BALANCED LEARNING FOR DOMAIN ADAPTIVE SEMANTIC SEGMENTATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Unsupervised domain adaptation (UDA) for semantic segmentation aims to transfer knowledge from a labeled source domain to an unlabeled target domain, improving model performance on the target dataset without additional annotations. Despite the effectiveness of self-training techniques in UDA, they struggle to learn each class in a balanced manner due to inherent class imbalance and distribution shift in both data and label space between domains. To address this issue, we propose Balanced Learning for Domain Adaptation (BLDA), a novel approach to directly assess and alleviate class bias without requiring prior knowledge about the distribution shift between domains. First, we identify over-predicted and under-predicted classes by analyzing the distribution of predicted logits. Subsequently, we introduce a post-hoc approach to align the positive and negative logits distributions across different classes using anchor distributions and cumulative density functions. To further consider the network's need to generate unbiased pseudo-labels during self-training, we couple Gaussian mixture models to estimate logits distributions online and incorporate logits correction terms into the loss function. Moreover, we leverage the resulting cumulative density as domain-shared structural knowledge to connect the source and target domains. Extensive experiments on two standard UDA semantic segmentation benchmarks demonstrate that BLDA consistently improves performance, especially for under-predicted classes, when integrated into existing methods. Our work highlights the importance of balanced learning in UDA and effectively mitigates class bias in domain adaptive semantic segmentation.

033

003

010 011

012

013

014

015

016

017

018

019

021

025

026

027

028

1 INTRODUCTION

Semantic segmentation is a fundamental computer vision task that assigns a semantic label to each
 pixel in an image, enabling comprehensive image understanding. Despite the remarkable progress in
 recent research (Long et al., 2015a; Chen et al., 2017b; Cheng et al., 2021; 2022), the performance
 of these methods often significantly drops when applied to new target datasets. This performance
 degradation stems from the differences between the source and target domains, challenging the
 generalization ability of these models. To address this issue, unsupervised domain adaptation (UDA)
 techniques have been extensively studied. UDA aims to bridge the domain gap by transferring
 knowledge from a labeled source domain to an unlabeled target domain, thereby improving the
 model's performance on the target dataset without requiring additional annotations.

In previous work, self-training techniques (Tranheden et al., 2021; Hoyer et al., 2022a) have been 043 naturally introduced into UDA tasks to fully utilize the large amount of unlabeled target domain data, 044 becoming a mainstream paradigm. This paradigm constructs a teacher network using a temporal 045 aggregation mechanism, treats its predictions on the target domain as pseudo-labels, and gradually 046 guides the student network's learning. Despite achieving remarkable results, these methods struggle 047 to learn each class in a balanced manner. Generally, the inherent class imbalance in segmentation 048 datasets (Cordts et al., 2016) (Fig.1(a)) leads networks to produce biased predictions towards head classes, often studied as the long-tail problem (Van Horn & Perona, 2017; Buda et al., 2018; Liu et al., 2019). However, in UDA, data and label distribution shifts between the training and test 051 data complicate the class bias. The network's bias towards classes does not entirely depend on the differences in class sample distribution. As shown in Fig.1(b), when a network trained on the 052 source domain is tested on the target domain, the performance degradation varies greatly across classes, distinguishing easy-to-transfer and hard-to-transfer classes. These factors jointly determine



Figure 1: Demonstration of factors that cause class bias. (a) The inherent class imbalance problem in segmentation datasets. (b) The differences in transfer difficulty across classes in cross-domain settings. "Oracle' represents the performance under full supervision, while "Src-only' represents training with the source domain and testing it on the target domain. (c) The differences in logits distributions predicted for each class by the network, including "positive distribution" and "negative distribution". (d) Bias assessment for different classes via Eq.4. The corresponding class IDs of (a), (b), and (c) are mapped in descending order onto this figure.

081

the network's different biases towards each class in target domain, resulting in over-prediction and
 under-prediction. Furthermore, confirmation bias (Guo et al., 2017) causes self-training techniques to
 exacerbate this phenomenon. Fig.2(a) shows the severe deterioration of classes like *rider* and *bicycle* after self-training, widening the performance gap across classes. Therefore, achieving balanced
 learning for each class in UDA is a challenging and worthwhile exploration.

Existing strategies to reduce model bias towards different classes can be broadly categorized into 880 re-weighting (Cui et al., 2019; Lin et al., 2017; Cao et al., 2019; Truong et al., 2023; Buda et al., 2018) and re-sampling (Hoyer et al., 2022a; Araslanov & Roth, 2021; He et al., 2008; 2021; Guan 089 et al., 2022). To compare these methods, we take self-training as the baseline method in Fig.2 090 and implement re-weighting (Cui et al., 2019) and re-sampling (Hoyer et al., 2022a) techniques, 091 respectively. We observe the class bias through class-wise accuracy (Fig.2(a)) and the frequency of 092 pseudo-labels generated on the target domain during training (Fig.2(b)). Loss re-weighting aims to 093 assign different weights to classes, making the model pay more attention to tail classes. Although 094 intuitive, the update frequency of each class to the network still varies greatly, with some classes remaining challenging to learn effectively, resulting in unstable performance in self-training. In 096 contrast, sample re-sampling proves more effective by directly adjusting the class sample distribution during training, significantly enhancing the performance of tail classes. Despite their empirical 098 solid performance, these methods are heuristic and rely on the assumption that the test and training 099 data share the same distribution in both data and label space. However, in the UDA setting, these assumptions are invalid because (1) the class distributions of the source and target domains differ, and 100 the target domain's prior class distribution is unavailable; (2) the data distributions also differ, leading 101 to varying transfer difficulties across classes in cross-domain settings. This raises the question: How 102 to assess and alleviate class bias directly without requiring prior knowledge about the distribution 103 shift between the two domains? 104

In this work, we propose to assess the degree of class bias by analyzing the distribution of logits predicted by the network (Sec.3.3.2). Fig.1(c) shows that the network exhibits differences in the predicted logits distributions for different classes, directly leading to class bias. Fig.1(d) illustrates that the ranking of class bias highly coincides with the ranking of logit distribution differences, i.e.,

126

127

128 129



Figure 2: In UDA, class bias can be expressed as over-predicted classes and under-predicted classes. (a) Class-wise accuracy under different training settings. (b) Frequency of pseudo-labels generated by the network for different classes during training.

over-predicted classes have larger logit values, while under-predicted classes have smaller logit values. 130 This assessment approach prompts us to propose BLDA, a method to achieve balanced learning for 131 domain adaptive semantic segmentation by balancing the logits distribution. First, we consider a post-132 hoc approach to adjust the logits (Sec.3.3.). We set shared anchor distributions for the positive and 133 negative logits distributions and align the class-wise logits distributions with the anchor distributions 134 based on the cumulative density function mapping. Furthermore, to generate unbiased pseudo-135 labels for classes during self-training, we propose an online logit adjustment method (Sec.3.3.4). 136 This strategy couples Gaussian mixture models to estimate the logits distributions online during 137 training and incorporates logit correction terms into the loss function to replace the post-hoc method. 138 Moreover, we find that the resulting cumulative density can measure the discrimination difficulty of 139 different sample points in each class, which is a domain-shared structural knowledge that can be used 140 as an auxiliary loss to connect the two domains, further enhancing domain adaptation performance (Sec.3.3.5). As shown in Fig.2, our method can be integrated into existing self-training-based UDA 141 paradigms and effectively balance the prediction bias across classes. 142

143 Our contributions can be summarized as follows: (1) We propose statistically analyzing the predicted 144 logits to directly assess the network's bias towards classes. (2) We propose BLDA to estimate logit 145 correction terms online during UDA training to achieve balanced learning for each class. (3) We demonstrate that BLDA can be easily integrated into existing UDA methods and consistently improves 146 performance on two standard UDA semantic segmentation benchmarks, significantly enhancing the 147 performance of under-predicted classes and confirming the effectiveness of our method. 148

149 150

2 **RELATED WORK**

151 152 153

2.1DOMAIN ADAPTIVE SEMANTIC SEGMENTATION

154 Unsupervised domain adaptation (UDA) aims to transfer semantic knowledge learned from labeled 155 source domains to unlabeled target domains. Due to the ubiquity of domain gaps, UDA methods 156 have been widely studied in various computer vision tasks, such as image classification, object 157 detection, and semantic segmentation. UDA is crucial for semantic segmentation to avoid laborious 158 pixel-wise annotation in new target scenarios. Recent UDA approaches for semantic segmentation 159 can be categorized into two main paradigms: adversarial training-based methods (Toldo et al., 2020; Tsai et al., 2018; Chen et al., 2018; Ganin & Lempitsky, 2015; Hong et al., 2018; Long et al., 2015b) 160 and self-training-based methods (Tranheden et al., 2021; Hoyer et al., 2022a; Araslanov & Roth, 161 2021; Zhang et al., 2021). Adversarial training-based methods learn domain-invariant representations

162 through a min-max optimization game, where a feature extractor is trained to confuse a domain 163 discriminator, aligning feature distributions across domains. Self-training-based methods, which 164 have come to dominate the field due to the domain-robustness of Transformers (Bhojanapalli et al., 165 2021), generate pseudo labels for target images based on a teacher-student optimization framework. 166 The success of this paradigm relies on generating high-quality pseudo labels, with strategies such as entropy minimization (Chen et al., 2019) and consistency regularization (Hoyer et al., 2023) being 167 developed for this purpose. However, due to the inherent class imbalance and distribution shift in 168 both data and label space between domains, networks often produce complicated class bias, which is further exacerbated by confirmation bias in the self-training paradigm. Our method focuses on 170 balanced learning in UDA training. 171

172 173

174

2.2 CLASS-IMBALANCED LEARNING

Class imbalance is a common problem in semantic segmentation, where the number of samples per 175 class varies significantly. Existing methods address this issue through re-weighting or re-sampling 176 techniques. Re-weighting methods assign different weights to classes during training, giving higher 177 importance to under-represented classes (Cui et al., 2019; Lin et al., 2017; Cao et al., 2019). Re-178 sampling techniques modify the class distribution in the training data by over-sampling minority 179 classes or under-sampling majority classes (He et al., 2008; 2021). In UDA for semantic segmentation, some approaches have introduced these strategies to alleviate class bias (Hoyer et al., 2022a; Araslanov 181 & Roth, 2021; Li et al., 2022). However, these methods are still empirical and focus on the single-182 domain setting, which follows the assumptions that the test data and training data share the same distribution in both data space and label space, without considering the additional challenges posed 183 by domain shift in UDA. In this work, we aim to access class bias directly and achieve balanced learning for each class with no prior knowledge about the distribution shift between domains. 185

186 187

188 189

190 191

195

197

198 199

200

201

202 203

204 205

3 METHOD

3.1 PROBLEM DEFINITION

In unsupervised domain adaptation for semantic segmentation, the network is simultaneously trained 192 on labeled source domain data and unlabeled target domain data. To be specific, the source domain 193 can be denoted as $D_s = \{(x_i^S, y_i^S)\}_{i=1}^{N_S}$, where $x_i^S \in X_S$ represents an image with $y_i^S \in Y_S$ as the corresponding pixel-wise one-hot label covering C classes. The target domain can be denoted as 194 $D_t = \{(x_i^T)\}_{i=1}^{N_T}$, which shares the same label space but has no access to the target label Y_T . 196

3.2 **REVISITING SELF-TRAINING IN UDA**

Self-training-based pipelines for UDA segmentation consist of a supervised branch for the source domain and an unsupervised branch for the target domain. For the supervised branch, loss \mathcal{L}^s can only be calculated on the source domain to train a neural network f_{θ} :

$$\mathcal{L}^{s} = \frac{1}{N_{S}} \sum_{i=1}^{N_{S}} \frac{1}{HW} \sum_{j=1}^{H \times W} \ell_{ce}(f_{\theta}(x_{ij}^{S}), y_{ij}^{S}),$$
(1)

where ℓ_{ce} denotes the cross-entropy loss. Unsupervised branch introduces teacher-student framework 206 to generate pseudo-labels $\hat{y}_{ij}^T = \operatorname{argmax}(g_{\phi}(x_{ij}^T))$ with the teacher model g_{ϕ} for target domain: 207

$$\mathcal{L}^{u} = \frac{1}{N_{T}} \sum_{i=1}^{N_{T}} \frac{1}{HW} \sum_{j=1}^{H \times W} q(p_{ij}) \ell_{ce}(f_{\theta}(x_{ij}^{T}), \hat{y}_{ij}^{T}),$$
(2)

210 211

208 209

212 where we define $q(p_{ij})$ as a quality estimate conditioned on confidence $p_{ij} = max(g_{\phi}(x_{ij}^{T}))$ for 213 pseudo labels, which gradually strengthens with increasing accuracy of models and can be implemented with threshold filtering or a weighting function. After each training step, the teacher model 214 g_{ϕ} is updated with the exponentially moving average of the weights of f_{θ} . Then, the overall objective 215 function is a combination of supervised loss and unsupervised loss as $\mathcal{L} = \mathcal{L}^s + \mathcal{L}^u$.

3.3 BALANCED LEARNING FOR DOMAIN ADAPTIVE SEMANTIC SEGMENTATION

218 3.3.1 OVERVIEW

In this section, we first propose to assess the network's prediction bias towards each class by statistically analyzing the distribution of logits (Sec.3.3.2). Based on the above analysis, we define a post-hoc method to balance the network's predictions (Sec.3.3.3). Furthermore, we introduce online logits adjustment tailored for the UDA training process (Sec.3.3.4). Finally, we introduce cumulative density estimation as domain-shared knowledge to bridge the two domains (Sec.3.3.5).

224 225 226

236 237

241 242

247

248

249

250

251

252 253

3.3.2 Assessing Prediction Bias from Logits Distribution

227 Given a label space $\mathcal{Y} = [C] = \{1, 2, ..., C\}$, the segmentation network can be seen as a scorer 228 $f_{\theta}: x_{ii} \to \mathcal{R}^{C}$ that assigns class-wise scores, also known as logits, to a pixel x_{ii} from image x_{i} . To 229 investigate the distribution of logits obtained by the network for different classes, we can analyze it from the perspective of the confusion matrix. The confusion matrix is a $C \times C$ matrix M, where 230 each element M_{cl} represents the number of pixels with ground truth label c predicted as class l. We 231 replace each element M_{cl} in the confusion matrix M with the corresponding set of logits, i.e., the 232 logits predicted for class l for all pixels with ground truth label c, to obtain the logits set matrix \mathcal{M} . 233 We then use \mathcal{M} to assess prediction bias. 234

Definition 1. Element in logits set matrix, \mathcal{M}_{cl} .

$$\mathcal{M}_{cl} = \{ f_{\theta}(x_{ij})[l] \mid y_{ij} = c \},$$
(3)

where $f_{\theta}(x_{ij})[l]$ represents the logit value predicted by the network for class l of pixel x_{ij} , and y_{ij} is the ground truth label of pixel x_{ij} . In the resulting $C \times C$ matrix \mathcal{M} , diagonal elements \mathcal{M}_{ll} represent the "positive logits distribution" for class l, while off-diagonal elements \mathcal{M}_{cl} ($c \neq l$) represent the "negative logits distribution" for class l with respect to class c.

Definition 2. Bias of the network towards class l, Bias(l).

$$\operatorname{Bias}(l) = \frac{1}{C} \sum_{c \in [C]} \mathbb{P}(\underset{c' \in [C]}{\operatorname{arg\,max}} f_{\theta}(x)[c'] = l|y = c) - \frac{1}{C},\tag{4}$$

where $\mathbb{P}(\arg \max_{c' \in [C]} f_{\theta}(x)[c'] = l|y = c)$ represents the probability that the network predicts a sample from class c as class l. This definition measures the average difference between the probability of predicting class l and the uniform probability 1/C across all classes. A positive bias indicates over-prediction, while a negative bias indicates under-prediction for class l.

