

# Understanding Multi-Granularity for Open-Vocabulary Part Segmentation

## (NeurIPS 2024)

Jiho Choi<sup>1\*</sup>, Seonho Lee<sup>1\*</sup>, Seungho Lee<sup>2</sup>, Minhyun Lee<sup>2</sup>, Hyunjung Shim<sup>1</sup>

<sup>1</sup> Graduate School of Artificial Intelligence, KAIST, Republic of Korea

<sup>2</sup> School of Integrated Technology, Yonsei University, Republic of Korea



# Open-Vocabulary Part Segmentation (OVPS)

- Task of segmenting different **parts** of an object in an image using **textual descriptions**
  - textual descriptions that are not constrained by a fixed set of predefined labels
- Recognizing parts that are more complex and diverse than object

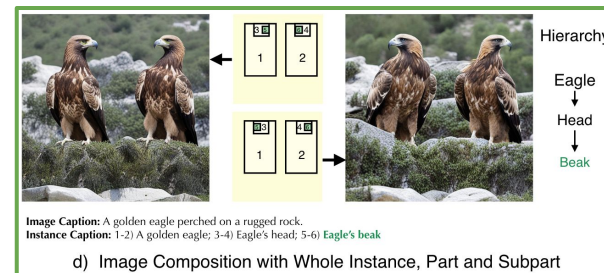
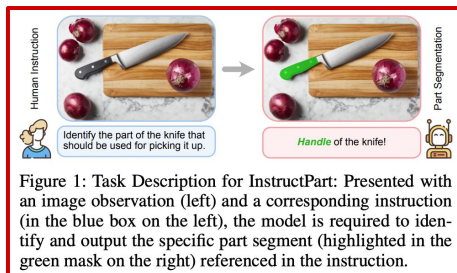
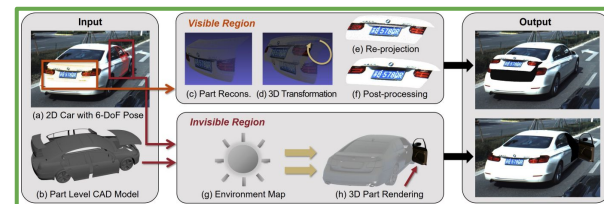
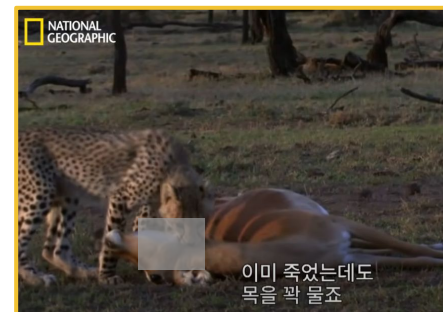


# 1. Introduction [2/5]

## Motivations of OVPS

### Why OVPS?

- Practical Usages
  - **Robot instruction:** Handle of tools (knife, ladle, pot)
  - Part Image **Editing / Generation**
- Multi-granularity Understanding
  - (Biomimicry) Mimicking and understanding **animal instincts**
    - e.g. When a “cheetah” hunts an “impala”, it can distinguish the “neck”



3D Part Guided Image Editing for Fine-grained Object Understanding (CVPR 2020)

InstructPart: Affordance-based Part Segmentation from Language Instruction (AAAIW 2024)

The Guardian - <https://www.theguardian.com/technology/2020/dec/06/the-robot-kitchen-that-will-make-you-dinner-and-wash-up-too>

National Geographic - <https://youtu.be/xVxMisFY3GY?si=kV92174pFiUsFUdx&t=180>

InstanceDiffusion: Instance-level Control for Image Generation (CVPR 2024)



# Challenges of OVPS

## Perspective on Semantic Segmentation

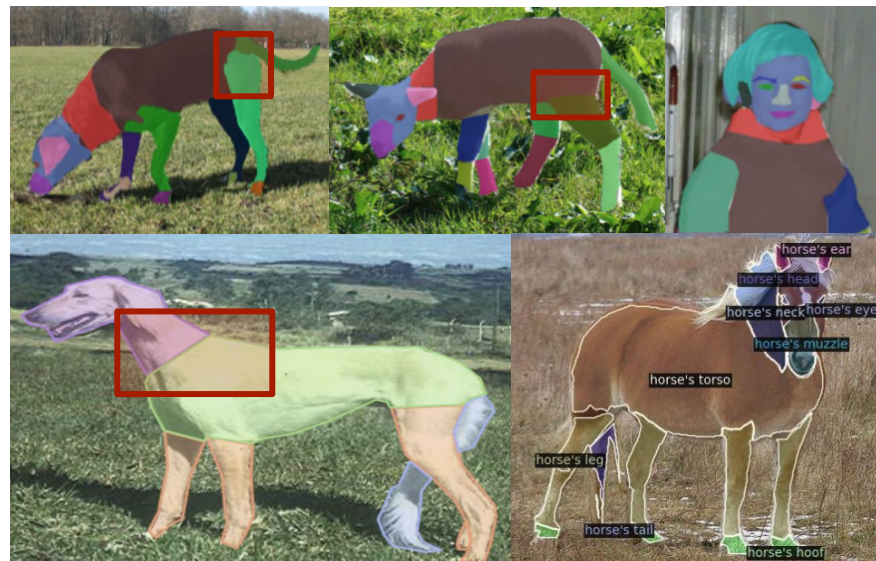
- Parts are **smaller** in size and **diverse** in category compared to objects

## 'Knowledge-based' characteristics of Parts

- Unlike Semantic, Instance, and Panoptic Segmentation, which can achieve clear answers solely from **visual information** present in nature, Parts reflect a knowledge-based nature, making the alignment with language crucial
- Parts are defined by **linguistic or social consensus**

