

Label Shift Adapter for Test-Time Adaptation under Covariate and Label Shifts

Sunghyun Park^{1,2} Seunghan Yang¹ Jaegul Choo² Sungrack Yun¹
¹Qualcomm AI Research* ²KAIST

¹{sunpar, seunghan, sungrack}@qti.qualcomm.com ²{psh01087, jchoo}@kaist.ac.kr

Abstract

Test-time adaptation (TTA) aims to adapt a pre-trained model to the target domain in a batch-by-batch manner during inference. While label distributions often exhibit imbalances in real-world scenarios, most previous TTA approaches typically assume that both source and target domain datasets have balanced label distribution. Due to the fact that certain classes appear more frequently in certain domains (e.g., buildings in cities, trees in forests), it is natural that the label distribution shifts as the domain changes. However, we discover that the majority of existing TTA methods fail to address the coexistence of covariate and label shifts. To tackle this challenge, we propose a novel label shift adapter that can be incorporated into existing TTA approaches to deal with label shifts during the TTA process effectively. Specifically, we estimate the label distribution of the target domain to feed it into the label shift adapter. Subsequently, the label shift adapter produces optimal parameters for target label distribution. By predicting only the parameters for a part of the pre-trained source model, our approach is computationally efficient and can be easily applied, regardless of the model architectures. Through extensive experiments, we demonstrate that integrating our strategy with TTA approaches leads to substantial performance improvements under the joint presence of label and covariate shifts.

1. Introduction

Despite the recent remarkable improvement of deep neural networks in various applications, the models still suffer from distribution shifts between source distribution and target distribution. One type of distribution shift, known as *covariate shift*, occurs when the target distribution $p_t(x)$ differs from the source distribution $p_s(x)$. In autonomous driving, for instance, models may degrade significantly during testing due to ambient factors such as weather and location. To design the models robust to covariate shifts, un-

*Qualcomm AI Research is an initiative of Qualcomm Technologies, Inc.

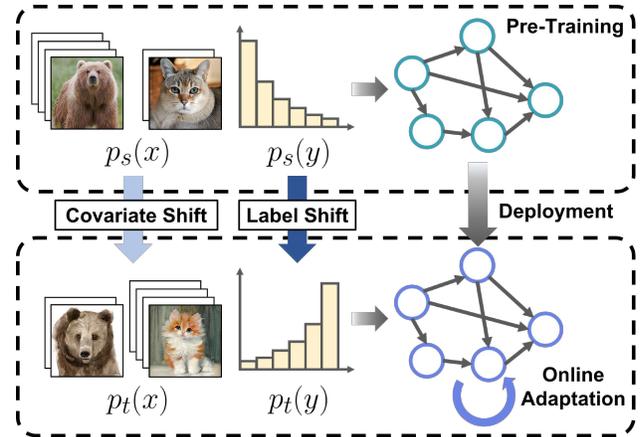


Figure 1. We consider the test-time adaptation scenario, when covariate and label shifts occur simultaneously. After deploying the pre-trained model, the model is adapted to the target domain. However, existing methods often suffer from the coexistence of covariate and label shifts. Employing our method enables online adaptation under shifted target label distributions.

supervised domain adaptation literature [9, 10, 23, 19] has explored the transfer of knowledge learned from labeled source data to unlabeled target data.

To be more practical in real-world scenarios, test-time adaptation (TTA) algorithms [46] have emerged to enhance practicality in real-world scenarios by adapting deep neural networks to the target domain during inference. Specifically, TTA approaches optimize the model parameters batch-by-batch using unlabeled test data, avoiding additional labeling costs. Previous TTA studies have mitigated the performance degradation caused by covariate shift by enhancing normalization statistics [41, 11, 24], optimizing model parameters with entropy minimization [46, 35], or utilizing pseudo labels [20].

Although the natural data encountered in practice often exhibits long-tailed label distribution, most previous TTA methods assume that the model is trained on class-balanced data. This assumption overlooks another type of distribution shift, known as *label shift*, in which label distribution varies between source $p_s(y)$ and target $p_t(y)$. Label shift has been studied extensively in the long-tailed recognition literature [18, 6, 8, 15]. Considering only one type

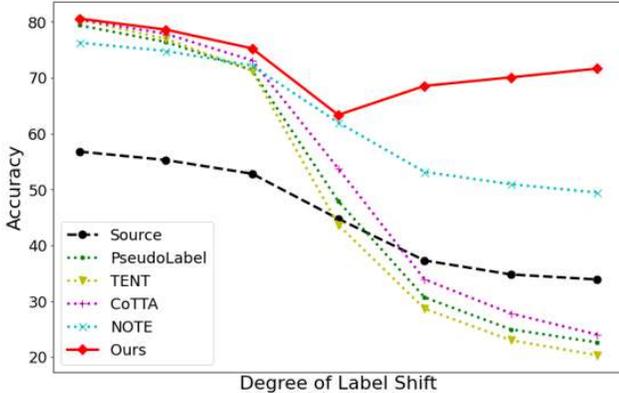


Figure 2. Plots of accuracy with different degrees of label shift in CIFAR-10-C [14]. Although we train source model using balanced softmax [39], most TTA baselines show lower performance than source model, when label shift is severe. On the other hand, our method is more robust to label shifts than other TTA methods.

of distribution shifts in real-world scenarios is infeasible, as both covariate and label shifts frequently occur simultaneously [29, 22]. For example, different object classes such as buildings and trees are more prevalent in certain environments, such as cities or forests.

We found that the majority of existing TTA approaches fail to adapt the pre-trained model when it is trained on a long-tailed dataset. This is because most TTA methods that employ entropy minimization [46, 35, 11, 7, 24] are flawed due to the model bias towards the dominant classes in source data. In other words, since the predictions on test samples are often biased toward the majority classes of training data, utilizing entropy minimization would lead the model to increase its confidence in predicting the dominant classes. While some recent TTA studies have addressed similar challenging issues, such as temporally correlated test data [11] and class-imbalanced test samples [53], they do not cover the situation where the source domain data has a long-tailed label distribution, which can lead to bias in the source model. Furthermore, we observed that training the source model on a long-tailed dataset with long-tailed recognition techniques, such as balanced softmax [39], is insufficient to stabilize TTA algorithms, as shown in Fig. 2.

To tackle such a non-trivial issue, we propose a novel label shift adapter, which is designed to adapt the pre-trained model according to the label distribution during inference. Before deploying the model to the server, we train the label shift adapter with a pre-trained source model, which takes the label distribution as input. The label shift adapter is optimized to produce the optimal parameters according to the label distribution. To make our method applicable to any model architecture, we design the label shift adapter to predict only the parameters associated with a part of the source model. After model deployment, we estimate the label distribution of target domain data for injecting the appropriate

input into the label shift adapter during inference. Moreover, our proposed method can be easily integrated with TTA methods such as TENT [46] and IABN [11] to adapt the model to the target domain. Combining TTA approaches with the proposed label shift adapter enables robust model adaptation to the target domain, even in the presence of covariate and label shifts simultaneously. Through extensive experiments, we demonstrate that our method outperforms the existing TTA methods when both source and target domain datasets have class-imbalanced label distributions.

In summary, the main contributions are as follows:

- We introduce a novel label shift adapter that produces the optimal parameters according to the label distribution. By utilizing the label shift adapter, we can develop a robust TTA algorithm that can handle both covariate and label shifts simultaneously.
- Our approach is easily applicable to any model regardless of the model architecture and pre-training process. It can be simply integrated with other TTA algorithms.
- Through extensive experiments on six benchmarks, we demonstrate that our method enhances the performance significantly when source and target domain datasets have class-imbalanced label distributions.

2. Related Work

Source-Free Domain Adaptation. Unsupervised domain adaptation (UDA) methods have been widely applied in cross-domain applications such as classification [9, 10], object detection [51], and semantic segmentation [55]. However, UDA approaches require access to both source and target domains simultaneously. This restriction makes these approaches frequently impractical due to computational costs and data privacy concerns. On the other hand, source-free domain adaptation (DA) [23, 19] overcomes this limitation by adapting a pre-trained model to the target domain using only unlabeled target data. However, existing source-free DA methods barely consider label shifts, which limits their applicability in real-world scenarios.

Domain Adaptation for Label Shift. Several methods [3, 28, 43, 17, 29, 22] have been developed to investigate a more realistic scenario of domain adaptation in which covariate and label shifts co-occur. To alleviate label shift, previous studies employ the re-weighting method [3, 28] by estimating target domain label distribution. Recent approaches employ an alternative training scheme [29] and a secondary pseudo label [22] to alleviate label shifts in unsupervised domain adaptation. However, it is challenging to adapt the model to the target domain in real time using these methods. Therefore, we focus on addressing such co-existence of covariate and label shifts in the TTA setting.

Test-Time Adaptation. Fully test-time adaptation [46] aims to improve model performance on target domain data

through adaptation with unlabeled test samples during inference. Previous work [41] improves the robustness under covariate shift by using the statistics of test batch in normalization layer. TENT [46] further optimizes affine parameters of batch normalization layers using entropy minimization. Before the pre-trained model deployment to the server, several approaches train additional modules to appropriately interpolate training and test statistics [56, 24] or regularize the model parameters [7] for TTA. Recent several studies [4, 11, 53] address the model to be more robust under non-i.i.d or class-imbalanced test samples. However, they assume that the source domain datasets are balanced, where the pre-trained model is not biased towards the dominant classes due to the class-imbalanced label distribution. Different from the existing studies, our research tackles the cases in which both source and target domain datasets are class-imbalanced, which is more challenging and practical.

Long-Tailed Recognition. It is natural that datasets have long-tailed distributions in the real world. Previous studies addressed this issue by altering loss functions [25, 8, 6, 40, 21], adjusting logits [33, 39], and employing multiple experts that are specialized in different label distributions [49, 54, 52]. Recently, several studies such as LADE [15], SADE [52], and BalPoE [1] have introduced test-agnostic long-tailed recognition, where the training label distribution is long-tailed while the test label distribution is agnostic. However, LADE needs the true test label distribution to adjust the logits, and SADE and BalPoE require multiple expert architectures. Due to these aspects, they are difficult to apply the TTA algorithms. Inspired by such studies, we design a novel label shift adapter, which has the capability to handle unknown test label distribution using training long-tailed distribution. In contrast to previous methods, our method is applicable to any model regardless of its architecture and can be employed without true test label distribution. In this paper, we focus on handling the label shifts in TTA setting.

Predicting Weights. A hypernetwork [12] is a deep neural network to produce the weights of another neural network. Hypernetworks have been developed for federated learning [42], multi-task learning [34, 26, 31], and continuous learning [45, 5]. Moreover, adapter layers between existing layers of the model have been proposed for fine-tuning [27, 16]. The key functional difference is that our method produces the parameters to handle the label shifts.

3. Method

3.1. Problem Formulation

In the TTA task, labeled samples from the source domain $\mathcal{D}_s = \{(x, y) \sim p_s(x, y)\}$ and unlabeled samples from the target domain $\mathcal{D}_t = \{x \sim p_t(x)\}$. TTA aims to predict the labels of target domain samples by updating the source

model to the target model during inference. Specifically, under the TTA scheme, the model receives a mini-batch x_t of test samples in the t -th inference step.