Let P_{cl} denote the distribution of \mathcal{M}_{cl} . Assuming each distribution P_{cl} is independent, we can estimate the $\mathbb{P}(l|c)$ by comparing the logit values:

$$\mathbb{P}(\underset{y'\in[C]}{\arg\max} f_{\theta}(x)[y'] = l|c) \approx \int_{-\infty}^{\infty} P_{cl}(z) \prod_{y'\neq l} \left(\int_{-\infty}^{z} P_{cy'}(t) dt \right) dz.$$
(5)

258

259

260

261

Combining Eq.4 and Eq.5, for an unbiased network, i.e., Bias(l) = 0 for all $l \in [C]$, a sufficient condition is that they have the same positive and negative distributions. This means the network's prediction performance is consistent across all classes. Fig.1(d) shows a direct correlation between logit distribution differences and class bias, indicating that variations in logit distributions lead to class bias in the network's predictions.

262 263 3.3.3 Post-Hoc Class Balancing

Generally, the network tends to produce larger logits for over-predicted classes and smaller logits
 for under-predicted classes, as shown in Fig.3 (a). Reweighting/resampling strategies can alleviate
 this gap by making the network pay more attention to tail classes and reducing the emphasis on head
 classes during training, as illustrated in Fig.3(b). However, as shown in Fig.1, class bias does not fully
 correlate with the inherent class imbalance problem, especially in the UDA setting, where different
 distribution shifts exist in both data and label space between domains. Furthermore, these methods are empirical and lack generalization capability across various scenarios.



Figure 3: Illustration of proposed post-hoc class balancing. (a) The logits distributions of overpredicted and under-predicted classes. (b) Reweighting/resampling strategies alleviate class imbalance by adjusting the training emphasis on different classes. (c) Our post-hoc logits adjustment method aligns the logits distributions of all classes with anchor distributions to achieve balanced prediction.

Based on the above analysis, to balance the network's prediction capabilities across classes, we adjust the network's predictions in a post-hoc manner. Specifically, we define an anchor distribution P_p for the positive logits distribution and an anchor distribution P_n for the negative logits distribution. We then align all the logits distributions with corresponding anchor distributions, as shown in Fig.3 (c). To preserve the relative ordering of logits, i.e., structural information within each distribution, we align them in a point-wise way via the cumulative distribution function (CDF). Let $F_{cl}(z)$, $F_p(z)$ and $F_n(z)$ be the CDFs of P_{cl} , P_p and P_n , respectively. We align the logit value z from P_{cl} to P_p or P_n as follows:

$$z' = \begin{cases} F_p^{-1}(F_{cl}(z)), & \text{if } c = l \\ F_n^{-1}(F_{cl}(z)), & \text{if } c \neq l \end{cases}$$
(6)

where z' is the aligned value of z with respect to the anchor distribution. For brevity, we define an offset for logit as $\Delta_{cl}(z) = z' - z$. Considering the probability estimate that $p(y_{ij} = c|x_{ij}) \propto \exp(f_{\theta}(x_{ij})[c])$, we can obtain revised prediction results for each pixel x_{ij} by:

$$\widetilde{y}_{ij} = \operatorname{argmax}_{c \in [C]} \frac{\exp(f_{\theta}(x_{ij})[c] + \tau \Delta_{cc}(f_{\theta}(x_{ij})[c]))}{\exp(f_{\theta}(x_{ij})[c] + \tau \Delta_{cc}(f_{\theta}(x_{ij})[c])) + \sum_{c' \neq c} \exp(f_{\theta}(x_{ij})[c'] + \tau \Delta_{cc'}(f_{\theta}(x_{ij})[c']))},$$
(7)

where τ is a scaling factor (the derivation of Eq.7 is detailed in Appendix A). When $\tau = 1$, the model produces balanced predictions. However, to achieve optimal performance on specific evaluation metrics (e.g., mIoU), we need to adjust the value of τ . We discuss the choice of τ in detail in the experimental section. In this way, the network generates balanced predictions for different classes.

307 Discussion about anchor distribution. In our experiments, we use the global positive and negative 308 logits distributions on the source domain to estimate the anchor distribution for both source and target 309 domains. This choice is based on two key considerations: (1) The network tends to produce larger logits for over-predicted classes and smaller logits for under-predicted classes. By using the global 310 logits distribution, we can effectively measure the average learning degree of the network across all 311 classes. Aligning each class-specific distribution with this global distribution can help neutralize the 312 class bias of the network, ensuring a more balanced learning process. (2) When estimating the logits 313 distributions for each class, there exist varying degrees of statistical errors. According to Bernstein 314 inequalities, estimating the global logits distribution can reduce estimation errors to a certain extent 315 and accelerate convergence rates. In our case, the global distribution, being more robust and stable, 316 serves as a reliable anchor distribution for subsequent alignment.

317 318

319

270

271

272

273

274

275

276 277

282

283

284

285

287

288

289

290

291

292

293

295 296

297

298

3.3.4 ONLINE LOGITS ADJUSTMENT FOR UDA

While the post-hoc logits adjustment method introduced in Sec.3.3.3 can effectively balance the network predictions across different classes, it is performed after training the model. In the UDA setting, it is significant to incorporate the logits balancing mechanism directly into the training process.
By doing so, the model can learn to make more balanced predictions while adapting to the target domain through pseudo labels.

To achieve this, we propose an online logit adjustment method tailored for UDA training. The key to this method lies in the online estimation of the logits distributions. We employ Gaussian Mixture Models (GMMs) to model these distributions. Considering that the source and target domains have inherently different logits distributions, we maintain two sets of GMMs separately, with each set containing $C \times C \times K$ Gaussian components, where C denotes the number of classes and K denotes the number of Gaussian components per element in \mathcal{M} . Formally, we define:

$$P_{cl}^{s} = \sum_{k=1}^{K} \pi_{clk}^{s} \mathcal{N}(\mu_{clk}^{s}, \sigma_{clk}^{s}), \quad P_{cl}^{t} = \sum_{k=1}^{K} \pi_{clk}^{t} \mathcal{N}(\mu_{clk}^{t}, \sigma_{clk}^{t}), \tag{8}$$

333 where P_{cl}^s and P_{cl}^t represent the estimated logits distributions for class l when the ground truth label 334 or pseudo label is c in the source and target domains, respectively. The parameters π_{clk}^s , μ_{clk}^s , σ_{clk}^s and π_{clk}^t , μ_{clk}^t , σ_{clk}^t denote the mixing coefficient, mean, and standard deviation of the k-th Gaussian 335 336 component in the corresponding GMM. We update the GMM parameters during each training 337 iteration using the logits obtained from the current mini-batch via the Expectation-Maximization 338 (EM) algorithm (McLachlan & Krishnan, 2007). Since the network f_{θ} gradually evolves during 339 training, we adopt a momentum-based EM update strategy. Specifically, we directly use the GMM parameters $\hat{\phi}_{cl}$ estimated in the latest iteration as the initialization $\phi_{cl}^{(0)}$ for the current iteration. After T EM loops, the current iteration is completed, and a momentum update is adapted with 340 341 342 $\phi_{cl}^{(T)} \leftarrow (1 - \tilde{\tau}^n)\phi_{cl}^{(T)} + \tilde{\tau}^n \hat{\phi}_{cl}$, where n represents the number of iterations that have not been 343 updated since the last parameter update for ϕ_{cl} . We also implement GMMs to estimate the anchor 344 distributions P_p and P_n using the source domain data's global positive and negative logits. The algorithm flow is detailed in Appendix D. 345

After estimating the logits distributions using GMMs, we compute the adjusted logits offset Δ^S and Δ^T for source and target domains, respectively. These offsets are then used to adjust the cross-entropy loss for both domains:

$$\widetilde{\mathcal{L}}^{s} = -\frac{1}{N_{S}} \sum_{i=1}^{N_{S}} \sum_{j=1}^{H \times W} \log \frac{\exp(f_{\theta}(x_{ij}^{S})[y_{ij}^{S}] - \tau \Delta_{y_{ij}^{S},y_{ij}^{S}}^{S}(f_{\theta}(x_{ij}^{S})[y_{ij}^{S}]))}{\sum_{c=1}^{C} \exp(f_{\theta}(x_{ij}^{S})[c] - \tau \Delta_{y_{ij}^{S},c}^{S}(f_{\theta}(x_{ij}^{S})[c]))},$$

$$\widetilde{\mathcal{L}}^{u} = -\frac{1}{N_{T}} \sum_{i=1}^{N_{T}} \sum_{j=1}^{H \times W} q(p_{ij}) \log \frac{\exp(f_{\theta}(x_{ij}^{T})[\hat{y}_{ij}^{T}] - \tau \Delta_{\hat{y}_{ij}^{T},\hat{y}_{ij}^{T}}^{T}(f_{\theta}(x_{ij}^{T})[\hat{y}_{ij}^{T}]))}{\sum_{c=1}^{C} \exp(f_{\theta}(x_{ij}^{T})[c] - \tau \Delta_{\hat{y}_{ij}^{T},c}^{S}(f_{\theta}(x_{ij}^{T})[\hat{y}_{ij}^{T}]))}.$$
(9)

350 351

330 331

332

354 355

356 357

358

359

360

361

362

363

364

In contrast to Eq.7, where the logits are adjusted post-hoc, Eq.9 directly incorporates the offset into the learning process of the logits. This approach is equivalent to learning a scorer of the form $g(x)[y] = f(x)[y] - \tau \Delta_y(x)$. Consequently, we have $\operatorname{argmax} f(x)[y] = g(x)[y] + \tau \Delta_y(x)$, which can be seen as analogous to the post-hoc adjustment. By employing these adjusted losses, we achieve a two-fold benefit. Firstly, the anchor distribution can serve as a reference distribution to balance the learning progress between classes within both domains. Secondly, since the pseudo-label-based loss in the target domain has a gradually increasing weight, using a shared anchor distribution allows the logits distribution of the target domain to gradually align with that of the source domain, thus establishing a connection between the two domains.

366 3.3.5 BRIDGING DOMAINS THROUGH CUMULATIVE DENSITY ESTIMATION

Furthermore, for each sample pixel, we can query the corresponding positive cumulative distribution value F_{cc} based on its label c, which ranges from 0 to 1. The positive distribution measures the discriminative ability of a class, and we find that this cumulative distribution value indicates the difficulty of the sample pixel belonging to that class. This structural knowledge depends only on the context of the pixel and is not affected by the image style, making it domain-invariant. To further bridge the two domains, we add an extra regression head to the network to predict this value as an additional auxiliary task.

Specifically, for each sample point x_{ij}^S in the source domain, we can query the corresponding positive cumulative distribution value based on its true label y_{ij}^S : $d_{ij}^S = F_{y_{ij}^S, y_{ij}^S}(f_\theta(x_{ij}^S)[y_{ij}^S])$, where $F_{y_{ij}^S, y_{ij}^S}$ is the positive cumulative distribution function for class y_{ij}^S . Similarly, for each sample point x_{ij}^T in the target domain, we can query the corresponding positive cumulative distribution value based on its

378Table 1: UDA segmentation performance on GTA. \rightarrow CS. using the mIoU (%) evaluation metric,379where the improvement is marked as **bold**. The results are acquired based on CNN-based model380(He et al., 2016; Chen et al., 2017a), denoted as C, and Transformer-based model (Xie et al., 2021),381denoted as T.

| Method | Arch. | Road | Sidewalk | Building | Wall | Fence | Pole | Light | Sign | Veg | Terrain | Sky | Person | Rider | Car | Truck | Bus | Train | Motor | Bike | mIoU | std |
|--------------------------------|-------|------|----------|----------|----------|-------------|------------|--------|------|-----------|-------------|----------|--------|-------|------|-------|------|-------|-------|--------------|-----------|------|
| DACS (Tranheden et al., 2021) | C | 89.9 | 39.7 | 87.9 | 39.7 | 39.5 | 38.5 | 46.4 | 52.8 | 88.0 | 44.0 | 88.8 | 67.2 | 35.8 | 84.5 | 45.7 | 50.2 | 0.2 | 27.3 | 34.0 | 52.1 | 24.4 |
| ProDA (Zhang et al., 2021) | C | 87.8 | 56.0 | 79.7 | 46.3 | 44.8 | 45.6 | 53.5 | 53.5 | 88.6 | 45.2 | 82.1 | 70.7 | 39.2 | 88.8 | 45.5 | 50.4 | 1.0 | 48.9 | 56.4 | 57.5 | 21.2 |
| CPSL (Li et al., 2022) | C | 92.3 | 59.5 | 84.9 | 45.7 | 29.7 | 52.8 | 61.5 | 59.5 | 87.9 | 41.6 | 85.0 | 73.0 | 35.5 | 90.4 | 48.7 | 73.9 | 26.3 | 53.8 | 53.9 | 60.8 | 20.3 |
| TransDA (Chen et al., 2022) | Т | 94.7 | 64.2 | 89.2 | 48.1 | 45.8 | 50.1 | 60.2 | 40.8 | 90.4 | 50.2 | 93.7 | 76.7 | 47.6 | 92.5 | 56.8 | 60.1 | 47.6 | 49.6 | 55.4 | 63.9 | 18.5 |
| DAFormer (Hoyer et al., 2022a) | Т | 95.7 | 70.2 | 89.4 | 53.5 | 48.1 | 49.6 | 55.8 | 59.4 | 89.9 | 47.9 | 92.5 | 72.2 | 44.7 | 92.3 | 74.5 | 78.2 | 65.1 | 55.9 | 61.8 | 68.3 | 16.8 |
| +BLDA | Т | 95.4 | 78.3 | 88.3 | 54.0 | 55.2 | 55.7 | 60.3 | 65.2 | 89.2 | 47.3 | 91.1 | 71.4 | 44.8 | 91.6 | 74.3 | 83.4 | 73.2 | 59.3 | 67.1 | 70.7 | 15.5 |
| HRDA (Hoyer et al., 2022b) | Т | 96.5 | 74.4 | 91.0 | 61.6 | 51.5 | 57.1 | 63.9 | 69.3 | 91.3 | 48.4 | 94.2 | 79.0 | 52.9 | 93.9 | 84.1 | 85.7 | 75.9 | 63.9 | 67.5 | 73.8 | 15.4 |
| +BLDA | Т | 96.4 | 77.6 | 90.7 | 63.3 | 57.9 | 62.1 | 66.5 | 72.5 | 91.3 | 52.2 | 94.4 | 76.9 | 57.3 | 93.5 | 86.2 | 87.7 | 79.9 | 66.8 | 68.9 | 75.6 | 13.8 |
| MIC (Hoyer et al., 2023) | T | 97.4 | 80.1 | 91.7 | 61.2 | 56.9 | 59.7 | 66.0 | 71.3 | 91.7 | 51.4 | 94.3 | 79.8 | 56.1 | 94.6 | 85.4 | 90.3 | 80.4 | 64.5 | 68.5 | 75.9 | 14.8 |
| +BLDA | T | 97.1 | 82.6 | 91.6 | 64.7 | 61.0 | 64.9 | 68.0 | 74.8 | 91.2 | 56.6 | 92.4 | 80.0 | 54.7 | 95.7 | 87.3 | 88.8 | 82.6 | 64.2 | 72.0 | 77.1 | 13.5 |
| Image | DAF | orm | er | | D/ +p | AFoi ost | me -hoo | r c | | DAI +I | Forr BLD | ner A | | | C | DE | | | C | Grou Trut | ınd th | |



Figure 4: Qualitative results. Note that the yellow boxes mark regions improved by BLDA.

pseudo label \hat{y}_{ij}^T : $d_{ij}^T = F_{\hat{y}_{ij}^T, \hat{y}_{ij}^T}(f_{\theta}(x_{ij}^T)[\hat{y}_{ij}^T])$, where $F_{\hat{y}_{ij}^T, \hat{y}_{ij}^T}$ is the positive cumulative distribution function for the class corresponding to the pseudo label \hat{y}_{ij}^T . We then add an extra regression head h_{ϕ} to the network to predict the cumulative distribution value for each sample point. The corresponding regression losses for the source and target domains can be defined as:

$$\mathcal{L}_{reg}^{S} = \frac{1}{N_{S}} \sum_{i=1}^{N_{S}} \sum_{j=1}^{H \times W} |h_{\phi}(\tilde{f}_{\theta}(x_{ij}^{S})) - d_{ij}^{S}|^{2}, \\ \mathcal{L}_{reg}^{T} = \frac{1}{N_{T}} \sum_{i=1}^{N_{T}} \sum_{j=1}^{H \times W} q(p_{ij}) |h_{\phi}(\tilde{f}_{\theta}(x_{ij}^{T})) - d_{ij}^{T}|^{2},$$
(10)

where $|\cdot|^2$ denotes the L2 loss, and $\tilde{f}_{\theta}(x_{ij})$ denote the features extracted by the network f_{θ} . Finally, the overall training objective can be expressed as $\mathcal{L} = \tilde{\mathcal{L}}^s + \tilde{\mathcal{L}}^u + \lambda(\mathcal{L}_{reg}^S + \mathcal{L}_{reg}^T)$, where λ is a hyperparameter balancing the cumulative density estimation loss.

