

## Contents

<b>A Limitations &amp; Future Work</b>	<b>1</b>
<b>B Experimental Details</b>	<b>1</b>
B.1. Code & Reproduction . . . . .	1
B.2. Device Information . . . . .	1
B.3. Implementation Details . . . . .	1
B.4. Evaluation Details . . . . .	2
B.5. Hyperparameters . . . . .	2
B.6. Baselines . . . . .	2
<b>C Additional Quantitative Evaluation</b>	<b>2</b>
C.1. Evaluation on Recall . . . . .	2
C.2. DINO Structural Guidance . . . . .	2
<b>D Additional Qualitative Results</b>	<b>2</b>
D.1. Pascal-Part-116 (Oracle-Obj Setting) . . . . .	2
D.2. PartImageNet (Pred-All Setting) . . . . .	3
D.3. PartImageNet (Oracle-Obj Setting) . . . . .	3
D.4. Cost Visualization . . . . .	3
<b>E Additional Ablation Study</b>	<b>3</b>
E.1. Compositional Loss . . . . .	3
<b>F. Datasets Details</b>	<b>4</b>
F.1. Pascal-Part-116 . . . . .	4
F.2. ADE20K-Part-234 . . . . .	5
F.3. PartImageNet . . . . .	5
F.4. PartImageNet (OOD) . . . . .	6

## A. Limitations & Future Work

Our PartCATSeg framework has shown strong performance but also has some limitations and potential directions for future work. Specifically, it struggles with handling extremely fine-grained part distinctions, where subtle differences between parts are challenging to capture accurately. Additionally, the framework is designed based on semantic segmentation, as it continues to demonstrate superior performance compared to instance segmentation in many scenarios. However, this design choice inherently limits the ability to distinguish individual parts as separate instances. For example, distinguishing between the left and right handles of a bicycle requires an instance-level understanding that the current framework lacks. Achieving this would necessitate the generation or inclusion of instance object masks, which are not currently supported. Expanding the framework to incorporate instance segmentation could address this limitation, enabling more precise and versatile part segmentation for applications that demand detailed instance-level information.

Future work aims to address these limitations by improving fine-grained differentiation through advanced attention mechanisms [23, 47, 68] and adaptive structural priors tailored to specific datasets. Enhancing the framework’s ability to distinguish subtle part-level differences could significantly improve its performance in complex scenarios. Additionally, integrating the framework with off-the-shelf open-vocabulary instance segmentation modules offers a promising solution for overcoming the inability to individual parts as separate instances. This integration could enable the model to assign unique labels to similar parts across different instances, further extending its applicability and effectiveness.

## B. Experimental Details

### B.1. Code & Reproduction

Details can be found in the publicly available code. For additional details, refer to the GitHub repository available at <https://github.com/kaist-cvml/part-catseg>

### B.2. Device Information

All experiments were conducted using eight NVIDIA A6000 GPUs and PyTorch 2.2 for training and evaluation.

### B.3. Implementation Details

Our model is based on CAT-Seg [13], a state-of-the-art open-vocabulary semantic segmentation (OVSS) method, redefined to suit open-vocabulary part segmentation (OVPS) tasks in OV-PARTS [70]. It employs a CLIP [57] encoder built on CLIP ViT-B/16 and leverages DINOv2 [54], a pre-trained model, for structural guidance.

We begin by utilizing the pre-trained object-level OVSS models from CAT-Seg and fine-tune them with the datasets described in Appendix F. The model undergoes training with the AdamW [49] optimizer, starting with an initial learning rate of 0.0001, over 20,000 iterations, and a batch size of 8. During training, model checkpoints are saved every 1,000 iterations. The final model is selected based on the highest validation performance. For instance, the best validation score on the Pascal-Part-116 dataset in

the Oracle-Obj setting comes from the checkpoint saved at 12,000 iterations.

## B.4. Evaluation Details

For the evaluation protocol, the Pred-All setup of PartCLIPSeg [14] and the Oracle-Obj setup of OV-PARTS [70] were utilized. The Pred-All setup assumes a more challenging scenario in which predictions are made without any prior information. In contrast, the Oracle-Obj setup assumes the availability of object-level masks. As noted in OV-PARTS, the Oracle-Obj setup simulates results achievable when using off-the-shelf open-vocabulary semantic segmentation models.

## B.5. Hyperparameters

The model architecture incorporates layers and parameters from CAT-Seg [13]. Furthermore, as outlined in Equation (15), our method defines three key hyperparameters,  $\lambda_{\text{obj}}$ ,  $\lambda_{\text{part}}$ , and  $\lambda_{\text{comp}}$ , which are associated with two primary loss functions: the disentanglement loss  $\mathcal{L}_{\text{disen}}$  and the compositional loss  $\mathcal{L}_{\text{comp}}$ . These lambda parameters were fine-tuned through experimental validation on the training set to balance the contributions of the proposed loss functions. The final values were determined as  $\lambda_{\text{obj}} = 1.0$ ,  $\lambda_{\text{part}} = 1.0$ , and  $\lambda_{\text{comp}} = 1.0$ .

