

# Interactive Network Perturbation between Teacher and Students for Semi-Supervised Semantic Segmentation

Hyuna Cho, Injun Choi, Suha Kwak, Won Hwa Kim  
Pohang University of Science and Technology (POSTECH)  
{hyunacho, surung9898, suha.kwak, wonhwa}@postech.ac.kr

## 1. Additional Implementation Details

In this section, we provide additional details for our implementation which were not included in the main manuscript due to the page limit. Weight decay was set to  $1e-4$ , and Sync-BN [3] is used for batch normalization. For a single-scale evaluation, the crop size is set to  $512 \times 512$  for PASCAL VOC 2012 and  $800 \times 800$  for Cityscapes. We used Online Hard Example Mining (OHEM) loss [5] as supervision losses with the ground truth (i.e.,  $L_{sup}^t$  and  $L_{sup}^s$ ) on Cityscapes for all baselines and our method as in CPS [1],  $n$ -CPS [2], and PS-MT [4].

## 2. Trade-off among hyperparameters: $\alpha, \beta, \gamma$

As shown in Table 1, we investigated various weight combinations on the pseudo-supervision losses of our method. Ablation study for  $\alpha, \beta, \gamma$ , i.e., the hyperparameters for  $L_{gps}$ ,  $L_{cps}$  and  $L_{fps}$ , were conducted. The best performances on both PASCAL VOC 2012 and Cityscapes were obtained with a relatively high weight on the  $L_{cps}$ . However, our proposed model is insensitive to the intensity of the hyperparameters, and the existence of all the three losses is significant to our method as reported in the ablation study (Table 4) of the main paper.

(a) PASCAL VOC 2012				(b) Cityscapes			
$\alpha$	$\beta$	$\gamma$	mIOU(%)	$\alpha$	$\beta$	$\gamma$	mIOU(%)
0.5	1.0	0.5	73.72	1	3	1	73.77
0.5	1.0	1.0	73.52	1	4	1	74.10
0.5	1.0	1.5	73.31	<b>1</b>	<b>5</b>	<b>1</b>	<b>74.62</b>
<b>0.5</b>	<b>1.5</b>	<b>0.5</b>	<b>73.95</b>	1	6	1	73.87
0.5	1.5	1.5	73.78	2	3	2	73.29
0.5	1.5	1.0	73.71	2	4	2	73.41
1.0	0.5	1.0	72.85	2	6	2	73.96
1.0	1.0	1.0	73.27	4	4	4	73.53
1.0	1.0	0.5	73.54				
1.0	1.0	1.5	72.98				
1.0	1.5	1.0	72.83				
1.5	0.5	1.0	72.85				
1.5	1.0	0.5	73.16				

Table 1. The results are obtained under 1/8 labeled partition protocol with ResNet-50. The comparison is performed without Cut-Mix.

## 3. Visualization of Interactions between TN and SN

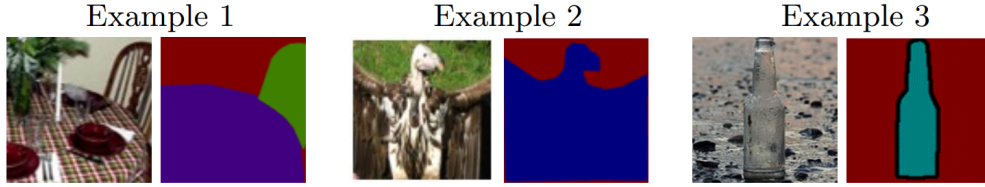
In one iteration, our method contains two losses (i.e.,  $L_{guide}$  and  $L_{fb}$ ) for a pair of Student Networks (SNs) and a Teacher Network (TN), respectively. To update the parameters of SNs, the  $L_{guide}$  uses pseudo-labels of TN to guide SNs, and concurrently uses pseudo-labels of SNs to teach each other. On the other hand, the  $L_{fb}$  takes the pseudo-labels from SNs to update the TN’s parameters.

