# LEARNING DYNAMICS OF LOGITS DEBIASING FOR LONG-TAILED SEMI-SUPERVISED LEARNING

**Anonymous authors**Paper under double-blind review

#### **ABSTRACT**

Long-tailed distributions are prevalent in real-world semi-supervised learning (SSL), where pseudo-labels tend to favor majority classes, leading to degraded generalization. Although numerous long-tailed SSL (LTSSL) methods have been proposed, the underlying mechanisms of class bias remain underexplored. In this work, we investigate LTSSL through the lens of learning dynamics and introduce the notion of baseline images to characterize accumulated bias during training. We provide a step-wise decomposition showing that baseline predictions are determined solely by shallow bias terms, making them reliable indicators of class priors. Building on this insight, we propose a novel framework, DyTrim, which leverages baseline images to guide data pruning. Specifically, we perform classaware pruning on labeled data to balance class distribution and label-agnostic soft pruning with confidence filtering on unlabeled data to mitigate error accumulation. Theoretically, we show that our method implicitly realizes risk reweighting, effectively suppressing class bias. Extensive experiments on public benchmarks show that DyTrim consistently enhances the performance of existing LTSSL methods by improving representation quality and prediction accuracy.

# 1 Introduction

Semi-supervised learning (SSL), exemplified by FixMatch (Sohn et al., 2020), has been proven to demonstrate significant generalization advantages over supervised learning, particularly in deep neural networks (Li et al., 2025). However, many existing SSL variants (e.g., FlexMatch; Zhang et al., 2021) implicitly assume that both labeled and unlabeled data are drawn from a balanced class distribution. In practice, datasets commonly exhibit a long-tailed label distribution, leading to biased pseudo-label toward majority classes. This discrepancy poses significant challenges to the effectiveness of SSL algorithms on real-world datasets.

Recent studies on long-tailed semi-supervised learning (LTSSL) have emerged to mitigate pseudo-label bias caused by class imbalance in both labeled and unlabeled data. These methods range from distribution alignment (Wei et al., 2021; Kim et al., 2020), data rebalancing (Fan et al., 2022; Lee et al., 2021), logit adjustment variants (Wei & Gan, 2023; Zhou et al., 2024), to foundation model-based methods (e.g., LADaS; Zheng et al., 2025). In particular, the approach employ baseline image was introduced as a simple yet efficient tool to quantify classifier bias by CDMAD (Lee & Kim, 2024), which has attracted considerable attention in the community (Xing et al., 2025). However, the underlying mechanisms of how class bias emerges and why existing approaches can mitigate it remain largely unexplored and poorly understood. That also prevents us from exploring a principle-based method to improve performance.

In this paper, we analyze the underlying mechanisms of class debiasing through an innovative lens of learning dynamics, investigating how an input affects the generation of biased pseudo-labels. We first point out that in the training processes of LTSSL, the logits of the baseline image serve as an indicator of the accumulated influence of the network's bias term. We further propose a framework that formalizes the learning dynamics of semi-supervised learning by decomposing the change of the model's prediction on the baseline image into three terms. Under this framework, many existing debiasing methods for class imbalance can be unified.

Furthermore, our analysis of bias accumulation dynamics motivates a pruning-based class debiasing framework. For labeled data, we compute class-wise pruning ratios to rebalance samples. For un-

labeled data, we apply a label-agnostic criterion that prunes low-confidence, inconsistent samples. Beyond empirical gains or ad-hoc analysis, DyTrim provide a principle-based theoretical guarantees that clarify how the proposed method can alleviate class biasing and why pruning enhances generalization. Extensive experiments confirm that DyTrim consistently enhances LTSSL performance across standard benchmarks.

# 2 PRELIMINARIES

**Notions.** We consider a labeled dataset  $\mathcal{X}=\{(x^n,y^n)\}_{n=1}^N$  with N samples and an unlabeled dataset  $\mathcal{U}=\{u^m\}_{m=1}^M$  with M samples, where  $x^n\in\mathbb{R}^d$  is the n-th labeled sample with label  $y^n\in[C]=\{1,\ldots,C\}$ , and  $u^m\in\mathbb{R}^d$  is the m-th unlabeled sample. Let  $N_c$  and  $M_c$  denote the number of labeled and unlabeled samples in class c, such that  $\sum_{c=1}^C N_c = N$  and  $\sum_{c=1}^C M_c = M$ . If classes are sorted by size, we have  $N_1\geq N_2\geq \cdots \geq N_C$ , and define the imbalance ratios as  $\gamma_l=N_l/N_c\geq 1$  and  $\gamma_u=M_l/M_c\geq 1$ , respectively. We denote the classifier by  $f_\theta:\mathbb{R}^d\mapsto 1,\ldots,C$  with parameters  $\theta$ , and its logits by  $g_\theta(x)\in\mathbb{R}^C$ , where  $f_\theta(x)=\arg\max_c g_\theta(x)_c$  and  $(\cdot)_c$  denotes the c-th component. For each iteration of training, we sample minibatches  $\mathcal{MX}=\{(x_b^n,y_b^n):b\in(1,\ldots,B)\}\subset\mathcal{X}$  and  $\mathcal{MU}=\{(u_b^m):b\in(1,\ldots,\mu B)\}\subset\mathcal{U}$  from the training set, where B denotes the minibatch size and  $\mu$  denotes the relative size of  $\mathcal{MU}$  to  $\mathcal{MX}$ . For brevity, when clear from context we drop the superscript on  $u_b^m(x_b^n)$  and simply write  $u_b(x_b)$ .

Base SSL algorithms. We use FixMatch (Sohn et al., 2020) as the base SSL algorithm, following other LTSSL studies. Specifically, FixMatch first predicts the class probability of a weakly augmented unlabeled data point  $\alpha(u_b)$  as  $q_b = \pi_\theta(y|\alpha(u_b))$  and then generates hard pseudo-label  $\hat{q}_b = \arg\max_c(q_{b,c})$ , where  $\pi_\theta(y|\cdot) = \text{Softmax}(g_\theta(\cdot))$ . For consistency regularization, FixMatch uses a hard pseudo-label  $\hat{q}_b$  only when  $\max_c(q_{b,c}) \geq \tau$ , where  $\tau$  denotes a predefined confidence threshold, to improve the quality of the pseudo-labels used for training. We express the training losses of FixMatch  $\mathcal L$  as:

$$\mathcal{L}(x_b, u_b, \hat{q}, \tau; \theta) = \mathcal{L}_{sup}(\alpha(x_b); \theta) + \mathcal{L}_{con}(\mathcal{A}(u_b), \hat{q}_b, \tau; \theta)$$
(1)

where  $x_b$  ( $u_b$ ) denotes the b-th labeled (unlabeled) samples in a minibatch  $\mathcal{MX}$  ( $\mathcal{MU}$ ).  $\mathcal{A}(u_b)$  denotes the strongly augmented version of  $u_b$ . The losses and other base SSL algorithms, *i.e.* Flex-Match (Zhang et al., 2021) and FreeMatch (Wang et al., 2023b), are detailed in Appendix B.

**Learning dynamics and its per-step decomposition.** Inspired by Ren & Sutherland (2024), we study how a single gradient update changes the model's confidence on an observation  $x_o$ . With  $\pi_{\theta}(y \mid x)$  denoting the predicted class probability distribution, the learning dynamics become,

$$\Delta \theta \triangleq \theta^{t+1} - \theta^t = -\eta \cdot \nabla \mathcal{L}(f_{\theta}(x_b), y_b); \quad \Delta \log \pi^t(y|x_o) \triangleq \log \pi_{\theta^{t+1}}(y|x_o) - \log \pi_{\theta^t}(y|x_o). \tag{2}$$

where the update of  $\theta$  during step  $t \to t+1$  is given by one gradient update on the sample pair  $(x_b, y_b)$  with learning rate  $\eta$ .  $\mathcal{L}$  is the loss function, we use the cross-entropy loss  $\mathbf{H}$  in our setting.

**Definition 1** (Per-step decomposition of learning dynamics). Let  $\pi = Softmax(\mathbf{z})$  with  $\mathbf{z} = g_{\theta}(x)$ . Then the one-step learning dynamics decompose as

$$\Delta \log \pi_{\theta}^{t}(y \mid x_{o}) = -\eta \mathcal{T}^{t}(x_{o}) \mathcal{K}^{t}(x_{o}, x_{b}) \mathcal{G}^{t}(x_{b}, y_{b}) + \mathcal{O}\left(\eta^{2} \|\nabla_{\theta} \mathbf{z}(x_{b})\|_{op}^{2}\right), \tag{3}$$

where  $\mathcal{T}^t(x_o) = \nabla_z \log \pi_{\theta^t}(x_o) = I - \mathbf{1}\pi_{\theta^t}^\top(x_o)$  only depends on the model's current predicted probability,  $\mathcal{K}^t(x_o, x_b) = (\nabla_\theta z(x_o)|_{\theta^t})(\nabla_\theta z(x_b)|_{\theta^t})^\top$  is the empirical neural tangent kernel (eNTK, Jacot et al. 2018) of the model, the product of the model's gradients with respect to  $x_o$  and  $x_b$ .  $\mathcal{G}^t(x_b, y_b) = \nabla_{\mathbf{z}} \mathcal{L}(x_b, y_b)|_{\mathbf{z}^t}$  is the loss gradient.  $\|\cdot\|_{op}^2$  denotes the spectral norm, which bounds the second-order remainder term.

This decomposition characterizes how each update at  $(x_b, y_b)$  influences predictions at  $x_o$ , forming the basis for our SSL analysis under class imbalance.

# 3 LEARNING DYNAMICS OF LONG-TAILED SEMI-SUPERVISED DEBIASING

#### 3.1 LEARNING DYNAMICS OF SEMI-SUPERVISED LEARNING

In this section, we characterize the learning dynamics of the semi-supervised version of gradient descent (GD) for the FixMatch algorithm Eq. (1),

$$\Delta \theta \triangleq \theta^{t+1} - \theta^t = -\eta \cdot (\nabla \mathcal{L}_{sup}(f_{\theta}(\alpha(x_b)), y_b) + \nabla \mathcal{L}_{con}(f_{\theta}(\alpha(u_b)), f_{\theta}(\mathcal{A}(u_b)));$$

$$\Delta f(x_o) \triangleq f_{\theta^{t+1}}(x_o) - f_{\theta^t}(x_o).$$
(4)

where  $x_o$  denotes the observation data point, the update of  $\theta$  during step  $t \to t+1$  is given by one gradient update on the labeled sample pair  $(x_b, y_b)$  and unlabeled sample  $(u_b)$  with learning rate  $\eta$ . Previous work (Ren & Sutherland, 2024) showed how a single gradient update influences model predictions in supervised learning. We now examine whether such characterization extends to the semi-supervised setting. Since FixMatch (Sohn et al., 2020) update naturally consists of a supervised part  $\mathcal{L}_{sup}$  and a consistency part  $\mathcal{L}_{con}$ , the gradient update can be decomposed accordingly. For an unlabeled sample  $u_b$  with target  $\hat{q}_b^t = \arg\max_c q_{b,c}^t$ , where  $q_b^t = \pi_{\theta^t}(\cdot \mid \alpha(u_b))$ . The per-step learning dynamics of semi-supervised learning become

$$\Delta \log \pi^{t}(y|x_{o}) \triangleq \Delta \log \pi_{\theta}^{t,\sup}(y \mid x_{o}; x_{b}) + \Delta \log \pi_{\theta}^{t, \operatorname{con}}(y \mid x_{o}; u_{b})$$

$$\tag{5}$$

where  $\Delta \pi_{\theta}^{t, \text{sup}}$  denotes the influence caused by  $x_b$  and  $\Delta \pi_{\theta}^{t, \text{con}}$  denotes the influence caused by  $u_b$ , respectively. We now state the decomposition of the per-step influence below:

**Proposition 1.** For an labeled (unlabeled) sample  $x_b$  ( $u_b$ ) with target  $y_b$  ( $\hat{q}_b^t = \arg \max_c q_{b,c}^t$ ), where  $q_b^t = \pi_{\theta^t}(y|\alpha(u_b))$ . The one-step learning dynamics of SSL decompose as

$$\Delta \log \pi_{\theta}^{t,\sup}(y \mid x_o; x_b) = -\eta \mathcal{T}^t(x_o) \mathcal{K}^t(x_o, \alpha(x_b)) \mathcal{G}_{\sup}^t(\alpha(x_b), y_b) + \mathcal{O}\left(\eta^2 \|\nabla_{\theta} \mathbf{z}(\alpha(x_b))\|_{op}^2\right)$$

$$\Delta \log \pi_{\theta}^{t,\operatorname{con}}(y \mid x_o; u_b) = -\eta \mathcal{T}^t(x_o) \mathcal{K}^t(x_o, \mathcal{A}(u_b)) \mathcal{G}_{\operatorname{con}}^t(\mathcal{A}(u_b), \hat{q}_b^t) + \mathcal{O}\left(\eta^2 \|\nabla_{\theta} \mathbf{z}(\mathcal{A}(u_b))\|_{op}^2\right)$$

$$(6)$$

where  $\mathcal{K}^t(x_o, \alpha(x_b))$  and  $\mathcal{K}^t(x_o, \mathcal{A}(u_b))$  are eNTK evaluations of the logit network  $\mathbf{z}(\cdot) = g_{\theta}(\cdot)$ , with different inputs.  $\mathcal{G}^t_{\sup}(\alpha(x_b), y_b) = \nabla_{\mathbf{z}} \mathcal{L}_{\sup}(\alpha(x_b), y_b)|_{\mathbf{z}^t}$  and  $\mathcal{G}^t_{\cos}(\hat{q}_b, \mathcal{A}(u_b)) = \nabla_{\mathbf{z}} \mathcal{L}_{\cos}(\hat{q}_b, \mathcal{A}(u_b))|_{\mathbf{z}^t}$ , respectively.

Accumulated influence and a demonstration on MNIST. As shown in Proposition 1, each update of  $\theta$  in FixMatch decomposes into a supervised part driven by  $(x_b, y_b)$  and a consistency part driven by  $(u_b, \hat{q}_b^t)$ . While this decomposition captures the per-step influence on  $\pi_{\theta}(y \mid x_o)$ , in practice training consists of many such steps, and the accumulated effect is governed by the iterative interaction between labeled and unlabeled updates.

