
OV-PARTS: Towards Open-Vocabulary Part Segmentation (*Supplementary Material*)

Coauthor
Affiliation
Address
email

1 The supplementary material is organized as follows:

- 2 • Implementation Details.(Sec. A)
- 3 • Details of Benchmark Datasets: Pascal-Part-116 and ADE20K-Part-234 (Sec. B).
- 4 • Qualitative Results of Three Benchmark Tasks (Sec. C).
- 5 • Future Works and Negative societal Impacts (Sec. D).

6 **A Implementation Details**

7 **Two-Stage Baselines.** Except for the Object Mask Prompt and Compositional Prompt Tuning designs,
8 we follow the default architecture in the original ZSseg paper. The number of part queries is set to 50.
9 All the two-stage baselines are trained with AdamW optimizer with the initial learning rate of $1e-4$
10 and weight decay of $1e-4$. A poly learning rate policy with a power of 0.9 is adopted. The total batch
11 size is 16 and the total training iteration is 20k. For the training of Compositional Prompt Tuning, a
12 SGD optimizer with the initial learning rate of $2e-2$ and weight decay of $5e-4$ is used. And we adopt
13 a warm-up cosine learning rate policy with 100 warm-up iterations. The total batch size is 32 and the
14 total training iteration is 3k. We sample 128 training samples for each object part class. The length
15 of the learnable object and part prompt tokens are 4. The object tokens are initialized from the text
16 embedding of the template “a photo of”. The initial value of the learnable fusion weight is 0.5.

17 **One-Stage Baselines.** We adopt the original architecture of both CATSeg and CLIPSeg as described
18 in their respective papers. For finetuning CATSeg, we utilize their pretrained model with a ResNet-101
19 backbone. However, while CATSeg achieves the best performance by finetuning the attention layers
20 of CLIP’s visual encoder in open vocabulary object segmentation, we observe poor performance
21 with the same finetuning strategy in OV-PARTS. In our experiments, we only finetune the backbone
22 with a backbone multiplier of 0.1 and the swin transformer based decoder. We employ an AdamW
23 optimizer with an initial learning rate of $2e-4$, weight decay of $1e-4$, and a cosine learning rate policy.
24 The total batch size is 8, and the training iterations amount to 40k. For CLIPSeg, we utilize the same
25 optimizer settings and learning rate policy as CATSeg. The training iterations are set to 20k for the
26 zero-shot/cross-dataset part segmentation setting and 3k for the few-shot part segmentation setting.
27 *More technical details about CLIPSeg.* CLIPSeg adds a parameter-efficient three-layered transformer
28 decoder to the original CLIP model for segmentation. It integrates visual features from the final
29 layer of the visual encoder and text features of all object part prompts from the text encoder through
30 the FiLM module, forming cross-modal input token embeddings for the decoder. Furthermore, the
31 features extracted from the 3rd, 6th, and 9th layers of CLIP’s visual encoder are projected and added
32 to the intermediate features of the corresponding decoder layers. It’s worth noting that, the visual

Table A.1: Model complexity analysis on Pascal-Part-116. ZSseg+ 1/2st is first/second stage.

Model	Backbone	Image Size	Trainable Params (M)	FLOPs (G)
ZSseg+ 1st	ResNet-101c	512 × 704	61.1	103.9
ZSseg+ 2st	ViT-B/16	384 × 384	0.004	58.9
CATSeg	ResNet-101 & ViT-B/16	384 × 384	30.9	139.0
CLIPSeg	ViT-B/16	352 × 352	1.5	97.9