## B.6. Baselines

- ZSSeg+ [70, 74]: ZSSeg is a two-stage framework for open-vocabulary semantic segmentation that uses CLIP to classify class-agnostic mask proposals, enabling segmentation of seen and unseen classes. ZSSeg+ extends ZSSeg to support part-level segmentation. We evaluate ZSSeg+ using a ResNet-50 [27] baseline, fine-tuned with Compositional Prompt Tuning based on CoOp [82].
- VLPART [61]: VLPART enables open-vocabulary part segmentation by training on data across multiple granularities (part-level, object-level, and image-level) and segments novel objects into parts through dense correspondences with base objects.
- CLIPSeg [50, 70]: CLIPSeg extends CLIP for segmentation tasks, using a transformer-based decoder to generate segmentation maps conditioned on text or image prompts, supporting tasks like referring expression segmentation and zero-shot segmentation. For evaluation, we fine-tune the FiLM layer, decoder, visual encoder, and language embedding layer in the text encoder, following the approach in [70].
- CAT-Seg [13, 70]: CAT-Seg adapts vision-language models like CLIP by aggregating cosine similarity between image and text embeddings to create cost volumes, enabling segmentation of seen and unseen classes. Additionally, CAT-Seg proposes learning the self-attention heads of CLIP’s encoders, achieving effective results.
- PartGLEE [38]: PartGLEE is a part-level segmentation model that uses a unified framework and the Q-Former to model hierarchical relationships between objects and parts, allowing segmentation at any granularity in open-world scenarios.
- PartCLIPSeg [14]: PartCLIPSeg leverages generalized parts and object-level contexts to enhance fine-grained part segmentation, incorporating competitive part relationships and attention mechanisms to improve segmentation accuracy and generalization to unseen vocabularies.

## C. Additional Quantitative Evaluation

### C.1. Evaluation on Recall

We evaluated various baselines across multiple datasets, as detailed in Table 1, Table 2, and Table 3. Here, we provide an additional analysis of recall performance in zero-shot evaluation on Pascal-Part-116 and ADE20K-Part-234. The recall metric measures a model’s ability to correctly identify less frequent or smaller objects, which are often more difficult to segment accurately. As shown in Table A1, our method achieves state-of-the-art recall scores, highlighting its superior capability in identifying and segmenting these challenging parts.

Method	Pascal-Part-116			ADE20K-Part-234		
	Seen	Unseen	h-Recall	Seen	Unseen	h-Recall
ZSSeg+ [74]	65.47	32.13	43.10	55.78	40.71	47.07
CLIPSeg [50, 70]	55.71	43.35	48.76	49.59	48.11	48.84
CAT-Seg [13, 70]	56.00	43.20	48.77	43.48	39.87	41.60
PartCLIPSeg [14]	58.46	47.93	52.67	53.31	51.52	52.40
PartCATSeg (Ours)	<b>67.15</b>	<b>61.02</b>	<b>63.94</b> (+11.27)	<b>64.81</b>	<b>64.22</b>	<b>64.52</b> (+12.12)

<sup>1</sup> The best score is **bold** and the second-best score is underlined.

Table A1. Comparison of zero-shot performance with state-of-the-art methods in terms of **Recall** for **Oracle-Obj** setting on Pascal-Part-116.

### C.2. DINO Structural Guidance

We confirmed the effectiveness of DINO’s structural guidance for PartCATSeg in Table 6 of the main text. Additionally, the supplementary material examines how DINO’s structural guidance impacts the original CAT-Seg. As shown in Table A2, while DINO’s structural guidance proves effective, its performance improvement is relatively modest compared to the proposed PartCATSeg. This highlights that the proposed framework for disentangling parts is more effective overall.

Method	Pred-All			Oracle-Obj		
	Seen	Unseen	h-IoU	Seen	Unseen	h-IoU
CAT-Seg	36.80	23.39	28.60	43.81	27.66	33.91
CAT-Seg w/ Structural Guidance	38.84	28.99	33.20	50.37	36.86	42.57
PartCATSeg w/o Structural Guidance	42.29	27.94	33.65	46.44	31.59	37.60
PartCATSeg	<b>52.62</b>	<b>40.51</b>	<b>45.77</b>	<b>57.49</b>	<b>44.88</b>	<b>50.41</b>

<sup>1</sup> The best score is **bold**.

Table A2. Comparison of zero-shot performance with state-of-the-art methods on Pascal-Part-116.

## D. Additional Qualitative Results

### D.1. Pascal-Part-116 (Oracle-Obj Setting)

The Figure A8 illustrates qualitative evaluations under the Oracle-Obj setting, which assumes that object masks are provided. This setting evaluates the fine-grained part segmentation results using object-level segmentation generated by other off-the-shelf OVSS models.

PartCATSeg consistently demonstrates superior segmentation performance compared to other baselines. Notably, it effectively segments small parts, such as arms and eyes, which are often missed by other models.

