to seeing smaller objects on average, or due to seeing more WebLI images per training step (a  $12 \times 12$  mosaic contains 144 WebLI images).

prompt templates to the pseudo-labels during self-training. During fine-tuning, we use the same augmentations as [26].

## A.6 Token Dropping

To improve training efficiency, we drop image patches based on their pixel variance (Section 3.2). Table A2 shows how the performance of a standard OWL-ViT model varies for different amounts of token dropping. Dropping up to 50% of tokens is within one standard deviation of the full performance. We therefore drop 50% of tokens during all of our experiments.

Table A2: Performance of standard OWL-ViT (L/14), trained on Objects365 and Visual Genome as in [26], for different token drop rates. For drop rate 0.0, the standard deviation over three runs is given.

Metric	Token drop rate				
	0.00	0.25	0.33	0.50	0.70
LVIS AP <sub>all</sub> <sup>val</sup>	33.3 ± 0.33	33.1	33.6	32.9	30.4
LVIS AP <sub>rare</sub> <sup>val</sup>	31.8 ± 1.16	31.0	32.6	30.8	28.2

To inject some stochasticity to the patch selection, we add a small amount of noise to the image before computing patch variance (uniformly distributed between 0.0 and 0.01 for images in the range [0.0, 1.0]). Figure A4 shows an example training image before and after token dropping.

Table A3: Hyperparameters of the models shown in Table 1. Only parameters that vary between models are shown; constant parameters are described in the text (Appendix A.8). For *Dropout rate* and *Droplayer rate*, the first number indicates the value used for the image encoder and the second for the text encoder. *Examples seen* includes both self-training and fine-tuning.

Method	Backbone	Image size	Learning rate	Dropout rate	Droplayer rate	Instance top $k$	Batch size (ST)	Batch size (FT)	Examples seen	
<b>Open vocabulary:</b>										
11	OWL-ST	CLIP B/16	960	$5 \times 10^{-5}$	.0/.0	.2/.1	256	256	–	$3.7 \times 10^8$
12	OWL-ST	CLIP L/14	1008	$2 \times 10^{-5}$	.0/.0	.2/.1	512	256	–	$2.3 \times 10^8$
13	OWL-ST	SigLIP G/14	1008	$2 \times 10^{-5}$	.0/.1	.2/.4	512	128	–	$1.6 \times 10^8$
14	OWL-ST+FT	CLIP B/16	960	$5 \times 10^{-5}$	.0/.0	.2/.1	256	256	256	$3.6 \times 10^8$
15	OWL-ST+FT	CLIP L/14	1008	$2 \times 10^{-5}$	.0/.0	.2/.1	512	256	128	$2.3 \times 10^8$
16	OWL-ST+FT	SigLIP G/14	1008	$2 \times 10^{-5}$	.0/.1	.2/.4	512	128	128	$1.6 \times 10^8$
<b>Human-curated vocabulary:</b>										
20	OWL-ST+FT	CLIP B/16	960	$5 \times 10^{-5}$	.0/.0	.2/.1	256	256	256	$8.2 \times 10^8$
21	OWL-ST+FT	CLIP L/14	1008	$2 \times 10^{-5}$	.0/.0	.2/.1	512	256	128	$3.6 \times 10^8$

## A.7 Further Efficiency Improvements

To further improve training efficiency beyond the methods described in Section 3.2, we also adopt previously proposed methods for large-scale Transformer training: To save memory, we use a variant [42] of the Adafactor optimizer [34] instead of Adam [15]. To avoid having to choose and optimize the total training duration ahead of time, we use the open-ended inverse square-root schedule [36, 42] with a fixed time-scale of 10’000 steps for all experiments and linearly “cool down” checkpoints along the way for evaluation (see Section 3.3).

## A.8 Model Hyperparameters

We use the following hyperparameters for all of our models. Hyperparameters that vary between models are listed in Table A3.

- Optimizer: Adafactor variant as in [42]
- Learning rate schedule: Inverse square-root [36] with timescale 10’000 steps
- Learning rate for the text encoder:  $2 \times 10^{-6}$
- Token dropping rate during training: 0.5
- Pseudo-annotation confidence score threshold: 0.3 (except for Figure 3)
- Augmentations: See Appendix A.5
- All remaining hyperparameters are as in [26].

