# Enhancing Sample Utilization through Sample Adaptive Augmentation in Semi-Supervised Learning

Guan Gui[1], Zhen Zhao[2], Lei Qi[3], Luping Zhou[2], Lei Wang[4], Yinghuan Shi[1*]

[1]Nanjing University   [2]University of Sydney   [3]Southeast University   [4]University of Wollongong

guiguan@smail.nju.edu.cn, {zhen.zhao, luping.zhou}@sydney.edu.au

qilei@seu.edu.cn, leiw@uow.edu.au, syh@nju.edu.cn

## Abstract

*In semi-supervised learning, unlabeled samples can be utilized through augmentation and consistency regularization. However, we observed certain samples, even undergoing strong augmentation, are still correctly classified with high confidence, resulting in a loss close to zero. It indicates that these samples have been already learned well and do not provide any additional optimization benefits to the model. We refer to these samples as "naive samples". Unfortunately, existing SSL models overlook the characteristics of naive samples, and they just apply the same learning strategy to all samples. To further optimize the SSL model, we emphasize the importance of giving attention to naive samples and augmenting them in a more diverse manner. Sample adaptive augmentation (SAA) is proposed for this stated purpose and consists of two modules: 1) sample selection module; 2) sample augmentation module. Specifically, the sample selection module picks out naive samples based on historical training information at each epoch, then the naive samples will be augmented in a more diverse manner in the sample augmentation module. Thanks to the extreme ease of implementation of the above modules, SAA is advantageous for being simple and lightweight. We add SAA on top of FixMatch and FlexMatch respectively, and experiments demonstrate SAA can significantly improve the models. For example, SAA helped improve the accuracy of FixMatch from 92.50% to 94.76% and that of FlexMatch from 95.01% to 95.31% on CIFAR-10 with 40 labels. The code is available at https://github.com/GuanGui-nju/SAA.*

(a) An example of *naive sample*.



(b) Model performance during training.

Figure 1: (a) shows an example of *naive sample*. Its augmented versions are correctly classified with high confidence, resulting in the loss close to 0. (b) shows the model performance during FixMatch training. Performance improvements are slow or even stagnant for a period of time.

## 1. Introduction

For the sake of reducing the cost of manual labeling, semi-supervised learning (SSL), which focuses on how to learn from unlabeled data, is a longstanding yet significant research topic in vision applications. Recently, data augmentation techniques and consistency regularization have been proven to be effective ways of utilizing unlabeled data. For example, FixMatch [39] encourages consistency in predictions between the weakly and strongly augmented versions, and it achieves an accuracy of 92.50% on the CIFAR-10 task with only 40 labels.

However, not all unlabeled samples are effectively utilized even with strong augmentation. In Figure 1a, if the strongly augmented versions are correctly classified with

1

high confidence, leading to a loss close to zero, it indicates that the sample has already been learned well and cannot further improve the model's performance. In other words, the sample was not effectively utilized to benefit model training, and we call this sample "*naive sample*". When the training process contains a large number of *naive samples*, it can cause slow or even stagnant model performance improvements, as shown in Figure 1b.

Unfortunately, existing SSL models [34] overlook the critical point of whether all samples are effectively utilized. Typically, these models apply the same fixed strong augmentation strategy to all samples, resulting in some strongly augmented versions that do not benefit the model train.

We emphasize that the key to alleviating this problem lies in how to further explore the value of the *naive samples* through new learning strategies. A natural idea that reminds us is to develop sample adaptive augmentation (SAA) to identify *naive samples* and increase their diversity after augmentation. Our proposed SAA is simple yet effective, which consists of two modules: 1) sample selection module and 2) sample augmentation module. The former is responsible for picking out *naive samples* in each epoch, and the latter applies a more diverse augmentation strategy for *naive samples*. Specifically, in the sample selection module, we first update the historical loss of the samples with exponential moving average (EMA) in each epoch, then these samples will be divided into two parts. The part of the samples with a smaller historical loss is considered to be the *naive sample*. Since historical loss captures the impact of the sample on model training, this approach allows us to identify samples that are not effectively utilized and would benefit from more diverse augmentation. While in the sample augmentation module, the more diverse augmented version of *naive sample* will be obtained by regrouping multiple strong augmented images, and the remaining samples are applied with the original strong augmentation.

Our proposed SAA is simple to implement, requiring only a few lines of code to add our proposed modules to the FixMatch or FlexMatch in PyTorch. It is also lightweight in terms of memory and computation, *i.e.*, SAA only needs to add two additional vectors and update them in each epoch, making it an efficient solution for improving SSL models.

We extended FixMatch and FlexMatch with SAA and conducted experiments on SSL benchmarks. The results of the experiments demonstrate that SAA can significantly improve performance. In summary, our contribution can be summarized as follows:

- **We identify ineffectively utilized samples and emphasize that they should be given more attention.** Under the consistency regularization based on data augmentation, some strongly augmented versions are not beneficial to model training, which results in the values of these samples not being fully exploited and

makes the model performance slow to improve. We refer to them as *"naive sample"*, and emphasize that they should be learned with a new learning strategy.

- **We propose SAA to make better use of *naive sample*.** To increase the probability that the augmented versions can benefit the model training, a simple yet effective method, sample adaptive augmentation (SAA), is proposed for identifying the *naive samples* and augmenting them in a more diverse manner.

- **We verify the validity of SAA on SSL benchmarks.** Using FixMatch and FlexMatch as the base framework, we proved that our approach can achieve state-of-the-art performance. For example, on CIFAR-10 with 40 labels, SAA helps FixMatch improve its accuracy from 92.50% to 94.76%, and helps FlexMatch improve its accuracy from 95.01% to 95.31%.

