

The GIST and RIST of Iterative Self-Training for Semi-Supervised Segmentation

Eu Wern Teh^{1,2,3} Terrance DeVries^{1,2} Brendan Duke^{3,4} Ruowei Jiang³ Parham Aarabi^{3,4} Graham W. Taylor^{1,2}
¹University of Guelph ²Vector Institute ³Modiface, Inc. ⁴University of Toronto

Abstract—We consider the task of semi-supervised semantic segmentation, where we aim to produce pixel-wise semantic object masks given only a small number of human-labeled training examples. We focus on iterative self-training methods in which we explore the behavior of self-training over multiple refinement stages. We show that iterative self-training leads to performance degradation if done naïvely with a fixed ratio of human-labeled to pseudo-labeled training examples. We propose Greedy Iterative Self-Training (GIST) and Random Iterative Self-Training (RIST) strategies that alternate between training on either human-labeled data or pseudo-labeled data at each refinement stage, resulting in a performance boost rather than degradation. We further show that GIST and RIST can be combined with existing semi-supervised learning methods to boost performance.

Keywords—semi-supervised learning; semantic segmentation; self-training

I. INTRODUCTION

Semantic segmentation is the task of producing pixel-wise semantic labels over a given image. This is an important problem that has many useful applications such as medical imaging, robotics, scene-understanding, and autonomous driving. Supervised semantic segmentation models are effective, but they require tremendous amounts of pixel-wise labels, typically provided by a time consuming human annotation process. To overcome the need of collecting more pixel-wise labeled data, there has been an increase in interest in semi-supervised semantic segmentation in recent years [1], [2], [3], [4], [5], [6], [7], [8].

Self-training is a classic semi-supervised learning method that uses pseudo-labels to guide its learning process. We define pseudo-labels as predictions generated by a given model in contrast to human-provided annotations. Self-training means using a model’s own predictions as pseudo-labels in its loss during training. Recently, there has been a resurgence of self-training methods in semi-supervised learning [9], [10], [11], [12], [13]. Despite the recent comeback of self-training methods, most recent self-training works are confined to only one refinement stage.

Iterative self-training consists of multiple refinement stages, each consisting of K training iterations. At the beginning of each stage the model is initialized with weights from the previous stage, and pseudo-labels are regenerated (see Sec. III-B). We aim to investigate the behaviour of self-training after many (i.e. > 3) refinement stages for semi-supervised semantic segmentation.

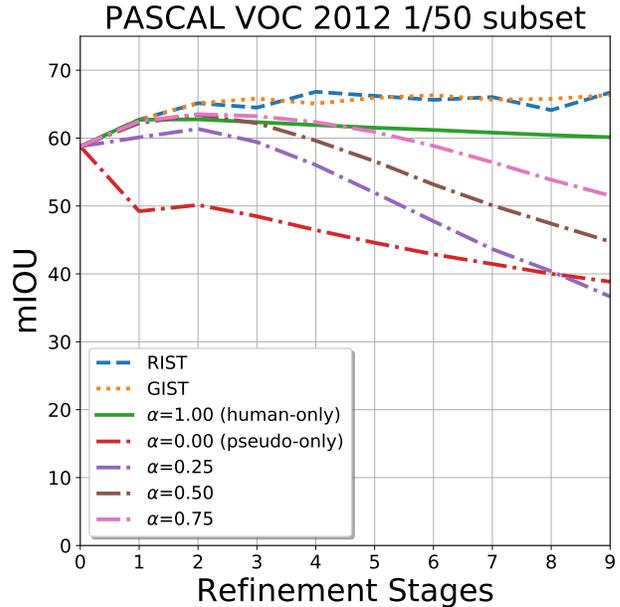


Figure 1. The performance of iterative self-training with various ratios of human-labels to pseudo-labels α on the PASCAL VOC 2012 validation datasets. Models are refined iteratively by bootstrapping on weights trained on a previous refinement stage, with only 2% of human-labels. A development set is used to select the best refinement stage.²

Iterative refinement on a small number of human-labels¹ may cause over-fitting on the training set, as no new information is introduced. Iterative refinement on pseudo-labels does introduce new information which can improve performance. However, it also results in a feedback loop that repeatedly reinforces and compounds incorrect predictions from previous iterations, ultimately resulting in “pseudo-label bloat”, where a single dominant class prediction spreads to cover an entire image eventually (see Figure 2).

The naïve solution of combining both human-labels and pseudo-labels in each batch slows the rate at which pseudo-label bloat occurs but does not combat it entirely. Instead, we find that alternating training on only human-labels or only pseudo-labels results in a more controlled training dynamic where pseudo-labels help expand predictions to regions that may have been missed, while human-labels prevent pseudo-labels from drifting too far away from the expected

¹For convenience, we refer to human-labeled training examples as “human-labels” and pseudo-labeled training examples as “pseudo-labels”.

annotations. By switching between these two extremes in a greedy fashion (GIST) or random fashion (RIST), our model enjoys the benefits of both label types, ultimately yielding better performance (see Figure 1). Our contributions are the following:

- We show that naïve application of iterative self-training to the problem of semi-supervised segmentation via a fixed human-labels to pseudo-labels ratio results in significant performance degradation when $\alpha < 1$.
- We introduce Greedy Iterative Self-Training (GIST) and Random Iterative Self-Training (RIST) to overcome performance degradation through iteratively training on either human- or pseudo-labels.
- We demonstrate that both RIST and GIST can improve existing semi-supervised learning methods, yielding performance boost in both the PASCAL VOC 2012 and City- scapes datasets across all eight subsets.

