# Appendix for DetCLIPv3: Towards Versatile Generative Open-vocabulary Object Detection

## A. Additional Implementation Details

**Training.** The training of DetCLIPv3 involves data from various sources. Table 1 summarizes detailed data information used in different training phases. Since the training process varies for different data type, (*e.g.*, the object captioner accepts only image-text pair data as input), we design each iteration's global batch to contain only one type of data.

For the training of the the open-vocabulary detector, following previous DetCLIP works [21, 23], we initialize the text encoder with the parameters of FILIP's [22] language model, and reduce its learning rate by 0.1 during training to preserve the knowledge obtained via FILIP's pre-training. To improve training efficiency, we set the maximum text token length for the text encoder to 16.

For the training of the object captioner, we initialize the captioner with the pre-trained weights of Qformer [8], whereas the deformable [26] cross-attention layers are randomly initialized. To preserve the knowledge acquired during Qformer [8] pretraining, the object captioner utilizes the same BERT [5] tokenizer for processing text input, different to the text encoder which employs the CLIP [15] tokenizer. The maximum text token length for the object captioner is set to 32.

In each training stage, to conserve GPU memory, automatic mixed-precision [13] and gradient checkpointing [2] are employed. Table 2 summarizes the detailed training settings for each training stage.

| Training Stage | Dataset | Total Volume |
|---|---|---|
| Stage 1 | O365 [17](0.66M), GoldG [9](0.77M), V3Det [20](0.18M) | 1.61M |
| Stage 2 | GranuCap50M | 50M |
| Stage 3 | O365 [17](0.66M), GoldG [9](0.77M), V3Det [20](0.18M), GranuCap600K (0.6M) | 2.21M |

Table 1. Dataset information for each training stage of Det-CLIPv3. O365 refers to the Objects365 v2, from which we sample 0.66M data with balanced class for training, similar to previous DetCLIP v1/v2 [21, 23] works. GranuCap50M is developed with our proposed auto-annotation pipeline, using 50M image-text pairs sampled from a collection of CC3M [18], CC12M [1], YFCC100M [19] and LAION400M [16].

**Inference.** Inference process of DetCLIPv3's OV detector follows DINO [25], where the results for each image are derived from the predictions of 300 object queries with highest confidence scores. For the *fixed* AP [4] evaluation on the LVIS [6] dataset, it is required that each category within the entire validation set has at least 10,000 predictions. To ensure an sufficient number of predictions per image, we adopt an inference process similar to that of GLIP [9]. Specifically, during inference for each data sample, the 1203 categories are split into 31 chunks, with a chunk size of 40 categories. We conduct inference separately for each chunk and retain the top 300 predictions based on their confidence scores.

For the inference process of DetCLIPv3's object captioner, as described in the main paper, for each image, we utilize the most frequent 15k concepts from our developed noun concept corpus as text queries to extract top 100 foreground regions with highest similarity. After the generation of descriptive labels for these regions by the object captioner, their confidence scores are re-calibrated using the OV detector. A class-agnostic non-maximum suppression (NMS) operation is then performed for regions with re-calibrated scores higher than 0.05, the results of which are output as predictions. We set beam search's beam size equal to 1 for inference of object captioner.

**Finetuning.** We fine-tune DetCLIPv3 on 2 datasets, *i.e.*, LVIS [6] and ODinW13 [9]. Table 3 and 4 summarize the detailed fine-tuning settings for LVIS and ODinW13, respectively. For LVIS, when fine-tuning with base categories, we exclude novel categories while sampling negative concepts. For ODinW13, similar to DetCLIPv2[23], we employ an auto-decay learning rate schedule. Specifically, when performance reaches a plateau and persists for a tolerance period $t_1$, we reduce the learning rate by a factor of 0.1. If there is no improvement in performance for a tolerance period $t_2$, we then terminate the training process.

## B. Additional Data Pipeline Details

Figure 1 illustrates an overview of DetCLIPv3's auto-annotation data pipeline.

**Prompts.** Here we provide the prompts used in each step, including those for the VLLMs as well as for GPT-4.
1. **Recaptioning with VLLM**: We employ Instruct-

| Config | Stage1 | Stage2 | Stage3 |
|---|---|---|---|
| GPUs (V100) | 32 (Swin-T)/64 (Swin-L) | | |
| training module | OV detector | object captioner | all modules |
| training objective | $\mathcal{L}_{det} = \mathcal{L}_{align} + \mathcal{L}_{box} + \mathcal{L}_{iou}$ | $\mathcal{L}_{lm} = \mathcal{L}_{lm}^{obj} + \mathcal{L}_{lm}^{img}$ | $\mathcal{L} = \mathcal{L}_{det} + \mathcal{L}_{lm}$ |
| training epochs | 12 | 3 | 5 |
| input resolution | $320 \times 240 \sim 1333 \times 800$ | $320 \times 320$ | $1333 \times 600 \sim 1333 \times 800$ |
| batch size | 128 | 2048 | 128 |
| learning rate | 2.8e-4 | 1e-4 | 1e-4 |
| text encoder lr reduce factor | 0.1 | – | 0.1 |
| numper of concepts (grounding/image-text pairs) | 4800 | – | 4800 |
| optimizer | AdamW [11] | | |
| optimizer momentum | $\beta_1 = 0.9, \beta_2 = 0.999$ | | |
| weight decay | 0.05 | | |
| warmup iters | 1000 | | |
| learning rate schedule | cosine annealing | | |
| text token length (text encoder) | 16 | | |
| text tokenizer (text encoder) | CLIP [15] Tokenizer | | |
| text token length (object captioner) | 32 | | |
| text tokenizer (object captioner) | BERT [5] Tokenizer | | |

Table 2. Detailed pre-training settings of DetCLIPv3. $\mathcal{L}_{lm}^{obj}$ and $\mathcal{L}_{lm}^{img}$ represent for language modeling training objective for object-level captioning and image-level captioning, respectively.
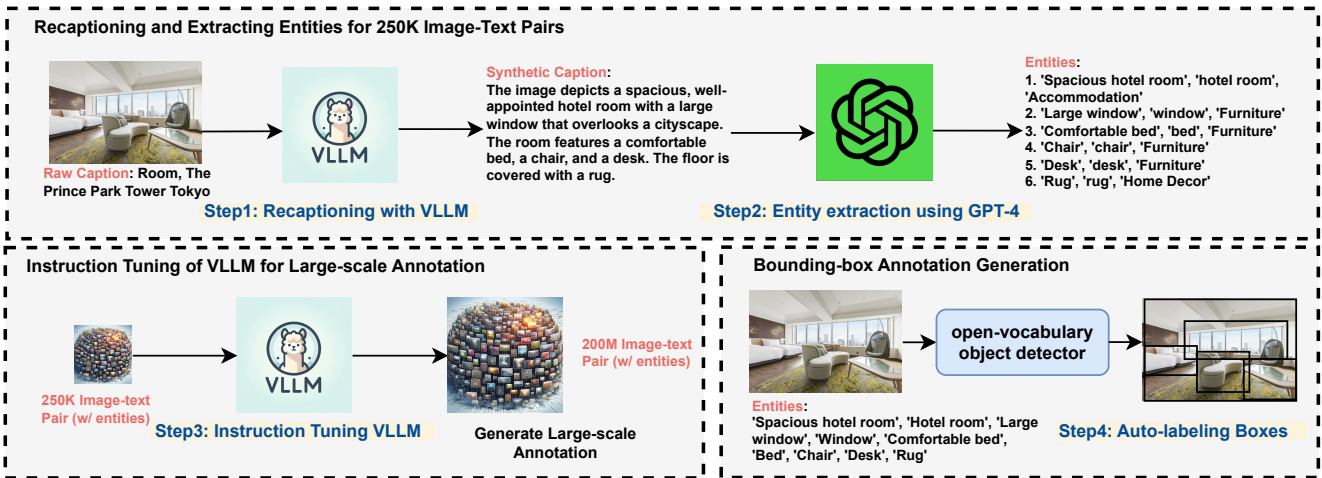


