

---

# Supplementary Materials

---

Anonymous Author(s)  
Affiliation  
Address  
email

1	<b>Contents</b>	
2	<b>1 Computational Complexity of NPR</b>	<b>1</b>
3	<b>2 Additional Experiments and Analysis</b>	<b>1</b>
4	2.1 Dataset and Models . . . . .	1
5	2.2 VOC results . . . . .	2
6	2.3 COCO results . . . . .	2
7	2.4 Context results . . . . .	2
8	2.5 LVIS & ImageNet-S results . . . . .	2
9	2.6 Two benefits . . . . .	2
10	<b>3 Broader Impacts</b>	<b>3</b>
11	<b>4 Limitations</b>	<b>3</b>

## 12 **1 Computational Complexity of NPR**

13 Recall that we claim the computational complexity of the prototype generation process in NPR is  
14  $\mathcal{O}(q \times m \times d)$ , which is acceptable to implement. Furthermore, the experimental results in Section 5.3  
15 also show that our proposed PGSeg indicates a reasonable level of computational efficiency. Here we  
16 aim to make a detailed deduction on this. The computational complexity of the prototype generation  
17 process is led by the EM process in Eq.(1) and (2), which could be divided into four stages:

- 18 1. The first stage could be regarded as a matrix multiplication operation between the prototypes  
19  $\mathbf{P} \in \mathbb{R}^{q \times d}$  and prototypical source feature  $\mathbf{V} \in \mathbb{R}^{m \times d}$ , and the computational complexity for that  
20 is  $\mathcal{O}(q \times m \times d)$ .
- 21 2. The second stage is the normalization process to obtain the estimated  $\mathbf{Y}$ , which could be treated as  
22 a SoftMax operation as presented in Eq.(1), and the complexity of such this operation is  $\mathcal{O}(q \times m)$ .
- 23 3. The third stage is to calculate the inner product between  $y_{ij}$  and  $v_j$ , namely  $\mathbf{YV}^\top$ , resulting in  
24 the complexity of  $\mathcal{O}(q \times m \times d)$ . Based on an extra normalization operation, the final complexity  
25 reaches  $\mathcal{O}(q \times m \times d) + \mathcal{O}(q \times m)$
- 26 4. The final step is to reconstruct the prototypical source by multiplying  $\mathbf{Y}^\top$  with  $\mathbf{P}$ , which is  
27 specifically used in I-NPR, leading to the computational complexity of  $\mathcal{O}(q \times m \times d)$ .

28 Overall, the total computational complexity of the prototype generation process is the summation of  
29 the results above, i.e.,  $T(\mathcal{O}(q \times m \times d) + \mathcal{O}(q \times m) + \mathcal{O}(q \times m \times d) + \mathcal{O}(q \times m) + \mathcal{O}(q \times m \times d))$ ,  
30 where  $T$  is the total iterations. Consequently, the estimated computational complexity is  $\mathcal{O}(q \times m \times d)$   
31 if neglecting the higher-order terms and constant factors, which concludes the proof.

## 32 **2 Additional Experiments and Analysis**

### 33 **2.1 Dataset and Models**

34 **Datasets.** In Section 5, we evaluate our PGSeg on five prevalent benchmarks, which are PASCAL  
35 VOC12 2012 [6], COCO [9], PASCAL Context [10], ImageNet-S [7], and LVIS [8]. Here is the  
36 detailed introduction of these five datasets as follows:

- 37 • **PASCAL VOC2012 [6]:** The PASCAL VOC12 dataset consists of a diverse collection of images  
38 spanning 21 different object categories (including one background class), such as a person, car,  
39 dog, and chair. The dataset provides annotations for both training and validation sets, with around  
40 1,464 images in the training set and 1,449 images in the validation set. We use the validation set for  
41 the downstream evaluation. During the inference, we set the background score as 0.95.
- 42 • **COCO [9]:** The COCO Object dataset covers a wide range of 80 object categories, such as cars,  
43 bicycles, people, animals, and household items. For semantic segmentation, it has 118,287 training  
44 images and 5,000 images for validation. Correspondingly, we merely use the validation set, and the  
45 background score is 0.85.
- 46 • **Context [10]:** The dataset contains a diverse set of images taken from various scenes, including  
47 indoor and outdoor environments. It covers 59 common object classes, such as a person, car,  
48 bicycle, and tree, as well as 60 additional stuff classes, including sky, road, grass, and water. It  
49 has 118,287 training images and 5,000 images for validation. Here we merely consider the object  
50 dataset part, and use the validation set. The background score during the inference is set as 0.36.
- 51 • **ImageNet-S [7]:** ImageNet-S, distilled from the ImageNet [4], is a human-annotated pixel-level  
52 dataset specifically used for semantic segmentation. ImageNet-S has three versions based on the  
53 category amount, and we use the version with the maximum number of classes, which contains 919  
54 classes and 12,419 validation samples. For a simple validation, we reduce the number of validation  
55 samples to 5,000. The background score is set to 0.11.
- 56 • **LVIS [8]:** The LVIS (Large Vocabulary Instance Segmentation) dataset is a large-scale dataset  
57 specifically constructed from COCO. It focuses on instance-level segmentation, where the goal is to  
58 identify and segment individual objects within an image. It includes a comprehensive vocabulary of  
59 over 1,200 object categories, making it one of the largest instance segmentation datasets available.  
60 It has 5,000 samples for evaluation. We ignore the instance-level annotation in the samples to  
61 proceed with semantic segmentation. The background score is set to 0.35.

62 **Models.** The PGSeg comprises an image encoder and a text encoder. The text encoder follows  
63 from [12] and consists of 12 transformer layers, each with a hidden dimension of 256. The text  
64 encoder adopts a *lower-cased byte pair encoding* (BPE) to encode the text with a vocabulary of  
65 49,512 words. For the image encoder, we turn to the ViT-S with 12 transformer layers, each having a  
66 multi-head (6) self-attention and an MLP. We use layer normalization [1] to the input of each PG  
67 Unit. 3 Transformer layers are added to the final output of the PG Unit. In T-NPR, we select the text  
68 embedding before the mapping MLP, used to align with image embedding, as the prototypical source.  
69 In I-NPR, we use the input token, fed before the transformer layers, after a layer normalization as  
70 the image-level prototypical source. During the training stage, any pre-trained model is not used  
71 for both the image and text encoder. Several data augmentation approaches are employed, such as  
72 Random Flip and Color Normalization. During the inference, we directly adopt the model without  
73 any fine-tuning or training. Besides, we strictly follow [16, 13, 17] to set the input image size as  
74  $448 \times 448$ , and adopt a stride strategy to generate the segmentation mask. The whole implementation  
75 is built on PyTorch [11] and MMSegmentation [3].

## 76 2.2 VOC results

77 Figure 1 presents more results of our PGSeg in VOC12. It is found that our PGSeg shows powerful  
78 grouping capability when segmenting the object-centric images. Besides, the learned group tokens  
79 could help segment objects in a compact and dense manner, which means there is less redundancy  
80 and noise in objects.

## 81 2.3 COCO results

82 Figure 2 presents some visualized results of COCO Object. Clearly, it has been observed that our  
83 PGSeg is able to perform fine-grained segmentation in the multi-object case. However, PGSeg is  
84 unable to completely capture some small objects, such as the bottle and plate in the image of the  
85 fourth row. This is essentially due to the wrongly-segmented group tokens, leading to some noise.

## 86 2.4 Context results

87 Figure 3 shows several visualized results of Context. Compared to COCO, the results of Context is  
88 comparably promising: the group tokens are distinctive from each other, capturing the whole object  
89 in a complete manner. Nevertheless, the semantic output seems to be not satisfying. On its face, such  
90 an issue is caused by a fix score of the background. Therefore, some areas in the group tokens with  
91 low scores, are directly recognized as the background. The dog in the third row could be an intuitive  
92 illustration. Besides, the over-segmentation phenomenon is quite severe in the multi-objects case,  
93 which is also a huge challenge in WOVSS.

## 94 2.5 LVIS & ImageNet-S results

95 Figure 4 shows the visualized results of both LVIS and ImageNet-S. In ImageNet-S, we find that  
96 the group token in our PGSeg is powerful to finely cluster the object-centric object. However, due  
97 to the limited vocabulary during the training stage, our PGSeg is unable to match the input text  
98 with the corresponding semantic groups. In LVIS, it has been observed that over-segmentation and  
99 the presence of noisy regions are inherent issues that have emerged, similar to those observed in  
100 the Context dataset. Besides, the confidence scores of some complex group areas are still not high,  
101 leading to an object-level under-segmentation.