II. RELATED WORK

Semi-supervised semantic segmentation has gained ground in recent years. Souly et al. [1] extend a typical Generative Adversarial Network (GAN) network by designing a discriminator that accepts images as input and produces pixel-wise predictions, consisting of confidence maps for each class and a fake label. Hung et al. [2] improve GAN-based semantic segmentation by using a segmentation network as a conditional generator. They also redesign the discriminator to accept the segmentation mask as input and restrict it to produce pixel-wise binary predictions. Mittal et al. [3] extend Hung et al.’s network by making the discriminator produce an image-level binary prediction. They also add a feature matching loss and self-training loss in their training pipeline. They further propose a separate semi-supervised classification network [14] to clean up segmentation masks’ prediction.

Recently, a few works propose non-GAN based solutions for semi-supervised semantic segmentation. Mendel et al. [4] propose an error-correcting network that aims to fix the predictions of the main segmentation network. This error-correcting network is also applied to unlabeled images to correct the generated mask. French et al. [5] adapt CutMix [15] to augment and regularize a segmentation network by creating composite images via mixing two different images and their corresponding segmentation masks. Similarly, Olsson et al. [6] also aim to regularize the segmentation network by mixing segmentation masks from two images by selecting half of the classes from one image and the other half from another image. Ouali et al. [7] minimize consistency loss between features of multiple auxiliary decoders, in which the perturbed encoder outputs are fed. Alonso et al. [8] improve the supervision

²A development set is the set of additional images used for meta-parameter selection (Sec 4 paragraph 1).

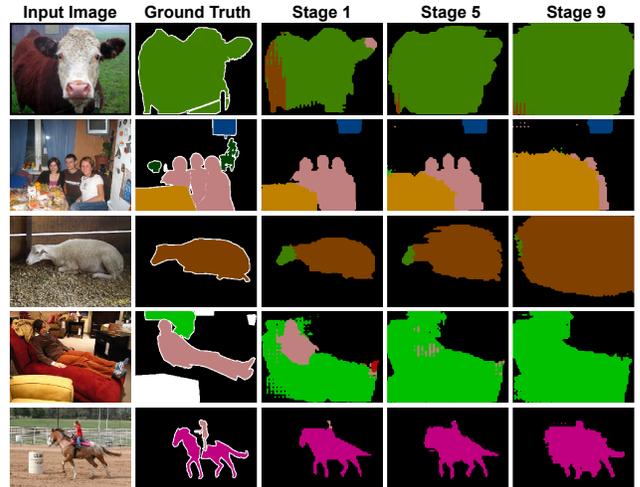


Figure 2. Pseudo-label degradation when a model is trained iteratively with a fixed human-labels to pseudo-labels ratio ($\alpha = 0.75$). The first column consists of input images. The second column consists of ground truth labels. The third to fifth columns show pseudo-labels generated after at refinement stage 1, 5 and 9.

signals from their teacher model by storing more samples in a separate memory banks.

We explore the iterative self-training method to tackle semi-supervised semantic segmentation. Self-training or pseudo-labeling is a classic semi-supervised learning recipe that can be traced back to 1996, where it was used in an NLP application [16]. A major benefit of self-training is that it allows easy extension from an existing supervised model without discarding any information. In the deep learning literature, Lee et al. [9] popularized self-training in semi-supervised classification. After its reappearance, it has gained traction in recent years. Zhai et al. [11] show the effectiveness of mixing self-supervised and semi-supervised learning along with pseudo-labeling in their training regime. Zoph et al. [12] show that a randomly initialized model with self-training via a joint-loss can yield better performance than a model initialized with a pre-trained model without self-training. Xie et al. [13] show that iterative self-training with noisy labels improves a classification model’s accuracy and robustness. Radosavovic et al. [10] use self-training in Omni-supervised learning, where they generate pseudo-labels by taking the average prediction of multiple perturbations of a single unlabeled image. Most of the recent self-training works [9], [10], [11], [12] are confined to only one refinement stage, with the exception of Xie et al.’s work [13], where they benefit from iterative self-training by repeating self-training for three stages of refinement. In this work, we aim to extend self-training for semi-supervised semantic segmentation and investigate self-training behavior under many (i.e. > 3) refinement stages.

III. METHODS

Self-training is a semi-supervised learning method that uses pseudo-labels to guide its learning process. As we improve the model, we also improve the quality of pseudo-labels. Self-training typically consists of the following steps: (1) training a model using human-labeled data; (2) generating pseudo-labels using the trained model; and (3) finetuning the trained model with the combination of human-labeled data and pseudo-labels.

A. Combining both human-labels and pseudo-labels in training

Given a set of images $\{(x_i, y_i)\}$, where x_i represents the image and y_i represents the corresponding human-label, we begin training a segmentation model, SEG for K iterations using a 2D Cross-Entropy Loss, ENT on available human-labeled data. The loss for iteration j is computed as:

$$L_j = \frac{1}{B} \sum_i^B \text{ENT}(o_i, y_i), \quad (1)$$

$$o_i = \text{SEG}(x_i).$$

where B represents the batch size, i indexes examples in a batch, and o_i represents a model's output.

After a model is trained on available human-labeled data, we can now use this pre-trained model to generate pseudo-labels on unlabeled data. Given another set of images $\{(x_i^p, y_i^p)\}$, where x_i^p represents the unlabeled image and y_i^p represents the corresponding pseudo-label, we can now combine human-labels and pseudo-labels by a linear combination of respective losses computed by Eq. 1.

$$L_j^\alpha = \frac{1}{B} \sum_i^B (\text{ENT}(o_i, y_i) * \alpha + \text{ENT}(o_i^p, y_i^p) * (1 - \alpha)). \quad (2)$$

We define α as the ratio of human-labels to pseudo-labels. Eq. 1 is a small modification to the classic self-training loss by Lee et al. [9], who apply a coefficient only to the unlabeled loss term, where that coefficient is iteration-dependent.

