

Local Context-Aware Active Domain Adaptation

Tao Sun
Stony Brook University
tao@cs.stonybrook.edu

Cheng Lu
XPeng Motors
luc@xiaopeng.com

Haibin Ling
Stony Brook University
hling@cs.stonybrook.edu

Abstract

Active Domain Adaptation (ADA) queries the labels of a small number of selected target samples to help adapting a model from a source domain to a target domain. The local context of queried data is important, especially when the domain gap is large. However, this has not been fully explored by existing ADA works. In this paper, we propose a Local context-aware ADA framework, named LADA, to address this issue. To select informative target samples, we devise a novel criterion based on the local inconsistency of model predictions. Since the labeling budget is usually small, fine-tuning model on only queried data can be inefficient. We progressively augment labeled target data with the confident neighbors in a class-balanced manner. Experiments validate that the proposed criterion chooses more informative target samples than existing active selection strategies. Furthermore, our full method clearly surpasses recent ADA arts on various benchmarks. Code is available at <https://github.com/tsun/LADA>.

1. Introduction

Unsupervised Domain Adaptation (UDA) [6, 18] adapts a model from a related source domain to an unlabeled target domain. It has been widely studied in the past decade. Despite its success in many applications, UDA is still a challenging task, especially when the domain gap is large [35]. In practical scenarios, it is often allowable to annotate a small number of unlabeled data. The new paradigm of Active Domain Adaptation (ADA), which queries the label of selected target samples to assist domain adaptation, draws increasing attention recently due to its promising performance with minimal labeling cost [32, 22, 42, 41].

Traditional Active Learning (AL) methods select unlabeled samples that are *uncertain* to the model [10] or *representative* to the data distribution [30]. Some works combine these two principles to design *hybrid* criteria [10]. These AL strategies, however, may be less effective in ADA due to the availability of labeled source data and the distribution shift between source and target domains. Recent ADA

works seek to use density weighted entropy [32], combine transferable criteria [5], focus on hard examples [42] or exploit free energy biases [41] to select target samples that are beneficial to domain adaptation.

Despite recent progress in ADA, the local context of queried data has not been fully explored. Different from Semi-Supervised Domain Adaptation (SSDA) [29, 15] where all labeled target data are given at the beginning of training and fixed afterwards, active querying to obtain labeled target data and model update interleave during the training of ADA. The local context can guide the selection of target samples that are uncertain and locally representative. It can also be utilized to update models and reduces the tendency to only memorize these newly queried data during fine-tuning [46]. Thus later training rounds can focus on harder cases. In the particular situation when the domain gap is large, it is also safer to trust the neighbors of queried data than other confident but distant samples.

Thus motivated, in this paper, we propose a novel framework of Local context-aware Active Domain Adaptation (LADA). A Local context-aware Active Selection (LAS) module is first designed, with a novel criterion based on the local inconsistency of model predictions. During the active selection stage, a diverse subset from the uncertain regions measured by our criterion is selected for querying labels. Then we design a Progressive Anchor-set Augmentation (PAA) module to overcome issues from the small size of queried data. Since the labeling budget for each round is usually small, it only requires a small fraction of training epoch before the model predicts well on the newly queried data. Another issue is that the labeled data can be imbalanced when the target data are long-tailed or the active selection focuses on partial classes. Our PAA handles the above-mentioned issues through augmenting labeled target data. Specifically, during each training epoch, we initialize an anchor set with all available queried data and progressively augment it with pseudo-labeled confident target data in a class-balanced manner. Target training batches are sampled from the anchor set instead to fine-tune the model. We show that choosing confident samples in the neighborhood of queried data overcomes the adversarial effects from false