Generally, previous TTA literature assumes only covariate shift, where $p_s(x) \neq p_t(x)$. In other words, the existing TTA methods only consider class-balanced datasets in training and testing. Different from previous works, we assume the joint presence of covariate and label shifts [17, 43, 22]: $p_s(x) \neq p_t(x)$ and $p_s(y) \neq p_t(y)$, which is more practical and natural in real-world scenarios. In particular, when $p_s(y)$ has long-tailed label distribution, TTA methods are flawed, despite leveraging long-tailed recognition methods. It is due to the fact that most TTA methods are not able to reduce the model’s bias toward the majority classes. Our goal is to design a novel method for TTA that can perform stably regardless of $p_s(y)$ and $p_t(y)$, while the model can be adapted during inference.

3.2. Label Shift Adapter

In this paper, we propose a novel label shift adapter for TTA to handle label distribution shifts in TTA. Entropy minimization [46, 35, 7, 11, 24] is widely used for TTA to optimize the model with unlabeled test samples during inference. Intuitively, entropy minimization makes individual predictions confident. During test time, entropy minimization loss is utilized as follows:

$$\mathcal{L}_{ent} = - \sum_{x \sim p_t(x)} f(x) \log f(x), \quad (1)$$

where f denotes a model $f : x \rightarrow y$. However, if the pre-trained source model $f(x)$ is biased toward the majority classes due to the long-tailed label distribution of \mathcal{D}_s , the predictions on test samples also would be biased towards the majority classes in \mathcal{D}_s regardless of label distribution of \mathcal{D}_t . In other words, it is not appropriate to minimize under the shifted label distribution because the model prediction estimates $p_s(y|x)$, which is strongly coupled with $p_s(y)$ and may differ from $p_t(y)$, as explained by the Bayes’ rule [15]:

$$p_s(y|x) = \frac{p_s(y)p_s(x|y)}{p_s(x)} = \frac{p_s(y)p_s(x|y)}{\sum_c p_s(c)p_s(x|c)}, \quad (2)$$

where c denotes the class index.

To address diverse $p_t(y)$ trained with a long-tailed distribution $p_s(y)$, recent long-tailed recognition methods [52, 1] have proposed the training strategy utilizing the multiple diverse experts, which are specialized in handling different label distributions, such as long-tailed and uniform label distributions. However, these approaches are not directly applicable to TTA setting, as they are designed for multiple-expert model architectures. Inspired by this strategy, we aim to develop the model f that is suitable for TTA by dynamically adapting the model f to diverse target label distributions $p_t(y)$. Therefore, we introduce a novel label shift

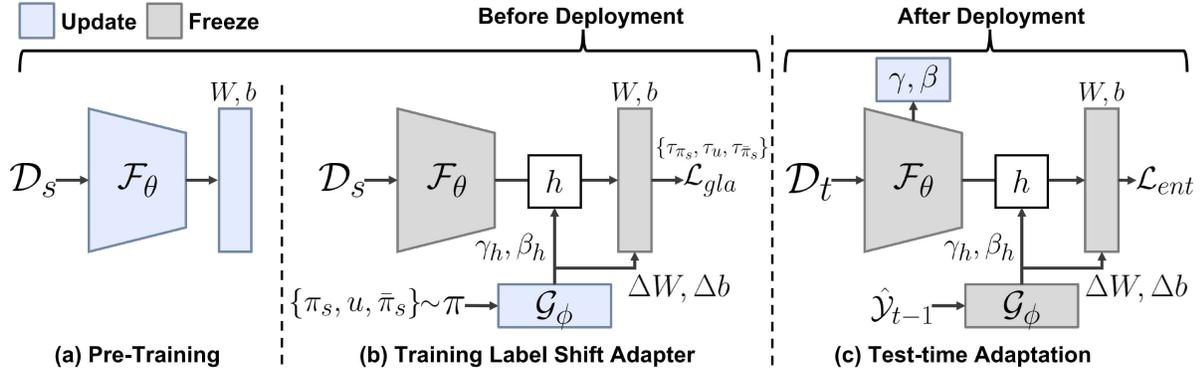


Figure 3. **Overview of the proposed method.** (a) We take a frozen pre-trained source model trained with \mathcal{D}_s . (b) Before model deployment, we train a label shift adapter with a frozen source pre-trained model. (c) After model deployment, we adapt the model to the target domain by integrating the label shift adapter into other TTA algorithms.

adapter that can produce the optimal parameters according to label distribution during inference while it is applicable to any model regardless of model architecture.

Fig. 3 shows the overview of the proposed method. Our method consists of three stages. First, our method takes the pre-trained source model in an off-the-shelf manner. Before model deployment, we train the label shift adapter with the frozen pre-trained model, which produces optimal parameters depending on the label distribution. Then, our label shift adapter can be integrated into other TTA algorithms, such as TENT [46], after model deployment. In specific, we optimize the affine parameters in the normalization layers of a feature extractor while adapting the model to the target label distribution $p_t(y)$ by estimating the label distribution of \mathcal{D}_t during inference.

Label Shift Adapter. Before training a label shift adapter, we pre-train a model $f : x \rightarrow y$ using a source domain data \mathcal{D}_s , where the model consists of a feature extractor \mathcal{F}_θ and a classifier weights $W \in \mathbb{R}^{d \times C}$, $b \in \mathbb{R}^{1 \times C}$. C and d denote the number of classes and channel of the output $h \in \mathbb{R}^{1 \times d}$ of the feature extractor, respectively. As several TTA methods [7, 56, 24] include an additional stage for training extra components, the label shift adapter is trained with the frozen pre-trained model before model deployment.

The label shift adapter \mathcal{G}_ϕ receives the label distribution $\pi \in \mathbb{R}^C$ as conditional input. For applicability and efficiency, we design the label shift adapter to generate the parameters for a part of the model. With π , the label shift adapter \mathcal{G}_ϕ predicts affine parameters $\gamma_h \in \mathbb{R}^{1 \times d}$, $\beta_h \in \mathbb{R}^{1 \times d}$ and the weight difference $\Delta W \in \mathbb{R}^{d \times C}$, $\Delta b \in \mathbb{R}^{1 \times C}$ for the classifier weights W, b . The affine parameters γ_h and β_h are applied to the hidden feature map h , which is the output of the feature extractor: $h = \mathcal{F}_\theta(x)$. Then, we compute the output \hat{y} using $W + \Delta W$ and $b + \Delta b$. Formally, \hat{y} is computed in the classifier layer as follows:

$$\hat{y} = (\gamma_h h + \beta_h) \cdot (W + \Delta W) + (b + \Delta b). \quad (3)$$

Objective Function. The label shift adapter aims to cre-

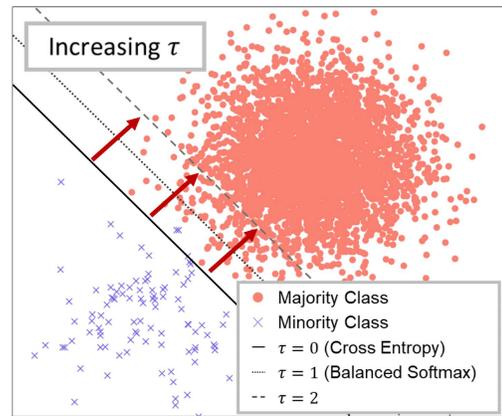


Figure 4. Visualization for understanding τ in a generalized logit adjusted loss. By adjusting τ , we can control the bias during training. As τ increases, the decision boundary

ate the optimal parameters depending on the label distribution π . During training the label shift adapter, we sample π based on π_s , which is the label distribution of \mathcal{D}_s .

To optimize the label shift adapter, we employ a generalized logit adjusted loss [33, 1], which incorporates a controlled bias during training:

$$\mathcal{L}_{gla} = - \sum_{(x_i, y_i) \sim \mathcal{D}_s} y_i \log \sigma(\hat{y}_i + \tau \log \pi_s), \quad (4)$$

where σ denotes a softmax function. \hat{y}_i indicates the output logits before computing the softmax function. $\tau \in \mathbb{R}^1$ is a scalar value for controlling the bias towards different parts of the label distribution. We sample π at each iteration during the label shift adapter training. Specifically, the label distribution π is sampled from three types of label distribution $\{\pi_s, u, \bar{\pi}_s\}$, where $\bar{\pi}_s$ indicate inverse label distribution that is obtained by inverting the order of training label distribution π_s . u denotes uniform label distribution. We select the appropriate τ for sampled label distribution π , where we set the hyperparameter τ matching each $\pi \subset \{\pi_s, u, \bar{\pi}_s\}$. For example, if π_s is sampled for π , τ is set to 0, resulting in

the use of a cross-entropy loss. On the other hand, τ is set to 1 and 2 for u and $\bar{\pi}_s$, respectively, which correspond to the balanced softmax [39] and inverse softmax [52] that simulate uniform and inverse label distributions. Intuitively, as τ increases, the decision boundary moves away from minority classes and towards the majority classes [33, 1], as shown in Fig. 4. Technically, it is possible to sample π in continuous label space instead of discrete label distributions. However, we found that sampling three distinctive label distributions is empirically sufficient to train the label shift adapter.

Test-Time Adaptation. After model deployment, we adapt the pre-trained model f using test samples x_t at t -th inference step. In specific, our method updates affine parameters γ, β in normalization layers of feature extractor F_θ . With t -th test samples x_t , the model minimizes the prediction entropy, following the previous work [46].

However, minimizing the entropy of the predictions \hat{y}_t is insufficient to adapt the model when label shift occurs. Therefore, the label shift adapter creates a part of parameters (i.e., $\gamma_h, \beta_h, \Delta W$, and Δb) of f depending on estimated label distribution $\hat{\mathcal{Y}}$. To estimate the label distribution $\hat{\mathcal{Y}}$ of D_t , we employ an exponential moving average. To be specific, the estimated target label distribution $\hat{\mathcal{Y}}_t$ at t -th step is updated recursively:

$$\hat{\mathcal{Y}}_t = \begin{cases} u, & \text{if } t = 0 \\ \alpha \bar{y}_t + (1 - \alpha) \hat{\mathcal{Y}}_{t-1}, & \text{if } t > 0 \end{cases} \quad (5)$$

where $\alpha \in [0, 1]$ denotes the momentum hyper-parameter, $\bar{y}_t = \frac{1}{n_t} \sum_{i=1}^{n_t} \hat{y}_t^i$ is the average model prediction on test samples x_t at the t -th step. We initialize $\hat{\mathcal{Y}}$ as $u \in \mathbb{R}^C$, which is a uniform distribution vector: $u[c] = \frac{1}{C}, c = 1, 2, \dots, C$. Based on estimated target label distribution $\hat{\mathcal{Y}}_{t-1}$, the label shift adapter produces the optimal parameters. With $\hat{\mathcal{Y}}_{t-1}$, we can adapt the model to the target domain using the refined entropy minimization as follows:

$$\mathcal{L}_{ent} = - \sum_{x \sim p_t(x)} f(x; \mathcal{G}_\phi(\hat{\mathcal{Y}})) \log f(x; \mathcal{G}_\phi(\hat{\mathcal{Y}})), \quad (6)$$

where the estimated label distribution $\hat{\mathcal{Y}}$ is fed into the label shift adapter \mathcal{G}_ϕ to handle label distribution shift, and the generated parameters from \mathcal{G}_ϕ adjusts the model f as shown in Eq. 3. Intuitively, while our strategy adapts the model to the target label distribution, entropy minimization boosts the confidence of correct classes.

