

Active Semi-Supervised Learning by Exploring Per-Sample Uncertainty and Consistency

Jaeseung Lim^{1,2}, Jongkeun Na¹, and Nojun Kwak²

¹SNUAILAB

²Seoul National University

Abstract

Active Learning (AL) and Semi-supervised Learning are two techniques that have been studied to reduce the high cost of deep learning by using a small amount of labeled data and a large amount of unlabeled data. To improve the accuracy of models at a lower cost, we propose a method called Active Semi-supervised Learning (ASSL), which combines AL and SSL. To maximize the synergy between AL and SSL, we focused on the differences between ASSL and AL. ASSL involves more dynamic model updates than AL due to the use of unlabeled data in the training process, resulting in the temporal instability of the predicted probabilities of the unlabeled data. This makes it difficult to determine the true uncertainty of the unlabeled data in ASSL. To address this, we adopted techniques such as exponential moving average (EMA) and upper confidence bound (UCB) used in reinforcement learning. Additionally, we analyzed the effect of label noise on unsupervised learning by using weak and strong augmentation pairs to address data-inconsistency. By considering both uncertainty and data-inconsistency, we acquired data samples that were used in the proposed ASSL method. Our experiments showed that ASSL achieved about 5.3 times higher computational efficiency than SSL while achieving the same performance, and it outperformed the state-of-the-art AL method.

1. Introduction

In recent years, deep learning has drastically innovated computer vision and other research areas. However, for a practical deep learning solution, a large amount of data is inevitable. For this reason, collecting and labeling a large amount of data has become a major barrier in the development and commercialization of deep learning solutions, requiring significant time and cost. To address this issue, several methods have been researched for decades. Semi-

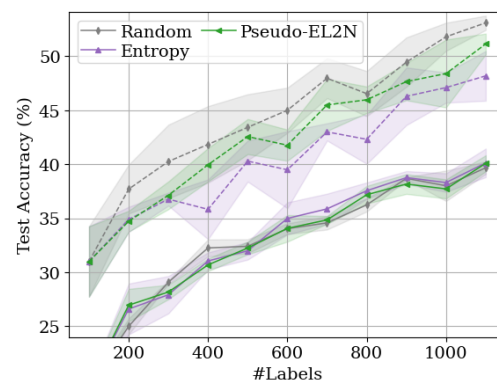


Figure 1: Accuracy comparison of random sampling and representative uncertainty-based methods on CIFAR-10. Solid and dashed lines are the results of AL and ASSL, respectively. In the AL case, the difference among the methods is almost negligible. However, in the ASSL case, even though the accuracies are better than those of AL, we can see that the improvement of accuracy for the two uncertainty-based methods is significantly lower than that of random sampling.

supervised learning (SSL) and active learning (AL) are the two representative methods among them. AL is a method of selecting a small number of data samples from a large data pool for labeling and using the newly labeled samples as well as the original labeled samples for further learning. Its acquisition function tries to select the most informative samples. This method can increase the performance of the model more efficiently than labeling the entire data pool by iteratively updating the model. In deep learning literature, one of the representative methods of AL is uncertainty-based methods that express the informativeness of data samples as scores [44, 50, 33, 12], while diversity-based methods select data samples that generalize the domain represen-

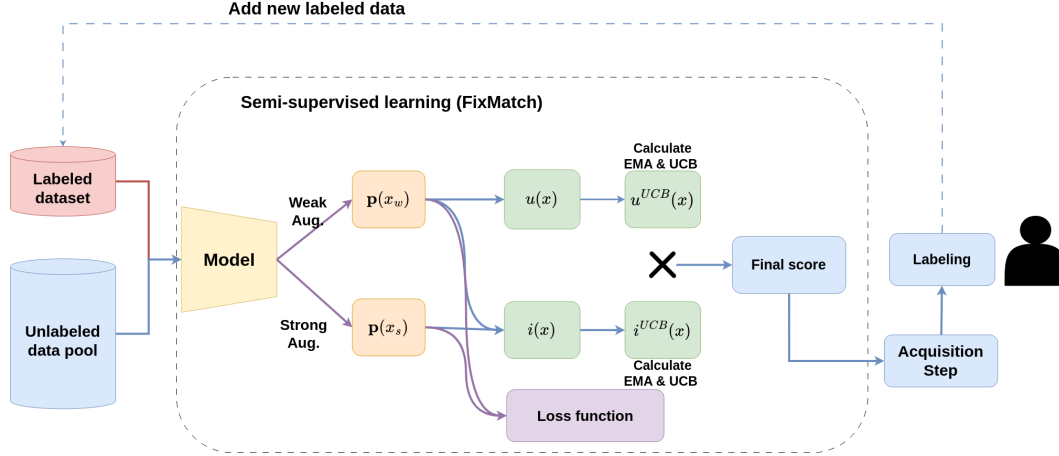


Figure 2: The overview and scoring flow of our proposed method. red line is flow for labeled data, blue line is flow for unlabeled data and purple line is flow for both. Our proposed method calculate uncertainty and data-inconsistency of unlabeled data and utilize the exponential moving average and upper confidence bound to achieve final active learning score, which was derived from the probability vector of the unlabeled data pool obtained during the SSL process.

tation of the entire data pool well in the feature embedding space [48, 36]. Recently, research has been conducted on combining these two methods [1, 3, 29] and utilizing methods such as adversarial networks [19, 38, 37].

Unlike AL, SSL directly utilizes large unlabeled data without an additional labeling process to improve the accuracy of a model with limited labeled data. Consistency regularization methods [42, 5, 6, 39, 26, 46] are the most popular methods in SSL, that enforce the outputs of pairs of differently augmented unlabeled (and labeled) samples to be similar. FixMatch [39] is the most representative method of this kind, which uses weak and strong augmentation pairs and pseudo-labeling to generate supervisory signals for unlabeled data.