B. Iterative Self-training

Self-training can be repeated through multiple refinement stages, where each refinement stage consists of a pseudo-label generation step and a finetuning step. A naïve solution is to fix α throughout all refinement stages, which we call Fixed Iterative Self-Training (FIST). Given that α can take a floating-point number between zero and one, there are infinitely many possible α values at each stage of refinement. By making α binary, we turn the problem into discrete path selection. In this setting, our goal is to find the sequence of stages which yields the best solution. We explore two different selection strategies: Greedy Iterative Self-Training (GIST) and Random Iterative Self-Training (RIST).

We define S as the maximum number of refinement stages and K as the maximum number of training iterations at a given stage. The number of possible paths in the search space (2^S) is exponential in the number of refinement stages.

Algorithm 1 GIST

```

1:  $\theta_{list} \leftarrow [\theta_0]$ ;  $\alpha_{list} \leftarrow [0, 1]$ 
2: for  $s \leftarrow 1 \dots S$  do
3:    $\theta_{list}^* \leftarrow []$ ;  $R_{list} \leftarrow []$ 
4:   for  $\theta_c$  in  $\theta_{list}$  do
5:     for  $\alpha$  in  $\alpha_{list}$  do
6:       for  $m \leftarrow 1 \dots M$  do
7:          $y^p \leftarrow \arg \max(\text{SEG}(x_m^p))$ 
8:       end for
9:       for  $j \leftarrow 1 \dots K$  do
10:         $L_j^\alpha \leftarrow \frac{1}{B} \sum_i^B (\text{ENT}(o_i, y_i) * \alpha + \text{ENT}(o_i^p, y_i^p) * (1 - \alpha))$ 
11:         $\theta_c \leftarrow \theta_c - \lambda \frac{\partial L_j^\alpha}{\partial \theta_c}$ 
12:      end for
13:      append  $\theta_c$  to  $\theta_{list}^*$ 
14:      append EVAL( $\theta_c$ ) to  $R_{list}$ 
15:    end for
16:  end for
17:  sort  $\theta_{list}^*$  based on  $R_{list}$  in descending order
18:  keep top  $G$  items in  $\theta_{list}^*$  and assign it to  $\theta_{list}$ 
19: end for

```

Algorithm 2 RIST

```

1: for  $s \leftarrow 1 \dots S$  do
2:    $\alpha \leftarrow (\text{Uniform}(0, 1) > 0.5)$ 
3:    $\theta_s \leftarrow \theta_{s-1}$ 
4:   for  $m \leftarrow 1 \dots M$  do
5:      $y^p \leftarrow \arg \max(\text{SEG}(x_m^p))$ 
6:   end for
7:   for  $j \leftarrow 1 \dots K$  do
8:      $L_j^\alpha \leftarrow \frac{1}{B} \sum_i^B (\text{ENT}(o_i, y_i) * \alpha + \text{ENT}(o_i^p, y_i^p) * (1 - \alpha))$ 
9:      $\theta_s \leftarrow \theta_s - \lambda \frac{\partial L_j^\alpha}{\partial \theta_s}$ 
10:  end for
11: end for

```

GIST works by finding the best stage of refinement by relying on the development set. Algorithm 1 describes the GIST algorithm, which is a beam search strategy. In line 3, we keep a list of candidates θ , and the corresponding evaluation results, R_{list} , using the EVAL function on the development set. These lists are updated in lines 13 and 14. In line 17, we sort θ_{list}^* in descending order based on R_{list} . In line 18, we keep the top G θ in the θ_{list}^* , and assign it to θ_{list} , where it will be used as the initial model for the next stage of refinement. In Line 6 to 8, we generate our pseudo-labels for M unlabeled images. In Line 9 to 12, we finetune a segmentation model for K iterations with learning rate λ . Figure 3 shows a hypothetical α path selection scenario using GIST.

One weakness of GIST is its smaller search space when compared to a random search (RIST). Bergstra et al. [17] show that in low dimensions, random search is an effective search strategy compared to a grid search. One can see that a greedy search is a subset of a grid search where we only explore the top performing branches. With a beam size of 1, the search space for a greedy search is $\log(S)$, and the search space for a random search is 2^S . At first glance, RIST

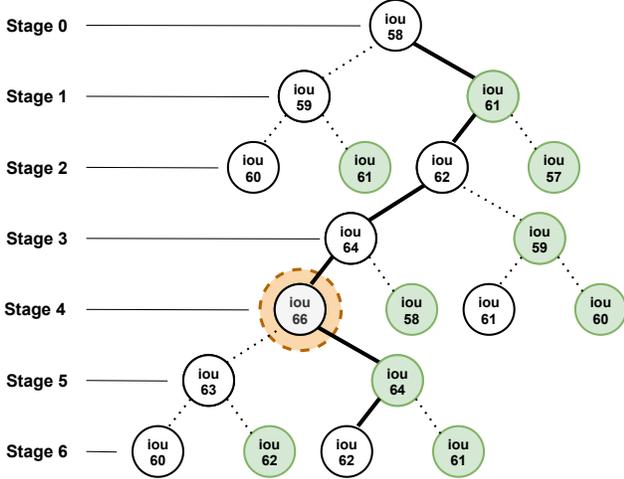


Figure 3. A hypothetical α path selection scenario with six refinement stages in self-training using the greedy approach (GIST). α indicates the ratio of human-labels to pseudo-labels. The open nodes indicate that only human-labels ($\alpha = 1$) are being used for training and the shaded nodes indicate only pseudo-labels ($\alpha = 0$) are being used. The number of possible paths is exponential in the number of refinement stages (2^6). At each stage of refinement, we evaluate the mean intersection over union (mIoU) of a model using a development set. Here, the optimum value is found at stage four of the refinement process.

may seem counter-intuitive, but RIST’s performance spread is relatively small (± 0.93), and this spread can be reduced by eliminating obvious degenerate solutions and increasing search paths (see Sec. IV-A for details). Algorithm 2 describes the RIST algorithm. In line 2, we randomly set α to either zero or one. Line 7 to 10 is similar to the GIST algorithm (line 9 to line 12).