## 'Open Granularity' characteristics of Parts:

- A pixel can be labeled as 'nose,' 'face,' or 'head' depending on the annotation
  - can yield **different ground truth** answers
  - analogous to the "ambiguity" discussed in SAM
- Parts are based on **relative and competitive concepts**



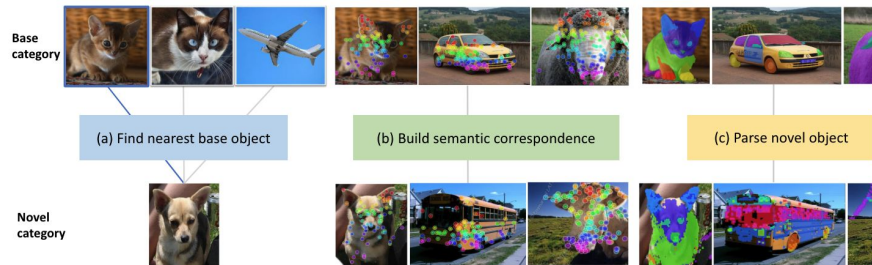
[Top] PASCAL-Part (2014)

[Bottom] Left: PartImageNet (2022) | Right: OV-PARTS (2023)

# Previous OVPS Methods

## VLPart (ICCV 2023.10)

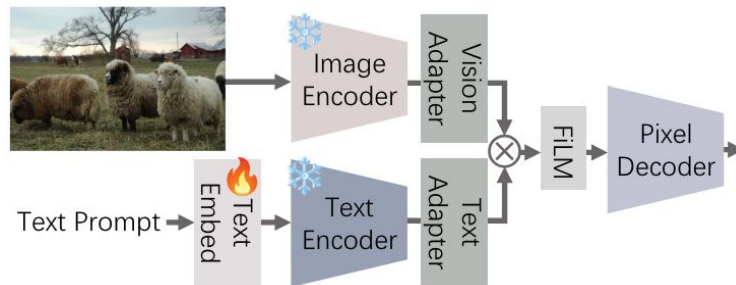
- **DINO** features to **map correspondences** between **base** and **novel** classes and creates **pseudo labels**
- → analogy



Pipeline of VLPart

## OV-PARTS (NeurIPS B&D 2023.12)

- introduces **object mask** prompts and transferring knowledge of base class with **few-shot** approach
- Image-Text Alignment
  - FiLM: Feature-wise Linear Modulation



Overall Framework of OV-PARTS

# Limitations of Previous Methods

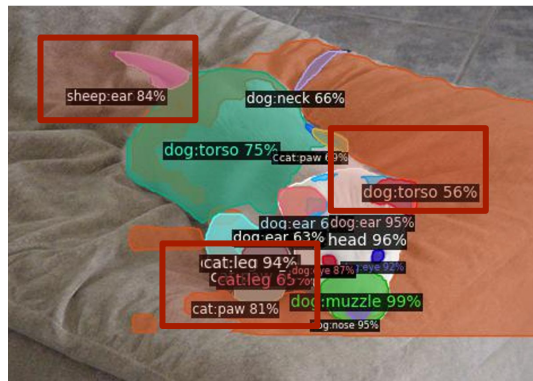
## Zero-shot Part Segmentation

Problem Phenomenon 1)

**Object-level / Part-level Label  
Misclassification**

Assumed Cause 1)

**Lack of Object-level Context /  
Lack of Part-level Generalization**



(a) Lack of generalization

Problem Phenomenon 2)

**Incomplete / Ambiguous  
Boundaries**

Assumed Cause 2)

**Lack of Competitive / Relative  
Partitioning Inductive Bias**



(b) Ambiguous boundaries

Problem Phenomenon 3)

**Missing Labels**

Assumed Cause 3)