## 102 2.6 Two benefits

103 Here we incorporate more results on the analysis of two benefits, i.e., *compactness* and *richness*, to  
104 better understand NPR. Figure 5 illustrates the visualized t-SNE results of input patch tokens. The  
105 label IDs for these tokens are provided by the group token. Upon observation, it is evident that in  
106 comparison to GroupViT, the group tokens in PGSeg demonstrate a better ability to form a compact  
107 foundation. This aids in densely clustering the input patch tokens while being free from noise, thereby  
108 reducing redundancy. For *richness*, Figure 6 reports the dimensional distribution of the group token  
109 in level 2, which has 8 group tokens in total. Clearly, the dimensional variance of our PGSeg is larger  
110 than that of GroupViT, leading to a diverse feature representation.

111 To investigate the effects of each NPR strategy in terms of visual and textual parts, we have added  
112 additional analysis to examine the influence of I-NPR and T-NPR on compactness and richness  
113 in Figure 7 and 8. We find that both components contribute to compactness and richness, but  
114 As shown in Figure 7, I-NPR exerts a more pronounced impact on amplifying compactness. This  
115 observation aligns with the intuitive understanding that the image information embedded within I-NPR  
116 is inherently structured, thereby leading to more cohesive clustering. As shown in Figure 8, T-NPR  
117 has a more significant effect on improving richness, since it leads to a larger dimensional variance  
118 compared to I-NPR. We conjecture that the prototypes, which are sourced from text embeddings,  
119 could impose stronger semantic regularization on the group tokens. The results also underscore the  
120 complementary roles of I-NPR and T-NPR, further substantiating their importance in PGSeg.

### 121 **3 Broader Impacts**

122 Note that our training datasets, CC12M [2] and RedCaps12M [5], are sourced from the Internet.  
123 Consequently, the collection of these datasets raises concerns regarding privacy if not appropriately  
124 regulated. Additionally, text supervision typically relies on human annotators, which can introduce  
125 biases, intentional or unintentional, if the annotators are not impartial. It is key to address these  
126 issues through proper data regulation, privacy protection measures, and meticulous selection on the  
127 annotated information to ensure fairness and relieve potential biases.

### 128 **4 Limitations**

129 While our proposed PGSeg model can be applied to segment various downstream datasets, it shows  
130 relatively poor performance compared to methods that use additional supervision, such as segmenting  
131 masks. Particularly, PGSeg exhibits subpar performance in datasets like LVIS and ImageNet-S,  
132 indicating its limited capability for fine-grained segmentation in real-world scenarios. Moreover,  
133 our model is trained from scratch using training datasets that primarily consist of common object  
134 categories like dogs, people, etc. As a result, its applicability may be limited in specific domains such  
135 as medical [15] and LiDAR [14] images. Therefore, further investigation is warranted to assess the  
136 segmenting ability and potential application scope of our model in the future.