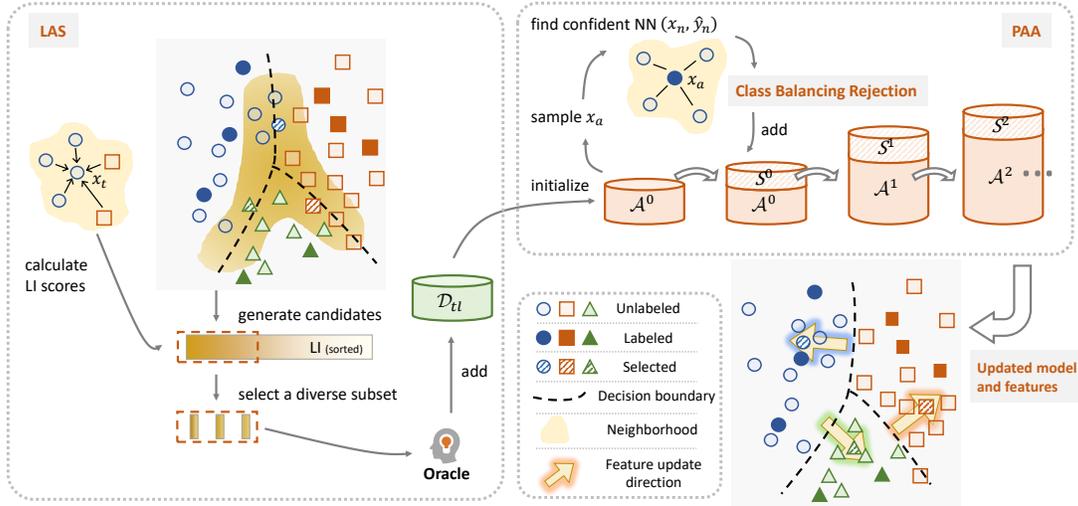


Figure 1: The framework of LADA. For active querying, the LAS (Local context-aware Active Selection) module selects informative target samples based on the Local Inconsistency (LI) of model predictions. For model adaptation, the PAA (Progressive Anchor set Augmentation) module exploits all queried data and their confident neighbors to fine-tune the model. The anchor set \mathcal{A} is expanded progressively during each training epoch and in a class balanced manner.

confident samples that are distant to labeled data.

To demonstrate the effectiveness of exploiting local context in ADA, we conduct extensive experiments on various domain adaptation benchmarks. We implement several representative active selection strategies, and compare them with LAS under the same configurations. Either with the simple model fine-tuning or a strong SSDA method named MME [29], LAS consistently selects more informative samples, leading to higher accuracies. Equipped with PAA, the full LADA method outperforms state-of-the-art ADA solutions on various benchmarks on both standard datasets and datasets with class distribution shift.

In summary, we make the following contributions:

- We advocate to utilize the local context of queried data in ADA, which may guide active selection and improve model adaptation.
- We propose a LAS module with a novel active criterion based on the local inconsistency of class probability predictions. It selects more informative samples than existing active selection criteria.
- We design a PAA module to overcome issues from the small size of queried data. It progressively supplements labeled target data with confident samples in a class-balanced manner.
- Extensive experiments show that the full LADA method outperforms state-of-the-art ADA solutions.

2. Related Work

Domain Adaptation. Domain Adaptation aims to adapt a model from a source domain to a related but different target domain. Early works seek to align the feature distribu-

tions through matching statistics [38] or adopt adversarial learning [6, 34]. Recent studies show their intrinsic limitations when label distribution shift (LDS) exists [14, 23]. Aligning feature distributions could also fail to learn discriminative class boundaries in SSDA where a few labeled target data are available [29, 44]. Self-training [18, 35] that learns from model predictions has become a promising approach in UDA with LDS [14], SSDA [44], source-data free UDA [17, 33], *etc.* Exploiting the structure of target data to improve adaptation performance is also an important research topic [37, 45, 15]. Since UDA can be challenging when the domain gap is large, some recent works try to explore Transformer [4] architecture for stronger feature representation [35], and leverage prior knowledge of target class distribution to rectify model predictions [33].

Active Learning. Active Learning (AL) selects a small number of unlabeled samples for annotation. There are two major paradigms: *uncertainty*-based methods select uncertain samples with criteria like entropy, confidence [13], prediction margin [2], *etc.*; *representativeness*-based methods select samples that can represent the entire data distribution. Popular methods include clustering [43] and CoreSet [30]. While most traditional AL algorithms use a one-by-one query method, deep AL [26] adopts batch-based query to avoid extra computation and model overfitting [1, 11]. Recently, the ideas of Bayesian learning [12], adversarial learning [31], and reinforcement learning [19] have also been introduced into deep AL.

Active Domain Adaptation. Active Domain Adaptation queries the label of selected target samples to assist domain

adaptation. It was first studied in the task of sentiment classification from text data [24], where samples are selected with uncertainty and a domain separator. More recently, AADA [32] studies ADA in vision tasks. It combines active learning with adversarial domain adaptation and selects samples with a diversity cue from importance weight plus an uncertainty cue. CLUE [22] performs uncertainty-weighted clustering to select samples that are both uncertain and diverse. S³VAADA [25] uses a submodular criterion combining three scores to select a subset. TQS [5] adopts an ensemble of transferable committee, transferable uncertainty and transferable domainness as the selection criterion. EADA [41] trains an energy-based model, and selects samples based on energy values. SDM [42] focuses on hard cases. It consists of a maximum margin loss and a margin sampling function. LAMDA [8] addresses the issue of label distribution mismatch in ADA by seeking target data that best approximate the entire target distribution.