In GIST, the time complexity of a beam search is $O(S * G)$. A beam search is only as efficient as a random search if we can parallelize the training for each α at each stage, and this condition requires a beam search to be trained on consecutive numbers of GPU resources. Unlike a beam search, each random search can be trained independently on a single GPU resource. It is easier and cheaper to obtain N independent GPUs rather than N consecutive GPU resources, making RIST computationally cheaper and faster in practical settings.

C. Additional Add-ons

In addition to cross entropy loss in Eq 2, we include three additional add-ons to boost the performance of FIST, GIST and RIST: Consistency Loss (CL), Label Erase (LE), and Temperature Scaling (TS). Table V shows the ablation study of the add-ons for both RIST and GIST for models trained on the Pascal VOC 1/50 subset.

Consistency Loss (CL): Consistency Loss (CL) is a common loss in semi-supervised learning where the goal is to minimize features between two perturbations of the same

input. CL is a common loss in semi-supervised learning [3], [6]. Mittal et al. [3] use CL via a feature matching loss to minimize the discrepancy between predicted features and ground truth features. Olsson et al. [3] use CL via mean-teacher [14] method. We also use the mean-teacher as our CL loss. We use the following equations (Equations 3 to 7) to calculate CL for both x_i and x_i^p .

$$\theta^* = \theta^* \beta + \theta * (1 - \beta) \quad (3)$$

$$o_i, f_i = \text{SEG}(x_i) \quad (4)$$

$$o_i^*, f_i^* = \text{SEG}^*(x_i) \quad (5)$$

$$f_i = \text{drop}(\text{pool}(f_i)) \quad (6)$$

$$f_i^* = \text{drop}(\text{pool}(f_i^*)) \quad (7)$$

$$L_{\text{feature}} = |f_i - f_i^*| \quad (8)$$

Following [3], [2], we use DeepLabV2 [18] as our segmentation model; therefore, f_i represents features before an Atrous Spatial Pyramid Pooling (ASPP) layer. Equation 3 describes the weight update rule for the teacher model, SEG^* , where an exponential moving average rule is applied to its weights, θ^* , which is controlled by β , ($0 \leq \beta \leq 1$). Equations 4 and 5 illustrate the extraction of features before the ASPP layer in both the student model, SEG , and the teacher model, SEG^* . In Equations 6 and 7, we apply global average pooling followed by a dropout perturbation to these features. After that, we compute the absolute differences between these two features. Figure 4 summarizes the feature consistent loss in our model.

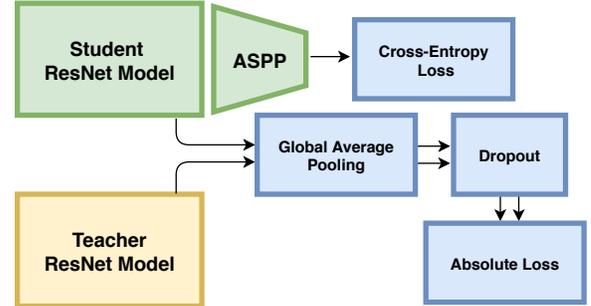


Figure 4. Depiction of the Consistency Loss (CL). At each batch of training, there are two copies of Resnet Models. The student model is the main model that is used to predict the segmentation mask, and the teacher model is a copy of the student model with an exponential moving weight update. The goal of CL is to minimize the differences between the features of the student and teacher models.

Label Erase (LE): Pseudo-labels are often very noisy, especially when we have a small amount of human-labeled data. To remove noise from pseudo-labels, we only keep predictions that pass a certain confidence threshold, ϕ . Label Erase (LE) is also used in [3], [6], [10], [12].

Equation 9 illustrates the process of flagging low confidence prediction regions so that they are ignored by the loss

function in Equation 1. We define pixel-wise confidence, $c_{i,j}^p$ as the pixel softmax output of $o_{i,j}^p$, where i represents the image index and j represents the pixel index.

$$y_{i,j}^p = \begin{cases} \arg \max(c_{i,j}^p), & \text{if } \max(c_{i,j}^p) \geq \phi, \\ \text{ignore label}, & \text{if } \max(c_{i,j}^p) < \phi, \end{cases} \quad (9)$$

Temperature Scaling (TS): Temperature scaling (TS) is introduced by Hinton et al. [19], where it is used to create a softer probability distribution for knowledge distillation. The formula for TS is defined as $q_i = \frac{\exp(y_i * T)}{\sum_j \exp(y_j * T)}$. As T becomes smaller, the output of the softmax function will tend towards uniform distribution. We employ TS [19] to overcome over-confident prediction in our model, where the activation values after softmax are highly skewed towards 100% (see Figure 5).

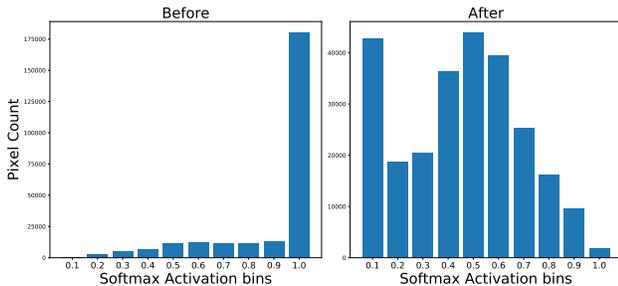


Figure 5. The before and after effects of applying temperature scaling to the output activation of a single image. A temperature scale of 0.2 is applied to the output activation.

