000		054
001		055
002		050
004	Supplemental Materials to	058
005		059
006	" SIM: Semantic-aware Instance Mask Generation for	060
007	Box-Supervised Instance Segmentation "	061
800	i b	062
009		063
010	Anonymous CVPR submission	064
011	•	065
012	Paper ID 532	066
013		067
014		800
015	In this supplemental file, we provide the following materials:	009
017	In this supplemental me, we provide the following materials.	070
018	• Details about the online weakly-supervised Copy-Paste operation, and some visualization results (cf. Sec3.3-Online	072
019	Weakly-Supervised Copy-Paste in the main paper);	073
020	• Implementation details on Mask2Former [2] (cf. Sec4.1-Implementation Details in the main paper):	074
021	• Implementation details on Wask21 of the [2] (c). See4.1-Implementation Details in the main paper),	075
022	• Analysis on positive weighting strategy and more visualizations (cf. Sec3.2.2-Positive Mask Weighting and Sec4.4-	076
023	Visualizations of Weights in the main paper);	077
024	• More qualitative results (of Sec/ 3 Qualitative Results in the main paper):	078
025	• More quantative results (c). See4.5-Quantative Results in the main paper),	079
026	• More parameter analyses (cf. Sec4.4-Ablation Study in the main paper).	080
027		081
028	A. Online Weakly-Supervised Copy-Paste	082
029	Copy-paste is a simple yet effective way to improve the data efficiency of instance segmentation models. By pasting	003
031	objects of various categories and scales to different images, copy-paste could achieve solid performance improvements on	085
032	strong baseline models [3–6]. In addition, Ghiasi et al. [6] demonstrated the efficacy of copy-paste under the semi-supervised	086
033	learning setting. However, Copy-Paste has rarely been explored for weakly-supervised instance segmentation. In this work,	087
034	we use copy-paste to create new training data for better handling the object occlusions and rare object categories. Following	088
035	the work [6], we adopt a simple strategy of randomly picking objects and pasting them on the target image, which could	089
036	provide a considerable boost on top of baselines.	090
037	Importance Sampling. The quality of pseudo mask varies significantly across different instances within an image. A	091
038	poor mask may produce an unconvincing paste, while a good mask could result in a convincing one. Therefore, we adopt an	092
039	importance sampling strategy to select instances with high-quality masks. Specifically, we calculate the averaged mask score	093
040	S that reflects the importance of each instance by:	094
041	$\sum^{N_k} M^{k,i} \cdot \mathbb{1}(\hat{M}^{k,i} = 1)$	095
042	$S_k = \frac{\sum_i \frac{M_{prob} - \mathbb{I}(M_i - 1)}{\sum_{i=1}^{N_k} \mathbb{I}(\widehat{\mathcal{I}}_k) - 1)},\tag{1}$	096
043	$\sum_{i}^{i < \kappa} \mathbb{I}(M^{\kappa, i} = 1)$	097
045	where N_k denotes the number of pixels from the k-th instance map of M_{rreck} and S_k is the score of the k-th instance. We	090
046	also store S into the memory bank \mathcal{M} .	100
047	For each training iteration, we randomly sample an image $\{X', Y', B', \hat{M}', S'\}$ from \mathcal{M} and randomly extract a subset of	101
048	instances according to their mask scores S' , so that the instances with higher-quality masks are more likely to be selected.	102
049	Examples of online weakly-supervised Copy-Paste are shown in Fig. S1, we cut objects which have significant contrast	103
050	with the background, and randomly paste them on the target image. In this way, more challenging data with various occlusion	104
051	patterns can be created, which could effectively improve the model ability to handle the occlusions between objects. In	105
052	addition, since small objects are the majority (more than 40%) in the COCO dataset, more small objects will be created so	106
053	that the Copy-Paste strategy will benefit small objects much (please refer to the Tab. 1 in our main paper).	107

CVPR 2023 Submission #532. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

108 109 110 111 112 113					
114 115 116 117 118 119				•	
120 121 122 123 124 125				*	
126 127 128 129 130 131			•	•	
132 133 134 135 136 137				4	
138 139 140 141 142				Τ	
143 144 145 146 147 148					-
149 150 151 152 153 154			COSS -		2
155 156 157 158 159 160	x	<i>x'</i>	Xpaste	Â,	Â'

Figure S1. Examples of online weakly-supervised Copy-Paste. We use \hat{M}' to extract instances from X' and paste them onto X, resulting in new training data X_{paste} .