Although various ADA methods have been proposed, the local context of queried data has not been fully explored. The closest to ours is [40] that proposes a Minimum Happy (MH) points learning for Active Source Free Domain Adaptation. The MH points are selected to have low neighbor purity and high neighbor affinity, where purity is measured with the entropy of hard neighbor labels. Different from theirs, we use the local inconsistency of class probability predictions to keep more information. Furthermore, they adopt one-shot querying in order to determine source-like samples as source data are unavailable, while we follow previous ADA works to conduct several rounds of active selection. We show that our active learning strategy is more compact and effective under ADA setting.

3. Local Context-Aware ADA

Problem Formulation. In ADA, there is a source domain with labeled data $\mathcal{D}_s = \{(x_i^s, y_i^s)\}$ and a target domain with unlabeled data $\mathcal{D}_t = \{(x_i^t)\}$. Meanwhile, we can actively select a few target samples to query their labels under a labeling budget B , which is usually much smaller than the total amount of target data (*i.e.*, $B \ll |\mathcal{D}_t|$). Let the obtained labeled target data (*a.k.a.* queried data) be $\mathcal{D}_{tl} = \{(x_i^{tl}, y_i^{tl})\}$, with $|\mathcal{D}_{tl}| = B$. The remaining unlabeled target data are $\mathcal{D}_{tu} = \mathcal{D}_t \setminus \mathcal{D}_{tl}$. The goal of ADA is to learn a model $h = g \circ f$ that generalizes well on the target domain, through querying most informative samples. f is a feature extractor and g is a task predictor. This work focuses on C -way classification. Hence labels $y \in \mathcal{Y} = \{0, 1, \dots, C - 1\}$ are categorical variables. In practice, it is more efficient to conduct active querying in R rounds, with a label budget of $b = B/R$ for each round.

It is worth mentioning that some ADA methods [5, 22, 42] train the model using only labeled data $\mathcal{D}_s \cup \mathcal{D}_{tl}$, while others [25, 8, 40] also require unlabeled data. When the la-

Algorithm 1 Local context-aware ADA.

Input: Source data \mathcal{D}_s , target data \mathcal{D}_t , training epochs E , iterations per epoch I

Initialization: $\mathcal{D}_{tl} = \emptyset, \mathcal{D}_{tu} = \mathcal{D}_t$

```

1: for  $e = 0$  to  $E$  do
2:   if need active querying then
3:     Obtain query set  $\mathcal{D}_q$  with ground-truth labels
       from Oracle as described in Sec. 3.1
4:      $\mathcal{D}_{tl} \leftarrow \mathcal{D}_{tl} \cup \mathcal{D}_q, \mathcal{D}_{tu} \leftarrow \mathcal{D}_{tu} \setminus \mathcal{D}_q$ 
5:   end if
6:   Initialize  $\mathcal{A} = \mathcal{D}_{tl}, \mathcal{S} = \emptyset$ 
7:   for  $i = 0$  to  $I$  do
8:     Sample mini-batches  $\mathcal{B}_s$  from  $\mathcal{D}_s, \mathcal{B}_a$  from  $\mathcal{A}$ 
9:     Update model parameters with Eq. 5
10:    Update  $\mathcal{S}$  as described in Alg. 2.
11:    if iteration over  $\mathcal{A}$  finishes then
12:       $\mathcal{A} \leftarrow \mathcal{A} \cup \mathcal{S}, \mathcal{S} = \emptyset$ 
13:    end if
14:  end for
15: end for

```

belonging budget is small or there exists label distribution shift, the unlabeled data are critical to the model performance. For our work, we present results under both settings.

Method Overview. We exploit the local context of queried data with the aim to select more informative target samples and improve model adaptation. The framework of our Local context-aware Active Domain Adaptation (LADA) is illustrated in Fig. 1. The LAS (Local context-aware Active Selection) module sorts all unlabeled target samples by their Local Inconsistency of model predictions (LI scores), and then selects a diverse subset with the largest scores for querying their ground-truth labels. The PAA (Progressive Anchor set Augmentation) module progressively supplements query set with confident target samples. A Class Balancing strategy is adopted to add samples from different classes evenly. The two modules run alternatively for several rounds until the labeling budget is used up. Alg. 1 summarizes the training procedure.

