

Adaptive Spatial Augmentation for Semi-supervised Semantic Segmentation

Lingyan Ran, Yali Li, Tao Zhuo, Shizhou Zhang, Yanning Zhang, *Fellow, IEEE*

Abstract—In semi-supervised semantic segmentation (SSSS), data augmentation plays a crucial role in the weak-to-strong consistency regularization framework, as it enhances diversity and improves model generalization. Recent strong augmentation methods have primarily focused on intensity-based perturbations, which have minimal impact on the semantic masks. In contrast, spatial augmentations like translation and rotation have long been acknowledged for their effectiveness in supervised semantic segmentation tasks, but they are often ignored in SSSS. In this work, we demonstrate that spatial augmentation can also contribute to model training in SSSS, despite generating inconsistent masks between the weak and strong augmentations. Furthermore, recognizing the variability among images, we propose an adaptive augmentation strategy that dynamically adjusts the augmentation for each instance based on entropy. Extensive experiments show that our proposed Adaptive Spatial Augmentation (ASAug) can be integrated as a pluggable module, consistently improving the performance of existing methods and achieving state-of-the-art results on benchmark datasets such as PASCAL VOC 2012, Cityscapes, and COCO.

Index Terms—Semi-supervised learning, Semantic Segmentation, Data Enhancement

I. INTRODUCTION

SEMANTIC segmentation focuses on assigning semantic labels to every pixel in an image, and it has been widely used in various domains, including analyses of natural images [1], [2], medical images [3], remote sensing [4], autonomous driving [5], etc. The success of supervised semantic segmentation methods depends on a large collection of high-quality pixel-level annotated images for training. However, the process of pixel-wise labeling is both labor-intensive and time-consuming. This challenge has led to the rise of semi-supervised learning (SSL) [6], [7] in semantic segmentation, which utilizes both annotated and unannotated data to train models, providing a more efficient and scalable approach.

Recent semi-supervised semantic segmentation (SSSS) approaches [8]–[10] often adopt the weak-to-strong consistency regularization framework [11]–[13]. These methods train a teacher model on weakly augmented data and a student model on strongly augmented data, ensuring consistent outputs despite the perturbations introduced by the augmentations. Weak augmentations typically include operations like image flipping, cropping, and scaling, while strong augmentations

This work is supported in part by the National Natural Science Foundation of China (62476226). (*Corresponding author: Tao Zhuo.*)

Lingyan Ran, Yali Li, Shizhou Zhang, and Yanning Zhang are with the Shaanxi Provincial Key Laboratory of Speech and Image Information Processing, and the National Engineering Laboratory for Integrated Aerospace-Ground-Ocean Big Data Application Technology, School of Computer Science, Northwestern Polytechnical University, Xi’an 710072, China. Tao Zhuo is in the College of Information Engineering, Northwest A&F University, Yangling, 712100, China.

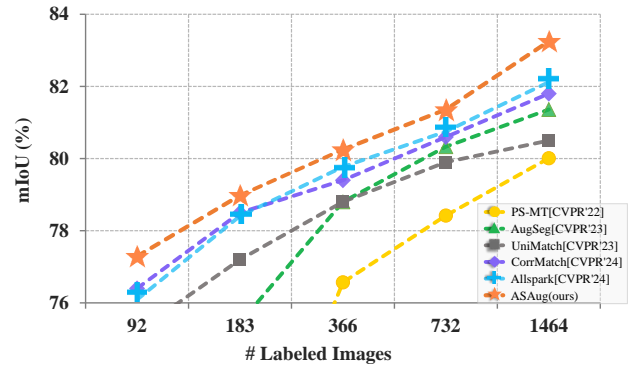


Fig. 1: Comparison with SOTA methods on the Pascal VOC 2012 dataset. Notably, our method outperforms other approaches in all partitioning scenarios.

are generally intensity-based [14], [15], such as identity, auto-contrast, Gaussian blur, equalize, sharpness, brightness, hue, color jitter, posterize, and solarize. As highlighted in [15], strong augmentations are a crucial component that enhances model training.

This study finds that commonly employed spatial augmentations in supervised semantic segmentation, like rotation and translation [16]–[18], are often not utilized for strong augmentation in semi-supervised scenarios. Spatial augmentations differ from intensity-based ones by altering pixel positions, thus changing the mask and potentially causing discrepancies between weak and strong augmentations. This observation prompts the following inquiry: can spatial augmentations enhance SSSS comparably to their impact on the supervised methods?

To address this question, we first conducted experiments by replacing intensity-based augmentations with spatial-based ones, specifically rotation and translation (see Table IV). Unexpectedly, our findings reveal that the inconsistent masks produced by these spatial augmentations enhance generalization. This improvement likely stems from the segmentation model’s need to accurately identify target regions under varying conditions, as spatial augmentations create a distinct gap between weak and strong augmentation scenarios. Furthermore, recognizing the significant variability among individual instances, we introduce an adaptive spatial augmentation (ASAug) technique that leverages entropy to modulate augmentation strength, thereby reinforcing consistency regularization. This approach draws inspiration from methods such as autoaugment [16], randaugment [17], and trivlaugment [18], all of which aim to automatically discover more effective data augmentation policies for image classification. Since our

TABLE I: Comparison of recent SSSS solutions in terms of “More Supervision”, “Augmentation”, and “Hybrid Techniques” which are combined with other algorithms (in order of publication year). We explain the acronyms as follows. “CATP”: co-training/auxiliary trainable-network/pseudo-rectifying, including dual-model co-training and trainable auxiliary networks, “MTS”: multiple training stages. “SDA”: strong data augmentation, “MIE”: multi-stream/multi-branch interaction/evaluation. “CR”: consistency regularisation, here mainly refers to introducing perturbations. “CL”: contrast learning. ASAug has a simple enough architecture and better performance.

Method	More Sup.		Aug.		Hybrid-tec.	
	CATP	MTS	SDA	MIE	CR	CL
CCT [19]	✓				✓	
C3-semi. [20]	✓		✓		✓	✓
DMT [21]			✓			
RoCo [22]			✓			✓
CPS [23]	✓		✓			
ST++ [24]		✓	✓			
ELN [25]	✓	✓	✓			
PS-MT [26]			✓		✓	
U2PL [27]			✓			✓
UCC [28]			✓	✓	✓	
UniMatch [13]			✓	✓	✓	
CCVC [29]	✓		✓	✓	✓	
iMAS [30]			✓	✓		
CISC-R [31]		✓	✓			
CorrMatch [32]	✓		✓			
Allspark [33]			✓	✓		
ASAug			✓			

method is pluggable, it can be easily integrated into existing weak-strong consistency regularization frameworks. Extensive experiments on three benchmark datasets demonstrate the effectiveness of our process to improve performance and achieve SOTA performance compared to other methods.

As shown in Table I, unlike previous strategies, we break away from traditional enhancement techniques by focusing on the essence of the segmentation task and practically analyzing the importance and substantial effect of enhancement.

Our contributions can be summarized in three folds:

- To the best of our knowledge, this is the first study to demonstrate that spatial augmentations, often overlooked in existing SSSS methods, can boost generalization, even though they produce inconsistent masks.
- We introduce a pluggable ASAug module that adaptively applies strong augmentations to each instance with entropy-based adaptive weight (EAW), improving consistency regularization.
- Extensive experiments on three well-known benchmarks, PASCAL VOC 2012 [34], Cityscapes [35], and COCO [36], demonstrate that our method surpasses the state-of-the-art methods by a noticeable margin.