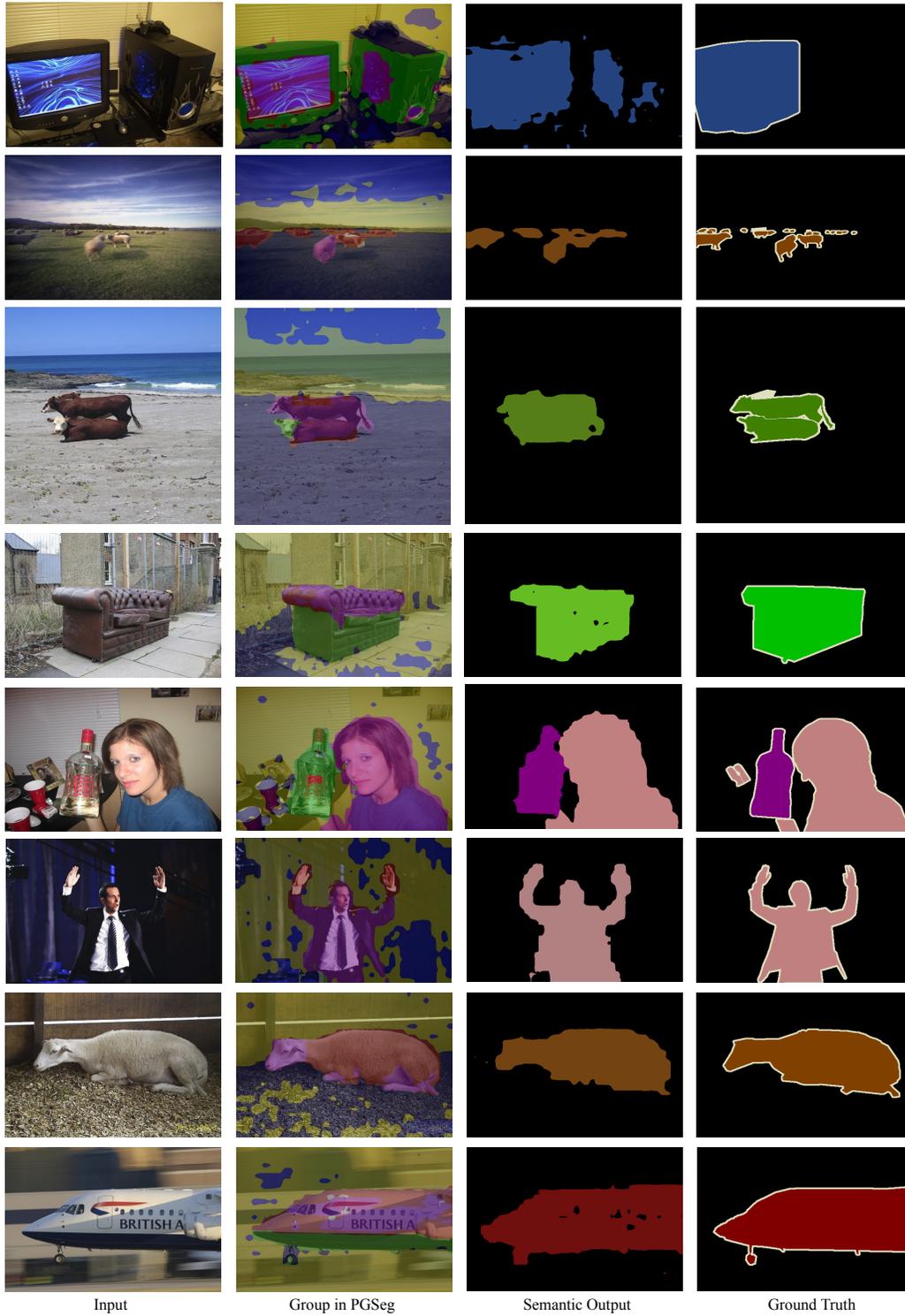


Figure 1: Qualitative results on PASCAL VOC12.

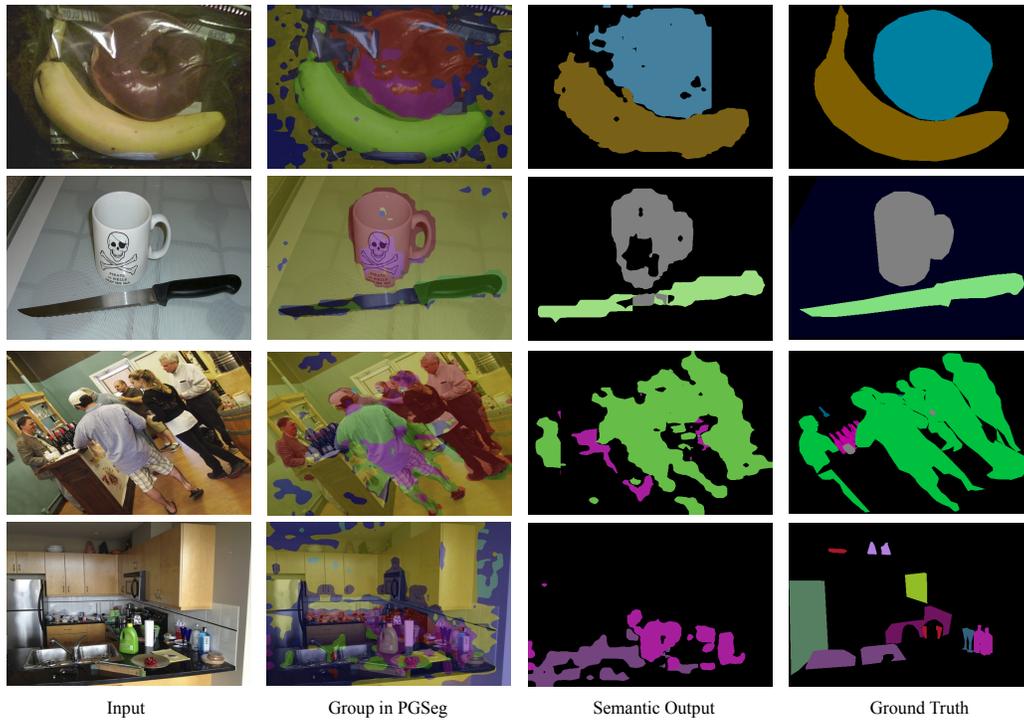


Figure 2: Qualitative results on COCO Object.