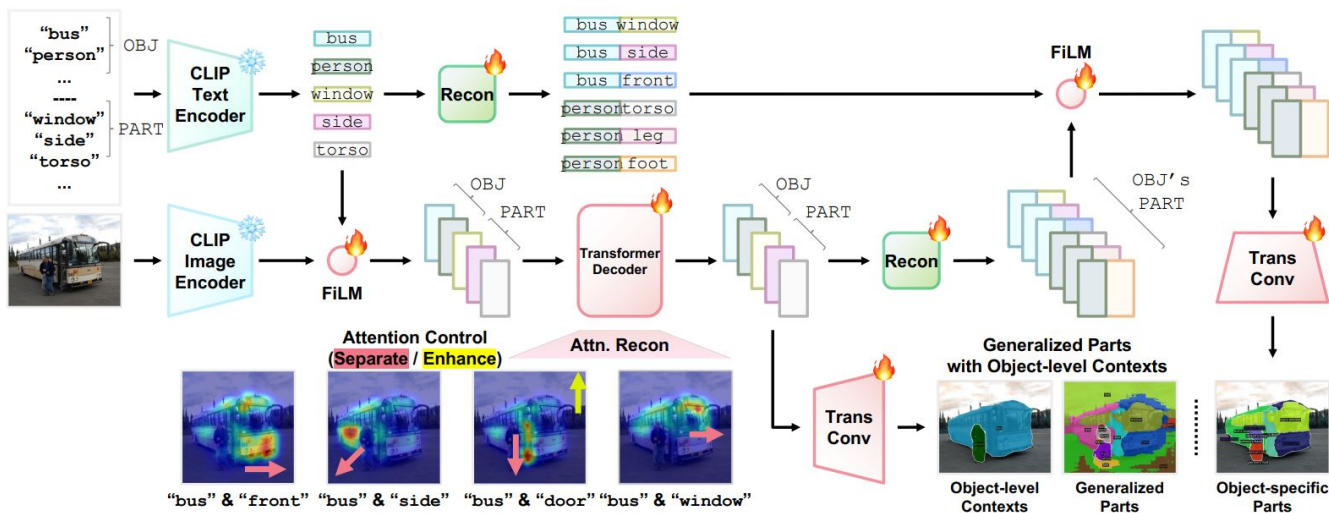
**Small Size / Low-frequency Labels**



(c) Missing underrepresented part

## PartCLIPSeg [1/1]

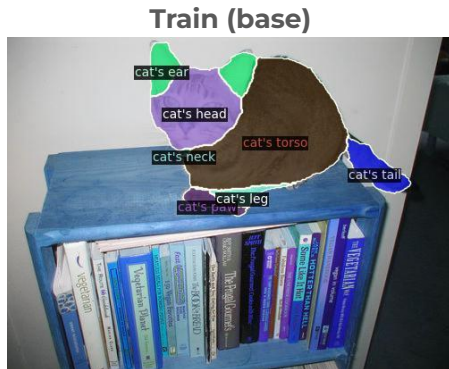
- modified the architecture from CLIPSeg
- Generalized Parts with Object-level Contexts
- Attention Control for Ambiguity and Omission



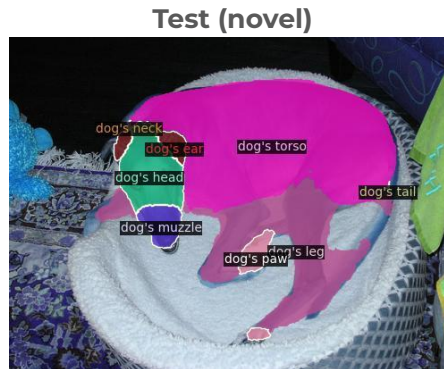
# Generalized Parts with Object-level Contexts [1/2]

## Object-level (Global) Context / Part-level Generalization

- considers **parts** as **common structural components** across object categories and integrates **object contexts**
  - → Part Information, Boundary Information, Holistic Understanding
- previous approaches
  - VLPART (base/novel correspondence, pseudo-label)
  - OV-PARTS (few-shot samples)



**object:** cat, cow, horse etc.  
**part:** head, ear, neck, torso, leg, paw etc.



**common structural components**

**object:** dog etc.  
**part:** head, ear, neck, torso, leg, paw etc.



# Generalized Parts with Object-level Contexts [2/2]

## Object and Part Embedding Generation

- feature extraction with CLIP image, text encoders
  - $\mathbf{e}_{[\text{obj} | \text{part}] }^{\mathcal{T}} = \text{CLIP}_{\mathcal{T}}^*(\mathbf{c}_{[\text{obj} | \text{part}]})$ ,  $\mathbf{e}^{\mathcal{I}} = \text{CLIP}_{\mathcal{I}}^*(\mathcal{I})$
- conditioning image features with FiLM
  - FiLM: an adaptive affine transformation
  - $\mathbf{e}_{[\text{obj} | \text{part}] }^{\mathcal{I}} = \mathbf{e}^{\mathcal{I}} \oplus \text{FiLM}(\mathbf{e}_{[\text{obj} | \text{part}] }^{\mathcal{T}})$

## Object-specific Part Construction

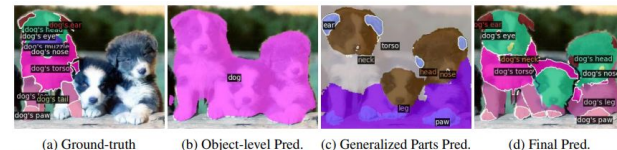
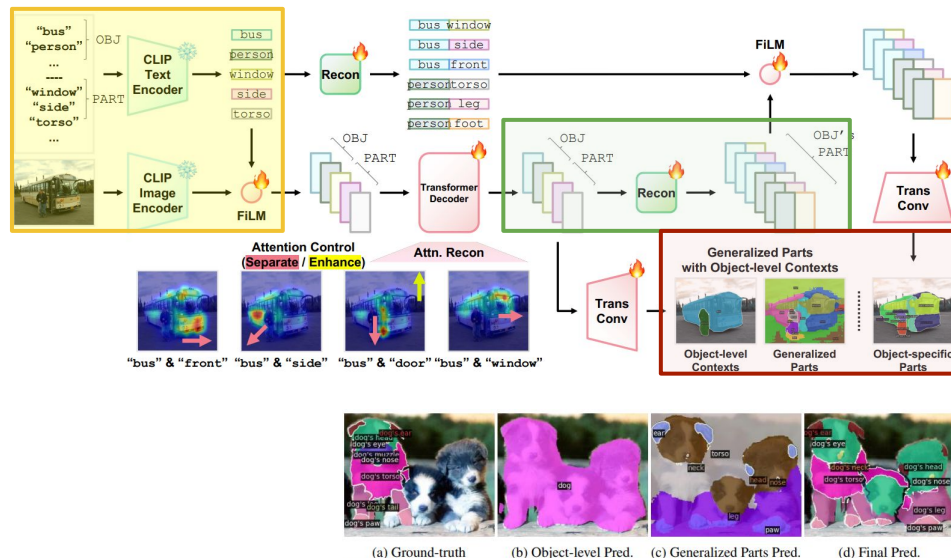
- reconstruction
  - $\mathbf{e}_{\text{obj-part}}^{[\mathcal{T} | \mathcal{I}]} = \text{Proj}(\left[ \mathbf{e}_{\text{obj}}^{[\mathcal{T} | \mathcal{I}]} \mid \mathbf{e}_{\text{part}}^{[\mathcal{T} | \mathcal{I}]} \right])$

