

Diverse Cotraining Makes Strong Semi-Supervised Segmentor

Yijiang Li^{1†} Xinjiang Wang² Lihe Yang^{3†} Litong Feng² Wayne Zhang² Ying Gao⁴

¹Johns Hopkins University ²SenseTime Research ³The University of Hong Kong ⁴South China University of Technology
yli556@jhu.edu, wangxinjiang@sensetime.com

Abstract

Deep co-training has been introduced to semi-supervised segmentation and achieves impressive results, yet few studies have explored the working mechanism behind it. In this work, we revisit the core assumption that supports co-training: multiple compatible and conditionally independent views. By theoretically deriving the generalization upper bound, we prove the prediction similarity between two models negatively impacts the model’s generalization ability. However, most current co-training models are tightly coupled together and violate this assumption. Such coupling leads to the homogenization of networks and confirmation bias which consequently limits the performance. To this end, we explore different dimensions of co-training and systematically increase the diversity from the aspects of input domains, different augmentations and model architectures to counteract homogenization. Our Diverse Co-training outperforms the state-of-the-art (SOTA) methods by a large margin across different evaluation protocols on the Pascal and Cityscapes. For example, we achieve the best mIoU of 76.2%, 77.7% and 80.2% on Pascal with only 92, 183 and 366 labeled images, surpassing the previous best results by more than 5%. Code will be available at <https://github.com/williamium3000/diverse-cotraining>.

1. Introduction

Deep learning has demonstrated impressive success in various applications [41, 32, 66, 35, 92, 93]. However, such SOTA performance heavily depends on supervised learning which requires expensive annotations. Particularly, labeling images for semantic segmentation is much more laborious and time-consuming compared with that of image classification [13]. Therefore, how to leverage the unlabeled images available in much larger quantities to improve the segmentation performance becomes crucial. Semi-supervised segmentation is thus proposed to alleviate the expensive an-

notation and is attracting growing attention [68].

One line of research in semi-supervised segmentation is co-training, which was first proposed by [5] with two compatible and independent views to guarantee the theoretical learnability [57]. It was first introduced to semi-supervised segmentation for cross pseudo supervision (CPS) [12]. Since most computer vision tasks provide only one view (RGB image) for each sample, CPS adopts two networks with identical architectures and different initializations to provide different opinions. Later researchers improve upon CPS through one or multiple extra networks [20], additional consistent constraints [38], multiple heads with a shared backbone [19, 60] and EMA teachers [83]. Despite these new variants of co-training, few studies have discussed the working mechanism behind the remarkable performance of co-training in semi-supervised segmentation. In this paper, we revisit the assumptions behind co-training: *two or multiple independent views compatible with the target function*. By deriving the generalization upper bound of Co-training, we theoretically show that the homogenization of networks accounts for the generalization error of Co-training methods. Empirically, we examine the existing co-training architectures and discover that they provide insufficient diversity (only by different initialization). We argue that similar decision boundaries or predictions will further lead to confirmation biases [1, 59] as no additional information or correction is induced. Given this problem, a natural question emerges, *how to create two or multiple views that are mutually independent?* To answer this question, we systematically explore different dimensions of the Co-training process to increase diversity. Specifically, both RGB and frequency domain are adopted as two inputs that cater to different properties of an image. Different augmentations of the same image first provide distinct views for each model. Different architectures including CNN and Transformer-based networks are then demonstrated to be quite effective in co-training due to the diverse inductive biases. Combining these findings, we propose our holistic approach: *Diverse Co-training*. To summarize, we make the following contributions:

- We theoretically prove that the homogenization of networks accounts for the generalization error of Co-

[†]Work done during internship at SenseTime.

training and discover the lack of diversity in current co-training methods that violate the assumptions.

- We comprehensively explore the different dimensions of co-training to promote diversity including the input domains, augmentations and architectures and demonstrate the significance of diversity in co-training.
- We propose a holistic framework combining the above three techniques to increase diversity and discuss two variants with high empirical performances.

2. Related Work

Semi-supervised Learning. Early works introduce self-training to tackle semi-supervised learning with an iterative EM algorithm [58, 14, 76]. Instead of labeling the unlabeled data before training, consistency regularization typically enforce invariance to perturbations on the unlabeled data in an online manner [65, 44, 85, 4, 55, 3, 67, 81]. Along this line of studies, researchers notice the significance of strong data augmentation. Combining with EMA Teacher [72], the "Teacher-Student" paradigm emerges. However, this framework suffers from the coupling problem since the teacher comes from the aggregation of student parameters [39], which fails to transfer meaningful knowledge and further leads to confirmation biases [72, 59]. We refer to [39] for detailed analysis. To this end, Deep Co-training is proposed [63, 39]. Initially, co-training is proposed to solve the semi text classification problem with two models and two views [5]. The author proves the learning ability with PAC framework [75] with the assumption that two views are compatible and independent given the class [57]. Unfortunately, most tasks in computer vision provides only single view for each sample (*i.e.* RGB image) [48, 15, 45, 13]. To simulate the condition required by co-training, [101, 39] propose to use different initialization while Tri-training [16] creates diverse training sets by injecting noise to true labels. Deep Co-training proposes a novel view difference constraint and adversarial examples as an additional pseudo view [63]. Other methods such as resampling [2], bagging [71] or bootstrapping [22] is also used to generate pseudo views. Deep Co-training is closely related to our work. Despite the two views increases diversity, adversarial examples are generated from the original sample leading to large dependence. Moreover, models trained with adversarial examples suffers from a degraded performance [74, 70] resulting to unequal roles where the original model serves as teacher and the model trained on adversarial examples the student (since original model is better with higher performance).

Semi-supervised Segmentation evolves from the early GAN-based methods [24, 68, 36, 54] which leverage the discriminator [25] to provide an auxiliary supervision for unlabeled images to simpler training paradigms of consis-

tency regularization [23, 38, 40, 34, 82, 61, 103] and entropy minimization [37, 33, 94, 91]. CPS first introduces co-training to provide cross supervision by using two identical networks with different initialization [12]. n-CPS builds upon CPS and propose to add additional models with different initialization [20] while [83] leverages the EMA of each model to teach the other model acting as teacher. Another variant [19, 60] of Co-training is the form of shared backbone and two heads, which is parameter and computation efficient. The co-training between CNN and transformer is also explored [98, 53] but focuses mainly on the powerful representation ability brought by the transformer. [98] propose to distill the feature maps between the CNN and transformer models leading to more coupled networks. Our work, on the contrary, explores the diversity presented in different architectures and demonstrate that such diversity along can achieve the SOTA performance.

Dense Vision Transformer. Recent works starting with ViT [17] prove the transformer's adaptability in CV. Later works such as Swin Transformer [52, 51] and PVT [78, 79] proves its superiority over CNN with different inductive biases trained [28, 30, 95, 73]. Recent work also demonstrates that transformer outperforms traditional CNN on dense predictions tasks such as object detection [6, 102, 9, 47] and semantic segmentation [69, 96, 84, 97]. SETR first introduce transformer to extract feature for segmentation [97]. Segmenter [69] leverages mask transformer to dynamically generate class masks while SegFormer [84] designs a novel transformer backbone with pyramid structures and simple multi-layer perceptron (MLP) head for aggregation of the multi-scale features.

3. Theoretical Analysis and Motivation

In this section, we first provide preliminary knowledge on co-training. Then a theoretical analysis on Co-training is provided to show the relation between homogenization and generalization bound. Thirdly, we conduct a thorough investigation of the existing variants of co-training and discover that they suffer significantly from the homogenization problem, which is mainly caused by a lack of two independent views. These findings further motivate our work: *how to approximate these assumptions by introducing more diversity into co-training?*

3.1. Preliminary

We describe Co-training in the context of segmentation. Segmentation network typically possesses an encoder f and a decoder head h . As aforementioned, co-training in semi-supervised segmentation generally utilizes two segmentation networks $f_i(\cdot, \theta_1^e)$, $h_i(\cdot, \theta_1^d)$ parameterized by θ_i^e, θ_i^d , $i = 1, 2$. Co-training simultaneously trains two models and the confident prediction from one model is used to supervise the other in sense of mutual teaching. Specif-

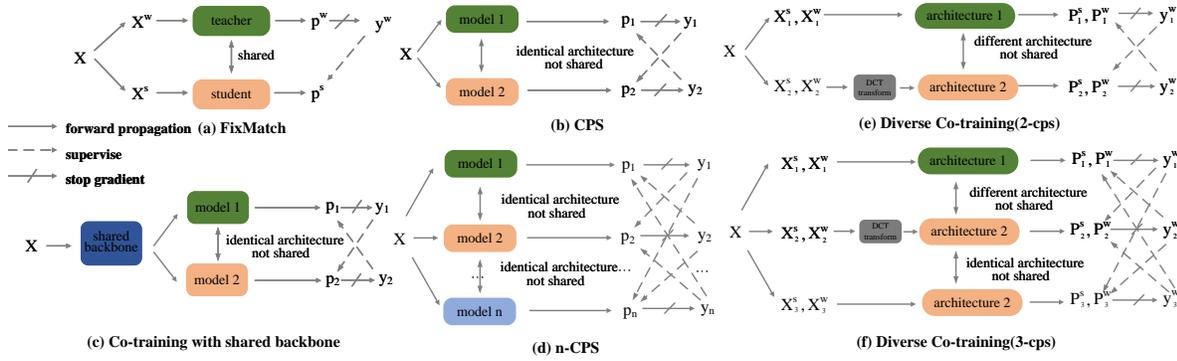


Figure 1. Illustrating the architectures for (a) FixMatch [67, 90], (b) CPS [12], (c) cross heads with shared backbone [19, 60] (d) n-CPS [20], (e) *Diverse Co-training (2-cps)* and (f) *Diverse Co-training (3-cps)*.

ically, each network generates pseudo segmentation confidence map p_i after softmax operation and then the one-hot pseudo maps y_i by argmax for each unlabeled data x : $p_i = h_i(f_i(x, \theta_i^e), \theta_i^d)$. The generated pseudo maps are used to supervise the other network on unlabeled data.

$$L^u = \frac{1}{|D_u|} \sum_{x \in D_u} \frac{1}{W \times H} \sum_{k=0}^{W \times H} CE(p_{1k}, y_{2k}) + CE(p_{2k}, y_{1k})$$

where D_u denotes the unlabeled data set, W and H denotes the size of input image and k denotes the k^{th} pixel. Cross entropy loss CE is adopted here and for the rest of the paper. For labeled data, each network is trained in a standard fully supervised manner:

$$L^s = \frac{1}{|D_l|} \sum_{x \in D_l} \frac{1}{W \times H} \sum_{k=0}^{W \times H} CE(p_{1k}, y_k^*) + CE(p_{2k}, y_k^*)$$

where D_l denotes the labeled data set and y_k^* is the ground truth label for pixel k . Then, the overall objective function is a combination of the above two losses with a balance term λ : $L = L^s + \lambda L^u$

3.2. Theoretical Analysis

We provide a generalization bound in the PAC learning framework following [5] on the Co-training method with two models. We first give the definition of homogenization. We simplify notations from above and denote the model $f_i^j, i = 1, 2$ as the i th model of j th iteration and the optimal model as f^* .