To demonstrate this, we train a WRN-28-2 on MNIST and visualize the accumulated influence in Figure 1. In the top row, when  $\hat{q}_b$  is correct, the consistency term reinforces the supervised signal, gradually pulling the prediction of  $x_o$  toward the correct class, consistent with the constructive dynamics implied by Eq. (6). In contrast, when  $\hat{q}_b$  is incorrect (top right), the consistency update exerts the opposite effect, systematically reducing the correct probability of  $x_o$ . This illustrates how pseudo-label errors, even if small at each step, can accumulate across iterations into a negative reinforcement

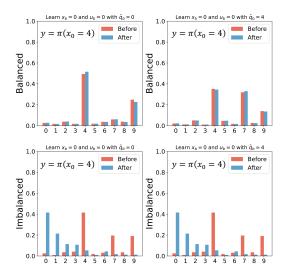


Figure 1: The per-step semi-supervised learning dynamics and the accumulated influence in an MNIST experiment.

loop. The bottom row shows that under class imbalance, such accumulated influence can drive the classifier to consistently predict the majority class (here 0), regardless of the true label. This confirms the implication of our dynamics analysis: in SSL, the effect of labeled data is mediated through pseudo-labels, so local errors can be amplified rather than averaged out, leading to catastrophic bias.

165

166

167

168

169

170

171

172

173

174 175

176

177

178

179

181

182

183

185

186

187

188

189

190

191

192 193

194

195 196

197

199

200

201

202

203 204 205

206

207

208

209 210

211

212

213

214

215

#### 3.2 LEARNING DYNAMICS ANALYSIS OF ACCUMULATED BIAS UNDER CLASS IMBALANCE

The aforementioned phenomenon, together with the learning dynamics of the semi-supervised framework, illustrates how class imbalance accumulates into systematic bias. While per-update dynamics capture the influence of individual samples on predictions, they fall short of reflecting the global effect of imbalance. This motivates the search for an indicator that bridges class-imbalance bias with the underlying learning dynamics. We propose to use a baseline image  $\mathcal{I}$  as such an indicator. To justify this choice, we analyze its theoretical properties in both linear and deep settings, and then incorporate it into the per-step influence decomposition.

Baseline image and its invariance property. For simplicity, we first consider a two-layer MLP with no bias in the first layer and a bias vector  $\boldsymbol{b} \in \mathbb{R}^C$  in the output layer  $h(x) = h^{(2)} \circ h^{(1)}(x)$ , where  $h^{(1)}(x) = \sigma(W_1 x)$  and  $h^{(2)} = W_2 x + b$ . This setting allows us to isolate and examine the predicted class probability  $\pi_{\theta}(\mathcal{I})$  of a baseline image. For a baseline image  $\mathcal{I} \in \mathbb{R}^d$ , we have

$$h(\mathcal{I}) = \mathbf{W}_2 h^{(1)}(\mathcal{I}) + \mathbf{b}. \tag{7}$$

In modern neural networks, the explicit bias term b is typically absorbed into the normalization layer, e.g., BatchNorm, LayerNorm, while other layers are usually set without bias. Without loss of generality, we take BatchNorm as an example for analysis. Since the BatchNorm transformation can be equivalently viewed as an affine linear layer with learnable parameters, we may replace  $h^{(2)}$ with a BatchNorm( $\cdot$ ) layer, i.e.,

$$h(\mathcal{I}) = \operatorname{BatchNorm}(h^{(1)}(\mathcal{I})) = \frac{h^{(1)}(\mathcal{I}) - \mathbb{E}[h^{(1)}(\mathcal{I})]}{\sqrt{\operatorname{Var}[h^{(1)}(\mathcal{I})]}} \cdot \mathbf{W}_2 + \mathbf{b}. \tag{8}$$

This replacement highlights that the baseline image prediction  $\pi_{\theta}(\mathcal{I})$  is directly governed by the BN bias b, thus allowing us to focus on its role in encoding and accumulating class-imbalance bias. We now state the main results regarding the  $\pi_{\theta}(\mathcal{I})$  below:

**Proposition 2** (Invariance of baseline image under affine normalization). Let  $\mathcal{I} = k \cdot \mathbf{1}_d$  be a baseline image, where  $k \in \{0, 1, \dots, 255\}$  and  $\mathbf{1}_d \in \mathbb{R}^d$  is an all-one vector. Suppose the output of the first hidden transformation is normalized by a normalization layer (e.g., BatchNorm, LayerNorm, InstanceNorm, or GroupNorm) with affine parameters  $(W_2, b)$ . Then the logits  $h(\mathcal{I})$  are independent of k and reduce to

$$h(\mathcal{I}) = \mathbf{b}, \quad \pi_{\theta}(\mathcal{I}) = Softmax(\mathbf{b}).$$
 (9)

One can immediately notice that  $\pi_{\theta}(\mathcal{I})$  in Eq. (9) does not contain any term related to the pixel value k of  $\mathcal{I}$ . This observation implies that the representation  $\pi_{\theta}(\mathcal{I})$  of a baseline image is entirely determined by the BatchNorm bias term b, and is invariant to the actual pixel value k.

Building upon this invariance, we now establish a connection between the baseline image and the underlying class distribution. Specifically, for the classifier formulation in Eq. (8), we show that the logits of the baseline image encode the class-imbalance ratio of the training data, thereby providing a direct bridge between  $\pi_{\theta}(\mathcal{I})$  and the long-tailed class prior.

**Theorem 1** (Bias as the conditional distribution prior). Assume the model h(x) which characterized in Eq. (8), is trained by cross-entropy,

$$\mathcal{L} = \mathbb{E}_{(x,y)} \left[ -y^{\top} \log Softmax(h(x)) \right]. \tag{10}$$

At a population risk minimizer 
$$(\mathbf{W}_{2}^{\star}, \mathbf{b}^{\star})$$
 we have
$$\hat{p}^{\star}(x) = P(y \mid x), \qquad \hat{p}^{\star}(\mathcal{I}) = Softmax(\mathbf{b}^{\star}) = P(y \mid \frac{h^{(1)}(\mathcal{I}) - \mathbb{E}[h^{(1)}(\mathcal{I})]}{\sqrt{\operatorname{Var}[h^{(1)}(\mathcal{I})] + \epsilon}} = \mathbf{0}). \tag{11}$$

In particular, for solid-color  $\mathcal{I}$  satisfying Proposition 2, the baseline prediction equals the conditional class distribution at the "normalized-zero" feature state, thereby encoding the class prior induced during training.

Therefore,  $\pi_{\theta}(\mathcal{I})$  can naturally serve as a proxy for the accumulated bias of the model, providing a bridge between class imbalance and learning dynamics.

**Per-step influence decomposition of the baseline image.** Let the estimate of the underlying class prior  $P_{\theta}(y|\cdot)$  be denoted by  $\pi$ . Then we can track the change in the model's confidence by observing  $\log \pi_{\theta}(y|\cdot)$ . Then learning dynamics become,

$$\Delta \log \pi^{t}(y|\mathcal{I}) \triangleq \log \pi_{\theta^{t+1}}(y|\mathcal{I}) - \log \pi_{\theta^{t}}(y|\mathcal{I}). \tag{12}$$

**Proposition 3.** Let  $\pi = \mathtt{Softmax}(z)$  and  $z = g_{\theta}(x)$ . The one-step dynamics decompose as

$$\Delta \log \pi^{t}(y \mid \mathcal{I}) = -\eta \mathcal{T}^{t}(\mathcal{I}) \mathcal{K}^{t}(\mathcal{I}, x) \mathcal{G}^{t}(x, y) + \mathcal{O}(\eta^{2} \|\nabla_{\theta} z(x)\|_{\text{op}}^{2}), \tag{13}$$

where  $\mathcal{T}^t(\mathcal{I}) = \nabla_z \log \pi^t(\mathcal{I}) = I - \mathbf{1}\pi_{\theta^t}^T(\mathcal{I})$ ,  $\mathcal{K}^t(\mathcal{I}, x) = (\nabla_{\theta} z(\mathcal{I})|_{\theta^t}) (\nabla_{\theta} z(x)|_{\theta^t})^T$  is the empirical Neural Tangent Kernel (NTK) of the logit network z, and  $\mathcal{G}^t(x, y) = \nabla_z \mathcal{L}(x, y)|_{z^t}$  (see Appendix F for details).

#### 3.3 EFFECT OF THE BASELINE IMAGE FOR GUIDING DATA PRUNING

The training objective can be interpreted as the minimization of the empirical risk  $\mathcal{L}$ . Assuming that all labeled samples  $x_b^n$  from  $\mathcal{X}$  and unlabeled samples  $u_b^m$  from  $\mathcal{U}$  are drawn from continuous distributions  $\rho^l(x_b^n)$  and  $\rho^u(u_b^m)$ , respectively, the training objective can be formulated as:

$$\underset{\theta \in \Theta}{\operatorname{arg\,min}} \underset{x_b^n \in \mathcal{X}, u_b^m \in \mathcal{U}}{\mathbb{E}} [\mathcal{L}(x_b^n, u_b^m; \theta)] = \int_{x_b^n} \mathcal{L}_{sup}(x_b^n, \theta) \rho^l(x_b^n) dx_b^n + \int_{u_b^m} \mathcal{L}_{con}(u_b^m, \theta) \rho^l(u_b^m) du_b^m.$$
(14)

After applying a data pruning policy, we sample  $x_b^n$  and  $u_b^m$  to obtain the labeled pruned subset  $\mathcal{S}_t^l$ , according to the labeled pruning probabilities  $\mathcal{P}_t^l(x_b^n)$  and unlabeled pruning probabilities  $\mathcal{P}_t^l(u_b^m)$ , respectively. For the labeled samples, we directly optimize over the pruned subset  $\mathcal{S}_t^l$  without reweighting the loss terms. Notably, the class-aware pruning probability  $r_c = \pi_\theta(\mathcal{I})_c$  inherently adjusts  $\mathcal{S}_t^l$  toward an asymptotically balanced class distribution. By retaining more samples from minority classes (lower  $r_c$ ) and pruning more samples from majority classes (higher  $r_c$ ), the pruned subset  $\mathcal{S}_t^l$  naturally mitigates class imbalance. As a result, even without explicit rescaling, the empirical risk over  $\mathcal{S}_t^l$  approximates:

$$\underset{\theta \in \Theta}{\arg\min} \ \underset{x_b^n \in \mathcal{S}_t^l}{\mathbb{E}} \left[ \mathcal{L}_{sup}(x_b^n, \theta) \right] \propto \frac{1 - \mathcal{P}_t^l(x_b^n)}{c_t^l} \int_{z} \mathcal{L}_{sup}(x_b^n, \theta) \rho_l(x_b^n) dx_b^n, \tag{15}$$

where  $c_t^l = \mathbb{E}_{x_b^n \sim \rho_l}[1 - \mathcal{P}_t^l(x_b^n)]$ . The term  $\frac{1 - \mathcal{P}_t^l(z)}{c_t^l}$  acts as an *implicit reweighting* due to the class-aware pruning policy. For unlabeled samples, pruning with uniform probability r and rescaling losses by  $\gamma_t(u) = \frac{1}{1 - \mathcal{P}_t^u(u)}$  yields

$$\underset{\theta \in \Theta}{\arg\min} \underset{u_b^m \in \mathcal{S}_t^u}{\mathbb{E}} \left[ \gamma_t(u_b^m) \mathcal{L}_{con}(u_b^m, \theta) \right] \propto \frac{1}{c_t^u} \int_z \mathcal{L}_{con}(u_b^m, \theta) \rho^l(u_b^m) du_b^m, \tag{16}$$

where  $c_t^u = \mathbb{E}_{u_b^m \sim \rho_u}[1 - \mathcal{P}_t^u(u_b^m)]$ . Crucially, even with uniform pruning rates, the interplay of consistency regularization and confidence thresholding ensures  $\mathcal{S}_t^u$  to be implicitly balanced, thus training on  $\mathcal{S}_t^u$  with rescaled factor  $\gamma_t(u_b^m)$  could achieve a better result as training on the  $\mathcal{U}$ .

# 4 DyTrim: A Baseline Image Guided Data Pruning Framework for CISSL

The theoretical results in Section 3 suggest that the distribution of the baseline image's logits is affected by imbalanced data and directly acts on the bias term in the shallow layers. This imbalance directly causes the model to produce bias. Fortunately, we indicated that this bias can be effectively reduced if the data is pruned to be more balanced. Based on these insights, we propose the algorithm DyTrim, which extends the data pruning by incorporating guidance from the baseline image's logits to select a balanced subset for the training of CISSL, as illustrated in Figure. 2.