## Mask Supervision

- Object, Part, and Object-specific Part

$$\mathcal{L}_{\text{mask}} = \sum_{i=1}^{|\mathcal{C}_{\text{obj-part}}|+1} \underbrace{\{1 - \text{BCE}(s_i, \hat{s}_i)\}}_{\text{object-specific part}} + \lambda_1 \sum_{i=1}^{|\mathcal{C}_{\text{obj}}|+1} \underbrace{\{1 - \text{BCE}(s_i^o, \hat{s}_i^o)\}}_{\text{object guidance}} + \lambda_2 \sum_{i=1}^{|\mathcal{C}_{\text{part}}|} \underbrace{\{1 - \text{BCE}(s_i^p, \hat{s}_i^p)\}}_{\text{generalized part guidance}}$$

\* FiLM: Feature-wise Linear Modulation



# Attention Control for Ambiguity and Omission [1/1]

## Self-attention activation maps

- visual tokens belonging to the same object-specific part mask should exhibit inter-similarity characteristics
- average self-attention map
  - $$\mathcal{A}_{\mathcal{M}_c} = \frac{1}{|\mathcal{M}_c|} \sum_{(h,w) \in \mathcal{M}_c} (\mathcal{A}_{c_{\text{obj}}}[h, w, :, :] + \mathcal{A}_{c_{\text{part}}}[h, w, :, :])$$

## Minimizing Part Overlaps for Ambiguity

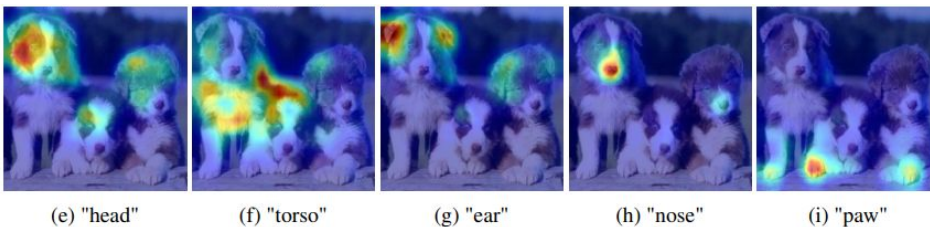
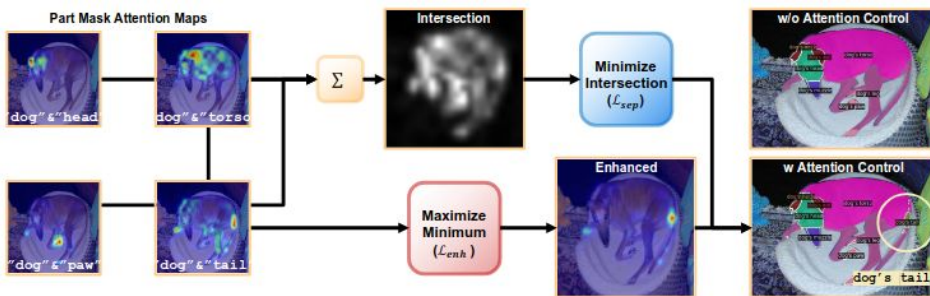
- mitigates the ambiguity issue in part boundaries
- parts with minimized intersection (binarize)

$$\mathcal{L}_{\text{sep}} = \frac{1}{|\mathbf{C}|} \left| \frac{\{(h, w) \mid \sum_{c \in \mathbf{C}} \mathcal{B}_{\mathcal{M}_c}(h, w) > 1\}}{\{(h, w) \mid \sum_{c \in \mathbf{C}} \mathcal{B}_{\mathcal{M}_c}(h, w) \geq 1\}} \right|$$

## Enhancing Part Activation for Omission

- the minimum activation of the part with the maximum value is enhanced

$$\mathcal{L}_{\text{enh}} = 1 - \min_{c \in \mathbf{C}} \left( \max_{(h,w) \in \mathcal{M}_c} \mathcal{A}_{\mathcal{M}_c}[h, w] \right)$$



## 4. Experiments [1/5]

# Experimental Setups [1/1]

### Datasets:

- Pascal-Part-116 (9.2k)
- ADE20K-Part-234 (8.3k)
- PartImageNet (40 classes from 158 categories)

### Tasks:

- Zero-Shot Part Segmentation (Cross-category Part Segmentation)
  - e.g. Unseen: bird, car, dog, sheep, motorbike
- Cross-dataset Part Segmentation

### Evaluation Protocols:

- Pred-All (w/o object mask)
- Oracle-Obj (w/ object mask)

# Performance Evaluation [1/3]

## Quantitative Results

- Zero-shot Part Segmentation
  - Pascal-Part-116
  - ADE20K-Part-234
  - PartImageNet

Table 4: Cross-dataset performance.

Method	Pred-All	Oracle-Obj
PartImageNet → Pascal-Part-116		
CLIPSeg [32, 46]	11.72	14.87
PartCLIPSeg (Ours)	<b>14.74</b> (+3.02)	<b>19.86</b> (+4.99)
ADE20K-Part-234 → Pascal-Part-116		
CLIPSeg [32, 46]	5.41	17.82
PartCLIPSeg (Ours)	<b>10.37</b> (+4.96)	<b>17.94</b> (+0.12)

Table 1: Comparison of zero-shot performance with state-of-the-art methods on Pascal-Part-116.