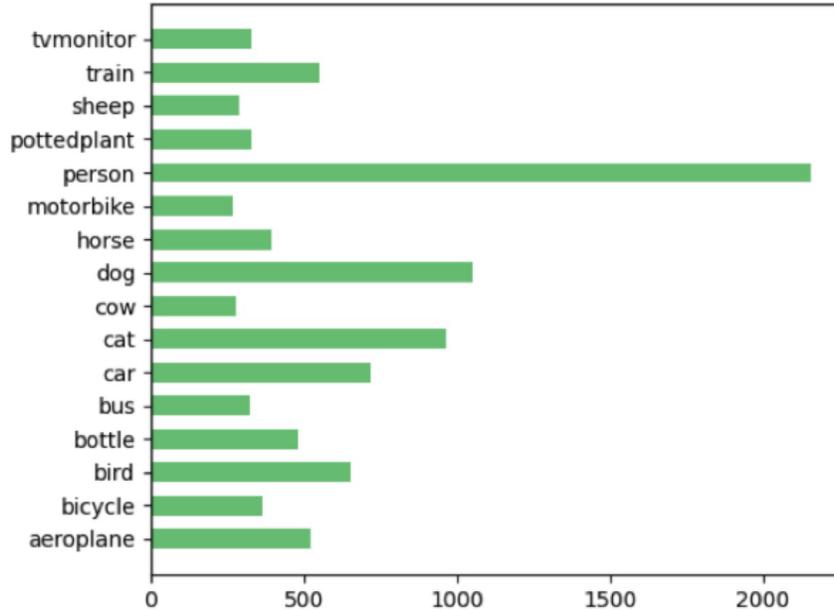


Figure A.1: The number of object masks that have corresponding part masks in Pascal-Part-116.

33 features extracted from the frozen CLIP visual encoder first pass through the added visual adapter,
 34 which consists of a two-layered MLP, before reaching the decoder. Finally, we replace the original
 35 binary segmentation head with a multi-class one to output the semantic segmentation map.

36 **Model Complexity.** We analyze the computational complexity of these two types of baselines
 37 and summarize the number of trainable parameters and FLOPs in Table A.1. The complexity
 38 is evaluated on Pascal-Part-116. It’s evident that the one-stage CLIPSeg, which solely refines a
 39 lightweight transformer decoder and employs parameter-efficient modules, showcases the fewest
 40 trainable parameters and the lowest FLOPs. In contrast, the two-stage ZSseg+ approach, involving
 41 the training of a complete maskformer with a larger resolution, leads to the highest count of trainable
 42 parameters and FLOPs.

43 B Benchmark Datasets Details

44 B.1 Pascal-Part-116

45 Pascal-Part-116 contains 8431 training images and 850 testing images. Compared to the original
 46 version of Pascal-Part, we discard the directional indicator for some part classes and merge them to
 47 avoid overly intricate part definitions. The category vocabulary and merging details are as follows:

48 aeroplane [body, stern, left/right wing, tail, engine, wheel]

49 bicycle [front/back wheel, saddle, handlebar, chainwheel, headlight]

50 bird [left/right wing, tail, head, left/right eye, beak, torso, neck,
 51 left/right leg, left/right foot]

