

039

040

041

Dictionary-based Framework for Interpretable and Consistent Object Parsing

Anonymous CVPRW submission

Paper ID *****

Abstract

001 In this work, we present CoCal, an interpretable and consistent object parsing framework based on dictionary-002 003 based mask transformer. Designed around Contrastive Components and Logical Constraints, CoCal rethinks ex-004 isting cluster-based mask transformer architectures used 005 in segmentation; Specifically, CoCal utilizes a set of dic-006 007 tionary components, with each component being explicitly 008 linked to a specific semantic class. To advance this concept, CoCal introduces a hierarchical formulation of dic-009 tionary components that aligns with the semantic hierarchy. 010 This is achieved through the integration of both within-level 011 012 contrastive components and cross-level logical constraints. 013 Concretely, CoCal employs a component-wise contrastive algorithm at each semantic level, enabling the contrast-014 ing of dictionary components within the same class against 015 those from different classes. Furthermore, CoCal addresses 016 logical concerns by ensuring that the dictionary compo-017 018 nent representing a particular part is closer to its corresponding object component than to those of other objects 019 through a cross-level contrastive learning objective. To fur-020 ther enhance our logical relation modeling, we implement 021 a post-processing function inspired by the principle that a 022 023 pixel assigned to a part should also be assigned to its corre-024 sponding object. With these innovations, CoCal establishes a new state-of-the-art performance on both PartImageNet 025 and Pascal-Part-108, outperforming previous methods by 026 a significant margin of 2.08% and 0.70% in part mIoU, re-027 spectively. Moreover, CoCal exhibits notable enhancements 028 029 in object-level metrics across these benchmarks, highlighting its capacity to not only refine parsing at a finer level but 030 also elevate the overall quality of object segmentation. 031

032 1. Introduction

Human perception involves the ability to decompose an
object into its semantically meaningful components (*i.e.*,
parts). For instance, when observing a dog, humans not
only identify it as a dog but also simultaneously discover
its head, torso, and other components, facilitating a more



: Part Components ↔: Logical Constraints ↔: Contrastive Components A: Object Components

Figure 1. Illustration of the proposed component-wise contrastive objectives. CoCal establishes two discriminative dictionaries at the part and object levels. Within the same semantic level, part/object components of the same classes are pulled closer $(\rightarrow \leftarrow)$, while those of different classes are pushed apart $(\leftarrow \rightarrow)$ (*i.e.*, contrastive components). At the cross-semantic level, part components and their corresponding object components are pulled closer and vice versa (*i.e.* logical constraints).

interpretable and resilient understanding of real-world scenarios. More specifically, humans can estimate the pose of a dog by considering the spatial arrangement of its parts, even in instances where some parts may be missing.

By contrast, emulating this innate human visual capa-042 bility presents a big challenge for modern computer vision 043 models. The predominant focus within the field has been 044 on addressing semantic segmentation at the object level, 045 with minimal attention given to intermediate part represen-046 tations. Notable works [15, 34, 48, 49] in object parsing 047 primarily extend algorithms designed for general segmenta-048 tion, overlooking the fact that parts, being at a lower seman-049

108

109

110

111

112

123

124

125

26

127

128

129

130

131

132

133

134

135

136

137

138

139

tic level, can be captured more efficiently and interpretably 050 through clustering. As a result, these works often adhere 051 052 to frameworks tailored for object segmentation without incorporating specialized designs for handling parts. More-053 054 over, even though certain studies [18, 22, 61] highlight the mutual benefit between object parsing and object segmen-055 tation, they typically treat these semantic levels separately. 056 disregarding the logical relationship between them. Conse-057 058 quently, the optimization objectives for these two levels are disjoint, leading to sub-optimal predictions. 059

060 In this work, we propose CoCal, a dictionary-based framework built on top of an off-the-shelf cluster-based 061 062 mask transformer, utilizing a set of dictionary components where each component is explicitly associated with a spe-063 cific semantic class to facilitate the grouping of pixels be-064 065 longing to that class. This enables CoCal to conduct inference in a straightforward parameter-free manner through 066 nearest neighbor search on the pixel feature maps within 067 the class dictionary. Taking this concept further, CoCal 068 069 introduces a hierarchical formulation of dictionary compo-070 nents, aligning with the semantic hierarchy, which naturally forms the logical paths within the structure (e.g., bird-head 071 \rightarrow bird). CoCal advances the learning of the above for-072 mulation through two simple yet effective targets: learn-073 ing contrastive objectives for obtaining discriminative dic-074 075 tionary components and exploring logical relations for consistent predictions. Specifically, as depicted in Fig. 1, at 076 each semantic level, CoCal employs a component-wise con-077 trastive algorithm to pull closer the dictionary components 078 079 withing the same class while pushing away those from dif-080 ferent classes. Then to model the cross-level logical rela-081 tions, CoCal further contrasts the positive pair between dictionary component representing a particular part and its cor-082 responding object dictionary components against the nega-083 tive pairs involving the part component and all other object 084 components. In addition, CoCal applies a post-processing 085 086 step enforcing that any pixel predicted as a certain part must also be labeled under its corresponding object class. Specif-087 ically, CoCal calculates a path probability by multiplying 088 part- and object-level similarities, then assigns pixels to the 089 path with the highest score. This mechanism captures cross-090 091 level semantics and resolves inconsistencies at inference.

092 2. Method

093 In this section, we begin by introducing the core principles behind mask transformer segmentation frameworks, 094 095 focusing on cluster-based methods and showing how our proposed dictionary-based formulation evolves from these 096 ideas. We then describe how CoCal leverages hierarchy, 097 contrastive components, and logical constraints to achieve 098 interpretable and consistent object parsing. Finally, we 099 present the meta-architecture of our method, detailing all 100 101 major components and their interactions.