Figure 1. The illustration for DetCLIPv3's auto-annotation data pipeline. It initially utilizes a VLLM to recaption 240K image-text pairs, followed by the use of GPT-4 to extract object entities, formatted as {*phrase, category, parent category*}. Subsequently, these image-text pairs are used for instruction tuning of a VLLM, with the trained model providing annotations for a larger-scale of 200M image-text pairs. Finally, an OV detector is employed to provide pseudo-label bounding boxes for these data, and after applying a confidence score filtering, 50M data are sampled to form GranuCap50M.

BLIP [3] to recaption 240K image-text pairs. To leverage information from the original caption texts, we use the following prompt:

*"Given a noisy caption of the image: {raw caption}, write a detailed clean description of the image."*

2. **Entity extraction using GPT-4**: In this step, we first utilize GPT-4 to filter out non-entity descriptions from the captions generated by the VLLM. The prompt used is:

*"Here is a caption for an image: {caption}. Extract the part of factual description related to what*

*is directly observable in the image, while filtering out the parts that refer to inferred contents, description of atmosphere/appearance/style and introduction of history/culture/brand etc. Return solely the result without any other contents. If you think there is no factual description, just return 'None'."*

Subsequently, we extract information about object entities from the filtered captions using the prompt:

*"You are an AI tasked with developing an open-set object detection dataset from a large number of image captions, without access to the actual images. Your mission*

| Config | Value |
|---|---|
| GPUs (V100) | 16 |
| training epochs | 16 |
| optimizer | AdamW [11] |
| optimizer momentum | $\beta_1 = 0.9, \beta_2 = 0.999$ |
| lr for image encoder | 2e-5 |
| lr for text encoder | 2e-6 |
| weight decay | 0.05 |
| warmup iters | 1000 |
| learning rate schedule | cosine decay |
| batch size | 64 |
| input resolution | $1333 \times 800$ |
| number of concepts per sample | 150 |
| augmentation | multi-scale training, random flip |

Table 3. Detailed fine-tuning settings for LVIS [6].

| Config | Value |
|---|---|
| GPUs (V100) | 8 |
| maximum training epochs | 250 |
| optimizer | AdamW [11] |
| optimizer momentum | $\beta_1 = 0.9, \beta_2 = 0.999$ |
| lr for image encoder | 4e-5 |
| lr for text encoder | 4e-7 |
| weight decay | 0.05 |
| warmup iters | 500 |
| learning rate schedule | auto-step decay |
| lr decay tolerance $t_1$ (epochs) | 10 |
| training terminate tolerance $t_2$ (epochs) | 15 |
| minimum lr to stop decay | 1e-8 |
| batch size | 32 |
| input resolution | $1333 \times 800$ |
| augmentation | multi-scale training, random flip |

Table 4. Detailed fine-tuning settings for ODinW13 [9].

*is to accurately identify and extract 'objects' from these captions, following the principles below:*
*1. 'Objects' are physically tangible: They must be concrete entities that can be visually represented in an image. They are NOT (1) abstract concepts (like 'history', 'culture') or feelings (like 'sorrow', 'happiness'), (2) meta-references to the image itself (e.g., 'image', 'picture', 'photo') or the camera (e.g. something is facing the 'camera'), unless they are specifically referring to physical elements within the image. (3) any descriptors (like 'appearance', 'atmosphere', 'color'), (4) events/activities and processes (like 'game', 'presentation', 'performance') and specific event types (like 'country style wedding', 'film festival'), (5) compositional aspects (like 'perspective', 'focus', 'composition') or viewpoint/perspective (like 'bird's eye view').*
*2. 'Objects' are visually distinct: They are standalone entities that can be visually isolated from their environ-*

*ment. They do not include environmental characteristics (like 'colorful environment') and general location/scene descriptors (e.g., 'scene set indoors', 'country setting', 'sunny day', 'black and white illustration')*
*Adhere to these guidelines for the extraction process:*
*1. Consolidate duplicates: If multiple extracted 'objects' refer to the same entity in the caption, merge them into one while retaining conceptual diversity.*
*2. Categorize the descriptive variants: For 'objects' described with adjectives, provide both versions - with and without the adjective.*
*3. Identify the broader category: Assign a 'parent category' that each 'object' belongs to.*
*Present your results as a numbered list in this format: id. 'object with adjective', 'object without adjective', 'parent category'. Your response should consist exclusively of results, with no superfluous content.*
*Here's the caption: {caption}"*

3. **Instruction tuning of VLLM for large-scale annotation**: In this phase, we use the caption texts and object entity information obtained from the above steps to fine-tune the LLaVA [10] model. Here, we combine the aforementioned information into a new concise prompt, and the question-answer pair is constructed as:

***Question**:*
*"From the noisy caption of the image: {raw caption}, generate a refined image description and identify all visible 'objects' – any visually and physically identifiable entity in the image. Keep the following guidelines in mind:*
*1. Merge similar 'objects' from the caption, preserving conceptual diversity.*
*2. For adjective-described 'objects', provide versions both with and without the adjective.*
*3. Assign a 'parent category' for each 'object'.*
*Present results as:*
*Caption: {caption}*
*Objects: {id. 'object with adjective', 'object without adjective', 'parent category'}.*
*<image tokens>"*
***Answer**:*
*Caption: {refined caption}*
*Objects: {entity information}*
Here, the VLLM receives image tokens, *i.e.*, <image tokens>, along with their original captions, *i.e.*, {raw caption}, as inputs, and learns to generate refined captions and extract information about object entities.

**Visualizations.** Figure 2-a and 2-b depict refined caption and extracted entity information obtained via our proposed data pipeline. Additionally, Figure 3 displays the bounding box pseudo-labels generated by our Swin-L-based model after stage-1 training.