Recently, to present a more powerful and efficient learning methodology by combining AL and SSL, active semi-supervised learning (ASSL) [13, 11, 21, 40] has been researched. These methods use the data pool of AL as the unlabeled data for SSL and have shown outstanding performance in object classification and detection tasks. However, as Oliver *et al.* [28] pointed out, training an SSL algorithm typically requires a very large number of training steps ($\sim 500k$), and if SSL is simply used in each iteration of AL, too much computational resource is required, which can prohibit ASSL from practical usage. However, since SSL can still yield better results than supervised learning using a small number of labeled data, if utilized appropriately, it can be combined with AL resulting in an efficient ASSL scheme. Therefore, we tried to combine AL’s acquisition step with a smaller number of training steps than typical SSL training to achieve more accurate and faster per-

formance than traditional AL or SSL.

However, there are two main challenges in directly combining AL and SSL. First, calculating the uncertainty of unlabeled samples is quite difficult. In the training of SSL, unlabeled data cause more dynamic model updates compared to supervised learning since the target for a sample is not fixed during training. This can lead to catastrophic forgetting [20] to unlabeled data. In other words, at the moment a specific unlabeled mini-batch is used for learning, the prediction of other unlabeled samples is affected, and this means that the uncertainty score of the unlabeled samples changes according to the acquisition time of ASSL. This phenomenon also occurs even after a significant number of training steps of ASSL, making it difficult for ASSL to view the observed uncertainty value as truly informativeness of the data. For example, as can be seen in Fig. 1, uncertainty-based methods show poorer performance than random sampling in ASSL. We name this problem as ‘*temporal-instability*’ which highlights that the prediction of unlabeled data suffers from time-dependent changes. To tackle this issue, we calculate the uncertainty for each unlabeled mini-batch and apply a moving average to minimize the effect of temporal-instability in data scoring, as shown in Fig. 2. Furthermore, we borrow the concept of *upper confidence bound* (UCB) from reinforcement learning to acquire more informative samples based on the variance of uncertainty score.

Second, SSL uses two types of loss functions, supervised loss and unsupervised loss, for model updates through back-propagation. However, uncertainty only focuses on informative samples for the supervised loss, and does not con-

sider the impact of unlabeled data on the unsupervised loss. Therefore, we measure the *data-inconsistency* of the weak and strong augmentation pairs of unlabeled samples, similarly as in [11, 13], to analyze the impact of unlabeled data on unsupervised loss. Through this, we identify data samples that give noisy supervisory signals during the consistency regularization process and incorporate this into the acquisition function to improve not only the performance of supervised learning but also that of unsupervised learning. Our main contributions can be summarized as follows:

- We propose a practical ASSL workflow to reduce the training cost, which speeds up training more than 6 times compared to full SSL training that takes around $500k$ steps.
- We observed the problem of instability in measuring the uncertainty of a sample in ASSL caused by temporal instability and resolved it by exploiting the exponential moving average (EMA) and its upper confidence bound (UCB) to acquire the informative samples in ASSL.
- We utilize data-inconsistency of weak & strong augmentation pair and detect unlabeled samples that make noisy supervisory signal in ASSL.
- We propose a new acquisition function in ASSL by comprehensively considering samples' uncertainty in prediction and data-inconsistency and the effectiveness of the proposed method is demonstrated through thorough experiments.

2. Related works

2.1. Active learning

Active learning is traditionally divided into two categories: uncertainty-based and diversity-based approaches. Uncertainty-based methods involve measuring the uncertainty of the unlabeled data pool and selecting the N samples with the highest uncertainty. Uncertainty can be determined by computing the entropy of the probability of a data sample [44], the margin of the first and second predicted probability [33], or by using Bayesian deep neural networks with Monte Carlo (MC) dropout [12]. More recent studies, such as LLAL [50], predict the loss of a data sample. Although uncertainty-based methods are relatively simple and have demonstrated excellent results, they tend to repeatedly select data samples with specific features, leading to poor generalization performance on the actual data pool since they do not consider the distribution of the data pool.

Diversity-based methods select data samples that represent the overall distribution of the data pool. For instance, studies have investigated clustering the unlabeled data pool to increase diversity [49] and creating a CoreSet of the data

pool [35]. These methods do not suffer from the limitations of uncertainty-based approaches, but they are influenced by the density of the data pool and have limited success in improving the decision boundary. Furthermore, since they require comparing the distance between feature embeddings, they have the disadvantage of relatively high computation.

Recently, there has been an emergence of hybrid methods that combine both uncertainty-based and diversity-based approaches. CDAL [1], for example, introduces a weighted probability distribution to define the context of data samples that applies the uncertainty of data samples and aims to maximize diversity between the contexts of data samples. BADGE [3], on the other hand, proposes gradient embedding by combining the gradients of data samples with feature embeddings and uses a strategy of clustering and selecting center points. ALFA-Mix [29] measures the sensitivity of a data sample by interpolating labeled and unlabeled sample feature embeddings and selecting center points of sensitive candidates.

2.2. Semi-supervised learning

Consistency regularization [34] is one of the most widely used methods for semi-supervised learning. Recent approaches of this kind have focused on regularizing consistency between weak and strong augmentations of unlabeled data [46, 5], and FixMatch [39] have combined consistency regularization with pseudo-labeling [24]. It serves as a base framework for various methods. For example, Semi-ViT [7] utilizes the ViT [10] model with an EMA framework [8] and pseudo-Mixup based on FixMatch. Other methods, such as FlexMatch [51] and FreeMatch [45], use curriculum-pseudo-labeling and self-adaptive thresholding, respectively, to obtain a pseudo-labeling with adaptive threshold. Lastly, CR [23] proposes contrastive regularization with contrastive learning [16] of the class cluster, which involves both confident and unconfident pseudo-labels.

