# Supplementary Material: Open-vocabulary Semantic Segmentation with Frozen Vision-Language Models

Chaofan Ma[*1]
chaofanma@sjtu.edu.cn

Yuhuan Yang[*1]
yangyuhuan@sjtu.edu.cn

Yanfeng Wang[†1,2]
wangyanfeng@sjtu.edu.cn

Ya Zhang[1,2]
ya_zhang@sjtu.edu.cn

Weidi Xie[1,2]
weidi@sjtu.edu.cn

[1] Coop. Medianet Innovation Center, Shanghai Jiao Tong University, China

[2] Shanghai AI Laboratory

# Contents

In this supplementary material, we start by giving a detailed description about our experimental setup, including datasets split, and prompt templates. The results for ablation study on the architecture of cross-modality fusion module are then presented. Lastly, more qualitative visualization results are displayed.

# 1 Datasets Settings

In Table 1, we provide the detailed datasets split settings used in our experiments. Datasets PASCAL-$5^i$ [1, 2] and COCO-$20^i$ [4] follow the split settings proposed in LSeg [4]. COCO-$20^i$ $\rightarrow$ PASCAL VOC refers to the transferability experiment mentioned in Section 4.6 of the main paper.

| Dataset | fold 0 | fold 1 | fold 2 | fold 3 |
|---|---|---|---|---|
| PASCAL-$5^i$ | aeroplane, bicycle, bird, boat, bottle | bus, car, cat, chair, cow | diningtable, dog, horse, motorbike, person | pottedplant, sheep, sofa, train, tv/monitor |
| COCO-$20^i$ | person, aeroplane, boat, parkingmeter, dog, elephant, backpack, suitcase, sportsball, skateboard, wineglass, spoon, sandwich, hotdog, chair, diningtable, mouse, microwave, refrigerator, scissors | bicycle, bus, trafficlight, bench, horse, bear, umbrella, frisbee, kite, surfboard, cup, bowl, orange, pizza, sofa, toilet, remote, oven, book, teddybear | car, train, firehydrant, bird, sheep, zebra, handbag, skis, baseballbat, tennisracket, fork, banana, broccoli, donut, pottedplant, tvmonitor, keyboard, toaster, clock, hairdrier | motorbike, truck, stopsign, cat, cow, giraffe, tie, snowboard, baseballglove, bottle, knife, apple, carrot, cake, bed, laptop, cellphone, sink, vase, toothbrush |
| COCO-$20^i$ $\rightarrow$ PASCAL VOC | aeroplane, boat, chair, diningtable, dog, person | bicycle, bus, horse, sofa | bird, car, pottedplant, sheep, train, tvmonitor | bottle, cat, cow, motorbike |

**Table 1.** Data split for PASCAL-$5^i$, COCO-$20^i$, and COCO-$20^i$ $\rightarrow$ PASCAL VOC. PASCAL-$5^i$ and COCO-$20^i$ use 4-fold cross-validation, and categories in each row are for test and the rest classes are used for training. In COCO-$20^i$ $\rightarrow$ PASCAL VOC, categories in fold $i$ are novel classes in PASCAL VOC after removing the seen classes in corresponding training split on fold $i$ of COCO-$20^i$.

## 2 Prompt Templates

Before feeding the semantic category into text encoder, we use multiple prompt templates as decorations to generate text embeddings for one category, and then ensemble these embeddings by averaging. The following are the prompt templates we used:

```
'a bad photo of the {category}.',
'a photo of the large {category}.',
'a photo of the small {category}.',
'a cropped photo of a {category}.',
'This is a photo of a {category}',
'This is a photo of a small {category}',
'This is a photo of a medium {category}',
'This is a photo of a large {category}',
'This is a masked photo of a {category}',
'This is a masked photo of a small {category}',
'This is a masked photo of a medium {category}',
'This is a masked photo of a large {category}',
'This is a cropped photo of a {category}',
'This is a cropped photo of a small {category}',
'This is a cropped photo of a medium {category}',
'This is a cropped photo of a large {category}',
'A photo of a {category} in the scene',
'a bad photo of the {category} in the scene',
'a photo of the large {category} in the scene',
'a photo of the small {category} in the scene',
'a cropped photo of a {category} in the scene',
'a photo of a masked {category} in the scene',
'There is a {category} in the scene',
'There is the {category} in the scene',
'This is a {category} in the scene',
'This is the {category} in the scene',
'This is one {category} in the scene',
'There is a masked {category} in the scene',
'There is the masked {category} in the scene',
'This is a masked {category} in the scene',
'This is the masked {category} in the scene',
'This is one masked {category} in the scene',
```