4 EXPERIMENTS

4.1 EXPERIMENTAL SETUP

Datasets. Following standard UDA protocols, we evaluate our method on two widely used bench-marks that involve transferring knowledge from a synthetic domain to a real domain in a street scene setting. Specifically, we use GTAv/SYNTHIA (Ros et al., 2016; Richter et al., 2016) as the labeled source domain and Cityscapes (Cordts et al., 2016) as the unlabeled target domain. GTAv contains 24,966 synthetic images with a resolution of 1914 × 1052, while SYNTHIA consists of 9,400 synthetic images with a resolution of 1280 × 960.

Implementation Details. Our method can be built with different self-training-based frameworks. For thorough evaluation, we apply BLDA to three strong baseline methods, i.e., DAFormer (Hoyer et al., 2022a), HRDA (Hoyer et al., 2022b), and MIC (Hoyer et al., 2023), with MiT-B5 (Xie et al., 2021) pretrained on ImageNet-1k (Deng et al., 2009) as the backbone. BLDA is implemented based on MMS egmentation (Contributors, 2020). All experiments are trained for 40K iterations and a batch size of 2, with one or two RTX-3090 (24 GB memory) GPUs, depending on the complexity of used UDA frameworks. We train the network with an AdamW optimizer with learning rates of 6×10^{-5} for the encoder and 6×10^{-4} for the decoder, a weight decay of 0.01, and linear learning rate warm-up for the first 1.5K iterations. The input images are rescaled and randomly cropped to

| Method | Arch. | Road | Sidewalk | Building | Wall | Fence | Pole | Light | Sign | Veg | Terrain | Sky | Person | Rider | Car | Truck | Bus | Train | Motor | Bike | mIoU | std |
|--------------------------------|-------|------|----------|----------|------|-------|------|-------|------|------|---------|------|--------|-------|------|-------|------|-------|-------|------|------|------|
| DACS (Tranheden et al., 2021) | C | 80.6 | 25.1 | 81.9 | 21.5 | 2.9 | 37.2 | 22.7 | 24.0 | 83.7 | - | 90.8 | 67.6 | 38.3 | 82.9 | - | 38.9 | - | 28.5 | 47.6 | 48.3 | 27.5 |
| ProDA (Zhang et al., 2021) | C | 87.8 | 45.7 | 84.6 | 37.1 | 0.6 | 44.0 | 54.6 | 37.0 | 88.1 | - | 84.4 | 74.2 | 24.3 | 88.2 | - | 51.1 | - | 40.5 | 45.6 | 55.5 | 25.6 |
| CPSL (Li et al., 2022) | C | 87.2 | 43.9 | 85.5 | 33.6 | 0.3 | 47.7 | 57.4 | 37.2 | 87.8 | - | 88.5 | 79.0 | 32.0 | 90.6 | - | 49.4 | - | 50.8 | 59.8 | 57.9 | 25.5 |
| TransDA (Chen et al., 2022) | Т | 90.4 | 54.8 | 86.4 | 31.1 | 1.7 | 53.8 | 61.1 | 37.1 | 90.3 | - | 93.0 | 71.2 | 25.3 | 92.3 | - | 66.0 | - | 44.4 | 49.8 | 59.3 | 26.5 |
| DAFormer (Hoyer et al., 2022a) | T | 84.5 | 40.7 | 88.4 | 41.5 | 6.5 | 50.0 | 55.0 | 54.6 | 86.0 | - | 89.8 | 73.2 | 48.2 | 87.2 | - | 53.2 | - | 53.9 | 61.7 | 60.9 | 22.1 |
| +BLDA | Т | 80.7 | 44.9 | 85.6 | 45.1 | 9.6 | 54.3 | 60.2 | 58.7 | 87.7 | - | 92.3 | 75.7 | 51.1 | 87.3 | - | 62.7 | - | 59.9 | 65.8 | 64.0 | 20.6 |
| HRDA (Hoyer et al., 2022b) | Т | 85.2 | 47.7 | 88.8 | 49.5 | 4.8 | 57.2 | 65.7 | 60.9 | 85.3 | - | 92.9 | 79.4 | 52.8 | 89.0 | - | 64.7 | - | 63.9 | 64.9 | 65.8 | 21.4 |
| +BLDA | Т | 83.9 | 54.9 | 87.5 | 53.1 | 11.5 | 63.2 | 69.4 | 64.4 | 87.2 | - | 93.1 | 79.1 | 54.7 | 88.3 | - | 69.1 | - | 64.2 | 65.7 | 67.9 | 19.3 |
| MIC (Hoyer et al., 2023) | T | 86.6 | 50.5 | 89.3 | 47.9 | 7.8 | 59.4 | 66.7 | 63.4 | 87.1 | - | 94.6 | 81.0 | 58.9 | 90.1 | - | 61.9 | - | 67.1 | 64.3 | 67.3 | 21.0 |
| +BLDA | Т | 86.1 | 61.2 | 89.8 | 47.2 | 10.2 | 62.6 | 70.3 | 67.1 | 90.0 | - | 94.4 | 81.4 | 56.4 | 90.5 | - | 67.2 | - | 64.8 | 66.3 | 69.1 | 20.4 |

Table 2: UDA segmentation performance on SYN. \rightarrow CS. using the mIoU (%) evaluation metric, where the improvement is marked as **bold**. Note that the mIoUs on are calculated over 16 classes.



Table 3: BLDA ablation study of different components built with DAFormer.

| - | | 80.0 | None Post-hoc | OLA_S | OLA | CDE | mIoU | mAcc |
|---------|-------------------------|----------------------------|---------------|--------------|--------------|--------------|------|------|
| 5 - | | 78.0 | \checkmark | | | | 68.3 | 77.8 |
| 1 | 0.30 0.60 0.90 | | \checkmark | | | | 69.2 | 80.3 |
| 0.00 | scaling factor (τ) | scaling factor (τ) | | \checkmark | | | 68.9 | 79.5 |
| | 5 0 1 6 1 | 1.00 / 1 / | | | \checkmark | | 70.2 | 81.8 |
| Figi | ure 5: Study of the | e different evalution | | | \checkmark | \checkmark | 70.7 | 82.0 |
| me | arics with respect | to scaling factor τ . | | | | | | |

 512×512 following the same data augmentation in DAFormer (Hoyer et al., 2022a), and the EMA coefficient for updating the teacher net is set to be 0.999. We set temperature coefficient $\tau = 0.1$ and loss weight $\lambda = 0.2$ respectively. For more details, please refer to the Appendix.

4.2 COMPARISON WITH EXISTING METHODS

Comparative Evaluation. We compare BLDA to existing state-of-the-art UDA approaches on the GTA. \rightarrow CS. and SYN. \rightarrow CS. benchmarks. In addition to the widely adopted mean intersection-over-union (mIoU) metric, we also report the mean accuracy (mAcc) metric, which is equivalent to measuring the balanced error (Menon et al., 2020), i.e., the average of each class's error rate, and is more suitable to be used to assess the balance among classes. We discuss these two metrics in detail in Appendix C. Additionally, we calculate the standard deviation of IoU and Acc for each class to reflect the balanced degree of performance across classes.

Evaluation Results. Tab.1-2 shows BLDA consistently improves the performance of all baseline methods on two benchmarks by a large margin, ranging from 1.2% to 3.1%. Furthermore, significant improvements are obtained for under-predicted classes, such as *sidewalk*, *fence*, *pole*, *light*, and *sign*, which demonstrates that BLDA can mitigate the class bias with decreased standard deviation and thus bring the performance gains. In Table 4, the same phenomenon is observed, and the improvement in mAcc is more significant, ranging from 2.9% to 4.2%.

ABLATION STUDY 4.3

We conduct a series of ablation studies on the $GTA \rightarrow CS$. benchmark built with DAFormer (Hoyer et al., 2022a). Please refer to the Appendix for further analysis, where we provide a deeper study of parameter settings and more visualization results.

Influence of scaling factor. As illustrated in Fig. 5, we explore the influence of different τ values on evaluation metrics. The mAcc metric gradually increases as τ grows, reaching its peak performance when $\tau = 1$. Interestingly, the mIoU metric does not demonstrate a perfectly positive correlation with the rise in mAcc. This discrepancy arises from the fact that the mIoU calculation is heavily impacted by the imbalanced distribution of the test set, whereas mAcc serves as a class-balanced metric. We comprehensively explain this phenomenon in Appendix C. Our method, which models class-balanced learning, effectively boosts mAcc. However, to achieve improvements in mIoU, selecting a smaller scaling factor τ is necessary.



Table 4: UDA segmentation performance on $GTA \rightarrow CS$. using the mAcc (%) evaluation metric, where the improvement is marked as **bold**.

Figure 6: Comparison of Logits Distribution. We choose {building (2), vegetation (8)} as overpredicted classes, and {*traffic sign* (7), *bicycle* (18)} as under-predicted classes for visualization. Note that the anchor distribution is counted separately at baseline and in our method.

Effectiveness of Components. In Tab.10, we delve into the various components of BLDA. By solely applying the post-hoc method to adjust the predictions, we observe a minor performance improvement of 0.9%. When introducing online logit adjustment exclusively during source domain image training (OLA_S) , the improvement is comparatively modest at 0.6%. However, by simultaneously performing adjustments in both domains (OLA), we witness a significant performance boost of 1.9%, suggesting that this strategy effectively captures the disparity in learning degrees between the domains. Lastly, the extra supervison from cumulative distribution estimation (CDE) further models the shared structural information across the domains (shown in Fig.4), producing additional performance gains.

Qualitative Results. Fig.4 presents the qualitative results. We observe that for under-predicted classes, such as *sidewalk* and *pole*, the baseline method struggles to recognize them accurately. While the post-hoc method can slightly improve the performance, our proposed BLDA approach significantly enhances the ability to predict these classes.

Comparison of Logits Distribution. In Fig.6, we visualize the positive distribution and negative distribution corresponding to the over-predicted classes { *building* (2), *vegetation* (8)} and under-predicted classes { traffic sign (7), bicycle (18) } on the Cityscapes Val. set. In the baseline method, the positive and negative logits of classes *building* and *vegetation* are larger than anchor distribution, while this phenomenon is reversed in classes *traffic light* and *bicycle*, which leads to class bias. Our method reduces this distribution difference by aligning with the anchor distribution and achieves class-balanced learning.

CONCLUSION

In this work, we present Balanced Learning for Domain Adaptation (BLDA), a novel approach to ad-dress class bias in unsupervised domain adaptation (UDA) for semantic segmentation. BLDA analyzes logits distributions to assess prediction bias and introduces an online logits adjustment mechanism to balance class learning in both source and target domains. Our method effectively mitigates class bias, promotes balanced learning, and enhances generalization to the target domain. Experimental results demonstrate consistent performance improvements on standard UDA benchmarks.

540 REFERENCES

- Nikita Araslanov and Stefan Roth. Self-supervised augmentation consistency for adapting semantic
 segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15384–15394, 2021.
- Srinadh Bhojanapalli, Ayan Chakrabarti, Daniel Glasner, Daliang Li, Thomas Unterthiner, and
 Andreas Veit. Understanding robustness of transformers for image classification. In *Proceedings* of the IEEE/CVF international conference on computer vision, pp. 10231–10241, 2021.
- Mateusz Buda, Atsuto Maki, and Maciej A Mazurowski. A systematic study of the class imbalance problem in convolutional neural networks. *Neural networks*, 106:249–259, 2018.
- 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.
- Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille.
 Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):834–848, 2017a.
- Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*, 2017b.
- Minghao Chen, Hongyang Xue, and Deng Cai. Domain adaptation for semantic segmentation with
 maximum squares loss. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 2090–2099, 2019.
- Runfa Chen, Yu Rong, Shangmin Guo, Jiaqi Han, Fuchun Sun, Tingyang Xu, and Wenbing Huang.
 Smoothing matters: Momentum transformer for domain adaptive semantic segmentation. *arXiv* preprint arXiv:2203.07988, 2022.
- Yuhua Chen, Wen Li, and Luc Van Gool. Road: Reality oriented adaptation for semantic segmentation
 of urban scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
 pp. 7892–7901, 2018.
- Bowen Cheng, Alex Schwing, and Alexander Kirillov. Per-pixel classification is not all you need for semantic segmentation. *Advances in Neural Information Processing Systems*, 34:17864–17875, 2021.
- Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked attention mask transformer for universal image segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 1290–1299, 2022.
- Peng Chu, Xiao Bian, Shaopeng Liu, and Haibin Ling. Feature space augmentation for long-tailed data. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIX 16*, pp. 694–710. Springer, 2020.
- 581
 582 MMSegmentation Contributors. Mmsegmentation: Openmmlab semantic segmentation toolbox and benchmark, 2020.
- Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo
 Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban
 scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3213–3223, 2016.
- 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.
- 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.

| 594 595 | Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In <i>International conference on machine learning</i> , pp. 1180–1189. PMLR, 2015. |
|--------------------------|--|
| 597 598 599 | Dayan Guan, Jiaxing Huang, Aoran Xiao, and Shijian Lu. Unbiased subclass regularization for semi-supervised semantic segmentation. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 9968–9978, 2022. |
| 600 601 | Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In <i>International conference on machine learning</i> , pp. 1321–1330. PMLR, 2017. |
| 603 604 | Haibo He and Edwardo A Garcia. Learning from imbalanced data. <i>IEEE Transactions on knowledge and data engineering</i> , 21(9):1263–1284, 2009. |
| 605 606 607 | Haibo He, Yang Bai, Edwardo A Garcia, and Shutao Li. Adasyn: Adaptive synthetic sampling approach for imbalanced learning. In 2008 IEEE international joint conference on neural networks (IEEE world congress on computational intelligence), pp. 1322–1328. Ieee, 2008. |
| 608 609 610 611 | Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 770–778, 2016. |
| 612 613 614 | Ruifei He, Jihan Yang, and Xiaojuan Qi. Re-distributing biased pseudo labels for semi-supervised semantic segmentation: A baseline investigation. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pp. 6930–6940, 2021. |
| 615 616 617 | Weixiang Hong, Zhenzhen Wang, Ming Yang, and Junsong Yuan. Conditional generative adversarial network for structured domain adaptation. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 1335–1344, 2018. |
| 619 620 621 | Lukas Hoyer, Dengxin Dai, and Luc Van Gool. Daformer: Improving network architectures and training strategies for domain-adaptive semantic segmentation. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 9924–9935, 2022a. |
| 622 623 624 | Lukas Hoyer, Dengxin Dai, and Luc Van Gool. Hrda: Context-aware high-resolution domain-adaptive semantic segmentation. In <i>European Conference on Computer Vision</i> , pp. 372–391. Springer, 2022b. |
| 625 626 627 628 | Lukas Hoyer, Dengxin Dai, Haoran Wang, and Luc Van Gool. Mic: Masked image consistency for context-enhanced domain adaptation. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 11721–11732, 2023. |
| 629 630 631 | Xinyue Huo, Lingxi Xie, Wengang Zhou, Houqiang Li, and Qi Tian. Focus on your target: A dual teacher-student framework for domain-adaptive semantic segmentation. <i>arXiv preprint arXiv:2303.09083</i> , 2023. |
| 632 633 634 | Ying Jin, Ximei Wang, Mingsheng Long, and Jianmin Wang. Minimum class confusion for versatile domain adaptation. In <i>Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXI 16</i> , pp. 464–480. Springer, 2020. |
| 636 637 638 | Jaehyung Kim, Jongheon Jeong, and Jinwoo Shin. M2m: Imbalanced classification via major-to- minor translation. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern</i> <i>recognition</i> , pp. 13896–13905, 2020. |
| 639 640 641 | Ruihuang Li, Shuai Li, Chenhang He, Yabin Zhang, Xu Jia, and Lei Zhang. Class-balanced pixel- level self-labeling for domain adaptive semantic segmentation. In <i>Proceedings of the IEEE/CVF</i> <i>Conference on Computer Vision and Pattern Recognition</i> , pp. 11593–11603, 2022. |
| 642 643 644 645 | Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In <i>Proceedings of the IEEE international conference on computer vision</i> , pp. 2980–2988, 2017. |
| 646 647 | Ziwei Liu, Zhongqi Miao, Xiaohang Zhan, Jiayun Wang, Boqing Gong, and Stella X Yu. Large-scale long-tailed recognition in an open world. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pp. 2537–2546, 2019. |