2.3. Active semi-supervised learning

Active learning (AL) and semi-supervised learning (SSL) both aim to leverage the unlabeled data pool to achieve efficient learning. Therefore, recent research has sought to combine the two algorithms. Song *et al.* [40] propose combining MixMatch [6] and margin [33], while Kong *et al.* [21] proposed Neural Pre-Conditioning with the gradient vectors of labeled and unlabeled samples, based on the FixMatch-DARP [17] approach. According to Gao *et al.* [13], in ASSL, assigning labels to highly inconsistent data samples based on data-inconsistency score is a more effective way to minimize unsupervised loss than acquisition based on uncertainty alone. Elezi *et al.* [11] demonstrate ASSL in an object detection task. Their method is based on CSD [15] and uses both uncertainty and data-inconsistency

to identify over-confident samples that are mispredicted. This approach produces better acquisition results than using only uncertainty and also uses pseudo-labels to low-score data samples as newly labeled data.

3. Methods

3.1. Problem setting

In this part, we define the ASSL problem for multi-class classification. Given the labeled dataset L_n and a large unlabeled data pool U_n at round n of active learning, we use both L_n and U_n as labeled and unlabeled data for SSL training. After SSL training, we apply an acquisition function to select K samples \hat{L}_n from U_n for labeling. The labeled dataset at round $n + 1$ is then updated as $L_{n+1} = L_n \cup \hat{L}_n$. Likewise, the unlabeled data pool is updated as $U_{n+1} = U_n \setminus \hat{L}_n$.

Let $C = \{1, \dots, k\}$, be the set of classes, and let x be an input image. For $y \in C$, we represent the final softmax output of the neural network f for x as $p(y|x; \theta) = \text{SOFTMAX}(f(x; \theta))$. For simplicity, we denote the vector $[p(y = 1|x; \theta), \dots, p(y = k|x; \theta)]^T$ as $\mathbf{p}(x)$. The predicted label for x is denoted as $\hat{y} = \arg \max(\mathbf{p}(x))$, and the corresponding one-hot vector is written as $\mathbf{1}_{\hat{y}} \in \mathbb{R}^k$.

3.2. Uncertainty score

As mentioned in previous works [3, 21], in deep neural networks using stochastic gradient descent (SGD), larger gradient magnitudes can result in larger updates to the network parameters. Following [30], which approximated the norm of the gradient with the Error L2-Norm (EL2N), which is the L2-norm of the difference between $\mathbf{p}(x)$ and the corresponding ground truth one-hot vector, we also approximate the gradient for unlabeled data using the predicted label \hat{y} . More specifically, we define the uncertainty score $u(x)$ for input image x as the pseudo-EL2N score between the one-hot vector of \hat{y} and $\mathbf{p}(x)$ as

$$u(x) = \|\mathbf{p}(x) - \mathbf{1}_{\hat{y}}\|_2. \quad (1)$$

3.3. Temporal-instability

As shown in Fig. 1, when using SSL in the framework of uncertainty-based AL approaches (entropy [44], pseudo-EL2N), the acquisition quality is significantly reduced compared to random sampling, unlike AL with supervised learning. To analyze this problem, we observed the prediction changes in the unlabeled data during SSL training using 100 randomly labeled data. We calculated the uncertainty $u(x)$ in Eq.(1) for unlabeled data U every 5k training step and observed how consistent the uncertainty was over time. Surprisingly, as shown in Fig. 3, the uncertainties of consecutive time-stamps showed almost no correlation. This phe-

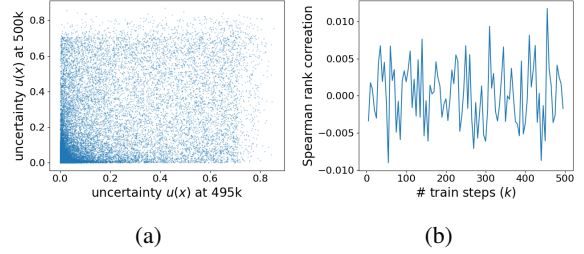


Figure 3: (a) The distribution of uncertainty(pseudo-EL2N) within 495k training step and 500k training step from unlabeled data pool. (b) Spreaman rank correlation coefficient between the uncertainty scores obtained by current and privious results of every $5k$ training step, there is almost no correlation of two training step.

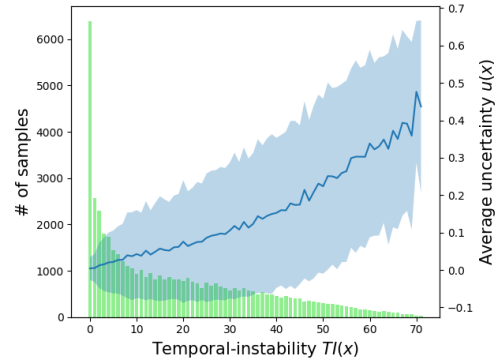


Figure 4: Mean and variance of uncertainty score according to temporal-instability. The green bars indicate the number of data for each temporal-instability value.

nomenon was consistent throughout all stages of SSL training including even after the model performance sufficiently converged. In other words, the weight parameters of the model continued to be dynamically updated during training SSL. We called this phenomenon as *temporal-instability* in SSL.