2.1. Dictionary-based Mask Transformers

Problem StatementSemantic segmentation aims to par-
tition an image $\mathbf{I} \in \mathbb{R}^{H \times W \times 3}$ into non-overlapping masks,
where each mask is associated with a semantic label. Con-
cretely, the ground-truth set of masks is:103103104104105105105

$$\{y_i\}_{i=1}^M = \{(d_i, c_i)\}_{i=1}^M,$$
(1) 107

where $d_i \in \{0, 1\}^{H \times W}$ is a binary mask denoting the pixels of the *i*-th region, c_i is its class label, and M represents the number of ground-truth masks. A typical segmentation model outputs a set of N predicted masks $(N \ge M)$ and their corresponding classes:

$$\{\hat{y}_i\}_{i=1}^N = \{(\hat{m}_i, \hat{c}_i)\}_{i=1}^N.$$
 (2) 113

Recap of Cluster-based Mask Transformers Cluster-114 based mask transformers [37, 71, 72] differ from standard 115 query-based transformers in their use of a one-hot argmax 116 assignment instead of a softmax for updating the query 117 features. Specifically, let $\mathbf{O} \in \mathbb{R}^{N \times D}$ denote the N ob-118 ject queries and $\hat{\mathbf{O}}$ the updated queries. Similarly, let 119 $\mathbf{Q}^{o}, \mathbf{K}^{p}, \mathbf{V}^{p}$ be the linearly projected features for queries, 120 keys, and values. Then, instead of the softmax cross-121 attention, 122

$$\hat{\mathbf{O}} = \mathbf{O} + \underset{HW}{\operatorname{softmax}} \left(\mathbf{Q}^{o} \times (\mathbf{K}^{p})^{\mathrm{T}} \right) \times \mathbf{V}^{p}, \qquad (3)$$

the cluster-based mask transformer adopts a one-hot argmax:

$$\hat{\mathbf{O}} = \mathbf{O} + \operatorname*{argmax}_{N} \left(\mathbf{Q}^{o} \times (\mathbf{K}^{p})^{\mathrm{T}} \right) \times \mathbf{V}^{p},$$
 (4) 1

so each pixel is clustered to the single query with which it has the highest affinity. The prediction set $\{\hat{y}_i\}_{i=1}^N$ is then matched with $\{y_i\}_{i=1}^M$ through Hungarian Matching [30], which guides mask and classification loss computation.

Dictionary-based Formulation While cluster-based methods typically employ N object queries (often larger than M), our dictionary-based approach seeks a more direct mapping between queries and classes. We replace the learnable queries **O** with a dictionary $\mathbf{C} \in \mathbb{R}^{P \times D}$, where each of the P dictionary components serves as a *cluster center* for a specific class. Hence, there is a one-to-one correspondence between each component \mathbf{C}_i and one of the P classes.

During training, the dictionary C updates in the same 140 spirit as Eqs. 2-4, except that we now have a *fixed* alignment 141 of dictionary elements to classes (so Hungarian Matching 142 is unnecessary). At inference time, the dictionary-based 143 mask transformer is fully parameter-free in its final assign-144 ment step, as each pixel's feature vector is classified by 145 whichever dictionary component is closest in feature space. 146 This streamlining removes the need for redundancies like 147 'void' labels and affords an inherently interpretable design: 148 each dictionary component explicitly encodes a particular 149 semantic class. 150

193

CVPRW 2025 Submission #*****. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.



Figure 2. **Meta-architecture of the proposed CoCal.** CoCal builds on top of an off-the-shelf clustering-based mask transformer, incorporating dictionary components that function as the cluster centers for each semantic class. Throughout training, the dictionary components in CoCal are updated via both mask-wise objectives from the transformer and contrastive objectives from the dictionary. During testing, CoCal adopts a straightforward inference approach by executing nearest neighbor search of the pixel features on the dictionary components.

151 2.2. CoCal: Interpretable and Consistent Object 152 Parsing

Although dictionary-based mask transformers already pro-153 154 vide a concise per-class representation, object parsing often 155 requires hierarchical reasoning (e.g., relating part classes to an object-level class) and advanced regularization to en-156 sure consistent outputs. CoCal addresses these issues with 157 three key ideas: (1) hierarchical dictionary components, (2) 158 contrastive learning of dictionary elements, and (3) logical 159 160 constraints at both training and inference time.

161 2.2.1. Hierarchical Structure of Dictionaries Across Mul 162 tiple Levels

Part labels naturally imply rich logical relationships (e.g., 163 dog-head is more semantically related to dog-torso than 164 fish-tail). To exploit this hierarchy, CoCal extends the dic-165 tionary with an additional tier for object-level classes. Con-166 cretely, we have a part-level dictionary $\mathbf{C} \in \mathbb{R}^{P \times D}$ for part 167 segmentation and an object-level dictionary $\widetilde{\mathbf{C}} \in \mathbb{R}^{\widetilde{P} \times D}$ 168 for object classes, where \widetilde{P} is the number of object classes. 169 170 These two dictionaries are learned jointly to reflect the inherent relationship between objects and their parts. 171