Definition 1 We define homogenization as the similarity between two networks, which can be approximated by the percentage of the agreement (agree rate) of all pixels.

$$H = Pr_{x \in D} [f_1(x) = f_2(x)] = \frac{1}{HW} \sum_{i=1}^{HW} \mathbb{1}(p_{1i} = p_{2i})$$

With different architectures, direct measures in parameter space are meaningless, we thus consider target space to

quantify homogenization. Diversity is exactly the opposite $d(f_1, f_2) = Pr_{x \in D} [f_1(x) \neq f_2(x)]$, which can be used to quantify the difference between any model. For instance, we can compute the generalization error with $d(f, f^*)$. We simplify the Co-training pipelines for easy theoretical derivation.

Assumption 1 At each step of optimization, pseudo labels of all unlabeled data are updated prior to optimization instead of online pseudo labeling.

Assumption 2 At each step of optimization, all unlabeled data is used except for the first step where only labeled data is used to get the initial model f_1^0, f_2^0 .

With the PAC learning framework and the two assumptions, we can extend the generalization bound of [80] to iterative Co-training instead of one-step optimization and obtain the following.

Theorem 1 Given hypothesis class \mathcal{H} and labeled data set D_l of size l that are sufficient to learn an initial segmentor f_i^0 with an upper bound of the generalization error of b_i^0 with probability δ (i.e. $l \geq \max\{\frac{1}{b_i^0} \ln \frac{|\mathcal{H}|}{\delta}\}$), we use empirical risk minimization to train f_i^0 on the combination of labeled and unlabeled set σ^i where pseudo label are provided by the other model f_{3-i}^0 . Then we have

$$Pr[d(f_i^k, f^*) \geq b_i^k] \leq \delta$$

if $lb_i^0 \leq e^{\sqrt{M}} - M$, where $M = ub_{3-i}^0$ and $b_i^k = \max\{\frac{lb_i^0 + ub_{3-i}^0 - ud(f_{3-i}^{k-1}, f_i^k)}{l}, 0\}$.

Theorem 1 shows that the bigger the difference between the two models f_{3-i}^{k-1} and f_i^k , the smaller the upper bound of the generalization error. Thus from Theorem 1 and Definition 1, we can conclude.

Remark 1 Homogenization negatively impacts the generalization ability of the Co-training method leading to sub-optimal performance.

With the condition that the difference between the two models is large enough $d(f_{3-i}^{k-1}, f_i^k) \geq b_{3-i}^0$, we can see that the larger the u the smaller the upper bound of the generalization error. Then we have Remark 2.

Remark 2 Given a large difference between the two models, more unlabeled data decreases the generalization error of Co-training.

This remark is consistent with empirical results that more unlabeled data leads to better performance. Further, with this remark, we provide theoretical guarantees for strong augmentations used in our method.

3.3. Homogenization in Co-training

Given the statement that homogenization negatively impacts performance, we now investigate the existing Co-training methods. Other than CPS ((b) of Figure 1), we summarize two co-training paradigms. co-training with cross heads and shared backbone is widely adopted in previous works [19, 60], as shown in (c) of Figure 1. n-CPS [20] leverages multiple models to perform co-training, which can be seen as a generalized CPS ((d) of Figure 1). We also display the paradigms used in FixMatch in (a) of Figure 1. Co-training is first introduced for its merit of providing decoupled models for cross-supervision, which can alleviate confirmation biases and generate additional information for its counterpart [39]. We demonstrate the superiority of co-training over FixMatch in Figure 2, from which co-training outperforms FixMatch consistently on all partitions and thresholds.

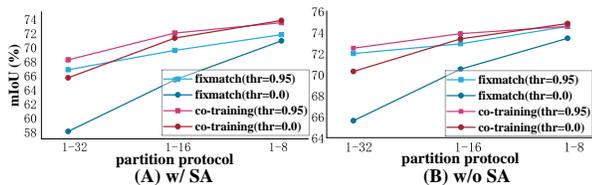


Figure 2. Comparison between FixMatch and co-training under different partition protocols and thresholds (a) with or (b) without strong augmentation (SA). Best viewed in color.

Despite the benefit, current deep co-training strongly violates the second assumption of co-training, as discussed in Section 1, since only a single view is used. The second assumption can weekly fulfilled through different initializations but it contributes too little difference for the two models to learn distinct decision boundaries. As shown in Figure 3, we can observe that all three paradigms have a severe homogenization issue. We also provide rigorous analysis in logits and prediction space with L2 distance and KL Divergence demonstrating similar phenomena in Appendix B.

Specifically, we find that co-training with a shared backbone has the most severe homogenization as the shared

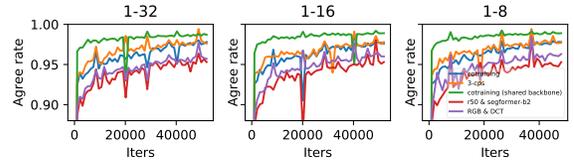


Figure 3. Homogenization problem (measured by agree rate) with different partitions. 3-cps refers to n-cps when $n = 3$.

backbone delivers the same features for each head. n-CPS also suffers from a severe homogenization problem than CPS because cross-supervision of a stack of models further enforces these models to predict similarly. We emphasize that agree rate alone is not sufficient to judge the effectiveness of a co-training method, as we cannot simply say a method is worse because they provide similar predictions. Thus, we empirically show in Table 1 that **less similar models bring performance benefits** (*i.e.* co-training outperforms the other two consistently over all settings). On the other hand, co-training with diverse input domains or different architectures achieves a much lower prediction similarity providing much more information to complement the counterpart. We show in Section 4 and evaluate in Section 5.2 that they both bring significant empirical benefits and thus can be viewed as two relatively more independent views compatible with the target function.

4. Diverse Co-training

After analyzing the limitation of current co-training paradigms, we provide a comprehensive investigation of co-training to i) promote the diversity between models and ii) provide a relatively more independent pseudo view that better fits the assumption in the PAC framework. We first introduce a better co-training baseline by adopting the strong-weak augmentation. Then, we propose and analyze three techniques to better increase the diversity between models.

Method	1/32 (331)	1/16 (662)	1/8 (1323)
Sup Baseline	61.2	67.3	70.8
Co-training	65.66	71.28	73.77
w/o SA shared backbone	58.97	65.94	71.25
3-cps	65.41	70.81	72.84
Co-training	70.28(+4.62)	73.36(+2.08)	74.82(+1.05)
w/ SA shared backbone	69.48	70.16	73.47
3-cps	69.68	71.83	74.36

Table 1. Co-training methods on ResNet50 with or without strong augmentation (SA).

Strong-Weak Augmentation. Strong-weak augmentation paradigm supervises a strongly perturbed unlabeled image x^s with the pseudo label provided by its corresponding weakly perturbed version x^w . A better pseudo label can be obtained with x^w while more efficient learning can be

conducted on x^s since x^s expands the knowledge [86, 90], alleviates confirmation biases [1] and enforce models with a decision boundary in the low-density regions [59]. Theoretically, we can also see the positive effect of strong augmentations through Remark 2 by showing that strong augmentations can potentially increase the size of unlabeled data. We argue that the improvements brought by strong augmentation are orthogonal to that of co-training with mutual benefits. In light of this statement, we combine the strong-weak augmentation with co-training to provide a better baseline for semi-supervised segmentation. Formally, we denote the weak augmentation O^w sampled from weak augmentation space S^w (i.e. random cropping and flipping) and strong augmentation O^s from S^w (detailed in Appendix G). For each image x , we obtain the strongly augmented image $x^s = O^s(O^w(x))$ and the weakly perturbed $x^w = O^w(x)$. To combine co-training with strong-weak augmentation, each model is fed with both x^w and x^s and predicts on x^w to supervise the other model as illustrated by (e) of Figure 1. This can be formulated by replacing the L^u with L_{st}^u , where the subscript st stands for "strong-weak".

$$L_{st}^u = \frac{1}{|D_u|} \sum_{x \in D_u} \frac{1}{W \times H} \sum_{i=0}^{W \times H} CE(p_{1i}^s, y_{2i}^w) + CE(p_{2i}^s, y_{1i}^w)$$

where $p_1^s = h_1(f_1(x^s, \theta_1^e), \theta_1^d)$, $p_1^w = h_1(f_1(x^w, \theta_1^e), \theta_1^d)$. p_2^s and p_2^w are similar and thus omitted. We empirically evaluate the effectiveness of strong-weak augmentation in combination with different co-training methods in Table 1. The improvement is significant and consistent on all co-training frameworks, demonstrating our argument that co-training and strong augmentation takes effect orthogonally and complements each other. We strongly suggest taking the improved co-training as the baseline for future studies.

Diverse Input Domains as Pseudo Views. Co-training methods build upon two independent views while most vision tasks provide only a single view. To relax the condition, the objective is to create pseudo views with the property of i) compatibility with the target function and ii) independent with the RGB view. To this end, we propose to learn two models from different input domains. We leverage the discrete cosine transform (DCT) coefficients to generate the frequency domain input, as illustrated in Figure 4. We refer to Appendix E for more details. The frequency domain, in its appearance, is extraordinarily different from the RGB domain. The compression and quantization process renders the DCT relatively more independent with the RGB image than, say, an image with augmentation or adversarial perturbations. The compressed representations in the frequency domain also contain rich patterns distinct from RGB domains [99, 89, 27] which not only provides additional information to the training process but also introduce additional inductive biases in data for different perspectives in the Co-training.

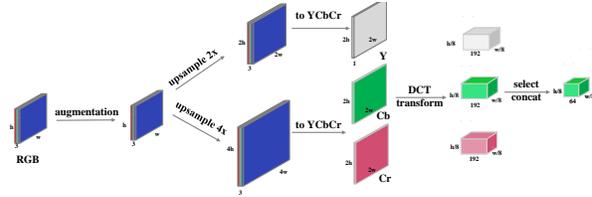


Figure 4. Illustration of the DCT transform process.