**Hyperparameter selection.** Most hyperparameters were either taken directly from [26] or technically constrained, e.g. we chose the largest batch size that fit into the memory of the available accelerators. Where hyperparameters were tuned, we ran short B/16-scale trial experiments and selected the parameters with the highest LVIS  $\text{mAP}_{\text{rare}}$  for our main runs.

**SigLIP G/14.** For the G/14 model, we started self-training with a learning rate of  $5 \times 10^{-5}$ , a droplayer rate of .1/.0, and no dropout. We found that the model overfit during fine-tuning with these settings, and switched to a learning rate of  $2 \times 10^{-5}$ , a droplayer rate of .2/.4, and a dropout rate of .0/.1 after 740’000 self-training steps. To save resources, we did not start training from the beginning. With the new settings, we observed no overfitting during fine-tuning, but it is possible that these settings are still not optimal.

Table A4: Open-vocabulary detection results on LVIS using the “fixed” AP metric [5]. Fixed AP is implemented as proposed in [5] by evaluating AP on the top 10’000 predictions per class over the entire validation set.

Method	Backbone	AP <sup>mini</sup> <sub>all</sub>		AP <sup>mini</sup> <sub>rare</sub>		AP <sup>val</sup> <sub>all</sub>		AP <sup>val</sup> <sub>rare</sub>		
		old	fixed	old	fixed	old	fixed	old	fixed	
<i>Open vocabulary:</i>										
1	RegionCLIP [46]	R50x4	–	–	–	–	32.3	–	22.0	–
2	OWL [26]	CLIP B/16	–	–	–	–	27.2	–	20.6	–
3	OWL [26]	CLIP L/14	–	–	–	–	34.6	–	31.2	–
4	GLIPv2 [45]	Swin-T	29.0	–	–	–	–	–	–	–
5	GLIPv2 [45]	Swin-B	48.5	–	–	–	–	–	–	–
6	GLIPv2 [45]	Swin-H	50.1	–	–	–	–	–	–	–
7	F-VLM [19]	R50x4	–	–	–	–	28.5	–	26.3	–
8	F-VLM [19]	R50x64	–	–	–	–	34.9	–	32.8	–
9	3Ways [1]	NFNet-F0	–	–	–	–	35.7	–	25.6	–
10	3Ways [1]	NFNet-F6	–	–	–	–	44.6	–	30.1	–
11	OWL-ST	CLIP B/16	31.8	34.4	35.4	38.3	27.0	28.6	29.6	30.3
12	OWL-ST	CLIP L/14	38.1	40.9	39.0	41.5	33.5	35.2	34.9	36.2
13	OWL-ST	SigLIP G/14	37.8	–	40.9	–	33.7	–	37.5	–
14	OWL-ST+FT	CLIP B/16	47.2	48.7	37.8	42.1	41.8	43.2	36.2	39.0
15	OWL-ST+FT	CLIP L/14	54.1	56.2	46.1	52.3	49.4	51.1	44.6	47.4
16	OWL-ST+FT	SigLIP G/14	51.3	–	50.9	–	47.0	–	47.2	–
<i>Human-curated vocabulary:</i>										
17	Detic [47]	R50	–	–	–	–	32.4	–	24.6	–
18	DetCLIPv2 [38]	Swin-T	–	40.4	–	36.0	–	32.8	–	31.0
19	DetCLIPv2 [38]	Swin-L	–	44.7	–	43.1	–	36.6	–	33.3
20	OWL-ST+FT	CLIP B/16	51.1	52.3	41.9	46.5	45.6	46.7	40.5	42.5
21	OWL-ST+FT	CLIP L/14	55.8	57.2	50.0	54.5	50.4	52.0	45.9	48.5

## A.9 Additional Results

### A.9.1 Fixed Average Precision

In the standard Average Precision metric (AP<sup>old</sup>), performance on one class depends on the performance on other classes. This dependence makes the metric “gameable” by re-scaling the scores of certain classes [5]. To avoid this issue, some prior work reports a “fixed” version of AP proposed in [5]. In Table 1, we report AP<sup>old</sup> for our models. For models from the literature, we report whichever AP version is available. Since AP<sup>fixed</sup> tends to produce higher values than AP<sup>old</sup>, Table 1 tends to underestimate the advantage of our method over prior work using AP<sup>fixed</sup>. We provide AP<sup>fixed</sup> for all of our models in Table A4. As proposed in [5], we implement AP<sup>fixed</sup> by evaluating AP on the top 10’000 predictions per class over the entire validation set. This ensures that classes do not compete with each other for inclusion in the evaluated predictions.