II. RELATED WORK

A. Semi-Supervised Learning

Semi-supervised learning(SSL) predominantly focuses on training models with a limited quantity of labeled data alongside a substantial volume of unlabeled data. Recent research

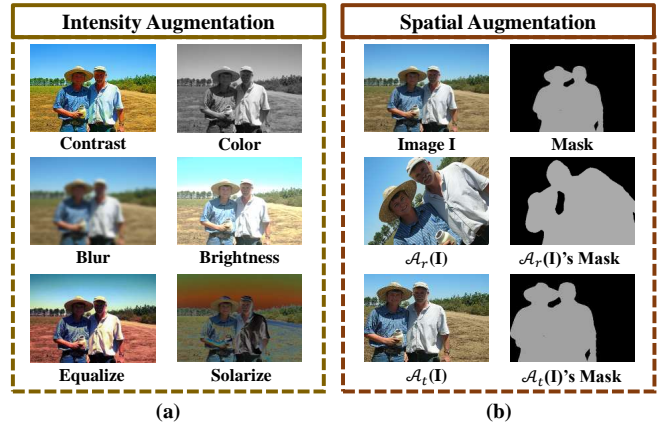


Fig. 2: Comparisons between the intensity and spatial augmentations. (a) intensity-based augmentations, like contrast, color, blur, brightness, and contrast adjustments, modify pixel appearance without changing spatial positions of the original mask; (b) spatial augmentations, such as rotation and translation, directly change the positions of pixels, leading to inconsistent masks between the original image and the augmented image.

branches into two principal categories: consistency regularization [19], [37] and pseudo-labeling [38], [39]. Consistency regularization entails training models to achieve high robustness by learning from both labeled and unlabeled data through perturbations, whereas pseudo-labeling aims to expand the dataset by having the model generate pseudo-labels for unlabeled data, subsequently utilizing both true and pseudo-labeled data for training. FixMatch [40] synthesizes these strategies by creating a hybrid method that leverages weakly augmented data for pseudo-label generation and employs cross-entropy for consistency regularization. Following this, FlexMatch [41] introduced curriculum pseudo-labels, which modify category thresholds without adding extra parameters or computations. To address the heavy reliance on labeled data, our proposed approach integrates SSL to significantly improve the performance of semantic segmentation models while reducing the need for annotated data.

B. Semi-Supervised Semantic Segmentation

SSL for semantic segmentation delivers outstanding performance and alleviates the challenge of limited labeled data for researchers. The survey [42] introduces a taxonomy that systematically categorizes these techniques into pseudo-labeling (PL) [8], [38], consistency regularization (CR) [43], contrastive learning (CL) [44], adversarial training (AT) [45], and hybrid methods (HM) [46].

Eary AT-related research utilized a common framework known as Generative Adversarial Networks (GANs) [47], which can be chiefly categorized into two types depending on their structure: those with Generators [48], [49] and those lacking Generators [50], [51]. Subsequently, the scientists explored leveraging CL to gain more meaningful representations within the embedding space to enhance segmentation results. RoCo [22] addresses the ambiguity among confusion classes

by employing contrast learning at the region level, while PRCL [52] suggests interpreting contrast learning through probability expressions. PL enhances the labeled dataset by incorporating pseudo-labels and is gaining traction because of its straightforward implementation and interpretability. The study ST++ [24] introduces a reliability retraining method that involves multiple stages, while CISC-R [31] recommends querying labeled images to refine inaccurate pseudo-labels. In the field of CR, there is a consensus on the smoothing assumption, which entails that identical predictions should be obtained for a pixel and its perturbed version. Building on FixMatch [40], both PseudoSeg [53] and UniMatch [13] apply strong and weak consistency to segmentation tasks. Additionally, the integration of these techniques has resulted in notable outcomes. U2PL [27] and DGCL [54] merge PL and CL to optimize pseudo-labels via feature similarity between samples, thereby steering the model toward learning more precise and robust feature representations. A prevalent approach [29], [40] nowadays involves integrating PL with CR under a consistency framework aimed at enhancing the confirmation bias and potentially boosting accuracy further. PC2Seg [55] uses feature space comparison learning and consistency training.

C. Data Augmentation Study in SSSS

Sufficient data is essential for tasks using deep learning methods, such as image classification, semantic segmentation, and change detection. With limited actuation, data augmentation [14], [16] is widely used in various scenarios. Methods like random shuffle or intensity change are the most common. Although many approaches have been developed in SSSS, ST++ [24] demonstrates that the simplest self-training pipeline in PL coupled with strong data augmentation, can achieve excellent performance without any modular improvement. This proves that data augmentation plays a crucial role with limited labeled data. Yuan et al. [14] introduced a simple yet efficient framework as a baseline for robust data enhancement techniques. CutMix [57] effectively exploited the regularization effect of training pixels and the loss of retained regions by cutting and pasting patches in the training image, and subsequent ClassMix [58] and ComplexMix [59] were improved on this basis. AugSeg [15] adaptively augmented samples based on confidence by randomly selecting various techniques. They also argue that most research ignores the variability between instances and proposes an instance-specific model adaptive supervision method called iMAS [30].

As shown in Fig. 2, traditional augmentation aims to artificially increase the diversity of the training data by performing various transformations or perturbations on the existing dataset. However, it may not adequately reflect the complexity and diversity required for a robust semantic segmentation model. Several researchers have noted the effect of positional factors on segmentation tasks and proposed geometric enhancements. Cao et al. [50] proposed a context-aware unsupervised strategy for micro-aggregate warping. M3L [60] introduces a robust perturbation model incorporating geometric warping and photometric variations. MR-PhTPS [61] relies on nonlinear geometric and photometric perturbations. Unlike the above methods, we address the issue of the importance

Algorithm 1 Pseudocode of ASAug in a PyTorch-like Style

```

# f/f_t: model / teacher_model
# aug_w: weak enhancements
# aug_spatial: spatial enhancement (A_r & A_t)
# k_r/k_t, d_r/d_t: scaling, offset parameters
# r_max/t_max: maximum rotation angle and
maximum translation ratio
# B: batch size, C: number of categories,
H/W: height/width of image

for each x in loader_u: do
  x_w = aug_w(x) # [B, H, W]
  x_w_probmap = f_t(x_w) # [B, C, H, W]
  # calculate the entropy
  p_w = x_w_probmap.argmax(dim=1)
  entropy = comp_entropy(x_w_probmap)
  norm_entropy = entropy.mean()

  # obtain spatial enhancement parameters
  A_r, A_t = spatial_by_entropy(norm_entropy,
  k_r, d_r, k_t, d_t)
  # apply A_r and A_t, get predictions
  x_spatial = aug_spatial(x, A_r, A_t)
  p_spatial = f(x_spatial)
  # apply the same enhancements to p_w
  p_w_spatial = aug_spatial(p_w, A_r, A_t)

  # loss function
  criterion = nn.MSELoss()
  loss_spatial = criterion(p_spatial,
  p_w_spatial)
end for

```

of geometric or positional factors on SSSS and accordingly propose an insertable adaptive mapping augmentation module that improves the efficiency of using unlabeled data through geometrically strong augmentation, which has been neglected in previous work.

III. METHOD

The proposed method is built upon the weak-to-strong consistency regularization framework [13], [30], [32], [33]. Unlike recent approaches that primarily focus on intensity-based augmentation, we propose that spatial-based augmentations can also enhance model generalization. Fourthmore, taking into account the variability among samples, we introduce an adaptive augmentation strategy. The details of our approach are introduced below.

In addition, to explain our ASAug in more detail, we show the overall architecture of our approach in Fig. 3, and the associated pseudo-code in Algorithm 1 in a PyTorch-like style for ease of understanding.

A. Preliminaries

The SSSS task focuses on training a segmentation model that optimally adapts to data distributions with a scarcity of labeled data and a surplus of unlabeled data. The labeled dataset $D^l = \{(x^l, y^l)\}^M$ consists of M examples with labels, while the unlabeled dataset $D^u = \{x^u\}^N$ includes N images without labels, where $N \gg M$. Our objective is to optimize a model utilizing the mean-teacher (MT) [62]