## 2. Related Work

### 2.1. Semi-Supervised Learning

Consistency regularization (CR) [2] is the main way to exploit unlabeled data in semi-supervised learning (SSL). The conventional implementation is to perturb the samples and then encourage the model to maintain a consistent prediction. The manner of perturbation has been studied in a variety of ways, *e.g.*, stochastic augmentation and drop out [25, 35], feature perturbations [24], adversarial perturbations [30], model perturbations [42]. [4, 44, 3] apply mixup to blend the images, also a perturbation of the image. With strong augmentation technique [7, 10], FixMatch applies a consistency regularization between the weakly augmented and strongly augmented versions, allowing the model to learn a greater diversity of images over long iterations. This approach has greatly simplified the framework and has led to breakthroughs in semi-supervised learning milestones. As we have previously analyzed, the framework applies the same fixed augmentation strategy to all images, which results in the *naive samples* not being fully utilized.

Due to the superiority of the FixMatch framework, a large number of SSL works [52, 12, 48, 55, 16, 13, 54] are now based on it for further optimization, but none of the work considers the effectiveness of utilization of *naive samples*. [12, 54] focuses on improving the quality of the pseudo labels by learning the distribution of unlabeled data. [55, 16] focus on learning the similarity relationship between samples or super-classes. [13, 48, 52] all emphasize the utilization of samples with low confidence. These works also allow the model to learn more samples within a certain number of iterations to some extent, but still ignore the issue of the validity of the augmentation, resulting in these augmented samples still possibly unhelpful to the training of the mode.

To the best of our knowledge, none of the SSL works has considered the utilization of *naive samples*. [52] treats each category differently and adjusts the threshold for each category, but we are considering treating each sample differently, with no relationship to the category.

## 2.2. Hard Example Mining

Our work is somewhat related to hard example mining, with the difference that we focus on *naive samples* that do not benefit model training, while they focus on hard samples that damage model training. A more common approach used to select hard samples is to rely on loss information between the sample and the ground truth [37, 46, 29, 38]. This is related to our approach, but the unlabeled sample is selected by its consistency loss due to the lack of ground-truth. In addition to this, distance metrics [49, 22] and false positive samples [20, 9, 11, 14] are common methods of hard sample selection. However, most of these methods for mining hard samples rely on labels, which is not practical under SSL. Both our proposed method and the field involve the selection of samples, the difference being that they focus on the selection of hard samples that are difficult to train, while we focus on *naive sample* that contributes no information of model training.

## 2.3. Date Augmentation

Data augmentation is an effective way of expanding the data space [36], which we roughly classify into the following categories: 1) Single perturbation [31, 21, 10]; 2) image blending based [43, 53, 50, 19]; 3) learning based [40, 28, 15]; and 4) search based [6, 27, 7]. Common operations for perturbing a single image include geometric transformations, color transformations, noise injection [31], random erasing [56], kernel flters [21], cutout [10]. [43, 53] direct mixing of the contents of two images, [50] mixes image patches, and [8, 19] Mix the content and style. For learning based strategy, adversarial training [40, 47] and GAN-based [28, 15] train the network to obtain augmented images. [6, 27] find the best combination in the perturbation space using a search strategy.

However, single perturbation and image blending-based methods are limited to enhancing the diversity of images. For learning and search-based methods, although they yield augmented images that facilitate model training, their time consumption is huge and therefore this is not suitable for training time-consuming CR-based SSL models. [7] combines random transformations to remove the search process, and is favored by the SSL model. [17] cuts the image and augments the patches, which enhances image diversity and has been validated to be effective on several tasks. Our augmentation is related to [17], but we augmented on the image, not on patches. We will further discuss it in the experimental section.

## 3. Preliminary and Background

### 3.1. Problem Setting

In semi-supervised learning, we denote labeled set $\mathcal{X} = \{(x_1, y_1), (x_2, y_2), \ldots, (x_M, y_M)\}$, where $y_i$ is the label of the $i$-th labeled sample $x_i$. We also denote unlabeled set $\mathcal{U} = \{u_1, u_2, \ldots, u_N\}$, where $u_i$ denotes $i$-th unlabeled sample, and typically $|\mathcal{X}| \ll |\mathcal{U}|$. In the implementation, the samples are provided on a per batch basis in each iteration, with a batch of labeled data $\mathcal{X}$ and batches of unlabeled data $\mathcal{U}$. How to use this unlabeled data for learning is the focus of SSL research. Generally, in consistency regularization (CR) based SSL models, unlabeled data generate different versions by perturbation, and then the model is encouraged to be consistent in its representations or predictions of these versions.

### 3.2. Preliminary for CR-based SSL Models

Strong augmentation is a good means of applying consistency regularization, and FixMatch [39] is representative of this idea. Many recent semi-supervised works [52, 48, 26, 55, 16] have also used FixMatch as a basis for further optimization, and to clearly introduce our approach, we also use FixMatch as a base framework.

We first review FixMatch, whose fundamental idea is that produce pseudo labels on weakly-augmented versions and use them as training targets for their corresponding strongly augmented versions. Of them, the weak augmentation $\alpha(\cdot)$ includes standard flip and shift operations, while the strong augmentation strategy $\mathcal{A}(\cdot)$ consists of RandAugment [7] and CutOut [10].

Let $p_i^w$ and $p_i^s$ represent the model's prediction on $\alpha(u_i)$ and $\mathcal{A}(u_i)$, respectively. Then this consistency regularization based unsupervised loss for unlabeled samples is,

$$\mathcal{L}_{unsup} = \frac{1}{|\mathcal{U}|} \sum_{i=1}^{|\mathcal{U}|} \mathbb{1}(\max(p_i^w) \geq \tau_c) H(p_i^w, p_i^s). \quad (1)$$

where $H(p_1, p_2)$ denotes the standard cross entropy between $p_1$ and $p_2$, and $\tau_c$ is a pre-defined threshold to retain only high-confidence pseudo-labels. As discussed in FixMatch [39], $\tau_c$ is commonly set as a large value to alleviate the confirmation bias [1] in SSL. Let $p_i$ denotes the model's prediction of $\alpha(x_i)$, then the supervised loss for labeled samples is,

$$\mathcal{L}_{sup} = \frac{1}{|\mathcal{X}|} \sum_{i=1}^{|\mathcal{X}|} H(q_i, y_i). \quad (2)$$

Finally, the total losses can be expressed as,

$$\mathcal{L} = \mathcal{L}_{sup} + \lambda \mathcal{L}_{unsup}. \quad (3)$$

where $\lambda$ is the weight of $\mathcal{L}_{unsup}$.