Techniques for Label Shift Adapter. In long-tailed recognition literature [18, 2], it is known that the classifier layer plays an essential role in resolving the label shifts. Based on this intuition, we have designed the label shift adapter to produce the parameters only for the parts associated with the classifier layer. Specifically, ΔW and Δb is the weight difference of W and b , and γ_h and β_h shift the feature vector h properly. By predicting only a small portion of the

model instead of predicting the entire model weights, our label shift adapter offers various benefits. Our label shift adapter is readily applicable to any pre-trained models, regardless of the model architecture. Moreover, this strategy is computationally efficient, as the label shift adapter requires negligible extra computational costs.

We discovered that it is more effective to utilize a mapping vector $m \in \mathbb{R}^C$ to make the label distribution π a scalar, instead of directly using label distribution π as the input of the label shift adapter. In specific, $m^\top \pi \in \mathbb{R}^1$ is fed into label shift adapter, instead of π . Here, the mapping vector is the class-wise coefficients that increase in proportion to the order of class frequency in the training set. Empirically, this strategy makes the label shift adapter training more stable. Instead of using a complex label space as a condition, this technique feeds the degree of imbalance to label shift adapter, which can easily interpret the condition.

Rationale behind Label Shift Adapter. Label shift adapter plays a crucial role in handling label shifts by adjusting the parameters based on the estimated label distribution at test time. Since the label distribution in the target domain is unknown, it is possible to train the label shift adapter by simulating various label distributions using the source domain dataset. By virtue of this process, the label shift adapter can learn to produce appropriate parameters based on the label distribution. Since the label distribution in the target domain is unknown, it is possible to train the adapter by simulating various label distributions using the source domain dataset. Moreover, this approach also allows for the effective handling of biases caused by the label distribution of the source domain dataset. This is why transferring the model from the source domain to the target domain proves to be an effective approach.

4. Experiments

4.1. Experiment Setup

We evaluate the effectiveness of our proposed method on six datasets widely used in domain adaptation literature. In contrast to the traditional TTA setting, we utilize imbalanced versions of the datasets for training and evaluation.

CIFAR-10&100 and ImageNet [6, 14]. For the training set of the source domain, we utilize CIFAR-10-LT, CIFAR-100-LT [6], and ImageNet-LT [30], which are the long-tailed version of CIFAR-10, CIFAR-100 and ImageNet, following the protocol in long-tailed recognition work. Note that the imbalance ratio ρ is defined as $\rho = \frac{\max_i n_i}{\min_i n_i}$, where n_i denotes the number of class i samples in the dataset. In our experiments, the imbalance ratio of CIFAR-10-LT and CIFAR-100-LT is set to 100. ImageNet-LT is obtained by sampling from ImageNet using a Pareto distribution with $\alpha = 6$ [30], following the previous work. The categories in the training set of ImageNet-LT contain between 5 and

Method	CIFAR-10-C								CIFAR-100-C							
	Forward-LT			Uni.	Backward-LT			Avg.	Forward-LT			Uni.	Backward-LT			Avg.
	50	25	10	1	10	25	50		50	25	10	1	10	25	50	
Source	56.78	55.27	52.82	44.74	37.28	34.78	33.88	45.08	33.53	31.95	29.36	22.14	14.98	12.37	11.12	22.21
BN Stats	78.60	76.21	71.77	53.19	35.00	28.62	25.18	52.65	49.03	46.61	42.92	31.30	20.07	16.09	13.86	31.41
ONDA	77.93	75.94	72.34	55.28	37.71	31.76	28.64	54.23	48.30	46.52	42.82	31.77	20.58	16.55	14.41	31.57
PseudoLabel	79.39	76.40	71.37	47.90	30.68	24.95	22.61	50.47	50.80	48.24	43.56	25.51	17.03	14.28	11.99	30.20
LAME	58.27	55.88	52.27	41.60	34.14	32.15	31.54	43.69	32.63	30.99	28.12	20.81	13.94	11.58	10.29	21.19
CoTTA	80.45	77.86	73.12	53.75	33.89	27.80	24.01	52.98	48.32	45.71	42.30	30.42	21.65	18.14	16.33	31.84
NOTE	76.27	74.79	72.18	61.98	53.13	50.94	49.45	62.68	44.52	43.15	40.52	34.25	23.92	20.60	19.11	32.30
TENT	80.36	76.99	71.28	43.65	28.64	23.01	20.32	49.18	51.74	49.24	44.05	21.30	15.86	13.40	11.25	29.55
+Ours	80.39	78.03	73.35	53.91	37.85	32.83	30.32	55.24	52.43	50.17	46.07	33.27	21.23	17.12	15.13	33.63
	+0.03	+1.04	+2.06	+10.25	+9.21	+9.82	+10.00	+6.06	+0.69	+0.93	+2.02	+11.98	+5.37	+3.72	+3.87	+4.08
IABN	76.23	74.84	72.22	62.22	53.29	50.88	49.69	62.77	44.79	43.24	40.63	34.01	23.95	20.67	19.16	32.35
+Ours	80.58	78.62	75.26	63.34	68.54	70.07	71.64	72.58	52.06	49.71	46.03	36.84	29.29	26.33	25.50	37.97
	+4.35	+3.78	+3.04	+1.12	+15.24	+19.20	+21.95	+9.81	+7.26	+6.47	+5.40	+2.83	+5.34	+5.67	+6.33	+5.62

Table 1. **Comparison of accuracy on CIFAR-10-C and CIFAR-100-C.** The source model is trained with CIFAR-10-LT and CIFAR-100-LT. We report the average accuracy of 15 corruption types on various test label distributions. Uni. indicates the uniform distribution. Numbers under Forward-LT and Backward-LT denote the imbalance ratio. We integrate our method into TENT [46] and IABN [11].

Method	Forward-LT			Uni.	Backward-LT			Avg.
	50	25	10	1	10	25	50	
Source	26.15	25.64	24.67	21.46	18.28	17.07	16.56	21.40
BN Stats	39.47	38.89	37.71	33.63	29.48	28.07	27.20	33.49
ONDA	39.45	38.83	37.71	33.56	29.33	28.01	26.96	33.41
PseudoLabel	41.46	40.78	39.31	33.49	29.79	28.36	27.67	34.41
LAME	26.08	25.57	24.58	21.37	18.20	17.01	16.48	21.33
CoTTA	40.22	39.81	39.10	35.40	30.21	28.72	27.65	34.44
NOTE	42.43	41.65	40.36	35.17	30.99	29.14	28.17	35.41
TENT	39.40	38.73	37.27	29.05	29.31	28.26	27.28	32.76
+Ours	44.52	43.03	40.86	34.18	31.32	31.21	31.28	36.63
	+5.12	+4.30	+3.59	+5.14	+2.01	+2.94	+3.99	+3.87
IABN	42.44	41.69	40.39	35.20	31.02	29.20	28.22	35.45
+Ours	46.88	45.16	42.68	35.72	33.18	32.91	33.17	38.53
	+4.44	+3.47	+2.29	+0.52	+2.16	+3.71	+4.95	+3.08

Table 2. **Comparison of accuracy on ImageNet-C.** We report the average accuracy of 15 corruption types on various test label distributions. Uni. indicates the uniform distribution. Numbers under Forward-LT and Backward-LT denote the imbalance ratio.

1280 samples, with an imbalanced ratio set to 256. For evaluation, we utilize three corrupted test sets: CIFAR-10-C, CIFAR-100-C, and ImageNet-C [14], which consist of 15 corruption types at five severity levels. The severity level is set to 5 and 3 for CIFAR-C and ImageNet-C, respectively.

Following the test-agnostic long-tailed recognition setting [15, 52], the models are evaluated on multiple subsets of test datasets that follow different label distributions. We construct three types of test label distributions as follows: (i) Forward distribution: as the imbalance ratio increases, it becomes similar to the training label distribution. (ii) Uniform distribution: a class-balanced test dataset. (iii) Backward distribution: the order of classes is reversed, causing it to deviate more from the training distribution, as the imbalance ratio increases. Note that the degree of label shifts increases from Forward to Backward.

VisDA-C [38]. VisDA-C is a challenging large-scale bench-

mark whose training data is synthesized through 3D model rendering, and its test data is sampled from the real world. The dataset contains 12 categories. We utilize an imbalanced dataset VisDA-C (RSUT), where source and target domains are subject to two reverse Pareto distributions, following the previous work [43]. Here, RSUT denotes the combination of Reversely-unbalanced Source (RS) and Unbalanced Target (UT) distribution. The label distributions of RSUT are described in the supplementary.

OfficeHome [44]. This dataset comprises four domains, each consisting of 65 categories. We also employ OfficeHome (RSUT) [43], which is created by the same protocol as VisDA-C (RSUT). Since the artistic domain in OfficeHome is too small to sample an imbalanced subset, we only utilize the remaining three distinct domains (*e.g.*, Clip Art, Product, and Real-World), as prior work [22].

DomainNet [37]. We employ a subset of DomainNet [43],

Method	VISDA-C	Method	C→P	C→R	P→C	P→R	R→C	R→P	Avg.
Source	51.45	Source	45.39	44.53	32.94	64.33	40.22	68.92	49.39
BN Stats	49.33	BN Stats	44.30	48.27	35.63	62.17	40.73	62.20	48.88
ONDA	50.68	ONDA	44.84	47.57	35.20	62.09	40.61	63.83	49.02
PseudoLabel	49.50	PseudoLabel	47.98	49.34	37.71	62.42	39.38	63.21	50.01
LAME	50.72	LAME	41.68	42.27	32.40	63.57	37.92	66.94	47.46
CoTTA	49.88	CoTTA	44.46	48.19	35.63	62.34	40.73	62.20	48.92
NOTE	49.37	NOTE	43.02	42.38	38.64	61.69	41.40	64.33	48.58
TENT	48.68	TENT	49.60	49.51	38.96	63.08	41.25	64.52	51.15
+Ours	72.97	+Ours	49.60	53.13	37.81	66.45	41.35	68.35	52.78
	+24.29		0.00	+3.62	-1.15	+3.37	+0.10	+3.83	+1.63

Table 3. Comparison of accuracy on VisDA-C (RSUT).

Table 4. Comparison of accuracy on three domains of Officehome (RSUT): C: Clipart, P: Product, R: Realworld.