3.1. Local Context-aware Active Selection

Uncertainty and *Representativeness* are two general principles in traditional active learning. Uncertainty-based strategies aim to select samples whose model predictions are less confident, while representativeness-based strategies aim to select samples to better represent the entire data distribution. In ADA, the source domain provides auxiliary labeling information, making it cost-inefficient to sample in the well-aligned regions of the two domains. Our empirical experiments indicate that it is more beneficial to focus on the uncertain target samples in ADA tasks.

Many criteria (*e.g.*, *entropy* and *margin*) have been pro-

Algorithm 2 Progressive Anchor set Augmentation.

Input: Confidence threshold τ , reject probability $q_{\mathcal{A}}$, supplementary set \mathcal{S} , labeled mini-batch \mathcal{B}_a , target data \mathcal{D}_t , neighborhood size K

- 1: **for** $x_a \in \mathcal{B}_a$ **do**
 - 2: Randomly sample x_n from the K nearest neighbors of x_a for LAA (or from \mathcal{D}_t for RAA)
 - 3: $p_n = \sigma(h(x_n))$, $\hat{y}_n = \arg \max p_n$
 - 4: $\xi \sim U(0, 1)$ $\triangleright U$ is uniform distribution
 - 5: **if** $p_n[\hat{y}_n] \geq \tau$ and $\xi \geq q_{\mathcal{A}}[\hat{y}_n]$ **then**
 - 6: $\mathcal{S} \leftarrow \mathcal{S} \cup \{(x_n, \hat{y}_n)\}$
 - 7: **end if**
 - 8: **end for**
 - 9: **Return** \mathcal{S}
-

posed to select uncertain samples. However, these measures only rely on the sample-wise model prediction while ignoring its local context, which may include some outliers during active selection. To address this issue, we design a novel criterion based on the local inconsistency of model predictions. Let $p(\cdot) = \sigma(h(\cdot))$ denote the class probability after softmax operation $\sigma(\cdot)$. For an unlabeled target sample x_t , its Local Inconsistency (LI) is measured as

$$\text{LI}(x_t) = -\frac{1}{K} \sum_{k=1}^K w_k p(x_{tk})^\top p(x_t) \quad (1)$$

where $\{x_{tk}\}$ are K nearest neighbors of x_t based on the cosine similarity $\langle f(x_{tk}), f(x_t) \rangle$, and $w_k \propto f(x_{tk})^\top f(x_t)$ is the normalized weight. Here we use the soft probability p , as pseudo labels are unreliable for uncertain samples near decision boundaries. When the neighbors are very close to the center sample, $\text{LI}(x_t) \approx -p(x_t)^\top p(x_t) = p(x_t)^\top (1 - p(x_t)) - 1$, which is essentially the Gini Impurity. In real ADA benchmarks, the data distribution is not dense and the neighborhood size K can be adjusted accordingly. To remove outlier values and enhance the clustering structure of candidate uncertain regions, we additionally smooth the scores among neighborhood by

$$\text{LI}(x_t) \leftarrow \text{LI}(x_t) + \frac{1}{K} \sum_{k=1}^K w_k \text{LI}(x_{tk}) \quad (2)$$

Data with high LI scores generally have large sample-wise uncertainty and inconsistent predictions to neighbors. An adjoint effect is that those highly scored target samples tend to co-occur in specific local regions. Active selection solely based on LI scores would include highly similar sample pairs, leading to a waste of labeling budgets. Fortunately, we can select a diverse subset from an oversampled candidate set. Let b be the labeling budget for current round and M be a predefined ratio. A clustering process is run on the top $(1 + M)b$ samples with highest LI scores, and then the b cluster centroids are selected as the

query set. The diverse sampling can also be implemented with determinant point process, though clustering is simple and works good enough. Note that unlike CLUE [22] that runs an uncertainty-weighted clustering on all unlabeled target data, ours focuses only on a small portion of uncertain samples and runs more efficient.

3.2. Progressive Anchor Set Augmentation

After each round of active querying, the model can be updated for 1-2 epochs. A common way is to fine-tune the model using all labeled data [5, 42, 41] with the objective¹

$$\mathcal{L} = \mathbb{E}_{(x,y) \sim \mathcal{D}_s} \ell_{ce}(h(x), y) + \mathbb{E}_{(x,y) \sim \mathcal{D}_{tl}} \ell_{ce}(h(x), y) \quad (3)$$

where ℓ_{ce} is the cross entropy loss.