| 648 649 650 | Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 3431–3440, 2015a. |
|---------------------------------|--|
| 651 652 653 654 | Mingsheng Long, Yue Cao, Jianmin Wang, and Michael Jordan. Learning transferable features with deep adaptation networks. In <i>International conference on machine learning</i> , pp. 97–105. PMLR, 2015b. |
| 655 656 | Geoffrey J McLachlan and Thriyambakam Krishnan. <i>The EM algorithm and extensions</i> . John Wiley & Sons, 2007. |
| 657 658 659 | Aditya Krishna Menon, Sadeep Jayasumana, Ankit Singh Rawat, Himanshu Jain, Andreas Veit, and Sanjiv Kumar. Long-tail learning via logit adjustment. <i>arXiv preprint arXiv:2007.07314</i> , 2020. |
| 660 661 662 | Duo Peng, Ping Hu, Qiuhong Ke, and Jun Liu. Diffusion-based image translation with label guidance for domain adaptive semantic segmentation. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pp. 808–820, 2023. |
| 663 664 665 | Xingchao Peng, Ben Usman, Neela Kaushik, Judy Hoffman, Dequan Wang, and Kate Saenko. Visda: The visual domain adaptation challenge. <i>arXiv preprint arXiv:1710.06924</i> , 2017. |
| 666 667 668 | Harsh Rangwani, Sumukh K Aithal, Mayank Mishra, Arihant Jain, and Venkatesh Babu Radhakrish- nan. A closer look at smoothness in domain adversarial training. In <i>International conference on</i> <i>machine learning</i> , pp. 18378–18399. PMLR, 2022. |
| 669 670 671 | Stephan R Richter, Vibhav Vineet, Stefan Roth, and Vladlen Koltun. Playing for data: Ground truth from computer games. In <i>Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part II 14</i> , pp. 102–118. Springer, 2016. |
| 672 673 674 675 676 | German Ros, Laura Sellart, Joanna Materzynska, David Vazquez, and Antonio M Lopez. The synthia dataset: A large collection of synthetic images for semantic segmentation of urban scenes. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 3234–3243, 2016. |
| 677 678 679 | Fengyi Shen, Akhil Gurram, Ziyuan Liu, He Wang, and Alois Knoll. Diga: Distil to generalize and then adapt for domain adaptive semantic segmentation. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 15866–15877, 2023. |
| 680 681 682 683 | Marco Toldo, Umberto Michieli, Gianluca Agresti, and Pietro Zanuttigh. Unsupervised domain adaptation for mobile semantic segmentation based on cycle consistency and feature alignment. <i>Image and Vision Computing</i> , 95:103889, 2020. |
| 684 685 686 | Wilhelm Tranheden, Viktor Olsson, Juliano Pinto, and Lennart Svensson. Dacs: Domain adaptation via cross-domain mixed sampling. In <i>Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision</i> , pp. 1379–1389, 2021. |
| 687 688 689 | Thanh-Dat Truong, Ngan Le, Bhiksha Raj, Jackson Cothren, and Khoa Luu. Fredom: Fairness domain adaptation approach to semantic scene understanding. In <i>IEEE/CVF Computer Vision and Pattern Recognition (CVPR)</i> , 2023. |
| 690 691 692 693 | Yi-Hsuan Tsai, Wei-Chih Hung, Samuel Schulter, Kihyuk Sohn, Ming-Hsuan Yang, and Manmohan Chandraker. Learning to adapt structured output space for semantic segmentation. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 7472–7481, 2018. |
| 694 695 | Grant Van Horn and Pietro Perona. The devil is in the tails: Fine-grained classification in the wild. <i>arXiv preprint arXiv:1709.01450</i> , 2017. |
| 696 697 698 699 | Kaihong Wang, Donghyun Kim, Rogerio Feris, and Margrit Betke. Cdac: Cross-domain attention consistency in transformer for domain adaptive semantic segmentation. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pp. 11519–11529, 2023. |
| 700 701 | Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. <i>Advances in Neural</i> <i>Information Processing Systems</i> , 34:12077–12090, 2021. |

Pan Zhang, Bo Zhang, Ting Zhang, Dong Chen, Yong Wang, and Fang Wen. Prototypical pseudo label denoising and target structure learning for domain adaptive semantic segmentation. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 12414–12424, 2021. Dong Zhao, Shuang Wang, Qi Zang, Dou Quan, Xiutiao Ye, Rui Yang, and Licheng Jiao. Learning pseudo-relations for cross-domain semantic segmentation. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 19191–19203, 2023.

A DERIVATION OF EQ.7

Considering probabilities estimate that $p(y_{ij} = c | x_{ij}) \propto \exp(f_{\theta}(x_{ij})[c])$, the discriminant probability for pixel x_{ij} can be presented as:

$$p(y_{ij} = c | x_{ij}) = \frac{\exp(f_{\theta}(x_{ij})[c])}{\sum_{c'} \exp(f_{\theta}(x_{ij})[c'])}.$$
(11)

Given offset for logits as $\Delta_{cl}(z) = z' - z$ and the label c for pixel x_{ij} , we can obtain revised class-conditional discriminant probability for each pixel x_{ij} by:

$$\widetilde{p}(y_{ij} = c | x_{ij}, c) = \frac{\exp(f_{\theta}(x_{ij})[c] + \Delta_{cc}(f_{\theta}(x_{ij})[c]))}{\sum_{c'} \exp(f_{\theta}(x_{ij})[c'] + \Delta_{cc'}(f_{\theta}(x_{ij})[c']))}.$$
(12)

768 Then, following Bayes rule, the revised posterior is derived as:

$$\widetilde{p}(y_{ij} = c | x_{ij}) = \frac{p(c)\widetilde{p}(y_{ij} = c | x_{ij}, c)}{\sum_{c'} p(c')\widetilde{p}(y_{ij} = c | x_{ij}, c')}.$$
(13)

So we can obtain revised prediciton results for each pixel x_{ij} by:

$$\widetilde{y}_{ij} = \operatorname{argmax}_{c \in [C]} p(c) \widetilde{p}(y_{ij} = c | x_{ij}, c).$$
(14)

Since the class probabilities p(c) are typically set as a uniform prior, i.e., $p(c) = \frac{1}{c}$, Eq.14 can be rewritten as: $\widetilde{u}_{ij} = \operatorname{argmax}_{clcr} p(u_{ij} = c | x_{ij}, c)$

$$\widetilde{y}_{ij} = \operatorname{argmax}_{c \in [C]} p(y_{ij} = c | x_{ij}, c) = \operatorname{argmax}_{c \in [C]} \frac{\exp(f_{\theta}(x_{ij})[c] + \Delta_{cc}(f_{\theta}(x_{ij})[c]))}{\exp(f_{\theta}(x_{ij})[c] + \Delta_{cc}(f_{\theta}(x_{ij})[c])) + \sum_{c' \neq c} \exp(f_{\theta}(x_{ij})[c'] + \Delta_{cc'}(f_{\theta}(x_{ij})[c']))}.$$
(15)

For Eq.6, we add a scaling factor τ to be adjusted for specific evaluation metrics.

B DISCUSSION ABOUT EQ.9

For a deeper understanding of the loss function in Eq.9, we can rewrite it as:

$$\widetilde{\mathcal{L}}^{s} = \frac{1}{N_{S}} \sum_{i=1}^{N_{S}} \sum_{j=1}^{H \times W} \log \left(1 + \sum_{c \neq y_{ij}^{S}} \left(\frac{\exp\left(\Delta_{y_{ij}^{S}, y_{ij}^{S}}(f_{\theta}(x_{ij}^{S})[y_{ij}^{S}])\right)}{\exp\left(\Delta_{y_{ij}^{S}, c}^{S}(f_{\theta}(x_{ij}^{S})[c])\right)} \right)^{\tau} \exp\left(f_{\theta}(x_{ij}^{S})[c] - f_{\theta}(x_{ij}^{S})[y_{ij}^{S}]\right)\right),\tag{16}$$

which can be interpreted as a standard cross-entropy loss with an adaptive margin. Specifically, if the class y_{ij}^S is an over-predicted class, it is reasonable to assume that for most other classes c, $\Delta_{y_{ij}^S, y_{ij}^S}^S(f_{\theta}(x_{ij}^S)[y_{ij}^S]) < \Delta_{y_{ij}^S, c}^S(f_{\theta}(x_{ij}^S)[c])$. This implies that:

$$\left(\frac{\exp\left(\Delta_{y_{ij}^S, y_{ij}^S}^S(f_\theta(x_{ij}^S)[y_{ij}^S])\right)}{\exp\left(\Delta_{y_{ij}^S, c}^S(f_\theta(x_{ij}^S)[c])\right)}\right)^{\tau} < 1$$

$$(17)$$

In the loss function Eq.16, this ratio is used to scale the term $\exp\left(f_{\theta}(x_{ij}^S)[c] - f_{\theta}(x_{ij}^S)[y_{ij}^S]\right)$. When y_{ij}^S is an over-predicted class, the ratio is less than 1, which reduces the weight of the term $\exp\left(f_{\theta}(x_{ij}^S)[c] - f_{\theta}(x_{ij}^S)[y_{ij}^S]\right)$. Consequently, for over-predicted classes, the loss function imposes a smaller penalty for misclassification. Conversely, if y_{ij}^S is an under-predicted class, it can be assumed that for most other classes c, $\Delta_{y_{ij}^S, y_{ij}^S}^S(f_{\theta}(x_{ij}^S)[y_{ij}^S]) > \Delta_{y_{ij}^S, c}^S(f_{\theta}(x_{ij}^S)[c])$. This leads to a ratio greater than 1, which increases the weight of the term $\exp\left(f_{\theta}(x_{ij}^S)[c] - f_{\theta}(x_{ij}^S)[y_{ij}^S]\right)$, thereby imposing a larger penalty for misclassification of under-predicted classes. In summary, by introducing the margin-based ratio term, the loss function Eq.16 adaptively adjusts the strength of the penalty based on the difficulty of the classes. This approach helps to mitigate the class bias problem and enhances the model's performance on under-predicted classes, leading to a more balanced and accurate classification. Moreover, since the logit offset term Δ is updated online during the training process, it aligns well with the self-training paradigm in UDA.

810 C DISCUSSION ABOUT EVALUATION METRICS 811

819 820 821

823

824

825

827 828 829

843

844

845 846

In this section, we investigate the impact of class imbalance on the evaluation metrics mIoU and mAcc in the context of semantic segmentation. We consider a multi-class problem with C classes, where the number of samples in the *i*-th class is denoted as N_i . Let P_{ij} represent the probability of classifying a sample from the *i*-th class as the *j*-th class in the confusion matrix. For the purpose of this analysis, we assume that the probabilities P_{ii} and P_{kk} are balanced across all classes, i.e., $P_{ii} = P_{kk} = p, \forall i, k \in 1, 2, ..., C$, and focus solely on the effect of class imbalance in terms of sample numbers. Under this assumption, the calculation formula for mAcc can be written as:

$$mAcc = \frac{1}{C} \sum_{i=1}^{C} \frac{N_i \cdot P_{ii}}{\sum_{j=1}^{C} N_i \cdot P_{ij}} = \frac{1}{C} \sum_{i=1}^{C} \frac{p}{\sum_{j=1}^{C} P_{ij}}$$
(18)

As evident from the equation above, the sample numbers N_i cancel out in the calculation of mAcc, making it independent of the class imbalance in terms of sample numbers. Therefore, mAcc remains unaffected by class imbalance under the given assumptions. On the other hand, the calculation formula for mIoU is given by:

$$mIoU = \frac{1}{C} \sum_{i=1}^{C} \frac{N_i \cdot P_{ii}}{\sum_{j=1}^{C} N_i \cdot P_{ij} + \sum_{j=1}^{C} N_j \cdot P_{ji} - N_i \cdot P_{ii}}$$
(19)

830 In contrast to mAcc, the sample numbers N_i do not cancel out in the calculation of mIoU. Con-831 sequently, when the classes are imbalanced, i.e., the sample numbers N_i vary significantly across 832 classes, the IoU of classes with larger sample numbers will dominate the overall mIoU result. To 833 illustrate the impact of class imbalance on mIoU, let us consider a case where class k is a head class with a significantly larger sample number N_k compared to a tail class i with sample number N_i . The 834 performance of class i will be greatly affected by class k through the term $N_k \cdot P_{ki}$ in the denominator 835 of mIoU. For class i to have a fair contribution to mIoU, the probability P_{ki} needs to be very small. 836 This implies that a balanced classifier may not be optimal for maximizing mIoU under class imbal-837 ance. Therefore, our proposed method implements a scaling factor τ to modulate the contributions of 838 introduced logits offset. In contrast, mAcc is inherently unaffected by class imbalance, as the sample 839 numbers N_i cancel out in its calculation formula. This means that a balanced classifier is indeed 840 optimal for maximizing mAcc, and our method, which aims to balance the contributions of different 841 classes, aligns well with this objective and can achieve consistent improvements. 842

D IMPLEMENTATION DETAILS OF ONLINE LOGITS DISTRIBUTION ESTIMATION

In this section, we provide the pseudo label to explain implementation details of online logits 847 distribution estimation, as shown in Alg.1. For computational efficiency, in each iteration, we sample 848 N_{sample} logits for each element in M_{cl}^S and M_{cl}^T , where N_{sample} is the minimum sample number of 849 classes in the minibatch (shoule be greater than or equal to N_{min} , where we set it as 100). This way, 850 we obtain $C \times C \times N_{sample}$ logits, and then update the $C \times C$ GMMs simultaneously in a parallel 851 manner. Since not all classes may be updated in each iteration, we maintain a variable n for each 852 GMM that is not updated, to record the number of iterations since its last update. This n is used to 853 adjust the momentum factor using $\tilde{\tau}$ in the EMA update, in order to match the update speed of the 854 network.