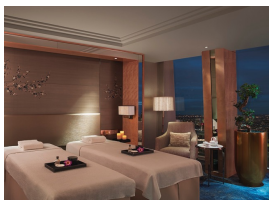| | Recaption text | Extracted entities |
|---|---|---|
|  | A man and a woman are sitting at a casino table, surrounded by chips and cards. The man is wearing a vest, while the woman is dressed in a black dress. They appear to be engrossed in a game of roulette. | 1. 'Man in a vest' \| 'Man' \| 'People'<br>2. 'Woman in a black dress' \| 'Woman' \| 'People'<br>3. 'Casino table' \| 'Table' \| 'Furniture'<br>4. 'Chips' \| 'Chips' \| 'Gaming Equipment'<br>5. 'Cards' \| 'Cards' \| 'Gaming Equipment'<br>6. 'Vest' \| 'Vest' \| 'Clothing'<br>7. 'Black dress' \| 'Dress' \| 'Clothing'<br>8. 'Game of roulette' \| 'Roulette' \| 'Gaming Equipment' |
|  | The image depicts a beautiful outdoor setting, featuring a wooden table and chairs placed in a garden area surrounded by lush greenery. The table is adorned with candles and lanterns. A walkway leads to the table. | 1. 'Beautiful outdoor setting' \| 'outdoor setting' \| 'Locations'<br>2. 'Wooden table' \| 'table' \| 'Furniture'<br>3. 'Chairs' \| 'chairs' \| 'Furniture'<br>4. 'Garden area' \| 'garden' \| 'Locations'<br>5. 'Lush greenery' \| 'greenery' \| 'Nature'<br>6. 'Candles' \| 'candles' \| 'Decorations'<br>7. 'Lanterns' \| 'lanterns' \| 'Decorations'<br>8. 'Walkway' \| 'walkway' \| 'Locations' |
|  | The image features a set of cute sea animals, including sharks, fish, piranha, octopus, manatee, whale, dolphin, narwhal, corals, and algae. | 1. 'Cute sea animals' \| 'sea animals' \| 'Animals'<br>2. 'Sharks' \| 'sharks' \| 'Sea Animals'<br>3. 'Fish' \| 'fish' \| 'Sea Animals'<br>4. 'Piranha' \| 'piranha' \| 'Sea Animals'<br>5. 'Octopus' \| 'octopus' \| 'Sea Animals'<br>6. 'Manatee' \| 'manatee' \| 'Sea Animals'<br>7. 'Whale' \| 'whale' \| 'Sea Animals'<br>8. 'Dolphin' \| 'dolphin' \| 'Sea Animals'<br>9. 'Narwhal' \| 'narwhal' \| 'Sea Animals'<br>10. 'Corals' \| 'corals' \| 'Sea Animals'<br>11. 'Algae' \| 'algae' \| 'Plants' |
|  | The image depicts a luxurious spa room in the Shangri-La Hotel at The Shard, located in London. The room features two massage beds, each with a white cover and pillows. There are also two chairs in the room, one of which is placed next to one of the massage beds. A potted plant can be seen in the corner of the room. | 1. 'Luxurious spa room' \| 'Spa room' \| 'Rooms'<br>2. 'Shangri-La Hotel' \| 'Hotel' \| 'Buildings'<br>3. 'The Shard' \| 'Building' \| 'Buildings'<br>4. 'Two massage beds' \| 'Massage beds' \| 'Furniture'<br>5. 'White cover' \| 'Cover' \| 'Textiles'<br>6. 'Pillows' \| 'Pillows' \| 'Textiles'<br>7. 'Two chairs' \| 'Chairs' \| 'Furniture'<br>8. 'Potted plant' \| 'Plant' \| 'Plants' |
|  | The image features a black bowl of Miso Noodle Soup, which contains tofu puffs, hard-cooked eggs, rice noodles, Asian green vegetables, roasted seaweed, and soup. | 1. 'Black bowl of Miso Noodle Soup' \| 'Bowl' \| 'Kitchenware'<br>2. 'Miso Noodle Soup' \| 'Soup' \| 'Food'<br>3. 'Tofu puffs' \| 'Tofu' \| 'Food'<br>4. 'Hard-cooked eggs' \| 'Eggs' \| 'Food'<br>5. 'Rice noodles' \| 'Noodles' \| 'Food'<br>6. 'Asian green vegetables' \| 'Vegetables' \| 'Food'<br>7. 'Roasted seaweed' \| 'Seaweed' \| 'Food' |
|  | The image features a dining table with a variety of food items, including croissants, bananas, grapes, and apples. There is also a bowl of fruit on the table. A bottle of water can be seen in the background. | 1. 'Dining table' \| 'table' \| 'Furniture'<br>2. 'Variety of food items' \| 'food items' \| 'Food'<br>3. 'Croissants' \| 'croissants' \| 'Food'<br>4. 'Bananas' \| 'bananas' \| 'Food'<br>5. 'Grapes' \| 'grapes' \| 'Food'<br>6. 'Apples' \| 'apples' \| 'Food'<br>7. 'Bowl of fruit' \| 'bowl' \| 'Kitchenware'<br>8. 'Bowl of fruit' \| 'fruit' \| 'Food'<br>9. 'Bottle of water' \| 'bottle' \| 'Kitchenware'<br>10. 'Bottle of water' \| 'water' \| 'Beverage' |

Figure 2-a. Examples of refined captions and extracted object entities yield by DetCLIPv3's data pipeline.