Different Augmentation Provides Different Views.

Given a single view (i.e. RGB image), augmentation is the most straightforward way to generate pseudo views [31, 11, 26, 10]. With different strong augmentation applied to the same image, distinct views can be generated so that predictions of the two models are not too similar. Particularly, random cropping has been proven effective to produce different views [8]. By randomly cropping images to different view crops, we inject view differences into the input data and Co-training is performed only on overlapping regions. Besides, random augmentations such as color jitter and greyscale also help prevent homogenization. We provide the details of augmentations used in Appendix G. Besides the diversity concern, different augmentations for each model also potentially increase the unlabeled data, which we show in Remark 2 can provide a better upper bound for generalization error.

Diverse Architecture Provides Different Inductive Biases.

Due to the lack of two views, deep co-training utilizes different initialization to relax the condition. By taking one more step, we propose to utilize models instantiated from different architectures, i.e. $f_1 \neq f_2$. In addition to different weights, diverse architectures provide different inductive biases to model the input. For the independence assumption, we provide an intuitive explanation to demonstrate why diverse architecture provides diverse views. Given the model as a function f parameterized by θ , it's essentially a composite function with each layer l as a function $z^l = f^l(\cdot, \theta^l)$ with input from the previous layer z^{l-1} . By freezing the layers below some i^{th} layer $f^l(\cdot, \theta^l)$, $l = i - 1, i - 2, \dots, 1$, we can view the i^{th} layers and above as a trainable function that maps the output of the $i - 1^{th}$ layer to the target class. Intuitively, the output of the two models from the $i - 1^{th}$ layer is different and less dependent compared with the input (i.e. same image) as different architectures or initialization are used. This intuitively explains why different initialization can relax the requirement of two independent views for co-training. It also explains why different architectures are better: due to the different inductive biases by different architectures, the output of every layer i is much more different compared with that of different initialization, thus better fulfilling the assumption. Practically, one can leverage different CNN architectures to instantiate the two models (e.g. ResNet and

ResNeXT [87]). But to promote diversity, the co-training of CNN and transformer can provide a distinct set of inductive biases that benefit each other (*i.e.* CNN with local modeling and transformer with the long-range dependence [17, 62, 64, 77, 46]).

Holistic Approach: Diverse Co-training. Following the spirit of the above sections, we combine the three proposed techniques to promote a holistic framework for diverse co-training. We provide two variants of *Diverse Co-training*, termed by 2-cps and 3-cps following [20], which leverage two models and three models to co-training respectively. Concretely, we leverage CNN and transformer as the two different architectures to maximize the discrepancy with one model trained on RGB and the other on DCT domain, as illustrated in (e) and (f) of Figure 1. The semi-supervised nature of co-training brings noise into pseudo labels for unlabeled data [67, 59, 90], thus we also provide confidence thresholding following FixMatch to filter out noisy pseudo labels in which the model has low confidence. Intuitively, each model should supervise the other model with the pseudo labels it’s most confident in. We introduce n ($n = 2, 3$) additional scalar hyperparameters τ_i ($i = 1, 2, \dots, n$) denoting the threshold above which we retain a pseudo-label. Formally, we reformulate the unlabeled term as followed. We omit the sum over $H \times W$ for simplicity.

$$L_{st}^u = \frac{1}{|D_u|} \sum_{x \in D_u} \mathbb{1}(\max(p_{2i}^w) > \tau_2) l_{ce}(p_{1i}^s, y_{2i}^w) + \mathbb{1}(\max(p_{1i}^w) > \tau_1) l_{ce}(p_{2i}^s, y_{1i}^w)$$

Unlike other methods [19, 94, 98] which either inserts modules or utilize the output of intermediate layers, we emphasize that our holistic approach leverages off-the-shelf segmentation networks without changing or probing its inside components, which can be quickly incorporated with any new SOTA segmentation methods and even be easily adapted to other fields such as semi-supervised classification and object detection.

	Input Domain	1/32 (331)	1/16 (662)
w/o SA	RGB	65.66	71.28
	DCT	65.33	67.37
	RGB & DCT	69.45 / 69.03	72.46 / 72.03
	RGB & HSV	69.65 / 67.05	71.74 / 69.89
w/ SA	RGB	70.28	73.36
	DCT	70.65	73.26
	RGB & DCT	71.88 / 72.00	74.10 / 73.94
	RGB & HSV	70.40 / 68.30	72.64 / 70.91

Table 2. Performance of co-training with different domains. For cells with two numbers, the left one is the result of RGB model.

5. Experiment

We conduct experiments in this section. The objective is to i) demonstrate our argument that diversity plays a crucial role in co-training ii) by illustrating that the three proposed techniques can effectively improve the performance and prevent the model from being tightly coupled with each other and iii) demonstrate the effectiveness by comparing with other state-of-the-arts.

5.1. Experiment Setup

Datasets. We leverage two datasets for evaluating the effectiveness of our idea. *PASCAL VOC 2012* [29] is constructed by a combination of the Pascal dataset [18] with high-quality train and validation images and the coarsely annotated SBD dataset [29], resulting in a total of 10582 training images. Following prior arts, we randomly sample labeled images from i) the original high-quality training images, and ii) the mixed 10582 images. *Cityscapes* [13] is an urban scene dataset with 19 classes and a total of 2975 high-resolution (2048×1024) training images as well as 500 validation images. We follow prior arts and divide the dataset by randomly sub-sampling 1/4, 1/8 and 1/30 of the total training images as labeled set and the rest as the unlabeled set. We crop each image to 769x769 during training.

Evaluation. We evaluate the segmentation performance with the mean Intersection-over-Union (mIoU) metric. For all partition protocols, we report the results on the *PASCAL VOC val* set with single scale testing on origin resolution and *Cityscapes val* set with single scale sliding window evaluation with crop size of 769 following [90, 82]

Implementation Details. For fair comparisons, we leverage the widely adopted DeepLabv3+ with ResNet as CNN and the SegFormer as the transformer architecture. The backbones of both architectures are pre-trained on ImageNet 1K. We utilize the pre-trained weights on ImageNet 1K for frequency domain from [89]. During training, we leverage a batch size of 16 for Pascal and 8 for Cityscapes with a labeled-unlabeled ratio of 1. We train Pascal and Cityscapes for 80 and 240 epochs with an initial learning rate of 0.001 and 0.005 respectively and polynomial learning rate decay following [12].

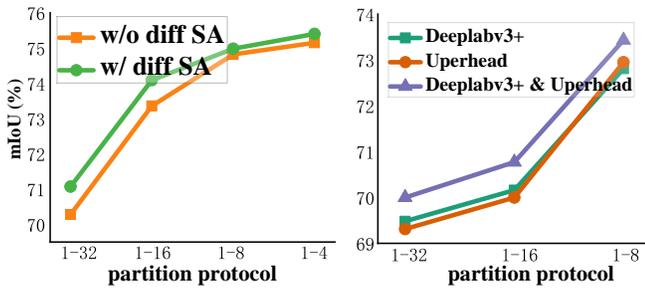
5.2. Analysis on How to Promote Diversity

We empirically analyze the three techniques proposed with ResNet50. All experiments and figures in Section 3.3 and this sections are conducted on *PASCAL VOC* with partition protocol 1/32 and 1/16.

Different Input Domains provides two views that are relatively independent of each other. We demonstrate that the DCT domain better fulfills the two assumptions of co-training. Firstly, we show in Table 2 that training only on the DCT domain achieves similar accuracy as the RGB domain, demonstrating that the DCT view is compatible with

	Backbone	1/32 (331)	1/16 (662)
w/o SA	R50	65.66	71.28
	mit-b2	71.01	74.53
	R50 & mit-b2	71.58 / 71.03	74.94 / 74.84
w/ SA	R50	70.28	73.36
	mit-b2	74.51	75.29
	R50 & mit-b2	<u>74.85 / 74.87</u>	<u>75.12 / 75.85</u>
	ResNeSt50	70.92	<u>75.58</u>
	ResNeXt50	71.18	72.77
	R50 & ResNeST50	<u>72.70 / 73.56</u>	<u>73.41 / 75.65</u>
	R50 & ResNeXT50	<u>72.15 / 72.39</u>	<u>74.41 / 74.56</u>

Table 3. Performance of co-training with different architectures. We refer ResNet50 as R50 and SegFormer-b2 as mit-b2 [84]. For cells with two numbers, the left one is the result of ResNet50.



(a) Different strong augmentation (SA) (b) Different head architecture

Figure 5. (a) Empirical evaluation of co-training with or without different strong augmentation (SA); (b) Empirical evaluation of co-training with different head structures.

the target function of the RGB domain and sufficient to train a quality segmentor. As per Figure 3, co-training with different domains significantly lowers the similarity between models, illustrating that different domains increase the diversity of the two models. We further evaluate empirically on *VOC PASCAL* where the RGB & DCT outperforms the baselines with single view consistently on all settings. Notice that both models of RGB & DCT obtain a significant improvement, demonstrating that diversity benefits two models mutually instead of a unidirectional teacher-student one. We also conduct experiments on HSV, a color space different from RGB and also observe similar but fewer improvements over baseline compared with DCT. We hypothesize this is due to the discrepancy between DCT and RGB being larger than that of HSV and RGB. **Different Augmentation** further promotes diversity in co-training. We empirically show in (a) of Figure 5 that different augmentations for each model are effective and show consistent improvement over the baseline.

Different Architectures provides different inductive biases and thus better independent views as discussed in Section 4. From Figure 3, we already demonstrate that less similar decision boundaries can be obtained with different architectures. We here empirically evaluate the performance

Method	Resolution	92	183	366	732	1464
ResNet50						
Sup Baseline	513x513	39.1	51.3	60.3	65.9	71.0
PseudoSeg [103]	512x512	54.9	61.9	64.9	70.4	-
PC ² Seg [100]	512x512	56.9	64.6	67.6	70.9	-
Ours (2-cps)	513x513	<u>71.8</u>	<u>74.5</u>	77.6	<u>78.6</u>	<u>79.8</u>
Ours (3-cps)	513x513	73.1	74.7	<u>77.1</u>	78.8	80.2
ResNet101						
Sup Baseline	321x321	44.4	54.0	63.4	67.2	71.8
ReCo [49]	321x321	64.8	72.0	73.1	74.7	-
ST++ [91]	321x321	65.2	71.0	74.6	77.3	79.1
ours (2-cps)	321x321	<u>74.8</u>	77.6	<u>79.5</u>	<u>80.3</u>	81.7
ours (3-cps)	321x321	75.4	<u>76.8</u>	79.6	80.4	<u>81.6</u>
Sup Baseline	512x512	42.3	56.6	64.2	68.1	72.0
MT [72]	512x512	48.7	55.8	63.0	69.16	-
CPS[12]	512x512	64.1	67.4	71.7	75.9	-
U ² PL [82]	512x512	68.0	69.2	73.7	76.2	79.5
PS-MT [50]	512x512	65.8	69.6	76.6	78.4	80.0
ours (2-cps)	513x513	76.2	76.6	80.2	80.8	81.9
ours (3-cps)	513x513	<u>75.7</u>	77.7	<u>80.1</u>	80.9	82.0

Table 4. Comparison with state-of-the-art methods on the *Pascal* dataset. Labeled images are from the high-quality training set.