In the algorithm, the cumulative distribution functions (CDFs) F_{cl}^s , F_{cl}^t , F_p , and F_n are computed using the estimated Gaussian mixture models (GMMs). These CDFs describe the cumulative probability distribution of the corresponding GMMs. For the source and target domain GMMs P_{cl}^s and P_{cl}^t , their CDFs can be represented as: $F_{cl}^s(z) = \sum_{k=1}^K \pi_{cl,k}^s \cdot \Phi(\frac{z-\mu_{cl,k}^s}{\sigma_{cl,k}^s}) F_{cl}^t(z) = \sum_{k=1}^K \pi_{cl,k}^t \cdot \Phi(\frac{z-\mu_{cl,k}^s}{\sigma_{cl,k}^s})$, where $\Phi(\cdot)$ is the CDF of the standard normal distribution, and $\pi_{cl,k}^s, \mu_{cl,k}^s$, and $\sigma_{cl,k}^s$ are the weight, mean, and standard deviation of the k-th component of the source domain GMM P_{cl}^s , respectively. $\pi_{cl,k}^t, \mu_{cl,k}^t$, and $\sigma_{cl,k}^t$ are the corresponding parameters for the target domain GMM P_{cl}^t . Similarly, the CDFs for the anchor GMMs P_p and P_n are: $F_p(z) = \sum_{k=1}^K \pi_{p,k} \cdot \Phi(\frac{z-\mu_{p,k}}{\sigma_{p,k}})$

| : Source domain $D_s = (x_i^S, y_i^S)_{i=1}^{N_S}$, target domain $D_t = (x_i^T)_{i=1}^{N_T}$, number of classes C , imber of Gaussian components K , momentum factor $\tilde{\tau}$, scaling factor τ , minimum number of ements N_{min} , number of EM Loop T . ut : Model parameters θ . itialize model parameters θ , source GMMs P_{cl}^s , target GMMs P_{cl}^t , anchor GMMs P_p and P_n or all $c, l \in [C]$. hile not converged do Sample a mini-batch of source data $(x_i^S, y_i^S)_{i=1}^{B_S}$ and target data $(x_i^T)_{i=1}^{B_T}$. Compute logits $f_{\theta}(x_{ij}^S)$ and $f_{\theta}(x_{ij}^T)$ for source and target samples. Compute pseudo-labels for target samples: $\hat{y}_{ij}^T = \operatorname{argmax}_{c \in [C]} f_{\theta}(x_{ij}^T)[c]$). Compute matrices M_{cl}^S and M_{cl}^T based on the source labels and target pseudo-labels: $M_{cl}^T = \{f_{\theta}(x_{ij}^S)[l] \mid \hat{y}_{ij}^T = c\}$ Update source GMMs P_{cl}^s , target GMMs P_{cl}^t , anchor GMMs P_p and P_n using the momentum- based EM algorithm: for $c, l \in [C]$ do if $ M_{cl}^S > N_{min}$ then Initialize $\phi_{cl}^{s,(0)} \leftarrow \hat{\phi}_{cl}^s$ from the latest iteration. for $t = 1, \dots, T$ do Update $\phi cl^{s,(t)}$ using the EM algorithm with current logits. end for $\hat{\phi} cl^s \leftarrow (1 - \tilde{\tau}^n)\phi_{cl}^{s,(T)} + \tilde{\tau}^n \hat{\phi} cl^s$, where n is the number of iterations since the last undate. |
|--|
| Imber of Gaussian components K , momentum factor $\tilde{\tau}$, scaling factor τ , minimum number of ements N_{min} , number of EM Loop T . It: Model parameters θ . itialize model parameters θ , source GMMs P_{cl}^s , target GMMs P_{cl}^t , anchor GMMs P_p and P_n or all $c, l \in [C]$. hile not converged do Sample a mini-batch of source data $(x_i^S, y_i^S)_{i=1}^{B_S}$ and target data $(x_i^T)_{i=1}^{B_T}$. Compute logits $f_{\theta}(x_{ij}^S)$ and $f_{\theta}(x_{ij}^T)$ for source and target samples. Compute pseudo-labels for target samples: $\hat{y}_{ij}^T = \operatorname{argmax}_{c \in [C]} f_{\theta}(x_{ij}^T[c])$. Compute matrices M_{cl}^S and M_{cl}^T based on the source labels and target pseudo-labels: $M_{cl}^T = \{f_{\theta}(x_{ij}^S)[l] \mid \hat{y}_{ij}^S = c\}$ $M_{cl}^T = \{f_{\theta}(x_{ij}^T)[l] \mid \hat{y}_{ij}^T = c\}$ Update source GMMs P_{cl}^s , target GMMs P_{cl}^t , anchor GMMs P_p and P_n using the momentum- based EM algorithm: for $c, l \in [C]$ do if $ M_{cl}^S > N_{min}$ then Initialize $\phi_{cl}^{s,(0)} \leftarrow \hat{\phi}_{cl}^s$ from the latest iteration. for $t = 1, \dots, T$ do Update $\phi cl^{s,(t)}$ using the EM algorithm with current logits. end for $\hat{\phi}cl^s \leftarrow (1 - \tilde{\tau}^n)\phi_{cl}^{s,(T)} + \tilde{\tau}^n \hat{\phi}cl^s$, where n is the number of iterations since the last undate. |
| ements N_{min} , number of EM Loop T . ut: Model parameters θ . itialize model parameters θ , source GMMs P_{cl}^s , target GMMs P_{cl}^t , anchor GMMs P_p and P_n or all $c, l \in [C]$. hile not converged do Sample a mini-batch of source data $(x_i^S, y_i^S)_{i=1}^{B_S}$ and target data $(x_i^T)_{i=1}^{B_T}$. Compute logits $f_{\theta}(x_{ij}^S)$ and $f_{\theta}(x_{ij}^T)$ for source and target samples. Compute pseudo-labels for target samples: $\hat{y}_{ij}^T = \operatorname{argmax}_{c \in [C]} f_{\theta}(x_{ij}^T[c])$. Compute matrices M_{cl}^S and M_{cl}^T based on the source labels and target pseudo-labels: $M_{cl}^S = \{f_{\theta}(x_{ij}^S)[l] \mid y_{ij}^S = c\}$ $M_{cl}^T = \{f_{\theta}(x_{ij}^T)[l] \mid \hat{y}_{ij}^T = c\}$ Update source GMMs P_{cl}^s , target GMMs P_{cl}^t , anchor GMMs P_p and P_n using the momentum based EM algorithm: for $c, l \in [C]$ do if $ M_{cl}^S > N_{min}$ then Initialize $\phi_{cl}^{s,(0)} \leftarrow \hat{\phi}_{cl}^s$ from the latest iteration. for $t = 1, \dots, T$ do Update $\phi cl^{s,(t)}$ using the EM algorithm with current logits. end for $\hat{\phi}cl^S \leftarrow (1 - \tilde{\tau}^n)\phi_{cl}^{s,(T)} + \tilde{\tau}^n\hat{\phi}cl^s$, where n is the number of iterations since the last undate. |
| It: Model parameters θ . itialize model parameters θ , source GMMs P_{cl}^s , target GMMs P_{cl}^t , anchor GMMs P_p and P_r , r all $c, l \in [C]$. hile not converged do Sample a mini-batch of source data $(x_i^S, y_i^S)_{i=1}^{B_S}$ and target data $(x_i^T)_{i=1}^{B_T}$. Compute logits $f_{\theta}(x_{ij}^S)$ and $f_{\theta}(x_{ij}^T)$ for source and target samples. Compute pseudo-labels for target samples: $\hat{y}_{ij}^T = \operatorname{argmax}_{c \in [C]} f_{\theta}(x_{ij}^T[c])$. Compute matrices M_{cl}^S and M_{cl}^T based on the source labels and target pseudo-labels: $M_{cl}^S = \{f_{\theta}(x_{ij}^S)[l] \mid y_{ij}^S = c\}$ $M_{cl}^T = \{f_{\theta}(x_{ij}^T)[l] \mid \hat{y}_{ij}^T = c\}$ Update source GMMs P_{cl}^s , target GMMs P_{cl}^t , anchor GMMs P_p and P_n using the momentum based EM algorithm: for $c, l \in [C]$ do if $ M_{cl}^S > N_{min}$ then Initialize $\phi_{cl}^{s,(0)} \leftarrow \hat{\phi}_{cl}^s$ from the latest iteration. for $t = 1, \dots, T$ do Update $\phi cl^{s,(t)}$ using the EM algorithm with current logits. end for $\hat{\phi}cl^s \leftarrow (1 - \tilde{\tau}^n)\phi_{cl}^{s,(T)} + \tilde{\tau}^n \hat{\phi}cl^s$, where n is the number of iterations since the lass update. |
| Initialize model parameters θ , source GMMs P_{cl}^{s} , target GMMs P_{cl}^{t} , anchor GMMs P_{p} and P_{n} or all $c, l \in [C]$. hile not converged do Sample a mini-batch of source data $(x_{i}^{S}, y_{i}^{S})_{i=1}^{B_{S}}$ and target data $(x_{i}^{T})_{i=1}^{B_{T}}$. Compute logits $f_{\theta}(x_{ij}^{S})$ and $f_{\theta}(x_{ij}^{T})$ for source and target samples. Compute pseudo-labels for target samples: $\hat{y}_{ij}^{T} = \operatorname{argmax}_{c \in [C]} f_{\theta}(x_{ij}^{T}[c])$. Compute matrices M_{cl}^{S} and M_{cl}^{T} based on the source labels and target pseudo-labels: $M_{cl}^{S} = \{f_{\theta}(x_{ij}^{S})[l] \mid y_{ij}^{S} = c\}$ $M_{cl}^{T} = \{f_{\theta}(x_{ij}^{T})[l] \mid \hat{y}_{ij}^{T} = c\}$ Update source GMMs P_{cl}^{s} , target GMMs P_{cl}^{t} , anchor GMMs P_{p} and P_{n} using the momentum based EM algorithm: for $c, l \in [C]$ do if $ M_{cl}^{S} > N_{min}$ then Initialize $\phi_{cl}^{s,(0)} \leftarrow \hat{\phi}_{cl}^{s}$ from the latest iteration. for $t = 1, \dots, T$ do Update $\phi cl^{s,(t)}$ using the EM algorithm with current logits. end for $\hat{\phi}cl^{s} \leftarrow (1 - \tilde{\tau}^{n})\phi_{cl}^{s,(T)} + \tilde{\tau}^{n}\hat{\phi}cl^{s}$, where n is the number of iterations since the lase undate. |
| hile not converged do Sample a mini-batch of source data $(x_i^S, y_i^S)_{i=1}^{B_S}$ and target data $(x_i^T)_{i=1}^{B_T}$. Compute logits $f_{\theta}(x_{ij}^S)$ and $f_{\theta}(x_{ij}^T)$ for source and target samples. Compute pseudo-labels for target samples: $\hat{y}_{ij}^T = \operatorname{argmax}_{c \in [C]} f_{\theta}(x_{ij}^T[c])$. Compute matrices M_{cl}^S and M_{cl}^T based on the source labels and target pseudo-labels: $M_{cl}^S = \{f_{\theta}(x_{ij}^S)[l] \mid y_{ij}^S = c\}$ $M_{cl}^T = \{f_{\theta}(x_{ij}^T)[l] \mid \hat{y}_{ij}^T = c\}$ Update source GMMs P_{cl}^s , target GMMs P_{cl}^t , anchor GMMs P_p and P_n using the momentum based EM algorithm: for $c, l \in [C]$ do if $ M_{cl}^S > N_{min}$ then Initialize $\phi_{cl}^{s,(0)} \leftarrow \hat{\phi}_{cl}^s$ from the latest iteration. for $t = 1, \dots, T$ do Update $\phi cl^{s,(t)}$ using the EM algorithm with current logits. end for $\hat{\phi}cl^s \leftarrow (1 - \tilde{\tau}^n)\phi_{cl}^{s,(T)} + \tilde{\tau}^n\hat{\phi}cl^s$, where n is the number of iterations since the last undate. |
| Sample a mini-batch of source data $(x_i^S, y_i^S)_{i=1}^{B_S}$ and target data $(x_i^T)_{i=1}^{B_T}$. Compute logits $f_{\theta}(x_{ij}^S)$ and $f_{\theta}(x_{ij}^T)$ for source and target samples. Compute pseudo-labels for target samples: $\hat{y}_{ij}^T = \operatorname{argmax}_{c \in [C]} f_{\theta}(x_{ij}^T[c])$. Compute matrices M_{cl}^S and M_{cl}^T based on the source labels and target pseudo-labels: $M_{cl}^S = \{f_{\theta}(x_{ij}^S)[l] \mid y_{ij}^S = c\}$ $M_{cl}^T = \{f_{\theta}(x_{ij}^T)[l] \mid \hat{y}_{ij}^T = c\}$ Update source GMMs P_{cl}^s , target GMMs P_{cl}^t , anchor GMMs P_p and P_n using the momentum based EM algorithm: for $c, l \in [C]$ do if $ M_{cl}^S > N_{min}$ then Initialize $\phi_{cl}^{s,(0)} \leftarrow \hat{\phi}_{cl}^s$ from the latest iteration. for $t = 1, \dots, T$ do Update $\phi cl^{s,(t)}$ using the EM algorithm with current logits. end for $\hat{\phi}cl^s \leftarrow (1 - \tilde{\tau}^n)\phi_{cl}^{s,(T)} + \tilde{\tau}^n\hat{\phi}cl^s$, where n is the number of iterations since the last undate. |
| Sample a mini-batch of source data $(x_i^r, y_i^r)_{i=1}$ and target data $(x_i^r)_{i=1}$. Compute logits $f_{\theta}(x_{ij}^S)$ and $f_{\theta}(x_{ij}^T)$ for source and target samples. Compute pseudo-labels for target samples: $\hat{y}_{ij}^T = \operatorname{argmax}_{c \in [C]} f_{\theta}(x_{ij}^T[c])$. Compute matrices M_{cl}^S and M_{cl}^T based on the source labels and target pseudo-labels: $M_{cl}^S = \{f_{\theta}(x_{ij}^S)[l] \mid y_{ij}^S = c\}$ $M_{cl}^T = \{f_{\theta}(x_{ij}^T)[l] \mid \hat{y}_{ij}^T = c\}$ Update source GMMs P_{cl}^s , target GMMs P_{cl}^t , anchor GMMs P_p and P_n using the momentum based EM algorithm: for $c, l \in [C]$ do if $ M_{cl}^S > N_{min}$ then Initialize $\phi_{cl}^{s,(0)} \leftarrow \hat{\phi}_{cl}^s$ from the latest iteration. for $t = 1, \dots, T$ do Update $\phi cl^{s,(t)}$ using the EM algorithm with current logits. end for $\hat{\phi}cl^s \leftarrow (1 - \tilde{\tau}^n)\phi_{cl}^{s,(T)} + \tilde{\tau}^n\hat{\phi}cl^s$, where n is the number of iterations since the last undate. |
| Compute logits $f_{\theta}(x_{ij}^{*})$ and $f_{\theta}(x_{ij}^{*})$ for source and target samples. Compute pseudo-labels for target samples: $\hat{y}_{ij}^{T} = \operatorname{argmax}_{c \in [C]} f_{\theta}(x_{ij}^{T}[c])$. Compute matrices M_{cl}^{S} and M_{cl}^{T} based on the source labels and target pseudo-labels: $M_{cl}^{S} = \{f_{\theta}(x_{ij}^{S})[l] \mid y_{ij}^{S} = c\}$ $M_{cl}^{T} = \{f_{\theta}(x_{ij}^{T})[l] \mid \hat{y}_{ij}^{T} = c\}$ Update source GMMs P_{cl}^{s} , target GMMs P_{cl}^{t} , anchor GMMs P_{p} and P_{n} using the momentum based EM algorithm: for $c, l \in [C]$ do if $ M_{cl}^{S} > N_{min}$ then Initialize $\phi_{cl}^{s,(0)} \leftarrow \hat{\phi}_{cl}^{s}$ from the latest iteration. for $t = 1, \dots, T$ do Update $\phi cl^{s,(t)}$ using the EM algorithm with current logits. end for $\hat{\phi}cl^{s} \leftarrow (1 - \tilde{\tau}^{n})\phi_{cl}^{s,(T)} + \tilde{\tau}^{n}\hat{\phi}cl^{s}$, where n is the number of iterations since the last undate. |
| Compute pseudo-labels for target samples: $y_{ij}^{t} = \operatorname{argmax}_{c \in [C]} f_{\theta}(x_{ij}^{t} c)$. Compute matrices M_{cl}^{S} and M_{cl}^{T} based on the source labels and target pseudo-labels: $M_{cl}^{S} = \{f_{\theta}(x_{ij}^{S})[l] \mid y_{ij}^{S} = c\}$ $M_{cl}^{T} = \{f_{\theta}(x_{ij}^{T})[l] \mid \hat{y}_{ij}^{T} = c\}$ Update source GMMs P_{cl}^{s} , target GMMs P_{cl}^{t} , anchor GMMs P_{p} and P_{n} using the momentum based EM algorithm: for $c, l \in [C]$ do if $ M_{cl}^{S} > N_{min}$ then Initialize $\phi_{cl}^{s,(0)} \leftarrow \hat{\phi}_{cl}^{s}$ from the latest iteration. for $t = 1, \dots, T$ do Update $\phi cl^{s,(t)}$ using the EM algorithm with current logits. end for $\hat{\phi}cl^{s} \leftarrow (1 - \tilde{\tau}^{n})\phi_{cl}^{s,(T)} + \tilde{\tau}^{n}\hat{\phi}cl^{s}$, where n is the number of iterations since the last undate. |
| Compute matrices M_{cl}^{S} and M_{cl}^{T} based on the source labels and target pseudo-labels: $M_{cl}^{S} = \{f_{\theta}(x_{ij}^{S})[l] \mid y_{ij}^{S} = c\}$ $M_{cl}^{T} = \{f_{\theta}(x_{ij}^{T})[l] \mid \hat{y}_{ij}^{T} = c\}$ Update source GMMs P_{cl}^{s} , target GMMs P_{cl}^{t} , anchor GMMs P_{p} and P_{n} using the momentum based EM algorithm: for $c, l \in [C]$ do if $ M_{cl}^{S} > N_{min}$ then Initialize $\phi_{cl}^{s,(0)} \leftarrow \hat{\phi}_{cl}^{s}$ from the latest iteration. for $t = 1, \dots, T$ do Update $\phi cl^{s,(t)}$ using the EM algorithm with current logits. end for $\hat{\phi}cl^{s} \leftarrow (1 - \tilde{\tau}^{n})\phi_{cl}^{s,(T)} + \tilde{\tau}^{n}\hat{\phi}cl^{s}$, where n is the number of iterations since the last undate. |
| $\begin{split} M_{cl}^{3} &= \{f_{\theta}(x_{ij}^{s})[l] \mid y_{ij}^{s} = c\} \\ M_{cl}^{T} &= \{f_{\theta}(x_{ij}^{T})[l] \mid \hat{y}_{ij}^{T} = c\} \\ \text{Update source GMMs } P_{cl}^{s}, \text{ target GMMs } P_{cl}^{t}, \text{ anchor GMMs } P_{p} \text{ and } P_{n} \text{ using the momentum based EM algorithm:} \\ \text{for } c, l \in [C] \text{ do} \\ \text{ if } M_{cl}^{S} > N_{min} \text{ then} \\ \text{Initialize } \phi_{cl}^{s,(0)} \leftarrow \hat{\phi}_{cl}^{s} \text{ from the latest iteration.} \\ \text{ for } t = 1, \dots, T \text{ do} \\ \text{ Update } \phi cl^{s,(t)} \text{ using the EM algorithm with current logits.} \\ \text{ end for} \\ \hat{\phi}cl^{s} \leftarrow (1 - \tilde{\tau}^{n})\phi_{cl}^{s,(T)} + \tilde{\tau}^{n}\hat{\phi}cl^{s}, \text{ where } n \text{ is the number of iterations since the last undate.} \end{split}$ |
| $\begin{split} M_{cl}^{I} &= \{f_{\theta}(x_{ij}^{I})[l] \mid \hat{y}_{ij}^{I} = c\} \\ \text{Update source GMMs } P_{cl}^{s}, \text{ target GMMs } P_{cl}^{t}, \text{ anchor GMMs } P_{p} \text{ and } P_{n} \text{ using the momentum based EM algorithm:} \\ \text{for } c, l \in [C] \text{ do } \\ \text{if } M_{cl}^{S} > N_{min} \text{ then} \\ \text{Initialize } \phi_{cl}^{s,(0)} \leftarrow \hat{\phi}_{cl}^{s} \text{ from the latest iteration.} \\ \text{for } t = 1, \dots, T \text{ do } \\ \text{Update } \phi cl^{s,(t)} \text{ using the EM algorithm with current logits.} \\ \text{end for} \\ \hat{\phi} cl^{s} \leftarrow (1 - \tilde{\tau}^{n})\phi_{cl}^{s,(T)} + \tilde{\tau}^{n}\hat{\phi}cl^{s}, \text{ where } n \text{ is the number of iterations since the lass undate.} \end{split}$ |
| Update source GMMs P_{cl}^s , target GMMs P_{cl}^t , anchor GMMs P_p and P_n using the momentum based EM algorithm: for $c, l \in [C]$ do if $ M_{cl}^S > N_{min}$ then Initialize $\phi_{cl}^{s,(0)} \leftarrow \hat{\phi}_{cl}^s$ from the latest iteration. for $t = 1, \dots, T$ do Update $\phi cl^{s,(t)}$ using the EM algorithm with current logits. end for $\hat{\phi}cl^s \leftarrow (1 - \tilde{\tau}^n)\phi_{cl}^{s,(T)} + \tilde{\tau}^n\hat{\phi}cl^s$, where n is the number of iterations since the lass update. |
| based EM algorithm: for $c, l \in [C]$ do if $ M_{cl}^S > N_{min}$ then Initialize $\phi_{cl}^{s,(0)} \leftarrow \hat{\phi}_{cl}^s$ from the latest iteration. for $t = 1, \dots, T$ do Update $\phi cl^{s,(t)}$ using the EM algorithm with current logits. end for $\hat{\phi}cl^s \leftarrow (1 - \tilde{\tau}^n)\phi_{cl}^{s,(T)} + \tilde{\tau}^n\hat{\phi}cl^s$, where n is the number of iterations since the last update. |
| for $c, t \in [C]$ do if $ M_{cl}^{s} > N_{min}$ then Initialize $\phi_{cl}^{s,(0)} \leftarrow \hat{\phi}_{cl}^{s}$ from the latest iteration. for $t = 1,, T$ do Update $\phi cl^{s,(t)}$ using the EM algorithm with current logits. end for $\hat{\phi}cl^{s} \leftarrow (1 - \tilde{\tau}^{n})\phi_{cl}^{s,(T)} + \tilde{\tau}^{n}\hat{\phi}cl^{s}$, where <i>n</i> is the number of iterations since the last update. |
| In $ M_{cl} > N_{min}$ then Initialize $\phi_{cl}^{s,(0)} \leftarrow \hat{\phi}_{cl}^s$ from the latest iteration. for $t = 1,, T$ do Update $\phi cl^{s,(t)}$ using the EM algorithm with current logits. end for $\hat{\phi}cl^s \leftarrow (1 - \tilde{\tau}^n)\phi_{cl}^{s,(T)} + \tilde{\tau}^n\hat{\phi}cl^s$, where <i>n</i> is the number of iterations since the last update. |
| Initialize $\phi_{cl} \leftarrow \phi_{cl}^s$ from the latest iteration. for $t = 1,, T$ do Update $\phi cl^{s,(t)}$ using the EM algorithm with current logits. end for $\hat{\phi}cl^s \leftarrow (1 - \tilde{\tau}^n)\phi_{cl}^{s,(T)} + \tilde{\tau}^n \hat{\phi}cl^s$, where <i>n</i> is the number of iterations since the last update. |
| In $t = 1,, 1$ do Update $\phi c l^{s,(t)}$ using the EM algorithm with current logits. end for $\hat{\phi} c l^s \leftarrow (1 - \tilde{\tau}^n) \phi_{cl}^{s,(T)} + \tilde{\tau}^n \hat{\phi} c l^s$, where <i>n</i> is the number of iterations since the last update. |
| end for $\hat{\phi}cl^s \leftarrow (1 - \tilde{\tau}^n)\phi_{cl}^{s,(T)} + \tilde{\tau}^n\hat{\phi}cl^s$, where <i>n</i> is the number of iterations since the last update. |
| $\hat{\phi}cl^s \leftarrow (1 - \tilde{\tau}^n)\phi_{cl}^{s,(T)} + \tilde{\tau}^n \hat{\phi}cl^s$, where <i>n</i> is the number of iterations since the last update. |
| update. |
| |
| end if |
| if $ M_{cl}^S > N_{min}$ then |
| Initialize $\phi_{cl}^{t,(0)} \leftarrow \hat{\phi}_{cl}^{t}$ from the latest iteration. |
| for $t = 1, \ldots, T$ do |
| Update $\phi c l^{t,(t)}$ using the EM algorithm with current logits. |
| end for \hat{t} |
| $\phi_{cl}^{\iota} \leftarrow (1 - \tilde{\tau}^n) \phi_{cl}^{\prime,(1)} + \tilde{\tau}^n \phi_{cl}^{\iota}$, where n is the number of iterations since the last update |
| ena n ord for |
| Undate anchor GMMs $P_{\rm c}$ and $P_{\rm c}$ using the global positive and negative logits from the source |
| domain: |
| Initialize $\phi_p^{(0)} \leftarrow \hat{\phi}_p, \phi_n^{(0)} \leftarrow \hat{\phi}_n$ from the latest iteration. |
| for $t = 1, \dots, T$ do |
| Update $\phi_p^{(t)}$, $\phi_n^{(t)}$ using the EM algorithm with current global logits. |
| end for |
| $\phi_p \leftarrow (1 - \tilde{\tau}) \phi_p^{(1)} + \tilde{\tau} \phi_p$ |
| $\hat{\phi}_n \leftarrow (1 - \tilde{\tau})\phi_n^{(T)} + \tilde{\tau}\hat{\phi}_n$ |
| Compute cumulative distributions F_{cl}^s , F_{cl}^t , F_p , F_n using the estimated GMMs. |
| Compute logits offsets for source domain: $\Delta^{S}(z) = \int F_{p}^{-1}(Fcl^{s}(z)) - z$, if $c = l$ |
| Compute logits offsets for source domain: $\Delta_{cl}(z) = \left\{ F_n^{-1}(Fcl^s(z)) - z, \text{ if } c \neq l \right\}$ |
| $\int F_n^{-1}(F_{cl}^t(z)) - z, \text{if } c = l$ |
| Compute logits offsets for target domain: $\Delta_{cl}^{-}(z) = \begin{cases} F_n^{-1}(F_{cl}^t(z)) - z, & \text{if } c \neq l \end{cases}$ |
| Compute the adjusted losses $\tilde{\mathcal{L}}^s$ and $\tilde{\mathcal{L}}^u$ using Eq. equation 8 with Λ^S , and Λ^T . |
| Undate model parameters \hat{A} by minimizing $\tilde{C}^s + \tilde{C}^u$ using an optimizer (e.g. SGD or Adam) |
| id while |
| eturn Model parameters θ . |
| • |
| |
| |
| |