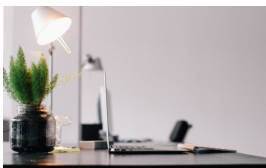| Recaption text | Extracted entities |
|---|---|
|  The image features a plush doll of Santa Claus sitting in a box, reading "The Night Before Christmas" book. | 1. 'Plush doll of Santa Claus' \| 'Plush doll' \| 'Toys'<br>2. 'Santa Claus' \| 'Santa Claus' \| 'Fictional Characters'<br>3. 'Box' \| 'Box' \| 'Containers'<br>4. '"The Night Before Christmas" book' \| 'Book' \| 'Literature' |
|  The image features a colorful and vibrant illustration of the planets of the solar system, including Earth, Mars, Venus, Jupiter, Saturn, Uranus, Neptune, and Pluto. These planets are arranged in a circular pattern on a dark background. | 1. 'Colorful and vibrant illustration' \| 'illustration' \| 'Artwork'<br>2. 'Planets of the solar system' \| 'planets' \| 'Space Objects'<br>3. 'Earth' \| 'Earth' \| 'Planets'<br>4. 'Mars' \| 'Mars' \| 'Planets'<br>5. 'Venus' \| 'Venus' \| 'Planets'<br>6. 'Jupiter' \| 'Jupiter' \| 'Planets'<br>7. 'Saturn' \| 'Saturn' \| 'Planets'<br>8. 'Uranus' \| 'Uranus' \| 'Planets'<br>9. 'Neptune' \| 'Neptune' \| 'Planets'<br>10. 'Pluto' \| 'Pluto' \| 'Planets' |
|  The image features a dining table set up on a balcony overlooking the ocean. The table is adorned with a white plate, which holds a delicious breakfast meal consisting of eggs, bacon, fruit, and orange juice. A fork, knife, and spoon are also present on the table, ready to be used for the meal. | 1. 'Dining table set up on a balcony' \| 'Dining table' \| 'Furniture'<br>2. 'Balcony overlooking the ocean' \| 'Balcony' \| 'Architectural elements'<br>3. 'White plate' \| 'Plate' \| 'Tableware'<br>4. 'Delicious breakfast meal' \| 'Meal' \| 'Food'<br>5. 'Eggs' \| 'Eggs' \| 'Food'<br>6. 'Bacon' \| 'Bacon' \| 'Food'<br>7. 'Fruit' \| 'Fruit' \| 'Food'<br>8. 'Orange juice' \| 'Juice' \| 'Beverages'<br>9. 'Fork' \| 'Fork' \| 'Tableware'<br>10. 'Knife' \| 'Knife' \| 'Tableware'<br>11. 'Spoon' \| 'Spoon' \| 'Tableware' |
|  The image features a collection of ingredients for making fried rice, arranged on a black surface. These ingredients include various vegetables such as peas, carrots, cauliflower, and onions, as well as eggs, soy sauce, and sesame seeds. | 1. 'Collection of ingredients' \| 'ingredients' \| 'Food items'<br>2. 'Fried rice' \| 'rice' \| 'Food items'<br>3. 'Black surface' \| 'surface' \| 'Household items'<br>4. 'Various vegetables' \| 'vegetables' \| 'Food items'<br>5. 'Peas' \| 'peas' \| 'Food items'<br>6. 'Carrots' \| 'carrots' \| 'Food items'<br>7. 'Cauliflower' \|    'cauliflower' \| 'Food items'<br>8. 'Onions' \| 'onions' \| 'Food items'<br>9. 'Eggs' \| 'eggs' \| 'Food items'<br>10. 'Soy sauce' \| 'sauce' \| 'Food items'<br>11. 'Sesame seeds' \| 'seeds' \| 'Food items' |
|  The image depicts a well-organized and spacious desk in a virtual office space. The desk is adorned with various items, including a laptop, a lamp, a potted plant, and a vase. The lamp provides ample lighting for the workspace, while the potted plant adds a touch of greenery and life to the room. | 1. 'Well-organized and spacious desk' \| 'desk' \| 'furniture'<br>2. 'Virtual office space' \| 'office space' \| 'location'<br>3. 'Laptop' \| 'laptop' \| 'electronics'<br>4. 'Lamp' \| 'lamp' \| 'lighting equipment'<br>5. 'Potted plant' \| 'plant' \| 'flora'<br>6. 'Vase' \| 'vase' \| 'decorative item'<br>7. 'Workspace' \| 'workspace' \| 'location'<br>8. 'Room' \| 'room' \| 'location' |
|  The image features a small dog wearing a cowboy hat, sheriff's badge, and boots. | 1. 'Small dog' \| 'Dog' \| 'Animals'<br>2. 'Cowboy hat' \| 'Hat' \| 'Clothing'<br>3. 'Sheriff's badge' \| 'Badge' \| 'Accessories'<br>4. 'Boots' \| 'Boots' \| 'Footwear' |

Figure 2-b. Examples of refined captions and extracted object object entities yield by DetCLIPv3's data pipeline.
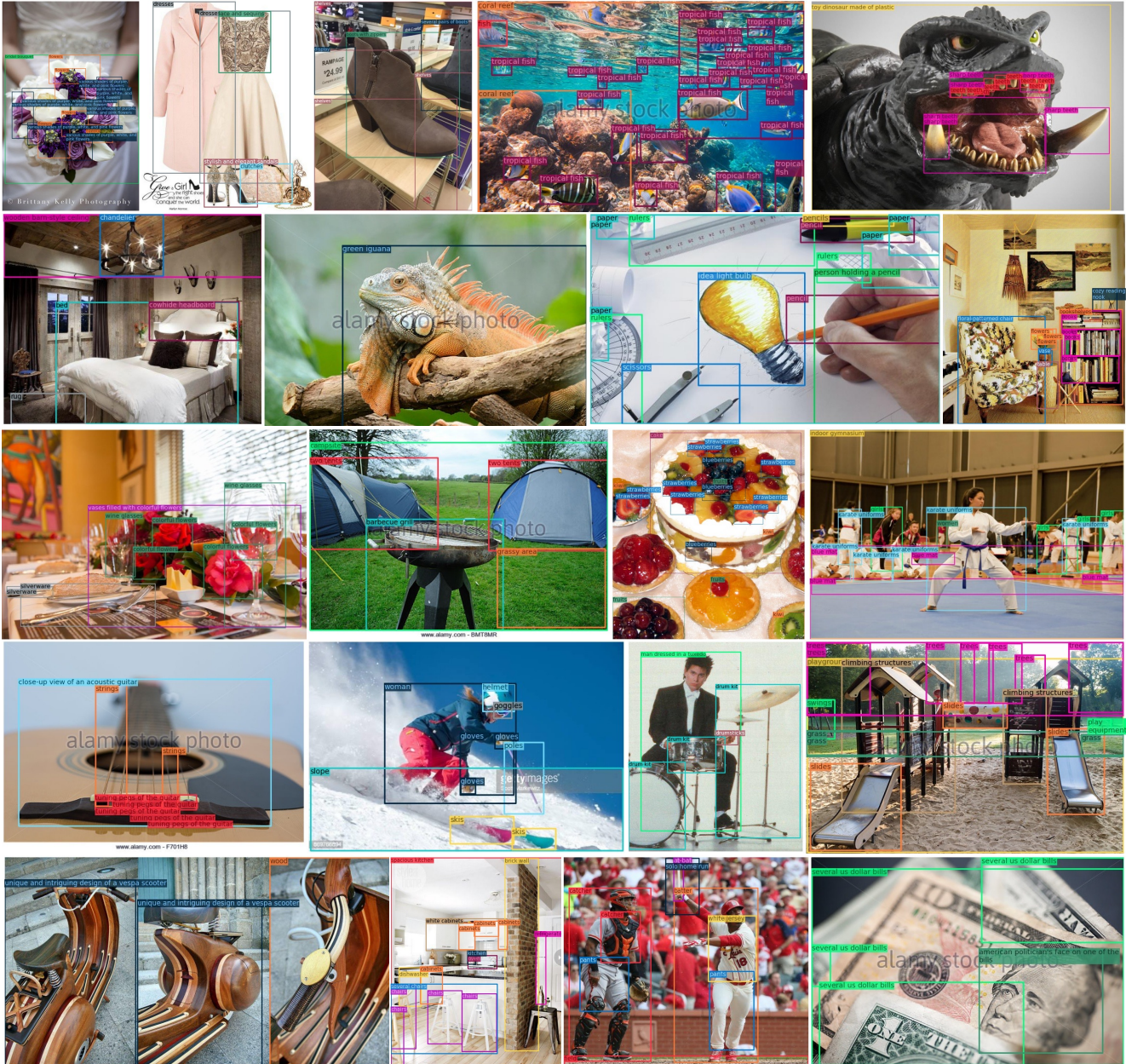
Figure 3. Examples of bounding box pseudo-labels generated by DetCLIPv3's Swin-L model after stage-1 training.