(a) Sample selection module in SAA.



(b) Sample augmentation module in SAA.

Figure 2: Overview of our method SAA. The core insight of SAA is that dynamically adjusts augmentation for samples, thus allowing *naive samples* to be used more effectively. In detail, SAA consists of two modules: sample selection module and sample augmentation module. (a). Each sample $u_i$ corresponds to a marker $\mathcal{F}_i$ and historical loss $\mathcal{H}_i$. In each epoch, samples' consistency losses are recorded and their historical losses $\mathcal{H}_i$ are updated with EMA. Then based on the historical losses, we divide these samples into two parts by OTSU. The part of the samples with a smaller historical loss are *naive samples*, and their markers are set to 1, then the rest of the markers are set to 0. (b). Sample $u_i$ is augmented in different ways depending on the marker $\mathcal{F}_i$, *i.e.*, if $\mathcal{F}_i = 0$, it is strongly augmented once, if $\mathcal{F}_i = 1$, it is strongly augmented twice, and the two augmented images are regrouped into one image. The regrouping may be in two parts top-bottom or two parts left-right, which is chosen randomly with a probability of 0.5.

## 3.3. Characteristics and Impact of Naive Samples

Taking FixMatch as an example, we tracked the loss of a *naive sample* and a *non-naive sample* separately, as shown in Figure 3. Note that we show the original loss without thresholds to exclude the interference of confidence thresholds on the loss values. It can be observed that in most epochs, the *naive*



Figure 3: loss of *non-naive sample* (bottom) and *naive sample* (top).

*sample's* cross-entropy losses are close to 0, indicating that the learning of $\mathcal{A}(u_i)$ does not contribute to the model's training progress.

With the same augmentation, the *non-naive sample* can encourage model optimization in the long term, whereas *naive samples* cannot. This confirms the necessity for dedicated attention to *naive samples* and the development of new augmentation strategies that can better exploit their potential value. Moreover, when there are too many *naive samples* in the training process, they can interfere with model performance improvement, as shown in Figure 1b. These findings highlight the importance of properly identifying and handling *naive samples* in SSL tasks.

We would highlight that there are several factors that cause the slow performance improvement, such as the number of high confidence pseudo-labels, *etc*. There are also multiple ways to solve the problem, such as adjusting the threshold [52], finding other learnable signals [55], *etc*. In this work, we concentrate on data augmentation. It should be noted that our approach can be used together with the above ways and thus beneficial to SSL.

## 4. SAA: Using Sample Adaptive Augmentation to Aid Semi-Supervised Learning

SAA aims to address the issue of not effectively utilizing *naive samples* by providing them with more attention and exploration of their value in aiding model training. To achieve this goal, SAA designs two modules: the sample selection module and the sample augmentation module. The function of the first module is to identify *naive samples* in each epoch, while the function of the second module is to apply more diverse augmentation strategies to the *naive samples* to facilitate their effective learning.

### 4.1. Sample Selection

We introduce two vectors $\mathcal{H} = \{\mathcal{H}_1, \mathcal{H}_2, ..., \mathcal{H}_N\}, \mathcal{F} = \{\mathcal{F}_1, \mathcal{F}_2, ..., \mathcal{F}_N\}$, where $N$ is the number of unlabeled samples. $\mathcal{H}$ records historical consistency loss information for each unlabeled sample, and $\mathcal{F}$ marks whether the unlabeled sample is a *naive sample*. For unlabeled sample $u_i$, the model calculates the consistency loss $l_i^t$ between its weakly- augmented version and strongly augmented versions once in $t$-th epoch. Then we update $\mathcal{H}_i^t$ with exponential moving average (EMA), which can be expressed as:

$$\mathcal{H}_i^t = (1 - a)\mathcal{H}_i^{t-1} + \alpha l_i^t. \tag{4}$$

Note that the parameter $\alpha$ introduced is not an additional model parameter, as the model parameters are also updated with EMA [39, 52]. Since the historical loss information can reflect the magnitude of the impact of a strongly augmented version on the model, it becomes the basis for our decision on *naive sample*. OTSU [33] is a commonly used method of threshold segmentation because of its computa-

**Algorithm 1:** Equip FixMatch with SAA

---

**Input:** Labeled data batch $\mathcal{B}_x = \{(x_i, y_i)\}^M$, unlabeled data batch $\mathcal{B}_u = \{u_i\}^N$, unsupervised loss weight $\lambda$, pre-training epochs $T'$, total training epochs $T$, augmentation strategies $\alpha(\cdot), \mathcal{A}(\cdot), \mathcal{A}'(\cdot)$, historical loss $\mathcal{H} = \{\mathcal{H}_1, \mathcal{H}_2, ..., \mathcal{H}_N\}$, mark $\mathcal{F} = \{\mathcal{F}_1, \mathcal{F}_2, ..., \mathcal{F}_N\}$