Method	C→P	C→R	C→S	P→C	P→R	P→S	R→C	R→P	R→S	S→C	S→P	S→R	Avg.
Source	52.73	74.87	52.15	58.42	81.22	61.82	66.03	69.58	55.31	63.92	59.68	75.43	64.26
BN Stats	56.81	77.05	54.10	63.63	81.12	60.22	67.38	70.00	56.84	71.75	68.72	80.18	67.32
ONDA	56.82	78.32	54.81	63.99	81.79	61.86	67.14	70.09	58.11	71.60	69.34	80.77	67.89
PseudoLabel	61.81	77.43	56.25	62.56	81.64	62.04	71.06	73.89	58.49	71.81	70.38	80.12	68.96
LAME	49.20	72.45	48.69	57.81	80.09	60.85	65.25	68.19	53.97	61.00	55.66	73.25	62.20
CoTTA	56.88	77.33	54.18	63.69	81.31	60.26	67.44	70.07	57.14	71.69	68.85	80.56	67.45
NOTE	55.38	74.15	57.98	65.59	81.66	64.65	71.29	73.32	63.28	72.28	68.31	80.25	69.01
TENT	63.26	77.10	59.76	66.69	80.02	64.32	71.88	74.34	62.25	73.13	72.64	78.73	70.34
+Ours	63.26	81.11	60.39	67.38	82.99	67.23	71.88	74.83	64.40	71.88	71.56	82.67	71.63
	0.00	+4.01	+0.63	+0.69	+2.97	+2.91	0.00	+0.49	+2.15	-1.25	-1.08	+3.94	+1.29

Table 5. Comparison of accuracy on four domains of DomainNet: C: Clipart, P: Painting, R: Real, S: Sketch.

comprising 40 categories from four domains: Real, Clipart, Painting, and Sketch. As the label shift between these domains is inherent, we did not need to modify the label distribution of the dataset. The visualization of label distribution in DomainNet are illustrated in the supplementary.

Baseline Methods. Note that we utilize the pre-trained model trained using *balanced softmax* [39], which is a widely used long-tailed recognition approach. **Source** indicates the pre-trained model with the source data using balanced softmax. We compare our method with the following TTA baselines: BN stats [41], ONDA [32], PseudoLabel [20], LAME [4], CoTTA [47], TENT [46], IABN [11], and NOTE [11]. Note that IABN is a normalization layer introduced in the NOTE paper. Since NOTE has been proposed for temporally correlated test samples, IABN layer has the capability to handle the class-imbalance in a batch.

Implementation Details. We utilize ResNet-18 [13] as the backbone for CIFAR-10 and CIFAR-100, and ResNeXt [50] for ImageNet. We also employ ResNet-50 for OfficeHome and DomainNet, and ResNet-101 for VisDA-C, which are pre-trained on ImageNet. For fair comparisons, the same architecture and optimizer are utilized for all TTA baselines. In all experiments, the source domain model is trained using Balanced softmax [39], a representative long-tailed recognition method. During inference, the batch size is set to 64 in all experiments. Further details

Dataset	Prior	Forward-LT			Uni.	Backward-LT			Avg.
		50	25	10	1	10	25	50	
CIFAR-10-C	✗	80.58	78.62	75.26	63.34	68.54	70.07	71.64	72.58
	✓	81.58	79.05	75.17	69.55	68.65	70.51	72.86	73.91
CIFAR-100-C	✗	52.06	49.71	46.03	36.84	29.29	26.33	25.50	37.97
	✓	53.59	50.95	46.96	37.17	29.79	27.46	27.03	38.99
ImageNet-C	✗	46.88	45.16	42.68	35.72	33.18	32.91	33.17	38.53
	✓	47.71	45.69	42.98	36.02	32.95	32.65	33.23	38.75

Table 6. Comparison between estimated label distribution and target prior on CIFAR-10-C, CIFAR-100-C, and ImageNet-C. Prior indicates that the true target label distribution is utilized as prior knowledge for label shift adapter, instead of estimated label distribution. We conduct the experiments using IABN+Ours [11].

regarding the hyperparameters for each baseline and our method are described in the supplementary.

4.2. Results on Corruption Data

Table 1 reports the average accuracy on 15 corruption types in CIFAR-10-C and CIFAR-100-C. The results on each target domain are presented in the supplementary. The results demonstrate that the performance of previous TTA methods is significantly inferior in the backward long-tailed distributions compared to the forward settings. This issue arises because TTA models learn based on model predictions, which are biased toward the majority classes.

In contrast, it is noteworthy that our strategy consider-

Dataset	Method	Forward-LT			Uni.	Backward-LT			Avg.
		50	25	10	1	10	25	50	
CIFAR-10-C	Logit Adjust	78.89	77.20	73.58	58.86	46.12	41.95	39.72	59.47
	IM Loss	76.54	75.39	72.34	56.43	49.79	47.58	46.62	60.67
	Ours+IABN	80.58	78.62	75.26	63.34	68.54	70.07	71.64	72.58
CIFAR-100-C	Logit Adjust	43.23	41.16	37.55	30.94	25.44	23.67	23.38	32.19
	IM Loss	44.10	42.64	40.00	33.36	23.61	20.54	19.11	31.91
	Ours+IABN	52.06	49.71	46.03	36.84	29.29	26.33	25.50	37.97

Table 7. Comparison with logit adjustment [33, 15] and information maximization (IM) loss [23].

ably improves the accuracy when label shift is severe. As our method produces the optimal parameters depending on the target label distribution, the models can be adapted to the target domain stably, even in the presence of severe label shifts. Furthermore, when integrated with IABN [11], a normalization layer for addressing the class imbalance in a batch, our method yields the best performance. Nevertheless, it should be noted that relying solely on IABN may not be sufficient in managing severe label shifts.

As shown in Table 2, our method consistently outperforms the existing TTA approaches on ImageNet-C, a more challenging dataset. Integrating our proposed method consistently improves the performance in all test sets. Furthermore, our method also improves the accuracy on the uniform dataset: TENT (+5.14%) and IABN (+0.52%). The promising results demonstrate the practicality and effectiveness of our method under covariate and label shifts.

4.3. Results on DA Benchmarks

We evaluate the effectiveness of our method in comparison with TTA methods on three domain adaptation benchmarks: VisDA-C (RSUT), OfficeHome (RSUT), and DomainNet. In VisDA-C (RSUT) and OfficeHome (RSUT), we utilize the test datasets, including the reversed Pareto label distribution. Surprisingly, the existing TTA methods perform worse than the source pre-trained model for evaluation data of VisDA-C, as shown in Table 3. This result indicates that the baselines do not work when the label shift is severe. In contrast, TENT combined with our method improves the performance significantly for VisDA-C test data.

The results on OfficeHome (RSUT) are reported in Table 4. We discovered that since the number of test samples in OfficeHome (RSUT) is limited (e.g., Clipart: 1,017, Product: 1,985, Real: 1,235), TTA methods generally do not result in significant performance improvements compared to the source model. Nonetheless, our method exhibits a general improvement on the OfficeHome (RSUT) datasets, except for P→C.

Table 5 shows the results on the DomainNet dataset, where the label shifts between different domains already exist. This result demonstrates that our approach generally outperforms the baselines. Moreover, we confirmed that our model performs better when adapting to Real domain that

Method	MACs	Params
ResNet-18	557.93M	11.22M
+ Ours	558.04M	11.34M

Table 8. Computational costs. We measure MACs and the number of parameters (Params.).

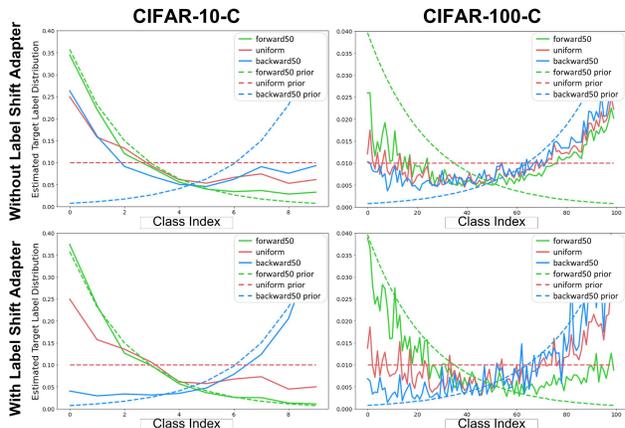


Figure 5. Visualization of estimated target label distributions and target priors on CIFAR-10-C and CIFAR-100-C. Applying the label shift adapter improves the target label distribution estimation.

contains a larger number of test samples compared to other domains (e.g., Real: 6,943, Clipart: 1,616, Painting: 2,909, Sketch: 2,399). It is because the target label distribution can be estimated more precisely with more test samples.

4.4. Analysis on Label Shift Adapter

Estimated Target Label Distribution. The proposed label shift adapter predicts the weight difference with the estimated target label distribution. Therefore, it is essential to estimate the target class prior accurately, which is the label distribution of target domain data. Different from previous TTA methods, our method adapts the pre-trained model according to the label distribution shift. As shown in Fig. 5, we visualize the estimated target label distribution of the model with and without our label shift adapter. This result shows that the model with the label shift adapter is superior for estimating various label distributions compared to the model without the label shift adapter. In addition, our method estimates the target label distribution similar to the target prior in CIFAR-10-C and CIFAR-100-C datasets.

As shown in Table 6, we compare the performances of the label shift adapter with estimated target label distribution and target prior. Obviously, the overall accuracy of the model increases if the target prior is utilized instead of the estimated target label distribution. Nevertheless, our

method shows comparable performances to the model using target prior as the input of label shift adapter.

Computational Costs. Table 8 shows the computational costs of our method. We measure the computational cost of ResNet-18, which is utilized for CIFAR-100-C. Since we design the architecture of the label shift adapter efficiently, our method requires a negligible amount of additional computational costs (MACs: +0.11M, Params: +0.12M). We describe the details of the label shift adapter architecture in the supplementary. Moreover, we report ablation study on the architecture of label shift adapter in the supplementary.

Effectiveness of Label Shift Adapter. We validate the effectiveness of our method over existing approaches for handling label shifts. Specifically, we apply the following non-trivial baselines with entropy minimization loss and IABN [11]: (i) Post-hoc logit adjustment [33, 15] modifies the logits using the estimated target label distribution. (ii) Information maximization loss [23] makes the outputs globally diverse by maximizing mean entropy, which can reduce the bias toward certain classes. Table 7 shows that these baselines are not sufficient to handle the label distribution shifts, particularly on inversely long-tailed distribution (e.g., Backward). In contrast, our approach leads to promising performance gains on various label distributions.

5. Discussions

Normalization Layers for Label Shifts. This paper introduces the label shift adapter for addressing the label shift problems in the test-time adaptation scenario, which can be applied to any model regardless of its architecture. In addition to our method, we found that managing the bias in batch statistics is also crucial for reducing performance degradation caused by label shifts. This observation is evidenced by the effectiveness of normalization layers, such as instance-aware batch normalization (IABN) layers, in dealing with the class imbalance in a batch. However, IABN layers are highly sensitive to hyperparameter selections such as the soft-shrinkage width α . Consequently, improving the robustness of normalization layers to label shifts is a promising future research direction.

Training Additional Component. One limitation of our work is that our method requires an additional stage for training the label shift adapter. In several recent test-time adaptation methods [7, 56, 24], an additional training process is often carried out before server deployment to train the additional components. Despite the drawback of requiring an additional training stage, we believe that these test-time adaptation models, including our method, are practically useful because they are more robust during inference. Furthermore, there are techniques [18, 6, 49, 48, 2] in the long-tailed recognition field that involves dividing the training process into two stages. Our approach can serve as an inspiration for methods that can be applied for handling

test-agnostic label distributions in long-tailed recognition, regardless of the model architecture.