**Dynamic data pruning for CISSL.** We use  $\mathcal{X} = \{(x^n, y^n)\}_{n=1}^N$  to denote the labeled set and  $\mathcal{U} = \{u^m\}_{m=1}^M$  for unlabeled set. Critically, the distribution mismatch between  $\mathcal{X}$  and  $\mathcal{U}$  necessitates separate scoring mechanisms for labeled and unlabeled samples—unlike conventional supervised dynamic pruning methods that assume identical data distributions. To this end, we define step-dependent scoring functions  $\mathcal{H}_t^l$  for labeled samples and  $\mathcal{H}_t^u$  for unlabeled samples, which dynamically quantify sample utility at training step t. For the dynamic pruning process, samples are discarded by the step-dependent pruning probabilities  $\mathcal{P}_t^l$  and  $\mathcal{P}_t^u$ :

$$\mathcal{P}_t^l(x;\mathcal{H}_t^l) = \mathbb{1}(\mathcal{H}_t^l(x),_{\prec r_c}\mathcal{H}_{c,t}^l); \quad \text{and} \quad \mathcal{P}_t^u(u;\mathcal{H}_t^u) = \mathbb{1}(\mathcal{H}_t^u(u),\bar{\mathcal{H}}_t^u), \tag{17}$$

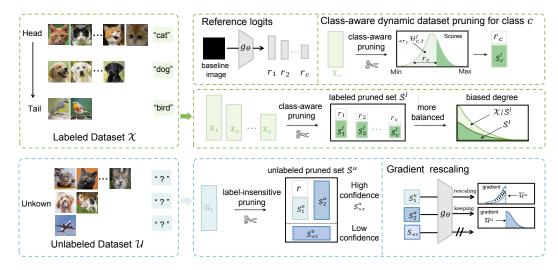


Figure 2: Illustration of the proposed DyTrim framework. DyTrim mainly consists of two operations, named labeled pruning and unlabeled pruning.  $_{\prec r_c}\mathcal{H}^l_{c,t}$  and  $\bar{\mathcal{H}}^u_t$  denote the adaptive thresholds of scores of labeled samples and unlabeled samples, with slight abuse of symbols.  $\mathcal{S}^u_{\prec \tau}$  denote the low confidence unlabeled sample which  $p^*(u^m_b) \geq \tau$ . Labeled pruning provides a class-aware pruning policy for each sample from class c. Unlabeled pruning provides a random pruning policy from the original data  $\mathcal{U}_1$  and uses a gradient rescaling strategy  $(\times 1/(1-r))$  for which sample from  $s^u_1$  is selected to prune) to keep the approximately same gradient expectation.

where  $_{\prec r_c} \mathcal{H}^l_{c,t}$  and  $\bar{\mathcal{H}}^u_t$  are adaptive thresholds,  $\mathbb{1}(\cdot,\cdot)$  is the indicator function. Thus, two dynamically pruned datasets  $\mathcal{S}^l_t$  and  $\mathcal{S}^u_t$  are formed for labeled and unlabeled datasets, respectively.

**Dynamic pruning for labeled data.** Since the labeled data follow a long-tailed class distribution, we design a class-aware pruning policy  $\mathcal{P}_t^l$  guided by  $\pi_{\theta}(\mathcal{I})$ . Critically, the classifier's pseudolabels are primarily influenced by the labeled samples, which introduce bias toward majority classes. Since Proposition 2 shows that the baseline image has invariance to solid-color intensity, from first principles, we leverage the logits from a **black image**  $\mathcal{I}$  to calibrate pruning probabilities. Given the labeled dataset  $\mathcal{X}$  in the t-th epoch, a class-aware pruning probability is assigned to each sample based on its score, which is formulated as:

$$\mathcal{P}_t^l(x_b^n) = \begin{cases} 1 & \mathcal{H}_t^l(x_b^n) \in \langle r_c \mathcal{H}_{c,t}^l, \\ 0 & \mathcal{H}_t^l(x_b^n) \notin \langle r_c \mathcal{H}_{c,t}^l, \end{cases}$$
(18)

where  $\forall r_c \mathcal{H}_{c,t}^l$  denotes the  $r_c \times N_c$  smallest scoring values of the class c and  $r_c = \pi_\theta(\mathcal{I})_c$  is the class-aware pruning probability. The labeled scoring function  $\mathcal{H}_{c,t}^l(x_b^n)$  is defined using the supervised loss  $\mathcal{L}_{sup}(x_b^n, y_b^n)$  to quantify sample utility. Specifically, we exploit the pruning policy to prune samples based on their scores. Then, for the pruned labeled samples, their scores remain unmodified as previously. For the remaining samples, their scores are updated by the losses in the current epoch. To ensure dynamic adaptation:

$$\mathcal{H}_{c,t+1}^{l}(x_b^n) = \begin{cases} \mathcal{H}_{c,t}^{l}(x_b^n) & x_b^n \in \mathcal{X}n\mathcal{S}^l, \\ \mathcal{L}_{sup}(x_b^n) & x_b^n \in \mathcal{S}^l. \end{cases}$$
(19)

**Dynamic pruning for unlabeled data.** While the distribution of the label of the unlabeled data and its imbalance ratio  $\gamma_u$  are unknown. To address the uncertainty and bias of pseudo-labels, we design a label-insensitive soft pruning policy  $\mathcal{P}^u_t$  inspired by (Qin et al., 2024), which introduces randomness and gradient scaling into the pruning process. Specifically, for an unlabeled dataset  $\mathcal{U}$  at the t-th epoch, a pruning probability is assigned to each sample based on its score, which is formulated as:

$$\mathcal{P}_t^u(u_b^m) = \begin{cases} r & \mathcal{H}_t^u(u_b^m) < \bar{\mathcal{H}}_t^m \text{ and } p^*(u_b^m) \ge \tau, \\ 0 & \mathcal{H}_t^u(u_b^m) \ge \bar{\mathcal{H}}_t^u \text{ or } p^*(u_b^m) < \tau, \end{cases}$$
(20)

where  $\bar{\mathcal{H}}^u_t$  is the adaptive threshold and r is a randomized pruning rate,  $\tau$  is the confidence threshold  $\tau$  and  $p^*(u^m_b) = \max(\text{softmax}(g^*_{\theta}(\alpha(u^m_b))))$  denote the debiased pseudo-label confidence. For

Table 1: Comparison of bACC/GM on CIFAR-10-LT.

Algorithm	CIFAR-10-LT (	$\gamma = \gamma_l = \gamma_u, \gamma_u$ is assum	ed to be known)
8	$\gamma = 50$	$\gamma = 100$	$\gamma = 150$
Vanilla	$65.2\pm0.05$ / $61.1\pm0.09$	$58.8 \pm 0.13 \ / \ 58.2 \pm 0.11$	55.6 ±0.43 / 44.0 ±0.98
Re-sampling	$64.3 \pm 0.48 / 60.6 \pm 0.67$	55.8 ±0.47 / 45.1 ±0.30	52.2 ±0.05 / 38.2 ±1.49
LDAM-DRW	$68.9 \pm 0.07$ / $67.0 \pm 0.08$	$62.8 \pm 0.17 / 58.9 \pm 0.60$	$57.9 \pm 0.20 / 50.4 \pm 0.30$
cRT	$67.8 \pm 0.13 / 66.3 \pm 0.15$	$63.2 \pm 0.45$ / $59.9 \pm 0.40$	$59.3\pm0.10$ / $54.6\pm0.72$
FixMatch	$79.2 \pm 0.33 \ / \ 77.8 \pm 0.36$	$71.5 \pm 0.72 / 66.8 \pm 1.51$	$68.4 \pm 0.15 / 59.9 \pm 0.43$
w/+DARP+cRT	$85.8 \pm 0.43$ / $85.6 \pm 0.56$	$82.4 \pm 0.26$ / $81.8 \pm 0.17$	$79.6 \pm 0.42 / 78.9 \pm 0.35$
w/+CReST+LA	$85.6 \pm 0.36 / 81.9 \pm 0.45$	$81.2 \pm 0.70$ / $74.5 \pm 0.99$	$71.9 \pm 2.24 / 64.4 \pm 1.75$
w/+ABC	$85.6 \pm 0.26$ / $85.2 \pm 0.29$	$81.1 \pm 1.14 / 80.3 \pm 1.29$	$77.3 \pm 1.25 / 75.6 \pm 1.65$
w/+CoSSL	$86.8 \pm 0.30$ / $86.6 \pm 0.25$	$83.2 \pm 0.49$ / $82.7 \pm 0.60$	$80.3 \pm 0.55$ / $79.6 \pm 0.57$
w/+SAW+LA	$86.2 \pm 0.15$ / $83.9 \pm 0.35$	$80.7 \pm 0.15$ / $77.5 \pm 0.21$	$73.7 \pm 0.06 / 71.2 \pm 0.17$
w/+Adsh	$83.4 \pm 0.06$ / $82.9 \pm 0.13$	$76.5 \pm 0.35 / 74.8 \pm 0.34$	$71.5\pm0.30$ / $68.8\pm0.35$
w/+DebiasPL	$85.6 \pm 0.20$ / $85.2 \pm 0.23$	$80.6 \pm 0.50$ / $79.9 \pm 0.57$	$76.6 \pm 0.12 / 75.8 \pm 0.71$
w/+UDAL	$86.5 \pm 0.29$ / $86.2 \pm 0.26$	$81.4 \pm 0.39$ / $80.6 \pm 0.38$	$77.9 \pm 0.33 / 75.8 \pm 0.71$
w/+L2AC	$86.6 \pm 0.31 / 86.7 \pm 0.30$	$82.1 \pm 0.57$ / $81.5 \pm 0.64$	$77.6 \pm 0.53$ / $75.8 \pm 0.71$
w/+CDMAD	$87.3 \pm 0.12 / 87.0 \pm 0.15$	$83.6 \pm 0.46$ / $83.1 \pm 0.57$	$80.8 \pm 0.86$ / $79.9 \pm 1.07$
w/+DYTRIM	<b>88.0</b> $\pm$ 0.31 / <b>87.8</b> $\pm$ 0.32	<b>84.8</b> $\pm$ 0.48 / <b>84.4</b> $\pm$ 0.51	<b>82.0</b> $\pm$ 0.09 / <b>81.3</b> $\pm$ 0.03
FlexMatch	$72.6 \pm 0.72 / 70.2 \pm 0.88$	$67.7 \pm 0.73 / 63.6 \pm 1.27$	$62.6 \pm 0.63 / 56.1 \pm 1.13$
w/+CDMAD	$74.4 \pm 0.82$ / $73.0 \pm 1.12$	$68.4 \pm 0.46$ / $66.8 \pm 0.53$	$67.0 \pm 0.52 / 63.2 \pm 0.44$
w/+DYTRIM	<b>77.2</b> $\pm$ 0.42 / <b>76.2</b> $\pm$ 0.44	<b>70.7</b> $\pm$ 0.49 / <b>67.8</b> $\pm$ 0.70	<b>68.6</b> $\pm$ 0.22 / <b>66.3</b> $\pm$ 0.07
FreeMatch	$71.9 \pm 0.24 / 69.4 \pm 0.61$	$65.7 \pm 0.18 / 60.9 \pm 0.69$	$62.5 \pm 0.12 / 57.3 \pm 0.53$
w/+CDMAD	$74.7 \pm 0.64 / 73.6 \pm 1.23$	$69.9 \pm 0.65$ / $68.2 \pm 0.74$	$66.2 \pm 0.27 / 63.2 \pm 0.44$
w/+DYTRIM	<b>76.9</b> $\pm$ 0.45 / <b>75.9</b> $\pm$ 0.52	<b>72.3</b> $\pm$ 0.12 / <b>71.4</b> $\pm$ 0.57	<b>69.4</b> ±0.35 / <b>67.5</b> ±0.63

a remaining sample with score  $\mathcal{H}^u_t(u^m_b) < \bar{\mathcal{H}}^m_t$ , whose corresponding pruning probability is r, its gradient is scaled to 1/(1-r) times of the original, otherwise the gradient remains unchanged. The score  $\mathcal{H}^u_{t+1}(u^m_b)$  is derived from the consistency regularization loss values  $\mathcal{L}_{con}(\alpha(u^m_b),\mathcal{A}(u^m_b))$  for unlabeled data points. To enhance pseudo-label reliability, we further apply a confidence threshold  $\tau$ , where only samples with  $p^*(u^m_b) > \tau$  contribute to  $\mathcal{L}_{con}$ , where  $\mathcal{L}_{con} = \frac{1}{B} \sum_{b=1}^B \mathbb{I}(p^*(u^m_b) > \tau) \mathbf{H}(P_\theta(y|\mathcal{A}(u^m_b),\hat{q}_b))$ . Thus, we formulate the update of  $\mathcal{H}^u_{t+1}(u^m_b)$  as:

$$\mathcal{H}_{t+1}^{u}(u_b^m) = \begin{cases} \mathcal{H}_t^{u}(u_b^m) & u_b^m \in \mathcal{U}n\mathcal{S}^u, \\ \mathcal{L}_{con}(u_b^m) & u_b^m \in \mathcal{S}^u. \end{cases}$$
 (21)

**Initialization:** at t = 0, scores  $\mathcal{H}_t^u$  and  $\mathcal{H}_t^l$  are all set to  $\{1\}$ , as no prior loss is available.

#### 5 EXPERIMENT

In this section, we conducted comprehensive experiments to verify the effectiveness of the proposed DyTrim on CIFAR10-LT, CIFAR100-LT (Cui et al., 2019), STL10-LT (Kim et al., 2020), and ImageNet-127 (Deng et al., 2009; Huh et al., 2016) datasets. Due to limited space, we defer the detailed experimental settings to the Appendix D.

# 5.1 Baselines

The classification performance of the DyTrim was compared with those of the following algorithms: 1. vanilla algorithm - Deep CNN trained with cross-entropy loss, 2. CIL algorithms - Resampling (JAPKOWICZ, 2000), LDAM-DRW (Cao et al., 2019), and cRT (Kang et al., 2020), 3. SSL algorithms - FixMatch (Sohn et al., 2020), and 4. CISSL algorithms - DARP, DARP+LA, DARP+cRT (Kim et al., 2020), CReST, CReST+LA (Wei & Gan, 2023), ABC (Lee et al., 2021), CoSSL (Fan et al., 2022), DASO (Oh et al., 2022), SAW, SAW+LA and SAW+cRT (Lai et al., 2022)

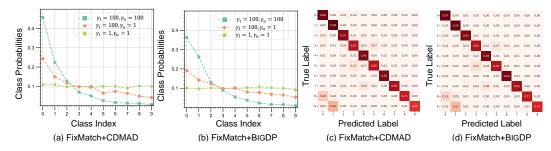


Figure 3: (a) and (b) present the  $\pi_{\theta}(\mathcal{I})$  using the CDMAD and DyTrim. (c) and (d) present the confusion matrices of the class predictions on test samples on CIFAR-10-LT ( $\gamma_l = \gamma_u = 100$ ). combined with FixMatch. Adsh(Guo & Li, 2022), DebiasPL (Wang et al., 2022), UDAL(Lazarow et al., 2023) and L2AC (Wang et al., 2023a) combined with FixMatch. We report the performance of the baseline algorithms reported in Tables of Lai et al. (2022) and Fan et al. (Fan et al., 2022) when it is reproducible; the performance measured using the uploaded code was reported otherwise.