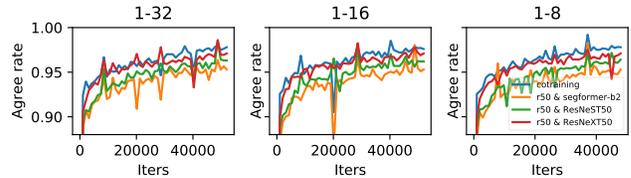


Figure 6. Homogenization (measured by agree rate) with different architectures.

of co-training with different architectures, as presented in Table 3. Both R50 and mit-b2 can be further improved through co-training compared with their corresponding individual baseline. For instance, we observe a remarkable improvement of ResNet50 over single R50 co-training (*e.g.* 5.37% and 3.56% on 1/32 and 1/16 w/o SA), which can be attributed to the high performance of mit-b2. However, ResNet50 of R50 & mit-b2 even surpasses the individual mit-b2 demonstrating that the cross-supervision between different architectures provides additional information other than the pseudo labels from mit-b2. We also provide co-training with different CNN architectures and illustrate in Figure 6 that different CNNs are more coupled compared with CNN and transformer. Empirically, despite baseline ResNeSt50 obtaining better performance than baseline mit-b2, R50 & mit-b2 outperforms R50 & ResNeSt50 in all settings. We also notice the improvement of ResNet50 in co-training with different CNN architectures is less than that of R50 & mit-b2, which further illustrates our point. To further prove the concept, we additionally conduct experiments on co-training with shared backbone and leverage different

decoder head structures, as shown in (b) of Figure 5. Co-training with DeepLabv3+ head [7] and UPerHead [88] is better than any baselines with single-head architecture consistently.

Method	Resolution	1/32	1/16	1/8	1/4
		(331)	(662)	(1323)	(2646)
Sup Baseline	321x321	55.8	60.3	66.8	71.3
CAC[43]	320x320	-	70.1	72.4	74.0
ST++[91]	321x321	-	72.6	74.4	75.4
Ours (2-cps)	321x321	75.2	76.0	76.2	76.5
Ours (3-cps)	321x321	74.9	76.4	76.3	76.6
Sup Baseline	513x513	54.1	60.7	67.7	71.9
CPS[12]	512x512	-	72.0	73.7	74.9
3-CPS [20]	512x512	-	72.0	74.2	75.9
ELN [42]	512x512	-	-	73.2	74.6
PS-MT [50]	512x512	-	72.8	75.7	76.4
U ² PL* [82]	513x513	-	72.0	75.1	76.2
Ours (2-cps)	513x513	75.2	76.2	77.0	77.5
Ours (3-cps)	513x513	74.7	76.3	77.2	77.7

Table 5. Comparison with state-of-the-art methods with ResNet50 on the *Pascal VOC 2012* dataset. Labeled images are sampled from the blended training set. The result of U²PL is reproduced with the same setting as ours.

5.3. Comparison with State-of-the-arts

Pascal VOC 2012. We only compare the most recent SOTA models due to limited space and a full comparison can be found in Appendix J. We first compare the performance of *Diverse Co-training* with SOTA on *PASCAL VOC* on two groups of data partition protocols described in 5.1. To ensure a fair comparison, we conduct training with resolutions of 321 and 513, which is reported with the results. On the first partition protocol, our *Diverse Co-training* outperforms the prior methods by a remarkable margin, as displayed in Table 4. For instance, we obtain an improvement of over 10% on ResNet50 compared with the PC²Seg on 92, 183 and 366 partitions. With ResNet101, our method surpasses the up-to-date SOTA (*i.e.* PS-MT) by a margin as large as 9% under scarce label conditions such as 92 and 183, and outperforms all prior arts consistently on other settings. On the second protocol, as indicated by Table 5, our method also gains remarkable improvements over most up-to-date studies. We emphasize that our performance on 1/32 (*i.e.* 75.2%) already outperforms other SOTA with 1/16 (around 72%) by 3%, which shows the effectiveness of our approach. We further report the comparison on ResNet101 and SegFormer-b3 in Appendix D. **Cityscapes.** We compare the SOTA with our method in Table 6. We report results on both ResNet50 & SegFormer-b2 and ResNet101 & SegFormer-b3. We can see that our method outperforms the current SOTA (*i.e.* U²PL) by more than 3% on 1/30 and 1/8 and 1.6% on 1/4 with ResNet50 & SegFormer-b2. Similar

trends can also be observed on ResNet101, where we contribute an improvement of 0.5% and 1.1% over PS-MT on 1/8 and 1/4 protocol.

Method	ResNet50			Method	ResNet101		
	1/30	1/8	1/4		1/16	1/8	1/4
	(100)	(372)	(744)		(186)	(372)	(744)
Sup Baseline	54.8	70.2	73.6	Sup Baseline	66.8	72.5	76.4
CAC [43]	-	69.7	72.7	CutMix [23]	67.9	73.5	75.4
CPS [12]	-	74.4	76.9	CPS [12]	70.5	75.7	77.4
ST++ [91]	61.4	72.7	73.8	U ² PL [82]	74.9	76.5	78.5
U ² PL* [82]	59.8	73.0	76.3	PS-MT [50]	-	76.9	77.6
Ours (2-cps)	64.5	76.3	77.1	Ours (2-cps)	75.0	77.3	78.7
Ours (3-cps)	65.5	76.5	77.9	Ours (3-cps)	75.7	77.4	78.5

Table 6. Comparison with state-of-the-art methods on the Cityscapes dataset. Result of U²PL on ResNet50 (marked with *) is from [90] which is reproduced with the same setting as ours.

5.4. Ablation Study and Analysis

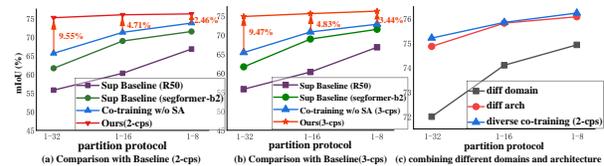


Figure 7. Comparison with baseline model with (a) *Diverse Co-training (2-cps)* (b) *Diverse Co-training (3-cps)*; (c) Ablation study of combining diverse domains and different architectures.

Comparison with Baseline. We compare our methods with baseline in Figure with ResNet50 and SegFormer-b2 on *PASCAL VOC*. We compare both variants of our method, *i.e.* 2-cps and 3-cps, to each corresponding baseline (*i.e.* CPS and n-CPS) and the labeled-only baseline. As per Figure 7, we achieve an improvement 9.55%, 4.71%, 2.46% and 9.47%, 4.83%, 3.44% over their corresponding baseline respectively. The improvement over the supervised baseline is larger. It’s worth mentioning that the n-CPS baseline with single architecture suffers from homogenization which consequently limits the performance compared with the CPS baseline (shown in Table 1). With the diversity boosted, *Diverse Co-training (3-cps)* now outperforms both baselines by a large margin, further demonstrating that diversity is crucial in co-training. Co-training is similar to knowledge distillation (KD) in the sense that they both possess a teaching process, the difference lies in that the teacher in KD is usually fixed and teaching is unidirectional. We provide a KD baseline comparison in Appendix F. **Combining Diverse domains and Different Architectures** renders model less coupled further increasing the diversity and encouraging the exploration at early stage of training. Empirically, we also demonstrate a improvement of the combination over the each individual, as shown in (c) of Figure 7. We

also provide a more complete ablation study on each one of the component and their combinations in Appendix K. **Confidence Threshold τ and Unlabeled Weight λ .** As proposed in Section 4, our holistic *Diverse Co-training* leverages confidence threshold τ to retrain confident samples as pseudo label and balance the losses on unlabeled data with weight λ . To be as simple as possible, we set the threshold $\tau = 0.0$ and $\lambda = 1.0$ as default for the experiments reported in this paper. But we also report the performance of different values in Figure 8. We emphasize that a better performance can be obtained if optimal hyperparameters is thoroughly searched on each setting as in [50], which is omitted due to computation reason. **Number of parameters.** We report the parameters of all backbones used in Section 5.2 in Appendix H. We also provide a rigorous analysis to show that our improvement is not trivial by adding more parameters.

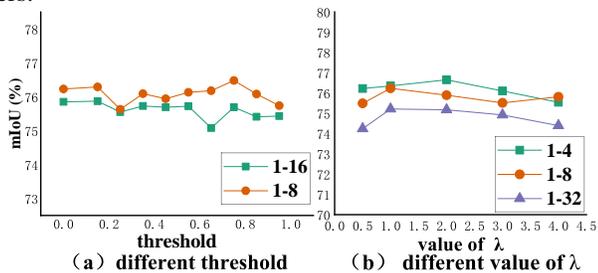


Figure 8. Performance on (a) different confidence thresholds τ and (b) different weight λ .

6. Conclusion

In this paper, we revisit the two core assumptions behind the deep co-training methods in semi-supervised segmentation and provide a theoretical upper bound over the generalization error that links with the homogenization of the two networks. We discover that the existing co-training paradigms suffer from severe homogenization problems and by exploring different dimensions of co-training and systematically increasing the diversity from three aspects, we propose a holistic framework: *Diverse Co-training* which achieves remarkable improvement over previous best results on all partitions of the *Pascal* and *Cityscapes* benchmarks.