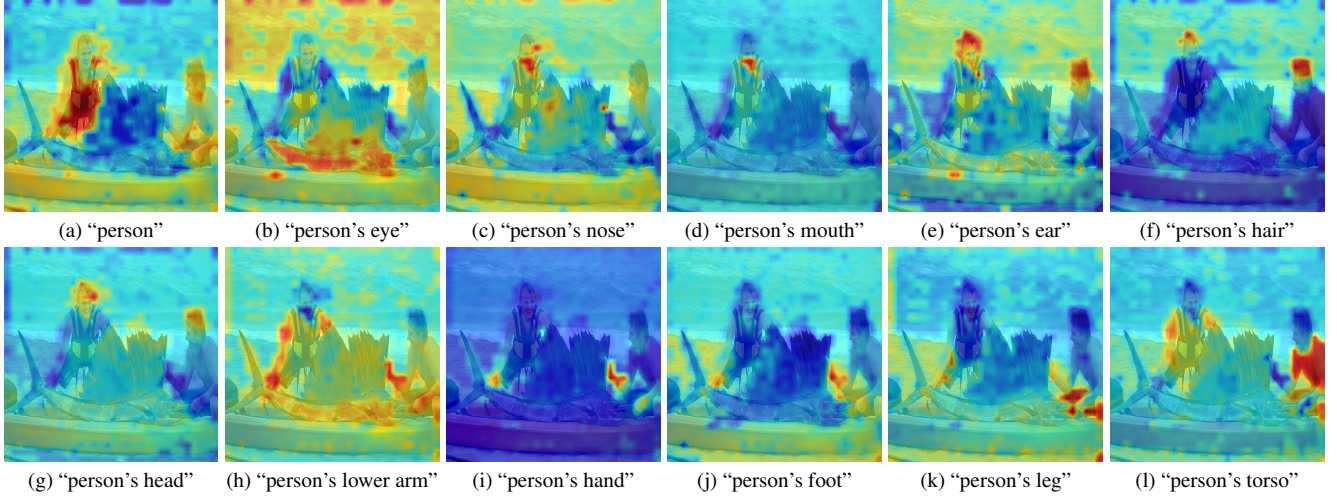


Figure A1. Image-Text Correspondence Visualization **Before** Training

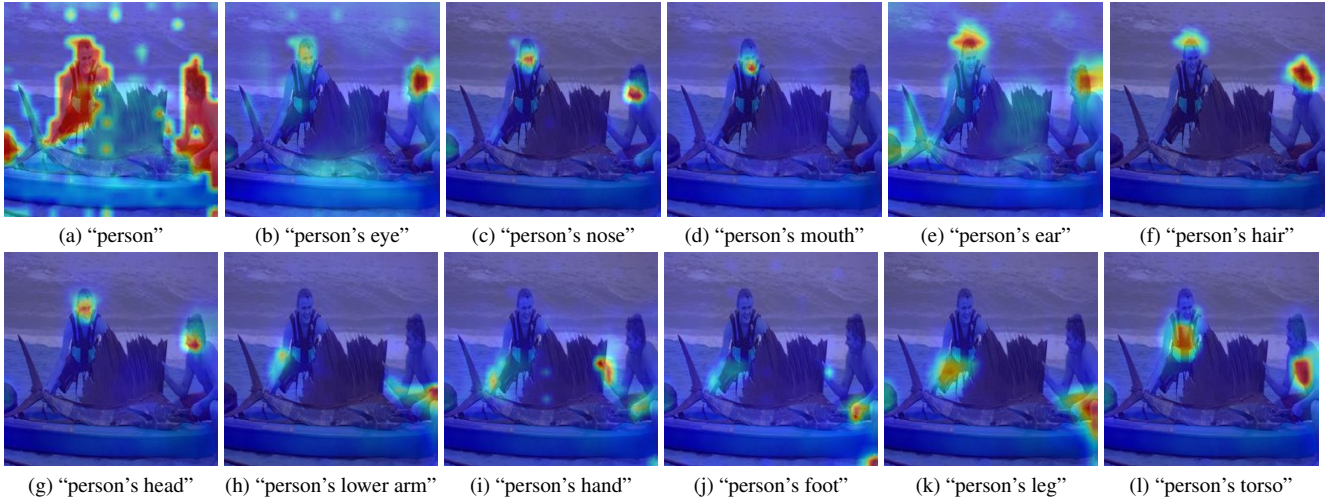


Figure A2. Image-Text Correspondence Visualization **After** Training

## D.2. PartImageNet (Pred-All Setting)

The following Figure A9 illustrates the PartImageNet prediction results of PartCATSeg in Pred-All setup. PartCATSeg demonstrates superior performance compared to other baselines, accurately predicting both appropriate classes and boundaries.

## D.3. PartImageNet (Oracle-Obj Setting)

The following Figure A10 illustrates the PartImageNet prediction results of PartCATSeg in Oracle-Obj setup. PartCATSeg demonstrates superior performance compared to other baselines, accurately predicting appropriate boundaries.