The number of queried data for each round is usually small, *e.g.*, only 1% target data are added to \mathcal{D}_{tl} after active selection when the total labeling budget is 5%. This causes two issues: it only requires a small fraction of training epoch before the model predicts well on those queried samples; \mathcal{D}_{tl} can be class imbalanced when the target data are imbalanced or the active selection focuses on a few classes. Consequently, it fails to fully utilize the training resources and to provide more informative target samples to the following rounds of active selection.

To supplement the limited queried data, we propose a Progressive Anchor set Augmentation (PAA) module to incorporate confident target data. An anchor set \mathcal{A} is initialized with current \mathcal{D}_{tl} . At each training iteration, a target mini-batch \mathcal{B}_a is sampled from \mathcal{A} instead of \mathcal{D}_{tl} for supervised training. Then some selected confident samples in the neighborhood of \mathcal{B}_a are added to a temporal set \mathcal{S} . Once the sampling iteration over \mathcal{A} concludes, we reinitialize the anchor set as $\mathcal{A} \leftarrow \mathcal{A} \cup \mathcal{S}$ and clear \mathcal{S} . The same procedure repeats until the end of this epoch. We name this as Local context-aware Anchor set Augmentation (LAA) as it exploits the local region of queried data. We also consider Random Anchor set Augmentation (RAA), where confident samples are randomly selected from the entire \mathcal{D}_t .

To obtain a class balanced anchor set, we reject a confident target sample with a higher probability if there already exist many data from the same class in current \mathcal{A} . Specifically, let $p_{\mathcal{A}}[c] = \sum_{(x_a, y_a) \in \mathcal{A}} \mathbb{I}[y_a = c] / |\mathcal{A}|$, $p_{\mathcal{A}}^{\max} = \max_c p_{\mathcal{A}}[c]$ and $p_{\mathcal{A}}^{\min} = \min_c p_{\mathcal{A}}[c]$. For a confident target sample with pseudo label c , we reject it with a probability

$$q_{\mathcal{A}}[c] = \frac{p_{\mathcal{A}}[c] - p_{\mathcal{A}}^{\min}}{p_{\mathcal{A}}^{\max}} \quad (4)$$

The detailed procedure of progressive anchor set augmentation is described in Alg. 2.

¹Some works sample training data from a joint labeled set $\mathcal{D}_l = \mathcal{D}_s \cup \mathcal{D}_{tl}$, which is more efficient but less generalizable.

Table 1: Accuracies (%) on **Office-Home** using 5%-budget. (The highest accuracies across all methods are **bolded** and within each section are underlined.)

AL method	DA method	Ar→Cl	Ar→Pr	Ar→Rw	Cl→Ar	Cl→Pr	Cl→Rw	Pr→Ar	Pr→Cl	Pr→Rw	Rw→Ar	Rw→Cl	Rw→Pr	Avg.
RAN		58.9	77.6	78.7	61.9	74.2	73.0	62.9	56.0	77.9	70.1	60.1	83.4	69.6
ENT		61.3	81.1	<u>82.4</u>	64.9	78.6	77.0	64.0	58.8	81.7	73.5	61.9	87.1	72.7
MAR		62.2	81.5	<u>82.4</u>	64.4	79.5	77.2	64.2	60.5	81.7	73.7	64.2	87.5	73.3
CoreSet		58.2	77.7	79.2	61.7	73.6	73.3	60.6	55.1	78.9	70.3	60.0	82.8	69.3
BADGE	ft w/ CE loss	63.4	81.9	81.5	65.1	79.9	77.0	64.4	61.0	82.0	73.8	63.5	88.0	73.5
AADA		60.4	81.5	82.3	64.5	79.0	76.9	63.2	58.9	81.9	73.3	63.2	87.4	72.7
CLUE		63.0	81.7	81.1	63.2	79.3	76.2	64.6	59.7	81.5	73.1	63.5	86.8	72.8
TQS		58.6	81.1	81.5	61.1	76.1	73.3	61.2	54.7	79.7	73.4	58.9	86.1	70.5
MHPL		63.5	81.7	82.1	65.0	79.2	77.2	65.0	61.7	82.2	73.8	65.1	87.7	73.7
LAS		<u>66.1</u>	<u>83.8</u>	<u>82.4</u>	<u>66.9</u>	<u>82.5</u>	<u>78.6</u>	<u>66.8</u>	<u>63.1</u>	<u>82.4</u>	<u>74.9</u>	<u>66.7</u>	<u>89.6</u>	<u>75.3</u>
RAN		64.0	82.3	81.2	68.4	80.8	77.3	68.2	63.5	82.0	76.0	65.1	86.2	74.6
AADA		64.8	84.2	84.7	71.8	82.3	80.4	71.7	62.7	85.1	79.2	67.2	88.6	76.9
CLUE	MME	65.5	84.9	83.8	71.1	82.4	79.5	71.4	62.9	85.2	79.0	66.9	88.4	76.7
MHPL		66.8	82.3	84.0	71.1	84.2	80.6	71.9	65.5	84.5	78.1	67.6	89.3	77.2
LAS		68.6	86.7	85.0	72.1	84.6	80.6	71.4	65.6	85.5	79.4	68.4	89.9	78.2