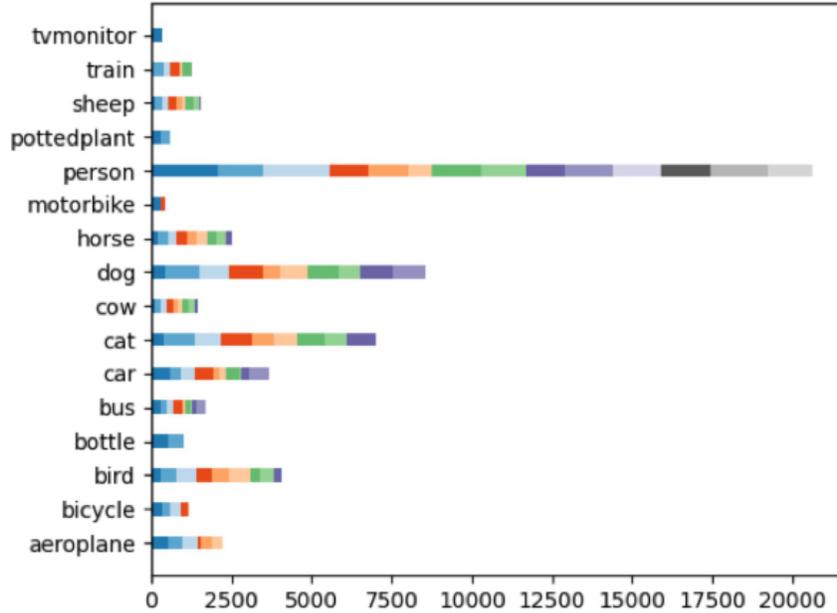


Figure B.2: The number of part masks for each object class in Pascal-Part-116. Each horizontal bar is color-coded to represent a specific part class belonging to the object. The colors of the bars are ordered from left to right based on the part sequence in the list of objects with parts.

```

52 bottle [body, cap]
53 bus [wheel, headlight, front, left/right side, back, roof, left/right
54 mirror, front/back license plate, door, window]
55 car [wheel, headlight, front, left/right side, back, roof, left/right
56 mirror, front/back license plate, door, window]
57 cat [tail, head, left/right eye, torso, neck, left-front/right-front
58 /left-back/right-back leg, nose, left-front/right-front/left-back/right-back
59 paw, left/right ear]
60 cow [tail, head, left/right eye, torso, neck, left-front-upper/left-front-lower
61 /right-front-upper/right-front-lower/left-back-upper/left-back-lower/right-back-upper
62 /right-back-lower leg, left/right ear, muzzle, left/right horn]
63 dog [tail, head, left/right eye, torso, neck, left-front/right-front
64 /left-back/right-back leg, nose, left-front/right-front/left-back/right-back
65 paw, left/right ear, muzzle]
66 horse [tail, head, left/right eye, torso, neck, left-front-upper/left-front-lower
67 /right-front-upper/right-front-lower/left-back-upper/left-back-lower/right-back-upper
68 /right-back-lower leg, left/right ear, muzzle, left-front/right-front/left-back
69 /right-back hoof]
70 motorbike [front/back wheel, saddle, handlebar, headlight]
71 person [head, left/right eye, torso, neck, left-lower/right-lower/left-upper/right-upper
72 leg, foot, nose, left/right ear, left/right eyebrow, mouth, hair,
73 left/right lower arm, left/right upper arm, left/right hand]
74 pottedplant [pot, plant]
75 sheep [tail, head, left/right eye, torso, neck, left-front-upper/left-front-lower
76 /right-front-upper/right-front-lower/left-back-upper/left-back-lower/right-back-upper
77 /right-back-lower leg, left/right ear, muzzle, left/right horn]

```

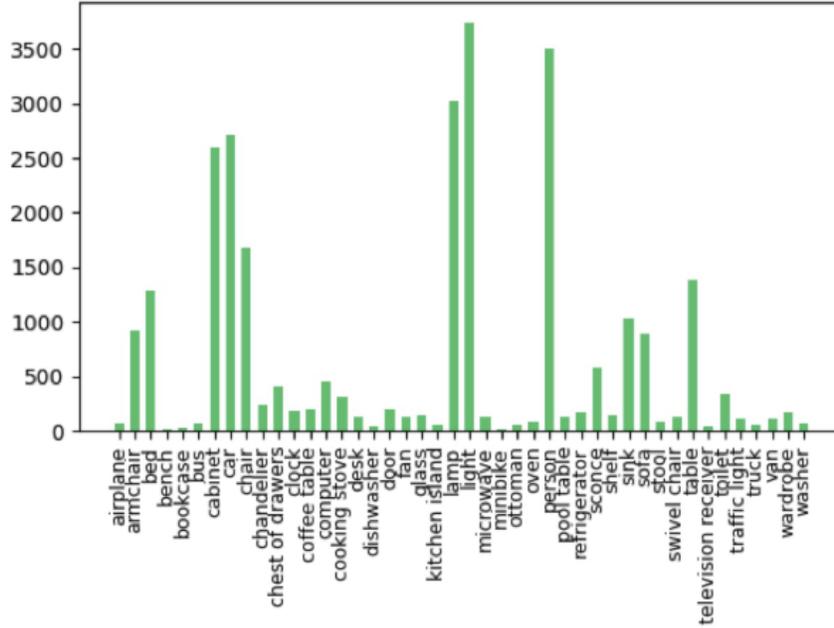


Figure B.3: The statistics for the number of object masks with part masks on ADE20K-Part-234.

78 train [headlight, head, front, left/right side, back, roof, coach]
 79 tvmonitor [screen]
 80 The unseen objects are colored blue and the removed terms are colored purple.