#### 5.2 RESULTS ON CIFAR10/100-LT AND STL-LT

Under the consistent condition where  $\gamma_u$  is known and matched to  $\gamma_l$ , the results in Table 1 show that CISSL algorithms consistently outperform their vanilla SSL counterparts by mitigating class imbalance while effectively exploiting unlabeled data. Among them, the proposed DyTrim achieves the best performance across all imbalance ratios. Compared with the state-of-the-art CDMAD, DyTrim improves bACC by 1.2% and GM by 1.4% on average, without incurring additional computational overhead. Furthermore, when integrated into FlexMatch and FreeMatch, DyTrim yields substantial improvements, boosting bACC/GM by 2–3% on average.

Table 2 further evaluates the methods on CIFAR-100-

Table 2: Comparison of bACC on CIFAR-100-LT.

Algorithm	CIFAR-100	-LT ( $\gamma = \gamma_l =$	$\gamma_u, \gamma_u$ is assumed to be known)
	$\gamma = 20$	$\gamma = 50$	$\gamma = 100$
FixMatch	$49.6 \pm 0.78$	$42.1 \pm 0.33$	$37.6 \pm 0.48$
w/+DARP	$50.8 \pm 0.77$	$43.1 \pm 0.54$	$38.3 \pm 0.47$
w/+DARP+cRT	$51.4 \pm 0.68$	$44.9 \pm 0.54$	$40.4 \pm 0.78$
w/+CReST	$51.8 \pm 0.12$	$44.9 \pm 0.50$	$40.1 \pm 0.65$
w/+CReST+LA	$52.9 \pm 0.07$	$47.3 \pm 0.17$	$42.7 \pm 0.70$
w/+ABC	$53.3 \pm 0.79$	$46.7 \pm 0.26$	$41.2 \pm 0.06$
w/+CoSSL	$53.9 \pm 0.78$	$47.6 \pm 0.26$	$43.0 \pm 0.61$
w/+UDAL	$54.1 \pm 0.23$	$48.0 \pm 0.56$	$43.7 \pm 0.41$
w/+CPE	$52.4 \pm 0.17$	$45.6 \pm 0.68$	$39.9 \pm 0.40$
w/+CDMAD	$54.3 \pm 0.44$	$48.8 \pm 0.75$	$44.1 \pm 0.29$
w/+DYTRIM	<b>55.5</b> $\pm 0.53$	$50.8 \pm 0.80$	<b>44.8</b> ±0.27
FlexMatch	36.5 ±0.51	29.6 ±0.35	25.8 ±0.79
w/+CDMAD	$39.2 \pm 0.47$	$31.9 \pm 0.46$	$27.0 \pm 0.66$
w/+DYTRIM	<b>40.9</b> $\pm 0.09$	$33.5 \pm 0.21$	<b>29.8</b> ±0.67
FreeMatch	35.9 ±0.69	31.3 ±0.65	24.5 ±0.66
w/+CDMAD	$36.9 \pm 0.96$	$32.8 \pm 0.93$	$28.0 \pm 0.68$
w/+DYTRIM	<b>39.0</b> ±0.61	$33.4 \pm 0.70$	29.8 ±0.09

LT, which involves more classes and stronger imbalance. The results demonstrate that DyTrim consistently outperforms all competing approaches under this more challenging setting. In particular, DyTrim delivers clear gains over CDMAD across all imbalance levels, confirming its scalability to large-scale, fine-grained datasets. Similar to the observations on CIFAR-10-LT, the integration of DyTrim continues to provide consistent improvements across different SSL frameworks.

Under the inconsistent condition where  $\gamma_u$  was unknown and mismatched to  $\gamma_l$ , the results in Table 3 show that DyTrim remains the most effective method overall. When the labeled and unlabeled data distributions deviate, DyTrim consistently outperforms CDMAD on both CIFAR-10-LT and STL-10-LT. The benefits are particularly notable when combined with FlexMatch and FreeMatch: on STL-10-LT, FlexMatch + DyTrim improves bACC by more than 2%, while FreeMatch + DyTrim achieves nearly 2% gains.

In addition, Table 4 highlights the performance of various algorithms under both consistent and inconsistent imbalance settings with ViT backbones. On CIFAR-10-LT, DyTrim yields the best results, improving bACC 0.6% over CDMAD and nearly 4% over FixMatch when  $\gamma_l = \gamma_u = 100$ . Under the inconsistent condition, DyTrim maintains a clear margin, surpassing CDMAD almost 2%. On CIFAR-100-LT, although the absolute accuracies are lower due to the increased difficulty, DyTrim still matches or slightly improves upon CDMAD, while consistently outperforming FixMatch. Additional experimental results are provided in Appendix E.

#### 5.3 SCALABILITY EVALUATION OF DYTRIM

DyTrim exhibited robust extensibility as a universal plug-in component, consistently boosting performance across diverse SSL frameworks (CDMAD/CCL), datasets (CIFAR/STL10-LT), and im-

Table 3: Comparison of bACC/GM on CIFAR-10-LT and STL-10-LT ( $\gamma_l \neq \gamma_u, \gamma_u$  is assumed to be unknown).

Algorithm	CIFAR-10-LT ( $\gamma_l =$	100, $\gamma_u$ = Unknown)	STL-10-LT ( $\gamma$	u = Unknown
8	$\gamma_u = 50$	$\gamma_u = 150$	$\gamma_l = 10$	$\gamma_l = 20$
FixMatch	$73.9 \pm 0.25$ / $70.5 \pm 0.52$	$69.6 \pm 0.60 / 62.6 \pm 1.11$	$72.9 \pm 0.09 / 69.6 \pm 0.01$	$63.4 \pm 0.21 / 52.6 \pm 0.09$
w/+DARP	$77.3 \pm 0.17 / 75.5 \pm 0.21$	$72.9 \pm 0.24$ / $69.5 \pm 0.18$	$77.8 \pm 0.33 / 76.5 \pm 0.40$	$69.9 \pm 1.77 / 65.4 \pm 3.07$
w/+DARP+LA	$82.3 \pm 0.32 / 81.5 \pm 0.29$	$78.9 \pm 0.23$ / $77.7 \pm 0.06$	$78.6 \pm 0.30$ / $77.4 \pm 0.40$	$71.9 \pm 0.49 / 68.7 \pm 0.51$
w/+DARP+cRT	$82.7 \pm 0.21$ / $82.3 \pm 0.25$	$80.7 \pm 0.44$ / $80.2 \pm 0.61$	$79.3 \pm 0.23$ / $78.7 \pm 0.21$	$74.1 \pm 0.61 / 73.1 \pm 1.21$
w/+ABC	$82.7 \pm 0.64$ / $82.0 \pm 0.76$	$78.4 \pm 0.87$ / $77.2 \pm 1.07$	$79.1 \pm 0.46$ / $78.1 \pm 0.57$	$73.8 \pm 0.15$ / $72.1 \pm 0.15$
w/+SAW	$79.8 \pm 0.25$ / $79.1 \pm 0.32$	$74.5 \pm 0.97$ / $72.5 \pm 1.37$	$78.3 \pm 0.25 \ / \ 77.0 \pm 0.19$	$71.9 \pm 0.81 / 69.0 \pm 0.81$
w/+SAW+LA	$82.9 \pm 0.38$ / $82.6 \pm 0.38$	$79.1 \pm 0.81 \ / \ 78.6 \pm 0.91$	$79.4 \pm 0.26$ / $78.4 \pm 0.17$	$73.9 \pm 0.91 / 71.8 \pm 0.99$
w/+SAW+cRT	$81.6 \pm 0.38$ / $81.3 \pm 0.32$	$77.6 \pm 0.40$ / $77.1 \pm 0.41$	$78.9 \pm 0.22$ / $77.8 \pm 0.14$	$72.3 \pm 0.86 / 69.5 \pm 0.83$
w/+CPE	$86.2 \pm 0.26$ / $85.9 \pm 0.33$	$82.4 \pm 0.49$ / $82.1 \pm 0.53$	$79.0\pm0.05$ / $78.7\pm0.54$	$77.0\pm0.73$ / $76.1\pm0.68$
w/+CDMAD	$85.7 \pm 0.36 / 85.3 \pm 0.38$	$82.3 \pm 0.23$ / $81.8 \pm 0.29$	$79.9 \pm 0.23$ / $78.9 \pm 0.38$	$75.2 \pm 0.40 / 73.5 \pm 0.31$
w/+DyTrim	<b>86.4</b> $\pm$ 0.43 / <b>86.0</b> $\pm$ 0.43	<b>83.8</b> $\pm$ 0.34 / <b>83.4</b> $\pm$ 0.33	<b>80.7</b> $\pm$ 0.64 / <b>79.8</b> $\pm$ 0.70	<b>77.9</b> $\pm$ 1.04 / <b>76.7</b> $\pm$ 1.26
FlexMatch	67.7 ±0.67 / 62.8 ±0.65	63.0 ±0.77 / 56.3 ±1.70	62.1 ±0.29 / 60.8 ±0.43	56.9 ±0.90 / 51.4 ±0.81
w/+CDMAD	$69.2 \pm 0.22$ / $67.0 \pm 0.11$	$67.0 \pm 1.69$ / $63.4 \pm 0.91$	$65.5 \pm 1.05$ / $63.7 \pm 1.02$	$62.4 \pm 1.05 / 60.5 \pm 0.99$
w/+DyTrim	<b>72.5</b> $\pm$ 0.39 / <b>70.7</b> $\pm$ 0.45	<b>70.3</b> $\pm$ 1.01 / <b>67.4</b> $\pm$ 0.21	<b>68.0</b> $\pm$ 0.94 / <b>66.4</b> $\pm$ 0.85	<b>63.9</b> $\pm$ 0.16 / <b>61.7</b> $\pm$ 0.28
FreeMatch	69.3 ±0.99 / 65.4 ±1.45	63.5 ±0.76 / 55.7 ±0.77	63.9 ±0.77 / 62.0 ±0.90	59.0 ±1.43 / 57.6 ±0.67
w/+CDMAD	$71.0\pm0.98$ / $69.0\pm1.05$	$67.1 \pm 0.96$ / $64.3 \pm 0.99$	$66.1 \pm 0.32$ / $63.8 \pm 0.97$	$61.5 \pm 0.47 / 59.5 \pm 0.63$
w/+DyTrim	<b>72.3</b> $\pm$ 0.69 / <b>71.1</b> $\pm$ 1.23	<b>69.9</b> $\pm$ 0.15 / <b>67.4</b> $\pm$ 0.37	<b>68.0</b> $\pm$ 0.64 / <b>66.5</b> $\pm$ 1.20	<b>64.6</b> ±0.77 / <b>62.7</b> ±1.16

Table 4: Comparison of bACC/GM on CIFAR-10-LT and CIFAR-100-LT with ViT.

Algorithm	CIFAR	CIFAR-100-LT	
8	$\gamma_l = \gamma_u = 100$	$\gamma_l = 100, \gamma_u = 150$	$\gamma_l = \gamma_u = 100$
FixMatch	$45.5 \pm 0.14 / 30.0 \pm 0.41$	$45.3 \pm 0.12 / 28.9 \pm 0.96$	$23.2 \pm 0.13 / 5.7 \pm 0.33$
w/+CDMAD	$48.7 \pm 0.49$ / <b>40.5</b> $\pm 0.26$	$45.4\pm0.13$ / $39.9\pm0.10$	$24.0\pm0.15$ / <b>9.0</b> $\pm0.77$
w/+DyTrim	<b>49.3</b> ±0.47 / 40.3 ±0.36	<b>47.3</b> $\pm$ 0.12 / 39.7 $\pm$ 0.57	<b>24.1</b> $\pm 0.22$ / $8.9 \pm 0.15$

Table 5: Comparison of bACC with two state-of-the-art CISSL algorithms with and without DyTrim on CIFAR-10, CIFAR-100, and STL-10.

Dataset		FixMatch+			FixMatch+		
zuusee	ituset		CDMAD CDMAD+DyTrim		CCL CCL+DyTrim		Gain
CIFAR10-LT	$\gamma_l = \gamma_u = 100$ $\gamma_l = \gamma_u = 150$ $\gamma_l = 100, \gamma_u = 1$	$83.6 \pm 0.46$ $80.8 \pm 0.86$ $87.5 \pm 0.46$	<b>84.8</b> ±0.48 <b>82.0</b> ±0.09 <b>88.9</b> ±0.88	↑1.2 ↑1.2 ↑1.4	$86.2 \pm 0.35$ $84.0 \pm 0.21$ $93.9 \pm 0.12$	<b>86.7</b> ±0.39 <b>84.0</b> ±0.26 <b>94.1</b> ±0.17	↑0.5 ↑0.0 ↑0.2
CIFAR100-LT	$\gamma_l = \gamma_u = 20$	54.3 ±0.44	<b>55.5</b> ±0.53	↑1.2	57.5 ±0.16	<b>58.1</b> ±0.49	↑0.6
STL10-LT	$ \gamma_l = 10 \\ \gamma_l = 20 $	$79.9 \pm 0.23$ $75.2 \pm 0.40$	<b>80.7</b> ±0.64 <b>77.9</b> ±1.04	↑0.8 ↑2.7	$84.8 \pm 0.15$ $83.1 \pm 0.18$	<b>85.1</b> ±0.33 <b>83.3</b> ±0.40	↑0.3 ↑0.2

balance ratios ( $\gamma=1\sim150$ ). Notably, it achieved up to +1.4% (CDMAD on CIFAR10-LT) and +2.7% (STL10-LT,  $\gamma_l$ =20) gains without architecture-specific tuning, validating its versatility in semi-supervised long-tailed scenarios.

#### 6 CONCLUSION

In this work, we introduce DyTrim, a simple yet effective method to address long-tailed semisupervised learning (LTSSL) by leveraging baseline image-guided dynamic data pruning for debiasing. We provide a theoretical framework that explains the role of baseline images in reducing classifier bias under class imbalance, showing how the class imbalance datasets influence the model prediction of the baseline image. Our method, DyTrim, offers an elegant solution to the problem of class imbalance in semi-supervised learning without requiring complex network branches or debiasing mechanisms. Through extensive experiments, we demonstrate that DyTrim significantly improves the performance of LTSSL models on public benchmarks. We believe that our approach can serve as a powerful tool for improving LTSSL in real-world datasets where class imbalance remains a major challenge.