1  **for** $t \leftarrow 1$ *to* $T'$ **do**
2  $\quad$ Run FixMatch; // `training as FixMatch in the first` $T'$ `epochs`
3  **end**
4  **for** $T' \leftarrow 1$ *to* $T$ **do**
5  $\quad$ Compute $\mathcal{L}_{sup} = 1/|\mathcal{B}_x| \sum_{i=1}^{|\mathcal{B}_x|} H(q_i, y_i)$; // `supervised loss in FixMatch`
6  $\quad$ **for** $i \leftarrow 1$ *to* $N$ **do**
7  $\quad\quad$ Apply perturbation $\alpha(\cdot)$ to $u_i$; // `weak augmentation in FixMatch`
8  $\quad\quad$ Apply augmentation $\mathcal{A}(\cdot)$ / $\mathcal{A}'(\cdot)$ to $u_i$ according to $\mathcal{F}_i$; // `sample augmentation in SAA`
9  $\quad\quad$ Compute loss $l_i = H(\arg\max p_m(\alpha(u_i)), p_m(\mathcal{A}(u_i)))$; // `consistency loss in FixMatch`
10 $\quad$ **end**
11 $\quad$ Compute $\mathcal{L}_{unsup} = 1/|\mathcal{B}_u| \sum_{i=1}^{|\mathcal{B}_u|} \mathbb{1}(\max p_m(\alpha(u_i)) \geq \tau)l_i$; // `unsupervised loss in FixMatch`
12 $\quad$ Update historical loss $\mathcal{H}_i^t = (1-a)\mathcal{H}_i^{t-1} + \alpha l_i^t$; // `sample selection in SAA`
13 $\quad$ Update mark $\mathcal{F}_i = \mathbb{1}(\mathcal{H}_i^t \leq \text{OTSU}(H_i^t))$; // `sample selection in SAA`
14 **end**

---

tional simplicity, stability, and strong self-adaptation. Inspired by this, we calculate the historical loss threshold in each epoch:

$$\tau_s = \text{OTSU}(\mathcal{H}_1, \mathcal{H}_2, ..., \mathcal{H}_N). \quad (5)$$

OTSU adaptively divides the sample into two parts based on the historical loss. Samples with small historical losses are considered as *naive samples* since they provide less help to the model. Then we update $\mathcal{F}$ by:

$$\mathcal{F}_i = \mathbb{1}(\mathcal{H}_i \leq \tau_s). \quad (6)$$

It can be seen that the decision of the *naive samples* is done at every epoch. In other words, whether a sample is a *naive sample* is related to the model performance. We should note that there may be multiple shifts in $\mathcal{F}$ in the training process. On the one hand, if the sample is regarded as a *naive sample*, we will apply a more diverse augmentation to it to avoid invalid learning. On the other hand, this more diverse augmentation may be too perturbing for the sample and negatively affect the model, so the augmentation strategy for these samples needs to be adjusted back to the original strategy in a timely manner.

### 4.2. Sample Augmentation

We apply different augmentations to the *non-naive sample* and the *naive sample*. The former will be applied with the original augmentation $\mathcal{A}$, while the latter will be applied with a new augmentation $\mathcal{A}'$, which increases the difficulty of the augmented version. This can be expressed as:

$$\text{Augmented}(u_i) = \begin{cases} \mathcal{A}(u_i), & \mathcal{F}_i = 1 \\ \mathcal{A}'(u_i), & \mathcal{F}_i = 0 \end{cases} \quad (7)$$

We implement this more diverse augmentation $\mathcal{A}'$ in a simple way, *i.e.*, by regrouping several $\mathcal{A}(u_i)$ into a new image. Formally, a new augmented image $\mathcal{A}'(u_i)$ can be expressed as:

$$\mathcal{A}'(u_i) = \text{Concat}(\text{Cut}(\mathcal{A}(u_i)_1), \text{Cut}(\mathcal{A}(u_i)_2)). \quad (8)$$

As shown in Figure 4, we have two strongly augmented images $\mathcal{A}(u_i)_1$ and $\mathcal{A}(u_i)_2$. To create a new augmented image $\mathcal{A}'(u_i)$, we randomly choose one of the following two options with equal probability: 1) Top-bottom concat: We take the top half of $\mathcal{A}(u_i)_1$ and the bottom half of $\mathcal{A}(u_i)_2$ and concatenate them to create a new image. 2) Left-right concat: We take the left half of $\mathcal{A}(u_i)_1$ and the right half of $\mathcal{A}(u_i)_2$ and concatenate them to create a new image.



Figure 4: Augmentation for *naive samples*.

Image regrouping is a technique to enhance image diversity by combining multiple augmented images into a new image. It is a simple and effective solution that has been shown to be effective in previous works such as "Cut-

| Method | CIFAR-10 | | | CIFAR-100 | | | SVHN | | | STL-10 |
| | 40 labels | 250 labels | 4000 labels | 400 labels | 2500 labels | 10000 labels | 40 labels | 250 labels | 1000 labels | 1000 labels |
|---|---|---|---|---|---|---|---|---|---|---|
| Mean-Teacher | 29.91±1.60 | 62.54±3.30 | 91.90±0.21 | 18.89±1.44 | 54.83±1.06 | 68.25±0.23 | 63.91±3.98 | 96.55±0.03 | 96.73±0.05 | - |
| MixMatch | 63.81±6.48 | 86.37±0.59 | 93.34±0.26 | 32.41±0.66 | 60.42±0.48 | 72.22±0.29 | 69.40±8.39 | 95.44±0.32 | 96.31±0.37 | 38.02±8.29 |
| ReMixMatch | 90.12±1.03 | 93.70±0.05 | 95.16±0.01 | 57.25±1.05 | 73.97±0.35 | **79.98±0.27** | 75.96±9.13 | 93.64±0.22 | 94.84±0.31 | 75.51±1.25 |
| Dash | 91.84±4.31 | 95.22±0.12 | 95.76±0.06 | 55.17±1.36 | 72.15±0.19 | 77.23±0.21 | 96.97±1.59 | 97.83±0.10 | 97.97±0.06 | 83.17±0.80 |
| CoMatch | 93.09±1.39 | 95.09±0.33 | 95.57±0.04 | - | - | - | - | - | - | 79.80±0.38 |
| SLA | 94.83±0.32 | 95.11±0.27 | 95.79±0.15 | 58.56±1.41 | 72.37±0.44 | 77.68±0.24 | 94.37±2.91 | 95.08±1.08 | 95.84±0.24 | - |
| NP-Match | 95.09±0.04 | 95.04±0.06 | 95.89±0.02 | 61.09±0.99 | 73.97±0.26 | 78.78±0.13 | - | - | - | - |
| SimMatch | 94.40±1.37 | 95.16±0.39 | 96.04±0.01 | **62.19±2.21** | 74.93±0.32 | 79.42±0.11 | - | - | - | 89.70±0.82 |
| FixMatch[†] | 92.50±0.67 | 95.10±0.04 | 95.81±0.05 | 53.17±0.51 | 72.64±0.17 | 77.60±0.09 | 96.24±0.98 | 97.54±0.04 | 97.98±0.02 | 85.27±1.15 |
| **FixMatch (w/SAA)** | 94.76±0.99 | 95.21±0.07 | 96.09±0.07 | 54.29±0.73 | 73.18±0.21 | 78.71±0.20 | **97.01±0.72** | **97.68±0.07** | **98.06±0.06** | 87.92±1.46 |
| FlexMatch[†] | 95.01±0.09 | 95.08±0.10 | 95.82±0.02 | 60.51±1.54 | 72.98±0.22 | 78.15±0.17 | 92.42±2.60 | 92.98±1.59 | 93.54±0.28 | 89.15±0.71 |
| **FlexMatch (w/SAA)** | **95.31±0.16** | **95.40±0.19** | **96.14±0.08** | 61.87±1.94 | **75.01±0.41** | 79.88±0.34 | 93.15±2.54 | 93.25±2.41 | 94.41±0.27 | **90.85±0.82** |
| Fully-supervised | | 95.38±0.05 | | | 81.70±0.09 | | | 97.87±0.02 | | - |