IV. EXPERIMENTS

We perform experiments on two semantic segmentation datasets: PASCAL VOC 2012 [20] and Cityscapes [21]. In each dataset, there are three subsets, where pre-defined ratios (1/50, 1/20, and 1/8) of training images are selected as images with human-labels. We also experiment on two additional subsets (1/30 and 1/4) for the Cityscapes dataset. These subsets of labeled images are selected using the same split as [2], [3], [6]. We treat the remaining images as unlabeled examples. We use the mean intersection-over-union (mIoU) as a performance metric. The validation images for both datasets are set aside and used for evaluation, which is consistent with [2], [3], [6]. We select 50 additional images from the training set as a modest “development set” for meta-parameter tuning.

For all of our experiments, we follow the same experimental setup as [2], [3], [6]. We use a DeepLabV2 [18] segmentation model that is initialized with MS-COCO pre-trained weights [22]⁵. We also freeze all the BatchNorm

⁵DeepLabV2 backbone is used in our experiments so that our method is comparable to S4GAN and ClassMix.

Layers in DeepLabV2. We optimize our model using the Stochastic Gradient Descent (SGD) optimizer with a base learning rate of 2.5e-4, a momentum of 0.9, and a weight decay of 5e-4. We use a polynomial learning rate decay policy, where we adjust the learning rate with the following equation: $\lambda_{\text{iter}} = \lambda_0(1 - \frac{\text{iter}}{\text{max_iter}})^{0.9}$ where λ_0 is a base learning rate. To augment the dataset, we use random-cropping (321 × 321 for PASCAL VOC 2012 and 256 × 512 for Cityscapes), horizontal-flipping (with a probability of 0.5), and random-resizing (with a range of 0.5 to 1.5) in all of our experiments.

For all subsets, we train our supervised models and stage-0 models for 25,000 iterations. Additionally, for the PASCAL VOC 2012 dataset, we use a batch-size of 8, and we refine our model for 3,000 iterations at each refinement stage (Stage 1 to 9). For the Cityscapes dataset, we use a batch-size of 6, and we refine our models for 4,000 iterations at each refinement stage (Stage 1 to 9). We use a search cost of two for both GIST and RIST and use the development set to select the best path for all experiments.

Table I shows the results of our experiments as well as relevant results that others have reported [2], [3], [5], [6]. We use the code provided by the respective authors for our experiments with S4GAN [3] and ClassMix [6]. For a fair comparison, we set the batch size for S4GAN and ClassMix to match with our experiments, and we also select the best iterations using our development set. We notice performance differences in S4GAN and ClassMix when compared to performances reported in the original papers. We speculate that the differences are caused by best iteration selection and batch sizes. For experiment with S4GAN+GIST/RIST and ClassMix+GIST/RIST, we first train the segmentation model with S4GAN and ClassMix algorithm (stage-0 models). After that, we further refine the segmentation model using the GIST/RIST algorithm by bootstrapping on the DeepLabV2 model trained with S4GAN/ClassMix.

A. Discussion

Figure 1 shows that a naïve application of iterative self-training leads to significant performance degradation in both datasets. Figure 2 shows the qualitative evidence of performance degradation in FIST at $\alpha = 0.75$. FIST suffers from over-confident pseudo-label predictions which spread to the surrounding pixel. Over multiple stages of refinement, the pseudo-labels expand and eventually engulf most of the image. We speculate that this may be why most of the recent self-training works are confined to one refinement stage.

Figure 1 also shows that both RIST and GIST overcome performance degradation. Additionally, both RIST and GIST generalize better than FIST in each successive refinement. Table I shows that both RIST and GIST can improve other semi-supervised segmentation techniques such as S4GAN and ClassMix. We show that both RIST and GIST can further refine S4GAN and ClassMix yielding performance

	VOC 2012 ($\approx 10k$ images)			Cityscapes ($\approx 3k$ images)				
# of labeled images	211	529	1,322	59	100	148	371	743
Subset	1/50	1/20	1/8	1/50	1/30	1/20	1/8	1/4
Supervised	54.15	62.94	67.44	49.68	53.96	54.71	59.90	62.21
Supervised + GIST	66.33	66.95	70.27	53.51	56.38	58.11	60.94	63.04
Supervised + RIST	66.71	68.28	69.90	53.33	56.28	57.81	61.38	63.92
S4GAN [3]	62.87	62.35	68.56	50.48	54.58	55.61	60.95	61.30
S4GAN [3] + GIST	67.21	68.50	70.61	52.36	57.18	57.40	61.27	64.24
S4GAN [3] + RIST	66.51	68.50	70.31	53.47	57.12	57.48	62.50	64.64
ClassMix [6]	63.63	66.74	66.14	52.14	57.02	58.77	61.56	63.90
ClassMix [6] + GIST	65.60	69.05	70.65	52.43	58.70	59.98	62.44	64.53
ClassMix [6] + RIST	66.30	69.40	70.76	53.05	58.55	59.54	62.57	65.14

Table I
SEMANTIC SEGMENTATION RESULTS (mIoU) ON THE PASCAL VOC 2012 AND CITYSCAPES VALIDATION DATASETS.

boosts across all subsets in the Pascal VOC 2012 and Cityscapes datasets. Figures 8 and 9 show RIST and GIST’s qualitative results that are trained with 2% of human-labels in both datasets.