Table 2: Accuracies (%) on **Office-Home** using 10%-budget. (†Using importance sampling)

AL method	DA method	Ar→Cl	Ar→Pr	Ar→Rw	Cl→Ar	Cl→Pr	Cl→Rw	Pr→Ar	Pr→Cl	Pr→Rw	Rw→Ar	Rw→Cl	Rw→Pr	Avg.
BADGE		71.3	88.6	86.3	<u>70.6</u>	86.7	82.0	71.6	69.7	<u>86.7</u>	<u>79.8</u>	72.4	91.9	79.8
AADA		69.0	87.5	<u>87.4</u>	70.1	85.3	81.8	70.3	67.2	<u>86.7</u>	79.1	69.6	90.6	78.7
CLUE	ft w/ CE loss	69.7	87.6	85.7	69.8	86.3	81.0	69.8	68.4	85.5	78.1	71.8	91.1	78.7
TQS		68.0	87.7	85.7	67.0	83.0	78.7	69.3	64.5	83.9	77.8	68.9	90.6	77.1
LAS		<u>73.6</u>	<u>90.0</u>	87.0	<u>70.6</u>	<u>88.7</u>	<u>82.6</u>	<u>72.4</u>	<u>71.8</u>	86.3	79.3	<u>73.8</u>	<u>92.2</u>	<u>80.7</u>
BADGE		68.8	86.4	84.4	75.3	88.2	83.6	75.2	71.3	87.1	79.6	73.5	91.3	80.4
AADA	CDAC	69.7	87.0	86.2	74.8	87.0	84.5	74.8	69.7	88.1	79.7	71.2	90.6	80.3
CLUE		72.3	87.6	85.8	74.0	88.6	84.3	74.4	72.6	87.2	79.4	73.3	91.1	80.9
LAS		<u>73.6</u>	<u>89.3</u>	<u>86.8</u>	<u>76.3</u>	<u>89.6</u>	<u>85.9</u>	<u>76.8</u>	<u>74.9</u>	<u>88.5</u>	<u>81.1</u>	<u>76.5</u>	<u>92.3</u>	<u>82.6</u>
TQS		68.7	80.1	83.1	64.0	83.1	76.9	67.7	71.0	84.4	76.4	72.7	90.0	76.5
S ³ VAADA	DANN†	65.5	79.6	80.0	65.4	82.2	75.5	68.4	68.1	84.0	73.5	70.7	88.6	75.1
LAMDA		74.8	88.5	86.9	73.8	88.2	83.3	74.6	75.5	86.9	80.8	77.8	91.7	81.9
LAS	MCC	77.2	91.0	88.6	77.1	90.7	86.8	76.4	76.4	89.1	81.9	77.7	93.3	83.9
LAS	RAA	77.8	91.8	88.4	77.7	91.5	87.7	78.1	79.1	89.5	83.4	79.8	94.1	84.9
LAS	LAA	77.2	91.9	88.1	76.9	91.1	86.8	76.6	78.1	88.3	82.0	79.0	93.8	84.2

The anchor set \mathcal{A} contains not only all labeled target data but also some pseudo-labeled confident data. With the class balancing rejection, \mathcal{A} is expected to contain adequate number of target samples from each class, thus effectively overcome the drawbacks from the small size of queried data. Later rounds of active selection may focus on more challenging samples. To further exploit \mathcal{A} , we incorporate a strong image transformation with RandAugment [3]. The objective becomes

$$\mathcal{L} = \mathbb{E}_{(x,y) \sim \mathcal{D}_s} \ell_{ce}(h(x), y) + \mathbb{E}_{(x,y) \sim \mathcal{A}} \ell_{ce}(h(\tilde{x}), y) \quad (5)$$

where $\tilde{x} = \beta x + (1-\beta)\alpha(x)$ for some random augmentation operator $\alpha(\cdot)$, and $\beta \sim \text{Beta}(0.2, 0.2)$ is a random variable. It can be viewed as a specific kind of *Mixup* between a weak augmentation and a strong augmentation of the same image. Comparing Eq. 5 to Eq. 3, we see that the memory usage and computation cost are not increased.