| | Method | Backbone | Pre-training data | Fine-tuning data | LVIS$^{\text{minival}}$ | | | | LVIS$^{\text{val}}$ | | | |
|---|--------|----------|-------------------|------------------|------|------|------|------|------|------|------|------|
| | | | | | AP$_{\text{all}}$ | AP$_{\text{r}}$ | AP$_{\text{c}}$ | AP$_{\text{f}}$ | AP$_{\text{all}}$ | AP$_{\text{r}}$ | AP$_{\text{c}}$ | AP$_{\text{f}}$ |
| 1 | OWL-ST [14] | CLIP B/16 | WebLI2B | – | 31.8 | 35.4 | – | – | 27.0 | 29.6 | – | – |
| 2 | OWL-ST [14] | CLIP L/14 | WebLI2B | – | 38.1 | 39.0 | – | – | 33.5 | 34.9 | – | – |
| 3 | DetCLIPv3 | Swin-T | O365,V3Det,GoldG,GranuCap50M | – | 43.7 | 39.3 | 44.5 | 43.7 | 36.7 | 34.2 | 34.9 | 39.9 |
| 4 | DetCLIPv3 | Swin-L | O365,V3Det,GoldG,GranuCap50M | – | 45.8 | 46.9 | 45.9 | 45.5 | 39.6 | 38.9 | 38.4 | 41.3 |
| 5 | OWL-ST+FT [14] | CLIP B/16 | WebLI2B | LVIS$_{\text{base}}$ | 47.2 | 37.8 | – | – | 41.8 | 36.2 | – | – |
| 6 | OWL-ST+FT [14] | CLIP L/14 | WebLI2B | LVIS$_{\text{base}}$ | 54.1 | 46.1 | – | – | 49.4 | 44.6 | – | – |
| 7 | DetCLIPv3+FT | Swin-T | O365,V3Det,GoldG,GranuCap50M | LVIS$_{\text{base}}$ | 54.4 | 46.7 | 56.1 | 54.3 | 48.2 | 40.2 | 48.5 | 51.3 |
| 8 | DetCLIPv3+FT | Swin-L | O365,V3Det,GoldG,GranuCap50M | LVIS$_{\text{base}}$ | 60.8 | 56.7 | 63.2 | 59.4 | 54.1 | 45.8 | 55.4 | 56.4 |

Table 5. Zero-shot and fine-tuning AP on LVIS val [6] and minival [7]. Results labeled without '+FT' represent zero-shot performance, whereas those with '+FT' indicate results of fine-tuning with LVIS base categories (LVIS$_{\text{base}}$). DetCLIPv3 significantly outperforms OWL-ST [14], which is pre-trained with 2 billion image-text pairs.

| Method | Backbone | Sketch | Weather | Cartoon | Painting | Tattoo | Handmake | Average |
|--------|----------|--------|---------|---------|----------|--------|----------|---------|
| DetCLIPv3 | Swin-T | 38.3 | 43.6 | 45.0 | 43.2 | 29.3 | 31.5 | 38.5 |
| DetCLIPv3 | Swin-L | 50.8 | 48.6 | 56.9 | 53.7 | 44.5 | 38.2 | 48.8 |

Table 6. Detailed performance on COCO-O [12] dataset. Zero-shot AP is reported.

| | Method | Backbone | PascalVOC | AerialDrone | Aquarium | Rabbits | EgoHands | Mushrooms | Packages | Raccoon | Shellfish | Vehicles | Pistols | Pothole | Thermal | Average |
|--|--------|----------|-----------|-------------|----------|---------|----------|-----------|----------|---------|-----------|----------|---------|---------|---------|---------|
| 1 | GLIP [9] | Swin-T | 62.3 | 31.2 | 52.5 | 70.8 | 78.7 | 88.1 | 75.6 | 61.4 | 51.4 | 65.3 | 71.2 | 58.7 | 76.7 | 64.9 |
| 2 | GLIPv2 [24] | Swin-T | 66.4 | 30.2 | 52.5 | 74.8 | 80.0 | 88.1 | 74.3 | 63.7 | 54.4 | 63.0 | 73.0 | 60.1 | 83.5 | 66.5 |
| 3 | GLIPv2 [24] | Swin-B | 71.1 | 32.6 | 57.5 | 73.6 | 80.0 | 88.1 | 74.9 | 68.2 | 70.6 | 71.2 | 76.5 | 58.7 | 79.6 | 69.4 |
| 4 | DetCLIPv2 [23] | Swin-T | 67.5 | 41.8 | 50.8 | 80.4 | 79.8 | 90.1 | 73.7 | 70.8 | 54.8 | 66.5 | 77.7 | 54.8 | 82.2 | 68.5 |
| 5 | DetCLIPv3 | Swin-T | 72.5 | 51.6 | 54.5 | 79.9 | 81.2 | 94.1 | 78.2 | 71.6 | 53.9 | 67.4 | 79.4 | 55.1 | 84.4 | **71.1** |
| 6 | GLIP [9] | Swin-L | 69.6 | 32.6 | 56.6 | 76.4 | 79.4 | 88.1 | 67.1 | 69.4 | 65.8 | 71.6 | 75.7 | 60.3 | 83.1 | 68.9 |
| 7 | GLIPv2 [24] | Swin-H | 74.4 | 36.3 | 58.7 | 77.1 | 79.3 | 88.1 | 74.3 | 73.1 | 70.0 | 72.2 | 72.5 | 58.3 | 81.4 | 70.4 |
| 8 | DetCLIPv2 [23] | Swin-L | 74.4 | 44.1 | 54.7 | 80.9 | 79.9 | 90 | 74.1 | 69.4 | 61.2 | 68.1 | 80.3 | 57.1 | 81.1 | 70.4 |
| 9 | DetCLIPv3 | Swin-L | 76.4 | 51.2 | 57.5 | 79.9 | 80.2 | 90.4 | 75.1 | 70.9 | 63.6 | 69.8 | 82.7 | 56.2 | 83.8 | **72.1** |

Table 7. Detailed fine-tuned AP on ODinW13 [9] dataset. DetCLIPv3 outperforms its counterparts by a large margin.

## C. More Qualitative Results

Figure 4-a, 4-b and 4-c present additional qualitative results showcasing multi-granular object labels generated by DetCLIPv3's object captioner. In the absence of candidate categories, DetCLIPv3's object captioner generates dense, fine-grained, multi-granular object labels, thus facilitating a more comprehensive image understanding.