# REFERENCES

- Jean-Michel Attendu and Jean-Philippe Corbeil. Nlu on data diets: Dynamic data subset selection for nlp classification tasks. *arXiv preprint arXiv:2306.03208*, 2023. 14
- Léon Bottou. Stochastic gradient descent tricks. In *Neural networks: tricks of the trade: second edition*, pp. 421–436. Springer, 2012. 17
  - Jiarui Cai, Yizhou Wang, and Jenq-Neng Hwang. Ace: Ally complementary experts for solving long-tailed recognition in one-shot. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 112–121, 2021. 14
  - Kaidi Cao, Colin Wei, Adrien Gaidon, Nikos Arechiga, and Tengyu Ma. Learning imbalanced datasets with label-distribution-aware margin loss. Advances in neural information processing systems, 32, 2019. 7, 17
  - Dingshuo Chen, Zhixun Li, Yuyan Ni, Guibin Zhang, Ding Wang, Qiang Liu, Shu Wu, Jeffrey Xu Yu, and Liang Wang. Beyond efficiency: Molecular data pruning for enhanced generalization. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. 14
  - Hao Chen, Ran Tao, Yue Fan, Yidong Wang, Jindong Wang, Bernt Schiele, Xing Xie, Bhiksha Raj, and Marios Savvides. Softmatch: Addressing the quantity-quality tradeoff in semi-supervised learning. In *The Eleventh International Conference on Learning Representations*, 2023. 14
  - Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie. Class-balanced loss based on effective number of samples. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 9268–9277, 2019. 7
  - Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009. 7
  - Guodong Du, Jia Zhang, Ning Zhang, Hanrui Wu, Peiliang Wu, and Shaozi Li. Semi-supervised imbalanced multi-label classification with label propagation. *Pattern Recognition*, 150:110358, 2024. 14
  - Yue Fan, Dengxin Dai, Anna Kukleva, and Bernt Schiele. Cossl: Co-learning of representation and classifier for imbalanced semi-supervised learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 14574–14584, June 2022. 1, 7, 8, 17
  - Qianhan Feng, Lujing Xie, Shijie Fang, and Tong Lin. Bacon: Boosting imbalanced semi-supervised learning via balanced feature-level contrastive learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 11970–11978, 2024. 14
  - Lan-Zhe Guo and Yu-Feng Li. Class-imbalanced semi-supervised learning with adaptive thresholding. In *International conference on machine learning*, pp. 8082–8094. PMLR, 2022. 8, 14, 17
  - Xiaoyu Guo, Xiang Wei, Shunli Zhang, Wei Lu, and Weiwei Xing. Dcrp: Class-aware feature diffusion constraint and reliable pseudo-labeling for imbalanced semi-supervised learning. *IEEE Transactions on Multimedia*, 26:7146–7159, 2024. 17
  - Yangyang Guo and Mohan Kankanhalli. Scan: Bootstrapping contrastive pre-training for data efficiency. *arXiv preprint arXiv:2411.09126*, 2024. 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*, pp. 770–778, 2016. 17
  - Chen Huang, Yining Li, Chen Change Loy, and Xiaoou Tang. Learning deep representation for imbalanced classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 5375–5384, 2016. 17

- Minyoung Huh, Pulkit Agrawal, and Alexei A Efros. What makes imagenet good for transfer learning? *arXiv preprint arXiv:1608.08614*, 2016. 7
- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. *Advances in neural information processing systems*, 31, 2018.
- N JAPKOWICZ. The class imbalance problem: Significance and strategies. In *Proc. 2000 International Conference on Artificial Intelligence*, volume 1, pp. 111–117, 2000. 7, 17
- Bingyi Kang, Saining Xie, Marcus Rohrbach, Zhicheng Yan, Albert Gordo, Jiashi Feng, and Yannis Kalantidis. Decoupling representation and classifier for long-tailed recognition. In *International Conference on Learning Representations*, 2020. 7, 17
- Kenji Kawaguchi and Haihao Lu. Ordered sgd: A new stochastic optimization framework for empirical risk minimization. In *International Conference on Artificial Intelligence and Statistics*, pp. 669–679. PMLR, 2020. 14
- Krishnateja Killamsetty, Sivasubramanian Durga, Ganesh Ramakrishnan, Abir De, and Rishabh Iyer. Grad-match: Gradient matching based data subset selection for efficient deep model training. In *International Conference on Machine Learning*, pp. 5464–5474. PMLR, 2021. 14
- 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. 1, 7, 17
- Miroslav Kubat. Addressing the curse of imbalanced training sets: one-sided selection. In *Proceedings of the 14th international conference on machine learning*, pp. 179–186. Morgan Kaufmann, 1997. 17
- Zhengfeng Lai, Chao Wang, Henrry Gunawan, Sen-Ching S Cheung, and Chen-Nee Chuah. Smoothed adaptive weighting for imbalanced semi-supervised learning: Improve reliability against unknown distribution data. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162, pp. 11828–11843, 2022. 7, 8, 17
- Justin Lazarow, Kihyuk Sohn, Chen-Yu Lee, Chun-Liang Li, Zizhao Zhang, and Tomas Pfister. Unifying distribution alignment as a loss for imbalanced semi-supervised learning. In 2023 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV), pp. 5633–5642. IEEE Computer Society, 2023. 8, 17
- 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*, pp. 3911–3920, June 2022. 14
- Hyuck Lee and Heeyoung Kim. Cdmad: Class-distribution-mismatch-aware debiasing for class-imbalanced semi-supervised learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 23891–23900, 2024. 1, 17
- Hyuck Lee, Seungjae Shin, and Heeyoung Kim. Abc: Auxiliary balanced classifier for class-imbalanced semi-supervised learning. In *Advances in Neural Information Processing Systems*, volume 34, pp. 7082–7094, 2021. 1, 7, 14, 17
- Jingyang Li, Jiachun Pan, Vincent Y. F. Tan, Kim chuan Toh, and Pan Zhou. Towards understanding why fixmatch generalizes better than supervised learning. In *The Thirteenth International Conference on Learning Representations*, 2025. 1
- Chengcheng Ma, Ismail Elezi, Jiankang Deng, Weiming Dong, and Changsheng Xu. Three heads are better than one: Complementary experts for long-tailed semi-supervised learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 14229–14237, 2024. 14
- Aditya Krishna Menon, Sadeep Jayasumana, Ankit Singh Rawat, Himanshu Jain, Andreas Veit, and Sanjiv Kumar. Long-tail learning via logit adjustment. In *Proceedings of the International Conference on Learning Representations*, 2021. 14

- Youngtaek Oh, Dong-Jin Kim, and In So Kweon. Daso: Distribution-aware semantics-oriented pseudo-label for imbalanced semi-supervised learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 9786–9796, 2022. 7, 17
- Ziheng Qin, Kai Wang, Zangwei Zheng, Jianyang Gu, Xiangyu Peng, Daquan Zhou, Lei Shang, Baigui Sun, Xuansong Xie, Yang You, et al. Infobatch: Lossless training speed up by unbiased dynamic data pruning. In *The Twelfth International Conference on Learning Representations*, 2024. 6, 14, 17, 18
- Ravi S Raju, Kyle Daruwalla, and Mikko Lipasti. Accelerating deep learning with dynamic data pruning, 2021. 14
- Yi Ren and Danica J Sutherland. Learning dynamics of llm finetuning. *arXiv preprint* arXiv:2407.10490, 2024. 2, 3, 21
- Yi Ren, Shangmin Guo, and Danica J. Sutherland. Better supervisory signals by observing learning paths. In *International Conference on Learning Representations*, 2022. 21
- Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. *arXiv preprint arXiv:1911.08731*, 2019. 14
- Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized experience replay. *arXiv* preprint arXiv:1511.05952, 2015. 14
- 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. *arXiv preprint arXiv:2001.07685*, 2020. 1, 2, 3, 7, 14, 17
- Truong Thao Nguyen, Balazs Gerofi, Edgar Josafat Martinez-Noriega, François Trahay, and Mohamed Wahib. Kakurenbo: adaptively hiding samples in deep neural network training. *Advances in Neural Information Processing Systems*, 36:37900–37922, 2023. 14
- Renzhen Wang, Xixi Jia, Quanziang Wang, Yichen Wu, and Deyu Meng. Imbalanced semisupervised learning with bias adaptive classifier. In *The Eleventh International Conference on Learning Representations*, 2023a. 8, 17
- Xudong Wang, Zhirong Wu, Long Lian, and Stella X Yu. Debiased learning from naturally imbalanced pseudo-labels. In 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 14627–14637. IEEE, 2022. 8, 14, 17
- 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. arXiv preprint arXiv:2205.07246, 2023b. 2, 14, 15
- Chen Wei, Kihyuk Sohn, Clayton Mellina, Alan Yuille, and Fan Yang. Crest: A class-rebalancing self-training framework for imbalanced semi-supervised learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10857–10866, 2021. 1, 14
- Tong Wei and Kai Gan. Towards realistic long-tailed semi-supervised learning: Consistency is all you need. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3469–3478, 2023. 1, 7, 14, 17
- Kan Wu, Jinnian Zhang, Houwen Peng, Mengchen Liu, Bin Xiao, Jianlong Fu, and Lu Yuan.Tinyvit: Fast pretraining distillation for small vision transformers. In *European conference on computer vision*, pp. 68–85. Springer, 2022. 17
- Liuyu Xiang, Guiguang Ding, and Jungong Han. Learning from multiple experts: Self-paced knowledge distillation for long-tailed classification. In *16th European Conference on Computer Vision*, *ECCV 2020*, pp. 247–263. Springer Nature, 2020. 17

- Weiwei Xing, Yue Cheng, Hongzhu Yi, Xiaohui Gao, Xiang Wei, Xiaoyu Guo, Yumin Zhang, and Xinyu Pang. Lcgc: Learning from consistency gradient conflicting for class-imbalanced semi-supervised debiasing. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 21697–21706, 2025. 1, 14, 17, 21
- Zhuoran Yu, Yin Li, and Yong Jae Lee. Inpl: Pseudo-labeling the inliers first for imbalanced semi-supervised learning. *arXiv preprint arXiv:2303.07269*, 2023. 14
- Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016. 16
- 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 *Advances in Neural Information Processing Systems*, volume 34, pp. 18408–18419, 2021. 1, 2, 14, 15
- Guibin Zhang, Haonan Dong, Yuchen Zhang, Zhixun Li, Dingshuo Chen, Kai Wang, Tianlong Chen, Yuxuan Liang, Dawei Cheng, and Kun Wang. Gder: Safeguarding efficiency, balancing, and robustness via prototypical graph pruning. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. 14
- Na Zheng, Xuemeng Song, Xue Dong, Aashish Nikhil Ghosh, Liqiang Nie, and Roger Zimmermann. Language-assisted debiasing and smoothing for foundation model-based semi-supervised learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 25708–25717, June 2025. 1
- Zi-Hao Zhou, Siyuan Fang, Zi-Jing Zhou, Tong Wei, Yuanyu Wan, and Min-Ling Zhang. Continuous contrastive learning for long-tailed semi-supervised recognition. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. 1, 14

# **APPENDIX**

# A RELATED WORK

#### A.1 MECHANISMS OF LONG-TAILED DEBIASING

This paper considers learning dynamics to study the debiasing mechanisms of SSL algorithms. We briefly introduce differences between the settings considered here and those in previous works. For debiasing on long-tailed learning, Menon et al. (2021) considered a unified framework for debiasing from the perspective of logits adjustment, which requires statistical label frequency. CCL (Zhou et al., 2024) considered debiasing from an information-theoretical lens. LCGC (Xing et al., 2025) used gradient flow to analyze the debiasing process. However, these methods only elucidate the model's behavior from an ad hoc perspective. We aim to develop a more comprehensive framework that enables a principle-based lens of the bias generation mechanisms inherent in long-tailed semi-supervised learning.

# A.2 SEMI-SUPERVISED LEARNING

Modern SSL methods typically integrate diverse strategies for exploiting unlabeled data, such as entropy minimization (Zhou et al., 2024), consistency regularization (Sohn et al., 2020), and contrastive learning (Zhou et al., 2024; Lee et al., 2022). Among them, most SSL approaches rely on selecting reliable pseudo-labels during training. FixMatch (Sohn et al., 2020) adopts a fixed confidence threshold of 0.95, whereas FlexMatch (Zhang et al., 2021) adapts thresholds per class based on learning difficulty and training progress. FreeMatch (Wang et al., 2023b) integrates global and local adjustments with a class-fairness regularizer to promote prediction diversity, while Soft-Match (Chen et al., 2023) employs a soft thresholding scheme that reweights samples to balance quantity and quality. In contrast, our method bypasses threshold tuning altogether and directly enforces class-balanced pseudo-labeling through dynamic pruning.

#### A.3 Long-tailed semi-supervised debiasing

Existing debiasing methods for LTSSL dominantly rely on consistent distribution assumptions (Guo & Li, 2022; Lee et al., 2021) and logit adjustment strategies (Wei & Gan, 2023). Notable approaches include CReST (Wei et al., 2021), which focuses on minority classes through selective self-training, and CoSSL (Cai et al., 2021), which balances representations using tail-class feature augmentation. Recent advances, like BaCon (Feng et al., 2024), utilize contrastive learning for balanced features, while SMCLP (Du et al., 2024) exploits collaborative label-instance correlations, and CPE (Ma et al., 2024) employs multiple expert classifiers. Innovative methods such as InPL (Yu et al., 2023) and DebiasMatch (Wang et al., 2022) move beyond traditional pseudo-labeling; InPL uses energy scores to detect reliable inliers, whereas DebiasMatch applies adaptive debiasing with a marginal loss to reduce long-tailed pseudo-label bias. Despite these advances, LTSSL techniques often demand intricate mechanisms or additional modules (Lee et al., 2021), posing challenges in minimizing bias while maintaining simplicity.