172 2.2.2. Enhancing Dictionary Discrimination via Con 173 trastive Objectives

To ensure discriminative power, CoCal applies contrastive learning on both part- and object-level dictionaries. Consider the part dictionary C. We maintain a *part memory bank* $\mathbf{B} \in \mathbb{R}^{P \times S \times D}$, where S is the number of stored "samples" per class. After retrieving the relevant components $\mathbf{C}(y)$ for classes present in a training sample, we compute the contrastive loss as:

$$\mathcal{L}_{n\,con}(\mathbf{C}(y)) =$$
182

$$\sum_{x \in M} \frac{-1}{|\mathbf{B}(x)|} \sum_{j \in \mathbf{B}(x)} \log \left(\frac{\exp(\mathbf{C}(y)_i \cdot B_j / \tau)}{\sum_{k \in \mathbf{B}} \exp(\mathbf{C}(y)_i \cdot B_k / \tau)} \right), \quad (5)$$

where $\mathbf{B}(x)$ are the stored features in \mathbf{B} belonging to the same class x, B_j is one such feature, and τ is the temperature. Inspired by [41, 52], we employ hard negative mining to focus on the most challenging negatives. An analogous memory bank $\widetilde{\mathbf{B}}$ and contrastive loss $\mathcal{L}_{o.con}$ are used for the object-level dictionary $\widetilde{\mathbf{C}}$: 189 190

$$\mathcal{L}_{o_con}(\widetilde{\mathbf{C}}(y)) =$$
 191

$$\sum_{x \in M} \frac{-1}{|\widetilde{\mathbf{B}}(x)|} \sum_{j \in \widetilde{\mathbf{B}}(x)} \log \left(\frac{\exp(\widetilde{\mathbf{C}}(y)_i \cdot \widetilde{B}_j / \tau)}{\sum_{k \in \widetilde{\mathbf{B}}} \exp(\widetilde{\mathbf{C}}(y)_i \cdot \widetilde{B}_k / \tau)} \right), \quad (6)$$

2.2.3. Logical Constraints for Consistent Predictions

Even with hierarchical dictionaries, inconsistent predictions 194 can arise-for example, labeling a pixel snake-head at the 195 part level but reptile at the object level. CoCal incorporates 196 two logical constraints to maintain cross-level consistency. 197 Cross-Level Contrastive Loss. First, parts belonging to 198 the same object should be closer to their object-level com-199 ponent than to other objects. We encode this by an addi-200 tional cross-level contrastive term: 201 202

$$\mathcal{L}_{logic}(\mathbf{C}(y)) =$$
 203

$$\sum_{x \in M} \frac{-1}{|\widetilde{\mathbf{B}}(x)|} \sum_{j \in \widetilde{\mathbf{B}}(x)} \log \left(\frac{\exp(\mathbf{C}(y)_i \cdot \widetilde{B}_j/\tau)}{\sum_{k \in \widetilde{\mathbf{B}}} \exp(\mathbf{C}(y)_i \cdot \widetilde{B}_k/\tau)} \right), \quad (7)$$

which brings part-level and object-level dictionary components closer when they belong to the same semantic object. 206

252

253

254

255

256

257



Figure 3. **Illustration of logical constraints at inference.** Here, a reptile-head and reptile-body are incorrectly predicted as snake-head and snake-body. CoCal corrects the wrong prediction by computing the logical path probability (multiplying part-level and object-level probabilities) and reassigning labels along the path to produce the correct part prediction.

207 Post-Processing at Inference. Second, we encode the
208 fact that if a pixel belongs to a certain part, it must also
209 belong to the corresponding object. CoCal thus multiplies
210 part- and object-level probabilities and selects the top joint
211 path (Fig. 3). This ensures if a pixel is labeled *dog-head*, it
212 cannot simultaneously be labeled *fish* at the object level.

213 2.2.4. Meta-Architecture Overview

Figure 2 shows the overall pipeline of CoCal. We start 214 with a cluster-based mask transformer backbone, which ex-215 tracts pixel features. On top of these features, two dictio-216 217 naries are learned: a part dictionary C and an object dictionary \tilde{C} . Two memory banks (B and \tilde{B}) store histori-218 cal dictionary components, enabling within-level and cross-219 level contrastive training. Finally, inference proceeds via a 220 parameter-free nearest neighbor search against both dictio-221 222 naries, augmented by a logical consistency check that re-223 assigns part labels according to the best object-part path. Therefore, CoCal delivers a highly interpretable and seman-224 tically coherent solution for object parsing. 225

3. Experiments

In this section, we first provide our main results on PartImageNet [21] and Pascal-Part-108 [49], followed by qualitative comparisons to highlight the effectiveness of CoCal.
For extended ablation studies and additional implementation details, we refer readers to the supplementary materials(including detailed experimental settings as well as ablation studies.).

3.1. Main Results

We summarize our core findings on PartImageNet [21] and 235 Pascal-Part-108 [49] in Table 1a and Table 1b, respec-236 tively. Notably, with a ResNet-50 [23] backbone on PartIm-237 238 ageNet, CoCal surpasses kMaX-DeepLab [72] by 2.08% in part mIoU. When upgraded to the ConvNeXt-Tiny [42] 239 backbone, CoCal continues to excel, reaching 70.31 part 240 mIoU-an additional 1.79% gain over kMaX-DeepLab un-241 242 der the same backbone. Beyond improving part mIoU, CoTable 1. PartImageNet *val* set and Pascal-Part-108 *test* set results. We report part-level and super-category/object-level mIoU (or mAvg), averaged over 3 runs.