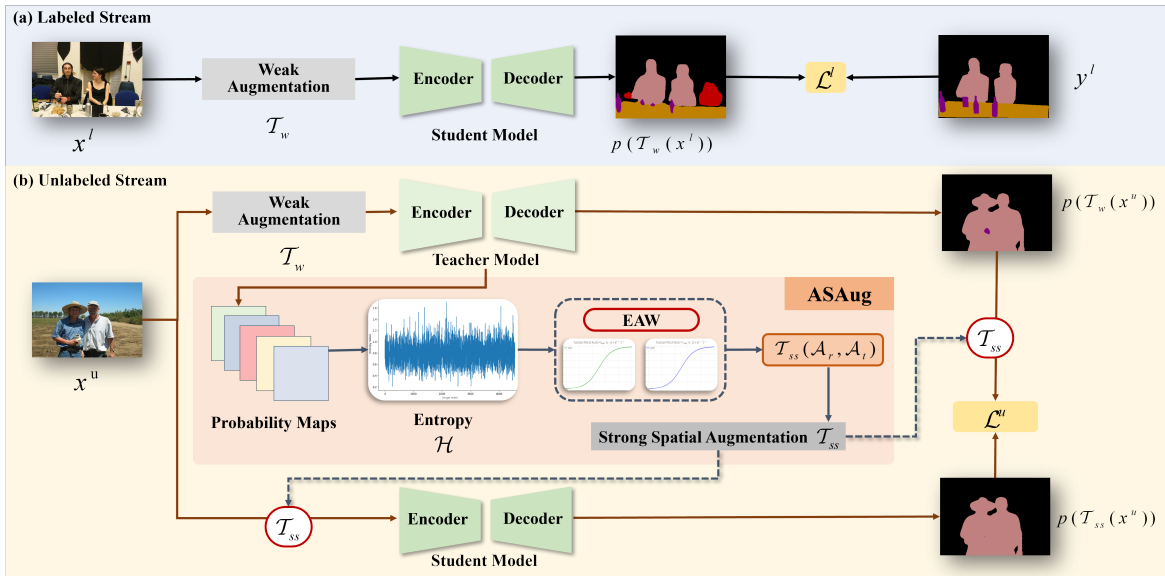


Fig. 3: Illustration of our ASAAug pipeline. Based on the teacher-student consistency training framework [56], we introduce an adaptive method that can selectively distort images as geometrically strong enhancements based on their reliability and importance. “EAW” denotes our entropy-based adaptive weight. Notably, geometric adjustment changes pixel point locations, and we apply the same operation to the teacher’s output to ensure the consistency of predictions before and after enhancement.

strategy, in which the prediction from a teacher model is regarded as ground truth for weakly augmented unlabeled instances. Initially, the student model undergoes training with labeled images, followed by optimization on unlabeled data. The teacher model begins as a copy of the student model and is subsequently refined using the EMA [62] protocol. The training loss \mathcal{L} can be minimized by utilizing both the labeled \mathcal{L}^l and unlabeled \mathcal{L}^u datasets.

$$\mathcal{L} = \mathcal{L}^l + \lambda \mathcal{L}^u, \quad (1)$$

where λ represents the hyper-parameter that balances the supervised and unsupervised trade-offs. It can be set to a constant value [40], [63] or dynamically modified [20], [64] throughout the training phase.

The supervised loss \mathcal{L}^l can be expressed as:

$$\mathcal{L}^l = -\frac{1}{M} \sum_{l=1}^M \sum_{c=1}^C y_c^l \log(p_c^l(\mathcal{T}_w(x^l); \theta)) \quad (2)$$

where $p_c^l(\cdot; \theta)$ is the model’s predicted probability that the sample x^l belongs to category c . C is the number of categories. $\mathcal{T}_w(\cdot)$ represents the weak augmentation, a commonly used method. y_c^l is the ground truth for sample x^l with category label c (one-hot).

As for the unsupervised loss \mathcal{L}^u , the key to SSL, is computed differently depending on the chosen method. Traditional weak-to-strong consistent regularization (WSCR) relies on strong and weak augmentation strategies to produce divergent predictions for the same input, obtaining the segmentation predictions $p_s^u = f(\mathcal{T}_s(x^u))$ for the student model on the strong augmentation $\mathcal{T}_s(\cdot)$ and $p_w^u = f(\mathcal{T}_w(x^u))$ for the teacher model on the weakly augmentation $\mathcal{T}_w(\cdot)$ to compute the loss, respectively:

$$\mathcal{L}^{\text{consistency}} = -\frac{1}{N} \sum_{u=1}^N \sum_{c=1}^C p_{w,c}^u \log(p_{s,c}^u) \quad (3)$$

where $\mathcal{T}_s(\cdot)$ denotes intensity enhancement transform applied to x^u (e.g., brightness). This is also the main technique recently used for WSCR in SSSS and has shown effectiveness. In contrast, widely used spatial enhancements (e.g., rotations and translations) are usually ignored. Strong spatial augmentation can effectively improve the model robustness, therefore, we propose an adaptive spatial augmentation (ASAAug), the details are described below.

B. Adaptive Spatial Augmentation

In particular, we substitute intensity-focused robust augmentations with strong spatial augmentation \mathcal{T}_{ss} , which involves both rotational \mathcal{A}_r and translational \mathcal{A}_t transformations. Compared to pixel level modifications like jittering and brightness change, \mathcal{A}_r and \mathcal{A}_t import more vibrant disturbance to unlabeled training samples.

- **Rotation \mathcal{A}_r :** Real-world objects and scenes are viewed from multiple perspectives. By rotating images, we simulate these different angles, helping the model identify patterns and characteristics across orientations. In ASAAug, we adjust the rotation within a restricted angle.
- **Translation \mathcal{A}_t :** In reality, cameras frequently adjust their positions horizontally or vertically to capture scenes from different angles. Translating images allows us to mimic these camera movements, which helps the model identify patterns and features in varied locations. In ASAAug, we achieve this by shifting an image sideways or up and down by a specified pixel count.

To tackle the challenges posed by increased uncertainty biases, we propose the adaptive spatial augmentation that modifies the augmentation intensity based on the sample’s quality. In WSCR works, the prediction of weak augmented image $\mathcal{T}_w(x^u)$ from the teacher model is treated as ground truth for x^u . We use the entropy of this prediction $p^u(\mathcal{T}_w(x^u))$ as the measurement for instance reliability. Specifically, we

TABLE II: Comparison with SOTA methods on the **Pascal classic** val set (Indicator: $mIoU$). Two crop sizes are tested.

Method	Size	1 / 16 (92)	1 / 8 (183)	1 / 4 (366)	1 / 2 (732)	Full (1464)
Supervised	513×513	50.43	56.20	64.73	67.37	70.35
CutMix-Seg [65]	513×513	55.58	63.20	68.36	69.84	76.54
PseudoSeg [53]	513×513	57.60	65.50	69.14	72.41	—
CPS [23]	513×513	64.07	67.42	71.71	75.88	—
PS-MT [26]	513×513	65.80	69.58	76.57	78.42	80.01
GTA [66]	513×513	70.00	73.20	75.60	78.40	80.50
AugSeg [15]	513× 513	71.09	75.45	78.80	80.33	81.36
CCVC [29]	513× 513	70.20	74.40	77.40	79.10	80.50
Allspark [33]	513×513	76.09	78.41	79.77	80.75	82.12
Allspark w/ ASAug	513×513	77.29	78.98	80.25	81.36	83.25
↑ Δ (%)	—	↑ 1.20	↑ 0.57	↑ 0.48	↑ 0.61	↑ 1.13
CorrMatch [~] [32]	513×513	76.41	77.83	79.57	80.66	81.93
CorrMatch w/ ASAug	513×513	77.06	78.75	80.09	81.03	82.40
↑ Δ (%)	—	↑ 0.65	↑ 0.92	↑ 0.52	↑ 0.37	↑ 0.47
Supervised	321×321	48.71	55.35	60.13	66.26	70.01
Mean Teacher [62]	321×321	52.72	58.93	65.92	69.54	72.42
UniMatch [13]	321×321	75.20	77.20	78.80	79.90	81.20
Allspark [~] [33]	321×321	72.36	76.31	77.07	79.56	80.62
Allspark w/ ASAug	321×321	76.52	78.76	79.92	81.11	83.00
↑ Δ (%)	—	↑ 4.16	↑ 2.45	↑ 2.85	↑ 1.55	↑ 2.38
CorrMatch [32]	321×321	76.40	78.50	79.40	80.60	81.80
CorrMatch w/ ASAug	321×321	76.91	79.31	79.98	80.94	82.28
↑ Δ (%)	—	↑ 0.51	↑ 0.81	↑ 0.58	↑ 0.34	↑ 0.48

[~] denotes the results we reproduced. We split the table with the different crop sizes.

assess the prediction’s entropy \mathcal{H} as an indicator of reliability and adjust the level of strong spatial augmentation $\mathcal{T}_{ss}(\mathcal{A}_r, \mathcal{A}_t)$ accordingly. For the predicted output probability $p^u(\cdot)$, the \mathcal{H} is defined as follows:

$$\mathcal{H} = - \sum_{c=1}^C p_c^u(\mathcal{T}_w(x^u)) \log(p_c^u(\mathcal{T}_w(x^u))) \quad (4)$$

where C denotes the number of categories, $p_c^u(\mathcal{T}_w(x^u))$ is the predicted probability of sample x^u on category c . The higher the value of \mathcal{H} , the more uncertain the model is about its prediction for that particular sample, and conversely.