#### A.4 DYNAMIC DATASET PRUNING

To reduce training cost on datasets, dynamic dataset pruning methods (Chen et al., 2024; Killamsetty et al., 2021; Sagawa et al., 2019; Schaul et al., 2015; Zhang et al., 2024) aim to reduce the number of training iterations while maintaining performance. Existing methods employ a variety of criteria to guide pruning, among which loss-based (Attendu & Corbeil, 2023; Kawaguchi & Lu, 2020; Thao Nguyen et al., 2023) method is the most popular. UCB (Raju et al., 2021) applies the cross-entropy loss with exponential moving average (EMA) smoothing to mitigate noise. Infobatch (Qin et al., 2024) randomly prunes low-loss samples and amplifies the gradients of retained ones to preserve the expected gradient. SCAN (Guo & Kankanhalli, 2024) categorizes samples as redundant or ill-matched based on their loss and gradually increases the pruning ratio using cosine annealing. While these methods effectively accelerate training and can yield nearly unbiased results, none have explored their potential to mitigate class imbalance in SSL by pruning.

# B MORE BASE SSL ALGORITHMS

#### B.1 MORE ABOUT TRAINING LOSSES OF FIXMATCH

Training losses of FixMatch on a minibatch for the labeled set  $\mathcal{MX}$  and a minibatch for the unlabeled set  $\mathcal{MU}$  can be expressed as follows:

$$\mathcal{L}_{sup}(x_b; \theta) = \frac{1}{B} \sum_{x_b \in \mathcal{MX}} \mathbf{H} \left( \pi_{\theta}(y | \alpha(x_b)), p_b \right)$$
 (22)

with

$$\mathcal{L}_{con}(u_b, \hat{q}, \tau; \theta) = \frac{1}{\mu B} \sum_{b=1}^{B} \mathbb{1}(\max(\hat{q}_b) \ge \tau) \mathbf{H}(P_{\theta}(y|\mathcal{A}(u_b), \hat{q}_b),$$
(23)

where  $\hat{q}$  denote the concatenations of  $\hat{q}_b$ .  $\mathcal{L}_{sup}$  denotes the supervised loss for weakly augmented labeled data points  $u_b$ .  $\mathcal{L}_{con}$  denotes the consistency regularization loss with the confidence threshold  $\tau$ .

#### B.2 MORE ABOUT FLEXMATCH

To overcome the limitation of FixMatch using a fixed threshold  $\tau$  across all classes, Flex-Match (Zhang et al., 2021) introduces the *Curriculum Pseudo Labeling (CPL)* strategy. The key idea is to dynamically adjust the confidence threshold according to the learning status of each class. Specifically, FlexMatch first predicts the class probability for a weakly augmented unlabeled sample  $u_b$  as  $q_b = \pi_\theta(y|\alpha(u_b))$ , and then estimates the learning effect of each class c by  $\sigma_t(c)$ , i.e., the number of samples predicted as class c that exceed the fixed threshold  $\tau$ . After normalization, a ratio coefficient  $\beta_t(c)$  is obtained, which defines the class-adaptive threshold:

$$T_t(c) = \beta_t(c) \cdot \tau. \tag{24}$$

In this way, hard-to-learn classes receive a lower threshold to include more samples in training, while easy-to-learn classes gradually increase their thresholds to ensure pseudo-label quality. The unsupervised loss is defined as:

$$\mathcal{L}_{con}(u_b, \hat{q}, T_t; \theta) = \frac{1}{\mu B} \sum_{b=1}^{\mu B} \mathbb{1}(\max(q_b) > T_t(\arg\max(q_b))) \ \mathbf{H}(\hat{q}_b, \pi_{\theta}(y|\mathcal{A}(u_b))), \tag{25}$$

where  $\hat{q}_b = \arg\max_c q_{b,c}$  denotes the hard pseudo-label, and  $\mathcal{A}(\cdot)$  is the strong augmentation function. The overall training objective is

$$\mathcal{L}_t = \mathcal{L}_{sup} + \lambda \mathcal{L}_{con}. \tag{26}$$

where  $\lambda$  is weighting hyperparameter.

#### B.3 MORE ABOUT FREEMATCH

Unlike FixMatch and FlexMatch, which rely on fixed or indirectly adjusted thresholds, FreeMatch (Wang et al., 2023b) proposes *Self-Adaptive Thresholding (SAT)* that dynamically determines thresholds based on the model's prediction confidence. Specifically, FreeMatch first estimates a global threshold  $\tau_t$  using an exponential moving average (EMA) of model confidence:

$$\tau_t = \rho \tau_{t-1} + (1 - \rho) \frac{1}{\mu B} \sum_{b=1}^{\mu B} \max(q_b), \tag{27}$$

and further refines it with class-specific local statistics  $\tilde{p}_t(c)$ :

$$\tau_t(c) = \frac{\tilde{p}_t(c)}{\max_{c'} \tilde{p}_t(c')} \cdot \tau_t.$$
 (28)

At the early stage of training, thresholds are low to encourage more unlabeled data utilization and faster convergence. As the model becomes more confident, thresholds increase to suppress incorrect pseudo-labels and reduce confirmation bias. The unsupervised loss at iteration t is thus:

$$\mathcal{L}_{con}(u_b, \hat{q}, \tau_t; \theta) = \frac{1}{\mu B} \sum_{b=1}^{\mu B} \mathbb{1}(\max(q_b) > \tau_t(\arg\max(q_b))) \ \mathbf{H}(\hat{q}_b, \pi_{\theta}(y|\mathcal{A}(u_b))). \tag{29}$$

In addition, FreeMatch introduces Self-Adaptive Fairness (SAF) regularization  $\mathcal{L}_f$ , which dynamically calibrates the prediction distribution to encourage diverse predictions and prevent class collapse during early training. Concretely, let  $h_t \in \mathbb{R}^C$  denotes the normalized class histogram of model predictions at iteration t, and let  $h^* \in \mathbb{R}^C$  denotes the target distribution (e.g., a uniform distribution). The SAF regularization is defined as

$$\mathcal{L}_f = D_{\mathrm{KL}}(h_t \parallel h^*) \,, \tag{30}$$

where  $D_{KL}(\cdot||\cdot)$  is the Kullback–Leibler divergence. The final training objective is:

$$\mathcal{L} = \mathcal{L}_{sup} + w_u \mathcal{L}_{con} + w_f \mathcal{L}_f, \tag{31}$$

where  $w_u$  and  $w_f$  are weighting hyperparameters.

# C PSEUDO CODE OF THE PROPOSED ALGORITHM

The pseudo-code that describes the DyTrim is presented in Algorithm 1 and Algorithm 2.

# Algorithm 1 DyTrim for Labeled Data Selection

**Input:** Labeled set of N samples  $\mathcal{X} = \{(x^n, y^n)\}_{n=1}^N$ , score set of the samples  $\mathcal{V}^l$ , number of classes  $n_c$ , biased degree b

**Output:** Labeled pruned set  $S^l$  ( $S^l \subseteq \mathcal{X}$ ,  $|S^l| <= |\mathcal{X}|$ )

```
1: \mathcal{S}^l \leftarrow \emptyset 
ightharpoonup Initialize the labeled pruned set
```

2: **for** c = 0 to  $n_c - 1$  **do** 

3: 
$$\mathcal{I}_c \leftarrow \{i \mid y_i = c\}$$
  
4:  $\mathcal{V}_c^l \leftarrow \{\mathcal{V}_i^l \mid i \in \mathcal{I}_c\}$   $\triangleright$  Select scores of class  $c$  samples

5: 
$$k_c \leftarrow \lfloor (1 - b_c) \cdot |\hat{\mathcal{X}}_c| \rfloor$$
  $\triangleright$  Compute target pruned set size of class  $c$  based on biased degree  
6:  $\mathcal{I}_c^{\text{top}} \leftarrow \text{TopK}(\mathcal{I}_c, \mathcal{V}_c^l, k_c)$   $\triangleright$  Select indices of top- $k_c$  scored samples

7: 
$$\mathcal{S}^l \leftarrow \mathcal{S}^l \cup \mathcal{I}_c^{\text{top}}$$

8: end for

9: **return**  $S^l$ 

#### Algorithm 2 DyTrim for Unlabeled Data Selection

**Input:** Unlabeled set of M samples  $\mathcal{U} = \{(u^m)\}_{m=1}^M$ , score set of the samples  $\mathcal{V}^u$ , pruning ratio r, weight of samples w

**Output:** Unlabeled pruned set  $S^l$  ( $S^l \subseteq U$ ,  $|S^u| <= |U|$ )

3:  $\mathcal{I}_{\neq 0} \leftarrow \{i \mid \mathcal{V}_i^u \neq 0\}$ 4:  $\mathcal{S}^u \leftarrow \mathcal{S}^u \cup \mathcal{I}_0$ 

5:  $\mu \leftarrow \text{Mean}(\{\mathcal{V}_i^u \mid i \in \mathcal{I}_{\neq 0}\})$ 

6: 
$$\mathcal{I}_{well} \leftarrow \{i \in \mathcal{I}_{\neq 0} \mid \mathcal{V}_i^u < \mu\}$$
  $\triangleright$  Select well-learned samples 7:  $\mathcal{I}_{poor} \leftarrow \mathcal{I}_{\neq 0} \setminus \mathcal{I}_{well}$   $\triangleright$  Select poorly-learned samples

7:  $\mathcal{I}_{\text{poor}} \leftarrow \mathcal{I}_{\neq 0} \setminus \mathcal{I}_{\text{well}}$ 8:  $\mathcal{S}^u \leftarrow \mathcal{S}^u \cup \mathcal{I}_{\text{poor}}$ 

9:  $\mathcal{I}_{\text{select}} \leftarrow \text{Randomly select } | (1-r) \cdot |\mathcal{I}_{\text{well}} | | \text{ samples from } \mathcal{I}_{\text{well}}$ 

10:  $S^u \leftarrow S^u \cup \mathcal{I}_{\text{select}}$ 

```
11: w_i \leftarrow 1, \ \forall i \in \{1, \dots, M\}
12: w_i \leftarrow \frac{1}{1-T}, \ \forall i \in \mathcal{I}_{\text{select}}

\triangleright Reset weights
\triangleright Rescaling
```

13: return  $\mathcal{S}^{u}$ 

#### D EXPERIMENTAL SETTINGS

#### D.1 MODELS

Unless otherwise specified, we adopt Wide ResNet (WRN) (Zagoruyko & Komodakis, 2016) as the default backbone following common practice in semi-supervised learning. Additionally, we also

evaluate Tiny Vision Transformers (TinyViT) (Wu et al., 2022) on CIFAR-10-LT and CIFAR-100-LT. For ImageNet-127, we employ ResNet-50 (He et al., 2016) as the backbone to ensure scalability on large-scale datasets.

#### D.2 BASELINES

 The classification performance of the DyTrim was compared with those of the following algorithms: 1. vanilla algorithm - Deep CNN trained with cross-entropy loss, 2. CIL algorithms - Resampling (JAPKOWICZ, 2000), LDAM-DRW (Cao et al., 2019), and cRT (Kang et al., 2020), 3. SSL algorithms - FixMatch (Sohn et al., 2020), and 4. CISSL algorithms - DARP, DARP+LA, DARP+cRT (Kim et al., 2020), CReST, CReST+LA (Wei & Gan, 2023), ABC (Lee et al., 2021), CoSSL (Fan et al., 2022), DASO (Oh et al., 2022), SAW, SAW+LA and SAW+cRT (Lai et al., 2022) combined with FixMatch. Adsh(Guo & Li, 2022), DebiasPL (Wang et al., 2022), UDAL(Lazarow et al., 2023) and L2AC (Wang et al., 2023a) combined with FixMatch. We report the performance of the baseline algorithms reported in Tables of Lai et al. (2022) and Fan et al. (Fan et al., 2022) when it is reproducible; the performance measured using the uploaded code was reported otherwise.

# D.3 IMPLEMENTATION DETAILS

All experiments are trained for 500 epochs with 500 steps per epoch, resulting in a total of 250,000 iterations. We use Stochastic Gradient Descent (SGD) (Bottou, 2012) with a fixed learning rate of  $\eta=0.0015$  and a batch size of 32. The pruning ratio of the unlabeled dataset is set to 0.7, and the parameter  $\delta$  is aligned with InfoBatch (Qin et al., 2024), fixed at 0.875. For CIFAR-10-LT, the largest labeled class contains 1,500 samples, while the largest unlabeled class contains 3,000 samples. For CIFAR-100-LT, the largest labeled and unlabeled classes contain 150 and 300 samples, respectively. For STL-10-LT, the largest labeled class contains 450 samples. To assess classification performance, we adopt balanced accuracy (bACC) (Huang et al., 2016) and geometric mean (GM) (Kubat, 1997) for CIFAR-10-LT and STL-10-LT. For CIFAR-100-LT and ImageNet-127, evaluation is conducted solely using bACC. Each experiment is repeated three times on RTX 4090 GPUs to ensure reproducibility, and we report both the mean and the standard error.

#### E ADDITIONAL EXPERIMENTAL RESULTS

#### E.1 ADDITIONAL RESULTS ON CIFAR-10-LT

Following prior works (Xing et al., 2025; Lee & Kim, 2024; Guo et al., 2024), we evaluate under a more challenging scenario where the unlabeled set is imbalanced in the reverse direction of the labeled set (Table 6). Across all settings, DyTrim delivers consistent gains by applying balanced pruning on the labeled data. Notably, when combined with FixMatch, DyTrim surpasses CDMAD by more than 1% in both bACC and GM. Similar benefits are observed for FlexMatch and FreeMatch: DyTrim improves FlexMatch by approximately 1.1–1.3% and FreeMatch by around 0.9–1.5%.

Table 6: Comparison of bACC/GM on CIFAR-10-LT( $\gamma_l = 100, \gamma_u = 100 \text{(reversed)})$ .