Figure S2. Visualizations of weights for different positive samples.

B. Implementation Details

We adopt CondInst [9] and Mask2Former [2] as our baselines. For Mask2former, we use Detectron2 and follow the commonly used settings for each dataset. More specifically, we use AdamW [7] and the *poly* [1] learning rate schedule with an initial learning rate of 10^{-4} . A learning rate multiplier of 0.1 is applied to CNN backbones. The more advanced multi-scale deformable attention transformer [11] is employed as the default pixel decoder.

C. Analysis on Positive Weights

In Sec. 3.2.2 of our main paper, we use the instance map to rectify the semantic map so that those falsely activated objects and regions will be suppressed while the correct ones are enhanced. Since the quality of masks produced by different positive samples varies significantly, we propose a positive mask weighting strategy to integrate different masks according to their quality. We visualize the weights of different positive samples in Fig. S2. One can observe that large weights tend to locate in the central regions of foreground objects, and the smaller weights tend to locate in ambiguous regions.

Effect of μ . The parameter μ in Eq. 6 of our manuscript controls the relative gaps between different weights. We investigate its effect in Tab. S1. Out method attains the best performance when μ is set to 5.0. When we equally treat different positive samples (*i.e.*, set μ to 0), the performance decreases by 0.8% AP. When we set μ to a larger value, *e.g.*, $\mu = 20$, the final result drops by 0.9% AP. This is because we rely on some certain positive samples while ignoring other useful ones. Actually, different positive samples provide complementary information and they should be fully integrated for producing better results.

CVPR
#532

	-											
		μ	AP	A	P ₅₀	AP_{75}	AP _S	AP_M	AP_L	,		
	-	0	31.4	4 52	2.9	31.6	14.3	33.7	46.9	1		
		1	31.6	5 53	3.4	32.0	14.5	34.1	47.2	r		
		5	32.2	2 54	1.0	33.0	15.8	34.5	48.3			
		10	31.8	3 53	3.9	32.3	15.0	34.8	46.8			
		20	31.3	3 53	3.1	31.9	15.1	33.9	46.1			
	Table S1. Effect of parameter μ .											
	τ_{low}	$ au_h$	iigh	AP	A	AP_{50} A	P ₇₅	AP_S	AP_M	AP_L		
	0.5	C).5	31.4	4	53.6 3	1.9	14.7	34.3	46.8		
	0.4	C).6	31.9	4	53.9 3	2.3	15.4	34.8	47.7		
	0.3	C	0.7	32.2	4	54.0 3	3.0	15.8	34.5	48.3		
	0.2	0).8	31.8	4	53.6 3	2.8	15.5	34.1	47.2		
			Table	52. E	ffect	t of thresh	olds $ au_{lo}$	τ_w and $ au_h$	igh.			
	tempe	eratur	$e(\tau)$	AP		AP ₅₀	AP_{75}	AP_S	AP_M	AP		
	0.1		0.1)	52.5	30.9	14.5	33.0	45.0		
		0.5 1.0 2.0 5.0		31.7	7	53.3	32.2	15.5	34.1	46.4		
				1.0		32.2	2	54.0	33.0	15.8	34.5	48.
				31.9)	54.1	32.1	15.6	34.0	47.′		
				31.6	5	53.8	32.5	15.6	34.2	47.:		
	Table S3. Effect of temperature τ .											
		λ_1	AP	A	P ₅₀	AP_{75}	AP _S	AP_M	AP _L	ŗ		
	-	0.1	31.9) 53	3.7	32.8	15.4	34.4	47.2	,		
		0.3	32.1	l 54	1.0	33.1	15.9	34.2	47.9)		
		0.5	32.2	2 54	1.0	33.0	15.8	34.5	48.3			
	0.7).7 31.4		3.6	32.2	15.2	33.9	46.9	1		
	Table S4. Effect of parameter λ_1 .											
	-											
Qualitative Results												
Za C2 above serves a 1	totime		ntot:		14	.f	411 .	. J D T	α «4 Γ10]	C(