| 918 919 920 | Table study c | 5: Parameter of K . | Table 6 study of | : Parameter $\tilde{\tau}$. | Table study o | 7: Parameter of T . | Table 8 study of | Parameter λ . |
|-------------------|------------------|-----------------------|--------------------------|------------------------------|------------------|-----------------------|----------------------|-----------------------|
| 921 | K | mIou (%) | $\frac{1}{\tilde{\tau}}$ | mIou (%) | \overline{T} | mIou (%) | $\overline{\lambda}$ | mIou (%) |
| 922 | 1 | 70.2 | 0 | 68.6 | 1 | 70.4 | 0.05 | 70.3 |
| 923 | 3 | 70.5 | 0.9 | 70.0 | 3 | 70.7 | 0.2 | 70.7 |
| 924 | 5 | 70.7 | 0.99 | 70.7 | 5 | 70.7 | 0.5 | 70.4 |
| 925 | 10 | 70.6 | 0.999 | 70.3 | 10 | 70.5 | 1 | 69.9 |
| 000 | | | - | | | | - | |

927 value of the variable corresponding to a given cumulative probability. For a given logit value z, by 928 computing $F_p^{-1}(F_{cl}^s(z))$ and $F_n^{-1}(F_{cl}^s(z))$, we obtain the corresponding logit values of the positive 929 anchor distribution P_p and the negative anchor distribution P_n at the cumulative probability $F_{cl}^s(z)$. 930 Then, the difference between these values and the original logit value z is used as the logits offset $\Delta_{cl}^{S}(z)$ for the source domain. Similarly, by computing $F_{p}^{-1}(F_{cl}^{t}(z))$ and $F_{n}^{-1}(F_{cl}^{t}(z))$, we obtain the 931 932 logits offset $\Delta_{cl}^T(z)$ for the target domain. To efficiently compute the CDF and its inverse function, 933 we use the Abramowitz-Stegun formula to approximate the CDF in the form of a polynomial and 934 employ interpolation methods to estimate the inverse function.

E INFLUENCE OF PARAMETERS SETTING

In this section, we further study the influence of parameters setting introduced in BLDA, i.e., number of Gaussian components K, momentum factor $\tilde{\tau}$, number of EM Loop T for GMM estimation and λ for cumulative density estimation loss. All experiments are conducted with DAFormer (Hoyer et al., 2022a) on GTA \rightarrow CS.

943 Number of Gaussian Components K. As shown in Tab.5, we find that BLDA can work well even 944 when K = 1, since the logit distribution is naturally close to Gaussian. The model achieves the best 945 performance when K = 5. A larger K allows for more flexibility in modeling the logits distributions 946 but may also introduce noise. We choose K = 5 as a balance between model capacity and robustness.

947 **Momentum Factor** $\tilde{\tau}$. The momentum factor $\tilde{\tau}$ controls the speed of updating the GMM parameters. 948 When $\tilde{\tau} = 0$, performance becomes erratic because the logits from the current iteration alone are 949 not sufficient to model the distribution of all logits. A larger $\tilde{\tau}$ leads to slower updates, retaining 950 the previously estimated distribution and making the estimation more stable but less adaptive. As 951 presented in Tab.6, setting $\tilde{\tau}$ to 0.99 yields the best performance, suggesting that a relatively stable 952 estimation of the logits distributions is beneficial for the adaptation process.

Number of EM Loop T. The number of EM loops T determines the number of iterations used to update the GMM parameters in each training step. Tab.7 shows that the model is not sensitive to the choice of T, since the convergence rate of GMM is faster than the rate of network update, and it can be estimated well even when T = 1. We choose T = 3 for stable performance while considering computational efficiency.

Cumulative Density Estimation Loss Weight λ . The weight λ balances the cumulative density estimation loss with the segmentation loss. A higher λ enforces stronger domain alignment through the cumulative density functions. As shown in Tab.8, $\lambda = 0.2$ provides the best performance gain. An overly large λ may distract the model from learning the primary segmentation task, leading to performance degradation.

963 964

935 936

937 938

965 966

F EXTENDED EXPERIMENT ON IMAGE CLASSIFICATION

To demonstrate the generality of BLDA, we implement BLDA based on MIC (Hoyer et al., 2023)
with ResNet-101 on the VisDA-2017 (Peng et al., 2017) UDA classification benchmark in Tab.9,
and our method still achieves improvements. In the classification task, the dataset does not have
severe class distribution differences like segmentation. However, as we point out in Fig.1, the transfer
difficulty differences between domains still lead to severe class bias in this task, and our method can
effectively alleviate this and achieve more balanced predictions.

Table 9: Image classification accuracy in % on VisDA-2017 for UDA, where the improvement is
 marked as **bold**.

| Method | Plane | Bcycl | Bus | Car | Horse | Knife | Mcyle | Persn | Plant | Sktb | Train | Truck | Mean |
|-----------------------------|---------|-------|------|------|-------|-------|-------|-------|-------|------|-------|-------|------|
| MCC (Jin et al., 2020) | 88.1 | 80.3 | 80.5 | 71.5 | 90.1 | 93.2 | 85.0 | 71.6 | 89.4 | 73.8 | 85.0 | 36.9 | 78.8 |
| SDAT (Rangwani et al., 2022 | 2) 95.8 | 85.5 | 76.9 | 69.0 | 93.5 | 97.4 | 88.5 | 78.2 | 93.1 | 91.6 | 86.3 | 55.3 | 84.3 |
| MIC (Hoyer et al., 2023) | 96.7 | 88.5 | 84.2 | 74.3 | 96.0 | 96.3 | 90.2 | 81.2 | 94.3 | 95.4 | 88.9 | 56.6 | 86.9 |
| +BLDA | 96.2 | 90.3 | 82.8 | 81.2 | 95.7 | 96.7 | 93.4 | 86.5 | 95.7 | 94.3 | 91.0 | 65.7 | 89.1 |

G VISUALIZATION OF ESTIMATED GMMS

In this section, we visualize the learned GMMs for target domain, i.e., M_{cl}^T for all $c, l \in [C]$. Fig. 7 presents the Estimated GMMs built with DAFormer, and Fig. 8 presents the estimated GMMs with introducing BLDA. We find that the estimated GMMs can accurately model logits distribution and our method reduces the difference in logits distribution across classes, thus achieving balanced learning.