To fine-tune the level of augmentation distortion, we introduce an entropy-based adaptive weight, utilizing information entropy as a measure of sample reliability. This means that the rotation angle and the degree of translation are modulated according to the entropy value of the weakly augmented prediction, denoted by \mathcal{H} . EAW dynamically adjusts the strength of spatial transformations. Samples with high entropy, where the model exhibits greater uncertainty, require more significant spatial transformations to explore a broader feature space. Conversely, samples with low entropy, where the model demonstrates greater certainty, are better suited for smaller transformations, thus preserving stable features. This EAW approach enables a smoother adjustment of augmentation across samples, enhancing the model’s capacity and resilience to spatial transformations.

We define our EAW used in adaptive spatial augmentation $\mathcal{A}_r(\mathcal{H})$ and $\mathcal{A}_t(\mathcal{H})$ as follows:

$$\mathcal{A}_r(\mathcal{H}) = r_{\max} \cdot k_r \cdot (1 + e^{d_r - \mathcal{H}})^{-1} \quad (5)$$

$$\mathcal{A}_t(\mathcal{H}) = t_{\max} \cdot k_t \cdot (1 + e^{d_t - \mathcal{H}})^{-1} \quad (6)$$

where k_r, k_t are the scaling parameters, r_{\max} and t_{\max} denote the maximum rotation angle ($^\circ$) and maximum translation ratio

($\times 100\%$), respectively. d_r and d_t are the offsets associated with the two spatial transformations to fine-tune the input entropy values. We combine \mathcal{A}_r and \mathcal{A}_t to form the strong spatial augmentation $\mathcal{T}_{ss}(\mathcal{A}_r, \mathcal{A}_t)$.

When \mathcal{H} is relatively small, the mapping outcome is likewise reduced, which decreases the magnitude of rotation and translation, thereby maintaining the stability of the applied enhancement.

C. Pixel-level Consistency Learning

When applying spatial transformations, the structure of the image changes, and so does the position of the mask. To ensure that the loss is a correct measure of the difference in prediction before and after enhancement, we apply the same spatial transformations to the weakly enhanced forecasts as we do to the strongly enhanced ones to align their spatial structure. We adopt the mean squared error (MSE) as the metric of consistency loss. Compared to other losses, like cross-entropy loss, MSE is smoother and more robust in measuring the variance of consecutive predicted values. It is particularly suitable for calculating prediction consistency that does not depend on category labels under geometric enhancement.

The pixel-level consistency loss \mathcal{L}^u for unlabelled samples is defined as follows:

$$\mathcal{L}^u = \frac{1}{N} \sum_{u=1}^N \frac{1}{HW} \sum_{i=1}^H \sum_{j=1}^W (\mathcal{T}_{ss}(p(\mathcal{T}_w(x_{i,j}^u))) - p(\mathcal{T}_{ss}(x_{i,j}^u)))^2 \quad (7)$$

where H and W are the height and width of the image, respectively. The difference in predicted values at each pixel location (i, j) is squared summed, and averaged to obtain a pixel-level MSE for that image. MSE directly measures spatially aligned pixel differences without sensitivity to the

TABLE III: Comparison with SOTA methods on the **Pascal blender** val set (Indicator: $mIoU$). The split table reports results by following U2PL [27] configurations.

Method	Size	1/16 (662)	1/8 (1323)	1/4 (2646)	1/2 (5291)
Supervised	321×321	65.58	70.37	72.46	73.61
CAC [63]	321×321	72.40	74.60	74.30	—
UniMatch [13]	321×321	76.50	77.00	77.20	—
CutMix-Seg [65]	513×513	72.56	72.69	74.25	75.89
CPS [23]	513×513	72.18	75.83	77.55	78.64
PS-MT [26]	513×513	75.50	78.20	78.72	79.76
ESL [67]	513×513	76.36	78.57	79.02	79.98
Allspark [~] [33]	321×321	76.89	78.18	80.02	79.86
w/ ASAug	321×321	77.60	80.22	80.97	80.37
↑ Δ (%)	—	↑ 0.71	↑ 2.04	↑ 0.95	↑ 0.51
CorrMatch [32]	321×321	77.60	77.80	78.30	—
w/ ASAug	321×321	78.52	78.27	78.87	—
↑ Δ (%)	—	↑ 0.92	↑ 0.47	↑ 0.57	—
U2PL [†] [27]	513×513	74.40	77.60	78.70	—
AugSeg [†] [15]	513×513	77.01	78.20	78.82	—
CCVC [†] [29]	513×513	77.20	78.40	79.00	—
iMAS [†] [30]	513×513	77.20	78.40	79.30	—
Allspark ^{†~}	321×321	80.44	81.00	79.34	78.91
w/ ASAug[†]	321×321	81.25	82.39	82.70	80.88
↑ Δ (%)	—	↑ 0.81	↑ 1.39	↑ 3.36	↑ 1.97
CorrMatch ^{†~}	321×321	78.44	80.64	79.08	78.14
w/ ASAug[†]	321×321	79.55	81.77	79.85	78.69
↑ Δ (%)	—	↑ 1.11	↑ 1.13	↑ 0.77	↑ 0.55

[~] denotes the results we reproduced.
[†] means the same split as U2PL [27], which are contained in split table.

TABLE IV: Ablation study on our ASAug. \mathcal{A}_r and \mathcal{A}_t represent the two main augmentation operations, respectively. Improvements to the baseline are highlighted in blue.

ASAug		mIoU(%)	
\mathcal{A}_r	\mathcal{A}_t	732	366
		77.85	76.77
✓		78.71 (↑ 0.86)	78.04 (↑ 1.27)
	✓	78.18 (↑ 0.33)	77.30 (↑ 0.53)
✓	✓	80.94 (↑ 3.09)	79.98 (↑ 3.21)

mask structure, improving the model’s robustness to the augmentation transform.

IV. EXPERIMENTS

This section starts with a thorough explanation of our experimental setup in subsection IV-A. Following that, subsection IV-B is dedicated to assessing ASAug in comparison with leading algorithms on two popular SSSS benchmarks. Next, in subsection IV-C, we perform ablation studies to further validate and confirm the robustness of ASAug. Finally, subsection IV-D offers a qualitative analysis, supported by visual results of Cityscapes.

A. Implementation Details

Datasets. We investigate the effects of ASAug on three benchmark segmentation datasets: Pascal VOC 2012 [68], Cityscapes [35], and COCO [36]. The Pascal VOC 2012

TABLE V: Comparison with SOTAs on the **Cityscapes** val set (Indicator: $mIoU$). Images are cropped to 801×801.

Method	1/16(186)	1/8(372)	1/4(744)	1/2(1488)
Supervised	64.39	72.14	74.83	77.93
CPS [23]	69.78	74.31	74.58	76.81
AEL [70]	74.45	75.55	77.48	79.01
PS-MT [26]	—	76.90	77.60	79.10
U2PL [27]	70.30	73.37	76.47	79.05
CCVC [29]	74.90	76.40	77.30	—
AugSeg [15]	75.22	77.82	79.56	80.43
DGCL [54]	73.18	77.29	78.48	80.71
UniMatch [13]	76.60	77.90	79.20	79.50
ESL [67]	75.12	77.15	78.93	80.46
CFCG [71]	77.28	79.09	80.07	80.59
Allspark [33]	78.33	79.24	80.56	81.39
w/ ASAug	78.91	79.68	81.06	81.95
↑ Δ (%)	↑ 0.58	↑ 0.44	↑ 0.50	↑ 0.56
CorrMatch [32]	77.30	78.50	79.40	80.40
w/ ASAug	77.97	79.22	79.86	80.98
↑ Δ (%)	↑ 0.67	↑ 0.72	↑ 0.46	↑ 0.58

TABLE VI: Comparison with SOTAs on the **COCO** val set (Indicator: $mIoU$). Images are cropped as 513×513. Corrmatch [32] did not report the COCO results.