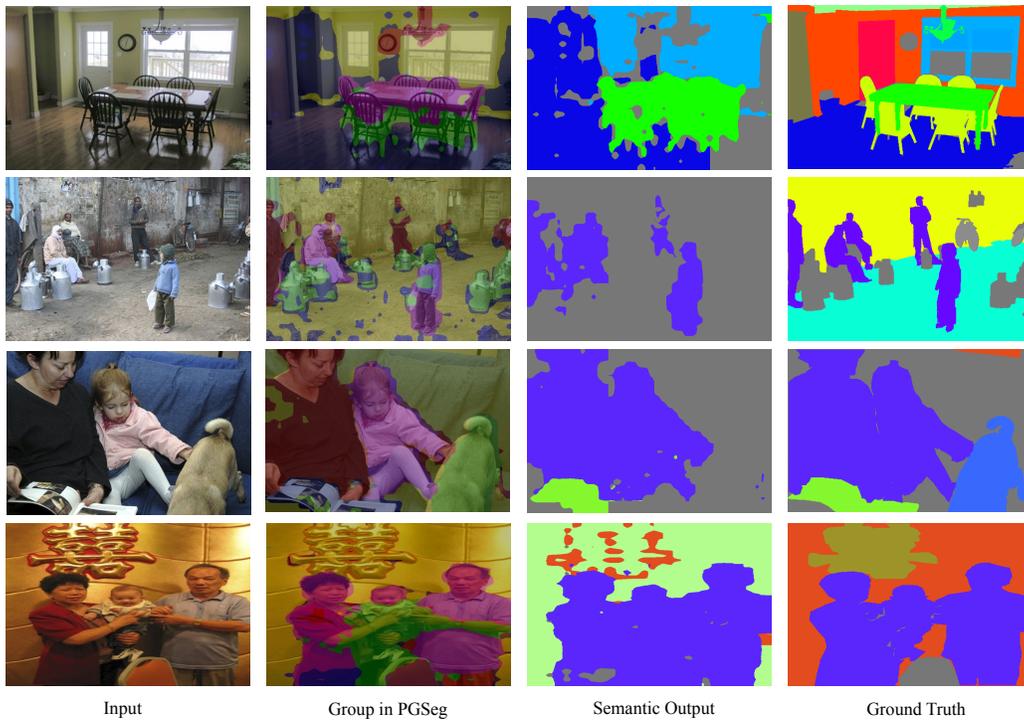


Figure 3: Qualitative results on Context.