6. Conclusion

This paper addresses the label shift problem in TTA, where both source and target domain datasets are class imbalanced. Existing TTA methods employing entropy minimization are often flawed due to the model bias toward the majority classes in source data. To address such a non-trivial issue, we propose a novel label shift adapter, which produces the optimal parameters according to the label distribution. Our label shift adapter is applicable to existing TTA methods regardless of the model architectures. Furthermore, we estimate the label distribution of target domain data to feed into the label shift adapter. Through extensive experiments, we demonstrate that our method outperforms the state-of-the-art TTA baselines. We believe that our work inspires future researchers to improve TTA methods under the joint presence of covariate and label shifts.

Acknowledgements. We would like to thank Kyuwoong Hwang, Simyung Chang, Hyunsin Park, Janghoon Cho, Juntae Lee, Hyoungwoo Park, Seokeon Choi, and Jungsoo Lee of Qualcomm AI Research team for their valuable discussions.

References

- [1] Emanuel Sanchez Aimar, Arvi Jonnarth, Michael Felsberg, and Marco Kuhlmann. Balanced product of experts for long-tailed recognition. *arXiv preprint arXiv:2206.05260*, 2022. 3, 4, 5, 15
- [2] Shaden Alshammari, Yu-Xiong Wang, Deva Ramanan, and Shu Kong. Long-tailed recognition via weight balancing. In *Proc. of the IEEE conference on computer vision and pattern recognition (CVPR)*, pages 6897–6907, 2022. 5, 9
- [3] Kamyar Azizzadenesheli, Anqi Liu, Fanny Yang, and Animeshree Anandkumar. Regularized learning for domain adaptation under label shifts. In *Proc. the International Conference on Learning Representations (ICLR)*, 2019. 2
- [4] Malik Boudiaf, Romain Mueller, Ismail Ben Ayed, and Luca Bertinetto. Parameter-free online test-time adaptation. In *Proc. of the IEEE conference on computer vision and pattern recognition (CVPR)*, pages 8344–8353, 2022. 3, 7, 13
- [5] Dhanajit Brahma, Vinay Kumar Verma, and Piyush Rai. Hypernetworks for continual semi-supervised learning. *arXiv preprint arXiv:2110.01856*, 2021. 3
- [6] Kaidi Cao, Colin Wei, Adrien Gaidon, Nikos Arechiga, and Tengyu Ma. Learning imbalanced datasets with label-distribution-aware margin loss. In *Proc. the Advances in Neural Information Processing Systems (NeurIPS)*, pages 1567–1578, 2019. 1, 3, 5, 9
- [7] Sungha Choi, Seunghan Yang, Seokeon Choi, and Sung-rack Yun. Improving test-time adaptation via shift-agnostic weight regularization and nearest source prototypes. In *Proc. of the European Conference on Computer Vision (ECCV)*, pages 440–458. Springer, 2022. 2, 3, 4, 9, 15

- [8] Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie. Class-balanced loss based on effective number of samples. In *Proc. of the IEEE conference on computer vision and pattern recognition (CVPR)*, pages 9268–9277, 2019. [1](#), [3](#)
- [9] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In *Proc. the International Conference on Machine Learning (ICML)*, pages 1180–1189. PMLR, 2015. [1](#), [2](#)
- [10] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *The journal of machine learning research*, 17(1):2096–2030, 2016. [1](#), [2](#)
- [11] Taesik Gong, Jongheon Jeong, Taewon Kim, Yewon Kim, Jinwoo Shin, and Sung-Ju Lee. Note: Robust continual test-time adaptation against temporal correlation. In *Proc. the Advances in Neural Information Processing Systems (NeurIPS)*, 2022. [1](#), [2](#), [3](#), [6](#), [7](#), [8](#), [9](#), [13](#), [15](#)
- [12] David Ha, Andrew Dai, and Quoc V Le. Hypernetworks. *arXiv preprint arXiv:1609.09106*, 2016. [3](#)
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proc. of the IEEE conference on computer vision and pattern recognition (CVPR)*, pages 770–778, 2016. [7](#)
- [14] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. In *Proc. the International Conference on Learning Representations (ICLR)*, 2018. [2](#), [5](#), [6](#)
- [15] Youngkyu Hong, Seungju Han, Kwanghee Choi, Seokjun Seo, Beomsu Kim, and Buru Chang. Disentangling label distribution for long-tailed visual recognition. In *Proc. of the IEEE conference on computer vision and pattern recognition (CVPR)*, pages 6626–6636, 2021. [1](#), [3](#), [6](#), [8](#), [9](#)
- [16] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. [3](#)
- [17] Xiang Jiang, Qicheng Lao, Stan Matwin, and Mohammad Havaei. Implicit class-conditioned domain alignment for unsupervised domain adaptation. In *International Conference on Machine Learning*, pages 4816–4827. PMLR, 2020. [2](#), [3](#)
- [18] Bingyi Kang, Saining Xie, Marcus Rohrbach, Zhicheng Yan, Albert Gordo, Jiashi Feng, and Yannis Kalantidis. Decoupling representation and classifier for long-tailed recognition. In *Proc. the International Conference on Learning Representations (ICLR)*, 2019. [1](#), [5](#), [9](#), [15](#)
- [19] Vinod K Kurmi, Venkatesh K Subramanian, and Vinay P Nambodiri. Domain impression: A source data free domain adaptation method. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 615–625, 2021. [1](#), [2](#)
- [20] 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. [1](#), [7](#), [12](#)
- [21] Jungsoo Lee, Jooyeol Yun, Sunghyun Park, Yonggyu Kim, and Jaegul Choo. Improving face recognition with large age gaps by learning to distinguish children. In *British Machine Vision Conference*, 2021. [3](#)
- [22] Xinhao Li, Jingjing Li, Lei Zhu, Guoqing Wang, and Zi Huang. Imbalanced source-free domain adaptation. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 3330–3339, 2021. [2](#), [3](#), [6](#), [12](#), [16](#)
- [23] Jian Liang, Dapeng Hu, and Jiashi Feng. Do we really need to access the source data? source hypothesis transfer for unsupervised domain adaptation. In *International Conference on Machine Learning*, pages 6028–6039. PMLR, 2020. [1](#), [2](#), [8](#), [9](#)
- [24] Hyesu Lim, Byeonggeun Kim, Jaegul Choo, and Sungha Choi. Ttn: A domain-shift aware batch normalization in test-time adaptation. In *Proc. the International Conference on Learning Representations (ICLR)*, 2023. [1](#), [2](#), [3](#), [4](#), [9](#)
- [25] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proc. of the IEEE international conference on computer vision (ICCV)*, pages 2980–2988, 2017. [3](#)
- [26] Xi Lin, Zhiyuan Yang, Qingfu Zhang, and Sam Kwong. Controllable pareto multi-task learning. *arXiv preprint arXiv:2010.06313*, 2020. [3](#)
- [27] Zhaojiang Lin, Andrea Madotto, and Pascale Fung. Exploring versatile generative language model via parameter-efficient transfer learning. *arXiv preprint arXiv:2004.03829*, 2020. [3](#)
- [28] Zachary Lipton, Yu-Xiang Wang, and Alexander Smola. Detecting and correcting for label shift with black box predictors. In *Proc. the International Conference on Machine Learning (ICML)*, pages 3122–3130. PMLR, 2018. [2](#)
- [29] Xiaofeng Liu, Zhenhua Guo, Site Li, Fangxu Xing, Jane You, C-C Jay Kuo, Georges El Fakhri, and Jonghye Woo. Adversarial unsupervised domain adaptation with conditional and label shift: Infer, align and iterate. In *Proc. of the IEEE international conference on computer vision (ICCV)*, pages 10367–10376, 2021. [2](#)
- [30] Ziwei Liu, Zhongqi Miao, Xiaohang Zhan, Jiayun Wang, Boqing Gong, and Stella X Yu. Large-scale long-tailed recognition in an open world. In *Proc. of the IEEE conference on computer vision and pattern recognition (CVPR)*, pages 2537–2546, 2019. [5](#)
- [31] Rabeeh Karimi Mahabadi, Sebastian Ruder, Mostafa Dehghani, and James Henderson. Parameter-efficient multi-task fine-tuning for transformers via shared hypernetworks. *arXiv preprint arXiv:2106.04489*, 2021. [3](#)
- [32] Massimiliano Mancini, Hakan Karaoguz, Elisa Ricci, Patric Jensfelt, and Barbara Caputo. Kitting in the wild through online domain adaptation. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1103–1109. IEEE, 2018. [7](#), [13](#)
- [33] Aditya Krishna Menon, Sadeep Jayasumana, Ankit Singh Rawat, Himanshu Jain, Andreas Veit, and Sanjiv Kumar. Long-tail learning via logit adjustment. In *Proc. the International Conference on Learning Representations (ICLR)*, 2020. [3](#), [4](#), [5](#), [8](#), [9](#), [15](#)
- [34] Aviv Navon, Aviv Shamsian, Ethan Fetaya, and Gal Chechik. Learning the pareto front with hypernetworks. In *Proc.*