81 B.2 ADE20K-Part-234

82 The original subset of SceneParse150 comprises 20,210 training images and 2,000 validation images.
 83 After filtering out less frequent parts, the subset is reduced to 7,347 training images and 1,016
 84 validation images. In ADE20K, most object parts have sparse mask annotations, and only a subset
 85 of object instances have part annotations. Hence, ADE20K-Part-234 provides the instance-level
 86 object mask annotations along with their part mask annotations. To maximize the use of labeled
 87 data and ensure authentic evaluations, different data splits are designed for the three task settings.
 88 (1) Generalized Zero-Shot Part Segmentation: Models are trained on the seen object instances from
 89 the 7,347 training images. Testing is performed on both unseen object instances from the same
 90 7,347 training images and all object instances from the 1,016 validation images. (2) Few-Shot
 91 Part Segmentation: For each object class, 16 training images are sampled following the approach
 92 in Pascal-Part-116. we adapt the validation set from the generalized zero-shot part segmentation
 93 setting by removing the images that occur in the sampled 16-shot training set. (3) Cross-Dataset Part
 94 Segmentation: The original data split (7347/1016 training/validation images) is used since we mainly
 95 test on the Pascal-Part-116 dataset. The annotated objects with their parts are listed as follows:

96 person [arm, back, foot, gaze, hand, head, leg, neck, torso]
 97 door [door frame, handle, knob, panel]
 98 clock [face, frame]
 99 toilet [bowl, cistern, lid]
 100 cabinet [door, drawer, front, shelf, side, skirt, top]
 101 sink [bowl, faucet, pedestal, tap, top]
 102 lamp [arm, base, canopy, column, cord, highlight, light source, shade, tube]
 103 sconce [arm, backplate, highlight, light source, shade]
 104 chair [apron, arm, back, base, leg, seat, seat cushion, skirt, stretcher]
 105 chest of drawers [apron, door, drawer, front, leg]
 106 chandelier [arm, bulb, canopy, chain, cord, highlight, light source, shade]

107 bed [footboard, headboard, leg, side rail]
108 table [apron, drawer, leg, shelf, top, wheel]
109 armchair [apron, arm, back, back pillow, leg, seat, seat base, seat cushion]
110 ottoman [back, leg, seat]
111 shelf [door, drawer, front, shelf]
112 swivel chair [back, base, seat, wheel]
113 fan [blade, canopy, tube]
114 coffee table [leg, top]
115 stool [leg, seat]
116 sofa [arm, back, back pillow, leg, seat base, seat cushion, skirt]
117 computer [computer case, keyboard, monitor, mouse]
118 desk [apron, door, drawer, leg, shelf, top]
119 wardrobe [door, drawer, front, leg, mirror, top]
120 car [bumper, door, headlight, hood, license plate, logo, mirror, wheel,
121 window, wiper]
122 bus [bumper, door, headlight, license plate, logo, mirror, wheel, window,
123 wiper]
124 oven [button panel, door, drawer, top]
125 cooking stove [burner, button panel, door, drawer, oven, stove]
126 microwave [button panel, door, front, side, top, window]
127 refrigerator [button panel, door, drawer, side]
128 kitchen island [door, drawer, front, side, top]
129 dishwasher [button panel, handle, skirt]
130 bookcase [door, drawer, front, side]
131 television receiver [base, buttons, frame, keys, screen, speaker]
132 glass [base, bowl, opening, stem]
133 pool table [bed, leg, pocket]
134 van [bumper, door, headlight, license plate, logo, mirror, taillight, wheel,
135 window, wiper]
136 airplane [door, fuselage, landing gear, propeller, stabilizer, turbine
137 engine, wing]
138 truck [bumper, door, headlight, license plate, logo, mirror, wheel,
139 windshield]
140 minibike [license plate, mirror, seat, wheel]
141 washer [button panel, door, front, side]
142 bench [arm, back, leg, seat]
143 traffic light [housing, pole]
144 light [aperture, canopy, diffusor, highlight, light source, shade]