References

- [1] Eric Arazo, Diego Ortego, Paul Albert, Noel E O’Connor, and Kevin McGuinness. Pseudo-labeling and confirmation bias in deep semi-supervised learning. In *2020 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2020. 1, 5
- [2] Eric Bauer and Ron Kohavi. An empirical comparison of voting classification algorithms: Bagging, boosting, and variants. *Machine learning*, 36(1):105–139, 1999. 2
- [3] David Berthelot, Nicholas Carlini, Ekin D Cubuk, Alex Kurakin, Kihyuk Sohn, Han Zhang, and Colin Raffel. Remixmatch: Semi-supervised learning with distribution alignment and augmentation anchoring. *arXiv preprint arXiv:1911.09785*, 2019. 2
- [4] David Berthelot, Nicholas Carlini, Ian Goodfellow, Nicolas Papernot, Avital Oliver, and Colin A Raffel. Mixmatch: A holistic approach to semi-supervised learning. *Advances in neural information processing systems*, 32, 2019. 2
- [5] Avrim Blum and Tom Mitchell. Combining labeled and unlabeled data with co-training. In *Proceedings of the eleventh annual conference on Computational learning theory*, pages 92–100, 1998. 1, 2, 3
- [6] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pages 213–229. Springer, 2020. 2
- [7] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *ECCV*, 2018. 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. 5
- [9] Ting Chen, Saurabh Saxena, Lala Li, David J Fleet, and Geoffrey Hinton. Pix2seq: A language modeling framework for object detection. *arXiv preprint arXiv:2109.10852*, 2021. 2
- [10] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020. 5
- [11] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15750–15758, 2021. 5
- [12] Xiaokang Chen, Yuhui Yuan, Gang Zeng, and Jingdong Wang. Semi-supervised semantic segmentation with cross pseudo supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2613–2622, 2021. 1, 2, 3, 6, 7, 8, 15, 16, 17, 18
- [13] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223, 2016. 1, 2, 6
- [14] Arthur P Dempster, Nan M Laird, and Donald B Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, 39(1):1–22, 1977. 2
- [15] 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, 2009. 2
- [16] WeiWang Dong-DongChen and Zhi-HuaZhou WeiGao. Tri-net for semi-supervised deep learning. In *Proceedings of twenty-seventh international joint conference on artificial intelligence*, pages 2014–2020, 2018. 2

- [17] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. [2](#), [6](#)
- [18] Mark Everingham, SM Eslami, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes challenge: A retrospective. *International journal of computer vision*, 111(1):98–136, 2015. [6](#)
- [19] Jiashuo Fan, Bin Gao, Huan Jin, and Lihui Jiang. Ucc: Uncertainty guided cross-head co-training for semi-supervised semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9947–9956, 2022. [1](#), [2](#), [3](#), [4](#), [6](#)
- [20] 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](#), [2](#), [3](#), [4](#), [6](#), [8](#), [15](#), [18](#)
- [21] Dominik Filipiak, Piotr Tempczyk, and Marek Cygan. n -cps: Generalising cross pseudo supervision to n networks for semi-supervised semantic segmentation. *arXiv preprint arXiv: Arxiv-2112.07528*, 2021. [17](#)
- [22] David A Freedman. Bootstrapping regression models. *The Annals of Statistics*, 9(6):1218–1228, 1981. [2](#)
- [23] Geoff French, Timo Aila, Samuli Laine, Michal Mackiewicz, and Graham Finlayson. Semi-supervised semantic segmentation needs strong, high-dimensional perturbations. 2019. [2](#), [8](#)
- [24] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020. [2](#)
- [25] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020. [2](#)
- [26] Jean-Bastien Grill, Florian Strub, Florent Althé, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent—a new approach to self-supervised learning. *Advances in neural information processing systems*, 33:21271–21284, 2020. [5](#)
- [27] Lionel Gueguen, Alex Sergeev, Ben Kadlec, Rosanne Liu, and Jason Yosinski. Faster neural networks straight from jpeg. *Advances in Neural Information Processing Systems*, 31, 2018. [5](#)
- [28] Kai Han, An Xiao, Enhua Wu, Jianyuan Guo, Chunjing Xu, and Yunhe Wang. Transformer in transformer. *Advances in Neural Information Processing Systems*, 34:15908–15919, 2021. [2](#)
- [29] Bharath Hariharan, Pablo Arbeláez, Lubomir Bourdev, Subhransu Maji, and Jitendra Malik. Semantic contours from inverse detectors. In *2011 international conference on computer vision*, pages 991–998. IEEE, 2011. [6](#)
- [30] 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. [2](#)
- [31] 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. [5](#)
- [32] 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. [1](#), [16](#)
- [33] Ruifei He, Jihan Yang, and Xiaojuan Qi. Re-distributing biased pseudo labels for semi-supervised semantic segmentation: A baseline investigation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6930–6940, 2021. [2](#)
- [34] Hanzhe Hu, Fangyun Wei, Han Hu, Qiwei Ye, Jinshi Cui, and Liwei Wang. Semi-supervised semantic segmentation via adaptive equalization learning. *Advances in Neural Information Processing Systems*, 34:22106–22118, 2021. [2](#), [15](#)
- [35] Siqian Huang, Yijiang Li, Chong Chen, Leyu Shi, and Ying Gao. Multi-metrics adaptively identifies backdoors in federated learning. *arXiv preprint arXiv:2303.06601*, 2023. [1](#)
- [36] Wei-Chih Hung, Yi-Hsuan Tsai, Yan-Ting Liou, Yen-Yu Lin, and Ming-Hsuan Yang. Adversarial learning for semi-supervised semantic segmentation. *arXiv preprint arXiv:1802.07934*, 2018. [2](#)
- [37] Rihuan Ke, Angelica I Aviles-Rivero, Saurabh Pandey, Saikumar Reddy, and Carola-Bibiane Schönlieb. A three-stage self-training framework for semi-supervised semantic segmentation. *IEEE Transactions on Image Processing*, 31:1805–1815, 2022. [2](#)
- [38] Zhanhan Ke, Di Qiu, Kaican Li, Qiong Yan, and Rynson WH Lau. Guided collaborative training for pixel-wise semi-supervised learning. In *European conference on computer vision*, pages 429–445. Springer, 2020. [1](#), [2](#), [15](#), [18](#)
- [39] Zhanhan Ke, Daoye Wang, Qiong Yan, Jimmy Ren, and Rynson WH Lau. Dual student: Breaking the limits of the teacher in semi-supervised learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6728–6736, 2019. [2](#), [4](#)
- [40] Jongmok Kim, Jooyoung Jang, and Hyunwoo Park. Structured consistency loss for semi-supervised semantic segmentation. *arXiv preprint arXiv:2001.04647*, 2020. [2](#)
- [41] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6):84–90, 2017. [1](#)
- [42] Donghyeon Kwon and Suha Kwak. Semi-supervised semantic segmentation with error localization network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9957–9967, 2022. [8](#), [15](#), [18](#)

- [43] Xin Lai, Zhuotao Tian, Li Jiang, Shu Liu, Hengshuang Zhao, Liwei Wang, and Jiaya Jia. Semi-supervised semantic segmentation with directional context-aware consistency. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1205–1214, 2021. 8, 15, 18
- [44] Samuli Laine and Timo Aila. Temporal ensembling for semi-supervised learning. *arXiv preprint arXiv:1610.02242*, 2016. 2
- [45] Kunpeng Li, Ziyang Wu, Kuan-Chuan Peng, Jan Ernst, and Yun Fu. Tell me where to look: Guided attention inference network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9215–9223, 2018. 2
- [46] Yijiang Li, Wentian Cai, Ying Gao, Chengming Li, and Xiping Hu. More than encoder: Introducing transformer decoder to upsampler. In *2022 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, pages 1597–1602. IEEE, 2022. 6
- [47] Yanghao Li, Hanzi Mao, Ross Girshick, and Kaiming He. Exploring plain vision transformer backbones for object detection. *arXiv preprint arXiv:2203.16527*, 2022. 2
- [48] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 2
- [49] Shikun Liu, Shuaifeng Zhi, Edward Johns, and Andrew J Davison. Bootstrapping semantic segmentation with regional contrast. *arXiv preprint arXiv:2104.04465*, 2021. 7, 18
- [50] Yuyuan Liu, Yu Tian, Yuanhong Chen, Fengbei Liu, Vasileios Belagiannis, and Gustavo Carneiro. Perturbed and strict mean teachers for semi-supervised semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4258–4267, 2022. 7, 8, 9, 15, 17, 18
- [51] Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda Xie, Yixuan Wei, Jia Ning, Yue Cao, Zheng Zhang, Li Dong, et al. Swin transformer v2: Scaling up capacity and resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12009–12019, 2022. 2
- [52] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021. 2
- [53] Xiangde Luo, Minhao Hu, Tao Song, Guotai Wang, and Shaoting Zhang. Semi-supervised medical image segmentation via cross teaching between cnn and transformer. *arXiv preprint arXiv:2112.04894*, 2021. 2
- [54] Sudhanshu Mittal, Maxim Tatarchenko, and Thomas Brox. Semi-supervised semantic segmentation with high-and low-level consistency. *IEEE transactions on pattern analysis and machine intelligence*, 43(4):1369–1379, 2019. 2
- [55] Takeru Miyato, Shin-ichi Maeda, Masanori Koyama, and Shin Ishii. Virtual adversarial training: a regularization method for supervised and semi-supervised learning. *IEEE transactions on pattern analysis and machine intelligence*, 41(8):1979–1993, 2018. 2
- [56] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of machine learning*. MIT press, 2018. 14
- [57] Kamal Nigam and Rayid Ghani. Analyzing the effectiveness and applicability of co-training. In *Proceedings of the ninth international conference on Information and knowledge management*, pages 86–93, 2000. 1, 2
- [58] Kamal Nigam, Andrew Kachites McCallum, Sebastian Thrun, and Tom Mitchell. Text classification from labeled and unlabeled documents using em. *Machine learning*, 39(2):103–134, 2000. 2
- [59] Yassine Ouali, Céline Hudelot, and Myriam Tami. An overview of deep semi-supervised learning. *arXiv preprint arXiv:2006.05278*, 2020. 1, 2, 5, 6
- [60] Yassine Ouali, Céline Hudelot, and Myriam Tami. Semi-supervised semantic segmentation with cross-consistency training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12674–12684, 2020. 1, 2, 3, 4
- [61] Yassine Ouali, Céline Hudelot, and Myriam Tami. Semi-supervised semantic segmentation with cross-consistency training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12674–12684, 2020. 2, 15, 18
- [62] Namuk Park and Songkuk Kim. How do vision transformers work? *arXiv preprint arXiv:2202.06709*, 2022. 6
- [63] Siyuan Qiao, Wei Shen, Zhishuai Zhang, Bo Wang, and Alan Yuille. Deep co-training for semi-supervised image recognition. In *Proceedings of the european conference on computer vision (eccv)*, pages 135–152, 2018. 2
- [64] Maithra Raghu, Thomas Unterthiner, Simon Kornblith, Chiyuan Zhang, and Alexey Dosovitskiy. Do vision transformers see like convolutional neural networks? *Advances in Neural Information Processing Systems*, 34:12116–12128, 2021. 6
- [65] Antti Rasmus, Mathias Berglund, Mikko Honkala, Harri Valpola, and Tapani Raiko. Semi-supervised learning with ladder networks. *Advances in neural information processing systems*, 28, 2015. 2
- [66] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28, 2015. 1
- [67] Kihyuk Sohn, David Berthelot, Nicholas Carlini, Zizhao Zhang, Han Zhang, Colin A Raffel, Ekin Dogus Cubuk, Alexey Kurakin, and Chun-Liang Li. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. *Advances in neural information processing systems*, 33:596–608, 2020. 2, 3, 6
- [68] Nasim Souly, Concetto Spampinato, and Mubarak Shah. Semi supervised semantic segmentation using generative adversarial network. In *Proceedings of the IEEE international conference on computer vision*, pages 5688–5696, 2017. 1, 2