In [43], forgetting events were proposed to observe the changes in prediction and catastrophic forgetting [20] during model updates in supervised learning. Inspired by this, we quantitatively measured the temporal-instability of unlabeled data by measuring the number of times the predicted label $\hat{y}_t(x)$ observed at time t was different from $\hat{y}_{t-1}(x)$ observed at time $t - 1$, i.e.,

$$TI_T(x) = \sum_{t=1}^T \mathbb{1}(\hat{y}_t(x) \neq \hat{y}_{t-1}(x)), \quad (2)$$

where $\mathbb{1}$ denotes the indicator function. Fig. 4 shows the strong positive correlation between the uncertainty score

$u(x)$ and temporal instability $TI(x)$. In the figure, the green bar at n denotes the number of samples with $TI(x) = n$, and the blue line is the average $u(x)$ of those samples accompanied by the corresponding standard deviation. The strong correlation indicates that samples with high-uncertainty are set on high values of $TI(x)$. And except for small $TI(x)$, data samples are likely to suffer from fluctuations in predictions and thus in its pseudo-labels. Therefore, unlike conventional AL that does not incorporate SSL, the acquisition quality of uncertainty-based methods has decreased in ASSL because we cannot accurately estimate the informativeness of a data sample based on its uncertainty at a specific time.

3.4. EMA and Upper Confidence Bound

To measure the uncertainty more accurately, methods such as Gao *et al.* [13], BALD [12] and ensemble-based AL [4] average the uncertainty of multiple inference results. Also, [31] proposed using exponential moving averages (EMA) and demonstrated the use of the Upper Confidence Bound (UCB) [41], a reinforcement learning technique, to alleviate temporal-instability and maximize acquisition effectiveness with labeled data in a scenario of dynamic model update. Because we can already obtain multiple predictions at different time steps for unlabeled data during the SSL training, we also measured the uncertainty of unlabeled mini-batches at each training step and calculated the UCB through the EMA method. The uncertainty of the unlabeled sample x at time t was measured using the weakly augmented image $\mathbf{p}(x_w)$, and the equation for EMA and UCB of the uncertainty are as follows:

$$\begin{aligned}\bar{u}_t(x) &= \alpha \cdot u_t(x_w) + (1 - \alpha)\bar{u}_{t-1} \\ \bar{v}_t^u(x) &= \alpha \cdot (u_t(x_w) - \bar{u}_t)^2 + (1 - \alpha)\bar{v}_{t-1}^u \\ u_t^{UCB}(x) &= \bar{u}_t(x) + c \cdot \sqrt{\bar{v}_t^u(x)}\end{aligned}\quad (3)$$

where $\bar{u}_t(x)$ and $\bar{v}_t(x)$ are EMA and exponential moving variance of uncertainty at time t , which are initialized as $\bar{u}_0(x) = 0$ and $\bar{v}_0(x) = 0$. Also α and c are EMA rate and confidence value of UCB. In addition, to adopt EMA and UCB functions to measure uncertainty, we no longer have to infer unlabeled data in the acquisition step and obtain the uncertainty score right after the end of the SSL training step.

3.5. Data-inconsistency

In ASSL, choosing data samples with a higher UCB value of uncertainty leads to the selection of samples that can generate a higher supervised loss in the next round. However, in FixMatch, acquisition functions that rely solely on supervised loss have restrictions due to the presence of unsupervised loss from consistency regularization. As a result, we explored the influence of unlabeled data on unsupervised loss. To measure data-inconsistency, Elezi

Sorted by	Data-inconsistency			Uncertainty		
	Top 1%	Top 5%	Top 10%	Top 1%	Top 5%	Top 10%
ratio	54.93%	47.46%	44.20%	15.67%	24.87%	28.14%

Table 1: Average ratio of pseudo-labeled samples among samples with top (1/5/10)% data-inconsistency and uncertainty in fully trained SSL (=500k training step)

et al. [11] used the Kullback-Leibler (KL) divergence between the original unlabeled image $\mathbf{p}(x)$ and its flipped version $\mathbf{p}(x_f)$. FixMatch-based methods [39, 7, 51, 45, 23, 17] use random weak and strong augmentation pairs [46, 5] in the consistency regularization process, making it difficult to precisely identify inconsistencies during the acquisition step. Therefore, like uncertainty $u(x)$, we measured data inconsistency for random weak and strong augmented images x_w and x_s of unlabeled mini-batches during the training step using the KL-divergence as shown in Eq. (4) and applied EMA as in Eq. (3):

$$\begin{aligned}i(x) &= \frac{KL(p(x_w), p(x_s)) + KL(p(x_s), p(x_w))}{2} \\ \bar{i}_t(x) &= \alpha \cdot i_t(x) + (1 - \alpha) \cdot \bar{i}_{t-1}(x).\end{aligned}\quad (4)$$

Gao *et al.* [13] pointed out that samples with high data-inconsistency are overly confident but the unsupervised loss based on these samples are hard to minimize. To measure this, similar to $TI(x)$, for each sample, we counted the number of times $\max(\mathbf{p}(x)) > \tau (= 0.95)$ at every 5k training step during SSL training. We called these samples ‘pseudo-labeled’, which have high confidence, and thus participate in the unsupervised loss. Table 1 reveals that the proportion of pseudo-labeled data with high data-inconsistency is approximately twice as high as that with high uncertainty. Interestingly, these highly data-inconsistent samples are utilized more frequently during training, despite not benefiting from consistency regularization. Consequently, such samples could introduce label noise to the unsupervised loss function.