Method	1/512 (232)	1/256 (463)	1/128 (925)	1/64 (1849)	1/32 (3697)
Supervised	22.94	27.96	33.60	37.80	42.24
PseudoSeg [53]	29.78	37.11	39.11	41.75	43.64
PC2Seg [55]	29.94	37.53	40.12	43.67	46.05
MKD [72]	30.24	38.04	42.32	45.50	47.25
UniMatch [13]	31.90	38.90	44.40	48.20	49.80
LogicDiag [73]	33.07	40.28	45.35	48.83	50.51
S4Former [74]	35.20	43.10	46.90	—	—
BRPG [75]	—	—	41.73	45.91	50.55
Allspark [33]	34.10	41.65	45.48	49.56	51.72
w/ASAug	35.78	43.28	48.42	52.31	55.25
↑ Δ (%)	↑ 1.68	↑ 1.63	↑ 2.94	↑ 2.75	↑ 3.53

[~] denotes the results we reproduced, which are missing in the paper.

dataset features 21 semantic categories and is divided into *classic* and *blender* subsets. The *classic* subset consists of 1,464 images with extensive labels for training and 1,449 images for validation. Conversely, the *blender* subset, as described in [24], [27], includes lower-resolution, roughly annotated images from the Segmentation Boundary Dataset (SBD) [69], thus enlarging the training pool to 10,582 images. We also apply the same settings as U2PL [27] to evaluate ASAug on the *blender* subset. The Cityscapes dataset, which encompasses 19 semantic categories within urban contexts, offers 2,975 training images with precise annotations and 500 images for validation. The COCO dataset contains 80 categories and provides 118,000 training images and 5,000 validation images covering indoor and outdoor scenes.

Parameter Settings. Following the conventions of most prior SSSS approaches [13], [15], [24], we evaluated our ASAug technique on both traditional CNNs and ViT structures. Specifically, for the CNN-based CorrMatch [32], we utilized DeepLabV3+ [76] integrated with ResNet-101 [77], pre-trained on ImageNet [78], and for the ViT-based Allspark [33], we employed SegFormer-B5 [79]. To evaluate ASAug’s performance, SGD was used as the optimizer with a polynomial decay learning rate policy: $lr_{init} \cdot \left(1 - \frac{iter}{total_iters}\right)^{power}$,

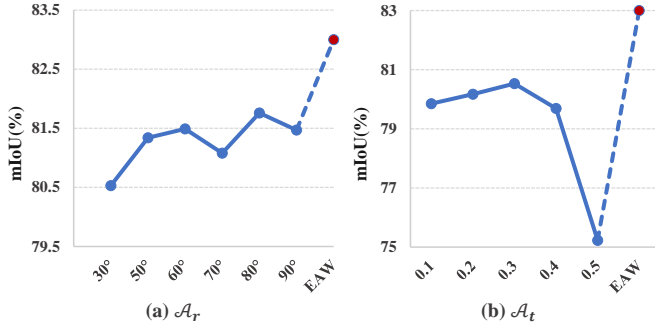


Fig. 4: Compare EAW with direct spatial augmentations. (a) EAW vs. fixed rotation angle, (b) EAW vs. same translation ratio (Based on Allspark [33]).

where lr_{init} is the starting learning rate, $iter$ is the current iteration number, $total_iters$ is the total number of iterations, with the power and weight decay set at 0.9 and $1e - 4$, accordingly. Loss weight λ is typically set to 0.5. For PASCAL, we began with a learning rate of 0.001, utilized crop sizes of 321×321 or 513×513 , had a batch size of 16, and trained for 80 epochs. For Cityscapes, the starting learning rate was 0.005, employed a crop size of 801×801 , used the same batch size, and trained for 240 epochs. All tests utilized the PyTorch deep learning framework and were executed on $4 \times$ NVIDIA V100 and $2 \times$ A800 GPUs for distinct numerical analyses. For COCO, we use the same initial learning rate and batch size as PASCAL, utilized crop sizes of 513×513 , and train for 10 epochs. Regarding data augmentation, the parameters were set to $r_{max} = 180$, $t_{max} = 0.5$. For PASCAL, parameters were $d_r = 1$, $d_t = 1$, $k_r = 11$, $k_t = 7$. Meanwhile, for Cityscapes, some parameters were adjusted to $d_r = 0.5$, $d_t = 0.5$, $k_r = 5.5$, and $k_t = 3$ based on the specific context.

Evaluations. In this experiment, we employ the mean intersection over union ($mIoU$) as our evaluation metric, which is a typical evaluation standard in SSSS. This metric is effective even when dealing with imbalanced classes, a frequent issue in pixel-level annotation tasks. Following the methodologies of CPS [23] and U2PL [27], we also utilize a sliding evaluation approach to assess model performance on the Cityscapes [35] dataset.

B. Comparison with SOTA Methods

Since our ASAug can be used as a pluggable module, we show its effectiveness by replacing the original strong augmentation with our method. In particular, we present the results of ASAug when combined with the highly competitive Allspark [33] and CorrMatch [32], comparing its performance to leading methods in both datasets using the diverse partition strategies.

Pascal VOC 2012 classic set. Table II reports the results of applying ASAug advancements to Allspark [33] and CorrMatch [32], along with comparisons with other leading methods in two dimension settings (321 and 513). The inclusion of ASAug consistently enhances the performance of existing techniques by a large margin. Besides, all methods show improved results with a larger crop size. With the support

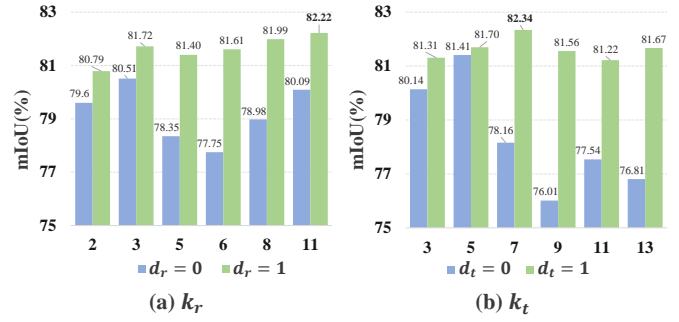


Fig. 5: Ablation study on EAW hyper-parameters k_r and k_t trained using the partitions of 1464, $d_t = d_r = 1.0$ (Based on Allspark [33]).

of ASAug, Allspark experiences an increase of $mIoU$ of 1.2% with a crop size of 513×513 and an increase of 4.16% with a crop size of 321×321 when trained on 92 labeled samples. Additionally, it can be seen that ASAug performs effectively in various scenarios, ranging from 1/16 partially labeled to fully labeled cases. This highlights ASAug’s effectiveness and establishes it as a versatile and valuable approach.

Pascal VOC 2012 blender set. Table III shows the numerical analysis for the larger *blender* set. Both Allspark [33] and CorrMatch [32] have demonstrated better performance with ASAug, which further verified the generalization of ASAug in all scenarios. At a 1/8 scale, Allspark’s performance sees a 2.04% improvement with ASAug, and CorrMatch results are similarly improved. An important finding is that our method significantly enhances performance, with results using a 321×321 crop size even outperforming some larger-scale 513×513 methods. Furthermore, we provide results within the same segmentation framework as U2PL [27], including fully annotated labels. Our method achieves a new state-of-the-art performance using Allspark with ASAug. These results verified the effectiveness and generalization of our ASAug.

Cityscapes. Table V evaluates the effectiveness of ASAug on the more challenging Cityscapes dataset, comparing it with other approaches. The evaluation uses a sliding window technique, showing that ASAug consistently delivers stable performance across different segments, even in complex urban scenes. We can easily see that ASAug can readily outperform other methods, especially with scarce labels. ASAug can improve Allspark by 4.16%, 2.45%, 2.85%, 1.55%, and 2.38% under 1/16, 1/8, 1/4, 1/2, and full partition protocols, respectively. This impressive improvement further demonstrates the effectiveness and importance of our claim that spatial data augmentation is effective and should be adopted for SSSS.

COCO. Table VI shows the performance of ASAug on the COCO dataset, and the results demonstrate that ASAug outperforms Allspark in all partitioning cases (1/512, 1/256, 1/128, 1/64, and 1/32), and in particular, in the case of 1/32, 1/64, and 1/128 partitions respectively, it improves by 3.53%, 2.75% and 2.94%, these experimental results further validate the effectiveness of spatial transformations on the multi-category datasets.