- [69] Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. Segmenter: Transformer for semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7262–7272, 2021. 2
- [70] Dong Su, Huan Zhang, Hongge Chen, Jinfeng Yi, Pin-Yu Chen, and Yupeng Gao. Is robustness the cost of accuracy?—a comprehensive study on the robustness of 18 deep image classification models. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 631–648, 2018. 2
- [71] Clifton D Sutton. Classification and regression trees, bagging, and boosting. *Handbook of statistics*, 24:303–329, 2005. 2
- [72] Antti Tarvainen and Harri Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. *Advances in neural information processing systems*, 30, 2017. 2, 7, 15, 18
- [73] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International Conference on Machine Learning*, pages 10347–10357. PMLR, 2021. 2
- [74] Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry. Robustness may be at odds with accuracy. *arXiv preprint arXiv:1805.12152*, 2018. 2
- [75] Leslie G Valiant. A view of computational learning theory. In *Foundations of Knowledge Acquisition*, pages 263–289. Springer, 1993. 2
- [76] Jean-Noël Vittaut, Massih-Reza Amini, and Patrick Gallinari. Learning classification with both labeled and unlabeled data. In *European Conference on Machine Learning*, pages 468–479. Springer, 2002. 2
- [77] Peihao Wang, Wenqing Zheng, Tianlong Chen, and Zhangyang Wang. Anti-oversmoothing in deep vision transformers via the fourier domain analysis: From theory to practice. *arXiv preprint arXiv:2203.05962*, 2022. 6
- [78] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 568–578, 2021. 2
- [79] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pvt v2: Improved baselines with pyramid vision transformer. *Computational Visual Media*, 8(3):415–424, 2022. 2
- [80] Wei Wang and Zhi-Hua Zhou. Analyzing co-training style algorithms. In *Machine Learning: ECML 2007: 18th European Conference on Machine Learning, Warsaw, Poland, September 17-21, 2007. Proceedings 18*, pages 454–465. Springer, 2007. 3
- [81] Xinjiang Wang, Xingyi Yang, Shilong Zhang, Yijiang Li, Litong Feng, Shijie Fang, Chengqi Lyu, Kai Chen, and Wayne Zhang. Consistent-teacher: Towards reducing inconsistent pseudo-targets in semi-supervised object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3240–3249, 2023. 2
- [82] Yuchao Wang, Haochen Wang, Yujun Shen, Jingjing Fei, Wei Li, Guoqiang Jin, Liwei Wu, Rui Zhao, and Xinyi Le. Semi-supervised semantic segmentation using unreliable pseudo-labels. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4248–4257, 2022. 2, 6, 7, 8, 15, 17, 18
- [83] Hui Xiao, Dong Li, Hao Xu, Shuibo Fu, Diquan Yan, Kangkang Song, and Chengbin Peng. Semi-supervised semantic segmentation with cross teacher training. *Neuro-computing*, 508:36–46, 2022. 1, 2, 18
- [84] 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. 2, 7
- [85] Qizhe Xie, Zihang Dai, Eduard Hovy, Thang Luong, and Quoc Le. Unsupervised data augmentation for consistency training. *Advances in Neural Information Processing Systems*, 33:6256–6268, 2020. 2
- [86] Qizhe Xie, Minh-Thang Luong, Eduard Hovy, and Quoc V Le. Self-training with noisy student improves imagenet classification. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10687–10698, 2020. 5
- [87] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1492–1500, 2017. 6
- [88] Yuwen Xiong, Renjie Liao, Hengshuang Zhao, Rui Hu, Min Bai, Ersin Yumer, and Raquel Urtasun. Upsnet: A unified panoptic segmentation network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8818–8826, 2019. 8
- [89] Kai Xu, Minghai Qin, Fei Sun, Yuhao Wang, Yen-Kuang Chen, and Fengbo Ren. Learning in the frequency domain. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1740–1749, 2020. 5, 6, 16
- [90] Lihe Yang, Lei Qi, Litong Feng, Wayne Zhang, and Yinghuan Shi. Revisiting weak-to-strong consistency in semi-supervised semantic segmentation. *arXiv preprint arXiv:2208.09910*, 2022. 3, 5, 6, 8, 16
- [91] Lihe Yang, Wei Zhuo, Lei Qi, Yinghuan Shi, and Yang Gao. St++: Make self-training work better for semi-supervised semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4268–4277, 2022. 2, 7, 8, 15, 18
- [92] Xingyi Yang, Daquan Zhou, Jiashi Feng, and Xinchao Wang. Diffusion probabilistic model made slim. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22552–22562, 2023. 1
- [93] Xingyi Yang, Daquan Zhou, Songhua Liu, Jingwen Ye, and Xinchao Wang. Deep model reassembly. *Advances in neural information processing systems*, 35:25739–25753, 2022. 1

- [94] Jianlong Yuan, Yifan Liu, Chunhua Shen, Zhibin Wang, and Hao Li. A simple baseline for semi-supervised semantic segmentation with strong data augmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8229–8238, 2021. [2](#), [6](#)
- [95] Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Zi-Hang Jiang, Francis EH Tay, Jiashi Feng, and Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch on imagenet. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 558–567, 2021. [2](#)
- [96] Yuhui Yuan, Xiaokang Chen, Xilin Chen, and Jingdong Wang. Segmentation transformer: Object-contextual representations for semantic segmentation. *arXiv preprint arXiv:1909.11065*, 2019. [2](#)
- [97] Sixiao Zheng, Jiachen Lu, Hengshuang Zhao, Xiatian Zhu, Zekun Luo, Yabiao Wang, Yanwei Fu, Jianfeng Feng, Tao Xiang, Philip HS Torr, et al. Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6881–6890, 2021. [2](#)
- [98] Xu Zheng, Yunhao Luo, Hao Wang, Chong Fu, and Lin Wang. Transformer-cnn cohort: Semi-supervised semantic segmentation by the best of both students. *arXiv preprint arXiv:2209.02178*, 2022. [2](#), [6](#)
- [99] Yijie Zhong, Bo Li, Lv Tang, Senyun Kuang, Shuang Wu, and Shouhong Ding. Detecting camouflaged object in frequency domain. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4504–4513, 2022. [5](#)
- [100] Yuanyi Zhong, Bodi Yuan, Hong Wu, Zhiqiang Yuan, Jian Peng, and Yu-Xiong Wang. Pixel contrastive-consistent semi-supervised semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7273–7282, 2021. [7](#), [18](#)
- [101] Yan Zhou and Sally Goldman. Democratic co-learning. In *16th IEEE International Conference on Tools with Artificial Intelligence*, pages 594–602. IEEE, 2004. [2](#)
- [102] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020. [2](#)
- [103] Yuliang Zou, Zizhao Zhang, Han Zhang, Chun-Liang Li, Xiao Bian, Jia-Bin Huang, and Tomas Pfister. Pseudoseg: Designing pseudo labels for semantic segmentation. *arXiv preprint arXiv:2010.09713*, 2020. [2](#), [7](#), [18](#)

A. Proof of Theorem 1.

First, we can show with PAC learning [56] that with labeled data set D_l of size l where $l \geq \max\{\frac{1}{b_i^0} \ln \frac{|\mathcal{H}|}{\delta}\}$, the generalization error of the initial segmentor f_i^0 is bounded by b_i^0 with probability δ , which is a standard PAC supervised learning problem. Then, without loss of generality, we show the probability that the generalization error of f_2^k denoted by $d(f_2^k, f^*)$ is larger than b_2^k is at most δ .

we analyze the prediction difference between the segmentor f_2^k and the total dataset which, at the k th iteration, contains the labeled set and the unlabeled set annotated by the previous segmentor f_1^{k-1} . We denote this dataset as σ_2 .

$$d(f^*, \sigma_2) = \frac{u \times d(f_1^{k-1}, f^*)}{l + u}$$

$$d(f_2^k, \sigma_2) = \frac{l \times d(f_2^k, f^*) + u \times d(f_2^k, f_1^{k-1})}{l + u}$$

Since the upper bound of the generalization error of the segmentor f_1^{k-1} is b_1^{k-1} , we have $d(f^*, \sigma_2) \leq \frac{ub_1^{k-1}}{l+u}$. Since σ_2 contains unlabeled data which may be incorrectly labeled, σ_2 must be sufficient to guarantee that if the difference of f_2^k and σ_2 is less than that of f^* which means f_2^k "learns" the mistake, then the probability that the generalization error of f_2^k is less than b_2^k is less than δ . Let $M = ub_1^{k-1}$, then the probability that f_2^k has a lower observed difference with σ_2 than f^* is less than

$$P = C_{l+u}^M d(f_2^k, \sigma_2)^M (1 - d(f_2^k, \sigma_2))^{l+u-M}$$

$$\text{Let } b_2^k = \max\left\{\frac{lb_2^0 + ub_1^0 - u \times d(f_1^{k-1}, f_1^k)}{l}, 0\right\},$$

$$d(f_2^k, \sigma_2) = \frac{l \times d(f_2^k, f^*) + u \times d(f_2^k, f_1^{k-1})}{l + u}$$

$$\geq \frac{lb_2^k + u \times d(f_2^k, f_1^{k-1})}{l + u}$$

$$\geq \frac{lb_2^0 + ub_1^0}{l + u}$$

As the function $C_s^t x^t (1-x)^{s-t}$ is monotonically decreasing in $\frac{t}{s} < x < 1$, it follows that

$$P \leq C_{l+u}^M \left(\frac{lb_2^0 + ub_1^0}{l + u}\right)^M \left(1 - \frac{lb_2^0 + ub_1^0}{l + u}\right)^{l+u-M}$$

We can approximate the RHS with Poisson Theorem.

$$C_{l+u}^M \left(\frac{lb_2^0 + ub_1^0}{l + u}\right)^M \left(1 - \frac{lb_2^0 + ub_1^0}{l + u}\right)^{l+u-M}$$

$$\approx \frac{(lb_2^0 + ub_1^0)^M}{M!} e^{-(lb_2^0 + ub_1^0)}$$

When $lb_2^0 \leq e^{\sqrt{M}} - M$,

$$\frac{(lb_2^0 + ub_1^0)^M}{M!} e^{-(lb_2^0 + ub_1^0)} \leq e^{lb_2^0}$$

We show at the beginning that $l \geq \frac{1}{b_2^0} \ln \frac{|\mathcal{H}|}{\delta}$, thus

$$P \leq e^{lb_2^0} \leq \frac{\delta}{|\mathcal{H}|}$$

Given at most $|\mathcal{H}| - 1$ (excluding the optimal f^*) segmentor with generalization error no less than b_2^k having a lower observed difference with σ_2 than f^* in hypothesis class \mathcal{H} , the probability that

$$Pr[d(f_2^k, f^*) \geq b_2^k] \leq \delta$$

. In order to let the above derivation holds, we need one more condition which is that the generalization error of f_1^{k-1} , which is the counterpart model in the last iteration, is bounded by b_1^{k-1} by probability δ . When $k = 0$, which is the initial segmentor that trains on the labeled set only, this condition is satisfied (by supervised PAC learning). When $k = 1$, the above holds as the the generalization error of f_1^0 is bounded by b_1^0 by probability δ . Then, by deduction, we can prove that the above holds for any k .