### A.9.2 Per-Dataset ODinW Results

Table A5 shows un-aggregated results on all 35 ODinW datasets for our main models. In addition, in the last row, we provide results for a weight-space ensemble of a self-trained and fine-tuned OWLv2 L/14 model (the same model is shown in Figure A2).

### A.9.3 Fine-Tuning Robustness Trade-Off for OWLv2 L/14

In Figure A2, we provide the same analysis of the robustness trade-off after fine-tuning for an L/14 model that we provided for a B/16 model in Figure 5.

## A.10 Qualitative Examples

In Figures A5 to A7, we provide qualitative examples of detection predictions from OWLv2 L/14 models. In each figure, the top image shows predictions obtained directly after self-training, and



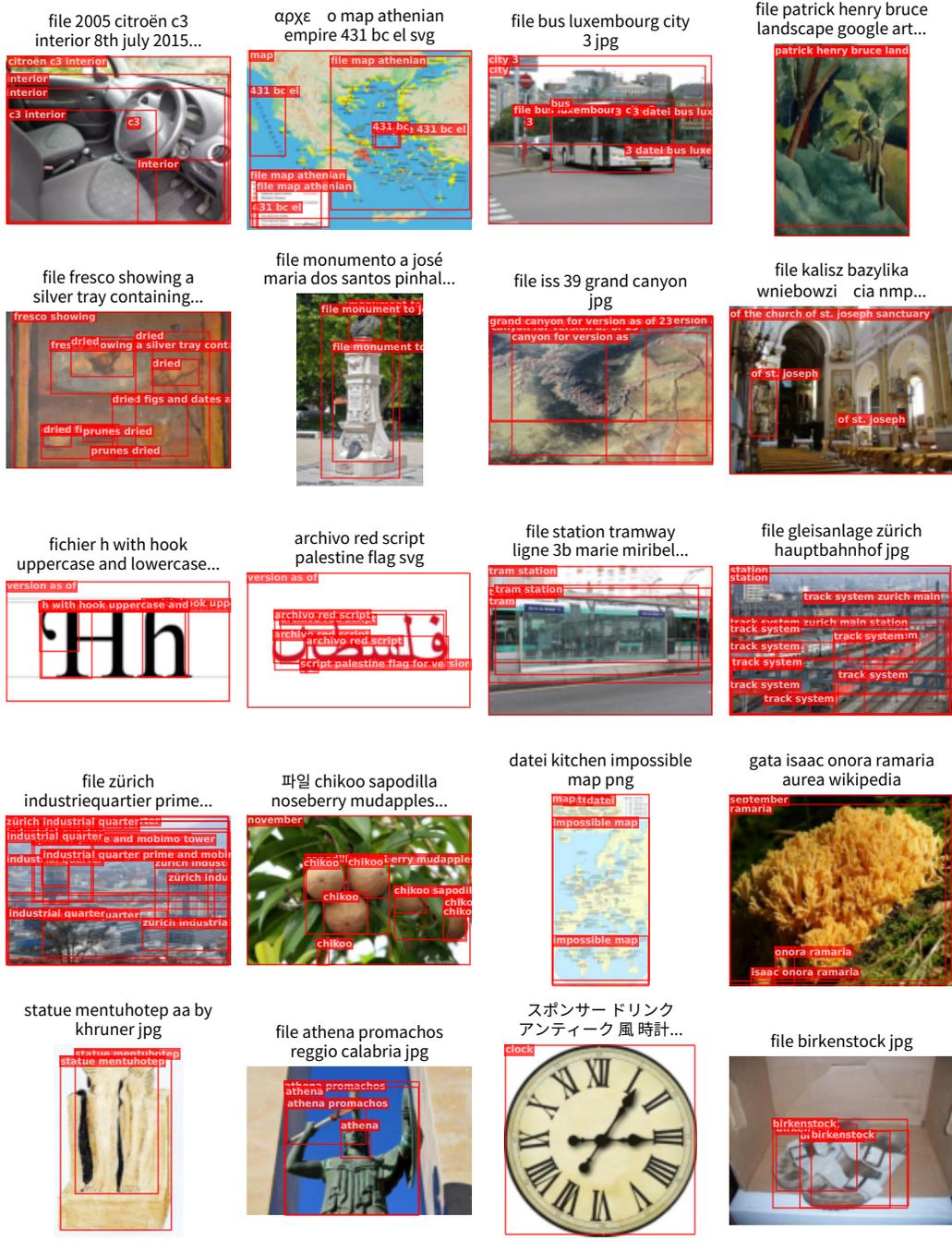


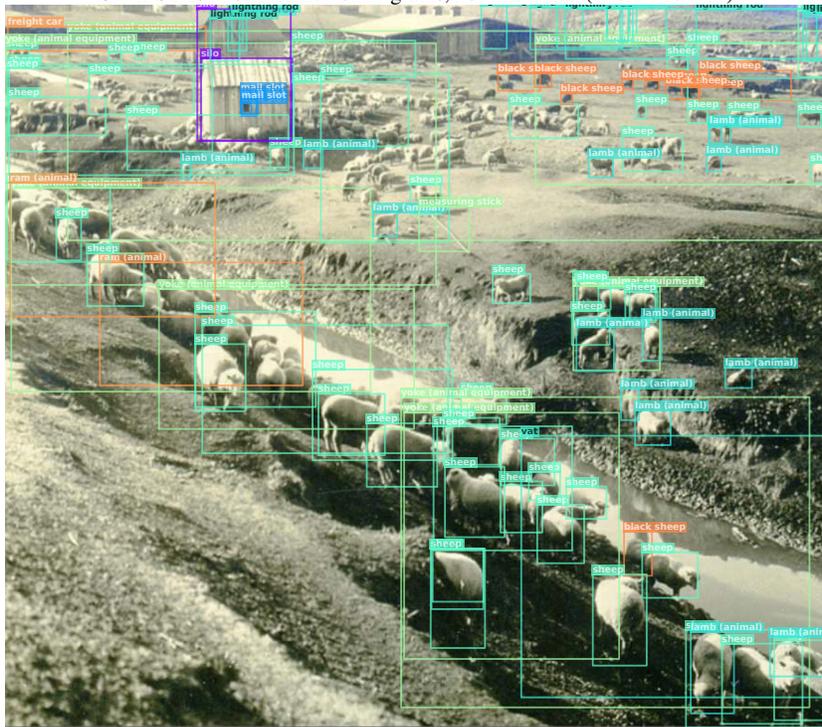
Figure A3: Example pseudo-annotations on WeBLI [4]. Image-associated text (from the HTML alt\_text tag) is shown above the images. If the text is not in English, an automatically generated translation is used. N-grams are extracted from these texts to generate queries for the annotator model. Pseudo-annotations were filtered as for our main experiments: To be included, boxes must have a score of at least 0.1, and images must have at least one box with a score above 0.3. All images from Wikimedia Commons.