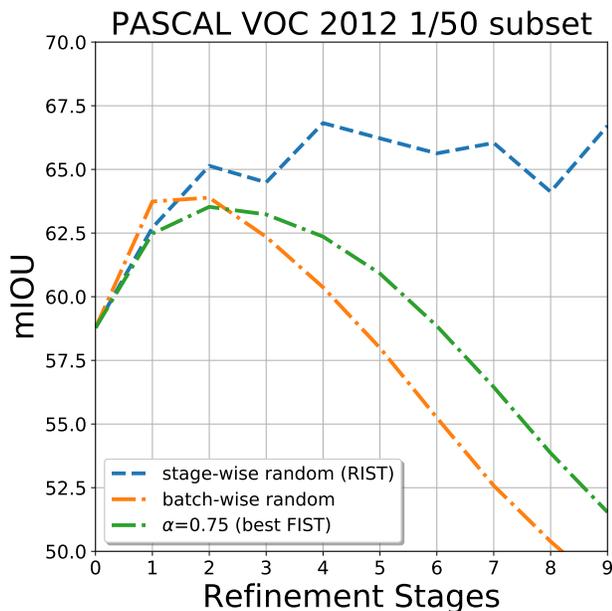


Figure 6. Self-training performance at various stages between batch-wise and stage-wise random selection strategies on the PASCAL VOC 2012 validation set.

Figure 6 demonstrates the importance of training on a single label type (i.e., human-label or pseudo-label) for an extended number of iterations. We find that randomly selecting the label type for each batch (batch-wise random) performs just as poorly as training with both label types in each batch (FIST). We speculate that the clean human-labels and the noisy pseudo-labels represent competing objectives, which are difficult for the model to satisfy simultaneously, resulting in it getting stuck at poor solutions. By applying stage-wise training, we allow the model to focus on a

single objective at a time, potentially escaping sub-optimal solutions from a previous stage.

Selection	mIoU (devel)	mIoU (val)	Group mIoU (val)
PPLLPLPL	54.35	66.14	67.03
LPLLLLPPL	55.05	66.71	
PPLLPPPLP	54.21	66.05	
LPLPPLLPL	55.12	67.03	
LLLLLLPPP	54.43	63.75	
PPPPPLLP	50.43	64.37	66.63
PLLLLPLPL	54.00	66.02	
LLPLPLLPL	55.19	66.24	
LPLPLLPLP	55.29	66.63	
LLPLPPLP	54.88	65.17	
PLPLPPPLP	54.64	66.40	66.62
LPPLPPLPP	54.46	66.83	
LPLLLLLPP	54.67	65.61	
LPPLLLLLP	54.82	66.62	
LPPPLLPLP	54.64	66.61	
Mean \pm 1 std. dev.		66.01 \pm 0.93	66.76 \pm 0.23

Table II
RIST SEMANTIC SEGMENTATION RESULTS ON THE PASCAL VOC 2012 VALIDATION SET. GROUP EXPERIMENT RESULTS ARE SELECTED BASED ON THE BEST DEVELOPMENT mIoU. RANDOM SELECTION CHOICES ARE P (PSEUDO-LABEL ONLY) OR L (HUMAN-LABEL ONLY). THERE ARE A TOTAL OF NINE REFINEMENT STAGES ORDERED SEQUENTIALLY FROM LEFT TO RIGHT.

Table II examines the stability of performance for RIST on the PASCAL VOC 2012 1/50 subset. We use the same subset for training and generate fifteen different permutations of binary α values uniformly at random. The degenerate solutions occur in row 5 and 6, where we have more than four consecutive numbers of the same α choice during training. If we were to perform a random search once, the standard deviation of mIoU is 0.93. If we were to remove the obvious degenerate solutions (row 5 and 6) using some heuristics, the standard deviation of mIoU is reduced to 0.52. Nevertheless, we can reduce this standard deviation further by selecting the best solution out of five different random solutions yielding a standard deviation of 0.23. Figure 7 shows that as we increase the number of random solutions,

the spread decreases.

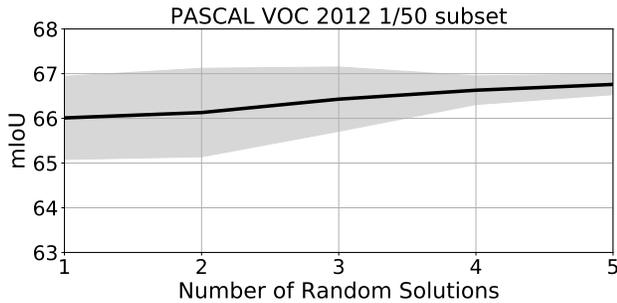


Figure 7. The effect of increasing stability of RIST when we increase the number of random solution. The best results of the random solutions are selected based on best development mIoU. The shaded area represent one standard deviation of uncertainty.

Sample Size	mIoU (RIST)	mIoU (GIST)
10	66.04	66.23
25	65.63	62.76
50	66.71	66.33
100	66.71	66.67
200	66.71	66.69
500	66.71	67.00

Table III

SEMANTIC SEGMENTATION RESULTS ON THE PASCAL VOC 2012 VALIDATION SET FOR RIST AND GIST BASED ON BEST EPOCH SELECTED WHILE VARYING THE NUMBER OF EXAMPLES IN THE DEVELOPMENT SET.

Beam Size	Search Cost	mIoU (devel)	mIoU (val)	Solution
1	2	55.17	66.33	LPLLPLLP
2	4	55.46	64.95	LLPLLPLLP
3	6	55.46	64.95	LLPLLPLLP

Table IV

SEMANTIC SEGMENTATION RESULTS ON THE PASCAL VOC 2012 VALIDATION SET FOR GIST ON DIFFERENT BEAM SIZE. SELECTION CHOICES AT EACH STAGE ARE P (PSEUDO-LABEL ONLY) OR L (HUMAN-LABEL ONLY).