(a) PartImagaNat val sat

(a) I artimagei (et vai set					
mathad	backbone –		mIoU		
method		Part	Supe	r-Category	
DeepLabv3+ [6]	ResNet-50	60.57	-		
MaskFormer [12]	ResNet-50	60.34	-		
Compositor [22]	ResNet-50	61.44		-	
kMaX-DeepLab [72]	ResNet-50	65.75	89.16		
CoCal	ResNet-50	67.83	90.41		
SegFormer [69]	MiT-B2	61.97	-		
MaskFormer [12]	Swin-T	63.96		-	
Compositor [22]	Swin-T	64.64		-	
kMaX-DeepLab [72]	ConvNeXt-T	68.52		91.34	
CoCal	ConvNeXt-T	70.31		92.65	
(b) Pascal-Part-108 test set					
method		Part r	nIoU mAvg		
SegNet [2]		18.6		20.8	
FCN [44]		31.6		33.8	
DeepLab [5]		31.6		40.8	
DeepLabv3+ [6]		46.5		48.9	
BSANet [74]		42.9		46.3	
GMNet [49]		45.8		50.5	
FLOAT [53]		48.0		53.0	
HSSN [31]		48.3		-	
DeepLabv3+ [6]+ LOGICSEG [32]		49.1		-	
kMaX-DeepLab [72]		48	.3	49.9	
CoCal		49	.8	52.0	

Cal also enhances super-category segmentation by over 1%, 243 demonstrating robust performance at multiple semantic lev-244 els. Turning to Pascal-Part-108, CoCal achieves 49.8 part 245 mIoU and 52.0 mAvg with a ResNet-101 backbone, thus 246 establishing new state-of-the-art results. Compared to LO-247 GISEG [32] and kMaX-DeepLab, CoCal offers improve-248 ments of 0.7% and 1.5% in part mIoU, respectively, along-249 side a notable 2.1% boost for object-level segmentation. 250

3.2. Qualitative Results

Figure 4 shows three representative examples on PartImageNet. Compared to kMaX-DeepLab [72], CoCal produces more accurate boundaries (see the first row) and detects parts that kMaX-DeepLab misses (rows 2 & 3), demonstrating its superior ability to capture fine structures.

4. Conclusion

In conclusion, this paper introduces CoCal, an innovative 258 model for object parsing that is rooted in a dictionary-based 259 framework. A key aspect of CoCal is its emphasis on eluci-260 dating the intrinsic relationships between parts and objects, 261 which significantly enhances the interpretability and con-262 sistency of parsing outcomes. This approach not only im-263 proves the accuracy of the parsing but also provides a deeper 264 understanding of the complex interplay between part and 265 object entities in images. 266

293

294

295

296

302

303

304

305

323

324

325

326

327

328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346

347

348

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

267 References

- [1] Benjamin Alt, Minh Da Nguyen, Andreas Hermann, Darko
 Katic, Rainer Jaekel, Ruediger Dillmann, and Eric Sax. Efficientpps: Part-aware panoptic segmentation of transparent
 objects for robotic manipulation. In *ISR Europe 2023; 56th International Symposium on Robotics*, pages 131–138. VDE,
 2023. 8
- [2] Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla.
 Segnet: A deep convolutional encoder-decoder architecture
 for image segmentation. *IEEE transactions on pattern anal ysis and machine intelligence*, 39(12):2481–2495, 2017. 4
- [3] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas
 Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-toend object detection with transformers. In *European confer- ence on computer vision*, pages 213–229. Springer, 2020. 8
- [4] Liang-Chieh Chen, Yi Yang, Jiang Wang, Wei Xu, and
 Alan L Yuille. Attention to scale: Scale-aware semantic image segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3640–3649,
 2016. 8
- [5] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos,
 Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image
 segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):834–848, 2017. 4
 - [6] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of the European conference on computer vision* (ECCV), pages 801–818, 2018. 4
- [7] Mu Chen, Zhedong Zheng, Yi Yang, and Tat-Seng Chua.
 Pipa: Pixel-and patch-wise self-supervised learning for domain adaptative semantic segmentation. In *Proceedings of the 31st ACM International Conference on Multimedia*,
 pages 1905–1914, 2023. 8
 - [8] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020. 8
- [9] Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad
 Norouzi, and Geoffrey E Hinton. Big self-supervised mod els are strong semi-supervised learners. *Advances in neural information processing systems*, 33:22243–22255, 2020. 8
- [10] Xianjie Chen, Roozbeh Mottaghi, Xiaobai Liu, Sanja Fidler,
 Raquel Urtasun, and Alan Yuille. Detect what you can: Detecting and representing objects using holistic models and
 body parts. In *CVPR*, pages 1971–1978, 2014. 8, 9
- [11] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He.
 Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020. 8
- [12] Bowen Cheng, Alex Schwing, and Alexander Kirillov. Perpixel classification is not all you need for semantic segmentation. *Advances in Neural Information Processing Systems*,
 320 34:17864–17875, 2021. 4, 8, 10
- [13] Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask

transformer for universal image segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1290–1299, 2022. 8, 10

- [14] Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le. Autoaugment: Learning augmentation policies from data. *arXiv preprint arXiv:1805.09501*, 2018.
 9
- [15] Daan de Geus, Panagiotis Meletis, Chenyang Lu, Xiaoxiao Wen, and Gijs Dubbelman. Part-aware panoptic segmentation. In *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition, pages 5485– 5494, 2021. 1, 8
- [16] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In 2009 IEEE conference on computer vision and pattern recognition, pages 248–255. Ieee, 2009. 8
- [17] Ye Du, Zehua Fu, Qingjie Liu, and Yunhong Wang. Weakly supervised semantic segmentation by pixel-to-prototype contrast. In *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition, pages 4320– 4329, 2022. 8
- [18] S Eslami and Christopher Williams. A generative model for parts-based object segmentation. Advances in Neural Information Processing Systems, 25, 2012. 2, 8
- [19] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2010 (VOC2010) Results. http://www.pascalnetwork.org/challenges/VOC/voc2010/workshop/index.html. 8
- [20] Bharath Hariharan, Pablo Arbeláez, Ross Girshick, and Jitendra Malik. Hypercolumns for object segmentation and fine-grained localization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 447–456, 2015. 8
- [21] Ju He, Shuo Yang, Shaokang Yang, Adam Kortylewski, Xiaoding Yuan, Jie-Neng Chen, Shuai Liu, Cheng Yang, Qihang Yu, and Alan Yuille. Partimagenet: A large, highquality dataset of parts. In *ECCV*, pages 128–145. Springer, 2022. 4, 8, 9
- [22] Ju He, Jieneng Chen, Ming-Xian Lin, Qihang Yu, and Alan L Yuille. Compositor: Bottom-up clustering and compositing for robust part and object segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11259–11268, 2023. 2, 4, 8
- [23] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 770–778, 2016. 4, 9
- [24] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020. 8
- [25] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16000– 16009, 2022. 8
 370