Table 1: Performance comparisons on CIFAR-10, CIFAR-100, SVHN, STL-10. We compare the performance with recent SSL works [42, 3, 4, 48, 26, 41, 55, 45]. We apply SAA on the top of FixMatch [39] and FlexMatch [52], respectively. For fair comparison, we re-ran FixMatch and FlexMatch under the exact same random seed, which is denoted by [†]. Fully-supervised comparisons follows FlexMatch [52], which is conducted with all labeled date with applying weak data augmentations. Experiments shows SAA provides a significant improvement to the SSL model. When we choose FlexMatch as the base framework, performance reaches SOTA for most settings.

Mix" [17]. In comparison to learning-based data augmentation methods, which can also yield augmented images suitable for model learning, image regrouping has lower memory and computational overheads. However, in the case of CutMix, augmentation is done on the cut images, which may result in some loss of information about the original image. In contrast, in our method, augmentation is applied to the whole image, which preserves more information about the original image. This is discussed further in the experimental section of the paper.

## 5. Experiments

We used FixMatch and FlexMatch as a base framework to verify the validity of SAA on SSL benchmark datasets: CIFAR-10, CIFAR-100 [23], SVHN [32] and STL-10 [5]. In section 5.1, we present the specific implementation details. In section 5.2, we first verify that SAA can help the model improve test accuracy and achieve SOTA on SSL tasks. In addition, we compare the performance for the same number of iterations and verify that SAA can accelerate the model's improvement speed.

### 5.1. Implementation Details

We adopt "WideResNet-28-2" [51] for CIFAR-10 and SVHN, "WideResNet-28-8" [51] for CIFAR-100 and "ResNet18" [18] for STL-10. For a fair comparison, we keep the same set of parameters as FixMatch and FlexMatch with $\{|\mathcal{B}_\mathcal{X}| = 64, |\mathcal{B}_\mathcal{U}| = 7|\mathcal{B}_\mathcal{X}|, \lambda = 1\}$. The test model is updated by EMA with a decay rate of 0.999. $\mathcal{H}$ is updated in the same way and with the same parameters ($\alpha = 0.999$). FixMatch and FlexMatch set the number of training iterations to $2^{20}$, and we keep this practice as well. In order to update the sample historical loss $\mathcal{H}$ and marker $\mathcal{F}$ in a

timely manner, we consider every 1024 iterations as one epoch, *i.e.*, a total of 1024 epochs are trained. As the augmentation $\mathcal{A}'$ is not suitable for the initial training of the model, we apply it only after the 100th epoch, while historical loss $\mathcal{H}$ is recorded from the beginning. We repeat the same experiment for five runs with different seeds to report the mean test accuracy and variance.

### 5.2. Main Results

**SAA improves the performance of baseline models.** As shown in Tabel 1, SAA successfully improves the test accuracy of FixMatch and FlexMatch under all settings. For instance, on CIFAR-10 with 40 labels, FixMatch and Flex-Match achieved a mean accuracy of 92.50% and 95.01%, while with SAA their average accuracy improved to 94.76% and 95.31%. For a challenging and realistic task STL-10, SAA helped FixMatch to improve its accuracy by 2.65% and FlexMatch by 1.70%. FixMatch and FlexMatch outperform even under fully-supervised in some settings, *e.g.*, Fix-Match achieved mean test accuracy of 95.81% on CIFAR-10 with 4000 labels and 97.98% on SVHN with 1000 labels. We can notice that the variance of the model becomes slightly larger after applying SAA, as the boosting effect of SAA on the model is different under different seeds. As shown in Figure 5a, SAA boosts FixMatch at all 5 seeds, but it can boost by 2.51% when the seed is 1 and 1.20% when the seed is 2.

**We achieve the SOTA performance.** With FixMatch as the base, SAA can help to bring its performance up to near or even beyond that of other SSL models. For example, on CIFAR-10 with 40 labels, FixMatch (w/SAA) achieves an accuracy of 94.76%, which is within 0.3% of NP-Match. While on CIFAR-10 with 250 labels, FixMatch (w/SAA)

(a) Accuracy under different seeds.

(b) Performance during training (with the same seed).