# 3    Ablations of Cross-modality Fusion Module

We conduct an ablation study on the architecture of the cross-modality fusion module, *i.e.*, the number of layers, heads, and width of the transformer. The width/head ratio is set to 64 in this ablation. In Table 2, we can see that, a lightweight fusion module (1-layer transformer) already works reasonably well, and a steady improvement can be further observed as the number of the transformer layer increases from 1 to 12. However, performance begins to drop when the number of layers exceeds 12. Likewise, with a fixed number of layer, the performance variation on the widths (or heads) of the model follows the same scheme. Based on this result, we set layers to 12, and heads to 8 in all experiments in our main paper.
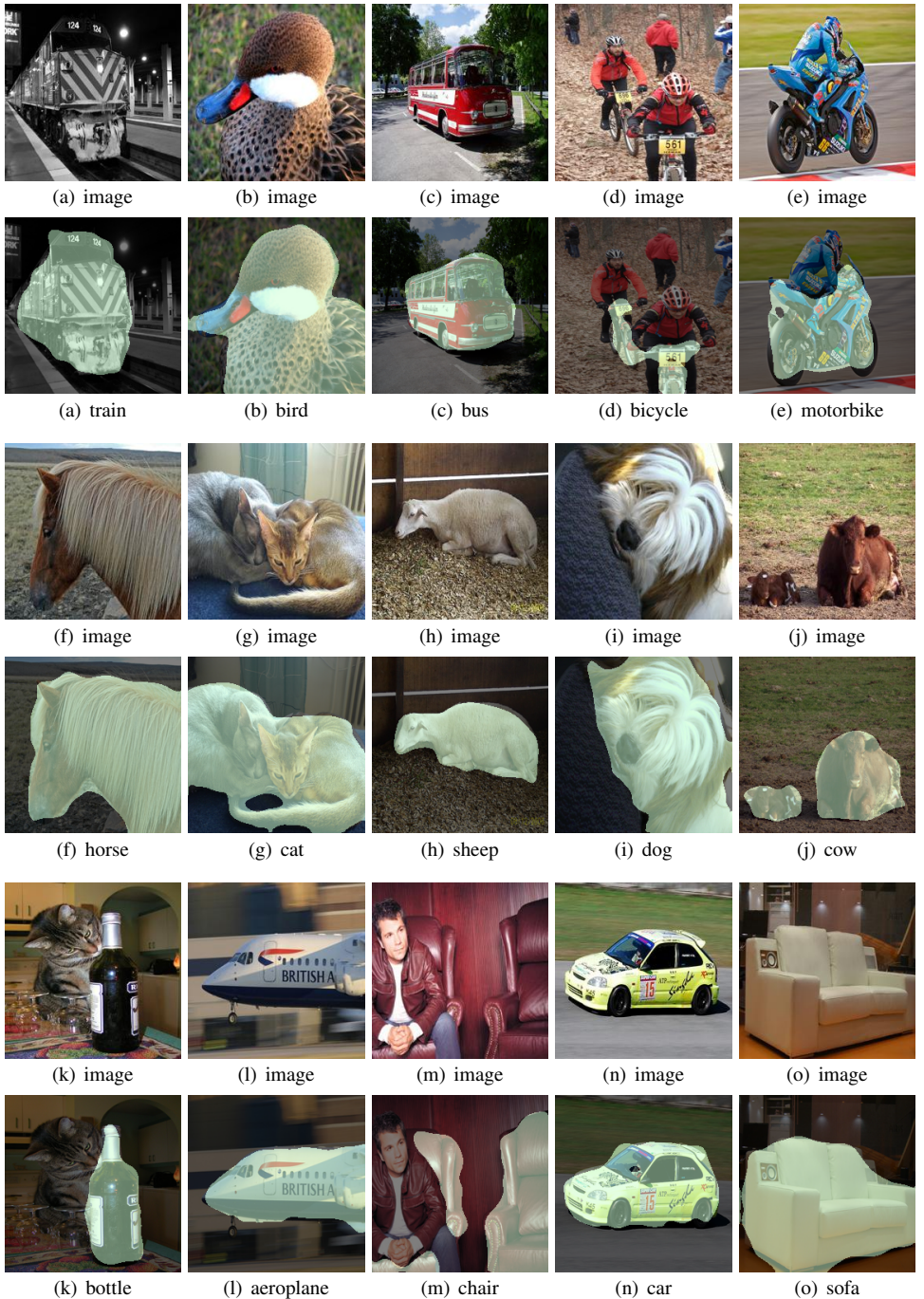
| layers | heads | width | $5^0$ | $5^1$ | $5^2$ | $5^3$ | mIoU |
|---|---|---|---|---|---|---|---|
| 1 | 8 | 512 | 52.2 | 59.2 | 45.0 | 44.6 | 50.2 |
| 3 | 8 | 512 | 57.1 | 68.0 | 46.7 | 50.0 | 55.4 |
| 6 | 8 | 512 | 58.9 | 68.9 | **51.3** | 51.5 | 57.7 |
| _12_ | _8_ | _512_ | **61.9** | 70.0 | 51.2 | **52.7** | **59.0** |
| 18 | 8 | 512 | 57.0 | 69.5 | 43.5 | 49.4 | 54.8 |
| 24 | 8 | 512 | 51.1 | **71.7** | 48.7 | 47.2 | 54.7 |
| 12 | 1 | 64 | 54.7 | 56.3 | 41.4 | 43.8 | 49.1 |
| 12 | 4 | 256 | 60.9 | 65.9 | 43.2 | 45.9 | 54.0 |
| _12_ | _8_ | _512_ | **61.9** | 70.0 | **51.2** | **52.7** | **59.0** |
| 12 | 10 | 640 | 54.1 | **70.5** | 49.1 | 47.7 | 55.3 |
| 12 | 12 | 768 | 55.5 | 65.8 | 47.8 | 46.8 | 54.0 |
| 12 | 16 | 1024 | 51.7 | 66.4 | 49.9 | 49.7 | 54.4 |