¹⁰⁸⁰ H MORE VISUALIZATION RESULTS OF LOGITS DISTRIBUTION

In this section, we provide more visualization results to compare logits distribution built with our method. As shown in Fig.9, for over-predicted classes, the network predicts larger positive logits and negative logits (column 1, 3), while for under-predicted classes, the network predicts smaller logits (column 2, 4). This difference in logits distribution leads to the class prediction bias. Our method reduces this difference through aligning with the anchor distribution and achieves class-balanced



1128Figure 9: Comparison of Logits Distribution. We choose {building (2), vegetation (8), car (13), sky1129(10)} as over-predicted classes, and {train (16), bicycle (18), traffic light (6), traffic sign (7) } as1130under-predicted classes for visualization. Note that the anchor distribution is counted separately at1131baseline and our method.

I MORE QUALITATIVE RESULTS

1135 1136

In this section, we provide more qualitative results between our method and other competitors on GTA \rightarrow CS. As shown in Fig.10, when previous methods fail to recognize the classes that are under-predicted and suffer severe performance decline in UDA (e.g. *sidewalk*, *pole*, *fence*, *terrain*, *bike*,*sign*), BLDA shows significant improvement on them, thereby demonstrating the effectiveness of our method.

1141



Figure 10: More qualitative comparison with DAFormer and MIC. The yellow boxes mark regions improved by BLDA.

1177 1178

1179

1180

1181 J LIMITATION AND SOCIETY IMPACT

1182

1183 Our method analyzes the class bias in domain adaptive semantic segmentation through logits dis-1184 tribution statistics and propose a method to implement online logits adjustment tailored for the 1185 UDA training process, which can be easily built with exiting methods and demonstrate consistent 1186 improvements. Although BLDA achieves remarkable performance, we balance the class under the 1187 assumption that each logits distribution in \mathcal{M} is independent, without considering correlation between 1188 classes. How to model this correlation and mitigate the class bias further is still to be resolved. Within this paper, we present an approach for domain adaptive semantic segmentation, a pivotal research area in the realm of computer vision, with no apparent negative societal implications known thus far.

Κ **EXTENDED EXPERIMENTS**

Table 10: UDA segmentation performance on SYN. \rightarrow CS. using the mAcc (%) evaluation metric, where the improvement is marked as **bold**.

| Method | Arch. | Road | Sidewalk | Building | Wall | Fence | Pole | Light | Sign | Veg | Terrain | Sky | Person | Rider | Car | Truck | Bus | Train | Motor | Bike | mAcc | std |
|-------------------------------|-------|------|----------|----------|------|-------|------|-------|------|------|---------|------|--------|-------|------|-------|------|-------|-------|------|------|------|
| DAFormer(Hoyer et al., 2022a) | Т | 89.8 | 90.2 | 96.2 | 33.8 | 8.3 | 51.9 | 63.1 | 57.8 | 95.1 | - | 98.4 | 86.3 | 63.6 | 96.7 | - | 83.4 | - | 55.4 | 61.4 | 70.7 | 25.1 |
| +BLDA | Т | 87.3 | 95.6 | 94.9 | 38.8 | 14.1 | 57.5 | 70.1 | 63.3 | 97.8 | - | 98.9 | 90.0 | 68.0 | 97.7 | - | 95.7 | - | 63.8 | 67.5 | 75.1 | 23.6 |
| HRDA (Hoyer et al., 2022b) | T | 90.3 | 85.2 | 96.0 | 69.4 | 8.0 | 70.6 | 81.2 | 69.6 | 94.7 | - | 99.1 | 88.0 | 68.9 | 97.4 | - | 93.8 | - | 71.3 | 74.0 | 78.6 | 21.2 |
| +BLDA | Т | 89.4 | 94.5 | 95.4 | 75.5 | 24.5 | 78.2 | 85.8 | 74.7 | 97.6 | - | 98.8 | 88.9 | 71.6 | 97.4 | - | 96.9 | - | 72.6 | 76.2 | 82.4 | 17.9 |
| MIC (Hoyer et al., 2023) | T | 89.9 | 87.4 | 96.2 | 71.3 | 8.4 | 66.7 | 81.5 | 69.7 | 95.6 | - | 98.5 | 89.1 | 73.0 | 97.3 | - | 93.5 | - | 78.2 | 71.5 | 79.2 | 21.2 |
| +BLDA | Т | 89.4 | 97.2 | 96.0 | 73.4 | 17.9 | 72.1 | 87.6 | 75.8 | 97.4 | - | 98.3 | 91.5 | 72.5 | 97.7 | - | 96.2 | - | 79.3 | 75.6 | 82.4 | 19.4 |
| | | | | | | | | | | | | | | | | | | | | | | |

Table 11: UDA segmentation performance on GTA. \rightarrow CS. using the mIoU (%) evaluation metric, where the improvement is marked as **bold**. The results are acquired based on CNN-based model (He et al., 2016; Chen et al., 2017a), denoted as C, and Transformer-based model (Xie et al., 2021), denoted as T. * denotes the reproduced result.

| } | | | ad | lewalk | ilding | IJ | nce | e | ght | E | 50 | rrain | ý | rson | ler | - | ıck | s | uin | otor | ée | | |
|---|---------------------------------|-------|------|--------|--------|------|------|------|------|------|------|-------|------|------|------|------|------|------|-----------------|------|------|------|------|
|) | Method | Arch. | 8 | Sic | Bu | W | Fei | Po | Γi | Sig | Ve | Te | Sk | Peı | Ric | Ca | Ę | Bu | Tr ₂ | M | Bil | mIoU | std |
| | DACS* (Tranheden et al., 2021) | С | 93.0 | 52.0 | 87.8 | 29.4 | 38.3 | 37.7 | 45.0 | 53.3 | 87.9 | 46.3 | 90.2 | 67.8 | 38.0 | 89.0 | 51.1 | 51.1 | 0.0 | 10.7 | 19.4 | 52.1 | 27.1 |
| | +BLDA | С | 92.9 | 67.5 | 87.1 | 36.3 | 39.3 | 41.2 | 50.7 | 58.5 | 87.3 | 45.5 | 87.7 | 69.1 | 40.4 | 88.3 | 45.3 | 53.5 | 1.2 | 10.8 | 36.3 | 54.7 | 25.6 |
| | DAFormer* (Hoyer et al., 2022a) | С | 95.7 | 69.9 | 87.2 | 35.6 | 36.7 | 37.0 | 49.4 | 52.8 | 87.3 | 44.1 | 87.9 | 69.0 | 42.2 | 86.5 | 40.0 | 51.7 | 0.2 | 41.1 | 54.0 | 56.2 | 24.0 |
| | +BLDA | C | 95.6 | 73.6 | 86.7 | 41.0 | 40.2 | 43.3 | 51.1 | 62.4 | 86.2 | 43.7 | 87.4 | 68.3 | 39.2 | 86.6 | 43.6 | 45.6 | 1.0 | 49.4 | 59.4 | 58.1 | 23.2 |
| | CDAC* (Wang et al., 2023) | Т | 96.1 | 72.8 | 90.5 | 55.2 | 48.0 | 51.8 | 57.1 | 61.8 | 90.8 | 50.4 | 91.9 | 73.2 | 46.9 | 93.6 | 80.9 | 78.6 | 58.2 | 56.9 | 64.5 | 69.2 | 16.7 |
| | +BLDA | Т | 96.6 | 78.1 | 90.0 | 57.9 | 52.5 | 55.1 | 58.7 | 64.5 | 90.1 | 50.8 | 90.9 | 73.3 | 47.5 | 93.2 | 75.1 | 80.0 | 65.3 | 60.7 | 68.9 | 71.0 | 15.4 |

Table 12: UDA segmentation performance on GTA. \rightarrow CS. using the mAcc(%) evaluation metric, where the improvement is marked as **bold**. The results are acquired based on CNN-based model (He et al., 2016; Chen et al., 2017a), denoted as C, and Transformer-based model (Xie et al., 2021), denoted as T. * denotes the reproduced result.

| Method | Arch. | Road | Sidewalk | Building | Wall | Fence | Pole | Light | Sign | Veg | Terrain | Sky | Person | Rider | Car | Truck | Bus | Train | Motor | Bike | mAcc | std |
|---------------------------------|-------|------|----------|----------|------|-------|------|-------|------|------|---------|------|--------|-------|------|-------|------|-------|-------|------|------|------|
| DACS* (Tranheden et al., 2021) | С | 97.9 | 58.9 | 94.9 | 40.7 | 49.4 | 46.1 | 51.7 | 57.2 | 94.5 | 66.9 | 98.6 | 82.0 | 62.2 | 93.3 | 82.9 | 79.0 | 0.0 | 74.4 | 20.0 | 65.8 | 26.5 |
| +BLDA | C | 97.1 | 80.5 | 93.4 | 54.1 | 39.2 | 50.5 | 64.6 | 67.4 | 93.7 | 71.5 | 99.0 | 83.4 | 58.8 | 95.0 | 73.0 | 76.2 | 1.0 | 73.5 | 39.2 | 69.0 | 24.2 |
| DAFormer* (Hoyer et al., 2022a) | С | 98.4 | 78.3 | 93.6 | 43.2 | 45.4 | 45.0 | 61.5 | 59.5 | 93.8 | 65.3 | 98.0 | 80.1 | 69.2 | 92.2 | 77.0 | 85.3 | 0.2 | 69.7 | 60.8 | 69.3 | 23.8 |
| +BLDA | С | 97.0 | 84.8 | 93.0 | 55.1 | 56.1 | 55.4 | 71.0 | 74.5 | 93.2 | 73.7 | 97.7 | 84.8 | 77.2 | 92.0 | 80.1 | 85.5 | 1.1 | 77.3 | 73.0 | 74.9 | 21.6 |
| CDAC* (Wang et al., 2023) | Т | 99.1 | 78.0 | 94.6 | 69.3 | 52.9 | 61.2 | 71.7 | 68.7 | 95.8 | 59.3 | 99.0 | 85.0 | 70.2 | 96.0 | 88.3 | 85.5 | 71.1 | 75.4 | 74.7 | 78.7 | 13.8 |
| +BLDA | Т | 98.1 | 82.5 | 94.0 | 74.8 | 66.7 | 68.0 | 76.8 | 74.1 | 95.5 | 67.9 | 98.3 | 87.1 | 72.9 | 95.6 | 90.8 | 88.8 | 71.8 | 80.8 | 82.9 | 82.5 | 11.5 |

Table 13: Performance comparison on GTA. \rightarrow CS. grouped by over-prediction and under-prediction classes. The results are acquired based on CNN-based model (He et al., 2016; Chen et al., 2017a), denoted as C, and Transformer-based model (Xie et al., 2021), denoted as T.

| 1230 | Metrics | | IoU: n | nean/std | Acc: n | nean/std |
|------|--------------------------------|-------|-----------------|------------------|-----------------|------------------|
| 1231 | Methods | Arch. | over-prediction | under-prediction | over-prediction | under-prediction |
| 1232 | DACS (Tranheden et al., 2021) | С | 68.7/28.9 | 37.0/13.1 | 80.3/29.2 | 52.3/14.5 |
| 1233 | +BLDA | C | 68.0/28.6 | 42.7/14.3 | 79.1/29.0 | 59.9/13.4 |
| 100/ | DAFormer (Hoyer et al., 2022a) | С | 67.3/29.6 | 46.3/10.0 | 79.8/29.0 | 59.8/11.2 |
| 1234 | +BLDA | C | 66.8/29.3 | 50.3/10.9 | 80.5/28.6 | 69.8/10.0 |
| 1235 | CDAC (Wang et al., 2023) | Т | 83.8/11.6 | 56.5/7.6 | 90.4/8.5 | 68.1/7.6 |
| 1236 | +BLDA | Т | 83.8/10.1 | 59.5/8.7 | 91.1/7.8 | 74.7/5.7 |
| 1237 | DAFormer (Hoyer et al., 2022a) | Т | 83.3/10.3 | 54.7/7.3 | 90.6/8.0 | 66.1/6.2 |
| 1238 | +BLDA | Т | 84.2/8.5 | 57.8/9.4 | 92.3/4.7 | 72.6/8.0 |
| 1020 | HRDA (Hoyer et al., 2022b) | Т | 87.8/6.8 | 61.0/8.0 | 93.0/5.6 | 72.5/10.1 |
| 1239 | +BLDA | Т | 88.6/6.2 | 64.5/7.2 | 93.9/4.0 | 77.0/8.2 |
| 1240 | MIC (Hoyer et al., 2023) | Т | 89.5/5.9 | 63.6/8.0 | 92.7/6.6 | 74.7/8.0 |
| 1241 | +BLDA | Т | 89.6/5.3 | 66.4/8.0 | 94.4/3.6 | 79.2/7.9 |

| Methods | GPU Memory (MB) | Time per iter (s) | Total Time (h) |
|--------------------------------|-----------------|-------------------|-------------------|
| DAFormer (Hoyer et al., 2022a) | 9,807 | 1.32 | 14.5 (40K iters) |
| +BLDA | 12,655 | 1.59 | 17.7 (40K iters) |
| DACS (Tranheden et al., 2021) | 11,078 | 0.52 | 35.5 (250K iters) |
| +BLDA | 14,354 | 0.81 | 56.1 (250K iters) |
| CDAC (Wang et al., 2023) | 35,443 | 1.66 | 18.7 (40K iters) |
| +BLDA | 35,443 | 1.95 | 22.3 (40K iters) |

Table 14: Computational resource requirements comparison

Table 15: Comparison with other class-imbalanced methods on GTA. \rightarrow CS. using the mIoU (%) evaluation metric, where baseline is based on DAFormer with uniform class sampling.

| Method | Arch. | Road | Sidewalk | Building | Wall | Fence | Pole | Light | Sign | Veg | Terrain | Sky | Person | Rider | Car | Truck | Bus | Train | Motor | Bike | mIoU | std |
|---------------------------|-------|------|----------|----------|------|-------|------|-------|------|------|---------|------|--------|-------|------|-------|------|-------|-------|------|------|------|
| baseline | T | 95.2 | 67.4 | 89.0 | 43.7 | 47.5 | 46.7 | 55.7 | 55.0 | 89.3 | 47.2 | 90.7 | 71.4 | 44.5 | 92.3 | 76.1 | 78.5 | 62.1 | 55.4 | 64.1 | 66.9 | 17.6 |
| [1] (Lin et al., 2017) | Т | 87.6 | 40.1 | 87.9 | 49.2 | 41.2 | 45.7 | 52.1 | 50.9 | 89.2 | 45.0 | 90.7 | 70.8 | 41.1 | 89.9 | 56.6 | 70.5 | 60.3 | 54.6 | 59.5 | 62.3 | 18.1 |
| [2] (Cui et al., 2019) | Т | 96.0 | 70.6 | 89.6 | 53.8 | 48.3 | 50.2 | 53.3 | 60.8 | 89.7 | 49.4 | 90.9 | 70.2 | 42.6 | 91.7 | 65.7 | 77.8 | 60.1 | 54.4 | 65.2 | 67.4 | 16.8 |
| [3] (Hoyer et al., 2022a) | T | 95.7 | 70.2 | 89.4 | 53.5 | 48.1 | 49.6 | 55.8 | 59.4 | 89.9 | 47.9 | 92.5 | 72.2 | 44.7 | 92.3 | 74.5 | 78.2 | 65.1 | 55.9 | 61.8 | 68.3 | 16.8 |
| [4] (Menon et al., 2020) | Т | 94.7 | 66.9 | 89.2 | 53.2 | 47.0 | 49.0 | 51.2 | 56.0 | 89.6 | 50.1 | 87.2 | 69.0 | 41.9 | 92.0 | 71.8 | 75.7 | 62.3 | 50.4 | 64.0 | 66.4 | 16.9 |
| +BLDA | Т | 94.5 | 72.1 | 88.4 | 48.3 | 51.6 | 52.9 | 59.5 | 64.5 | 88.1 | 51.5 | 90.5 | 72.4 | 47.6 | 94.0 | 75.9 | 80.8 | 65.6 | 53.7 | 67.7 | 69.5 | 15.9 |

Table 16: Comparison with other class-imbalanced methods on GTA. \rightarrow CS. using the mAcc (%) evaluation metric, where baseline is based on DAFormer with uniform class sampling.