Method	Backbone	Pred-All			Oracle-Obj		
		Seen	Unseen	Harmonic	Seen	Unseen	Harmonic
ZSSeg+ [52]	ResNet-50	<u>38.05</u>	3.38	6.20	<b>54.43</b>	19.04	28.21
VLPpart [40]	ResNet-50	35.21	9.04	14.39	42.61	18.70	25.99
CLIPSeg [32, 46]	ViT-B/16	27.79	13.27	17.96	48.91	27.54	<u>35.24</u>
CAT-Seg [11, 46]	ViT-B/16	28.17	<b>25.42</b>	<u>26.72</u>	36.20	<u>28.72</u>	32.03
PartCLIPSeg (Ours)	ViT-B/16	<b>43.91</b> <sub>±0.45</sub>	<u>23.56</u> <sub>±0.21</sub>	<b>30.67</b> <sub>±0.09</sub> (+3.94)	<u>50.02</u> <sub>±0.51</sub>	<b>31.67</b> <sub>±0.29</sub>	<b>38.79</b> <sub>±0.13</sub> (+3.55)

<sup>1</sup> The best score is **bold** and the second-best score is underlined. The standard error of an average of 5 results is reported. These are the same for all experiments.

- Cross-dataset Part Segmentation
  - PartImageNet → Pascal-Part-116

Table 2: Comparison of zero-shot performance with state-of-the-art methods on ADE20K-Part-234.

Method	Backbone	Pred-All			Oracle-Obj		
		Seen	Unseen	Harmonic	Seen	Unseen	Harmonic
ZSSeg+ [52]	ResNet-50	<b>32.20</b>	0.89	1.74	<b>43.19</b>	27.84	33.85
CLIPSeg [32, 46]	ViT-B/16	3.14	0.55	0.93	38.15	<u>30.92</u>	<u>34.15</u>
CAT-Seg [11, 46]	ViT-B/16	7.02	2.36	3.53	28.01	21.24	24.16
PartCLIPSeg (Ours)	ViT-B/16	<u>14.15</u> <sub>±0.51</sub>	<b>9.52</b> <sub>±0.13</sub>	<b>11.38</b> <sub>±0.10</sub> (+7.85)	<u>38.37</u> <sub>±0.14</sub>	<b>38.82</b> <sub>±0.31</sub>	<b>38.60</b> <sub>±0.08</sub> (+4.45)

Table 3: Comparison of zero-shot performance with state-of-the-art method on PartImageNet.

Method	Backbone	Pred-All			Oracle-Obj		
		Seen	Unseen	Harmonic	Seen	Unseen	Harmonic
CLIPSeg [32, 46]	ViT-B/16	32.39	12.27	17.80	53.91	37.17	44.00
PartCLIPSeg (Ours)	ViT-B/16	<b>38.82</b> <sub>±0.74</sub>	<b>19.47</b> <sub>±0.45</sub>	<b>25.94</b> <sub>±0.32</sub> (+8.14)	<b>56.26</b> <sub>±0.29</sub>	<b>51.65</b> <sub>±0.62</sub>	<b>53.85</b> <sub>±0.37</sub> (+9.85)

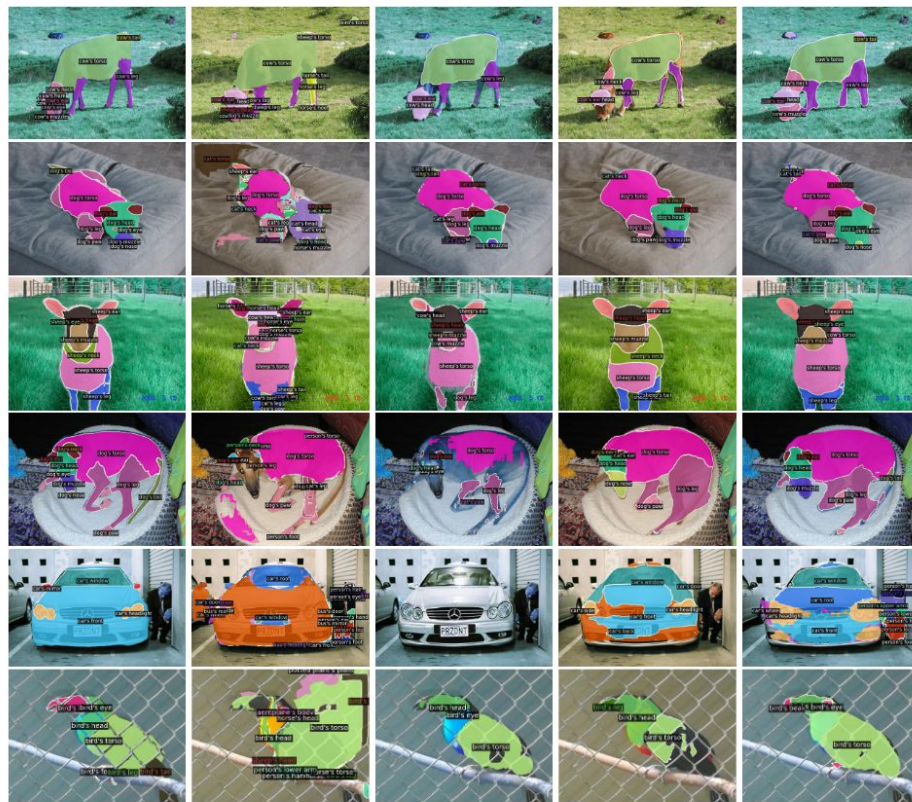


## 4. Experiments [3/5]

# Performance Evaluation [2/3]

## Qualitative Results [1/2]

- Pred-All



(a) Ground-truth (b) VLPART [40] (c) CLIPSeg [38, 46] (d) CAT-Seg [11, 46] (e) PartCLIPSeg (Ours)