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

- [26] Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q
 Weinberger. Deep networks with stochastic depth. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*, pages 646–661. Springer, 2016. 9
- [27] Tsung-Wei Ke, Jyh-Jing Hwang, Yunhui Guo, Xudong
 Wang, and Stella X Yu. Unsupervised hierarchical semantic
 segmentation with multiview cosegmentation and clustering
 transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2571–
 2581, 2022. 8
- [28] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna,
 Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and
 Dilip Krishnan. Supervised contrastive learning. Advances *in neural information processing systems*, 33:18661–18673,
 2020. 8
- [29] Diederik P Kingma and Jimmy Ba. Adam: A method for
 stochastic optimization. *arXiv preprint arXiv:1412.6980*,
 2014. 9
- 400 [30] Harold W Kuhn. The hungarian method for the assignment
 401 problem. *Naval research logistics quarterly*, 2(1-2):83–97,
 402 1955. 2
- [31] Liulei Li, Tianfei Zhou, Wenguan Wang, Jianwu Li, and Yi
 Yang. Deep hierarchical semantic segmentation. In *Proceed- ings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1246–1257, 2022. 4, 8
- 407 [32] Liulei Li, Wenguan Wang, and Yi Yang. Logicseg: Parsing
 408 visual semantics with neural logic learning and reasoning.
 409 In *Proceedings of the IEEE/CVF International Conference*410 on *Computer Vision*, pages 4122–4133, 2023. 4, 8
- [33] Tianyu Li, Subhankar Roy, Huayi Zhou, Hongtao Lu, and
 Stéphane Lathuilière. Contrast, stylize and adapt: Unsupervised contrastive learning framework for domain adaptive semantic segmentation. In *Proceedings of the IEEE/CVF Con- ference on Computer Vision and Pattern Recognition*, pages
 4868–4878, 2023. 8
- 417 [34] Xiangtai Li, Shilin Xu, Yibo Yang, Guangliang Cheng, Yunhai Tong, and Dacheng Tao. Panoptic-partformer: Learning a unified model for panoptic part segmentation. In *European Conference on Computer Vision*, pages 729–747. Springer, 2022. 1, 8
- [35] Xiangtai Li, Shilin Xu, Yibo Yang, Haobo Yuan, Guangliang
 Cheng, Yunhai Tong, Zhouchen Lin, Ming-Hsuan Yang, and
 Dacheng Tao. Panopticpartformer++: A unified and decoupled view for panoptic part segmentation. *arXiv preprint arXiv:2301.00954*, 2023. 8
- 427 [36] Zhiheng Li, Wenxuan Bao, Jiayang Zheng, and Chenliang
 428 Xu. Deep grouping model for unified perceptual parsing.
 429 In Proceedings of the IEEE/CVF Conference on Computer
 430 Vision and Pattern Recognition, pages 4053–4063, 2020. 8
- 431 [37] James Liang, Tianfei Zhou, Dongfang Liu, and Wenguan
 432 Wang. Clustseg: Clustering for universal segmentation.
 433 2023. 2, 8
- [38] Xiaodan Liang, Si Liu, Xiaohui Shen, Jianchao Yang, Luoqi
 Liu, Jian Dong, Liang Lin, and Shuicheng Yan. Deep human
 parsing with active template regression. *IEEE transactions on pattern analysis and machine intelligence*, 37(12):2402–
 2414, 2015. 8

- [39] Xiaodan Liang, Ke Gong, Xiaohui Shen, and Liang Lin.
 Look into person: Joint body parsing & pose estimation network and a new benchmark. *IEEE transactions on pattern analysis and machine intelligence*, 41(4):871–885, 2018. 8
- [40] Xiaodan Liang, Hongfei Zhou, and Eric Xing. Dynamicstructured semantic propagation network. In *Proceedings* of the IEEE Conference on Computer Vision and Pattern Recognition, pages 752–761, 2018. 8
- [41] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017. 3
- [42] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 11976–11986, 2022. 4, 9
- [43] Stuart Lloyd. Least squares quantization in pcm. *IEEE transactions on information theory*, 28(2):129–137, 1982. 8
- [44] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015. 4
- [45] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2018. 9
- [46] Wenhao Lu, Xiaochen Lian, and Alan Yuille. Parsing semantic parts of cars using graphical models and segment appearance consistency. arXiv preprint arXiv:1406.2375, 2014. 8
- [47] Lele Lv, Qing Liu, Shichao Kan, and Yixiong Liang. Confidence-aware contrastive learning for semantic segmentation. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 5584–5593, 2023. 8
- [48] Umberto Michieli and Pietro Zanuttigh. Edge-aware graph matching network for part-based semantic segmentation. *International Journal of Computer Vision*, 130(11):2797– 2821, 2022. 1
- [49] Umberto Michieli, Edoardo Borsato, Luca Rossi, and Pietro Zanuttigh. Gmnet: Graph matching network for large scale part semantic segmentation in the wild. In *Computer Vision– ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VIII 16*, pages 397–414. Springer, 2020. 1, 4, 8
- [50] Shishir Muralidhara, Sravan Kumar Jagadeesh, René Schuster, and Didier Stricker. Jppf: Multi-task fusion for consistent panoptic-part segmentation. *SN Computer Science*, 5(1):187, 2024. 8
- [51] Xuecheng Nie, Jiashi Feng, and Shuicheng Yan. Mutual learning to adapt for joint human parsing and pose estimation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 502–517, 2018. 8
- [52] Joshua David Robinson, Ching-Yao Chuang, Suvrit Sra, and Stefanie Jegelka. Contrastive learning with hard negative samples. In *International Conference on Learning Representations*, 2020. 3, 8
- [53] Rishubh Singh, Pranav Gupta, Pradeep Shenoy, and Ravikiran Sarvadevabhatla. Float: Factorized learning of object attributes for improved multi-object multi-part scene parsing.