In Fig 1, we visualize confidence maps and pseudo-label maps of a exemplary samples from the 2nd and the last epochs within a single iteration. The qualitative comparison is performed before and after the two losses are applied to update the SNs and the TN in turn. The pseudo-labels of SNs before computing  $L_{guide}$  are used for the pseudo supervision between the SNs. After minimizing  $L_{guide}$ , they serve as pseudo-labels for the TN working as the feedback pseudo supervision (FPS). With  $L_{guide}$ , we empirically observed significant improvement on the prediction qualities of the SNs, and the changes are notably shown in both of the confidence maps and pseudo-labels in the 2nd epoch.

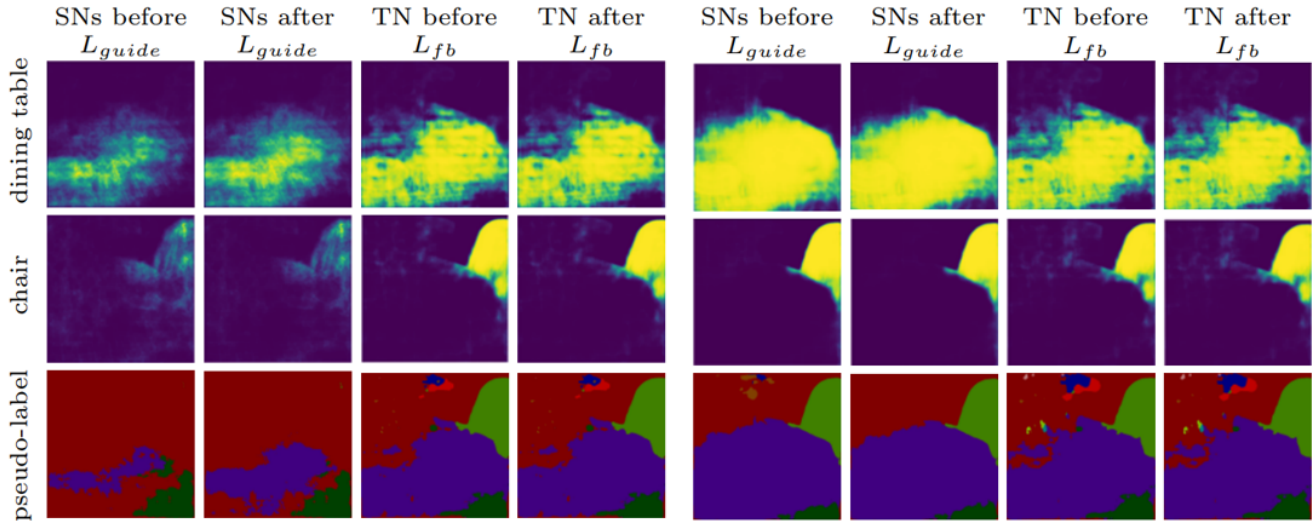
Regardless of  $L_{fb}$ , the TN’s pseudo-labels are employed as the guided pseudo supervision (GPS) to evaluate SNs’ predictions. The difference with the computation of  $L_{fb}$  is when they act as guides to SNs: before minimizing the  $L_{fb}$ , the TN’s pseudo-labels are used in the *current* iteration with  $L_{guide}$ , whereas the updated TN produces pseudo-labels in the *next* iteration. With Adaptive Ramp-Up, the performance of TN is maintained until the training ends, and thus the TN can consistently provide informative pseudo-labels to SNs. Due to our interactive loss operation between  $L_{guide}$  and  $L_{fb}$ , the results of SNs outperform the TN’s, and the superiority can be seen with clear and vivid confidence maps of SNs compared to the TN’s prediction at the last epoch. In summary, students become better than the teacher.