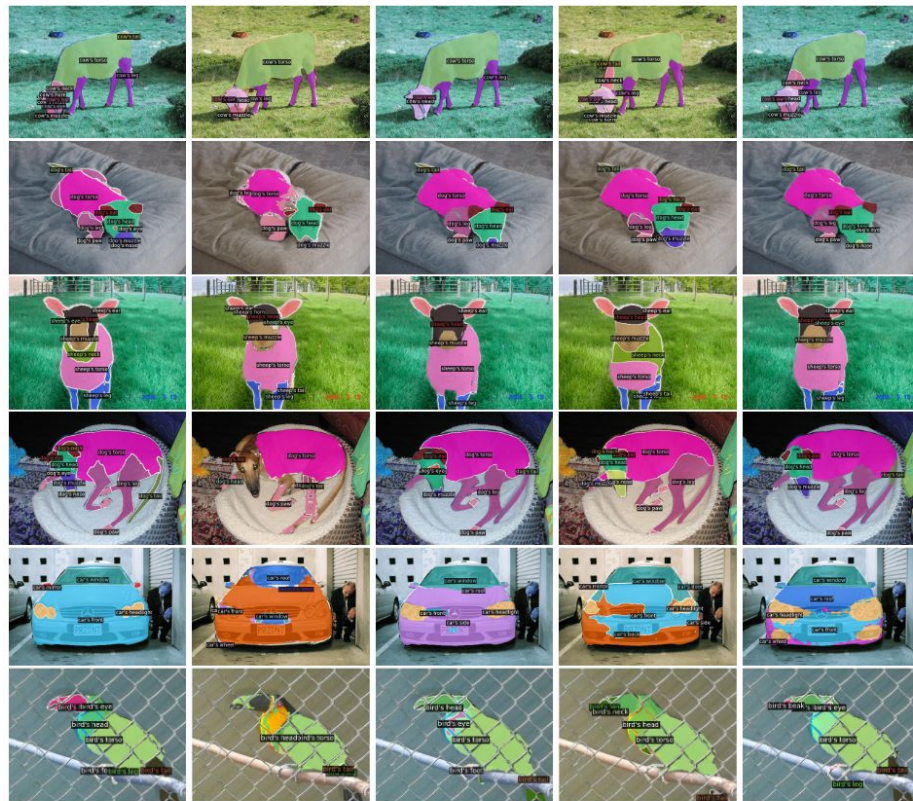
Figure 5: Qualitative results of zero-shot part segmentation on Pascal-Part-116 in **Pred-All** setting. Annotations for unseen categories (bird, car, dog, sheep, etc.) are not included in the train set.

## 4. Experiments [4/5]

# Performance Evaluation [3/3]

## Qualitative Results [2/2]

- Oracle-obj



(a) Ground-truth (b) VLPART [40] (c) CLIPSeg [38, 46] (d) CAT-Seg [11, 46] (e) PartCLIPSeg (Ours)

Figure 6: Qualitative results of zero-shot part segmentation on Pascal-Part-116 in **Oracle-Obj** setting.

# Ablation Study [1/1]

- Separation & Enhance Losses
  - both separation and enhance losses in improving performance
- Impact of PartCLIPSeg for Underrepresented Parts
  - (small parts)



Figure A2: Comparison of results using only  $\mathcal{L}_{sep}$  (top) with both  $\mathcal{L}_{sep}$  and  $\mathcal{L}_{enh}$  (bottom). The heatmap illustrates attention activation for the “sheep’s neck” class.

Table 5: Impact of attention control losses.

Loss		Pred-All			Oracle-Obj		
$\mathcal{L}_{sep}$	$\mathcal{L}_{enh}$	Seen	Unseen	Harmonic	Seen	Unseen	Harmonic
Pascal-Part-116							
✗	✗	43.86	21.89	29.20	49.09	31.26	38.20
✓	✗	<b>44.01</b>	<u>23.18</u>	<u>30.37</u>	<b>50.37</b>	<u>31.45</u>	<u>38.72</u>
✓	✓	<u>43.91</u>	<b>23.56</b>	<b>30.67</b>	<u>50.02</u>	<b>31.67</b>	<b>38.79</b>
ADE20K-Part-234							
✗	✗	10.86	8.33	9.43	37.39	<u>36.49</u>	36.93
✓	✗	<u>12.78</u>	9.38	<u>10.82</u>	<b>39.46</b>	36.04	<u>37.67</u>
✓	✓	<b>14.15</b>	<b>9.52</b>	<b>11.38</b>	<u>38.37</u>	<b>38.82</b>	<b>38.60</b>

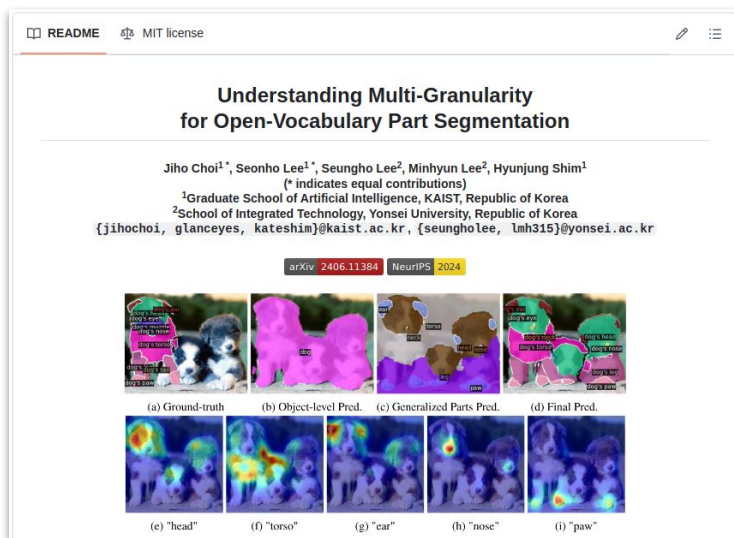
Table 7: Impact of PartCLIPSeg for small parts on Pascal-Part-116 in Oracle-Obj setting. (mIoU)