As same as uncertainty, we utilize UCB as described in Eq. (3) to determine the final data inconsistency score as shown in Eq. (5).

$$\begin{aligned}\bar{v}_t^i(x) &= \alpha \cdot (i_t(x) - \bar{i}_t)^2 + (1 - \alpha)\bar{v}_{t-1}^i \\ i_t^{UCB}(x) &= \bar{i}_t(x) + c \cdot \sqrt{\bar{v}_t^i(x)}.\end{aligned}\quad (5)$$

Consequently, we define the ultimate AL score by considering both uncertainty and inconsistency in the unlabeled data pool, as described in Eq. (6).

$$Score(x) = u^{UCB}(x) \times i^{UCB}(x).\quad (6)$$

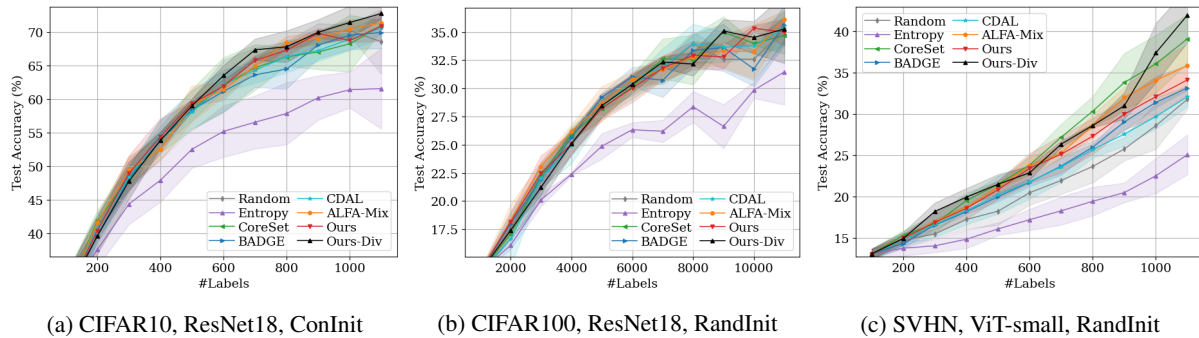


Figure 5: Comparison accuracy results with various setting.

4. Experiments

To evaluate our ASSL method, we used CIFAR-10 [22], CIFAR-100 [22], SVHN [27], and MiniImageNet [32]. For CIFAR-10 and SVHN, we randomly selected 100 images for an initial labeled dataset and acquire 100 additional images for labeling in each round. For CIFAR-100 and MiniImageNet, we started with 1,000 random initial labeled datasets and labeled 1,000 images in each round.

All experiments were conducted with 10 rounds of active learning cycles, and we tested two methods for initial weight condition: random initial weight (RandInit) and the weight from the previous round (ConInit). For RandInit, we used the same random initial weights in all rounds and the initial model. We used the ResNet-18 [14] and ViT-small [10] models with the same layer configuration as the model used in ALFA-Mix [29]. We used $1,024 \text{ steps} \times 5 \text{ epochs}$ for training in each round to achieve efficient SSL training, which is approximately 1% of the training step ($500k$) proposed by Oliver *et al.* [28]. We set the learning rate as $lr = 0.03$ for ResNet-18 and $lr = 0.01$ for ViT-small. We did not use learning rate scheduler because of smaller training steps. The other hyper-parameters used in SSL were the same as those in FixMatch [39], and we used RandAugment [9] for augmenting unlabeled data.

We implement our ASSL framework based on ALFA-Mix¹ and pytorch implementation of FixMatch². For the hyper-parameters of the acquisition function, we set the EMA rate α to 0.8, the confidence value of uncertainty c_u to 0.5, and the confidence value of data-inconsistency c_i to 2.0. We also tried to apply diversity-measure to our proposed method by simply multiplying the score value of each unlabeled sample with its feature embedding vector, and using K-means++ [2] for clustering. We denote the proposed method as ‘Ours’, and the method with diversity as ‘Ours-Div’. The results of all experiments are reported as the av-

	Random	Entropy	CoreSet	BADGE	CDAL	ALFA-Mix	Ours	Ours-Div
Random	-	11.1	3.6	2.9	1.9	2.3	3.3	2.8
Entropy	0.0	-	0.0	0.0	0.0	0.0	0.0	0.0
CoreSet	2.9	10.9	-	3.1	2.5	2.8	1.7	2.0
BADGE	3.1	11.1	3.1	-	1.8	2.0	2.7	1.8
CDAL	3.5	11.4	3.5	2.8	-	2.5	2.8	2.2
ALFA-Mix	3.9	11.4	3.5	2.9	2.1	-	2.7	1.6
Ours	3.5	11.3	2.1	3.0	3.1	2.6	-	1.8
Ours-Div	3.2	11.3	3.1	3.2	2.9	2.6	3.1	-
Average	2.87	11.21	2.7	2.56	2.04	2.11	2.33	1.74

Figure 6: Pairwise comparison matrix of SOTA methods. The overall performance(lower is better) is shown at the bottom row and maximum value is 12.

erage of three runs with different seeds.

4.1. Results

To evaluate our proposed method, we compared our methods with random sampling and state-of-the-art AL methods such as entropy [44], BADGE [3], CDAL [1], CoreSet [35], and ALFA-Mix [29]. We evaluate them by using Resnet-18 on four aforementioned datasets. Furthermore, to test if our method works on different architectures, we compare the result of ViT-small on CIFAR10 and SVHN. To comprehensively analyze the experimental results, we compute the pairwise comparison matrix [3] in Fig. 6 that each elements $e_{i,j}$ indicates that the number of i -th method outperformed j -th and its symmetric elements $e_{j,i}$ is the opposite case. the last row at bottom is average score by column that lower score means better. As shown in Fig. 6, ‘Ours-Div’ achieved the best performance in the