## D. More Experimental Results

**More results on LVIS.** To comprehensively evaluate the performance of DetCLIPv3, Table 5 provides the standard Average Precision (AP) on LVIS, comparing it with the state-of-the-art method OWL-ST [14], which is pretrained on 2 billion image-text pairs. Specifically, we assess two settings on the LVIS minival [7] and validation [6] datasets: the zero-shot performance and the performance after fine-tuning on LVIS base categories. Despite being pretrained with only 50M image-text pairs, DetCLIPv3 markedly outperforms OWL-ST, *e.g.*, DetCLIPv3 surpasses OWL-ST's counterparts by over 5 AP across all settings, demonstrating the superior learning efficiency of our method. Figure 5 provides the detection results on both zero-shot and LVIS$_{base}$ fine-tuning settings.

**Detailed performance on COCO-O.** Table 6 reports the detailed zero-shot AP performance on COCO-O [12]'s 6 domains, *i.e.*, sketch, weather, cartoon, painting, tattoo and handmake. Figure 6-a and 6-b visualizes the detection results, demonstrating DetCLIPv3's robust domain generalization capability.

**Detailed performance on ODinW13.** Table 7 reports the

| Method | Backbone | OV detector | Object captioner |
|--------|----------|-------------|------------------|
| GLIP [9] | Swin-T | 2.5 FPS | – |
| DetCLIP [21] | Swin-T | 2.3 FPS | – |
| DetCLIPv3 | Swin-T | **14.5** FPS | 1.2 FPS |
| GLIP [9] | Swin-L | 0.3 FPS | – |
| DetCLIPv3 | Swin-L | **8.2** FPS | 0.9 FPS |

Table 8. Inference speed. We test the speed with V100 GPU, using batch size=1 and FP16 inference. DetCLIPv3 can run significantly faster than previous methods like GLIP [9] and DetCLIP [21].

detailed fine-tuned performance on ODinW13 [9] dataset.

**Inference speed.** Table 8 reports the inference speed of DetCLIPv3 as well as a comparison with previous methods.
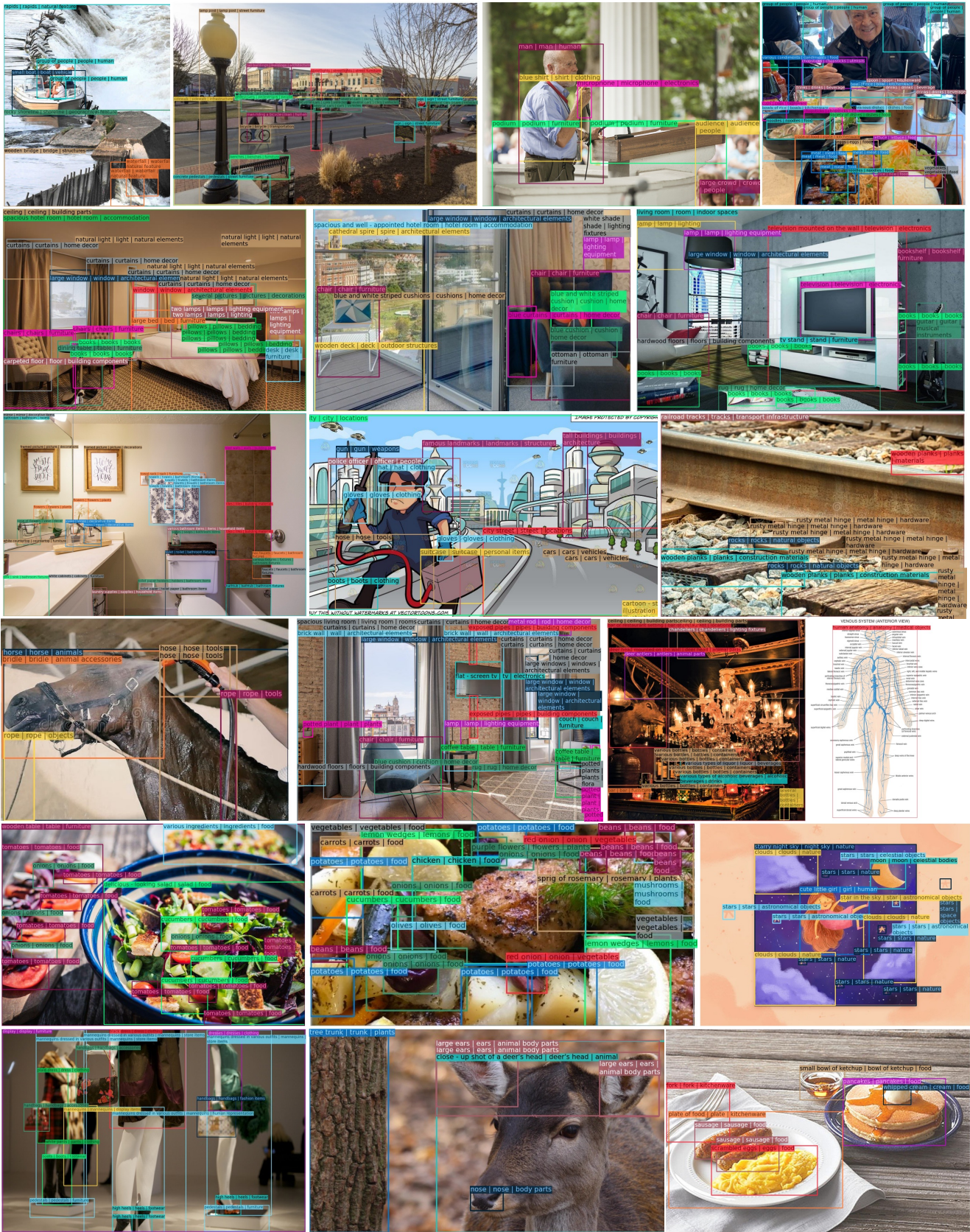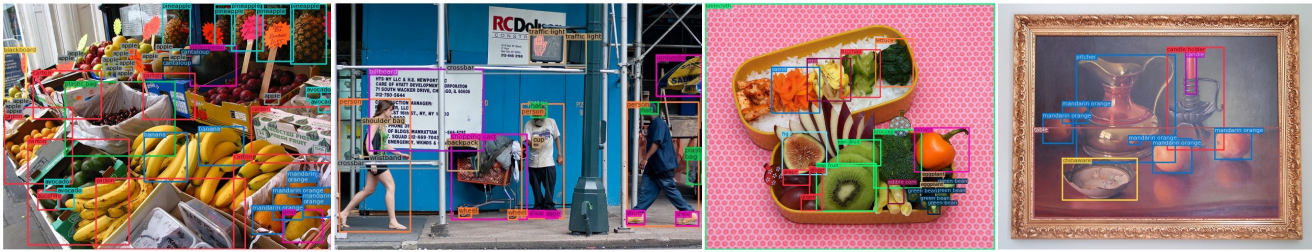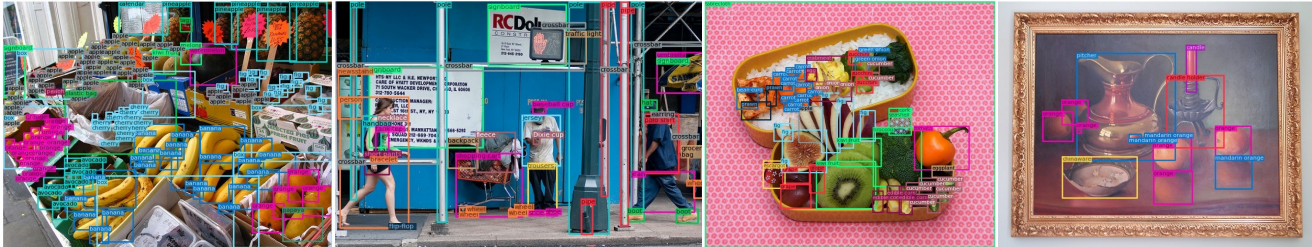
Figure 4-a. Qualitative results of multi-granular object labels generated by DetCLIPv3's object captioner. In the absence of candidate categories, DetCLIPv3's object captioner generates dense, fine-grained, multi-granular object labels, thus facilitating a more comprehensive image understanding.
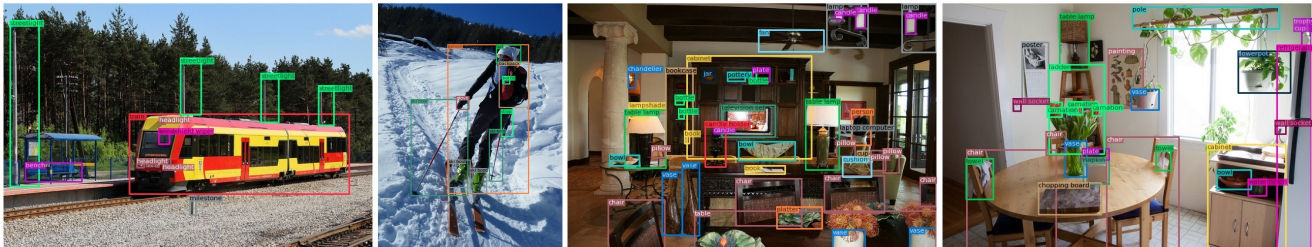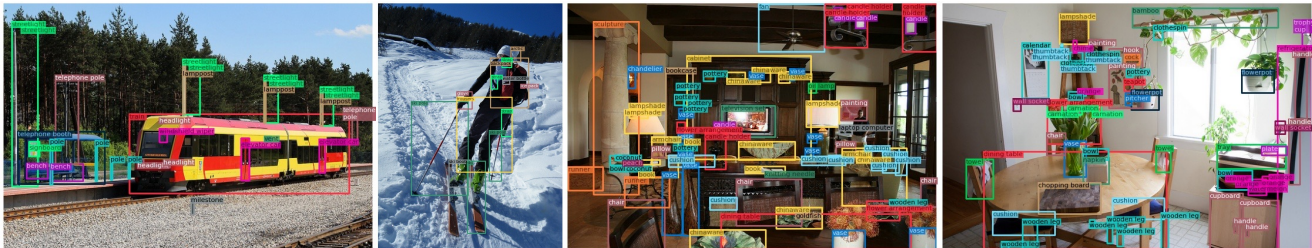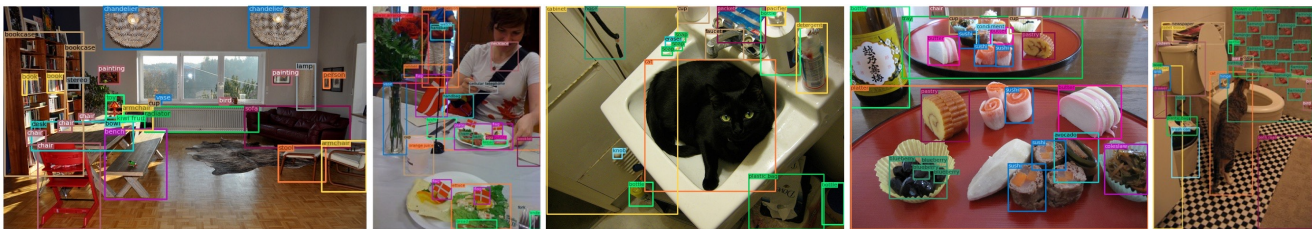
Figure 4-b. Qualitative results of multi-granular object labels generated by DetCLIPv3's object captioner. In the absence of candidate categories, DetCLIPv3's object captioner generates dense, fine-grained, multi-granular object labels, thus facilitating a more comprehensive image understanding.