B. Quantitative Analysis of Homogenization problem

To quantitatively analyze the homogenization problem of Co-training (or to quantify the diversity between two models in the Co-training), we further propose two metrics to measure the similarity in target space. As discussed in Section 3.3, we can only quantify in the target space since measures in parameter space of different architectures is meaningless. Specifically, we use L2 distance to measure the similarity of logits output by the two models in Co-training methods.

$$D_{l2} = \frac{1}{HWC} \sum_{i=0}^{HW} \sum_{j=0}^C \|\text{logit}_{1i}^j - \text{logit}_{2i}^j\|_2$$

As the model outputs probabilistic distributions, we can also measure the similarity of models by KL Divergence.

$$D_{kl} = \frac{1}{HW} \sum_{i=0}^{HW} \sum_{j=0}^C s_{1i}^j \log \frac{s_{1i}^j}{s_{2i}^j}$$

As shown in Figure 9, we can see that Co-training with a shared backbone suffers the most from the homogenization problem while different architecture and different input domains allow more diverse model in Co-training, which is consistent with the findings in Section 3.3.

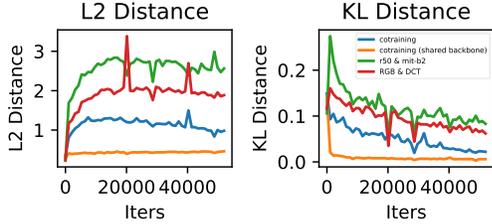


Figure 9. Demonstration of homogenization problem in Co-training

C. Quantification of Diversity in Different Techniques

After identifying the homogenization problem in Co-training methods, we provide three techniques to alleviate this problem. As discussed in Section 3.3, 5.2 and Appendix B, we first show that the three techniques can reduce the homogenization (measured by prediction similarity) and then empirically show the effectiveness of each technique individually and combined. Here, we are curious about how much diversity they each introduce, or more specifically, to compare the diversity they bring to the Co-training. We conjecture that with more diversity introduced, the empirical performance is better. The first and simplest approach is to directly and qualitatively analyze the homogenization plots. We can see that different architectures provide more diverse predictions than different input domains as well as Co-training, and Co-training (shared backbone). The second approach can quantify the diversity brought by leveraging one of the three metrics (agree rate, l2, or kld discussed in Appendix B). Due to the stochastic nature of SGD optimization, we can use an exponential moving average to estimate the metrics. An alternative can be a weighted average of the metrics at the last epoch over the whole dataset. However, we emphasize here that the three techniques tackle homogenization in three different perspectives in the training process and they mutually benefit each other as shown in the ablation study.

D. VOC PASCAL 2012 Results on ResNet101 and Comparison with SOTA

We provide comparison with ResNet101 and SegFormer-b3 on VOC PASCAL 2012 under the second partition protocol mentioned in Section 5.1. For *Diverse Co-training*, we use ResNet101 and SegFormer-b3 as backbones and compare two variants (*i.e.* 2-cps and 3-cps) with other methods with ResNet101 in Table 7. We further demonstrate the effectiveness of our *Diverse Co-training* by showing that the improvement over current SOTA methods with resolutions of 321 and 513. We outperform the previous best consistently by more than

2% with resolution of 321 and around 1% with resolution of 513. For instance, ours (3-cps) surpasses ST++ [91] by 2.8%, 2.0% and 2.0% on 1/16, 1/8 and 1/4 partition protocols respectively. We also compare with AEL [34], U²PL[82] and PS-MT [50] which obtains the best previous performance. We outperforms the best of them by 0.7%, 0.8% and 1.3% on 1/16, 1/8 and 1/4 partition protocols respectively. It’s worth mentioning that, our performance with resolution 321 already outperforms the previous SOTA with resolution 512. The remarkable performance of our *Diverse Co-training* illustrate the significance of diversity in co-training.

Method	Resolution	1/16 (662)	1/8 (1323)	1/4 (2646)
Sup Baseline	321x321	67.5	70.4	73.7
CAC [43]	321x321	72.4	74.6	76.3
CTT*	321x321	73.7	75.1	-
ST++ [91]	321x321	74.5	76.3	76.6
ours (2-cps)	321x321	77.6	78.3	78.7
ours (3-cps)	321x321	<u>77.3</u>	<u>78.0</u>	<u>78.6</u>
Sup Baseline	513x513	66.6	70.5	74.5
MT [72]	512x512	70.6	73.2	76.6
CCT [61]	512x512	67.9	73.0	76.2
GCT [38]	512x512	67.2	72.2	73.6
CPS [12]	512x512	74.5	76.4	77.7
CutMix [82]	512x512	72.6	72.7	74.3
3-CPS [20]	512x512	75.8	78.0	79.0
DSBN‡	769x769	-	74.1	77.8
ELN [42]	512x512	-	75.1	76.6
U ² PL [82]	513x513	74.4	77.6	78.7
PS-MT [50]	512x512	75.5	78.2	78.7
AEL [34]	513x513	77.2	77.6	78.1
ours (2-cps)	513x513	77.9	78.7	79.0
ours (3-cps)	513x513	<u>77.6</u>	<u>79.0</u>	<u>80.0</u>

Table 7. Comparison with state-of-the-art methods with ResNet101 on the *Pascal VOC 2012* dataset. Labeled images are sampled from the blended training set. Results of MT, CCT, GCT are from [12]. Results of CTT (denoted by *) is based on DeepLabv2 and results of DSBN (denoted by ‡) is based on Xception65

E. Detailed DCT Transform

We detailed the DCT transform in this section. As illustrated in Figure 4, we first transform images to YCbCr color space, consisting of one luma component (Y), representing the brightness, and two chroma components, Cb and Cr, representing the color. Since the spatial resolution of the Cb and Cr channel is reduced by a factor of two, we upsample the original image by two to obtain the same resolution as Y channel. The image is then converted to the frequency domain through DCT transform where each of the three Y, Cb, and Cr channels is split into blocks of 8x8

pixels and transformed to DCT coefficients of 192 channels. The two-dimensional DCT coefficients at the same frequency are grouped into one channel to form the three-dimensional DCT cubes. After the DCT transform, we obtain frequency domain input of 192 channels but with resolution downsampled by 8. Following [89], we select 64 channels (44, 10 and 10 channels each from Y, Cb and Cr components respectively) close to upper-left squares from the total 192 channels to reduce computation. We refer to [89] for more details regarding the channel selections.

Since the number of channels for frequency domain is different than the RGB domain (*i.e.* three), we have to modify the backbone to adapt it. We take ResNet [32] as an example. To be as simple as possible and further reduce training parameters and computation, we remove the stem layers at the beginning of ResNet and modify the first convolution layer in the first ResLayer to have 64 in channels.

Notice that, the above DCT transform are not contradictory to standard pre-processing techniques widely applied to RGB images it takes RGB images as input, requiring minimum modifications to the current pre-process pipeline and model architecture. To maintain the strong-weak augmentation proposed above, we first perform augmentations on RGB images and then transform it to DCT for training models on the frequency domain.

F. Comparison with Knowledge Distillation

As discussed in Section 5.4, Co-training is similar to knowledge distillation (KD) in the sense that they both possess a teaching process, the difference lies in that the teacher in KD is usually fixed and teaching is unidirectional while Co-training does not possess the "teacher" and "student" concept and the model teaches each other mutually. To demonstrate that the effectiveness of our method is not simply a knowledge transfer from one model to another but a mutually beneficial process, we compare the knowledge distillation with our method. Specifically, a Segformer with mit-b2 is trained alone and distills the knowledge to FixMatch with ResNet50. From Table 8, we show knowledge transfer do take effect improving the original FixMatch baseline by 3% 1%, which can be attributed to the diverse inductive bias and the high-quality pseudo label introduced by the transformer model. However, we show that our method still outperforms knowledge distillation by 1% consistently. This is because Co-training mutually benefits the two models while KD fails to enjoy this benefit. This can be demonstrated from Figure 3 that Co-training improves the mit-b2 by 1% while KD uses a trained and fixed model.

Method	Param	1/32	1/16	1/8	1/4
FixMatch	40.5M	70.28	73.36	74.0	74.3
FixMatch Distill	65.2M	74.1	74.9	75.6	75.8
Ours (2-cps)	65.2M	75.2	76.0	76.2	76.5

Table 8. Comparison with knowledge distillation. Labeled images are sampled from the original high-quality training set.

G. Detail of Strong Augmentation

We provide a full list of strong augmentations applied in Table 9.

CutMix is applied twice to the two different views individually. Notably, instead of batch-wise CutMix adopted by CPS [12, 90], we use in-batch CutMix which leverages the shuffled samples of the same batch to cutmix. We leverage the random cropped image directly as a weakly augmented view to generate labels. Despite CutMix is applied to each strong view individually, in-batch CutMix allows us to generate cutmixed pseudo labels by forwarding each model only once.

H. Number of Parameters

The objective of this section is to (1) demonstrate that our improvement is not trivial by simply adding more parameters and (2) facilitate a fair comparison with the SOTA method. We first report the parameters of the different architectures used in Table 3.