¹https://github.com/AminParvaneh/alpha_mix_active_learning

²<https://github.com/kekmodel/FixMatch-pytorch>

Methods	Accuracy (%)				Time (s)			
	RandInit(AL)	RandInit	ConInit	Total	CIFAR10	SVHN	CIFAR100	MiniImageNet
Random	27.99±0.84	37.36±1.32	42.49±0.90	39.93±1.11	-	-	-	-
Entropy	27.33±1.11	30.74±1.71	38.44±1.42	34.59±1.56	25.50	30.38	26.08	51.34
CoreSet	27.35±1.08	36.58±1.76	42.59±1.28	39.59±1.52	37.36	43.63	128.13	209.38
BADGE	28.28±1.02	37.15±1.43	42.01±1.59	39.58±1.51	204.65	245.87	>1 hour	>1 hour
CDAL	28.45±0.97	37.02±1.39	42.62±1.07	39.82±1.23	28.39	34.04	64.06	135.39
ALFA-Mix	28.34±1.04	37.53±1.43	42.47±0.90	40.00±1.17	83.17	104.76	490.95	1013.44
Ours	-	37.07±1.44	42.80±0.81	39.93±1.13	0.05	0.06	0.03	0.05
Ours-Div	-	37.68±1.49	42.62±1.04	40.15±1.27	36.05	43.39	129.14	209.10

Table 2: Average accuracy in terms of initial weights condition and acquisition times

combination of four datasets and two models. Also ‘Ours’ shows better results than random sampling in ASSL and outperformed CoreSet and BADGE that rely on diversity. For more details, as shown in Fig. 5,

‘Ours’ shows comparable performance to state-of-the-art (SOTA) methods, and in some cases, even shows better performance.

The left part of Table 2 is the average test accuracy of all rounds and datasets by two initial weight parameter conditions. ‘Our-Div’ shows the best average accuracy with Randinit and Total, and ‘Ours’ shows the best average accuracy with ConInit. Also, we tested conventional AL frameworks to compare the accuracies between AL and ASSL. We train the AL framework with the same hyperparameter setting that was used by ALFA-Mix [29] and average test accuracies as done in ASSL with RandInit condition. Our ASSL framework shows average accuracy improvements of 8.11%p compared to all state-of-the-art AL methods.

Moreover, diversity-based methods, such as CoreSet [35], show relatively low sensitivity to temporal-instability in the setting of ASSL. This means that even if $p(x)$ of data samples is temporally unstable, the representation of the entire domain is relatively preserved. In other words, $p(x)$ is expected to change within a specific range without random direction and instability. So adapting EMA and UCB helps to find the maximum expected value in the unstable area.

4.2. Training efficiency

In terms of training time, the proposed method was compared to FixMatch. Up to 10 rounds, the total number of training steps was 56,320, which represents about 11% of the 500k training steps. Fig. 7 compares the test accuracy of FixMatch and Ours during training. Using CIFAR-10 on the Resnet-18 model, FixMatch reached 70.82%p at 300k training steps which is the closest to 70.90%p obtained by Ours(ConInit) at around 56k training steps. In this case, our proposed method is about 5.3 times faster than SSL. In our Nvidia Titan X environment, it took approximately 20 hours for our ASSL to achieve the best accuracy, while SSL required about 4.4 days.



Figure 7: Compare accuracies of FixMatch and Ours according to training steps

In actual active learning scenarios, data labeling incurs an additional time cost. However, annotating approximately 100 images does not require significant effort. In our experiment, three people spent less than 5 minutes on this task. Therefore, in practice, the small amount of labeling work required by ASSL is a more efficient approach for achieving the target accuracy compared to fully converging the model through SSL.

In terms of acquisition time, as shown in the right section of Table 2, ‘Ours’ calculates the score during the training process, so it can be confirmed that there is almost no time cost in the acquisition step. In addition, similar to CoreSet [35], ‘Ours-Div’ also has a faster acquisition time compared to BADGE [3] or ALFA-Mix [29] whose time costs depend on the number of classes.

4.3. Ablation study

4.3.1 Uncertainty and its UCB

In the ablation study, we analyzed the effects of EMA and UCB of uncertainty (Pseudo-EL2N) with CIFAR-10, as shown in Fig. 8. We found that using EMA achieves a better average accuracy than vanilla Pseudo-EL2N, and we observed that the accuracy score was higher when additionally using UCB. The average accuracy is the highest when the

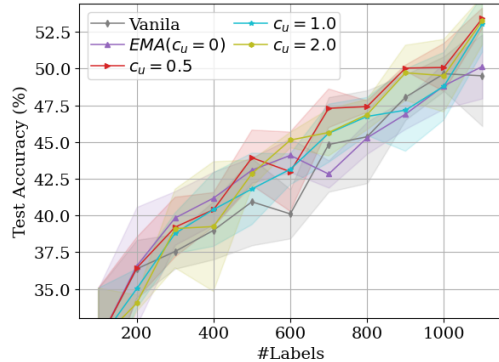


Figure 8: Pseudo-EL2N and its UCB with various confidence value in CIFAR10 and RandInit

confidence value of UCB is set to $c_u = 0.5$. Therefore, when the model updates dynamically like in SSL, scoring through moving average is an effective method, and applying additional UCB shows that samples with large supervised loss were selected.

4.3.2 Data-Inconsistency and its UCB

To analyze the effect of data inconsistency, as shown in Figure 9, we found that $c_i = 2.0$ provided the best performance when combined with u^{UCB} . This means that using i^{UCB} with high confidence values is most effective, as data-inconsistency is important to have high values while also having high variance. Reducing the instability of data-inconsistency by UCB improves accuracy than considering only EMA of data-inconsistency.

5. Discussion

Trade-off of labeling cost and training cost. In this paper, we conducted experiments using same number of data samples as AL framework for comparison. However, we believe that it is possible to achieve higher accuracy with less labeling cost by leveraging training steps of SSL. On the other hand, simply acquiring more labeled samples would decrease the training cost to achieve the same accuracy. So there would be a trade-off relationship between increasing labeling cost to use fewer training cost. Therefore, adjusting the balance between the two costs according to each task would make our proposed method more efficient.