## 4. Pseudo Code

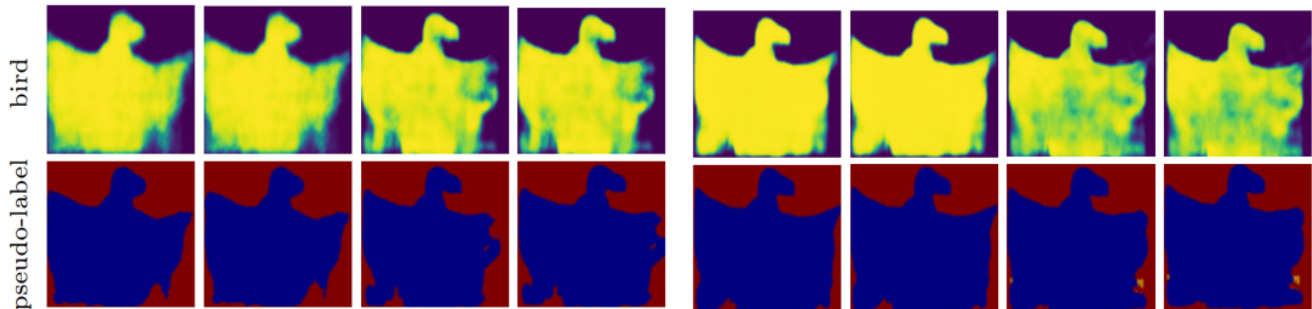
We present the pseudo code of our whole training process in Algorithm 1.



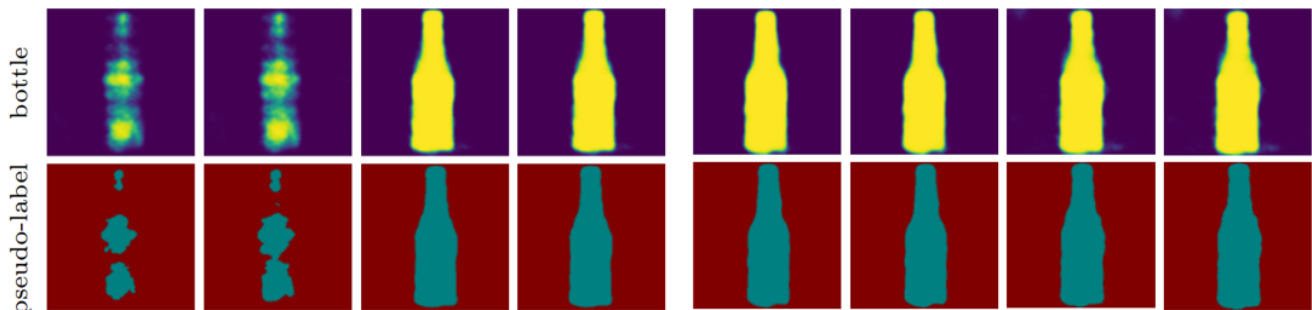
(a) Inputs and Ground Truths



Example 1



Example 2



Example 3

(b) 2-nd Epoch

(c) Last Epoch

Figure 1. **Visualization of Interactions between a TN and SNs** before and after applying loss functions in one iteration. The resultant pseudo-labels of SNs are forwarded from averaged confidence maps of the two SNs.