Figure 5: All experiments are conducted on task of CIFAR-10 with 40 labels. (a) shows the performance improvement of SAA on the model with 5 seeds. Although the magnitude of SAA's performance improvement on the model is different under different seeds, there is a steady improvement in general. (b) shows the performance growth during training. For the same iterations, SAA can significantly improve the performance of the model.

achieves an accuracy of 95.21%, which outperforms current SSL models. FlexMatch, which improves FixMatch by adjusting the threshold, can be further improved with the help of SAA. For example, on CIFAR-10 with 40 labels, SAA helped FlexMatch increase its mean accuracy to 95.31%. For more difficult tasks, SAA helped FlexMatch increase its mean accuracy to 75.01% and 90.85% on CIFAR-100 with 2500 labels and STL-10, which outperforms all current SSL models. Note that for unbalanced datasets, FlexMatch's threshold estimates for each class can produce large deviations, which is the reason for FlexMatch performs less favorably under the SVHN tasks. Since SVHN is a simple task, a fixed high threshold in FixMatch is more advantageous. When applying SAA to FixMatch, it is also possible to further improve its performance and outperform existing SSL models.

**SAA accelerates the improvement of model performance.** Figure 5b shows the performance curve of the model training with the same seed. For example, at the 200k-th and 400k-th iterations, SAA helps FixMatch improve its performance from 86.72% and 87.26% to 91.15% and 91.11%, respectively. FlexMatch can also improve the performance of FixMatch for the same iterations by adjusting the confidence threshold so that the model can learn more samples. SAA, on the other hand, allows *naive samples* to be learned more effectively, and therefore succeeds in further enhancing the learning of the model. For example, FixMatch reached an accuracy of 87.26%, FlexMatch helped FixMatch improve to 92.74%, and SAA helped FlexMatch further improve to 94.22%. More often, we can observe that FixMatch encountered performance stagnant between approximately 200k-th and 600k-th iterations. This is because during this period the model learns a mass of strongly augmented versions that are non-useful for model performance improvement, while our proposed SAA successfully avoids this phenomenon by changing the augmentation for *naive samples*.

| #Methods of selecting samples to apply $\mathcal{A}'$ | CIFAR-10 | STL-10 |
|---|---|---|
| Baseline-1: Applying $\mathcal{A}$ to all samples | 92.50±0.67 | 85.27±1.15 |
| Baseline-2: Applying $\mathcal{A}'$ to all samples | 92.98±2.94 | 83.19±3.98 |
| Baseline-3: Applying $\mathcal{A}'$ to 50% samples (random) | 94.05±2.00 | 85.98±2.98 |
| Setting the threshold on $\mathcal{H}$: | | |
| Fixed threshold (0.001) | 93.82±0.95 | 85.98±1.00 |
| Fixed threshold (0.002) | 94.10±1.22 | 85.22±1.87 |
| Fixed proportion threshold (25%) | 93.10±0.89 | 85.38±1.20 |
| Fixed proportion threshold (50%) | 93.87±1.52 | 86.08±1.97 |
| Fixed proportion threshold (75%) | 93.85±2.29 | 84.29±2.03 |
| OTSU threshold | **94.50±1.05** | **87.92±1.46** |

Table 2: Different methods of selecting samples that applied with $\mathcal{A}'$. Experiment on conducted on the base of Fix-Match. There are three ways to set the threshold on $\mathcal{H}$: 1) fixed value; 2) percentile of sorted $\mathcal{H}$; 3) automatic OTSU division.

## 6. Discussion

**The more diverse augmentation is not applicable to all samples.** Our approach differs from [17] in that we only apply the more diverse augmentation to a subset of samples (i.e., the *naive samples*). We experimentally validated this, as shown in Table 2 with Baseline-1 and Baseline-2. It can be clearly seen that applying diverse augmentation to all samples can lead to instability and reduce performance on some tasks. This indicates that some images have too much semantic information corrupted under augmentation $\mathcal{A}'$, leading to an accumulation of errors. To further explore this, we applied $\mathcal{A}'$ to a randomly selected sample of 50% at each epoch and the mean test accuracy of the model was slightly improved, but still unstable. This further gives us the sense that more diverse augmentation is necessary, but can only be used on the *naive samples* to work better.

**Adaptively dividing *naive samples*.** To identify the *naive samples*, we use the historical consistency loss of the sample. We tested this approach on CIFAR-10 and STL-10 with different threshold settings. As shown in Table 2, both

(a) Proportion of *naive sample*.

(b) Consistency loss of two *naive samples*.

Figure 6: (a) shows the proportion of *naive sample* divided by OTSU with model training. We recorded the proportion of *naive samples* on CIFAR-10 with 40 labels, SVHN with 40 labels and STL-10. The pictures on the left and right are FixMatch (w/SAA) and FlexMatch (w/SAA) respectively. (b) shows the consistency loss of two *naive samples*. The pictures on the left and right are FixMatch and FixMatch (w/SAA) respectively.

the fixed threshold and fixed scale approaches have a boosting effect on the model, although the effect is unstable and varies for different datasets. For example, the fixed threshold $\tau_s = 0.002$ outperforms better on CIFAR-10 task, while $\tau_s = 0.001$ outperforms better on STL-10 task, so the fixed threshold will be a more tricky hyperparameter for different datasets. Compared with the first two approaches, OTSU is not only better adapted to cross-dataset tasks, but can also be tuned as the model is trained. Figure 6a shows the proportion of *naive sample* divided by the OTSU method at different iterations. We can find that *naive samples* are not only related to task difficulty, but also to model performance. For simpler datasets, the proportion of *naive samples* is greater, and as model performance increases, more samples are also treated as *naive samples*.

**SAA prefers warm-up models.** We anaylze the model warm-up on of CIFAR-10 with 40 labels. The results in Figure 7 show that model warm-up with 100k iterations (10% of total iterations) performs the best. This is because more diverse augmented images are more difficult to recognize, which



Figure 7: Model warm-up. Experiments are conducted on CIFAR-10 with 40 labels.

can damage the initial training of the model. Therefore, SAA performs better on warm-up models. However, warming up the model too soon would reduce the action time of the SAA.