<b>Part: “eye”</b>	bird	cat	cow	dog	sheep	person
CLIPSeg [32, 46]	<b>3.33</b>	18.77	3.65	16.05	0.00	15.30
PartCLIPSeg (Ours)	1.95	<b>31.01</b>	<b>28.16</b>	<b>32.79</b>	<b>0.67</b>	<b>29.16</b>
<b>Part: “neck”</b>	bird	cat	cow	dog	sheep	person
CLIPSeg [32, 46]	19.09	6.57	0.78	8.12	8.47	30.93
PartCLIPSeg (Ours)	<b>32.51</b>	<b>12.00</b>	<b>2.75</b>	<b>16.37</b>	<b>18.80</b>	<b>50.71</b>
<b>Part: “leg”</b>	bird	cat	cow	dog	sheep	person
CLIPSeg [32, 46]	19.61	38.62	27.85	39.34	52.63	52.67
PartCLIPSeg (Ours)	<b>31.12</b>	<b>44.82</b>	<b>63.78</b>	<b>41.55</b>	<b>54.73</b>	<b>55.35</b>



# Conclusion

- introduced SOTA OVPS method: PartCLIPSeg
- utilizes generalized parts and object-level guidance
- separates parts by minimizing their overlaps in attention maps, handling ambiguous part boundaries
- enhanced loss function to improve the detection of underrepresented parts
- Code: <https://github.com/kaist-cvml/part-clipseg/>





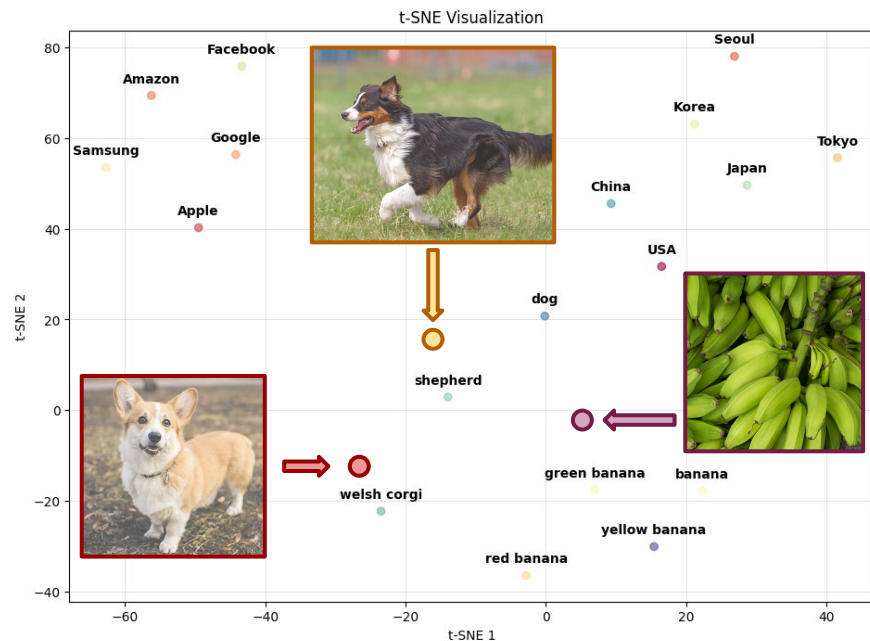
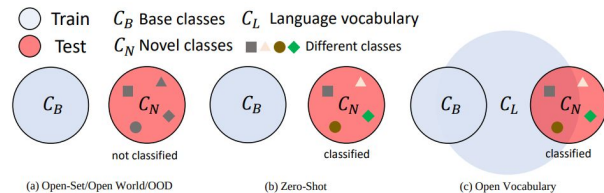
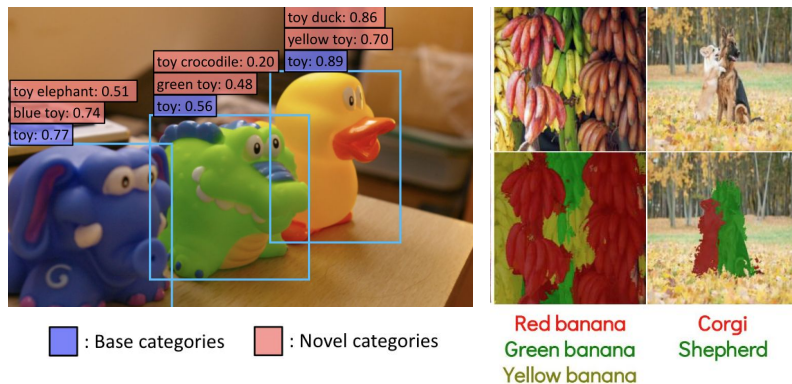
## **Appendix**

## 2. Related Work [1/3]

# VLM & Open-Vocabulary

## Open-Vocabulary & Zero-shot Learning

- Predicting **unseen (novel)** category names, not present in the train set (**seen, base**), is possible using the embedding space of Vision-Language Models (VLMs) like CLIP



Open-vocabulary Object Detection via Vision and Language Knowledge Distillation (ICLR 2022)

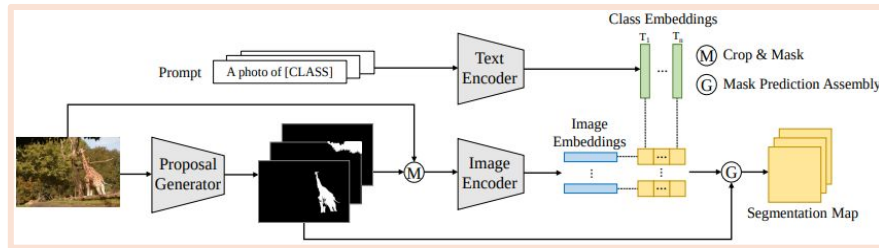
Learning to Generate Text-grounded Mask for Open-world Semantic Segmentation from Only Image-Text Pairs (CVPR 2023)