Figure A4: Training inputs after pre-processing. **Top:** A  $4 \times 4$  mosaic of randomly resized and padded images as used for self-training. **Bottom:** The same mosaic after dropping the 50% of patches with lowest pixel variance (image size:  $1008 \times 1008$ ; patch size:  $14 \times 14$ ). Most dropped patches belong to padding areas or uniform image backgrounds. All images from Wikimedia Commons.



OWL-ST L/14 self-trained on N-grams, not fine-tuned (Table 1 row 12)



OWL-ST+FT L/14 self-trained on N-grams and fine-tuned on LVIS<sub>base</sub> (Table 1 row 15)



Figure A6: Qualitative example for OWLv2 L/14 from the LVIS val set. For the visualization, all LVIS classes were used as prompts. LVIS<sub>rare</sub> classes are labeled in black. **Top:** OWL-ST self-trained on N-grams, not fine-tuned (Table 1 row 12). **Bottom:** OWL-ST+FT self-trained on N-grams and fine-tuned on LVIS<sub>base</sub> (Table 1 row 15). Boxes above score 0.08 (top) or 0.3 (bottom) are shown.

OWL-ST L/14 self-trained on N-grams, not fine-tuned (Table 1 row 12)



OWL-ST+FT L/14 self-trained on N-grams and fine-tuned on LVIS<sub>base</sub> (Table 1 row 15)



Figure A7: Qualitative example for OWLv2 L/14 from the LVIS val set. For the visualization, all LVIS classes were used as prompts. LVIS<sub>rare</sub> classes are labeled in black. **Top:** OWL-ST self-trained on N-grams, not fine-tuned (Table 1 row 12). **Bottom:** OWL-ST+FT self-trained on N-grams and fine-tuned on LVIS<sub>base</sub> (Table 1 row 15). Boxes above score 0.08 (top) or 0.3 (bottom) are shown.