Data-inconsistency. Through this study, we began to examine how unlabeled data works in SSL. By identifying data inconsistency, we observed unlabeled samples that are not trainable and create noisy supervisory signals. By considering data inconsistency in the acquisition steps, we obtained better results. However, we did not deeply investigate how and to what extent such data samples have a

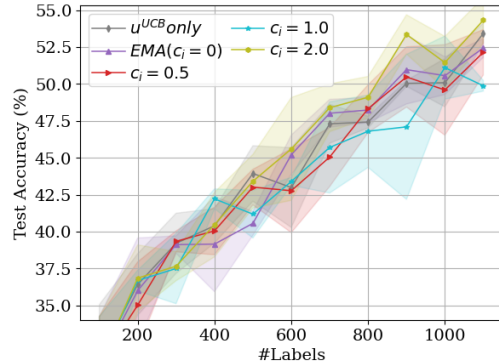


Figure 9: Adapt Data-inconsistency with various confidence value in CIFAR10, RandInit and $c_u = 0.5$

negative impact on SSL. For example, ignoring high data-inconsistency samples during the ASSL and SSL training, or simply removing those samples from unlabeled data pool during the acquisition steps. Additionally, we did not explore whether low data inconsistency data is useful for training. Therefore, we believe that further research and supplementation of these aspects can lead to the development of new methodologies in ASSL and SSL.

Apply to other tasks. As Elezi *et al.* [11] proposed, ASSL is tried to apply to object detection tasks. Also semi supervised learning methods based on consistency regularization has been studied continuously for many computer vision tasks, such as segmentation [25] and human pose estimation [18, 47] that suffer high labeling costs than object classification and object detection. We believe that if we apply the proposed ASSL method to the above methods. we can update model fast with more reasonable costs as object classification.

6. Conclusion

We present a novel approach to Active Semi-supervised Learning (ASSL) that effectively addresses temporal instability. ASSL faces a challenge in accurately determining the informativeness of unlabeled data using conventional uncertainty-based functions due to temporal instability. To tackle this challenge, we leverage the exponential moving average (EMA) and its upper confidence bound (UCB) of uncertainty to identify truly informative data samples. Additionally, we introduce data-inconsistency to identify samples that are not trainable and incorporate both factors comprehensively in the acquisition function. Our experimental results demonstrate that our method outperforms state-of-the-art active learning methods, with an 8.11%p improvement over existing AL methods [44, 35, 3, 1, 29], while also being approximately 5.3 times more efficient than semi-supervised learning. Going forward, we plan to focus on

developing more practical ASSL methods with a smaller number of labeled data and adapting them to other tasks.

References

- [1] Sharat Agarwal, Himanshu Arora, Saket Anand, and Chetan Arora. Contextual diversity for active learning. In *Computer Vision – ECCV 2020*, pages 137–153, 2020.
- [2] David Arthur and Sergei Vassilvitskii. K-means++: The advantages of careful seeding. In *Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA '07*, page 1027–1035, USA, 2007. Society for Industrial and Applied Mathematics.
- [3] Jordan T. Ash, Chicheng Zhang, Akshay Krishnamurthy, John Langford, and Alekh Agarwal. Deep batch active learning by diverse, uncertain gradient lower bounds. In *International Conference on Learning Representations*, 2020.
- [4] William H Beluch, Tim Genewein, Andreas Nürnberger, and Jan M Köhler. The power of ensembles for active learning in image classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 9368–9377, 2018.
- [5] David Berthelot, Nicholas Carlini, Ekin D Cubuk, Alex Kurakin, Kihyuk Sohn, Han Zhang, and Colin Raffel. Remixmatch: Semi-supervised learning with distribution alignment and augmentation anchoring. *arXiv preprint arXiv:1911.09785*, 2019.
- [6] David Berthelot, Nicholas Carlini, Ian Goodfellow, Nicolas Papernot, Avital Oliver, and Colin A Raffel. Mixmatch: A holistic approach to semi-supervised learning. *Advances in neural information processing systems*, 32, 2019.
- [7] Zhaowei Cai, Avinash Ravichandran, Paolo Favaro, Manchen Wang, Davide Modolo, Rahul Bhotika, Zhuowen Tu, and Stefano Soatto. Semi-supervised vision transformers at scale. In *NeurIPS*, 2022.
- [8] Zhaowei Cai, Avinash Ravichandran, Subhransu Maji, Charles Fowlkes, Zhuowen Tu, and Stefano Soatto. Exponential moving average normalization for self-supervised and semi-supervised learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 194–203, 2021.
- [9] 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*, pages 702–703, 2020.
- [10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021.
- [11] Ismail Elezi, Zhiding Yu, Anima Anandkumar, Laura Leal-Taixé, and Jose M. Alvarez. Not all labels are equal: Rationalizing the labeling costs for training object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14492–14501, June 2022.
- [12] Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep Bayesian active learning with image data. In *Proceedings of the 34th International Conference on Machine Learning*, pages 1183–1192, 06–11 Aug 2017.
- [13] Mingfei Gao, Zizhao Zhang, Guo Yu, Sercan Ö Arık, Larry S Davis, and Tomas Pfister. Consistency-based semi-supervised active learning: Towards minimizing labeling cost. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part X 16*, pages 510–526. Springer, 2020.
- [14] 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.
- [15] Jisoo Jeong, Seungeui Lee, Jeesoo Kim, and Nojun Kwak. Consistency-based semi-supervised learning for object detection. *Advances in neural information processing systems*, 32, 2019.
- [16] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschiot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *Advances in neural information processing systems*, 33:18661–18673, 2020.
- [17] Jaehyung Kim, Youngbum Hur, Sejun Park, Eunho Yang, Sung Ju Hwang, and Jinwoo Shin. Distribution aligning refinery of pseudo-label for imbalanced semi-supervised learning. *Advances in neural information processing systems*, 33:14567–14579, 2020.
- [18] JongMok Kim, Hwijun Lee, Jaeseung Lim, Jongkeun Na, Nojun Kwak, and Jin Young Choi. Pose-mum : Reinforcing key points relationship for semi-supervised human pose estimation, 2022.
- [19] Kwanyoung Kim, Dongwon Park, Kwang In Kim, and Se Young Chun. Task-aware variational adversarial active learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8166–8175, June 2021.
- [20] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.
- [21] Seo Taek Kong, Soomin Jeon, Dongbin Na, Jaewon Lee, Hong-Seok Lee, and Kyu-Hwan Jung. A neural pre-conditioning active learning algorithm to reduce label complexity. In *Advances in Neural Information Processing Systems*, 2022.
- [22] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [23] Doyup Lee, Sungwoong Kim, Ildoo Kim, Yeongjae Cheon, Minsu Cho, and Wook-Shin Han. Contrastive regularization for semi-supervised learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pages 3911–3920, June 2022.