| Method | Arch. | Road | Sidewalk | Building | Wall | Fence | Pole | Light | Sign | Veg | Terrain | Sky | Person | Rider | Car | Truck | Bus | Train | Motor | Bike | mAcc | std |
|---------------------------|-------|------|----------|----------|------|-------|------|-------|------|------|---------|------|--------|-------|------|-------|------|-------|-------|------|------|------|
| baseline | Т | 98.3 | 75.4 | 94.9 | 53.8 | 57.8 | 53.6 | 70.2 | 59.7 | 96.2 | 57.6 | 99.1 | 83.5 | 67.6 | 95.5 | 87.1 | 84.3 | 73.8 | 74.5 | 77.5 | 76.9 | 15.3 |
| [1] (Lin et al., 2017) | Т | 91.7 | 63.3 | 94.6 | 66.4 | 50.2 | 54.4 | 63.9 | 55.1 | 95.6 | 66.9 | 98.6 | 84.9 | 58.5 | 92.6 | 84.4 | 81.9 | 70.0 | 68.8 | 68.2 | 74.2 | 15.2 |
| [2] (Cui et al., 2019) | Т | 98.8 | 77.7 | 94.1 | 71.4 | 56.8 | 63.7 | 77.2 | 71.2 | 95.3 | 63.5 | 99.0 | 84.1 | 72.0 | 94.6 | 90.1 | 89.2 | 68.5 | 76.4 | 80.2 | 80.2 | 12.6 |
| [3] (Hoyer et al., 2022a) | Т | 99.1 | 74.8 | 95.2 | 61.0 | 53.1 | 59.8 | 70.7 | 68.6 | 96.0 | 63.4 | 98.7 | 84.2 | 69.9 | 95.5 | 88.9 | 84.1 | 74.0 | 70.0 | 69.7 | 77.8 | 14.2 |
| [4] (Menon et al., 2020) | Т | 97.9 | 79.9 | 95.0 | 64.4 | 64.7 | 66.6 | 74.3 | 72.6 | 96.5 | 66.1 | 99.7 | 87.9 | 72.5 | 96.3 | 90.5 | 91.3 | 85.7 | 74.3 | 72.4 | 81.5 | 12.2 |
| +BLDA | Т | 97.5 | 79.4 | 93.4 | 65.5 | 68.0 | 63.9 | 77.2 | 75.5 | 94.9 | 68.1 | 99.0 | 85.2 | 71.4 | 95.2 | 89.2 | 87.6 | 80.8 | 77.0 | 83.6 | 81.7 | 10.9 |

Table 17: The JS divergence between the P_{cl} and corresponding anchor distribution P_p and P_n , where the results on P_n is averaged over the C-1 negative components.

| Method | Anchor | Road | Sidewalk | Building | Wall | Fence | Pole | Light | Sign | Veg | Terrain | Sky | Person | Rider | Car | Truck | Bus | Train | Motor | Bike | mean | std |
|----------|-------------|-------|----------|----------|-------|-------|-------|-------|-------|-------|---------|-------|--------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| baseline | P_p | 0.412 | 0.289 | 0.381 | 0.286 | 0.293 | 0.137 | 0.325 | 0.142 | 0.429 | 0.273 | 0.596 | 0.130 | 0.217 | 0.300 | 0.225 | 0.164 | 0.190 | 0.145 | 0.180 | 0.270 | 0.119 |
| +BLDA | P_p | 0.367 | 0.246 | 0.261 | 0.153 | 0.212 | 0.221 | 0.146 | 0.234 | 0.281 | 0.173 | 0.310 | 0.244 | 0.070 | 0.273 | 0.284 | 0.181 | 0.125 | 0.207 | 0.125 | 0.217 | 0.072 |
| baseline | $\hat{P_n}$ | 0.517 | 0.327 | 0.315 | 0.394 | 0.305 | 0.280 | 0.422 | 0.288 | 0.412 | 0.404 | 0.433 | 0.278 | 0.396 | 0.279 | 0.387 | 0.382 | 0.401 | 0.430 | 0.394 | 0.371 | 0.064 |
| +BLDA | P_n | 0.231 | 0.165 | 0.153 | 0.135 | 0.209 | 0.118 | 0.224 | 0.119 | 0.248 | 0.246 | 0.240 | 0.186 | 0.158 | 0.170 | 0.199 | 0.208 | 0.176 | 0.162 | 0.136 | 0.183 | 0.041 |

Table 18: Ablation on different selection criteria for anchor distributions on GTA. \rightarrow CS. using the mIoU (%) evaluation metric. OLA denotes online logits adjustment.

| 288 | | | | | | | | | | | | | | | | | | | | | | | |
|------|--------------------------------|----------|------|--------|--------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|
| 1289 | | | ad | lewalk | ilding | Ξ | Jce | e | tht | E | 50 | rain | × | son | ler | L | ıck | s | ij. | otor | e | | |
| 1290 | Anchor | Strategy | R l | Sid | Bu | Wa | Fei | Pol | Ľ | Sig | Ve | Ter | Sk | Peı | Rić | Ca | Ę | Bu | Tr | M | Bil | mIoU | std |
| 1200 | DAFormer (Hoyer et al., 2022a) | - | 95.7 | 70.2 | 89.4 | 53.5 | 48.1 | 49.6 | 55.8 | 59.4 | 89.9 | 47.9 | 92.5 | 72.2 | 44.7 | 92.3 | 74.5 | 78.2 | 65.1 | 55.9 | 61.8 | 68.3 | 16.8 |
| 1291 | target (global) | post-hoc | 95.7 | 77.0 | 87.8 | 61.0 | 53.7 | 54.3 | 56.1 | 61.2 | 87.2 | 50.2 | 90.4 | 74.4 | 43.4 | 90.9 | 73.3 | 82.5 | 56.6 | 54.8 | 68.3 | 69.4 | 15.8 |
| 1000 | source (building) | post-hoc | 95.8 | 77.1 | 90.0 | 59.8 | 54.5 | 51.6 | 56.6 | 59.8 | 86.2 | 48.2 | 90.6 | 75.4 | 43.7 | 90.4 | 72.0 | 79.5 | 58.9 | 53.7 | 69.2 | 69.1 | 16.0 |
| 1292 | source (fence) | post-hoc | 95.6 | 77.5 | 88.0 | 58.4 | 55.0 | 51.2 | 56.9 | 57.7 | 88.7 | 48.3 | 90.5 | 75.2 | 43.8 | 90.1 | 71.7 | 80.6 | 61.1 | 53.9 | 65.5 | 68.9 | 16.0 |
| 1293 | source (global) | post-hoc | 95.7 | 77.1 | 88.6 | 60.5 | 55.3 | 48.5 | 57.3 | 60.9 | 89.5 | 47.2 | 91.0 | 72.9 | 43.7 | 91.3 | 73.7 | 80.8 | 61.1 | 55.3 | 63.8 | 69.2 | 16.2 |
| | source (building) | OLA | 95.5 | 77.9 | 88.5 | 60.4 | 59.1 | 53.3 | 57.9 | 60.4 | 88.3 | 49.2 | 90.3 | 76.1 | 43.6 | 90.3 | 76.4 | 85.7 | 56.2 | 57.9 | 68.8 | 70.3 | 15.8 |
| 1294 | source (fence) | OLA | 95.6 | 81.2 | 88.2 | 57.9 | 57.0 | 51.6 | 54.2 | 60.2 | 87.9 | 50.2 | 90.3 | 76.3 | 44.3 | 90.0 | 75.0 | 82.4 | 57.7 | 56.0 | 70.6 | 69.8 | 16.0 |
| 1295 | source (global) | OLA | 95.4 | 78.3 | 88.3 | 54.0 | 55.2 | 55.7 | 60.3 | 65.2 | 89.2 | 47.3 | 91.1 | 71.4 | 44.8 | 91.6 | 74.3 | 83.4 | 73.2 | 59.3 | 67.1 | 70.7 | 15.5 |

| 299 300 | Anchor | Strategy | Road | Sidewalk | Building | Wall | Fence | Pole | Light | Sign | Veg | Terrain | Sky | Person | Rider | Car | Truck | Bus | Train | Motor | Bike | mAcc | std |
|------------|--------------------------------|----------|------|----------|----------|------|-------|------|-------|------|------|---------|------|--------|-------|------|-------|------|-------|-------|------|------|------|
| 1301 | DAFormer (Hoyer et al., 2022a) | - | 99.1 | 74.8 | 95.2 | 61.0 | 53.1 | 59.8 | 70.7 | 68.6 | 96.0 | 63.4 | 98.7 | 84.2 | 69.9 | 95.5 | 88.9 | 84.1 | 74.0 | 70.0 | 69.7 | 77.8 | 14.2 |
| 1001 | target (global) | post-hoc | 98.7 | 85.6 | 92.2 | 65.9 | 67.4 | 67.3 | 75.1 | 68.6 | 90.8 | 66.8 | 95.5 | 91.8 | 67.3 | 93.1 | 92.7 | 95.7 | 62.0 | 84.5 | 80.2 | 81.1 | 12.4 |
| 1302 | source (building) | post-hoc | 98.7 | 89.6 | 92.2 | 73.5 | 67.5 | 66.3 | 76.5 | 67.0 | 90.8 | 63.8 | 95.5 | 90.7 | 68.3 | 93.1 | 92.6 | 95.6 | 60.3 | 79.9 | 80.9 | 81.2 | 12.5 |
| | source (fence) | post-hoc | 98.5 | 85.3 | 91.3 | 74.1 | 66.7 | 66.2 | 82.3 | 68.4 | 90.5 | 63.4 | 95.5 | 88.9 | 62.1 | 92.5 | 91.9 | 94.9 | 63.8 | 83.5 | 79.2 | 81.0 | 12.2 |
| 1303 | source (ours) | post-hoc | 98.6 | 90.0 | 91.6 | 71.0 | 66.5 | 64.1 | 78.5 | 68.2 | 91.2 | 65.1 | 96.4 | 91.9 | 64.8 | 92.7 | 92.5 | 95.1 | 67.6 | 82.2 | 77.4 | 81.3 | 12.4 |
| 130/ | source (building) | OLA | 98.6 | 84.7 | 92.6 | 71.4 | 66.2 | 68.8 | 76.6 | 71.0 | 91.2 | 67.3 | 95.8 | 94.6 | 69.7 | 93.5 | 93.2 | 96.3 | 59.7 | 79.5 | 80.2 | 81.6 | 12.2 |
| 1304 | source (fence) | OLA | 98.6 | 84.8 | 92.2 | 71.6 | 64.1 | 66.9 | 78.3 | 68.2 | 90.9 | 68.2 | 95.6 | 92.9 | 65.4 | 93.3 | 93.0 | 95.9 | 60.4 | 82.3 | 84.8 | 81.4 | 12.6 |
| 1305 | source (ours) | OLA | 98.3 | 86.7 | 93.4 | 70.9 | 61.9 | 66.0 | 74.0 | 78.9 | 95.4 | 59.0 | 98.6 | 85.7 | 73.1 | 95.3 | 89.8 | 89.3 | 85.4 | 77.2 | 78.6 | 82.0 | 11.9 |

1296Table 19: Ablation on different selection criteria for anchor distributions on GTA. \rightarrow CS. using the1297mAcc (%) evaluation metric. OLA denotes online logits adjustment.

L EXTENDED DISCUSSION

1310 L.1 DOMAIN ADAPTIVE SEMANTIC SEGMENTATION

Unsupervised domain adaptation (UDA) transfers semantic knowledge learned from labeled source
domains to unlabeled target domains. Due to the ubiquity of domain gaps, UDA methods have been
widely studied in various computer vision tasks, such as image classification, object detection, and
semantic segmentation. UDA is crucial for semantic segmentation to avoid laborious pixel-wise
annotation in new target scenarios.

Recent UDA approaches for semantic segmentation fall into two main paradigms: adversarial training-1317 based methods (Toldo et al., 2020; Tsai et al., 2018; Chen et al., 2018; Ganin & Lempitsky, 2015; 1318 Hong et al., 2018; Long et al., 2015b) and self-training-based methods (Tranheden et al., 2021; Hoyer 1319 et al., 2022a; Araslanov & Roth, 2021; Zhang et al., 2021). Adversarial training-based methods learn 1320 domain-invariant representations through a min-max optimization game, where a feature extractor 1321 is trained to confuse a domain discriminator, aligning feature distributions across domains. Self-1322 training-based methods, which have become dominant in the field due to the domain-robustness of 1323 Transformers (Bhojanapalli et al., 2021), generate pseudo labels for target images based on a teacher-1324 student optimization framework. The success of this paradigm depends on generating high-quality 1325 pseudo labels, with strategies such as entropy minimization (Chen et al., 2019) and consistency 1326 regularization (Hoyer et al., 2023) being developed for this purpose.

1327 Recently, several approaches have been proposed to tackle the challenges in UDA for semantic 1328 segmentation: DTS (Huo et al., 2023) employs a Dual Teacher-Student Framework, promoting 1329 the self-training paradigm to fully adapt models to the target domain by exploring different mix 1330 strategies. CDAC (Wang et al., 2023) introduces consistency constraints in attention to enhance 1331 the model's cross-domain performance. RTea (Zhao et al., 2023) defines proxy tasks based on 1332 structural information and incorporates them into the self-training paradigm as additional supervision 1333 signals. Peng et al. (2023) apply Diffusion-based image translation techniques to directly mitigate the distribution differences between the target and source domains. DiGA (Shen et al., 2023) integrates 1334 distillation strategies and self-training through multi-stage training. 1335

However, due to the inherent class imbalance and distribution shift in both data and label space
between domains, networks often exhibit complex class biases, which are further amplified by the
confirmation bias inherent in the self-training paradigm. Our method aims to achieve balanced
learning in UDA training to mitigate these issues and improve the overall performance of domain
adaptation for semantic segmentation.

1341

1298

1306 1307

1308 1309

1311

1342 L.2 CLASS-IMBALANCED LEARNING

Class imbalance is a prevalent issue in semantic segmentation, where the number of samples per class varies significantly. Existing methods tackle this problem through re-weighting or re-sampling techniques. Re-weighting methods assign different weights to classes during training, giving higher importance to under-represented classes (Cui et al., 2019; Lin et al., 2017; Cao et al., 2019; Liu et al., 2019). Re-sampling techniques modify the class distribution in the training data by over-sampling minority classes or under-sampling majority classes (He et al., 2008; 2021; He & Garcia, 2009; Kim et al., 2020; Chu et al., 2020).

Recent works have proposed various approaches to address class imbalance in different tasks. For object detection, Lin et al. (2017) propose a re-weighting method that assigns weights to samples based on their confidence, with harder samples receiving higher weights. In classification, Cui et al. (2019) define an effective number of samples for each class to weight different classes, while Menon et al. (2020) introduce a logit adjustment method that incorporates correction terms in logits, determined by prior knowledge of class distribution. For segmentation, Hoyer et al. (2022a) present a re-sampling technique that samples based on prior knowledge of class distribution, with more sampling for rare classes. Truong et al. (2023) considers class-imbalance as fairness problem and propose to model the context of structural dependency to tackle it.

In UDA for semantic segmentation, some approaches have introduced these strategies to alleviate
class bias (Hoyer et al., 2022a; Araslanov & Roth, 2021; Li et al., 2022). However, these methods are
still empirical and focus on the single-domain setting, which assumes that the test data and training
data share the same distribution in both data space and label space, without considering the additional
challenges posed by domain shift in UDA.

In contrast to these methods, our approach aims to directly address class bias and achieve balanced learning for each class without relying on prior knowledge about the distribution shift between domains. We evaluate class bias through logit sets in the form of a confusion matrix and explicitly balance components in the matrix using anchor distributions and cumulative density functions, implemented in an online self-training paradigm. Our method does not depend on prior knowledge of class distribution and instead directly models class bias based on the actual logit distribution, making it more adaptable to the challenges posed by domain shift in UDA.

1372 M CORE CODE

| 1374 | We provide our core code in BLDA to demonstrate our implementation. |
|------|---|
| 1375 | |
| 1376 | |
| 1377 | |
| 1378 | |

- 138?