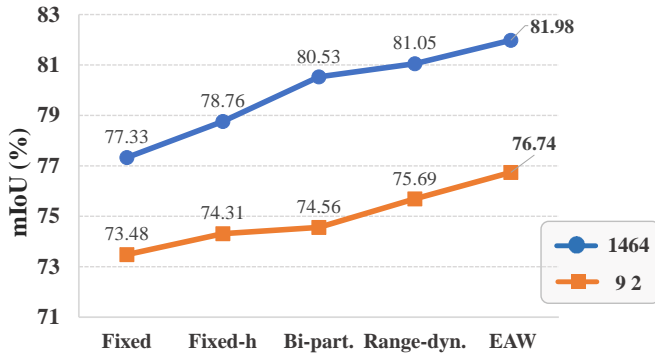


Fig. 6: Ablation study on mapping strategies. “EAW” denotes our entropy-based adaptive dynamic mapping.

C. Ablation Studies

We undertake multiple ablation studies to confirm the structure of the improvement approach introduced in ASAug. For consistency across studies, we used both CNN-based CorrMatch [32] method and ViT-based Allspark [33] method frameworks for ablation experiments. We present our findings using the Pascal VOC 2012 *classic* dataset, with a training dimension of 321×321 .

Effectiveness of ASAug. We begin by assessing the impact of individual components in ASAug, which are comprehensively presented in Table IV. The baseline results are from CorrMatch [32], and it is worth noting the point that the data presented in Table IV are from our reproducible ablation experiments using CorrMatch. Our streamlined augmentations, denoted as \mathcal{A}_r and \mathcal{A}_t , significantly enhance performance. The full implementation, $\mathcal{T}_{ss}(\mathcal{A}_r, \mathcal{A}_t)$, exceeds the performance of either \mathcal{A}_r or \mathcal{A}_t when applied independently. This indicates that strong spatial augmentations \mathcal{A}_r and \mathcal{A}_t both serve effectively as regularizers, with their combination yielding even superior results.

Effectiveness of EAW. We verified the need to dynamically adjust the weight of the spatial increase according to the variance of the instance, and Fig. 4 illustrates the effectiveness of EAW compared to other spatial transformations. (a) presents the results using a constant rotation angle for all samples, while (b) depicts those with a fixed translation ratio. A noticeable variance in model performance is observed at different angles, generally improving as the angle increases. Similarly, the augmentation effect varies across different translation ratios. In both scenarios, EAW’s adaptive adjustment of augmentation distortion degree demonstrates superior results.

Effects of hyper-parameters. We further analyze the performance of ASAug for \mathcal{A}_r and \mathcal{A}_t with various scaling parameters k_r, k_t and validate the importance of the offset parameters d_r, d_t , as displayed in Fig. 5. The findings demonstrate that incorporating d_r and d_t generally enhances model performance across different hyperparameter settings, notably in higher k_r and medium k_t scenarios (82.22% and 82.34% for $k_r = 11$ and $k_t = 7$, respectively). Fig. 5 also illustrates that no single parameter combination is consistently superior to others.

We assess the effect of the mapping method in Fig. 6, focusing solely on rotation mapping \mathcal{A}_r . Here, “EAW” stands for our entropy-driven adaptive weight. “Fixed” and “Fixed-h”

TABLE VII: Ablation study of the efficiency. We give a comparison of running time and results before and after applying ASAug, based on CorrMatch [32]. (Indicator: $mIoU$).

	1 / 16 (662)	1 / 8 (1323)	1 / 4 (2646)	1 / 2 (5291)
CorrMatch	32min25s	30min06s	25min48s	17min25s
w/ ASAug	39min02s	36min02s	30min31s	21min05s
Time diff.	+ 6min37s	+ 5min56s	+ 4min43s	+ 3min40s
$\uparrow \Delta$ (%)	\uparrow 1.11	\uparrow 1.13	\uparrow 0.77	\uparrow 0.55

refer to spatial rotation transformations applied to images for a standard fixed-angle view (30°) and a higher fixed angle (80°), respectively. “Bi-part.” signifies a two-part mapping (chosen between Fixed or Fixed-h based on the entropy value), while “Range-dyn.” involves splitting the angular range into three parts (low, medium, and high) based on experimental results, then conducting dynamic partition mapping based on entropy values. Fig. 6 shows that rotation angles significantly affect the effectiveness of ASAug. Our EAW strategy surpasses the other typical mappings, highlighting its notable robustness in handling complex features and limited samples.

Efficiency. The ASAug serves as a plug-and-play data augmentation tool that improves feature learning without requiring modifications to the existing training framework. Although it marginally extends the duration of the training, it significantly enhances performance and leaves the inference process unchanged. As demonstrated in the table below for the Pascal blender dataset (both involving 1085 iterations), our approach effectively balances enhanced performance with practical efficiency. (Indicator: $mIoU$)

D. Qualitative Analysis

Visualization on the Pascal dataset. We visually compare the ASAug results with Allspark [33] and CorrMatch [32] for segmenting the Pascal validation datasets in Fig. 7. The ASAug technique shows a higher level of segmentation accuracy, regularly surpassing the performances of both Allspark and CorrMatch across various complex conditions. It offers enhanced capabilities in identifying object classes and boundaries, preserving detailed features, and reducing incorrect background predictions. For instance, the outlines of aircraft ((a)-Row 1), animals ((a)-Rows 2/4 and (b)-Rows

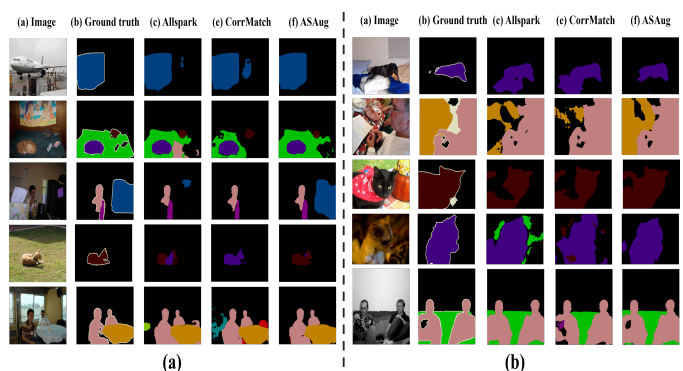


Fig. 7: Visualization of the segmentation results on Pascal validation set, compared with Allspark [33] and CorrMatch [32].

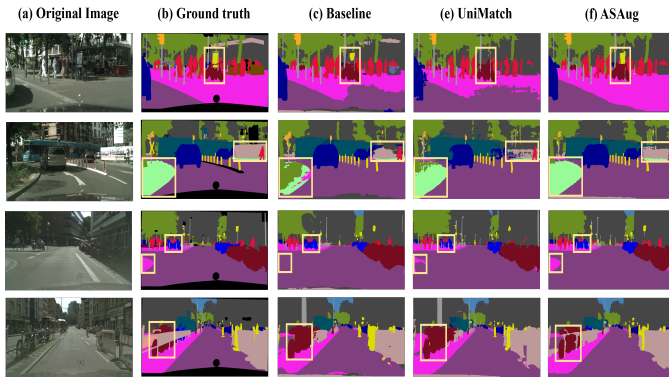


Fig. 8: Visualization of sample predictions from Cityscapes, highlighting the improved segmentation accuracy achieved with ASAug.

1/3/4), and people ((a)-Rows 3/5 and (b)-Row 5) are more precise and align closely with the actual situation. Additionally, ASAug reduces unnecessary background inclusions in the output, noticeable with the cat in the grass ((a)-Row 4) and the sleeping dog ((a)-Row 5). In scenes with multiple categories ((a)-Rows 2/5 and (b) Row 7), ASAug provides a more accurate segmentation of individual characters, whereas other methods generate some degree of noisy predictions, offering clearer outputs. These findings underscore the efficacy of ASAug’s spatial improvements and consistency-driven regularization strategy, which enhance ASAug’s performance across straightforward and challenging segmentation tasks.