145 **B.3 Data Statistics Analysis.**

146 We report the statistics for the number of object masks that have part annotations in Pascal-Part-116
147 (see Figure A.1) and ADE20K-Part-234 (see Figure B.3). The total number of part masks for each
148 object and the proportion of each part are shown in Figure B.2 (Pascal-Part-116) and Figure B.4
149 (ADE20K-Part-234). In Figure B.2, the color sequence from left to right corresponds to the part word
150 sequence as listed in Section B.1. In Figure B.4, the color sequence from bottom to up corresponds to
151 the part word sequence as listed in Section B.2. Additionally, we report the scale distribution for the
152 part masks of each object as shown in Figure B.8.

153 **C Qualitative Results**

154 The qualitative results on the comparison among ZSseg+, CATSeg and CLIPSeg for the challenging
155 case “bird” are shown in Figure B.5.

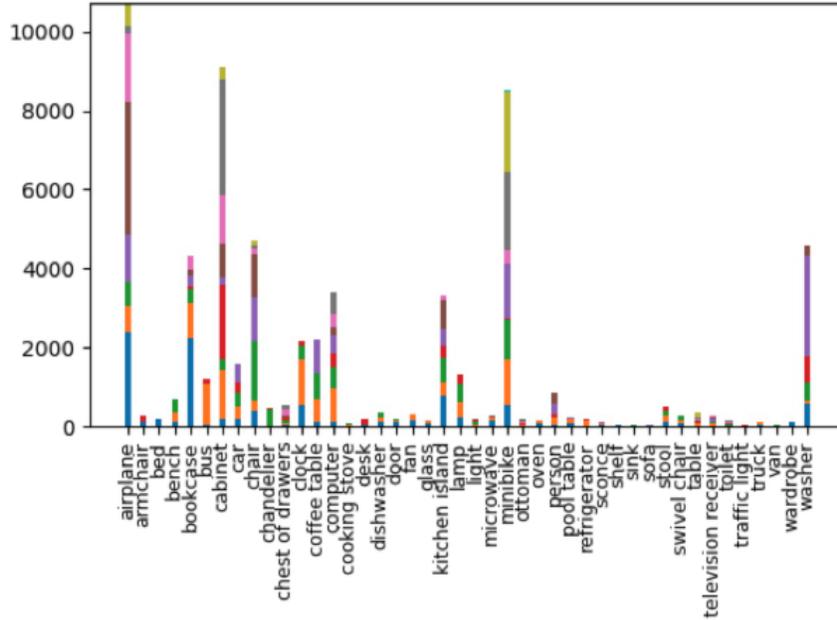


Figure B.4: The number of part masks for each object class in ADE20K-Part-234. Each horizontal bar is color-coded to represent a specific part class belonging to the object. The colors of each bar are ordered from bottom to top according to the part sequence in the list of objects with parts.



Figure B.5: Qualitative results on ZSseg+, CATSeg and CLIPSeg concerning the challenging unseen “bird” class in Pascal-Part-116, as shown in the first row. The second row shows the corresponding ground truth. We can observe that CATSeg and CLIPSeg can generalize to the more novel parts: “Bird’s Beak” and “Bird’s Wing”



Figure B.6: Qualitative results on CATSeg’s multi-granular generalization ability. From the left to the middle image, the model generalizes from “head” to the more fine-grained “beard”. From the middle to the right image, the model generalizes from [“hair”, “eyebrow”, “eye”] to the coarse-grained “head” and also from “neck” to “torso”.