- the International Conference on Learning Representations (ICLR)*, 2020. 3
- [35] Shuaicheng Niu, Jiaxiang Wu, Yifan Zhang, Yaofo Chen, Shijian Zheng, Peilin Zhao, and Mingkui Tan. Efficient test-time model adaptation without forgetting. In *Proc. the International Conference on Machine Learning (ICML)*, pages 16888–16905. PMLR, 2022. 1, 2, 3
- [36] Shuaicheng Niu, Jiaxiang Wu, Yifan Zhang, Zhiquan Wen, Yaofo Chen, Peilin Zhao, and Mingkui Tan. Towards stable test-time adaptation in dynamic wild world. 2023. 15, 16
- [37] Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment matching for multi-source domain adaptation. In *Proc. of the IEEE international conference on computer vision (ICCV)*, pages 1406–1415, 2019. 6
- [38] Xingchao Peng, Ben Usman, Neela Kaushik, Judy Hoffman, Dequan Wang, and Kate Saenko. Visda: The visual domain adaptation challenge. *arXiv preprint arXiv:1710.06924*, 2017. 6
- [39] Jiawei Ren, Cunjun Yu, Shunan Sheng, Xiao Ma, Haiyu Zhao, Shuai Yi, and Hongsheng Li. Balanced meta-softmax for long-tailed visual recognition. In *Proc. the Advances in Neural Information Processing Systems (NeurIPS)*, pages 4175–4186, 2020. 2, 3, 5, 7, 12
- [40] Dvir Samuel and Gal Chechik. Distributional robustness loss for long-tail learning. In *Proc. of the IEEE international conference on computer vision (ICCV)*, 2021. 3
- [41] Steffen Schneider, Evgenia Rusak, Luisa Eck, Oliver Bringmann, Wieland Brendel, and Matthias Bethge. Improving robustness against common corruptions by covariate shift adaptation. *Proc. the Advances in Neural Information Processing Systems (NeurIPS)*, 33:11539–11551, 2020. 1, 3, 7, 12
- [42] Aviv Shamsian, Aviv Navon, Ethan Fetaya, and Gal Chechik. Personalized federated learning using hypernetworks. In *Proc. the International Conference on Machine Learning (ICML)*, pages 9489–9502. PMLR, 2021. 3
- [43] Shuhan Tan, Xingchao Peng, and Kate Saenko. Class-imbalanced domain adaptation: an empirical odyssey. In *Computer Vision—ECCV 2020 Workshops: Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*, pages 585–602. Springer, 2020. 2, 3, 6, 12
- [44] Hemanth Venkateswara, Jose Eusebio, Shayok Chakraborty, and Sethuraman Panchanathan. Deep hashing network for unsupervised domain adaptation. In *Proc. of the IEEE conference on computer vision and pattern recognition (CVPR)*, pages 5018–5027, 2017. 6
- [45] Johannes von Oswald, Christian Henning, Benjamin F Grewe, and João Sacramento. Continual learning with hypernetworks. In *Proc. the International Conference on Learning Representations (ICLR)*. Proc. the International Conference on Learning Representations (ICLR), 2020. 3
- [46] Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, and Trevor Darrell. Tent: Fully test-time adaptation by entropy minimization. In *Proc. the International Conference on Learning Representations (ICLR)*, 2021. 1, 2, 3, 4, 5, 6, 7, 12, 13, 15
- [47] Qin Wang, Olga Fink, Luc Van Gool, and Dengxin Dai. Continual test-time domain adaptation. In *Proc. of the IEEE conference on computer vision and pattern recognition (CVPR)*, pages 7201–7211, 2022. 7, 13
- [48] Xudong Wang, Long Lian, Zhongqi Miao, Ziwei Liu, and Stella X Yu. Long-tailed recognition by routing diverse distribution-aware experts. In *Proc. the International Conference on Learning Representations (ICLR)*, 2020. 9
- [49] Liuyu Xiang, Guiguang Ding, and Jungong Han. Learning from multiple experts: Self-paced knowledge distillation for long-tailed classification. In *Proc. of the European Conference on Computer Vision (ECCV)*, pages 247–263, 2020. 3, 9
- [50] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *Proc. of the IEEE conference on computer vision and pattern recognition (CVPR)*, pages 1492–1500, 2017. 7
- [51] Jayeon Yoo, Inseop Chung, and Nojun Kwak. Unsupervised domain adaptation for one-stage object detector using offsets to bounding box. In *Proc. of the European Conference on Computer Vision (ECCV)*, pages 691–708. Springer, 2022. 2
- [52] Yifan Zhang, Bryan Hooi, HONG Lanqing, and Jiashi Feng. Self-supervised aggregation of diverse experts for test-agnostic long-tailed recognition. In *Proc. the Advances in Neural Information Processing Systems (NeurIPS)*, 2022. 3, 5, 6
- [53] Bowen Zhao, Chen Chen, and Shu-Tao Xia. Delta: degradation-free fully test-time adaptation. In *Proc. the International Conference on Learning Representations (ICLR)*, 2023. 2, 3, 16
- [54] Boyan Zhou, Quan Cui, Xiu-Shen Wei, and Zhao-Min Chen. Bbn: Bilateral-branch network with cumulative learning for long-tailed visual recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9719–9728, 2020. 3
- [55] Yang Zou, Zhiding Yu, BVK Kumar, and Jinsong Wang. Unsupervised domain adaptation for semantic segmentation via class-balanced self-training. In *Proc. of the European Conference on Computer Vision (ECCV)*, pages 289–305, 2018. 2
- [56] Yuliang Zou, Zizhao Zhang, Chun-Liang Li, Han Zhang, Tomas Pfister, and Jia-Bin Huang. Learning instance-specific adaptation for cross-domain segmentation. In *Proc. of the European Conference on Computer Vision (ECCV)*, pages 459–476. Springer, 2022. 3, 4, 9

Supplementary Material

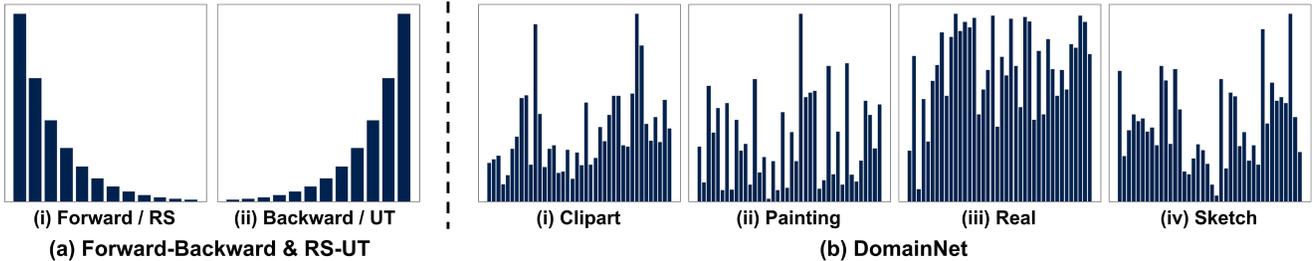


Figure 6. **Visualization of label distributions in datasets.** (a) shows the illustrations of forward or reversely-unbalanced source (RS) setting and backward or unbalanced target (UT) setting. In specific, forward and backward are used in CIFAR-10-C, CIFAR-100-C, and ImageNet-C. In addition, RS and UT are utilized in VisDA-C and OfficeHome. (b) shows the natural label shift of DomainNet.

A. Implementation Details

In this section, we introduce further information regarding the datasets, along with the implementation details for the baseline test-time-adaptation (TTA) methods and the label shift adapter.

A.1. Datasets

Fig. 6 illustrates the label distributions for the datasets utilized in our experiments. As depicted in Fig. 6 (a), ‘forward’ and ‘RS’ represent long-tailed label distributions, with class order corresponding to the training label distribution. Conversely, ‘backward’ and ‘UT’ indicate a reversed class order.

In the forward and backward settings, the imbalance ratios for CIFAR-10-C, CIFAR-100-C, and ImageNet-C are configured to 10, 25, and 100. We adjust the label distribution by reducing the number of images per class based on the specified imbalance ratio. For VisDA-C, The imbalance ratio is set to 100 for both training and test datasets. Furthermore, we utilize an imbalanced version of OfficeHome created by the previous research [43].

Fig. 6 (b) shows the label distributions of DomainNet, in which existing label shifts are significant enough. The superior performance of our method on DomainNet demonstrates its ability to handle label shifts that arise in real-world scenarios.

A.2. Details of Baselines

We carry out the experiments using the official implementations of the baseline models. We provide additional details regarding the implementation specifics, including

hyperparameters. Note that the batch size for test-time adaptation is configured to 64 for fair comparisons. For simplicity, we present the hyperparameters in the following sequence: **{CIFAR-10-C, CIFAR-100-C, ImageNet-C, VisDA-C, OfficeHome, DomainNet}** for test-time adaptation baselines. In instances where hyperparameters are not separately described for each dataset, the same values are employed across all datasets.

Source. Different from the previous TTA studies, we employ long-tailed datasets in our research. To mitigate model bias towards the majority classes, we utilize a balanced softmax [39], which is a prominent method for long-tailed recognition. Formally, the balanced softmax is expressed as:

$$\mathcal{L}_{\text{bal}} = - \sum_{(x_i, y_i) \sim \mathcal{D}_s} y_i \log \sigma(\hat{y}_i + \log(\pi_s)),$$

where π_s represents the frequency of the training classes, and σ denotes the softmax function.

Table 9 describes the hyperparameters utilized for training on source domain datasets. We select the hyperparameters for VisDA-C, OfficeHome, and DomainNet in accordance with the imbalanced source-free domain adaptation study [22]. As described in the main manuscript, we utilize pre-trained ResNet-50 and ResNet-101 on ImageNet, when conducting the experiments on VisDA-C, OfficeHome, and DomainNet. Moreover, the learning rate for the feature extractor and the classifier is set to $0.1 \times \text{LR}$ and LR, respectively, when training the model on VisDA-C, OfficeHome, and DomainNet. All experiments are conducted using NVIDIA RTX A5000 GPU.

BN Stats. BN stats [41] utilizes test batch statistics instead of running statistics within batch normalization layers.

PseudoLabel. In accordance with previous studies [20, 46],

Src Data	Tgt Data	Model	Optim.	Scheduler	Epoch	Batch	WD	Momentum	LR
CIFAR-10-LT	CIFAR-10-C	ResNet-18	SGD	CosineAnneal	200	128	5e-4	0.9	0.1
CIFAR-100-LT	CIFAR-100-C	ResNet-18	SGD	CosineAnneal	200	128	5e-4	0.9	0.1
ImageNet-LT	ImageNet-C	ResNeXt-50	SGD	Manual	90	64	2e-4	0.9	0.1
VisDA-C (RS)	VisDA-C (UT)	ResNet-101	SGD	-	15	40	1e-3	0.9	1e-3
OfficeHome (RS)	OfficeHome (UT)	ResNet-50	SGD	-	50	40	1e-3	0.9	1e-2
DomainNet	DomainNet	ResNet-50	SGD	-	20	40	1e-3	0.9	1e-2

Table 9. **Hyperparameters for training the model with source domain data.** Src Data and Tgt Data denote source domain and target domain datasets, respectively. Optim. indicates the optimizer. WD and LR denote the weight decay and learning rate for training. The manual scheduler for ImageNet-LT is to decay the learning rate at 60 and 80 epochs.

Model Architecture				Forward-LT			Uni.	Backward-LT			Avg.
γ_h	β_h	ΔW	Δb	50	25	10	1	10	25	50	
✓				51.20	49.30	46.06	37.36	27.28	23.75	21.84	36.69
	✓			48.79	42.45	32.10	15.21	22.52	24.09	25.14	30.04
		✓		51.32	49.36	46.11	<u>37.17</u>	27.06	23.56	21.71	36.61
			✓	52.50	50.19	46.64	37.51	28.95	25.56	24.06	37.92
✓	✓			<u>52.09</u>	49.48	45.43	35.92	29.60	26.82	26.25	37.94
		✓	✓	51.93	<u>49.78</u>	<u>46.43</u>	37.18	27.71	24.28	<u>22.57</u>	37.12
✓	✓	✓	✓	52.06	49.71	46.03	36.84	<u>29.29</u>	<u>26.33</u>	<u>25.50</u>	37.97

Table 10. **Ablation study on architecture design of label shift adapter using CIFAR-100-LT and CIFAR-100-C.**

we update the affine parameters in the batch normalization layers using the hard pseudo labels. The learning rate is set to $\{1e-3, 1e-3, 2.5e-4, 5e-5, 5e-5, 1e-3\}$ for each respective dataset, following the hyperparameters of TENT [46].

ONDA. Online domain adaptation (ONDA) [32] modifies the batch normalization statistics for target domains using a batch of target data through an exponential moving average. We set the update frequency $N = 10$ and the decay of the moving average $m = 0.1$, adhering to the default values of the original paper.

TENT. Test entropy minimization (TENT) [46] optimizes the affine parameters of batch normalization layers via entropy minimization. The learning rate is configured to $\{1e-3, 1e-3, 2.5e-4, 5e-5, 5e-5, 1e-3\}$ for each dataset. We referred to the official implementation for hyperparameter selection.

LAME. Laplacian adjusted maximum-likelihood estimation (LAME) [4] alters the output probability of the classifier. Following the authors’ implementation, we set the kNN affinity matrix with the value of k as 5.