## D.4. Cost Visualization

The following Figure A2 and Figure A4 present the cost (correspondence) visualization after training. It visualizes the object-specific part correspondence between the caption text and the im-

age. Unlike (pretrained) CLIP Image-Text Similarity Visualization (Figure 2) discussed in Section 1 and Figure A1, as well as Figure A3, the cost volume demonstrates significant improvement in fine-grained alignment after training, as illustrated in Figure A2 and Figure A4.

## E. Additional Ablation Study

### E.1. Compositional Loss

As detailed in Section 3.5, parts are not only compositional components that constitute an object but also maintain relationships with adjacent parts. Previous methodologies have proposed learning strategies that consider granularity at two levels—the object and its parts. However, they have not focused on the composition of object-specific parts within the object and the relationships between these parts. This limitation often results in small parts, such as “cat’s eye” and “cat’s neck”, being undetected within the



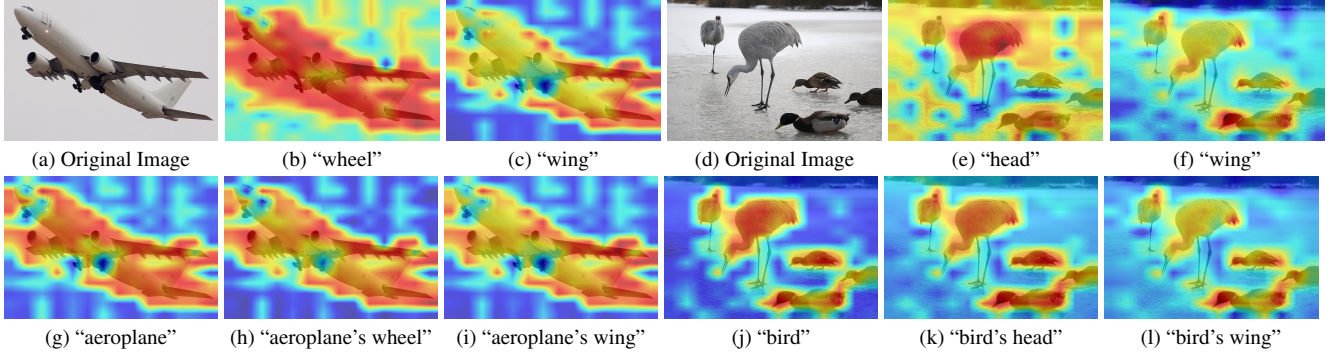


Figure A3. **CLIP Image-Text Similarity Visualization for Object-Level and Part-Level Text.** The visualization compares the frozen CLIP image-text similarity between object-level and part-level text descriptions. (a), (d) show the original images; (b), (c), (e), (f) depict the part-level similarities for terms such as "wheel" and "wing" while (g)-(l) show object-specific parts. The stronger activation for object-level text suggests a dominant focus on the entire object rather than individual parts in the image-text correspondence.

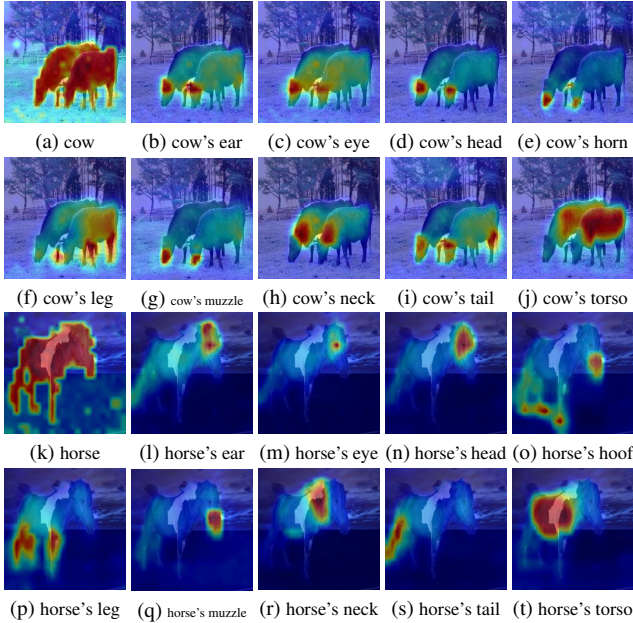


Figure A4. **Cost Volume Visualization After Additional Training.** Cost volume visualization showed that PartCATSeg significantly enhanced fine-grained alignment.

context of a larger object, as shown in Figure A5a. We hypothesize that this issue arises because certain parts become excessively dominant relative to other parts they encompass, leading to a failure to recognize their spatial and compositional relationships. By introducing compositional loss, our model better identifies and segments smaller or more discriminative parts. As illustrated in Figure A5, the inclusion of compositional loss resolves issues of overlapping or diffuse cost volumes, as visualized in Figure A5c and Figure A5d. These visualizations highlight how compositional loss sharpens the focus on specific parts, mitigating the spread of cost volume and ensuring better part-level segmentation.