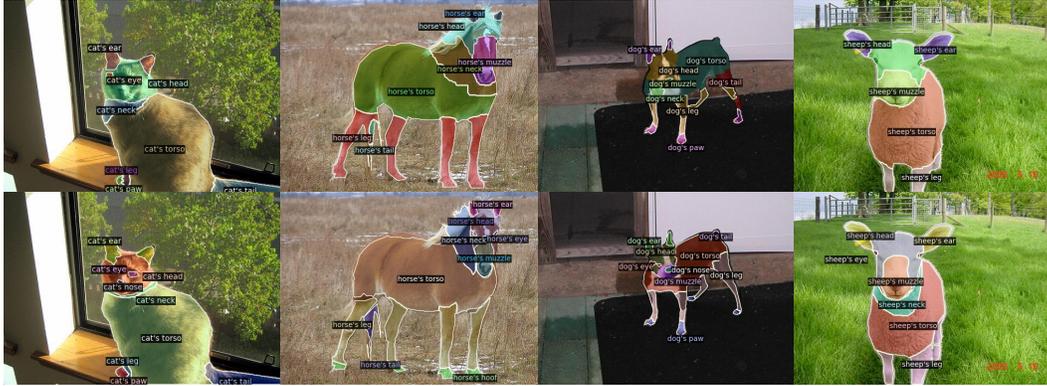


Figure B.7: More qualitative results of generalized zero-shot part segmentation on Pascal-Part-116 are in the first row. The ground truth is in the second row. The seen classes are “cat” and “horse” while the unseen classes are “dog” and “sheep”.

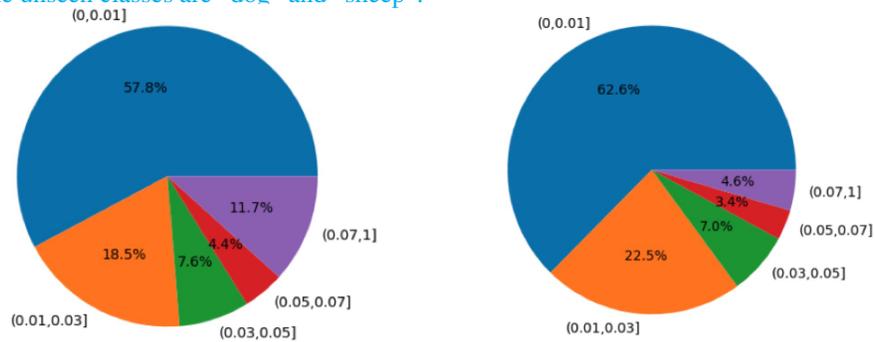


Figure B.8: The scale ratio (number of pixels in the object part mask out of all pixels in an image.) distribution of all part masks of Pascal-Part-116 (Left) and ADE20K-Part-234 (Right).

156 Figure B.6 shows the multi-granular generalization ability of the one-stage baselines. The adopted
 157 model is CATSeg. The visualization sample is from the “person” class in Pascal-Part-116.

158 We give more qualitative results on Pascal-Part-116 and ADE20K-Part-234 on the three proposed
 159 task settings. The adopted model is CLIPSeg with finetuning (VA+L+F+D). The visualization results
 160 for the **Generalized Zero-Shot Part Segmentation** on Pascal-Part-116 and ADE20K-Part-234 are
 161 shown in Figure B.7 and Figure B.9 respectively. We report the qualitative results for the **Few-Shot**
 162 **Part Segmentation** on Pascal-Part-116 in Figure B.10 and on ADE20K-Part-234 in Figure B.11. And
 163 the results for the **Cross-Dataset Part Segmentation** on Pascal-Part-116 are shown in Figure B.12.
 164 Furthermore, we present the qualitative results for models trained on Pascal-Part-116 and then tested
 165 on ADE20K-Part-234 are shown in Figure B.13.

166 D Future Works and Negative Societal Impacts

167 Although part-level OVSS indeed presents more challenges compared to object-level OVSS, the
 168 OV-PARTS benchmark datasets have lower quality than existing object-level OVSS benchmark
 169 datasets. The original version of Pascal-Part and ADE20K-Part are annotated without considering
 170 the open vocabulary scenario especially the analogical reasoning ability and open granularity ability
 171 that we care about in a part-level OVSS model. The benchmark datasets need to be continuously
 172 expanded and improved to encompass more diverse and complex object-part annotations.

173 There may be potential negative societal impacts associated with the OV-PART benchmark. The
 174 deployment of fine-grained part segmentation models in various real-world applications may lead to
 175 unintended consequences. We must ensure that the predictions be reliable and accurate in critical
 176 applications, such as medical diagnosis or autonomous vehicles. Also, there is a possibility of



Figure B.9: Qualitative results of generalized zero-shot part segmentation on ADE20K-Part-234. The first and second rows show the generalize from the seen classes [chair, armchair, sofa] to the unseen classes [swivel chair, ottoman, stool]. The third row shows the generalize from the seen classes [lamp, chandelier] to the unseen class [fan]. Notably, “fan’s blade” is novel at the object and part level.