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

In Proceedings of the IEEE/CVF Conference on Computer 498 Vision and Pattern Recognition, pages 1445-1455, 2022. 4, 499 8.9

- 500 [54] Yafei Song, Xiaowu Chen, Jia Li, and Qinping Zhao. Em-501 bedding 3d geometric features for rigid object part segmen-502 tation. In Proceedings of the IEEE international conference 503 on computer vision, pages 580-588, 2017. 8
- 504 [55] Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia 505 Schmid. Segmenter: Transformer for semantic segmenta-506 tion. In Proceedings of the IEEE/CVF international confer-507 ence on computer vision, pages 7262-7272, 2021. 8
- [56] BLE Verboeket and Gijs Dubbelman. A hierarchical ap-508 509 proach to part-aware panoptic segmentation. 2022. 8
- [57] Changqi Wang, Haoyu Xie, Yuhui Yuan, Chong Fu, and 510 Xiangyu Yue. Space engage: Collaborative space super-511 512 vision for contrastive-based semi-supervised semantic seg-513 mentation. In Proceedings of the IEEE/CVF International 514 Conference on Computer Vision, pages 931–942, 2023. 8
- 515 [58] Huiyu Wang, Yukun Zhu, Bradley Green, Hartwig Adam, 516 Alan Yuille, and Liang-Chieh Chen. Axial-deeplab: Stand-517 alone axial-attention for panoptic segmentation. In European 518 conference on computer vision, pages 108-126. Springer, 519 2020.8
- 520 [59] Huiyu Wang, Yukun Zhu, Hartwig Adam, Alan Yuille, and 521 Liang-Chieh Chen. Max-deeplab: End-to-end panoptic 522 segmentation with mask transformers. In Proceedings of 523 the IEEE/CVF conference on computer vision and pattern 524 recognition, pages 5463-5474, 2021. 8
- 525 [60] Jianyu Wang and Alan L Yuille. Semantic part segmentation using compositional model combining shape and appear-526 527 ance. In Proceedings of the IEEE conference on computer 528 vision and pattern recognition, pages 1788–1797, 2015. 8
- 529 [61] Peng Wang, Xiaohui Shen, Zhe Lin, Scott Cohen, Brian 530 Price, and Alan L Yuille. Joint object and part segmentation 531 using deep learned potentials. In Proceedings of the IEEE 532 International Conference on Computer Vision, pages 1573-533 1581, 2015, 2
- 534 [62] Wenguan Wang, Zhijie Zhang, Siyuan Qi, Jianbing Shen, 535 Yanwei Pang, and Ling Shao. Learning compositional neu-536 ral information fusion for human parsing. In Proceedings of 537 the IEEE/CVF international conference on computer vision, 538 pages 5703-5713, 2019. 8
- 539 [63] Wenguan Wang, Hailong Zhu, Jifeng Dai, Yanwei Pang, 540 Jianbing Shen, and Ling Shao. Hierarchical human pars-541 ing with typed part-relation reasoning. In Proceedings of 542 the IEEE/CVF conference on computer vision and pattern recognition, pages 8929-8939, 2020. 8 543
- 544 [64] Wenguan Wang, Tianfei Zhou, Fisher Yu, Jifeng Dai, En-545 der Konukoglu, and Luc Van Gool. Exploring cross-image 546 pixel contrast for semantic segmentation. In Proceedings 547 of the IEEE/CVF International Conference on Computer Vi-548 sion, pages 7303-7313, 2021. 8
- 549 [65] Fangting Xia, Peng Wang, Liang-Chieh Chen, and Alan L 550 Yuille. Zoom better to see clearer: Human and object parsing 551 with hierarchical auto-zoom net. In Computer Vision-ECCV 552 2016: 14th European Conference, Amsterdam, The Nether-553 lands, October 11-14, 2016, Proceedings, Part V 14, pages 554 648-663. Springer, 2016. 8

- [66] Fangting Xia, Peng Wang, Xianjie Chen, and Alan L Yuille. 555 Joint multi-person pose estimation and semantic part seg-556 In Proceedings of the IEEE conference on 557 mentation. computer vision and pattern recognition, pages 6769-6778, 558 2017.8 559
- [67] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing for scene understanding. In Proceedings of the European conference on computer vision (ECCV), pages 418-434, 2018. 8
- [68] Binhui Xie, Shuang Li, Mingjia Li, Chi Harold Liu, Gao Huang, and Guoren Wang. Sepico: Semantic-guided pixel contrast for domain adaptive semantic segmentation. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2023.8
- [69] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. Advances in Neural Information Processing Systems, 34:12077-12090, 2021. 4
- [70] Kota Yamaguchi, M Hadi Kiapour, Luis E Ortiz, and Tamara L Berg. Parsing clothing in fashion photographs. In 2012 IEEE Conference on Computer vision and pattern recognition, pages 3570-3577. IEEE, 2012. 8
- [71] Qihang Yu, Huiyu Wang, Dahun Kim, Siyuan Qiao, Maxwell Collins, Yukun Zhu, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. Cmt-deeplab: Clustering mask transformers for panoptic segmentation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 2560-2570, 2022. 2, 8
- [72] Qihang Yu, Huiyu Wang, Siyuan Qiao, Maxwell Collins, Yukun Zhu, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. k-means Mask Transformer. In ECCV, pages 288-307. Springer, 2022. 2, 4, 8, 9
- [73] Wenwei Zhang, Jiangmiao Pang, Kai Chen, and Chen Change Loy. K-net: Towards unified image segmentation. Advances in Neural Information Processing Systems, 34, 2021. 8
- [74] Yifan Zhao, Jia Li, Yu Zhang, and Yonghong Tian. Multiclass part parsing with joint boundary-semantic awareness. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 9177-9186, 2019. 4, 8, 9
- [75] Long Zhu, Yuanhao Chen, Chenxi Lin, and Alan Yuille. 596 Max margin learning of hierarchical configural deformable 597 templates (hcdts) for efficient object parsing and pose esti-598 mation. International journal of computer vision, 93:1-21, 599 2011.8 600

