



Figure 6: **Object Instances per Image:** histograms of the distribution of the number of object instances referenced in each image within the training sets of the studied datasets.

A CONSTRAINED GREEDY MATCHING ALGORITHM

The details of our constrained greedy matching algorithm are presented in [1](#). For the sake of simplicity of notation, the image index i is dropped as this matching only occurs for all objects within an image.

Algorithm 1 Constrained Greedy Matching

```

1: Input: mask choices  $\{\mathbf{m}_{j,k}^c\}_c$ , model predictions  $\hat{\mathbf{m}}_{j,k}$ 
2: Result: greedy matching  $\delta_{j,k,c}^*$ 
3:  $\delta_{j,k,c}^* \leftarrow 0$ 
4:  $\mathcal{C} \leftarrow \emptyset$  ▷ Initialize the exclusion set
5:  $\mathcal{M} \leftarrow \left\{ j, k, c : \ell_{\text{IoU}} \left( \hat{\mathbf{m}}_{j,k}, \mathbf{m}_{j,k}^c \right) \right\}$  ▷ Compute pseudo-IoU for all mask & prediction pairs
6: SORT( $\mathcal{M}$ ) ▷ Sort them in descending order
7: while  $\mathcal{M} \neq \emptyset$  do
8:    $j', k', c' \leftarrow \text{POP}(\mathcal{M})$  ▷ Get the next highest IoU mask choice
9:   if  $c \in \mathcal{C}$  or  $j \in \mathcal{C}$  then ▷ If either the object or mask has been matched, skip it
10:    continue
11:   end if
12:   for  $k$  do ▷ Match it for all expressions of the same object
13:      $\delta_{j',k,c'}^* \leftarrow 1$ 
14:   end for
15:    $\mathcal{C} \leftarrow \mathcal{C} \cup \{(j', c')\}$  ▷ Add object & mask to the exclusion set to avoid re-matching them
16: end while
17: return  $\delta_{j,k,c}^*$ 

```

B NUMBER OF OBJECT INSTANCES PER IMAGE

In Figure [6](#) we show the distribution of the number of objects per image in the training sets of RefCOCO, RefCOCO+ and RefCOCOg (umd and google splits). Note that this average is higher for RefCOCO and RefCOCO+ than for RefCOCOg. The higher the average number of object instances per image, the more effective we expect the loss from Stage 3 to be in correcting the zero-shot mistakes due to the matching mechanism. This intuition is backed by the experimental results presented in Section [4.1](#) of the main paper, where the improvement in RefCOCO and RefCOCO+ from S+S+C is higher than the one in RefCOCOg (umd and google partitions).

C PIXEL-LEVEL CONTRASTIVE ABLATION

On top of the constrained greedy matching results, we also experimented with a pixel-dense contrastive loss term. This contrastive term, $\mathcal{L}_{\leftrightarrow}$ has the goal of regularizing the output by explicitly leveraging the insight that references to the same object (positive examples) in an image should have the same output, and that other objects in the same image (negative examples) should have a different output.

Table 5: **Zero-shot and Weakly-Supervised Comparison with $\mathcal{L}_{\leftrightarrow}$** : oIoU (top) and mIoU (bottom) results on benchmark datasets for our zero-shot method (S+S), our full method (S+S+C), the zero-shot bootstrapped models from Stage 3 (ZSBOOTSTRAP), and an ablation with the proposed pixel-dense contrastive loss $\mathcal{L}_{\leftrightarrow}$, along with existing baselines and ablations. The first column refers to the type of method: zero-shot (ZS), weakly-supervised (WS) or fully-supervised (FS). Best zero-shot results are highlighted in **purple**, and the best weakly-supervised ones in **green**. For RefCOCOg, U refers to the UMD partition, and G refers to the Google partition.