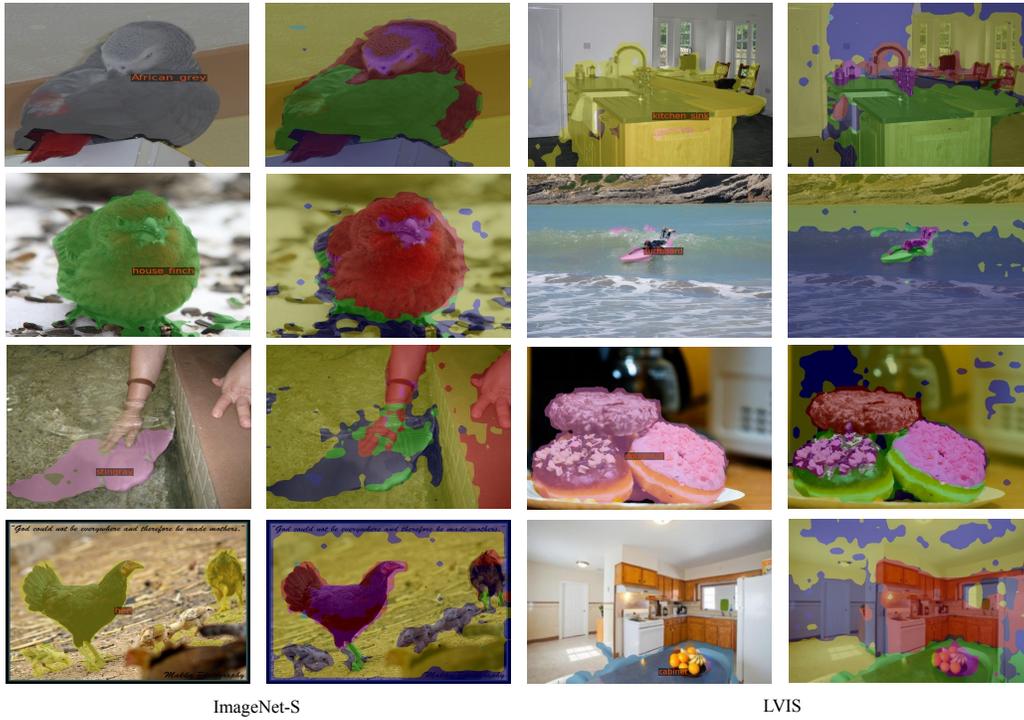


Figure 4: Qualitative results on ImageNet-S and LVIS.

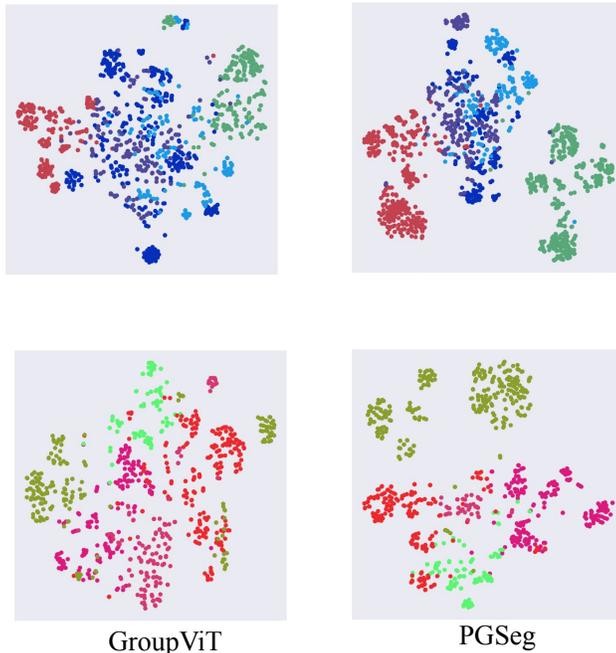


Figure 5: Compactness analysis. The results come from 5 clustered patch tokens based on the group tokens.

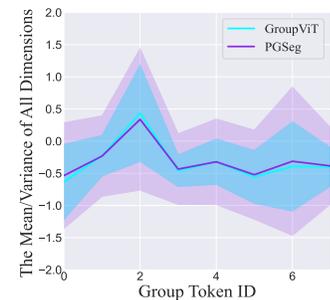


Figure 6: Dimension distributions of each group token.

137 **References**

138 [1] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint*  
 139 *arXiv:1607.06450*, 2016.

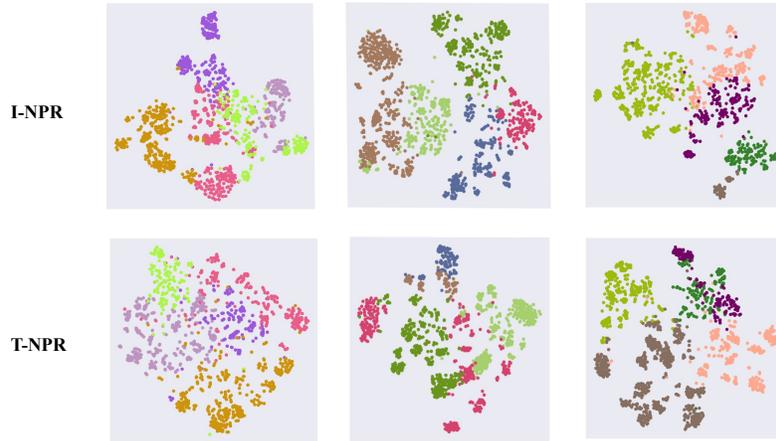


Figure 7: Comparison of *compactness* between I-NPR and T-NPR. I-NPR exerts a more pronounced impact on amplifying compactness, leading to more cohesive clustering.