Furthermore, as shown in Figure A6, compositional loss con-

sistently improves segmentation performance across datasets. For example, in Pascal-Part-116, it enhances the segmentation of challenging parts such as "cow's eye" or "cat's nose," which are often difficult to detect. Specifically, the comparison demonstrates that softmax normalization in the compositional loss ( $\mathcal{L}_{\text{comp-SM}}$ ) outperforms L1 normalization ( $\mathcal{L}_{\text{comp-L1}}$ ). Similarly, as shown in Figure A7, compositional loss demonstrates its effectiveness in capturing fine-grained part relationships in PartImageNet, such as "goose's tail," "tench's tail," and "killer whale's head." This improvement underscores the importance of explicitly modeling compositional relationships in part segmentation, particularly for smaller or less distinct parts.

The effectiveness of compositional loss is further validated through quantitative results, as shown in Table A3. The inclusion of  $\mathcal{L}_{\text{comp}}$  improves performance across both the Pred-All and Oracle-Obj settings on PartImageNet. Notably, it enhances the harmonic IoU for unseen parts, demonstrating its ability to better capture fine-grained compositional relationships and improve segmentation consistency for challenging parts.

Compositional Loss	Pred-All			Oracle-Obj		
	Seen	Unseen	h-IoU	Seen	Unseen	h-IoU
w/o $\mathcal{L}_{\text{comp}}$	<b>59.21</b>	50.75	54.66	72.17	68.42	70.24
w/ $\mathcal{L}_{\text{comp}}$	57.33	<b>53.07</b>	<b>55.12</b>	<b>73.83</b>	<b>71.52</b>	<b>72.66</b>

Table A3. Impact of Compositional Loss on PartImageNet

## F. Datasets Details

### F.1. Pascal-Part-116

Pascal-Part-116 [70] is a modified version of the Pascal-Part [8] dataset specifically designed for Open-Vocabulary Part Segmentation tasks. The dataset includes a mix of base and novel categories, focusing on diverse object-level classes. It features novel classes for the object categories "bird", "car", "dog", "sheep", and "motorbike". Additionally, based on object-specific part categories, the dataset comprises 74 base categories and 42 novel categories, offering a comprehensive benchmark for evaluating segmentation



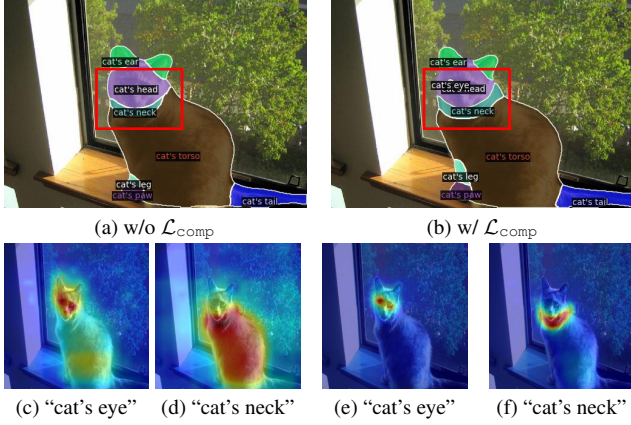


Figure A5. **Ablation on Compositional Loss.** (a) and (b) show segmentation results without and with  $\mathcal{L}_{comp}$ , respectively. (c) and (d) show less defined cost volumes without  $\mathcal{L}_{comp}$ , while (e) and (f) reveal more exclusive similarities in the cost volumes with  $\mathcal{L}_{comp}$ . Notably, "cat's eye" is successfully segmented with the inclusion of  $\mathcal{L}_{comp}$ .

models. Detailed information about the base and novel classes can be found in Table A4.

Pascal-Part-116 Object-specific Part Categories				
Base Categories (74)				
aeroplane's body	aeroplane's stern	aeroplane's wing	aeroplane's tail	aeroplane's engine
aeroplane's wheel	bicycle's wheel	bicycle's saddle	bicycle's handlebar	bicycle's chainwheel
bicycle's headlight	bottle's body	bottle's cap	bus's wheel	bus's headlight
bus's front	bus's side	bus's back	bus's roof	bus's mirror
bus's license plate	bus's door	bus's window	cat's tail	cat's head
cat's eye	cat's torso	cat's neck	cat's leg	cat's nose
cat's paw	cat's ear	cow's tail	cow's head	cow's eye
cow's torso	cow's neck	cow's leg	cow's ear	cow's muzzle
cow's horn	horse's tail	horse's head	horse's eye	horse's torso
horse's neck	horse's leg	horse's ear	horse's muzzle	horse's hoof
person's head	person's eye	person's torso	person's neck	person's leg
person's foot	person's nose	person's ear	person's eyebrow	person's mouth
person's hair	person's lower arm	person's upper arm	person's hand	pottedplant's pot
pottedplant's plant	train's headlight	train's head	train's front	train's side
train's back	train's roof	train's coach	tvmonitor's screen	
Novel Categories (42)				
bird's wing	bird's tail	bird's head	bird's eye	bird's beak
bird's torso	bird's neck	bird's leg	bird's foot	car's wheel
car's headlight	car's front	car's back	car's side	car's roof
car's mirror	car's license plate	car's door	car's window	dog's tail
dog's head	dog's eye	dog's torso	dog's neck	dog's leg
dog's nose	dog's paw	dog's ear	dog's muzzle	motorbike's wheel
motorbike's saddle	motorbike's handlebar	motorbike's headlight	sheep's tail	sheep's head
sheep's eye	sheep's torso	sheep's neck	sheep's leg	sheep's ear
sheep's muzzle	sheep's horn			

Table A4. List of object-specific classes in **Pascal-Part-116**.