It is worth mentioning that \mathcal{A} is reinitialized as \mathcal{D}_{tl} at each training epoch. This reduces cumulative biases from noisy pseudo labels. Another reason is that after each round of active selection, restarting from \mathcal{D}_{tl} ensures to train the model sufficiently on the newly queried target samples.

4. Experiments

We conduct experiments on four widely-used domain adaptation benchmarks: **Office-31** [28], **Office-Home** [39], **VisDA** [21] and **DomainNet**. **Office-Home RSUT** [36] is a subset of Office-Home created with the protocol of Reverse-unbalanced Source and Unbalanced Target to have a large label distribution shift.

All experiments are implemented with Pytorch. In consistent with previous works [5, 41, 42], we use a pre-trained ResNet-50 [7] backbone and perform 5 rounds of active selection with $B = 5\%$ or $B = 10\%$ of all target data. For a fair comparison, we reproduce the results of several traditional active learning criteria including random (RAN), least confidence (CONF), entropy (ENT), prediction margin (MAR) [10], CoreSet [30], BADGE [1], and ADA methods including AADA [32] and CLUE [22] in a unified framework. Results of recent state-of-the-art ADA methods like S³VAADA [25], TQS [5] and LAMDA [8] are borrowed from the corresponding papers whenever applicable.

We set the confidence threshold τ to 0.9 for all datasets. The neighborhood size K is 5 for Office-31 and 10 for

Table 3: Accuracies (%) on **Office-31** using 5%-budget.

AL	DA	A→D	A→W	D→A	D→W	W→A	W→D	Avg.
RAN		82.5	87.6	73.4	98.6	76.2	99.6	86.3
ENT		90.8	92.5	75.8	99.9	76.9	100.	89.3
MAR		89.4	92.5	78.5	99.8	79.0	99.9	89.9
CoreSet		83.8	85.9	74.8	97.6	76.1	99.6	86.3
BADGE	fit w/	88.0	90.9	79.2	99.8	80.5	99.9	89.7
AADA	CE loss	88.7	91.4	76.3	100.	77.3	100.	88.9
CLUE		89.8	91.8	79.1	100.	79.8	99.9	90.1
TQS		<u>92.8</u>	92.2	80.6	100.	80.4	100.	91.1
MHPL		90.0	91.0	78.3	99.4	78.7	99.9	89.5
LAS		91.6	<u>93.9</u>	<u>81.5</u>	99.7	<u>81.8</u>	99.6	<u>91.4</u>
<hr/>								
RAN		89.9	92.6	78.1	99.2	78.9	99.8	89.7
AADA		95.9	94.1	81.4	<u>99.5</u>	81.3	100.	92.0
CLUE	MME	96.1	<u>96.3</u>	82.1	99.3	82.0	100.	92.6
MHPL		95.2	95.6	82.2	99.3	82.1	100.	92.4
LAS		<u>96.7</u>	96.1	<u>84.8</u>	99.3	<u>84.8</u>	100.	<u>93.6</u>
<hr/>								
RAN		90.6	91.7	75.0	98.2	77.0	99.7	88.7
AADA		94.7	96.1	77.8	99.5	79.6	99.5	91.2
CLUE	CDAC	94.9	95.3	78.8	99.5	79.5	100.	91.4
MHPL		94.2	95.6	76.8	99.5	79.2	99.6	90.8
LAS		<u>96.2</u>	<u>96.6</u>	<u>80.8</u>	<u>99.6</u>	<u>82.2</u>	99.8	<u>92.5</u>
<hr/>								
	DANN	95.0	96.4	83.3	99.0	83.7	99.8	92.9
	MCC	94.2	95.5	85.1	100.	86.0	99.8	93.4
LAS	RAA	96.9	97.6	84.2	100.	86.0	100.	94.1
	LAA	97.8	98.5	82.8	100.	85.2	100.	94.0

other datasets. We set M with an empirical formulation $\lceil \frac{55}{100B} - 1 \rceil$, leading to a candidate set of about $\sim 12\%$ unlabeled target data. Since VisDA has a huge number of data in each class, we reduce the empirical value by a half. Additional implementation details and analyses can be found in the supplementary.