CoTTA. Continual test-time adaptation (CoTTA) [47] adapts the model to accommodate continually evolving target domains by employing a weight-averaged teacher model, data augmentations, and stochastic restoring. CoTTA incorporates three hyperparameters: augmentation confidence threshold p_{th} , restoration factor p , and the decay of EMA m . p and m are set to 0.01 and 0.999, respectively. Additionally, p_{th} is configured to $\{0.92, 0.72, 0.01, 0.01, 0.01, 0.01\}$. Given that the authors do not provide the hyperparameters for VisDA-C, OfficeHome, and DomainNet,

Algorithm 1 Training Process of Label Shift Adapter

Require: Dataset $\mathcal{D}_s = \{(x_i, y_i)\}_{i=1}^n$. A pre-trained model f . A label shift adapter \mathcal{G}_ϕ .

- 1: Initialize the parameters ϕ randomly
- 2: **for** $k = 1$ to K **do**
- 3: $\mathcal{B} \leftarrow \text{SampleMiniBatch}(\mathcal{D}, m)$
 \triangleright a mini-batch of m examples
- 4: $\pi, \tau \leftarrow \text{Sample}(\{\pi_s, u, \bar{\pi}_s\}, \{\tau_{\pi_s}, \tau_u, \tau_{\bar{\pi}_s}\})$
 \triangleright sample τ matching π
- 5: $\mathcal{L}(\mathcal{G}_\phi) \leftarrow \frac{1}{m} \sum_{(x,y) \in \mathcal{B}} \mathcal{L}_{glu}((x, y, \pi); f, \mathcal{G}_\phi)$
- 6: $\mathcal{G}_\phi \leftarrow \mathcal{G}_\phi - \eta \nabla_{\theta} \mathcal{L}(\mathcal{G}_\phi)$ \triangleright one SGD step
- 7: **end for**

we fine-tune the appropriate hyperparameters for them.

NOTE. Non-i.i.d. test-time adaptation (NOTE) [11] comprises two components: (i) Instance-aware batch normalization (IABN), and (ii) Prediction-balanced reservoir sampling (PBRS). In accordance with the original paper, we substitute the batch normalization layers with IABN layers before pre-training the source models. Two hyperparameters are associated with IABN: soft-shrinkage width α and EMA momentum m . The values of α are configured as $\{4, 4, 8, 8, 8, 8\}$, while m is set to $\{0.01, 0.01, 0.1, 0.1, 0.1, 0.1\}$. The memory size of PBRS is set to 64, equal to the batch size. In our experiments, we incorporate our label shift adapter into the models using IABN layers.

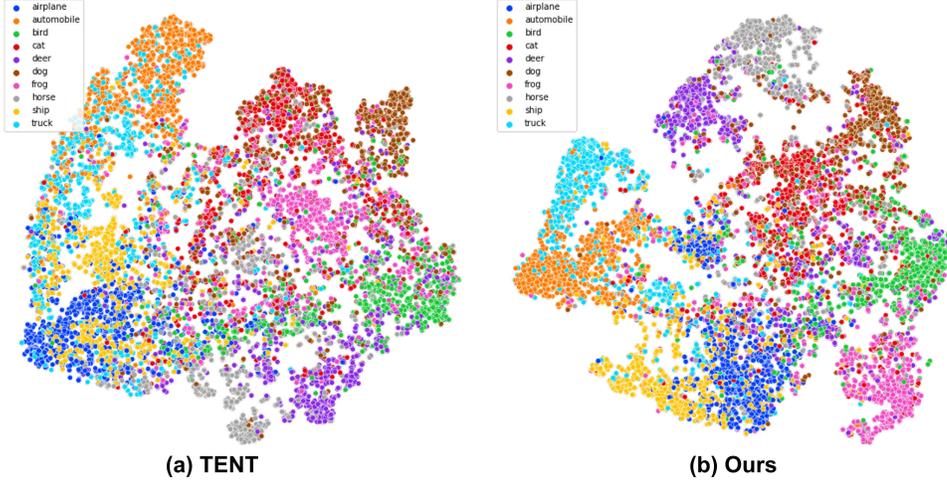


Figure 7. T-SNE visualizations of (a) TENT and (b) IABN+Ours. We visualize the feature map h obtained from the Gaussian noise corruption in the CIFAR-10-C uniform test dataset. The number of training samples is large in the order of the classes in the legend.

Src Method	TTA Method	Forward-LT			Uni.	Backward-LT			Avg.
		50	25	10	1	10	25	50	
<i>Cross Entropy</i>	Source	40.20	37.26	33.18	22.20	12.19	8.95	7.24	23.03
	BN Stats	48.12	45.44	40.17	26.17	13.67	9.69	7.80	27.30
	ONDA	48.49	45.80	41.16	27.66	15.17	11.02	8.98	28.33
	PseudoLabel	48.76	45.26	39.09	19.84	11.32	8.19	6.44	25.56
	TENT	49.17	45.21	38.42	15.32	9.99	7.44	5.67	24.46
	LAME	38.17	35.17	31.22	20.44	11.02	7.93	6.30	21.46
	CoTTA	32.83	29.35	25.72	14.75	7.38	4.94	3.67	16.95
	NOTE	46.41	44.10	40.49	29.30	15.77	11.87	9.84	28.26
	IABN	46.43	43.92	39.80	25.35	14.71	11.13	9.27	27.23
+Ours	53.20	50.77	46.26	32.34	18.56	14.07	11.74	32.42	
<i>Balanced Sampling</i>	Source	34.88	32.54	29.12	20.29	11.80	9.09	7.51	20.75
	BN Stats	45.07	42.93	39.10	28.67	17.82	13.97	11.88	28.49
	ONDA	44.79	42.60	39.01	29.20	18.62	14.65	12.70	28.80
	PseudoLabel	47.08	44.66	40.45	26.98	17.31	13.59	11.23	28.76
	TENT	48.28	45.64	41.07	24.04	16.88	13.50	11.30	28.67
	LAME	32.88	30.44	27.08	18.48	10.48	7.89	6.39	19.09
	CoTTA	30.38	28.68	25.58	17.32	10.74	7.96	6.49	18.16
	NOTE	46.62	44.27	40.64	30.42	16.99	12.77	10.32	28.86
	IABN	46.85	44.29	40.40	27.20	16.12	12.05	9.88	28.12
+Ours	51.32	49.18	44.99	31.92	19.64	15.11	13.00	32.16	
<i>Classifier Re-Training</i>	Source	40.17	37.47	33.59	23.23	13.42	10.18	8.51	23.80
	BN Stats	48.33	45.60	40.82	27.49	15.13	11.16	9.11	28.23
	ONDA	48.33	45.77	41.57	28.96	16.78	12.50	10.63	29.22
	PseudoLabel	49.10	45.89	40.19	21.01	12.74	9.50	7.66	26.58
	TENT	49.70	46.19	39.79	16.59	11.28	8.93	6.92	25.63
	LAME	38.22	35.41	31.60	21.48	12.28	9.13	7.52	22.23
	CoTTA	31.93	29.23	25.91	15.26	8.18	5.55	4.33	17.20
	NOTE	45.96	44.04	41.14	31.68	18.52	14.66	12.67	29.81
	IABN	46.11	44.04	40.53	27.61	17.28	13.71	11.94	28.75
+Ours	53.11	50.86	46.66	34.04	20.77	16.41	14.12	33.71	

Table 11. Ablation study on the source pre-trained model using CIFAR-100-LT and CIFAR-100-C.

A.3. Details of Label Shift Adapter

Model Architecture. We utilize the same model architecture for the label shift adapter across all datasets. The proposed label shift adapter consists of two fully-connected (FC) layers and a ReLU activation function, structured as FC-ReLU-FC. Furthermore, the label shift adapter is partitioned into two neural networks producing (γ_h, β_h) and $(\Delta W, \Delta b)$. As described in the main manuscript, the label shift adapter takes $m^\top \pi \in \mathbb{R}^1$ and produces (γ_h, β_h) and $(\Delta W, \Delta b)$ in each respective neural network. The hidden layer size in the label shift adapter is configured to 100.

Details of Label Shift Adapter. We provide the algorithm of the training process for the label shift adapter as a pseudo-code in Algorithm 1. The primary objective of the label shift adapter is to learn the relationship between π and adaptive parameters by selecting appropriate τ based on sampled π within generalized logit adjusted loss [33, 1] function. Increasing τ results in decision boundary shifting away from the minority class towards the majority class. Consequently, instead of sampling batches differently based on π , we sample π and τ iteratively, as described in Algorithm 1. This enables the label shift adapter to optimize its parameters in accordance with input label distributions (e.g., π and \hat{Y}_t), thereby producing suitable parameter adjustments.

During the training of the label shift adapter, we sample the label distribution π from three types of label distributions: $\{\pi_s, u, \bar{\pi}_s\}$. For each sampled label distribution, we select the appropriate $\tau \subset \{\tau_{\pi_s}, \tau_u, \tau_{\bar{\pi}_s}\}$, with the hyperparameter τ corresponding to each π . Different τ values are employed for each dataset. We set τ to $\{1, -1.5, 3\}$, $\{1, 0, -2\}$, $\{1, 0, -2\}$, $\{1, 0, -2\}$, $\{1, -1, -3\}$, and $\{1, 0, -2\}$, for CIFAR-10-LT, CIFAR-100-LT, ImageNet-LT, VisDA-C, OfficeHome, and DomainNet, respectively.

The mapping vector m maps the label distribution’s vector to the scalar of the imbalance degree. We set the range of m from -1 to 1, with the values increasing proportionally to the data count rank of each class. This technique enables the adapter to effectively utilize the degree of imbalance as an input, circumventing the challenges associated with complex label spaces encountered when using π directly.

While training the label shift adapter, we employ the same optimizer and batch size as those employed for training the source models. The learning rate is set to 1e-3 for all datasets. Moreover, we train the label shift adapter for $\{200, 200, 30, 15, 50, 20\}$ epochs.

During inference, the momentum hyperparameter α for target label distribution estimation is configured to 0.1. For learnable parameters in the test-time adaptation process, we only update affine parameters in normalization layers by following TENT [46] and IABN [11]. Unlike TENT, we freeze the top layers and update the affine parameters

of the layer in the remaining shallow layers, inspired by previous work [7, 36]. Specifically, for ResNet, including four layer groups (layer 1, 2, 3, 4), we only freeze layer4 in CIFAR-10-C, CIFAR-100-C, and ImageNet-C. In other datasets, there is no significant difference in performance, so all affine parameters are trained. When estimating the label distribution on ImageNet-C, we utilize only the top-3 probability to update the estimated label distribution \hat{Y}_t . Empirically, we discovered that it is effective to consider only top- k when the number of classes is particularly large.

B. Further Analysis on Label Shift Adapter

Ablation Study on Architecture Design. We examine the model architecture design for the proposed label shift adapter. The label shift adapter produces four types of outputs: $\gamma_h, \beta_h, \Delta W$, and Δb . Table 10 presents the ablation study for each component. Interestingly, even when only Δb is employed, the performance is quite good. However, we observed that as the degree of the label shift increases, the performance of using only Δb declines. Moreover, utilizing γ_h and β_h only also yields impressive results, indicating that appropriately shifting the feature map h is effective in addressing the label shifts. We choose the architecture design of the label shift adapter that achieves the best average accuracy, indicating that the final model generally performs well across a variety of label distributions.