Backbone	Param
R50	$2 \times 40.5M = 81M$
mit-b2	$2 \times 24.7M = 49.4M$
R50 & mit-b2	$40.5M + 24.7M = 65.2M$
ResNeSt50	$2 \times 42.3M = 84.6M$
ResNeXt50	$2 \times 39.8M = 79.6M$
R50 & ResNeSt50	$40.5M + 42.3M = 82.8M$
R50 & ResNeXt50	$40.5M + 39.8M = 80.3M$

Table 10. We show the parameters of each architecture.

As per Table 10, our R50 & mit-b2 possess 20M parameters less than CNN variants such as R50 & ResNeSt50 and R50 & ResNeXt50 but still achieve better performance. Then we compare FixMatch-Distill and FixMatch-Ensemble which uses exactly the same or more parameters than ours but a different learning paradigm. FixMatch-Distill uses a trained Segformer-b2 to distill knowledge to ResNet50 model as described in Appendix F. FixMatch-Ensemble is an ensemble of two ResNet50 model is uses 20M parameters more than ours. As shown in the first section of Table 11, our model outperforms both FixMatch-Distill and FixMatch-Ensemble consistently by a large margin. This demonstrates that the improvements by our *Diverse Co-training* is not trivially by adding more parameters. Finally, we also compare the parameters used in our

Weak Augmentation	
Random Rescale	Resizes randomly the image by [0.5, 2.0].
Random Flip	Flip the image horizontally with a probability of 0.5.
Random Crop	Randomly crop a region from the image.
Strong Augmentation	
Color Jitter	Randomly jitter the color space of the image with a probability of 0.8.
Gaussian Blur	Blurs the image with a Gaussian kernel with a probability of 0.5.
Random Grayscale	Turn the image to greyscale with a probability of 0.2.
Cutmix	Cut a patch from one image and paste the patch to another image. We always apply Cutmix to every image.

Table 9. List of various image transformations.

method and the previous SOTA methods. CPS [12] uses two models to perform Co-training while n-CPS (n=3) [21] uses three. Although PS-MT [50] uses only one architecture, they leverage two teachers (which are two different sets of parameters) and one student which equals three times the parameters of one model. U²PL [82] leverages the popular teacher-student framework which also leverages two sets of parameters. We show dominant performance with 20M parameters less which further demonstrates the effectiveness of our *Co-training*.

Method	Param	1/32	1/16	1/8	1/4
FixMatch Ensemble	81.0M	73.0	74.3	75.6	75.9
FixMatch Distill	65.2M	74.1	74.9	75.6	75.8
CPS [12]	81.0M	-	72.0	73.7	74.9
n-CPS (n=3) [21]	121.5M	-	72.0	74.2	75.9
PS-MT [50]	121.5M	-	72.8	75.7	76.4
U ² PL* [82]	81M	-	72.0	75.1	76.2
Ours (2-cps)	65.2M	75.2	76.0	76.2	76.5

Table 11. Comparison of parameters and performance with different learning paradigms and previous SOTA. Labeled images are sampled from the original high-quality training set.

I. Visualization

Figure 10 visualizes some segmentation results on *PASCAL VOC 2012* validation set. First, we can observe the better results obtained by co-training methods (*i.e.* (d) and (e)) as shown in the third and last row, where FixMatch is prone to under-segmentation (classifies many foreground pixels as background). Our *Diverse Co-training*, compared with co-training baseline, can better segments the small objects that FixMatch and co-training baseline tends to ignore (*e.g.* the forth and fifth row). The FixMatch and co-training baseline tends to ignore some foreground while our *Diverse Co-training* does not, such as the visualization of the second row. These visualization further demonstrate the remarkable performance of *Diverse Co-training* and proves the argument that diversity matters significantly in co-training.

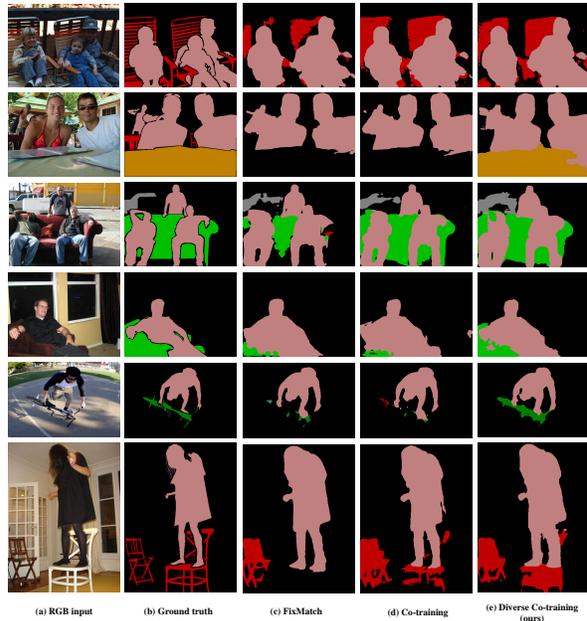


Figure 10. Example qualitative results from *PASCAL VOC 2012*. (a) RGB input; (b) ground truth; (c) FixMatch; (d) Co-training baseline; (e) Diverse Co-training (ours). (c) and (d) use DeepLabv3+ with ResNet50 as the segmentation network while (e) uses DeepLabv3+ with ResNet50 and SegFormerb2 (with MLP head) as the two segmentation networks.

J. Full Comparison with SOTA on Pascal VOC 2012

Due to limited space, we only compare the most recent SOTA in Section 5.3. We provide a full comparison here.

K. Full Ablation Study

We further provide a table to show the importance and performance gain of each component. As per table 14, we can see that all component is effective when incorporate into the holistic framework. The combination of diverse domains and different architecture provides the best result of 75.21%, 75.85% and 76.23% on 1/32, 1/16 and 1/8 labeled data.

Method	Resolution	92	183	366	732	1464
ResNet50						
Sup Baseline	513x513	39.1	51.3	60.3	65.9	71.0
PseudoSeg [103]	512x512	54.9	61.9	64.9	70.4	-
PC ² Seg [100]	512x512	56.9	64.6	67.6	70.9	-
Ours (2-cps)	513x513	<u>71.8</u>	<u>74.5</u>	77.6	<u>78.6</u>	<u>79.8</u>
Ours (3-cps)	513x513	73.1	74.7	<u>77.1</u>	78.8	80.2
ResNet101						
Sup Baseline	321x321	44.4	54.0	63.4	67.2	71.8
PseudoSeg [103]	321x321	57.6	65.5	69.1	72.4	73.2
PC ² Seg [100]	321x321	57.0	66.3	69.8	73.1	74.2
ReCo [49]	321x321	64.8	72.0	73.1	74.7	-
ST++ [91]	321x321	65.2	71.0	74.6	77.3	79.1
ours (2-cps)	321x321	<u>74.8</u>	77.6	<u>79.5</u>	<u>80.3</u>	81.7
ours (3-cps)	321x321	75.4	<u>76.8</u>	79.6	80.4	<u>81.6</u>
Sup Baseline	512x512	42.3	56.6	64.2	68.1	72.0
MT [72]	512x512	48.7	55.8	63.0	69.16	-
GCT [38]	512x512	46.0	55.0	64.7	70.7	-
CTT* [83]	512x512	64	71.1	72.4	76.1	-
CPS[12]	512x512	64.1	67.4	71.7	75.9	-
U ² PL [82]	512x512	68.0	69.2	73.7	76.2	79.5
PS-MT [50]	512x512	65.8	69.6	76.6	78.4	80.0
ours (2-cps)	513x513	76.2	<u>76.6</u>	80.2	<u>80.8</u>	<u>81.9</u>
ours (3-cps)	513x513	<u>75.7</u>	77.7	<u>80.1</u>	80.9	82.0

Table 12. Full Comparison with state-of-the-art methods on the *Pascal* dataset. Labeled images are from the original high-quality training set. Results of CTT (denoted by *) is based on DeeplabV2.

Method	Resolution	1/32 (331)	1/16 (662)	1/8 (1323)	1/4 (2646)
Sup Baseline	321x321	55.8	60.3	66.8	71.3
CAC[43]	320x320	-	70.1	72.4	74.0
ST++[91]	321x321	-	72.6	74.4	75.4
Ours (2-cps)	321x321	75.2	<u>76.0</u>	<u>76.2</u>	<u>76.5</u>
Ours (3-cps)	321x321	<u>74.9</u>	76.4	76.3	76.6
Sup Baseline	513x513	54.1	60.7	67.7	71.9
CutMix [82]	512x512	-	68.9	70.7	72.5
CCT [61]	512x512	-	65.2	70.9	73.4
GCT [38]	512x512	-	64.1	70.5	73.5
CPS[12]	512x512	-	72.0	73.7	74.9
3-CPS [20]	512x512	-	72.0	74.2	75.9
ELN [42]	512x512	-	-	73.2	74.6
PS-MT [50]	512x512	-	72.8	75.7	76.4
U ² PL* [82]	513x513	-	72.0	75.1	76.2
Ours (2-cps)	513x513	75.2	<u>76.2</u>	<u>77.0</u>	<u>77.5</u>
Ours (3-cps)	513x513	<u>74.7</u>	76.3	77.2	77.7

Table 13. Full Comparison with state-of-the-art methods with ResNet50 on the *Pascal VOC 2012* dataset. Labeled images are sampled from the blended training set. The result of U^2PL is reproduced with the same setting as ours.

Table 14. **Ablation study of different component combinations** on PASCAL VOC dataset with ResNet50 and SegFormer-b2. The results are obtained under 1/32, 1/16 and 1/8 partition protocols and the observations are consistent for other partition protocols. L^s represents the supervision loss on the labeled data. L^u represents the pseudo supervision loss on the unlabeled data. SA (strong augmentation) denotes strong-weak augmentation is used. Diff SA stands for different strong augmentation for each model. Diff domain means using RGB and frequency domain to train separate models with cross supervision. Diff arch means different architectures are used to instantiate the two models.

Components						PASCAL VOC		
L^s	L^u	SA	diff SA	diff domain	diff arch	1-32	1-16	1-8
✓						55.78	60.3	66.79
✓	✓					65.66	71.28	73.77
✓	✓	✓				70.28	73.36	74.82
✓	✓		✓			69.45	72.43	74.84
✓	✓			✓		71.58	74.94	75.97
✓	✓	✓	✓			71.07	74.09	74.98
✓	✓	✓		✓		72.00	74.10	74.93
✓	✓	✓			✓	74.89	75.82	76.08
✓	✓	✓		✓	✓	75.21	75.99	76.23