Towards Open Vocabulary Learning: A Survey (PAMI 2024)

# Open-Vocabulary Dense Prediction

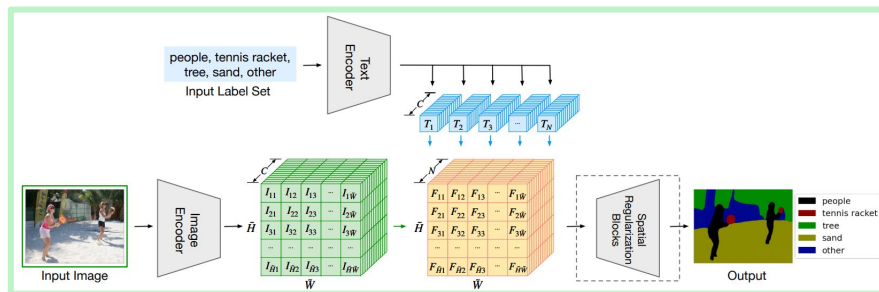
## Two-Stage OVSS / OVIS

- **ZSSeg** (ECCV 2022)
  - **class-agnostic mask proposal**
  - classify proposals with CLIP embedding



## One-Stage OVSS

- **LSeg** (ICLR 2022), **MaskCLIP** (ECCV 2022)
  - (language-driven semantic segmentation)
  - **image-text similarity (alignment)**
  - ViT attention map as saliency
    - image-text cross-attention



\* Open-vocabulary Semantic Segmentation (OVSS)

\*\* Open-vocabulary Instance Segmentation (OVIS)



ZSSeg: A Simple Baseline for Open-Vocabulary Semantic Segmentation with Pre-trained Vision-language Model (ECCV 2022)

LSeg: Language-driven Semantic Segmentation (ICLR 2022)

MaskCLIP: Extract Free Dense Labels from CLIP (ECCV 2022)

## 2. Related Work [3/3]

# Part Segmentation

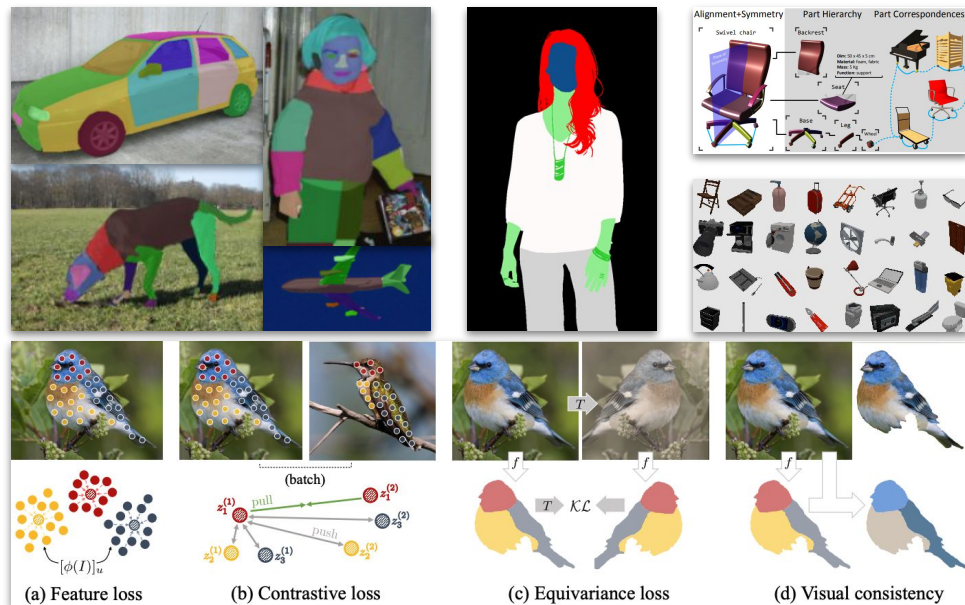
- Datasets

- CUB-200-2011 (2011)
- PASCAL-Part (CVPR 2014)
- DeepFashion (CVPR 2016)
- PartImageNet (ECCV 2022)
- OV-PARTS (NeurIPS B&D 2023)
- SubPartImageNet (ECCV 2024)
- ShapeNet (2015) / PartNet-Mobility (CVPR 2019)

- Tasks ( $\approx$ )

- Part Segmentation
- Fine-grained Segmentation
- Part Discovery (unsupervised, clustering)
- 3D Part Segmentation

- → fine-grained understanding



DeepFashion: Powering Robust Clothes Recognition and Retrieval with Rich Annotations (CVPR 2016)

PASCAL-Part: Detect What You Can: Detecting and Representing Objects using Holistic Models and Body Parts (CVPR 2014)

PartImageNet: A Large, High-Quality Dataset of Parts (ECCV 2022)

PACO: Parts and Attributes of Common Objects (CVPR 2023)

OV-PARTS: Towards Open-Vocabulary Part Segmentation (NeurIPS Datasets and Benchmarks Track 2023)

SPIN: Hierarchical Segmentation with Subpart Granularity in Natural Images (ECCV 2024)

Unsupervised Part Discovery from Contrastive Reconstruction (NeurIPS 2021)

Compositor: Bottom-up Clustering and Compositing for Robust Part and Object Segmentation (CVPR 2023)



### 3. Method

## PartCLIPSeg [1/1]