**Table 2.** Performance on layers, heads, and width of cross-modality fusion module on PASCAL-$5^i$. Underlined are the best parameters we used in the main paper.
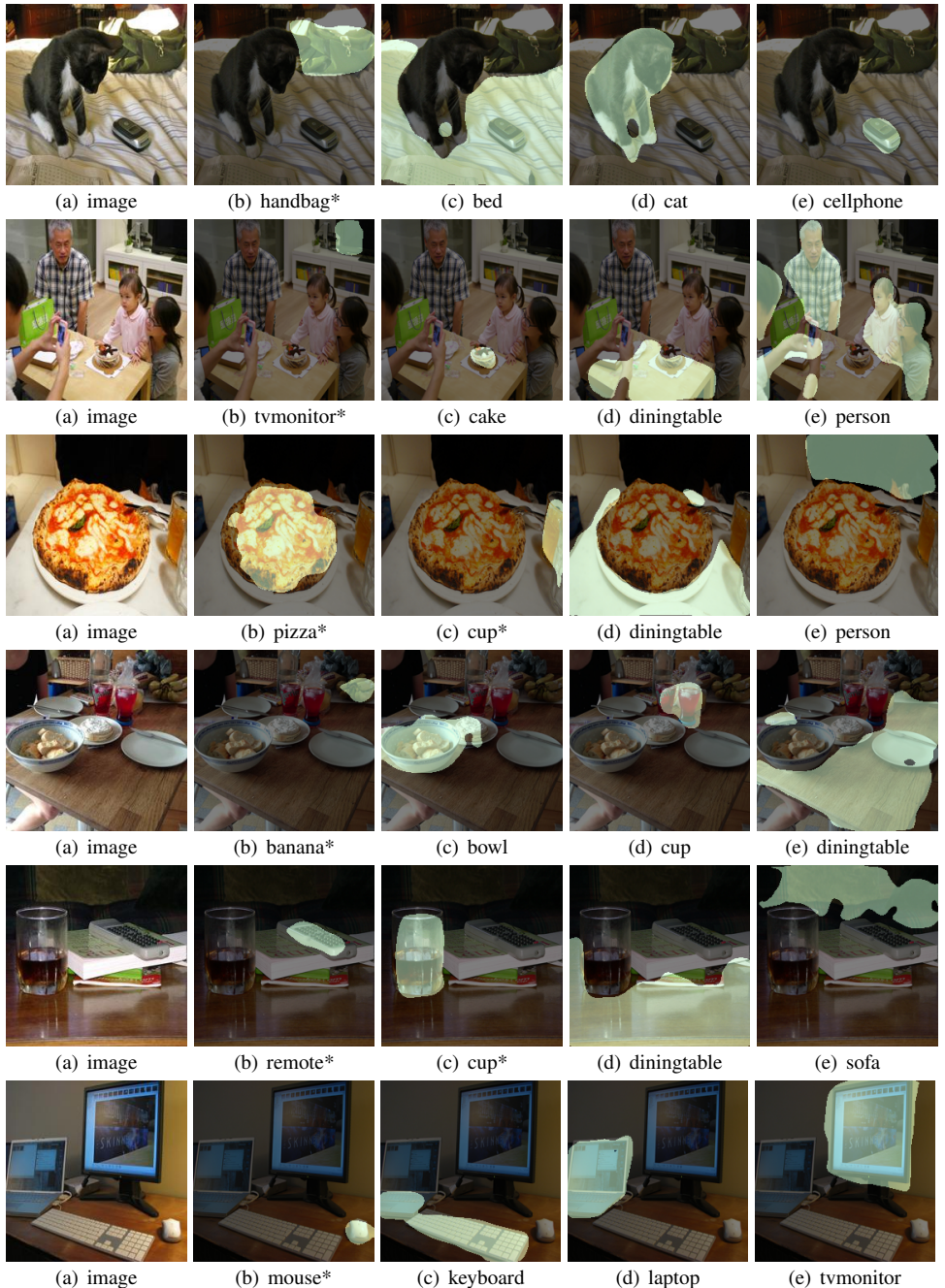
# 4    More Qualitative Visualization