Figure B.10: Qualitative results of few-shot part segmentation on Pascal-Part-116. We display the segmentation map of four classes: “bird”, “aeroplane”, “car” and “bicycle”.



Figure B.11: Qualitative results of few-shot part segmentation on ADE20K-Part-234. We display the segmentation map of four classes: “lamp”, “sink”, “toilet” and “cooking stove”.

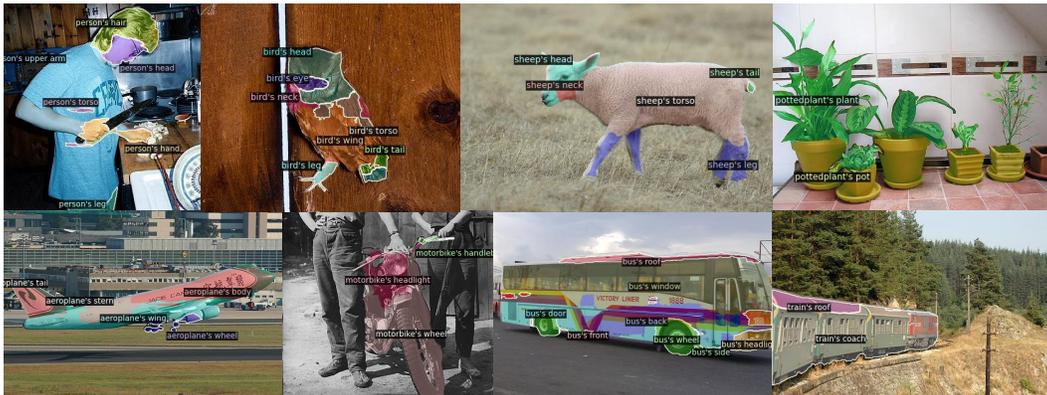


Figure B.12: Qualitative results of cross-dataset part segmentation on Pascal-Part-116. Pascal-Part-116 provides more fine-grained part annotations for the “person” category, such as “hair” and “upper arm”. The model trained on ADE20K-Part-234 demonstrates the ability to recognize “hair” but struggles to generalize from “arm” to “upper arm” and “lower arm” accurately. Moreover, the model exhibits potential in generalizing parts of the “airplane” category. Although ADE20K-Part-234 annotates the parts as “door”, “fuselage”, “landing gear”, “propeller”, “stabilizer”, “turbine engine”, and “wing”, the model can generalize them to Pascal-Part-116’s parts, including “body”, “stern”, “wing”, “tail”, “engine”, and “wheel”, despite the differences in vocabulary and granularity. Notably, ADE20K-Part-234 does not contain related classes to “bird”, “sheep”, and “potted plant”, but the model demonstrates a certain level of generalization ability to segmenting these categories.



Figure B.13: Qualitative results of cross-dataset part segmentation on ADE20K-Part-234. For the categories “car” and “bus”, the part annotations in Pascal-Part-116 are more coarse-grained. When tested on ADE20K-Part-234, the model trained on Pascal-Part-116 can predict novel parts like “logo”, “wiper”, “hood”, and “bumper”. However, the segments and part labels don’t align accurately. For example, the model still segments the “bus’s roof”, which is annotated in Pascal-Part-116, but wrongly assigns it to “bus’s bumper” in ADE20K-Part-234. This showcases the challenge of generalizing across different granular part definitions. For the novel object “swivel chair”, the model adeptly delineates part boundaries even without relevant objects in Pascal-Part-116. But the category errors are still present. In the case of the “person” category, the model only segments the “upper arm”, which demonstrates the difficulty of generalizing from “upper/lower arm” to “arm”.

177 misuse of part segmentation technology for malicious purposes, such as creating deepfake images or
178 spreading misinformation. Ensuring security measures and appropriate regulations to prevent such
179 misuse is vital in the development and deployment process.