---

**Algorithm 1** Training networks with GPS

---

```
1: procedure TRAINING(epoch, batch, dataloader)
2:    $\sigma(\cdot) \leftarrow \text{Softmax}(\cdot)$ 
3:    $f_{\theta_t}(\cdot), f_{\theta_{s_1}}(\cdot), f_{\theta_{s_2}}(\cdot) \leftarrow \text{TN}, \text{SN}_1, \text{SN}_2$ 
4:    $w(e) = \exp(-1 - \frac{e}{\text{epoch}})^2 \leftarrow \text{Adaptive Ramp-Up function}$ 
5:   for each epoch do ▷ Establishing a Teacher Network
6:     for each batch do
7:        $D^l = \{X_i^l\}_{i=1}^n, D^u = \{X_j^u\}_{j=1}^m \leftarrow \text{dataloader.iter}()$ 
8:       for  $X_i^l \in D^l$  do
9:          $P_i^{l,t} \leftarrow \sigma(f(X_i^l; \theta_t))$ 
10:      end for
11:      Compute  $L_t^l$  using  $P^{l,t}, Y^l$ 
12:      Update  $\theta_t$ 
13:    end for
14:  end for
15:  for  $e$  in each batch do
16:    for each batch do
17:       $D^l = \{X_i^l\}_{i=1}^n, D^u = \{X_j^u\}_{j=1}^m \leftarrow \text{dataloader.iter}()$ 
18:      for  $X_i^l \in D^l$  and  $X_j^u \in D^u$  do ▷ Guiding Step, Students Update
19:         $P_i^{l,s_1}, P_i^{l,s_2} \leftarrow \sigma(f(X_i^l; \theta_{s_1})), \sigma(f(X_i^l; \theta_{s_2}))$ 
20:         $P_j^{u,t}, P_j^{u,s_1}, P_j^{u,s_2} \leftarrow \sigma(f(X_j^u; \theta_t)), \sigma(f(X_j^u; \theta_{s_1})), \sigma(f(X_j^u; \theta_{s_2}))$ 
21:         $Y_j^t, Y_j^{s_1}, Y_j^{s_2} \leftarrow \max(P_j^{u,t}), \max(P_j^{u,s_1}), \max(P_j^{u,s_2})$ 
22:      end for
23:      Compute  $L_s^l$  using  $P^{l,s_1}, P^{l,s_2}, Y^l$ 
24:      Compute  $L_{gps}$  using  $P^{u,s_1}, P^{u,s_2}, Y_j^t$ 
25:      Compute  $L_{cps}$  using  $P^{u,s_1}, P^{u,s_2}, Y_j^{s_1}, Y_j^{s_2}$ 
26:       $L_{guide} \leftarrow L_s^l + \alpha L_{gps} + \beta L_{cps}$ 
27:      Update  $\theta_{s_1}$  and  $\theta_{s_2}$ 
28:      for  $X_i^l \in D^l$  and  $X_j^u \in D^u$  do ▷ Feedback Step, Teacher Update
29:         $P_i^{l,t} \leftarrow \sigma(f(X_i^l; \theta_t))$ 
30:         $P_j^{u,t}, P_j^{u,s_1}, P_j^{u,s_2} \leftarrow \sigma(f(X_j^u; \theta_t)), \sigma(f(X_j^u; \theta_{s_1})), \sigma(f(X_j^u; \theta_{s_2}))$ 
31:         $Y_j^s \leftarrow \frac{1}{2}(P_j^{u,s_1} + P_j^{u,s_2})$ 
32:      end for
33:      Compute  $L_t^l$  using  $P^{l,t}, Y^l$ 
34:      Compute  $L_{fps}$  using  $P^{u,t}, Y^s$ 
35:       $L_{fb} \leftarrow L_t^l + w(e) \cdot \gamma L_{fps}$ 
36:      Update  $\theta_t$ 
37:    end for
38:    if  $0.5 * \sum^{\text{batch}} L_s^l \leq \sum^{\text{batch}} L_t^l$  then
39:       $w(e) = 1$ 
40:    end if
41:  end for
42: end procedure
```

---

## References

- [1] Xiaokang Chen, Yuhui Yuan, Gang Zeng, et al. Semi-supervised semantic segmentation with cross pseudo supervision. In CVPR, 2021. 1
- [2] Dominik Filipiak, Piotr Tempczyk, and Marek Cygan.  $n$ -cps: Generalising cross pseudo supervision to  $n$  networks for semi-supervised semantic segmentation. arXiv preprint arXiv:2112.07528, 2021. 1
- [3] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In ICML, 2015. 1
- [4] Yuyuan Liu, Yu Tian, Yuanhong Chen, et al. Perturbed and strict mean teachers for semi-supervised semantic segmentation. In CVPR, pages 4258–4267, 2022. 1
- [5] Abhinav Shrivastava, Abhinav Gupta, and Ross Girshick. Training region-based object detectors with online hard example mining. In CVPR, pages 761–769, 2016. 1