T-SNE Visualization. To further substantiate the effectiveness of our method, we visualize the feature map h using t-SNE by extracting h during test-time adaptation. As illustrated in Fig. 7, our method shows a more well-separated representation space in a class-wise manner compared to TENT. Notably, it is evident that the minority classes (e.g., horse and truck) are not well divided in the representation space of TENT. In contrast, our method integrating into IABN layers enhances class-discriminability.

Ablation Study on Source Model. In the main manuscript, we employ the balanced softmax to reduce the model bias towards the majority classes. To further validate the effectiveness of the proposed method, we apply our method to several source pre-trained models utilizing different training strategies. We employ three types of techniques: (i) Cross-entropy loss, (ii) Balanced sampling, (iii) Classifier re-training [18], where the feature extractor is trained using cross-entropy loss, and then the classifier is randomly re-initialized and re-trained using class-balanced sampling. Table 11 demonstrates that our method effectively handles the label shifts, regardless of the source pre-trained models. Moreover, these results indicate that existing long-tailed recognition methods can be combined with our method to further reduce the model bias towards the majority classes in source domain data.

Ablation Study on π . As described in the main manuscript, we sampled three kinds of label distributions for π during

Num.	F50	F25	F10	U	B10	B25	B50	Avg.
3	52.06	49.71	46.03	36.84	29.29	26.33	25.50	37.97
5	50.91	48.82	45.41	36.90	29.28	26.37	25.51	37.60
7	51.13	48.99	45.52	36.96	29.24	26.25	25.45	37.64
∞	51.62	49.36	45.85	37.09	29.12	26.06	25.03	37.73

Table 12. Ablation study on the number of π for training label shift adapter using CIFAR-100-C. Num. denotes the number of π for training the adapter. F, U, and B indicate forward, uniform, and backward distributions, respectively. We chose three label distributions.

	DELTA	ISFDA	TENT+Ours
VISDA-C	50.10	61.02	72.97

Table 13. Comparison with additional baselines in test-time adaptation setting.

	Method	F50	F25	F10	U	B10	B25	B50	Avg.
CIFAR10	SAR+GN	57.22	57.20	57.07	57.12	61.84	63.06	64.37	59.70
	SAR+BN	78.63	76.28	71.82	53.28	34.99	28.60	25.18	52.68
	Ours+IABN	80.58	78.62	75.26	63.34	68.54	70.07	71.64	72.58
CIFAR100	SAR+GN	9.09	9.59	10.23	14.05	18.93	20.46	21.70	14.86
	SAR+BN	49.44	47.04	43.39	32.18	20.22	16.24	13.96	31.78
	Ours+IABN	52.06	49.71	46.03	36.84	29.29	26.33	25.50	37.97

Table 14. Comparison with SAR using CIFAR-10-C and CIFAR-100-C in TTA setting. F, U, and B denote forward, uniform, and backward, respectively. GN and BN indicate group and batch normalization, respectively.

training label shift adapter. Regarding the effect of sampling different numbers of π , Table 12 indicates that such variations have negligible impact on performance. Specifically, in this experiment, we interpolate three distributions (*i.e.*, π_s , u , $\bar{\pi}_s$) and τ to train the label shift adapter when different numbers of π are utilized.

C. Additional Experiments

Comparison with Baselines Related to Label Shifts.

We’ve compared two baselines in Table 7, which have the capability of handling label shifts. We compare additional baselines, DELTA [53] and ISFDA [22], which address covariate and label shifts simultaneously. Although ISFDA requires several epochs for adapting the source models, we conduct the experiments in the test-time adaptation setting for a fair comparison. Table 13 demonstrates that our method is superior to baselines significantly in the VISDA-C dataset. ISFDA, a domain adaptation model, exhibits limitations in its suitability for online learning during inference. Since DELTA only focuses on class imbalances in the target domain, it lacks the ability to handle imbalances in the source domain. In contrast, our method successfully addresses the imbalance in both source and target domains in the test-time adaptation setting.

Comparison with Recent TTA Baseline. We compare recent test-time adaptation baseline, sharpness-aware and

reliable entropy minimization (SAR) [36]. SAR proposes an optimizer and analyzes normalization layers to resolve imbalances in the target domain. However, it is important to note that our work addresses imbalances in both the source and target domains. Table 14 demonstrates that our method outperforms both SAR+GN and SAR+BN significantly. Moreover, it is a viable option to integrate our method with SAR method.

D. Domain-wise Results

Table 15, 16, 17 show the average classification accuracy on CIFAR-10-C, CIFAR-100-C, and ImageNet-C, shown per domain. To compute the accuracy of each domain, we calculate the average performances of Forward50, Forward25, Forward10, Uniform, Backward10, Backward25, and Backward50, as described in the main manuscript. These results demonstrate that our method consistently enhances performance across various domains.

Method	Gaussian	Shot	Impulse	Defocus	Glass	Motion	Zoom	Snow	Frost	Fog	Brightness	Contrast	Elastic	Pixelate	JPEG	Avg
Source	23.06	26.94	19.59	46.89	41.06	45.63	49.42	58.30	45.31	45.72	69.85	21.23	57.60	60.14	65.45	45.08
BN Stats	49.18	50.52	46.29	58.11	45.05	56.14	55.96	52.32	50.65	53.60	59.11	54.93	51.63	54.15	52.12	52.65
ONDA	50.70	51.66	47.68	59.80	46.49	57.45	57.78	53.81	52.36	55.14	61.52	54.76	53.60	56.49	54.19	54.23
PseudoLabel	46.87	48.76	44.47	55.96	43.87	53.89	53.46	49.96	48.70	50.91	56.34	52.57	49.29	51.88	50.16	50.47
LAME	17.97	22.75	15.43	44.74	40.36	42.40	47.08	61.30	48.62	45.20	67.84	20.33	55.40	61.36	64.59	43.69
CoTTA	51.69	53.11	50.31	55.80	47.28	54.31	54.88	52.60	52.18	52.81	58.30	49.66	52.12	55.32	54.39	52.98
NOTE	54.48	56.22	53.24	68.20	48.64	64.87	65.09	65.56	64.43	64.08	73.33	67.59	60.24	66.95	67.26	62.68
TENT	46.41	48.29	43.38	53.82	42.42	52.22	51.57	48.92	47.51	49.81	55.06	50.62	47.90	50.58	49.20	49.18
+ Ours	51.66	53.55	48.66	60.74	46.94	58.77	57.37	55.04	53.48	56.00	62.05	57.68	54.11	57.48	55.07	55.24
IABN	54.77	56.48	53.25	68.24	48.39	64.53	64.89	65.63	64.44	64.60	73.79	67.24	60.32	67.81	67.17	62.77
+ Ours	68.72	69.54	64.94	77.48	61.47	76.24	75.52	72.06	72.83	74.53	79.69	79.08	69.66	74.25	72.69	72.58

Table 15. Domain-wise results on CIFAR-10-C.

Method	Gaussian	Shot	Impulse	Defocus	Glass	Motion	Zoom	Snow	Frost	Fog	Brightness	Contrast	Elastic	Pixelate	JPEG	Avg
Source	14.42	16.45	9.08	22.55	23.14	25.65	25.82	27.37	22.98	18.57	34.51	5.84	33.03	17.47	36.19	22.21
BN Stats	27.59	28.51	26.20	36.89	28.61	34.53	36.44	29.40	28.81	28.98	36.66	29.22	33.37	33.78	32.18	31.41
ONDA	27.53	28.77	26.17	36.85	29.07	34.37	36.91	29.71	29.09	29.38	36.97	27.77	33.89	34.00	33.00	31.57
PseudoLabel	27.44	28.54	25.59	34.92	27.78	32.83	34.56	28.58	27.40	28.58	34.96	26.01	31.58	32.86	31.38	30.20
LAME	13.41	15.73	7.66	20.93	22.30	24.79	24.63	27.21	22.39	17.07	33.62	4.48	32.41	15.62	35.66	21.19
CoTTA	30.01	30.85	28.45	34.77	30.64	34.04	35.64	30.92	30.10	28.65	36.09	24.30	33.54	35.58	34.00	31.84
NOTE	24.17	25.64	18.62	35.73	28.08	36.89	37.48	34.91	33.95	29.13	41.47	33.93	36.29	32.01	36.13	32.30
TENT	27.50	28.49	25.28	33.98	26.89	32.12	33.36	28.00	26.79	27.78	34.05	25.20	31.01	32.25	30.53	29.55
+ Ours	30.95	32.21	28.77	38.60	30.58	36.06	38.24	32.30	30.74	31.52	38.34	30.36	34.75	36.19	34.86	33.63
IABN	24.54	25.79	18.92	35.50	28.00	36.80	37.58	34.97	34.20	29.02	41.39	33.99	36.17	32.09	36.29	32.35
+ Ours	33.65	34.37	28.17	41.86	33.62	41.08	41.79	37.82	38.36	34.77	43.51	42.54	39.18	39.98	38.76	37.97

Table 16. Domain-wise results on CIFAR-100-C.

Method	Gaussian	Shot	Impulse	Defocus	Glass	Motion	Zoom	Snow	Frost	Fog	Brightness	Contrast	Elastic	Pixelate	JPEG	Avg
Source	5.72	5.88	4.82	15.48	11.17	17.98	18.73	21.25	15.21	28.04	44.67	28.63	39.16	29.56	34.77	21.40
BN Stats	29.27	28.90	25.46	25.59	22.94	33.14	32.91	28.37	25.15	40.06	45.27	40.19	43.59	41.84	39.74	33.49
ONDA	29.15	28.88	25.33	25.49	22.74	32.73	32.96	28.18	25.04	40.12	45.46	39.69	43.60	42.02	39.72	33.41
PseudoLabel	31.25	30.93	29.00	27.72	25.80	34.36	33.64	30.20	25.39	40.33	44.09	39.56	42.78	41.43	39.65	34.41
LAME	5.56	5.73	4.67	15.33	11.03	17.92	18.67	21.19	15.16	28.01	44.64	28.61	39.12	29.49	34.75	21.33
CoTTA	30.93	30.36	27.47	27.28	24.95	34.42	33.56	30.01	26.44	40.62	44.79	40.58	43.32	41.83	40.05	34.44
NOTE	31.34	30.83	29.26	27.39	24.66	35.67	32.70	33.34	28.33	39.52	46.55	44.97	43.67	41.38	41.59	35.41
TENT	29.29	28.94	25.84	25.68	22.94	32.11	32.37	27.17	23.50	39.43	43.94	38.18	42.30	40.66	39.02	32.76
+ Ours	32.38	32.39	28.77	29.07	26.21	35.70	35.81	31.44	27.86	43.14	48.56	43.00	46.59	44.92	43.59	36.63
IABN	31.47	30.86	29.28	27.37	24.66	35.69	32.77	33.35	28.37	39.54	46.65	45.03	43.72	41.40	41.60	35.45
+ Ours	34.28	34.00	32.18	30.25	27.14	38.92	35.53	36.48	31.58	42.50	49.95	48.35	46.93	44.95	44.87	38.53

Table 17. Domain-wise results on ImageNet-C.