Visualization on the Cityscapes dataset. Fig. 8 presents visualization results on the Cityscapes dataset trained with 744 labeled samples, to illustrate the enhanced performance of ASAug. As our baseline, we utilized the original self-training model, which was trained solely with labeled data. The baseline results in some categories having vague edges, leading to lost details, such as in the boundary of lawns, trees, and buildings, with evident misclassifications and reduced detection accuracy. While UniMatch [13] somewhat enhances edge quality, it still struggles with misclassification in intricate regions like those with pedestrians and road objects. In contrast, ASAug’s outcomes are more aligned with the ground truth, particularly achieving higher segmentation accuracy for small objects like pedestrians and vehicles, showcasing its superior capability in detailing regions. This advancement may be attributed to ASAug’s unique focus on positional information, which enables it to adapt to different scene variations, thereby improving detail discrimination while maintaining overall consistency, allowing ASAug to outperform other methods in terms of edge preservation and complex regions, validating its effectiveness in improving model generalization ability and detail fidelity. Nonetheless, segmentation in complex category stacking remains challenging. We anticipate future improvements to ASAug for such scenarios.

V. CONCLUSION

In this paper, we present ASAug, a simple yet effective method designed to address the challenges of SSSS. Unlike recent approaches that primarily rely on intensity-based

augmentations, our findings demonstrate that strong spatial augmentations can significantly improve SSSS performance. ASAug utilizes translation and rotation as data augmentation techniques, promoting better generalization in the presence of dynamic disturbances. To ensure stability, we incorporate an entropy-based adaptive weight strategy to handle more challenging cases. ASAug achieves substantial improvements over state-of-the-art, without the need for complex architectures.

We believe that shifting the focus from intensity-based augmentations to spatial techniques, such as translation and rotation, will lay a strong foundation for future SSSS research. Additionally, combining strong spatial augmentation with other augmentation strategies may show great potential.

REFERENCES

- [1] Q. Ma, Z. Zhang, P. Qiao, Y. Wang, R. Ji, C. Liu, and J. Chen, “Dual-level masked semantic inference for semi-supervised semantic segmentation,” *TMM*, 2025.
- [2] S. Zhang, D. Kong, Y. Xing, Y. Lu, L. Ran, G. Liang, H. Wang, and Y. Zhang, “Frequency-guided spatial adaptation for camouflaged object detection,” *TMM*, vol. 27, pp. 72–83, 2025.
- [3] Y. Ding, L. Li, W. Wang, and Y. Yang, “Clustering propagation for universal medical image segmentation,” in *CVPR*, 2024.
- [4] L. Ran, L. Wang, T. Zhuo, Y. Xing, and Y. Zhang, “DDF: A novel dual-domain image fusion strategy for remote sensing image semantic segmentation with unsupervised domain adaptation,” *TGRS*, 2024.
- [5] D. Guo, D.-P. Fan, T. Lu, C. Sakaridis, and L. Van Gool, “Vanishing-point-guided video semantic segmentation of driving scenes,” in *CVPR*, 2024.
- [6] B. Chen, Z. Ye, Y. Liu, X. Fang, G. Lu, S. Xie, and X. Li, “Towards robust semi-supervised distribution alignment against label distribution shift with noisy annotations,” *TMM*, 2025.
- [7] X. Lu, L. Li, L. Jiao, X. Liu, F. Liu, W. Ma, and S. Yang, “Uncertainty-aware semi-supervised learning segmentation for remote sensing images,” *TMM*, 2025.
- [8] L. Ran, Y. Li, G. Liang, and Y. Zhang, “Pseudo labeling methods for semi-supervised semantic segmentation: A review and future perspectives,” *TCSVT*, vol. 35, no. 4, pp. 3054–3080, 2025.
- [9] K. Hu, X. Chen, Z. Chen, Y. Zhang, and X. Gao, “Multi-perspective pseudo-label generation and confidence-weighted training for semi-supervised semantic segmentation,” *TMM*, 2024.
- [10] H. Xiao, Y. Hong, L. Dong, D. Yan, J. Xiong, J. Zhuang, D. Liang, and C. Peng, “Multi-level label correction by distilling proximate patterns for semi-supervised semantic segmentation,” *TMM*, 2024.
- [11] Y. Zou, Z. Zhang, H. Zhang, C.-L. Li, X. Bian, J.-B. Huang, and T. Pfister, “PseudoSeg: Designing pseudo labels for semantic segmentation,” in *ICLR*, 2021.
- [12] J. Lee, E. Kim, and S. Yoon, “Anti-adversarially manipulated attributions for weakly and semi-supervised semantic segmentation,” in *CVPR*, 2021.
- [13] L. Yang, L. Qi, L. Feng, W. Zhang, and Y. Shi, “Revisiting weak-to-strong consistency in semi-supervised semantic segmentation,” in *CVPR*, 2023.
- [14] J. Yuan, Y. Liu, C. Shen, Z. Wang, and H. Li, “A simple baseline for semi-supervised semantic segmentation with strong data augmentation,” in *ICCV*, 2021.
- [15] Z. Zhao, L. Yang, S. Long, J. Pi, L. Zhou, and J. Wang, “Augmentation matters: A simple-yet-effective approach to semi-supervised semantic segmentation,” in *CVPR*, 2023.
- [16] E. D. Cubuk, B. Zoph, D. Mane, V. Vasudevan, and Q. V. Le, “AutoAugment: Learning augmentation strategies from data,” in *CVPR*, Jun. 2019.
- [17] E. D. Cubuk, B. Zoph, J. Shlens, and Q. V. Le, “RandAugment: Practical automated data augmentation with a reduced search space,” in *CVPRW*, 2020.
- [18] S. G. Müller and F. Hutter, “TrivialAugment: Tuning-free yet state-of-the-art data augmentation,” in *ICCV*, 2021.
- [19] Y. Ouali, C. Hudelot, and M. Tami, “Semi-supervised semantic segmentation with cross-consistency training,” in *CVPR*, 2020.
- [20] Y. Zhou, H. Xu, W. Zhang, B. Gao, and P.-A. Heng, “C3-semiseg: Contrastive semi-supervised segmentation via cross-set learning and dynamic class-balancing,” in *ICCV*, 2021.
- [21] Z. Feng, Q. Zhou, Q. Gu, X. Tan, G. Cheng, X. Lu, J. Shi, and L. Ma, “Dmt: Dynamic mutual training for semi-supervised learning,” *PR*, 2022.