Table III explores the sensitivity of RIST and GIST’s performance on the PASCAL VOC 2012 1/50 subset with respect to the size of the development set. In general, RIST and GIST are relatively stable. RIST performance remains the same for 50 to 500 sample sizes. GIST performance improves slightly as we increase the sample size from 50 to 500. On the Pascal VOC 2012 dataset, we show that our supervised+GIST and supervised+RIST trained with 211 human-labels (1/50 subset) outperform the supervised model that is trained with 529 human-labels (1/20 subset). This result shows that GIST and RIST improve the model,

not just because they got more supervised signals from the development set (see Table I).

Table IV explores GIST at various beam sizes. GIST can find a better solution for the development set; however, since there is a mismatch between the distribution of the development set and the original validation set due to the small sample size, the best solution of the development set is not the best solution for the original validation set. This study shows that GIST has a higher chance to overfit the development set when compared to RIST.

Method	RIST	GIST
no add-on	60.80	58.23
+CL	61.72	62.37
+CL +LE	64.69	64.56
+CL +LE +TS	66.71	66.33

Table V

ABLATION STUDY OF THE ADD-ONS FOR BOTH RIST AND GIST FOR MODELS TRAINED ON THE PASCAL VOC 1/50 SUBSET. RESULTS ARE REPORTED IN MIOU ON THE VALIDATION SET.

V. CONCLUSION

We show that iterative self-training with a fixed human-labels to pseudo-labels ratio (FIST) leads to performance degradation. This degradation can be overcome by alternating training on only human-labels or only pseudo-labels in a greedy (GIST) or random (RIST) fashion. A clear benefit of self-training is that it can easily extend existing architectures. We show that both GIST and RIST can further refine models trained with other semi-supervised techniques resulting in a performance boost.

REFERENCES

- [1] N. Souly, C. Spampinato, and M. Shah, “Semi supervised semantic segmentation using generative adversarial network,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 5688–5696.
- [2] W. C. Hung, Y. H. Tsai, Y. T. Liou, Y. Y. Lin, and M. H. Yang, “Adversarial learning for semi-supervised semantic segmentation,” in *29th British Machine Vision Conference, BMVC 2018*, 2019.
- [3] S. Mittal, M. Tatarchenko, and T. Brox, “Semi-supervised semantic segmentation with high-and low-level consistency,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019.
- [4] R. Mendel, L. A. de Souza Jr, D. Rauber, and J. Paulo, “Semi-supervised segmentation based on error-correcting supervision.”
- [5] G. French, T. Aila, S. Laine, M. Mackiewicz, and G. Finlayson, “Semi-supervised semantic segmentation needs strong, high-dimensional perturbations,” *arXiv preprint arXiv:1906.01916*, 2019.

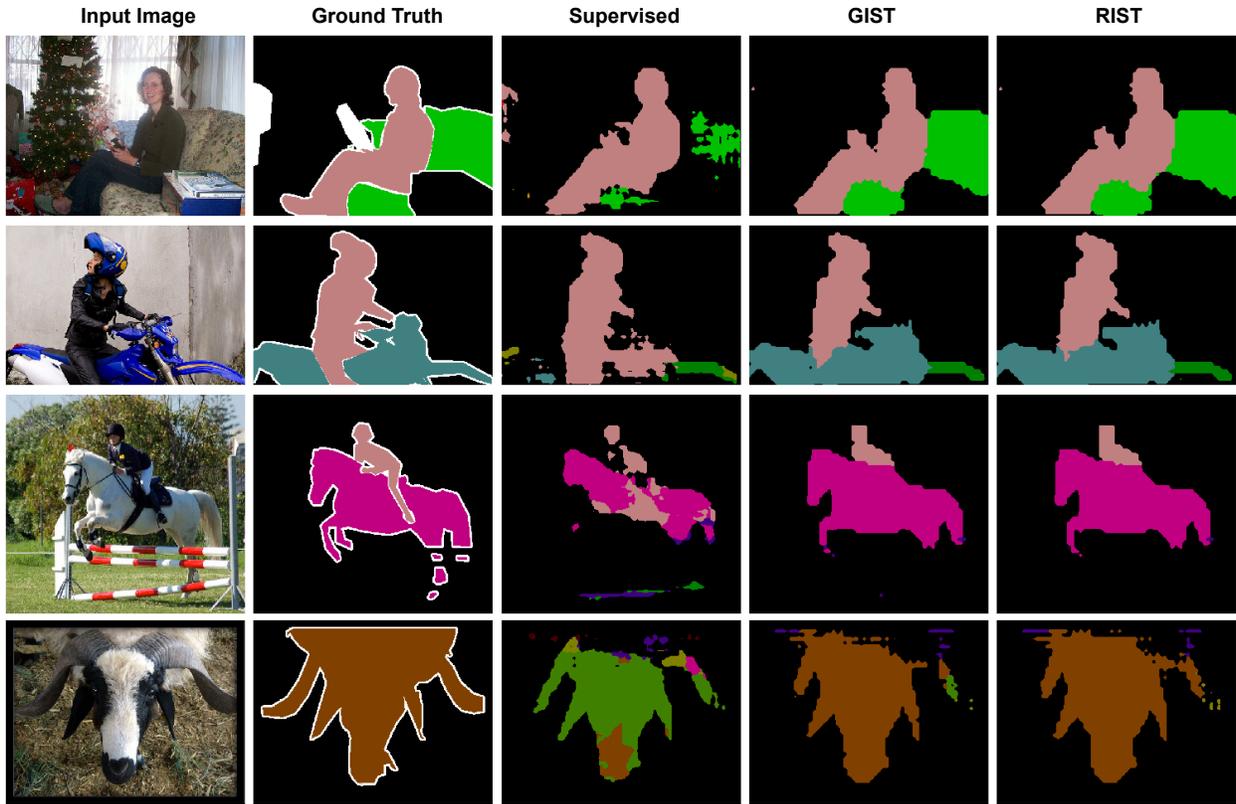


Figure 8. The qualitative results of our model train on 2% human-labels from the PASCAL VOC 2012 dataset.