		RefCOCO			RefCOCO+			RefCOCOg		
		val	testA	testB	val	testA	testB	val(U)	test(U)	val(G)
oIoU										
ZS	GL CLIP	24.88	23.61	24.66	26.16	24.90	25.83	31.11	30.96	30.69
	GL CLIP (SAM)	24.50	26.00	21.00	26.88	29.95	22.14	28.92	30.41	28.92
	S+S (<i>Ours</i>)	33.31	40.35	26.14	34.84	43.16	28.22	35.71	42.10	41.70
WS	TRIS	31.17	32.43	29.56	30.90	30.42	30.80	36.00	36.19	36.23
	ZSBOOTSTRAP (<i>Ours</i>)	33.61	42.20	26.12	34.13	42.03	26.60	38.27	40.09	37.03
	S+S+C (<i>Ours</i>)	50.13	60.70	43.46	40.61	49.68	29.54	41.96	42.59	42.18
	S+S+C + $\mathcal{L}_{\leftrightarrow}$ (<i>Ours</i>)	50.43	61.66	43.28	39.47	48.97	30.08	41.62	42.48	42.26
FS	LAVT	72.73	75.82	68.79	62.14	63.38	55.10	61.24	62.09	60.50
mIoU										
ZS	GL CLIP	26.20	24.94	26.56	27.80	25.64	27.84	33.52	33.67	33.61
	GL CLIP (SAM)	30.79	33.08	27.51	32.99	37.17	29.47	39.45	40.85	40.66
	S+S (<i>Ours</i>)	36.95	43.77	27.97	37.68	46.24	29.31	41.41	47.18	47.57
WS	TSEG	25.95	–	–	22.26	–	–	23.41	–	–
	Shatter&Gather	34.76	34.58	35.01	28.48	28.60	27.98	–	–	28.87
	ZSBOOTSTRAP (<i>Ours</i>)	37.29	44.18	28.43	38.84	46.13	29.60	43.41	44.78	42.72
	S+S+C (<i>Ours</i>)	56.03	64.73	38.64	46.89	55.45	33.88	48.18	48.61	49.41
	S+S+C + $\mathcal{L}_{\leftrightarrow}$ (<i>Ours</i>)	55.46	64.45	38.54	46.56	55.96	34.61	48.53	48.71	49.84
FS	LAVT	74.46	76.89	70.94	65.81	70.97	59.23	63.34	63.62	63.66

However, note that for negative examples considered pairwise, the output of the network should only be distinct in locations where the instance chosen masks from Equation 3 are active, as it should be 0 elsewhere for both. Given the full matching δ^* , we can obtain the chosen mask for each input as $\mathbf{m}_{i,j,k} = \sum_c \mathbf{m}_{i,j,k}^c \delta_{j,k,c}^*$. Now consider a pair of negative examples for an object j, k and j^-, k^- . The "active" pixels in either of the chosen, pseudo-ground-truth masks as $\mathbf{A}(i, j, k, j^-, k^-) = \mathbf{m}_{i,j,k} \cup \mathbf{m}_{i,j^-,k^-}$. These are then used to slice each of the masks such that $\tilde{\mathbf{m}}_{i,j,k} = \hat{\mathbf{m}}_{i,j,k} |_{\mathbf{A}(i,j,k,j^-,k^-)}$ (and similarly for j^-, k^-). We can then write the contrastive loss as:

$$\mathcal{L}_{\leftrightarrow}(i, j, k) = \sum_{k^+ \neq k} \text{KL}(\hat{\mathbf{m}}_{i,j,k} \| \hat{\mathbf{m}}_{i,j,k^+}) + \sum_{j^- \neq j, k^-} [\gamma \text{KL}(\tilde{\mathbf{m}}_{i,j,k} \| \tilde{\mathbf{m}}_{i,j^-,k^-})]^{-1}, \quad (4)$$

where $\gamma \in \mathbb{R}^+$ is a hyperparameter to tune the balance of negative and positive examples in the contrastive term.

We present the full results in Table 5. As can be observed, the $\mathcal{L}_{\leftrightarrow}$ loss does not lead to a clear improvement of the results, which is why we have decided to omit it from the main method.