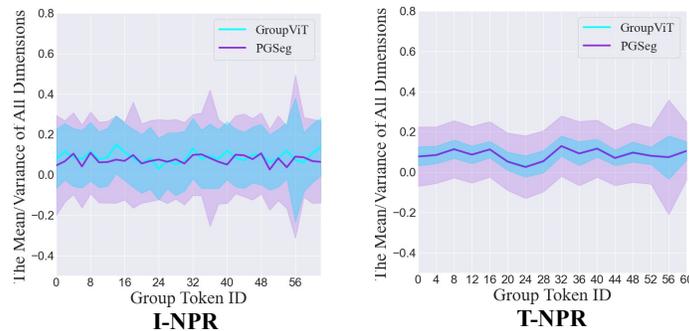


Figure 8: Comparison of *richness* between I-NPR and T-NPR. T-NPR has a more significant effect on improving richness, leading to a larger dimensional variance compared to I-NPR.

- 140 [2] Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12m: Pushing web-scale  
 141 image-text pre-training to recognize long-tail visual concepts. In *Proceedings of the IEEE/CVF Conference*  
 142 *on Computer Vision and Pattern Recognition*, pages 3558–3568, 2021.
- 143 [3] MMSegmentation Contributors. Mmsegmentation: Openmmlab semantic segmentation toolbox and  
 144 benchmark. <https://github.com/open-mmlab/>, 2020.
- 145 [4] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical  
 146 image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255.  
 147 Ieee, 2009.
- 148 [5] Karan Desai, Gaurav Kaul, Zubin Aysola, and Justin Johnson. Redcaps: Web-curated image-text data  
 149 created by the people, for the people. *arXiv preprint arXiv:2111.11431*, 2021.
- 150 [6] Mark Everingham, SM Ali Eslami, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew  
 151 Zisserman. The pascal visual object classes challenge: A retrospective. *International journal of computer*  
 152 *vision*, 111:98–136, 2015.
- 153 [7] Shanghua Gao, Zhong-Yu Li, Ming-Hsuan Yang, Ming-Ming Cheng, Junwei Han, and Philip Torr.  
 154 Large-scale unsupervised semantic segmentation. *IEEE Transactions on Pattern Analysis and Machine*  
 155 *Intelligence*, 2022.
- 156 [8] Agrim Gupta, Piotr Dollar, and Ross Girshick. Lvis: A dataset for large vocabulary instance segmentation.  
 157 In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5356–5364,  
 158 2019.
- 159 [9] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár,  
 160 and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014:*  
 161 *13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pages  
 162 740–755. Springer, 2014.

- 163 [10] Roozbeh Mottaghi, Xianjie Chen, Xiaobai Liu, Nam-Gyu Cho, Seong-Whan Lee, Sanja Fidler, Raquel  
164 Urtasun, and Alan Yuille. The role of context for object detection and semantic segmentation in the wild.  
165 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 891–898, 2014.
- 166 [11] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen,  
167 Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep  
168 learning library. *Advances in neural information processing systems*, 32, 2019.
- 169 [12] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish  
170 Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from  
171 natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR,  
172 2021.
- 173 [13] Pengzhen Ren, Changlin Li, Hang Xu, Yi Zhu, Guangrun Wang, Jianzhuang Liu, Xiaojun Chang, and  
174 Xiaodan Liang. Viewco: Discovering text-supervised segmentation masks via multi-view semantic  
175 consistency. *arXiv preprint arXiv:2302.10307*, 2023.
- 176 [14] Santiago Royo and Maria Ballesta-Garcia. An overview of lidar imaging systems for autonomous vehicles.  
177 *Applied sciences*, 9(19):4093, 2019.
- 178 [15] Dinggang Shen, Guorong Wu, and Heung-Il Suk. Deep learning in medical image analysis. *Annual review*  
179 *of biomedical engineering*, 19:221–248, 2017.
- 180 [16] Jiarui Xu, Shalini De Mello, Sifei Liu, Wonmin Byeon, Thomas Breuel, Jan Kautz, and Xiaolong Wang.  
181 Groupvit: Semantic segmentation emerges from text supervision. In *Proceedings of the IEEE/CVF*  
182 *Conference on Computer Vision and Pattern Recognition*, pages 18134–18144, 2022.
- 183 [17] Jilan Xu, Junlin Hou, Yuejie Zhang, Rui Feng, Yi Wang, Yu Qiao, and Weidi Xie. Learning open-vocabulary  
184 semantic segmentation models from natural language supervision. *arXiv preprint arXiv:2301.09121*, 2023.