- [22] S. Liu, S. Zhi, E. Johns, and A. J. Davison, “Bootstrapping semantic segmentation with regional contrast,” *ICLR*, 2022.
- [23] X. Chen, Y. Yuan, G. Zeng, and J. Wang, “Semi-supervised semantic segmentation with cross pseudo supervision,” in *CVPR*, 2021.
- [24] L. Yang, W. Zhuo, L. Qi, Y. Shi, and Y. Gao, “St++: Make self-training work better for semi-supervised semantic segmentation,” in *CVPR*, 2022.
- [25] D. Kwon and S. Kwak, “Semi-supervised semantic segmentation with error localization network,” in *CVPR*, 2022.
- [26] Y. Liu, Y. Tian, Y. Chen, F. Liu, V. Belagiannis, and G. Carneiro, “Perturbed and strict mean teachers for semi-supervised semantic segmentation,” in *CVPR*, 2022.
- [27] Y. Wang, H. Wang, Y. Shen, J. Fei, W. Li, G. Jin, L. Wu, R. Zhao, and X. Le, “Semi-supervised semantic segmentation using unreliable pseudo-labels,” in *CVPR*, 2022.
- [28] J. Fan, B. Gao, H. Jin, and L. Jiang, “Ucc: Uncertainty guided cross-head co-training for semi-supervised semantic segmentation,” in *CVPR*, 2022.
- [29] Z. Wang, Z. Zhao, X. Xing, D. Xu, X. Kong, and L. Zhou, “Conflict-based cross-view consistency for semi-supervised semantic segmentation,” in *CVPR*, 2023.
- [30] Z. Zhao, S. Long, J. Pi, J. Wang, and L. Zhou, “Instance-specific and model-adaptive supervision for semi-supervised semantic segmentation,” in *CVPR*, 2023.
- [31] L. Wu, L. Fang, X. He, M. He, J. Ma, and Z. Zhong, “Querying labeled for unlabeled: Cross-image semantic consistency guided semi-supervised semantic segmentation,” *TPAMI*, 2023.
- [32] B. Sun, Y. Yang, L. Zhang, M.-M. Cheng, and Q. Hou, “Corrmatch: Label propagation via correlation matching for semi-supervised semantic segmentation,” in *CVPR*, 2024.
- [33] H. Wang, Q. Zhang, Y. Li, and X. Li, “Allspark: Reborn labeled features from unlabeled in transformer for semi-supervised semantic segmentation,” in *CVPR*, 2024.
- [34] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, “The pascal visual object classes (voc) challenge,” *IJCV*, 2010.
- [35] 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 *CVPR*, 2016.
- [36] 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 *ECCV*. Springer, 2014, pp. 740–755.
- [37] Y. Xu, L. Shang, J. Ye, Q. Qian, Y.-F. Li, B. Sun, H. Li, and R. Jin, “Dash: Semi-supervised learning with dynamic thresholding,” in *ICLR*, 2021.
- [38] D.-H. Lee *et al.*, “Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks,” in *Workshop on challenges in representation learning, ICML*, 2013.
- [39] L. Ran, D. Wen, T. Zhuo, S. Zhang, X. Zhang, and Y. Zhang, “Adasemicd: An adaptive semi-supervised change detection method based on pseudo-label evaluation,” *TGRS*, vol. 63, pp. 1–14, 2025.
- [40] K. Sohn, D. Berthelot, N. Carlini, Z. Zhang, H. Zhang, C. A. Raffel, E. D. Cubuk, A. Kurakin, and C.-L. Li, “Fixmatch: Simplifying semi-supervised learning with consistency and confidence,” *NeurIPS*, 2020.
- [41] B. Zhang, Y. Wang, W. Hou, H. Wu, J. Wang, M. Okumura, and T. Shinozaki, “Flexmatch: Boosting semi-supervised learning with curriculum pseudo labeling,” in *NeurIPS*, 2021.
- [42] A. Peláez-Vegas, P. Mesejo, and J. Luengo, “A survey on semi-supervised semantic segmentation,” *arXiv:2302.09899*, 2023.
- [43] P. Bachman, O. Alsharif, and D. Precup, “Learning with pseudo-ensembles,” *NeurIPS*, 2014.
- [44] A. v. d. Oord, Y. Li, and O. Vinyals, “Representation learning with contrastive predictive coding,” *arXiv:1807.03748*, 2018.
- [45] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial nets,” *NeurIPS*, 2014.
- [46] J. Hou, X. Ding, and J. D. Deng, “Semi-supervised semantic segmentation of vessel images using leaking perturbations,” in *WACV*, 2022.
- [47] Y. Xu, F. He, B. Du, D. Tao, and L. Zhang, “Self-ensembling gan for cross-domain semantic segmentation,” *TMM*, vol. 25, pp. 7837–7850, 2022.
- [48] N. Souly, C. Spampinato, and M. Shah, “Semi supervised semantic segmentation using generative adversarial network,” in *ICCV*, 2017.
- [49] D. Li, J. Yang, K. Kreis, A. Torralba, and S. Fidler, “Semantic segmentation with generative models: Semi-supervised learning and strong out-of-domain generalization,” in *CVPR*, 2021.
- [50] C. Cao, T. Lin, D. He, F. Li, H. Yue, J. Yang, and E. Ding, “Adversarial dual-student with differentiable spatial warping for semi-supervised semantic segmentation,” *TCSVT*, 2022.
- [51] G. Jin, C. Liu, and X. Chen, “Adversarial network integrating dual attention and sparse representation for semi-supervised semantic segmentation,” *Inf. Process. Manag.*, 2021.
- [52] H. Xie, C. Wang, M. Zheng, M. Dong, S. You, C. Fu, and C. Xu, “Boosting semi-supervised semantic segmentation with probabilistic representations,” in *AAAI*, 2023.
- [53] Y. Zou, Z. Zhang, H. Zhang, C.-L. Li, X. Bian, J.-B. Huang, and T. Pfister, “PseudoSeg: Designing pseudo labels for semantic segmentation,” in *ICLR*, 2021.
- [54] X. Wang, B. Zhang, L. Yu, and J. Xiao, “Hunting sparsity: Density-guided contrastive learning for semi-supervised semantic segmentation,” in *CVPR*, 2023.
- [55] Y. Zhong, B. Yuan, H. Wu, Z. Yuan, J. Peng, and Y.-X. Wang, “Pixel contrastive-consistent semi-supervised semantic segmentation,” in *ICCV*, 2021.
- [56] Y. Ouali, C. Hudelot, and M. Tami, “An overview of deep semi-supervised learning,” *arXiv:2006.05278*, 2020.
- [57] 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 *ICCV*, 2019.
- [58] V. Olsson, W. Tranheden, J. Pinto, and L. Svensson, “Classmix: Segmentation-based data augmentation for semi-supervised learning,” in *WACV*, 2021.
- [59] Y. Chen, X. Ouyang, K. Zhu, and G. Agam, “Complexmix: Semi-supervised semantic segmentation via mask-based data augmentation,” in *ICIP*, 2021.
- [60] I. Grubišić, M. Oršić, and S. Šegvić, “Revisiting consistency for semi-supervised semantic segmentation,” *Sensors*, 2023.
- [61] ———, “A baseline for semi-supervised learning of efficient semantic segmentation models,” in *MVA*, 2021.
- [62] A. Tarvainen and H. Valpola, “Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results,” *NeurIPS*, 2017.
- [63] X. Lai, Z. Tian, L. Jiang, S. Liu, H. Zhao, L. Wang, and J. Jia, “Semi-supervised semantic segmentation with directional context-aware consistency,” in *CVPR*, 2021.
- [64] H. Kong, G.-H. Lee, S. Kim, and S.-W. Lee, “Pruning-guided curriculum learning for semi-supervised semantic segmentation,” in *WACV*, 2023.
- [65] G. French, S. Laine, T. Aila, M. Mackiewicz, and G. Finlayson, “Semi-supervised semantic segmentation needs strong, varied perturbations,” *arXiv:1906.01916*, 2019.
- [66] Y. Jin, J. Wang, and D. Lin, “Semi-supervised semantic segmentation via gentle teaching assistant,” *NeurIPS*, 2022.
- [67] J. Ma, C. Wang, Y. Liu, L. Lin, and G. Li, “Enhanced soft label for semi-supervised semantic segmentation,” in *ICCV*, 2023.
- [68] M. Everingham, S. A. Eslami, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, “The pascal visual object classes challenge: A retrospective,” *IJCV*, 2015.
- [69] B. Hariharan, P. Arbeláez, L. Bourdev, S. Maji, and J. Malik, “Semantic contours from inverse detectors,” in *ICCV*, 2011.
- [70] H. Hu, F. Wei, H. Hu, Q. Ye, J. Cui, and L. Wang, “Semi-supervised semantic segmentation via adaptive equalization learning,” *NeurIPS*, 2021.
- [71] S. Li, Y. He, W. Zhang, W. Zhang, X. Tan, J. Han, E. Ding, and J. Wang, “Cfcg: Semi-supervised semantic segmentation via cross-fusion and contour guidance supervision,” in *ICCV*, 2023.
- [72] J. Yuan, J. Ge, Z. Wang, and Y. Liu, “Semi-supervised semantic segmentation with mutual knowledge distillation,” in *ACM MM*, 2023.
- [73] C. Liang, W. Wang, J. Miao, and Y. Yang, “Logic-induced diagnostic reasoning for semi-supervised semantic segmentation,” in *ICCV*, 2023, pp. 16 197–16 208.
- [74] X. Hu, L. Jiang, and B. Schiele, “Training vision transformers for semi-supervised semantic segmentation,” in *CVPR*, 2024.
- [75] J. Dong, Z. Meng, D. Liu, J. Liu, Z. Zhao, and F. Su, “Boundary-refined prototype generation: A general end-to-end paradigm for semi-supervised semantic segmentation,” *EAAI*, 2024.
- [76] L.-C. Chen, G. Papandreou, F. Schroff, and H. Adam, “Rethinking atrous convolution for semantic image segmentation,” *arXiv:1706.05587*, 2017.
- [77] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *CVPR*, 2016.
- [78] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *CVPR*, 2009.
- [79] E. Xie, W. Wang, Z. Yu, A. Anandkumar, J. M. Alvarez, and P. Luo, “Segformer: Simple and efficient design for semantic segmentation with transformers,” *NeurIPS*, 2021.