## F.2. ADE20K-Part-234

ADE20K-Part-234 [70] is an adapted version of the ADE20K dataset [80] tailored for Open-Vocabulary Part Segmentation tasks. The dataset includes a mix of base and novel categories, with a focus on a diverse range of object-level classes. It features novel classes for the object categories "bench", "bus", "fan", "desk", "stool", "truck", "van", "swivel chair", "oven", "ottoman", and "kitchen island". Additionally, based on object-specific part categories, the dataset comprises 176 base categories and 58 novel categories, providing a robust benchmark for evaluating segmentation models. The dataset presents additional challenges due to its diverse categories and the frequent appearance of small parts, which require precise part segmentation.

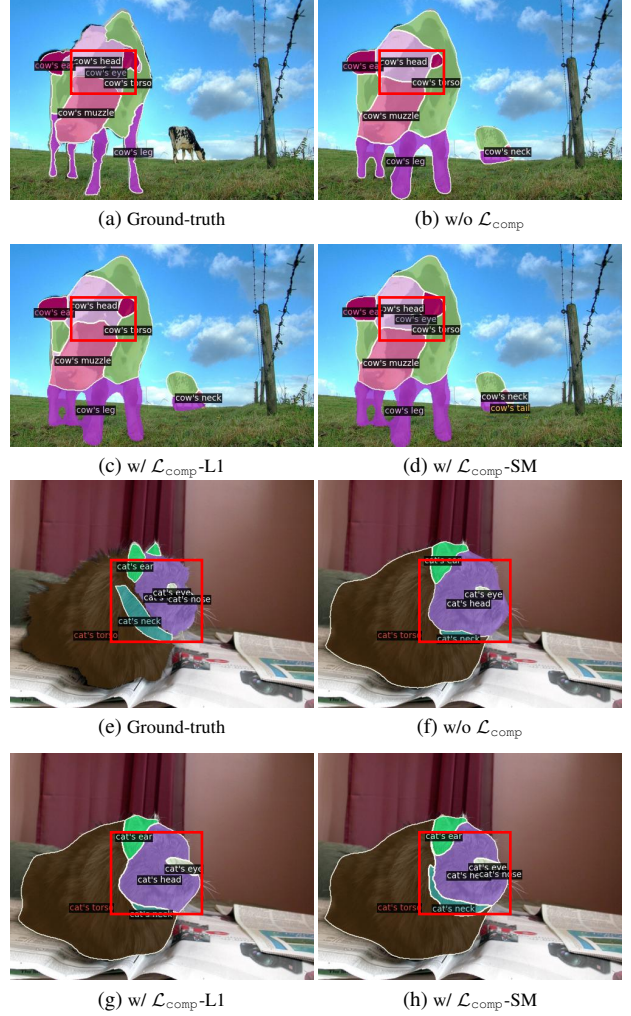


Figure A6. **Qualitative Ablation on Compositional Loss for Pred-All setting in Pascal-Part-116.** Applying compositional loss improves segmentation performance, particularly for challenging parts like "cow's eye" or "cat's nose". Furthermore, using softmax normalization in the compositional loss ( $\mathcal{L}_{comp-SM}$ ) outperforms L1 normalization ( $\mathcal{L}_{comp-L1}$ ) by better capturing these fine-grained parts such as "cat's neck".

## F.3. PartImageNet

PartImageNet [25] is a dataset adapted from ImageNet [17], comprising approximately 24,000 images across 158 categories, each annotated with detailed part information. These categories are grouped into 11 superclasses, based on the hierarchical taxonomy provided by WordNet [51]. In cross-dataset evaluation settings, PartImageNet categories are not pre-divided into base and novel classes. Therefore, for zero-shot evaluation, we select 40 representative object classes from the dataset, which we further split into 25 base object classes and 15 novel object classes, just as in PartCLIPSeg [14]. Each superclass contains corresponding part classes, allowing us to construct part categories by associating these object classes with their respective part classes within each superclass. This selective splitting ensures that the novel classes