Algorithm	CIFAR-10-LT, $\gamma_l = 100, \gamma_u = 100 \text{(reversed)}$					
8	ABC	SAW	SAW+LA	SAW+cRT	CDMAD	DyTrim
FixMatch+	69.5/66.8	72.3/68.7	74.1/72.0	75.5/73.9	77.1/75.4	78.2 / 76.7
FlexMatch+	-/-	-/-	-/-	-/-	67.2/65.1	68.3 / 66.4
FreeMatch+	-/-	-/-	-/-	-/-	68.5/66.4	69.4 / 67.9

We also compared the classification performance of CDMAD with ACR (Xiang et al., 2020) and BaCon, two recent CISSL algorithms. From Table. 7, we can observe that CDMAD outperforms both ACR and BaCon.

#### E.2 RESULTS ON SMALL-IMAGENET-127

Table 7: Comparison of bACC/GM on CIFAR-10-LT

Algorithm/CIFAR-10-LT	$\gamma_l = \gamma_u = 100$	$\gamma_l = \gamma_u = 1$
FixMatch+ACR	81.8 / 81.4	85.6 / 85.3
FixMatch+BaCon	84.4 / 84.0	82.0 / 81.5
FixMatch+CDMAD	83.6 / 83.1	87.5 / 87.1
FixMatch+DyTrim	84.8 / 84.4	87.9 / 87.5

ImageNet-127 is a naturally long-tailed dataset, widely used to evaluate class-imbalanced semi-supervised learning (CISSL) algorithms at scale. Following standard protocol, we downsample images to resolutions of  $32\times32$  and  $64\times64$  using the box interpolation method from the Pillow library, and randomly select 10% of the training samples as labeled data. Under such limited supervision and class imbalance, learning discriminative representations and a balanced classifier is particularly challenging. As reported in Table. 8, DyTrim achieves the highest balanced accuracy (bACC) at both resolutions, outperforming the strongest baseline CDMAD by 3.0% at  $32\times32$  and 1.2% at  $64\times64$ . These improvements demonstrate the robustness of our method, especially under

Table 8: Comparison of bACC on Small-ImageNet-127.

Algorithm	Small-ImageNet-127		
i ingoriumi	$32 \times 32$	$64 \times 64$	
FixMatch	29.7	42.3	
w/+DARP	30.5	42.5	
w/+DARP+cRT	39.7	51.0	
w/+CReST	32.5	44.7	
w/+CReST+LA	40.9	55.9	
w/+ABC	46.9	56.1	
w/+CoSSL	43.7	53.8	
w/+CPE	47.8	58.2	
w/+CDMAD	48.4	59.3	
w/+DyTrim	50.6	60.0	

low-resolution and low-resource conditions. The performance gain at lower resolutions suggests that DyTrim effectively handles the compounded difficulty of reduced visual fidelity and severe label scarcity. This makes it a promising solution for real-world applications where high-resolution data and abundant labels are often unavailable.

#### E.3 RESULTS ON DYNAMIC DATA PRUNING EXPERIMENT

Recently, Infobatch (Qin et al., 2024) provides a no-bias dynamic data pruning method. In this section, we compare it with DyTrim in the framework of CISSL. The experiment is conducted on the CIFAR-10-LT dataset, comparing the settings of  $\gamma_l = \gamma_u$  and  $\gamma_l \neq \gamma_u$ . Specifically, we directly apply the pruning policy of InfoBatch to labeled samples and unlabeled samples without distinction, and the results are shown in the Table. 9 and Table. 10. It can be seen that compared with the proposed method, the pruning policy directly combined with InfoBatch is not consistently effective in all settings. In particular, when  $\gamma_l \neq \gamma_u$ , it will cause a decrease in accuracy, which is caused by the mismatch in the distribution of labeled samples and unlabeled samples.

Table 9: Comparison of bACC/GM on CIFAR-10-LT.

	CIFAR-10-LT ( $\gamma = \gamma_l = \gamma_u, \gamma_u$ is assumed to be known)				
Method	$\gamma_l = 50, \gamma_u = 50$	$\gamma_l = 100, \gamma_u = 100$	$\gamma_l = 150, \gamma_u = 150$		
FixMatch	$79.2 \pm 0.33$ / $77.8 \pm 0.36$	$71.5\pm0.72$ / $66.8\pm1.51$	$68.4 \pm 0.15 / 59.9 \pm 0.43$		
w/+CDMAD	$87.3 \pm 0.12 / 87.0 \pm 0.15$	$83.6 \pm 0.46$ / $83.1 \pm 0.57$	$80.8 \pm 0.86$ / $79.9 \pm 1.07$		
w/+InfoBatch	$87.2 \pm 0.18$ / $86.9 \pm 0.19$	$84.1 \pm 0.61$ / $83.7 \pm 0.69$	$81.6 \pm 0.45$ / $80.9 \pm 0.59$		
w/+DyTrim	<b>88.0</b> $\pm$ 0.31 / <b>87.8</b> $\pm$ 0.32	<b>84.8</b> $\pm$ 0.48 / <b>84.4</b> $\pm$ 0.51	<b>82.0</b> $\pm$ 0.09 / <b>81.3</b> $\pm$ 0.03		

Table 10: Comparison of bACC/GM on CIFAR-10-LT ( $\gamma_l \neq \gamma_u, \gamma_u$  is assumed to be unknown).

	CIFAI	R-10-LT ( $\gamma_l = 100$ , $\gamma_u = \text{Unk}$	nown)
Method	$\gamma_u = 1$	$\gamma_u = 50$	$\gamma_u = 150$
FixMatch	$68.9 \pm 1.95$ / $42.8 \pm 8.11$	$73.9\pm0.25$ / $70.5\pm0.52$	$69.6 \pm 0.60 / 62.6 \pm 1.11$
w/+CDMAD	$87.5 \pm 0.46 / 87.1 \pm 0.50$	$85.7 \pm 0.36$ / $85.3 \pm 0.38$	$82.3 \pm 0.23$ / $81.8 \pm 0.29$
w/+InfoBatch	$86.4 \pm 0.63$ / $85.9 \pm 0.73$	$85.5 \pm 0.33$ / $85.1 \pm 0.37$	$83.3 \pm 0.08$ / $82.8 \pm 0.11$
w/+DyTrim	<b>88.9</b> $\pm 0.88$ / <b>88.6</b> $\pm 1.03$	<b>86.4</b> $\pm$ 0.43 / <b>86.0</b> $\pm$ 0.43	<b>83.8</b> $\pm$ 0.34 / <b>83.4</b> $\pm$ 0.33

Table 11: Ablation study for the proposed algorithm on CIFAR-10-LT.

Labeled	Unlabeled	Rescaling	$\gamma_l = \gamma_u$	= 50	$\gamma_l = \gamma_u$	= 100	$\gamma_l = \gamma_u$	= 150
Pruning	Pruning	researing	bACC		bACC	GM	bACC	GM
			87.3	87.0	83.6	83.1	80.8	79.9
$\checkmark$			87.5	87.2	84.4	84.0	81.3	80.6
	$\checkmark$	$\checkmark$	87.7	87.4	84.0	83.6	81.4	80.6
$\checkmark$	$\checkmark$		87.2	86.9	83.6	83.1	79.9	79.0
✓	$\checkmark$	$\checkmark$	88.0	87.8	84.8	84.4	82.0	81.3

Table 12: Comparison of bACC/GM on CIFAR-10-LT under different baseline images.

FixMatch+DyTrim	CIFA	AR-10-LT
Input	$\gamma_l = \gamma_u = 100$	$\gamma_l = 100, \gamma_u = 150$
Noise images	77.7 / 76.8	76.7 / 75.8
Dataset means	78.0 / 76.1	76.7 / 74.2
Red	83.5 / 83.2	82.2 / 81.7
Green	83.7 / 83.3	81.5 / 81.0
Blue	84.5 / 84.2	83.1 / 82.6
Gray	84.1 / 83.7	82.3 / 81.9
Black	84.2 / 83.8	82.4 / 82.0
White	84.8 / 84.4	83.8 / 83.4

#### E.4 ABLATION STUDY

Effectiveness of each component. We conducted ablation studies on CIFAR-10-LT to assess the contribution of each component in DyTrim, varying the hyperparameter  $\gamma = \gamma_l = \gamma_u$  across 50, 100, and 150. As shown in Table. 11, the best performance was achieved when both labeled and unlabeled pruning were combined with rescaling. Removing rescaling led to a bACC drop of 0.8–2.1 points across  $\gamma$  values. Excluding either pruning component also reduced performance (e.g., -0.5 and -0.3 at  $\gamma = 50$  without unlabeled or labeled pruning, respectively). Removing both pruning strategies resulted in the most significant degradation. These results highlighted the complementary benefits of pruning and rescaling.

Sensitivity of different baseline images  $\mathcal{I}$ . We further examined the sensitivity of DyTrim to the choice of baseline image by conducting ablation studies on CIFAR-10-LT with different types of inputs, including noise, dataset means, and solid colors. Table 12 shows that solid-color images consistently outperform noise or mean-based baselines. Among them, white and black images deliver the strongest results.

#### E.5 QUALITATIVE ANALYSES

Since the baseline image could implicitly reflect the bias of the classifier, we argued that by customizing dynamic data pruning methods for labeled and unlabeled data, DyTrim significantly reduced classifier bias while improving performance. To verify this claim, in Figure. 3 (a) and (b), we analyzed the class probabilities predicted on the baseline image using FixMatch+DyTrim, trained on CIFAR-10-LT under various settings. We observed that classifiers trained with DyTrim consistently produced more balanced predictions than CDMAD across all settings, with improved accuracy on tail classes. To further validate the balanced classification effect of DyTrim, we visualized the dynamics of baseline image logits during training as shown in Figure. 4 (a), (b) and (c). The results clearly showed that BiGDP significantly reduced classifier bias induced by class imbalance. We defined r as the probability of pruning an unlabeled sample  $u_b^m$  when  $\mathcal{H}_t^u(u_b^m) < \bar{\mathcal{H}}_t^m$  and  $\max(P_\theta(y|\alpha(u_b^m))) \geq \tau$ . In Figure. 5, we evaluated different pruning ratios for unlabeled samples on CIFAR-10-LT. Results showed that setting  $r \geq 0.1$  yields higher performance across both archi-

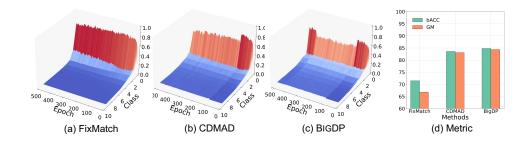


Figure 4: (a), (b) and (c) present the change of  $\pi_{\theta}(\mathcal{I})$  for the baseline image on CIFAR-10-LT with  $\gamma_l = \gamma_u = 100$  across different methods. (d) present the bACC and GM on those methods.

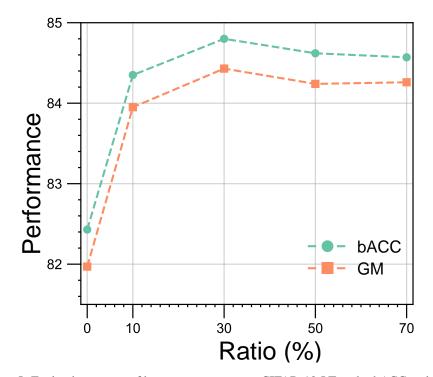


Figure 5: Evaluation curves of hyper-parameter r on CIFAR-10-LT under bACC and GM.

tectures, indicating that DyTrim was relatively robust with respect to the hyperparameter r, with the best performance achieved when r=0.3.

#### F Proof for Section 3

In this section, we present the technical details of Section 3. In particular, Section F.1 first discuss the relationship between classifier  $f_{\theta}$ , dataset  $\{\mathcal{X};\mathcal{U}\}$  and the baseline image  $\mathcal{I}$ . Then, Section ?? proves Theorem ?? for revealing the baseline image's intrinsic debiasing effect, and Section F.2 presents the details of the bias term and running statistics.

# F.1 Proof of Proposition 3

**Proposition 1.** Let  $\pi = \text{Softmax}(z)$  and  $z = g_{\theta}(x)$ . The one-step dynamics decompose as

$$\Delta \log \pi^{t}(y \mid \mathcal{I}) = -\eta \mathcal{T}^{t}(\mathcal{I}) \mathcal{K}^{t}(\mathcal{I}, x) \mathcal{G}^{t}(x, y) + \mathcal{O}(\eta^{2} \|\nabla_{\theta} z(x)\|_{\text{op}}^{2}), \tag{32}$$

where  $\mathcal{T}^t(\mathcal{I}) = \nabla_z \log_{\pi^t}(\mathcal{I}) = I - \mathbf{1} \pi_{\theta^t}^T(\mathcal{I})$ ,  $\mathcal{K}^t(\mathcal{I}, x) = (\nabla_{\theta} z(\mathcal{I})|_{\theta^t})(\nabla_{\theta} z(x)|_{\theta^t})^T$  is the empirical neural tangent kernel of the logit network z, and  $\mathcal{G}^t(x, y) = \nabla_z \mathcal{L}(x, y)|_{z^t}$ .