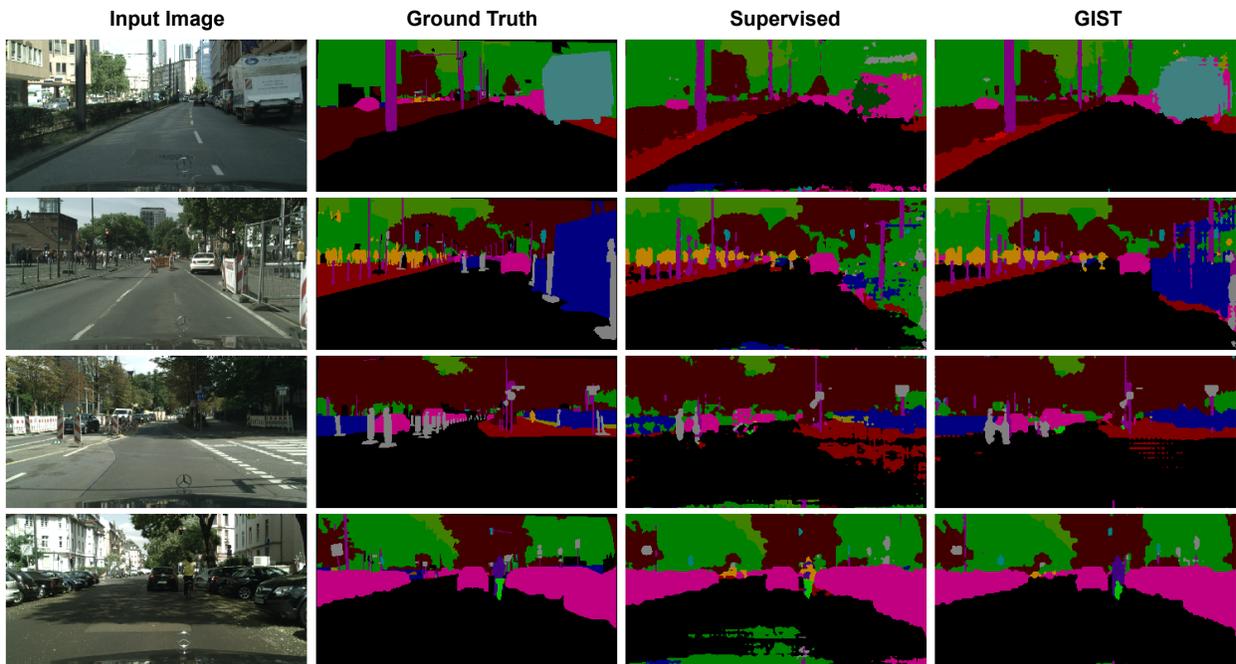


Figure 9. The qualitative results of our model trained on a 2% human-labels from the Cityscapes dataset.

- supervised learning,” 2020.
- [7] Y. Ouali, C. Hudelot, and M. Tami, “Semi-supervised semantic segmentation with cross-consistency training,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [8] I. Alonso, A. Sabater, D. Ferstl, L. Montesano, and A. C. Murillo, “Semi-supervised semantic segmentation with pixel-level contrastive learning from a class-wise memory bank,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 8219–8228.
- [9] D.-H. Lee, “Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks,” in *Workshop on challenges in representation learning, ICML. Vol. 3. No. 2. 2013*.
- [10] I. Radosavovic, P. Dollár, R. Girshick, G. Gkioxari, and K. He, “Data distillation: Towards omni-supervised learning,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 4119–4128.
- [11] X. Zhai, A. Oliver, A. Kolesnikov, and L. Beyer, “S4l: Self-supervised semi-supervised learning,” in *Proceedings of the IEEE international conference on computer vision*, 2019, pp. 1476–1485.
- [12] B. Zoph, G. Ghiasi, T.-Y. Lin, Y. Cui, H. Liu, E. D. Cubuk, and Q. V. Le, “Rethinking pre-training and self-training,” *arXiv preprint arXiv:2006.06882*, 2020.
- [13] Q. Xie, M.-T. Luong, E. Hovy, and Q. V. Le, “Self-training with noisy student improves imagenet classification,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10 687–10 698.
- [14] A. Tarvainen and H. Valpola, “Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results,” in *Advances in neural information processing systems*, 2017, pp. 1195–1204.
- [15] S. Yun, D. Han, S. J. Oh, S. Chun, J. Choe, and Y. Yoo, “Cutmix: Regularization strategy to train strong classifiers with localizable features,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 6023–6032.
- [16] D. Yarowsky, “Unsupervised word sense disambiguation rivaling supervised methods,” in *33rd annual meeting of the association for computational linguistics*, 1995, pp. 189–196.
- [17] J. Bergstra and Y. Bengio, “Random search for hyperparameter optimization.” *Journal of machine learning research*, vol. 13, no. 2, 2012.
- [18] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, “Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 4, pp. 834–848, 2017.
- [19] G. Hinton, O. Vinyals, and J. Dean, “Distilling the knowledge in a neural network,” *arXiv preprint arXiv:1503.02531*, 2015.
- [20] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman, “The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results,” <http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html>.
- [21] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, “The cityscapes dataset for semantic urban scene understanding,” in *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [22] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft coco: Common objects in context,” in *European conference on computer vision*. Springer, 2014, pp. 740–755.