- [24] Dong-Hyun 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*, page 896, 2013.
- [25] Yuyuan Liu, Yu Tian, Yuanhong Chen, Fengbei Liu, Vasileios Belagiannis, and Gustavo Carneiro. Perturbed and strict mean teachers for semi-supervised semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4258–4267, June 2022.
- [26] 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.
- [27] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bisacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011.
- [28] Avital Oliver, Augustus Odena, Colin A Raffel, Ekin Dogus Cubuk, and Ian Goodfellow. Realistic evaluation of deep semi-supervised learning algorithms. *Advances in neural information processing systems*, 31, 2018.
- [29] Amin Parvaneh, Ehsan Abbasnejad, Damien Teney, Gholamreza (Reza) Haffari, Anton van den Hengel, and Javen Qinfeng Shi. Active learning by feature mixing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12237–12246, June 2022.
- [30] Mansheej Paul, Surya Ganguli, and Gintare Karolina Dziugaite. Deep learning on a data diet: Finding important examples early in training. In *Advances in Neural Information Processing Systems*, 2021.
- [31] Ravi S Raju, Kyle Daruwalla, and Mikko Lipasti. Accelerating deep learning with dynamic data pruning. *arXiv preprint arXiv:2111.12621*, 2021.
- [32] Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. In *International conference on learning representations*, 2017.
- [33] Dan Roth and Kevin Small. Margin-based active learning for structured output spaces. In *Machine Learning: ECML 2006*, pages 413–424, Berlin, Heidelberg, 2006. Springer Berlin Heidelberg.
- [34] 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.
- [35] Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set approach. *arXiv preprint arXiv:1708.00489*, 2017.
- [36] Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set approach. In *International Conference on Learning Representations*, 2018.
- [37] Changjian Shui, Fan Zhou, Christian Gagné, and Boyu Wang. Deep active learning: Unified and principled method for query and training. In *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pages 1308–1318, Online, 26–28 Aug 2020. PMLR.
- [38] Samarth Sinha, Sayna Ebrahimi, and Trevor Darrell. Variational adversarial active learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.
- [39] Kihyuk Sohn, David Berthelot, Chun-Liang Li, Zizhao Zhang, Nicholas Carlini, Ekin D. Cubuk, Alex Kurakin, Han Zhang, and Colin Raffel. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. In *Advances in Neural Information Processing Systems*, 2020.
- [40] Shuang Song, David Berthelot, and Afshin Rostamizadeh. Combining mixmatch and active learning for better accuracy with fewer labels, 2019.
- [41] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press, second edition, 2018.
- [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*, 30, 2017.
- [43] Mariya Toneva, Alessandro Sordoni, Remi Tachet des Combes, Adam Trischler, Yoshua Bengio, and Geoffrey J Gordon. An empirical study of example forgetting during deep neural network learning. In *ICLR*, 2019.
- [44] Dan Wang and Yi Shang. A new active labeling method for deep learning. In *2014 International Joint Conference on Neural Networks (IJCNN)*, pages 112–119, 2014.
- [45] Yidong Wang, Hao Chen, Qiang Heng, Wenxin Hou, Yue Fan, , Zhen Wu, Jindong Wang, Marios Savvides, Takahiro Shinozaki, Bhiksha Raj, Bernt Schiele, and Xing Xie. Freematch: Self-adaptive thresholding for semi-supervised learning. In *International Conference on Learning Representations (ICLR)*, 2023.
- [46] Qizhe Xie, Zihang Dai, Eduard Hovy, Thang Luong, and Quoc Le. Unsupervised data augmentation for consistency training. *Advances in neural information processing systems*, 33:6256–6268, 2020.
- [47] Rongchang Xie, Chunyu Wang, Wenjun Zeng, and Yizhou Wang. An empirical study of the collapsing problem in semi-supervised 2d human pose estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 11240–11249, October 2021.
- [48] Yi Yang, Zhigang Ma, Feiping Nie, Xiaojun Chang, and Alexander G. Hauptmann. Multi-class active learning by uncertainty sampling with diversity maximization. *Int. J. Comput. Vision*, 113(2):113–127, jun 2015.
- [49] Yi Yang, Zhigang Ma, Feiping Nie, Xiaojun Chang, and Alexander G Hauptmann. Multi-class active learning by uncertainty sampling with diversity maximization. *International Journal of Computer Vision*, 113:113–127, 2015.
- [50] Donggeun Yoo and In So Kweon. Learning loss for active learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [51] Bowen Zhang, Yidong Wang, Wenxin Hou, Hao Wu, Jindong Wang, Manabu Okumura, and Takahiro Shinozaki. Flexmatch: Boosting semi-supervised learning with curriculum pseudo labeling. In *Neural Information Processing Systems (NeurIPS)*, 2021.