**SAA allows the augmented versions can further optimize the model.** Figure 6b compares the training loss of *naive samples* with and without SAA in FixMatch. The plot shows that without SAA, the loss of *naive samples* remains close to 0 most of the time after the 100th epoch, indicating that strongly augmented versions are not helpful for model training. However, with SAA, the number of times that strongly augmented versions aid model training

significantly increases due to the dynamic adjustment of the augmentation strategy for these samples.

**Augment on image, not on patches.** Previous work [17] proposed to first cut an image into crops and then apply augmentation on them, while we instead apply augmentation on the whole image and then cut it into crops. We conducted experiments to compare these two methods, and the results are shown in Figure 8. Our augmentation method performs better on both CIFAR-10 and STL-



Figure 8: Different more diverse augmentation. Experiments are conducted on CIFAR-10 with 40 labels and STL-10.

10 tasks. We attribute this to the fact that augmenting the whole image preserves more semantic information, which is safer for training SSL models.

**Limitations**. Since the augmentation we use is unlearnable, there is no guarantee that every augmented version is capable of contributing to model learning. Thus, the SAA serves to increase the likelihood that the augmented versions are useful to the model. In addition to this, if the model is already making good use of the sample, further augmentation may not be necessary or may even be detrimental to the model's performance, then the role of SAA is diminished.

# 7. Conclusion

In this paper, we first discuss the characteristics of *naive samples* and their impact on model training and highlight that these samples should receive attention to uncover more value. We propose SAA to achieve this goal, which identifies *naive samples* in real-time and dynamically adjusts their augmentation strategy so that they can contribute to model training. Our experimental results show that SAA significantly improves the performance of SSL methods, such as FixMatch and FlexMatch, on various datasets.

# 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 *International Joint Conference on Neural Networks*, pages 1–8. IEEE, 2020. 3

[2] Philip Bachman, Ouais Alsharif, and Doina Precup. Learning with pseudo-ensembles. *Advances in Neural Information Processing Systems*, 27, 2014. 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. *In Eighth International Conference on Learning Representations*, 2020. 2, 6

[4] David Berthelot, Nicholas Carlini, Ian Goodfellow, Avital Oliver, Nicolas Papernot, and Colin Raffel. Mixmatch: A holistic approach to semi-supervised learning. *Advances in Neural Information Processing Systems*, 32, 2019. 2, 6

[5] Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised feature learning. In *Proceedings of the fourteenth International Conference on Artificial Intelligence and Statistics*, pages 215–223. JMLR Workshop and Conference Proceedings, 2011. 6

[6] Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le. Autoaugment: Learning augmentation strategies from data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 113–123, 2019. 3

[7] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data augmentation with a reduced search space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2020. 2, 3

[8] Ali Dabouei, Sobhan Soleymani, Fariborz Taherkhani, and Nasser M Nasrabadi. Supermix: Supervising the mixing data augmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13794–13803, 2021. 3

[9] Navneet Dalal and Bill Triggs. Histograms of oriented gradients for human detection. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, volume 1, pages 886–893. Ieee, 2005. 3

[10] Terrance DeVries and Graham W Taylor. Improved regularization of convolutional neural networks with cutout. *In Fifth International Conference on Learning Representations*, 2018. 2, 3

[11] Piotr Dollár, Zhuowen Tu, Pietro Perona, and Serge Belongie. Integral channel features. 2009. 3

[12] Yue Duan, Lei Qi, Lei Wang, Luping Zhou, and Yinghuan Shi. Rda: Reciprocal distribution alignment for robust semi-supervised learning. In *European Conference on Computer Vision*, pages 533–549. Springer, Cham, 2022. 2

[13] Yue Duan, Zhen Zhao, Lei Qi, Lei Wang, Luping Zhou, Yinghuan Shi, and Yang Gao. Mutexmatch: semi-supervised learning with mutex-based consistency regularization. *Transactions on Neural Networks and Learning Systems*, 2022. 2

[14] Pedro F Felzenszwalb, Ross B Girshick, David McAllester, and Deva Ramanan. Object detection with discriminatively trained part-based models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32:1627–1645, 2010. 3

[15] 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. 3

[16] Guan Gui, Zhen Zhao, Lei Qi, Luping Zhou, Lei Wang, and Yinghuan Shi. Improving barely supervised learning by discriminating unlabeled samples with super-class. In *Advances in Neural Information Processing Systems*, 2022. 2, 3

[17] Junlin Han, Pengfei Fang, Weihao Li, Jie Hong, Mohammad Ali Armin, Ian Reid, Lars Petersson, and Hongdong Li. You only cut once: Boosting data augmentation with a single cut. *International Conference on Machine Learning*, 2022. 3, 6, 7, 8

[18] 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. 6

[19] Minui Hong, Jinwoo Choi, and Gunhee Kim. Stylemix: Separating content and style for enhanced data augmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14862–14870, 2021. 3

[20] SouYoung Jin, Aruni RoyChowdhury, Huaizu Jiang, Ashish Singh, Aditya Prasad, Deep Chakraborty, and Erik Learned-Miller. Unsupervised hard example mining from videos for improved object detection. In *Proceedings of the European Conference on Computer Vision*, pages 307–324, 2018. 3

[21] Guoliang Kang, Xuanyi Dong, Liang Zheng, and Yi Yang. Patchshuffle regularization. *arXiv preprint arXiv:1707.07103*, 2017. 3

[22] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *Proceedings of the IEEE International Conference on Computer Vision workshops*, pages 554–561, 2013. 3

[23] A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. *Handbook of Systemic Autoimmune Diseases*, 1(4), 2009. 6

[24] Chia-Wen Kuo, Chih-Yao Ma, Jia-Bin Huang, and Zsolt Kira. Featmatch: Feature-based augmentation for semi-supervised learning. In *European Conference on Computer Vision*, pages 479–495, 2020. 2

[25] Samuli Laine and Timo Aila. Temporal ensembling for semi-supervised learning. *International Conference on Learning Representations*, 2017. 2