Superclass	Base Object Categories (109)	Novel Object Categories (19)
Quadruped	impala, Egyptian cat, warthog, otter, Tibetan terrier timber wolf, polecat, water buffalo, ox, redbone English springer, tiger, American black bear, leopard, hartebeest vizsla, Brittany spaniel, giant panda, Boston bull, ram cairn, Arabian camel, fox squirrel, Eskimo dog, Irish water spaniel Saluki, Walker hound, cheetah, gazelle, soft-coated wheaten terrier bighorn, brown bear, chow, weasel	golden retriever, cougar, ice bear, mink, Saint Bernard
Snake	night snake, boa constrictor, green mamba, thunder snake, green snake hognose snake, sidewinder, horned viper, diamondback, rock python garter snake, vine snake	Indian cobra
Reptile	Gila monster, common newt, green lizard, bullfrog, American alligator leatherback turtle, spotted salamander, box turtle, tailed frog, African chameleon Komodo dragon, agama, frilled lizard, loggerhead	whiptail, alligator lizard
Boat	yawl, pirate	trimaran
Fish	goldfish, coho, tench, anemone fish, killer whale	barracouta, great white shark
Bird	albatross, spoonbill, black stork, dowitcher, American egret goose, ruddy turnstone, bee eater, kite	little blue heron, bald eagle
Car	garbage truck, minibus, ambulance, snowplow, golfcart police van, minivan, convertible, limousine, recreational vehicle go-kart, tractor, school bus, racer	beach wagon, cab
Bicycle	motor scooter, tricycle, mountain bike	unicycle
Biped	gorilla, gibbon, guenon, macaque, patas howler monkey, chimpanzee, proboscis monkey, spider monkey, baboon colobus, capuchin	siamang, marmoset
Bottle	beer bottle, pop bottle, pill bottle	water bottle
Aeroplane	airliner	

Table A7. **PartImageNet (OOD)**. List of selected object classes and their corresponding part classes per superclass from PartImageNet (OOD) [25]. Object categories are divided into base and novel object classes. Detailed associations of part classes with their respective superclasses are provided in Table A6.