Fig. S3 shows some qualitative segmentation results of our method and BoxInst [10] on COCO *val* split. On the one hand, our method could better segment foreground instances that heavily tangle with the background or other objects with similar appearances. On the other hand, our method is good at separating overlapping objects of the same semantics and can keep the integrity of objects.

Fig. S4 gives more results of SIM with the ResNet-101-FPN backbone and $3 \times$ training schedule. One can see that our method achieves precise predictions around the objects' boundaries.

E. Parameter analysis

Tab. S2 shows the effects of the two thresholds τ_{low} and τ_{high} . When we set $\tau_{low} = \tau_{high} = 0.5$, all pixels provide supervision, which inevitably introduces much noise. When we set $\tau_{low} = 0.2$ and $\tau_{high} = 0.8$, many pixels are neglected and hence limited supervision is provided and the performance is slightly degraded.

Tab. S3 shows the effect of temperature τ . We employ the instance map $M_{\rm I}$ as a weight map to online rectify the semantic map $M_{\rm S}$, where the temperature τ controls the modulation intensity. When $\tau \to 0$, the modulation intensity increases so that the final pseudo mask \hat{M} will rely more on $M_{\rm S}$. When $\tau \to \infty$, the modulation intensity decreases so that the final pseudo mask \hat{M} will rely more on $M_{\rm I}$.

Tab. S4 shows the segmentation results by using different weights λ_1 in pseudo loss. One can see that our method is insensitive to this parameter when $0.1 < \lambda_1 < 0.5$.

432 References

- [1] 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. 3
- [2] 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. 1, 3
 - [3] Nikita Dvornik, Julien Mairal, and Cordelia Schmid. Modeling visual context is key to augmenting object detection datasets. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 364–380, 2018. 1
 - [4] Debidatta Dwibedi, Ishan Misra, and Martial Hebert. Cut, paste and learn: Surprisingly easy synthesis for instance detection. In *Proceedings of the IEEE international conference on computer vision*, pages 1301–1310, 2017. 1
 - [5] Hao-Shu Fang, Jianhua Sun, Runzhong Wang, Minghao Gou, Yong-Lu Li, and Cewu Lu. Instaboost: Boosting instance segmentation via probability map guided copy-pasting. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 682– 691, 2019.
 - [6] Golnaz Ghiasi, Yin Cui, Aravind Srinivas, Rui Qian, Tsung-Yi Lin, Ekin D Cubuk, Quoc V Le, and Barret Zoph. Simple copy-paste is a strong data augmentation method for instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision* and Pattern Recognition, pages 2918–2928, 2021.
 - [7] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101, 2017. 3
 - [8] Tal Remez, Jonathan Huang, and Matthew Brown. Learning to segment via cut-and-paste. In Proceedings of the European conference on computer vision (ECCV), pages 37–52, 2018.
 - [9] Zhi Tian, Chunhua Shen, and Hao Chen. Conditional convolutions for instance segmentation. In *European conference on computer vision*, pages 282–298. Springer, 2020. **3**
 - [10] Zhi Tian, Chunhua Shen, Xinlong Wang, and Hao Chen. Boxinst: High-performance instance segmentation with box annotations. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 5443–5452, 2021. 4, 6
- [11] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end
 object detection. In *International Conference on Learning Representations*, 2020. 3

CVPR #532



Figure S3. Qualitative results of BoxInst [10] (in the red box) and our method (in the blue box) on COCO val2017.

CVPR #532



7