*Proof.* Inspired by the analysis of the learning dynamic of (Ren et al., 2022; Ren & Sutherland, 2024). In this work, we want to observe the classifier's prediction on the baseline image  $\mathcal{I}$ . Starting from Eq (12), we first approximate  $\log \pi^{t+1}(y \mid \mathcal{I})$  using first-order Talyor expansion, with slightly abused symbols, we use  $\pi^t$  to represent  $\pi_{\theta^{t+1}}$ , x to represent labeled sample  $x_b^n$  and u to represent unlabeled sample  $u_b^m$ :

$$\log \pi^{t+1}(y|\mathcal{I}) = \log \pi^t(y|\mathcal{I}) + \langle \nabla \log \pi^t(y|\mathcal{I}), \theta^{t+1} - \theta^t \rangle + \mathcal{O}(\|\theta^{t+1} - \theta^t\|^2)$$

Then, assuming the model updates its parameters using SGD calculated by an "updating labeled example" (x,y) or an "updating unlabeled example" u, we can rearrange the terms in the above equation to get the following expression:

$$\Delta \log \pi^t(y|\mathcal{I}) = \log \pi^{t+1}(y|\mathcal{I}) - \log \pi^{t+1}(y|\mathcal{I}) = \nabla_{\theta} \log \pi^t(y|\mathcal{I})|_{\theta^t}(\theta^{t+1} - \theta^t) + \mathcal{O}(\|\theta^{t+1} - \theta^t\|^2),$$

To evaluate the leading term, we first take a labeled sample as an example plug in the definition of SGD, and repeatedly use the chain rule:

$$\nabla_{\theta} \log \pi^{t}(y|\mathcal{I})|_{\theta^{t}}(\theta^{t+1} - \theta^{t}) = (\nabla_{z} \log \pi^{t}(y|\mathcal{I})|_{z^{t}})(-\eta \nabla_{\theta} \mathcal{L}(x)|_{\theta^{t}})^{T}$$

$$= (\nabla_{z} \log \pi^{t}(y|\mathcal{I})|_{z^{t}})(-\eta \nabla_{\theta} \mathcal{L}(x)|_{z^{t}} - \nabla_{\theta} z^{t}(x)|_{\theta^{t}})^{T}$$

$$= -\eta \nabla_{z} \log \pi^{t}(\mathcal{I})|_{z_{t}}[\nabla_{\theta} z(\mathcal{I})|_{\theta^{t}}(\nabla_{\theta} z(x)|_{\theta^{t}})^{T}](\nabla_{z} \mathcal{L}(x)|_{z^{t}})^{T}$$

$$= -\eta \mathcal{T}^{t}(\mathcal{I})\mathcal{K}^{t}(\mathcal{I}, x)\mathcal{G}^{t}(x, y)$$
(33)

For the higher-order term, using as above that

$$\theta^{t+1} - \theta^t = -\eta \nabla_{\theta} z^t(x) |_{\theta_t}^T \mathcal{G}^t(x, \hat{y})$$

and noting that, since the residual term  $\mathcal{G}^t$  is usually bouned, we have that

$$\mathcal{O}(\|\theta^{t+1} - \theta^t\|^2) = \mathcal{O}(\eta^2 \|(\nabla_{\theta} z^t(x)|\theta^t)^T\|_{op}^2 \|\mathcal{G}^t(x,\hat{y})\|_{op}^2) = \mathcal{O}(\eta^2 \|\nabla_{\theta} z(x)\|_{op}^2)$$

In the decomposition, we can write our  $\mathcal{T}^(t)$  as  $\mathcal{T}^t(\mathcal{I}) = \nabla_z \log_{\pi^t}(\mathcal{I}) = I - \mathbf{1}\pi_{\theta^t}^T(\mathcal{I})$ , this term is only related to the input  $\mathcal{I}$ , and reflects the model's correspondence to the baseline image, we will further analysis  $\log_{\pi^t}(\mathcal{I})$ . The second term in this decomposition,  $\mathcal{K}^t(\mathcal{I},x)$  is the product of gradients at  $\mathcal{I}$  and x or u. As shown in (Ren & Sutherland, 2024), if their gradients have similar directions, the Frobenius norm of this matrix is large, and vice versa. This matrix is known as the empirical neural tangent kernel, and it can change through the course of training as the network's notion of "similarity" evolves. The third term in this decomposition,  $\mathcal{G}^t$  is determined by the loss function  $\mathcal{L}$ , which provides the energy and direction for the model's adaptation. We have  $\mathcal{L} = \mathcal{L}_{sup}(x_b^n, y_b^n)$  for each labeled sample and  $\mathcal{L} = \mathcal{L}_{con}(\alpha(u_b^m), \mathcal{A}(u_b^m))$  for each unlabeled sample. According to the analysis of Xing et al. (2025),  $\mathcal{G}^t$  using the baseline image enhances the balance of the base SSL model implicitly utilizing the integrated gradient flow  $\nabla_{\theta}\mathcal{L}_{\text{Con}} = \sum_b \left(\sum_{i=1}^d \text{IntegratedGrads}_i(u_b^m)\right) \nabla g_b + \sum_b q_{\mathcal{A},b} \frac{\partial q_{\mathcal{A},b}}{\partial \theta}$ .

*Proof.* For the baseline image  $\mathcal{I}$  is a solid-balck image, i.e., the k=0, with  $h(\mathcal{I})=\beta$ 

$$g_{\theta}^{*}(x_{b}) = h(x_{b}) - h(\mathcal{I}) = \frac{\langle \mathbf{w}, x_{b} \rangle - \mathbb{E}[\langle \mathbf{w}, x_{b} \rangle]}{\sqrt{\text{Var}[\langle \mathbf{w}, x_{b} \rangle]}} * \gamma + \beta - \beta = \frac{\langle \mathbf{w}, x_{b} \rangle - \mathbb{E}[\langle \mathbf{w}, x_{b} \rangle]}{\sqrt{\text{Var}[\langle \mathbf{w}, x_{b} \rangle]}} * \gamma$$

# F.2 Detail of the bias term and running statistics

**Effects of bias term.** When the bias term  $\beta$  of the BN layer is frozen and equal to 0,  $h(\mathcal{I})$  becomes  $\gamma * (\langle \mathbf{w}, k \rangle - \mathbb{E}[\langle \mathbf{w}, k \rangle]) / \sqrt{\text{Var}[\langle \mathbf{w}, k \rangle]}$  which is the same as the Eq.(7) except for a bias term. Ignoring the running statistics strategy, the form of  $h(\mathcal{I})$  only depends on the  $\beta$ . As a result,  $h(\mathcal{I})$  becomes  $h(\mathcal{I}) \to 0$  during training and  $h(\mathcal{I}) \to -\gamma * \mathbb{E}_{mom}[\langle \mathbf{w}, x_b \rangle] / \sqrt{\text{Var}_{mom}[\langle \mathbf{w}, x_b \rangle]}$  during testing. This shows that the  $g_{\theta}^*$  operation has no effect in the training phase and only eliminates the impact of the unbalanced running means in the testing phase. This will affect the ability to benefit h from  $g_{\theta}^*$ , as shown in Table. 13.

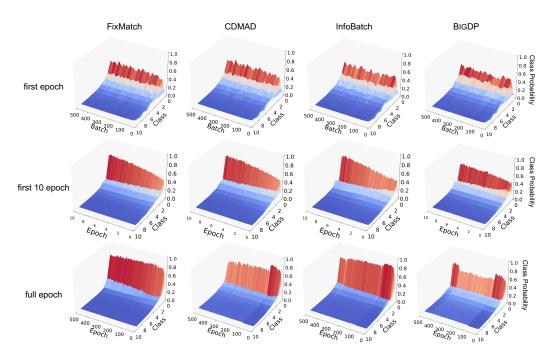


Figure 6: Comparison of the change of logits's probability distribution  $\pi_{\theta}(\mathcal{I})$  for the baseline image on CIFAR-10-LT with  $\gamma_l = \gamma_u = 100$  across different CISSL methods.

**Effects of running statistics.** When we do not keep running estimates, batch statistics are instead used during evaluation time as well. The form of  $h(\mathcal{I})$  becomes  $h(\mathcal{I}) \to \beta$  both training and testing. We can rewrite  $g_{\theta}^*(x_t) = \gamma * (\langle \mathbf{w}, x_t \rangle - \mathbb{E}[\langle \mathbf{w}, x_t \rangle]) / \sqrt{\text{Var}[\langle \mathbf{w}, x_t \rangle]}$ . On the other hand, as  $h(\mathcal{I}) \to 0$ , the benefit of  $g_{\theta}^*$  is also vanishes, also shown in Table. 13.

We then extend our results to a non-linear neural network, thus we have the following corollary:

Table 13: Comparison of bACC/GM on CIFAR-10-LT.

Metric	With original $g_{\theta}^*$	$g_{\theta}^*$ without $\beta$	$g_{ heta}^*$ without $\mathbf{x}_{mom}$	$g_{ heta}^*$ without $eta$ & $\mathbf{x}_{mom}$
bACC	$83.6 \pm 0.46$	$80.92 \pm 0.02 \downarrow 2.68$	$71.63 \pm 0.35 \downarrow 11.97$	$64.01 \pm 0.14 \downarrow 19.59$
GM	$83.1 \pm 0.57$	$80.37 \pm 0.23 \downarrow 2.73$	$67.85 \pm 0.51 \downarrow 15.25$	$54.48 \pm 0.36 \downarrow 28.62$

# G VISUALIZATION

#### G.1 Details of the change of logits's probability distribution

In this section, we conduct some visualization experiments to demonstrate the advantages of the DyTrim in debiasing and improving classifier performance. We first analyze the change of logits's probability distribution Softmax( $g_{\theta}(\mathcal{I})$ ) for the baseline image on CIFAR-10-LT with  $\gamma_l = \gamma_u = 100$  for fixmatch, CDMAD, and the DyTrim as shown in Figure. 6. It can be seen intuitively that in the first epoch, the classifier has bias due to the imbalance of categories in the data. This situation increases significantly with the number of network training times, as shown in the second column of the figure. However, we can see that DyTrim can effectively slow down the increase of this bias. Furthermore, after the model is fully trained for 500 epochs, it can be seen that after the 100th epoch, CDMAD starts to use the baseline image for post-hoc debiasing, which significantly reduces the representation of the model. However, by dynamically pruning the data set, DyTrim obtains a more distinct debias effect as shown in Figure. 7.

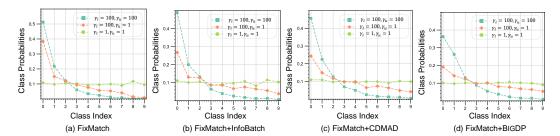


Figure 7: Class probabilities predicted on a baseline image using (a) FixMatch, (b) FixMatch+InfoBatch, (c) FixMatch+CDMAD, (d) FixMatch+DyTrim.

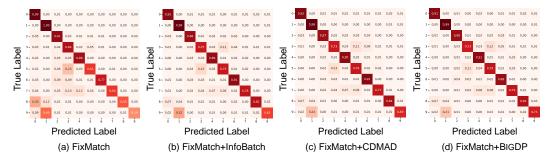


Figure 8: Confusion matrices of the class predictions on the test set of CIFAR-10 using (a) FixMatch, (b) FixMatch+InfoBatch, (c) FixMatch+CDMAD, and (d) FixMatch+DyTrim trained on CIFAR-10-LT under  $\gamma_l=100$  and  $\gamma_u=100$ .

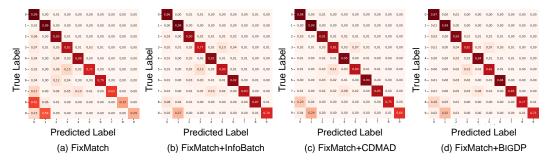


Figure 9: Confusion matrices of the class predictions on the test set of CIFAR-10 using (a) FixMatch, (b) FixMatch+InfoBatch, (c) FixMatch+CDMAD, and (d) FixMatch+DyTrim trained on CIFAR-10-LT under  $\gamma_l = 100$  and  $\gamma_u = 1$ .

#### G.2 Details of the change of logits's probability distribution

Figure. 8 and Figure. 9 compare the confusion matrices of the class predictions on the test set of CIFAR-10 using (a) FixMatch, (b) FixMatch+Infobatch, (c) FixMatch+CDMAD, and (d) FixMatch+DyTrim trained on CIFAR-10-LT under  $\gamma_l=100, \gamma_u=1,100$ . FixMatch+DyTrim made more balanced predictions across classes. Furthermore, we also conducted experiments under a balanced setting ( $\gamma=\gamma_1=\gamma_u=1$ ), as shown in Figure. 10. The results show that even under a balanced data distribution, BigDP can still achieve better results on the pruned dataset than methods such as CDMAD trained on the full dataset.

Similar to confusion matrices, we also compare t-distributed stochastic neighbor embedding (t-SNE) of representations obtained for the test set of CIFAR-10 using FixMatch, FixMatch+CDMAD, FixMatch+InfoBatch, and FixMatch+DyTrim trained on CIFAR-10 with  $\gamma_l=100$  and  $\gamma_u=1,100$  (unknown  $\gamma_u$ ), where different colors indicate different classes in CIFAR-10 Figure. 11, Figure. 12. We can observe that the representations obtained using FixMatch+DyTrim are separated into classes with clearer boundaries compared the those from FixMatch and CDMAD. This is prob-

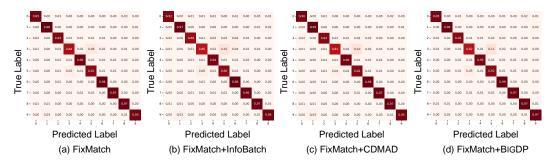


Figure 10: Confusion matrices of the class predictions on the test set of CIFAR-10 using (a) Fix-Match, (b) FixMatch+InfoBatch, (c) FixMatch+CDMAD, and (d) FixMatch+DyTrim trained on CIFAR-10-LT under  $\gamma_l = 1$  and  $\gamma_u = 1$ .

ably because CDMAD appropriately refined the biased pseudo-labels and used them for training, whereas FixMatch failed to learn the representations properly because they used the biased pseudo-labels for training. These results demonstrate that the quality of representations can be improved by using well-refined pseudo-labels for training.

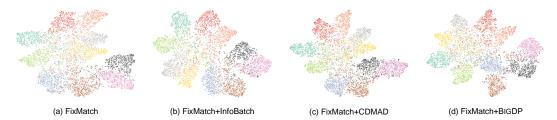


Figure 11: t-SNE of representations obtained for the test set of CIFAR-10 using (a) FixMatch, (b) FixMatch+InfoBatch, (c) FixMatch+CDMAD, and (d) FixMatch+DyTrim trained on CIFAR-10-LT under  $\gamma_l = 100$  and  $\gamma_u = 100$ .



Figure 12: t-SNE of representations obtained for the test set of CIFAR-10 using (a) FixMatch, (b) FixMatch+InfoBatch, (c) FixMatch+CDMAD, and (d) FixMatch+DyTrim trained on CIFAR-10-LT under  $\gamma_l = 100$  and  $\gamma_u = 1$ .

# H USE OF LLMS

This work did not involve the use of large language models (LLMs) at any stage. The design of experiments, data analysis, and manuscript preparation were conducted entirely by the authors through conventional computational methods and human expertise, without reliance on automated text generation or model-driven reasoning.