Figure 4-c. Qualitative results of multi-granular object labels generated by DetCLIPv3's object captioner. In the absence of candidate categories, DetCLIPv3's object captioner generates dense, fine-grained, multi-granular object labels, thus facilitating a more comprehensive image understanding.

Figure 5. Visualization of detection results on LVIS [6]. In each group, the first row represents the zero-shot results, while the second row indicates the results after fine-tuning on the LVIS base categories.
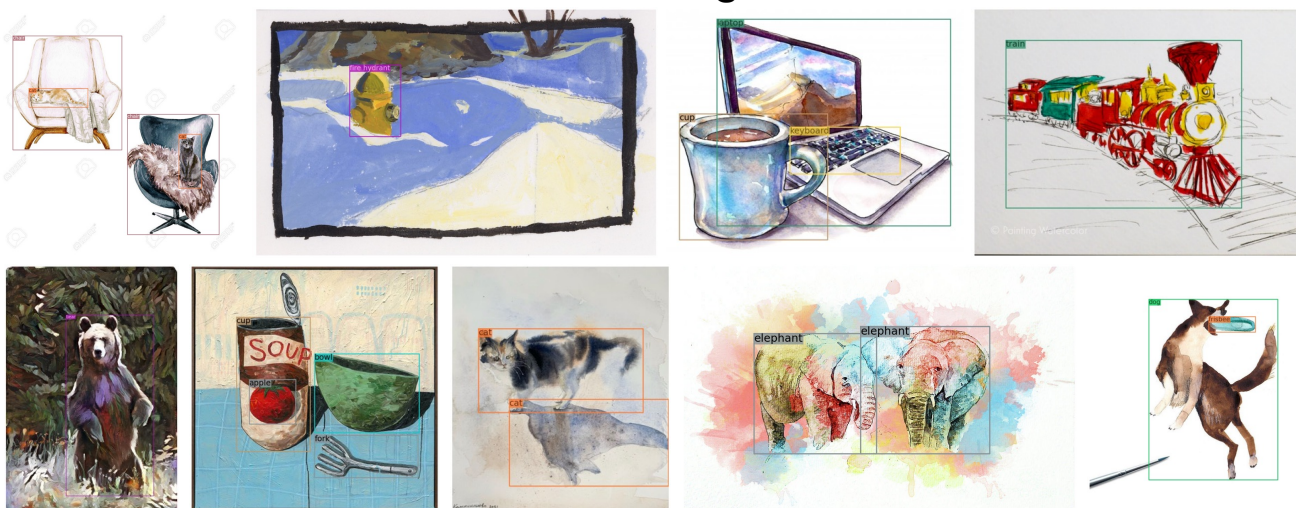
# Sketch

# Weather

# Painting

Figure 6-a. Zero-shot detection results on COCO-O[12] dataset. DetCLIPv3 exhibits a robust domain generalization capability.