In this section, we present additional qualitative results on PASCAL-$5^i$ and COCO-$20^i$ using our model with ViT-L backbone. Specifically, Figure 1 shows the results on PASCAL-$5^i$. All categories are novel (unseen) in their corresponding fold. Taking into account the variety of images, we choose 15 different categories to display. Figure 1(d) bicycle, Figure 1(e) motorbike, Figure 1(m) chair, and Figure 1(k) bottle show **Fusioner**'s ability to distinguish object of target semantic from other objects (distractor), *e.g.*, the person in Figure 1(e), and the cat in Figure 1(g). Besides, in Figure 1(d) bicycle, Figure 1(j) cow, Figure 1(m) chair, and Figure 1(g) cat, the model segments precisely even when the target instance contains more than one.

The visualization of COCO-$20^i$ is shown in Figure 2, with both seen and unseen categories are displayed. Facing a more noisy and complex scene, **Fusioner** is still able to recognize the desired (unseen) categories that are small and cornered, for example, tvmonitor, banana, remote and mouse in column of Figure 2(b).

**Figure 1.** Qualitative results on PASCAL5$^i$. All categories are novel (unseen) in their corresponding folds, and we show 15 different categories. For each two rows, top: images; bottom: segmentation masks for unseen categories.

**Figure 2.** Qualitative results on COCO-20$^i$. (a) input images, (b)-(e) segmentation masks for different categories. Unseen categories are marked as *.

# References

[1] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88(2):303–338, 2010.

[2] Bharath Hariharan, Pablo Arbeláez, Lubomir Bourdev, Subhransu Maji, and Jitendra Malik. Semantic contours from inverse detectors. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 991–998. IEEE, 2011.

[3] Boyi Li, Kilian Q Weinberger, Serge Belongie, Vladlen Koltun, and Rene Ranftl. Language-driven semantic segmentation. In *International Conference on Learning Representations*, 2022.

[4] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*, pages 740–755. Springer, 2014.