665

678

682

683

684

685

686

687

688

689

690

691

5. Related Work

602 5.1. Object Parsing

The extensive literature on object parsing can be divided 603 into single-object multi-part parsing [4, 20, 39, 51, 65, 66] 604 605 and multi-object multi-part parsing [22, 49, 53, 74]. Single-606 object multi-part parsing has primarily focused on specific 607 classes, such as humans [38, 70, 75], animals [60], and vehicles [18, 46, 54]. While the methodologies addressing 608 multi-object multi-part parsing mainly focus on employing 609 top-down or coarse-to-fine strategies. Specifically, Singh 610 611 et al. [53] proposed FLOAT, a factorized top-down parsing framework by first detecting the object followed with 612 zooming in for obtaining higher quality part masks. On the 613 contrary, He et al. [22] introduced Compositor, a bottom-up 614 architecture designed to iteratively learn objects by cluster-615 616 ing pixels to derive parts. Recently, there are also explo-617 rations in the closely related area of panoptic part segmentation within the research community. Notable works such 618 as [1, 15, 34, 35, 50, 56] have delved into the semantic pars-619 ing of objects while also distinguishing parts between dif-620 621 ferent instances. However, a common trend in these works, 622 whether focused on semantic object parsing or panoptic part segmentation, involves extending standard segmenta-623 tion models, often overlooking the nuanced semantic levels 624 of parts. In contrast, CoCal takes a novel approach by fo-625 626 cusing specifically on semantic object parsing. It redefines the paradigm of cluster-based mask transformers and intro-627 duces a novel dictionary-based framework meticulously tai-628 lored for object parsing. 629

630 5.2. Cluster-based Mask Transformer

631 With the recent progress in transformers [3], a new paradigm named mask classification [12, 13, 55, 58, 59, 73] 632 633 has been proposed, where segmentation predictions are represented by a set of binary masks with its class label, which 634 is generated through the conversion of object queries to 635 mask embedding vectors followed by multiplying with the 636 image features. The predicted masks are trained by Hun-637 638 garian matching with ground truth masks. Thus the essential component of mask transformers is the decoder which 639 takes object queries as input and gradually transfers them 640 into mask embedding vectors. Recently, cluster-based mask 641 642 transformers are introduced in [37, 71, 72], which rethinks the design of the decoder by replacing the cross-attention 643 with a k-means [43] attention. Building upon these ex-644 plorations, CoCal introduces a global class dictionary and 645 replaces the Hungarian matching with a fixed one-to-one 646 matching, thereby establishing an interpretable dictionary-647 648 based framework for part segmentation.

5.3. Contrastive Learning in Segmentation

CVPRW 2025 Submission #*****. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

Contrastive learning [8, 9, 11, 24, 25, 28, 52] has emerged 650 as a prominent technique in computer vision as an effec-651 tive method for learning feature representation for self-652 supervised models. The core idea lies in contrasting sim-653 ilar (positive) data pairs against dissimilar (negative) pairs. 654 Recently, Wang et al. [64] raise a pixel-to-pixel contrastive 655 learning method for semantic segmentation, which enforces 656 pixel embeddings belonging to the same semantic class to 657 be more similar than embeddings from different classes. 658 [7, 17, 33, 47, 57, 68] are built upon this concept, extend-659 ing it to various segmentation domains. Motivated by these 660 advancements, we propose a component-wise contrastive 661 learning method tailored for modern cluster-based mask 662 transformers, which effectively learns discriminative dictio-663 nary components within the clustering scheme. 664

5.4. Logical Constraints in Segmentation

Few segmentation models [27, 31, 32, 36, 40, 62, 63, 67] 666 consider the implicit logic rules inherent in structured la-667 bels. While the majority of them are dedicated to human 668 parsing, a few recent works [31, 32] tackle the general seg-669 mentation in a flexible function and avoid incorporating la-670 bel taxonomies into the network topology. Concretely, Li 671 et al. [31] enhance the logical consistency by modeling the 672 segmentation as a pixel-wise multi-label classification. Li 673 et al. [32] exploit neuro-symbolic computing for grounding 674 logical formulae onto data. In contrast to these efforts, Co-675 Cal introduces an object level on top of the part and models 676 logical rules as a contrastive objective during training. 677

6. Extended Experiments and Analyses

This section provides details on our experiments, including
ablation studies, dataset statistics, training parameters, and
an in-depth exploration of CoCal's design components.679680
681680

6.1. Experimental Setup and Datasets

Datasets We conduct experiments on two popular object parsing benchmarks: PartImageNet [21] and PASCAL-Part-108 [49]. We provide the detailed statistics of each dataset and the class definitions below:

- PartImageNet [21] contains 24095 elaborately annotated general images from ImageNet [16], which are split into 20481/1206/2408 for *train/val/test*. It is associated with 40 part classes, which are grouped into 11 object classes following the official class definition.
- Pascal-Part-108 [49] expands upon the part definition introduced in Pascal-Part-58 [10], providing a more intricate benchmark with finer part-level details. This extension maintains the original split of VOC [19] and encompasses a dataset of 10,103 images across 20 object classes and 108 part classes. Our experiments adhere to the orig-697



Figure 4. Qualitative comparison between CoCal and kMaX-DeepLab on PartImageNet. Our CoCal yields more precise part boundaries (row 1) and captures missed parts (rows 2 & 3).

inal split, utilizing 4,998 images for training and 5,105images for testing.