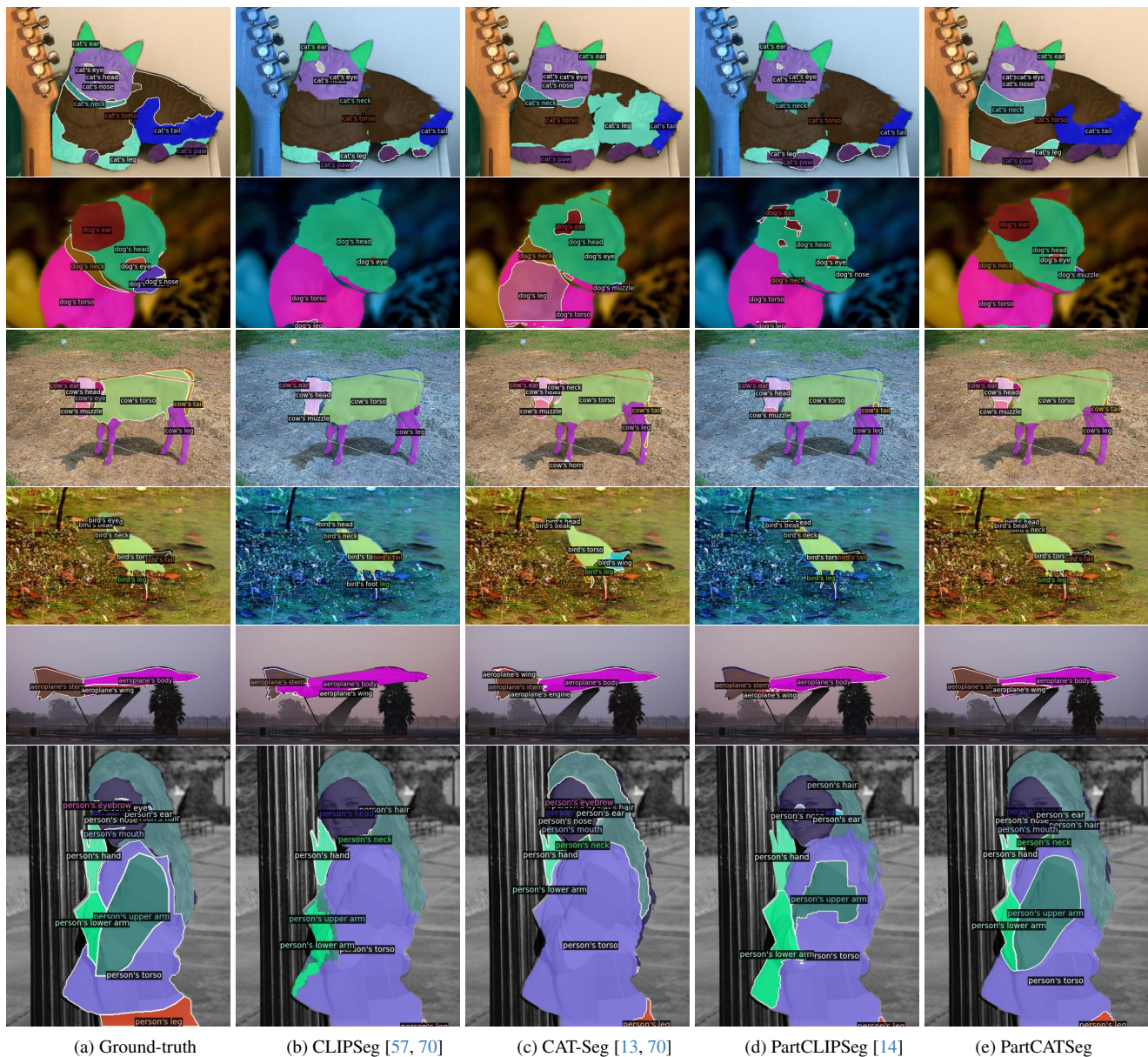


Figure A8. Qualitative evaluation of zero-shot part segmentation on Pascal-Part-116 in the **Oracle-Obj** configuration. Note that annotations for unseen categories (e.g., bird, cow, dog) are excluded from the training set.





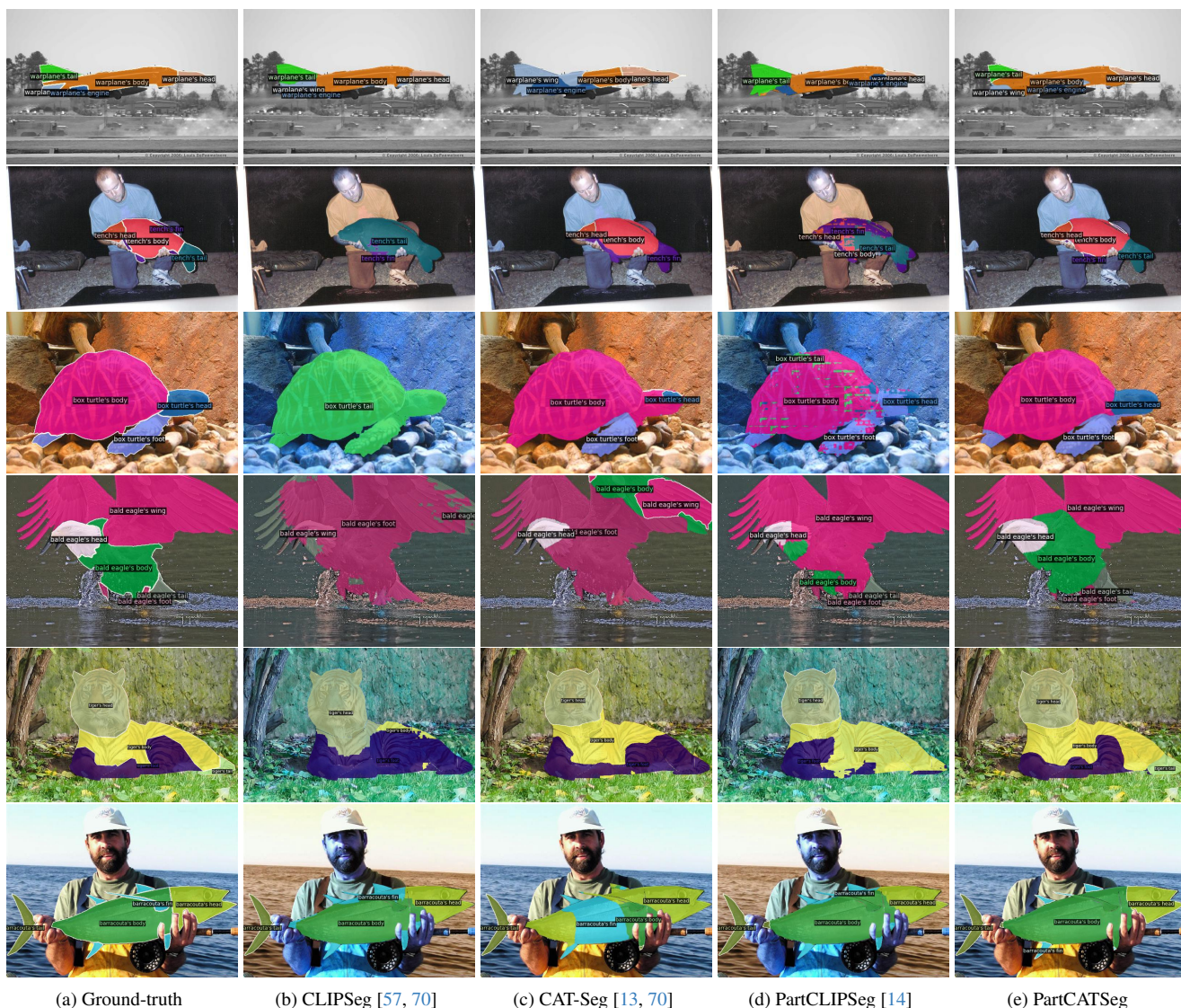


Figure A10. Qualitative evaluation of zero-shot part segmentation on PartImageNet in the **Oracle-Obj** configuration.