[26] Junnan Li, Caiming Xiong, and Steven CH Hoi. Comatch: Semi-supervised learning with contrastive graph regularization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021. 3, 6

[27] Sungbin Lim, Ildoo Kim, Taesup Kim, Chiheon Kim, and Sungwoong Kim. Fast autoaugment. *Advances in Neural Information Processing Systems*, 32, 2019. 3

[28] Zilong Lin, Yong Shi, and Zhi Xue. Idsgan: Generative adversarial networks for attack generation against intrusion detection. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, pages 79–91. Springer, 2022. 3

[29] Ilya Loshchilov and Frank Hutter. Online batch selection for faster training of neural networks. *International Conference on Learning Representations Workshop*, 2015. 3

[30] 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

[31] Francisco J Moreno-Barea, Fiammetta Strazzera, José M Jerez, Daniel Urda, and Leonardo Franco. Forward noise adjustment scheme for data augmentation. In *2018 IEEE Symposium Series on Computational Intelligence (SSCI)*, pages 728–734. IEEE, 2018. 3

[32] Yuval Netzer, Tao Wang, Adam Coates, Bo Bissacco, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. *In NIPS Workshop on Deep Learning and Unsupervised Feature Learning*, 2011. 6

[33] Nobuyuki Otsu. A threshold selection method from gray-level histograms. *IEEE Transactions on Systems, Man, and Cybernetics*, 9(1):62–66, 1979. 4

[34] Yassine Ouali, Céline Hudelot, and Myriam Tami. An overview of deep semi-supervised learning. *arXiv preprint arXiv:2006.05278*, 2020. 2

[35] Mehdi Sajjadi, Mehran Javanmardi, and Tolga Tasdizen. Regularization with stochastic transformations and perturbations for deep semi-supervised learning. *Advances in Neural Information Processing Systems*, 29, 2016. 2

[36] Connor Shorten and Taghi M Khoshgoftaar. A survey on image data augmentation for deep learning. *Journal of Big Data*, 6(1):1–48, 2019. 3

[37] Abhinav Shrivastava, Abhinav Gupta, and Ross Girshick. Training region-based object detectors with online hard example mining. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 761–769, 2016. 3

[38] Edgar Simo-Serra, Eduard Trulls, Luis Ferraz, Iasonas Kokkinos, and Francesc Moreno-Noguer. Fracking deep convolutional image descriptors. *arXiv preprint arXiv:1412.6537*, 2014. 3

[39] Kihyuk Sohn, David Berthelot, Chun-Liang Li, Zizhao Zhang, Nicholas Carlini, Alex Kurakin, Han Zhang, and Colin Raffel. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. *Advances in Neural Information Processing Systems*, 2020. 1, 3, 4, 6

[40] Jiawei Su, Danilo Vasconcellos Vargas, and Kouichi Sakurai. One pixel attack for fooling deep neural networks. *IEEE Transactions on Evolutionary Computation*, 23(5), 2019. 3

[41] Kai Sheng Tai, Peter Bailis, and Gregory Valiant. Sinkhorn label allocation: Semi-supervised classification via annealed self-training. *Workshop on Challenges in Representation Learning*, 2021. 6

[42] 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*, 2017. 2, 6

[43] Yuji Tokozume, Yoshitaka Ushiku, and Tatsuya Harada. Between-class learning for image classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5486–5494, 2018. 3

[44] Vikas Verma, Alex Lamb, Juho Kannala, Yoshua Bengio, and David Lopez-Paz. Interpolation consistency training for semi-supervised learning. In *International Joint Conference on Artificial Intelligence*, pages 3635–3641, 2019. 2

[45] Jianfeng Wang, Thomas Lukasiewicz, Daniela Massiceti, Xiaolin Hu, Vladimir Pavlovic, and Alexandros Neophytou. Np-match: When neural processes meet semi-supervised learning. In *International Conference on Machine Learning*, pages 22919–22934, 2022. 6

[46] Xiaolong Wang and Abhinav Gupta. Unsupervised learning of visual representations using videos. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2794–2802, 2015. 3

[47] Lingxi Xie, Jingdong Wang, Zhen Wei, Meng Wang, and Qi Tian. Disturblabel: Regularizing cnn on the loss layer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4753–4762, 2016. 3

[48] Yi Xu, Lei Shang, Jinxing Ye, Qi Qian, Yu-Feng Li, Baigui Sun, Hao Li, and Rong Jin. Dash: Semi-supervised learning with dynamic thresholding. In *Proceedings of the International Conference on Machine Learning*, 2021. 2, 3, 6

[49] Hong Xuan, Abby Stylianou, Xiaotong Liu, and Robert Pless. Hard negative examples are hard, but useful. In *European Conference on Computer Vision*, pages 126–142. Springer, 2020. 3

[50] Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6023–6032, 2019. 3

[51] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *In Proceedings of the British Machine Vision Conference*, 2016. 6

[52] Bowen Zhang, Yidong Wang, Wenxin Hou, Hao Wu, Jindong Wang, Manabu Okumura, and Takahiro Shinozaki. Flexmatch: Boosting semi-supervised learning with curriculum pseudo labeling. *Advances in Neural Information Processing Systems*, 2021. 2, 3, 4, 6

[53] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. Mixup: Beyond empirical risk minimization. *In Fifth International Conference on Learning Representations*, 2017. 3

[54] Zhen Zhao, Luping Zhou, Yue Duan, Lei Wang, Lei Qi, and Yinghuan Shi. Dc-ssl: Addressing mismatched class distribution in semi-supervised learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9757–9765, 2022. 2

[55] Mingkai Zheng, Shan You, Lang Huang, Fei Wang, Chen Qian, and Chang Xu. Simmatch: Semi-supervised learning with similarity matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14471–14481, 2022. 2, 3, 4, 6

[56] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random erasing data augmentation. In *Proceedings of the AAAI Conference on Artificial Untelligence*, volume 34, pages 13001–13008, 2020. 3