Evaluation Metrics We evaluate the performance of Co-700 701 Cal on the PartImageNet dataset [21] using the mean Intersection over Union (mIoU) on both part and super-category 702 703 levels. It's important to note that for PartImageNet, we choose to report performance on the super-category level 704 705 because the parts in PartImageNet are defined within the context of super-category. The hierarchy of super-category 706 707 is inherited for training CoCal on this dataset. In the case of Pascal-Part-108, our evaluation includes reporting part 708 709 mIoU, and additionally, we calculate the mAvg on the object level. The mAvg metric, as defined in the literature [74], 710 provides the average mIoU score of all parts belonging to an 711 object. We refer the reader to FLOAT [53] for a detailed ex-712 713 planation of these metrics.

Training details We implement CoCal based on the kMaX-DeepLab architecture [72], utilizing its official PyTorch re-implementation codebase. To ensure a fair comparison, we adopt the training settings from kMaXDeepLab. The backbone, pretrained on ImageNet [23, 42],
followed a learning rate multiplier of 0.1. For regularization and augmentations, we incorporate drop path [26]

and random color jittering [14]. The optimizer used is 721 AdamW [29, 45] with a weight decay of 0.05. Unless oth-722 erwise specified, we train all models with a batch size of 723 64 on a single A100 GPU, performing 40,000 iterations 724 on PartImageNet [21] and 10,000 iterations on Pascal-Part-725 108 [10]. The first 2,000 and 500 steps serve as the warm-up 726 stage, where the learning rate linearly increases from 0 to 727 5×10^{-4} . The training objective for CoCal includes the 728 combination of kMaX-DeepLab's original losses and the 729 proposed contrastive loss terms, as specified in Eq. 5, Eq. 6 730 and Eq. 7: 731

$$\mathcal{L} = \lambda_{kMaX} \mathcal{L}_{kMaX} + \lambda_{p_con} \mathcal{L}_{p_con} +$$
732

$$\lambda_{o_con} \mathcal{L}_{o_con} + \lambda_{logic} \mathcal{L}_{logic}.$$
733

Here, \mathcal{L}_{kMaX} represents the loss from kMaX-DeepLab [72], 734 and λ_{kMaX} follows the official setting. The weights for the 735 proposed loss terms are set to $\lambda_{p_con} = 2$, $\lambda_{o_con} = 2$, and 736 $\lambda_{logic} = 1$. CoCal uses the exact same number of part and 737 object queries corresponding to the part and object classes 738 in the dataset. Specifically, we set P to 41 and 109, and \tilde{P} to 739 12 and 21 (with one additional learnable component for rep-740 resenting the background at both the part and object levels) 741 in PartImageNet and Pascal-Part-108, respectively. This de-742

sign enables a straightforward and highly interpretable in-743 ference process, using nearest neighbor search for parts and 744

objects separately during inference. Afterward, we com-745

746 pute the top-scoring logical path and reassign the predicted classes based on that path. 747

6.2. Ablation Studies 748

Dictionary Components, Contrastive Objectives, and 749 Logical Constraints Table 2 (reproduced from the main 750 paper for convenience) shows an ablation on core design 751 choices. Simply switching kMaX-DeepLab to a dictionary-752 753 based version slightly degrades performance (65.75 to 64.31), but adding contrastive objectives and logical con-754 straints incrementally boosts part mIoU to 67.83, surpassing 755 the original baseline. 756

> Table 2. Ablation of CoCal components on PartImageNet val with ResNet-50.

method	Dict	\mathcal{L}_{p_con}	\mathcal{L}_{o_con}	\mathcal{L}_{logic}	Part mIoU
kMaX-DeepLab	X	X	X	X	65.75
Dictionary-based	\checkmark	×	X	×	64.31
CoCal (ours)	\checkmark	\checkmark	×	×	65.87
CoCal (ours)	\checkmark	\checkmark	\checkmark	×	66.53
CoCal (ours)	\checkmark	\checkmark	\checkmark	\checkmark	67.83

Memory Bank Size S In Table 3, we vary S and observe 757 that excessively small or large values degrade performance 758 due to insufficient or redundant samples. 759

> Table 3. Impact of memory bank size on PartImageNet val (ResNet-50).

$\# \operatorname{memory} \operatorname{bank} S$	Part mIoU
50	66.50
100	67.83
150	67.16
200	67.02

Number of Negative Samples k Table 4 shows that using 760 too few negatives (e.g., k = 50) reduces performance to 761 66.28 part mIoU, while too many (e.g., k = 200 or "all") 762 also hurts accuracy.

763

Table 4. Influence of negative samples on PartImageNet val (ResNet-50).

# negative sample k	Part mIoU
50	66.28
100	67.83
200	66.40
all	65.74

Generalizability of CoCal Finally, Table 5 demonstrates 764 that CoCal also boosts other transformer-based frameworks 765

such as MaskFormer [12] and Mask2Former [13], improv-766 ing part mIoU by 3.18 and 2.77, respectively. 767

Table 5. Generalizability to other baselines on PartImageNet val (ResNet-50).

method	mIoU		
	Part	Super-Category	
MaskFormer	60.34	-	
CoCal (MaskFormer)	63.52	86.67	
Mask2Former	63.62	87.20	
CoCal (Mask2Former)	66.39	88.72	