# Cartoon
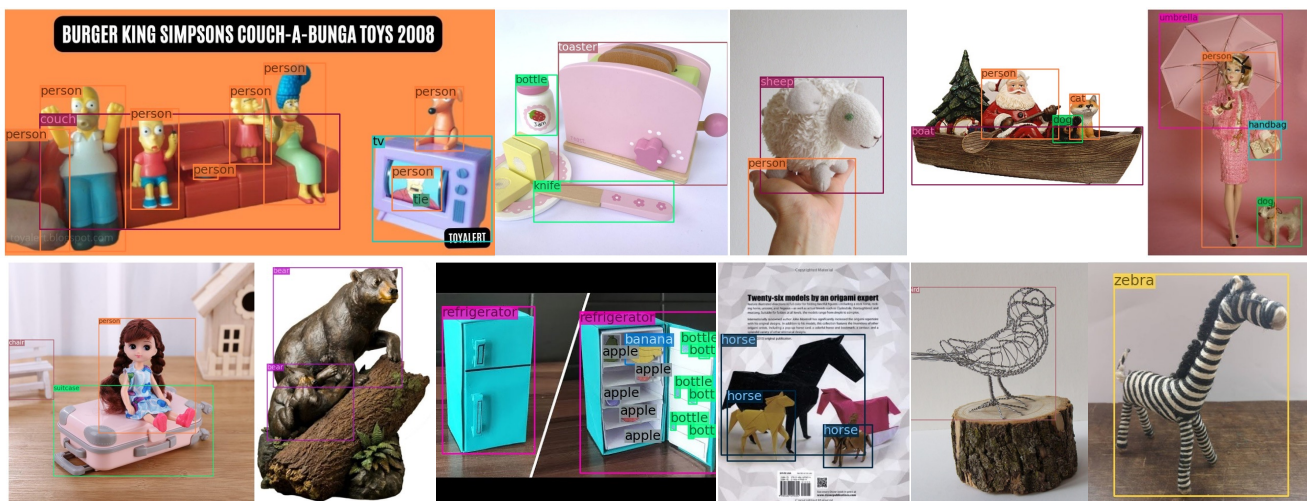


# Tattoo



# Handmake



Figure 6-b. Zero-shot detection results on COCO-O[12] dataset. DetCLIPv3 exhibits a robust domain generalization capability.

# References

[1] Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12M: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In *CVPR*, 2021. 1

[2] Tianqi Chen, Bing Xu, Chiyuan Zhang, and Carlos Guestrin. Training deep nets with sublinear memory cost. preprint arXiv:1604.06174, 2016. 1

[3] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. In *NeurIPS*, 2023. 2

[4] Achal Dave, Piotr Dollár, Deva Ramanan, Alexander Kirillov, and Ross Girshick. Evaluating large-vocabulary object detectors: The devil is in the details. *arXiv preprint arXiv:2102.01066*, 2021. 1

[5] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*, 2019. 1, 2

[6] Agrim Gupta, Piotr Dollar, and Ross Girshick. Lvis: A dataset for large vocabulary instance segmentation. In *CVPR*, 2019. 1, 3, 6, 7, 11

[7] Aishwarya Kamath, Mannat Singh, Yann LeCun, Gabriel Synnaeve, Ishan Misra, and Nicolas Carion. Mdetr-modulated detection for end-to-end multi-modal understanding. In *ICCV*, 2021. 6, 7

[8] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023. 1

[9] 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 *CVPR*, 2022. 1, 3, 7

[10] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023. 3

[11] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2019. 2, 3

[12] Xiaofeng Mao, Yuefeng Chen, Yao Zhu, Da Chen, Hang Su, Rong Zhang, and Hui Xue. Coco-o: A benchmark for object detectors under natural distribution shifts. In *ICCV*, 2023. 7, 12, 13

[13] Paulius Micikevicius, Sharan Narang, Jonah Alben, Gregory Diamos, Erich Elsen, David Garcia, Boris Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, et al. Mixed precision training. In *ICLR*, 2018. 1

[14] Matthias Minderer, Alexey Gritsenko, and Neil Houlsby. Scaling open-vocabulary object detection. *arXiv preprint arXiv:2306.09683*, 2023. 6, 7

[15] 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 *ICML*, 2021. 1, 2

[16] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021. 1

[17] 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 *ICCV*, 2019. 1

[18] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *ACL*, 2018. 1

[19] Bart Thomee, David A Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li. Yfcc100m: The new data in multimedia research. *Communications of the ACM*, 2016. 1

[20] Jiaqi Wang, Pan Zhang, Tao Chu, Yuhang Cao, Yujie Zhou, Tong Wu, Bin Wang, Conghui He, and Dahua Lin. V3det: Vast vocabulary visual detection dataset. In *ICCV*, 2023. 1

[21] Lewei Yao, Jianhua Han, Youpeng Wen, Xiaodan Liang, Dan Xu, Wei Zhang, Zhenguo Li, Chunjing Xu, and Hang Xu. Detclip: Dictionary-enriched visual-concept paralleled pre-training for open-world detection. In *NeurIPS*, 2022. 1, 7

[22] Lewei Yao, Runhui Huang, Lu Hou, Guansong Lu, Minzhe Niu, Hang Xu, Xiaodan Liang, Zhenguo Li, Xin Jiang, and Chunjing Xu. Filip: Fine-grained interactive language-image pre-training. In *ICLR*, 2022. 1

[23] Lewei Yao, Jianhua Han, Xiaodan Liang, Dan Xu, Wei Zhang, Zhenguo Li, and Hang Xu. Detclipv2: Scalable open-vocabulary object detection pre-training via word-region alignmen. In *CVPR*, 2023. 1, 7

[24] Haotian Zhang, Pengchuan Zhang, Xiaowei Hu, Yen-Chun Chen, Liunian Li, Xiyang Dai, Lijuan Wang, Lu Yuan, Jenq-Neng Hwang, and Jianfeng Gao. Glipv2: Unifying localization and vision-language understanding. In *NeurIPS*, 2022. 7

[25] Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun Zhu, Lionel M Ni, and Heung-Yeung Shum. Dino: Detr with improved denoising anchor boxes for end-to-end object detection. In *ICLR*, 2023. 1

[26] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. In *ICLR*, 2021. 1