4.1. Main Results

Comparison with ADA methods on standard datasets. Results on Office-Home using 5% budget is listed in Tab. 1. Among the active selection criteria, uncertainty-based criteria (e.g., ENT and MAR) generally obtain higher accuracies than representativeness-based criteria (e.g., CoreSet), indicating that it is inefficient to select target data that are well-aligned with the source domain. When the semi-supervised solver MME [29] is used, the performance gaps among these criteria become smaller. With either strategy, our proposed LAS consistently obtains the best scores, showing that it can select more informative target samples. Table 2 lists the results using 10%-budget. When using fine-tuning or CDAC [15], LAS obtains better accuracies than other active selection criteria. Compared with the recent LAMDA method, our LAS w/ LAA as a unified solution boosts the accuracy by +2.3%. LAS w/ RAA is slightly better.

Table. 3 shows results on Office-31 using 5%-budget. When fixing the adaptation method as fine-tuning, MME or CDAC, LAS consistently performs better than other active selection criteria. When fixing the active selection method as LAS, the proposed RAA/LAA performs better than MME, MCC [9], and CDAC. Table 5 lists results on VisDA using 10%-budget. LAS w/ LAA outperforms LAMDA by +1.3%.

Table 4: Accuracies (%) on **Office-Home RSUT** using 10%-budget. (†Using importance sampling)

AL	DA	C→P	C→R	P→C	P→R	R→C	R→P	Avg.
S ³ VAADA	VAADA	73.0	63.0	50.7	69.6	52.6	78.3	64.5
TQS	DANN†	67.6	61.4	54.8	74.7	53.6	77.6	64.9
LAMDA	DANN†	81.2	75.7	64.1	81.6	65.1	87.2	75.8
<hr/>								
	DANN	80.2	69.4	58.4	76.9	60.8	85.3	71.8
	MME	75.6	68.9	56.6	76.9	57.9	87.2	70.5
	MCC	82.3	72.0	62.6	81.2	65.9	88.3	75.4
LAS	CDAC	79.6	71.9	62.1	81.4	63.7	85.6	74.1
	RAA	83.8	73.6	64.0	82.6	65.2	88.6	76.3
	LAA	83.2	77.2	63.8	83.0	65.4	88.1	76.8

Table 5: Accuracies (%) on **VisDA** and **DomainNet** using 10%-budget. (†Using importance sampling)

AL	DA	VisDA	DomainNet				Avg.
			R→C	C→S	S→P	C→Q	
TQS	DANN†	87.7	59.3	50.9	52.4	41.5	51.0
LAMDA	DANN†	91.8	65.3	56.1	58.1	48.3	57.0
<hr/>							
	DANN	91.3	62.1	53.1	53.1	42.4	52.7
	MME	92.2	65.9	54.1	55.9	42.9	54.7
LAS	RAA	93.0	70.3	58.1	60.3	48.9	59.4
	LAA	93.1	69.4	57.5	59.9	49.1	59.0

Comparison with ADA methods on label-shifted datasets. Since label distribution mismatch between source and target domains raises a critical issue in ADA, we compare with the LAMDA [8] devised to address this issue on label-shifted datasets. Following their settings, we use 10%-budget. Tables 4,5 list the comparison results. LAMDA applies importance sampling on source data to match the label distribution of source and target domains. This is particularly useful to domain adversarial methods like DANN. Our RAA/LAA belong to self-training methods. We use Class Balancing Rejection to create class-balanced training data. Given LAS, RAA/LAA achieve higher accuracies than other adaptation methods. The full LAS w/ LAA method improves over LAMDA by +1%.

4.2. Analysis on LAS

Uncertainty measures in LAS. LAS discovers uncertain regions with the LI score. To analyze its properties, we replace the LI score in LAS with other uncertainty measures while keeping the same diverse sampling procedure. Comparison measures include *prediction margin*, *entropy*, *prediction confidence* and *NAU score*[40]. Fig. 4a and Fig. 4b plots the performances under different oversampling rate M . From the figure, using a large M is beneficial to all uncertainty measures, indicating the necessity to balance between uncertainty and diversity. CLUE and BADGE are also hybrid methods, but use a different way to LAS. CLUE runs a clustering on all target samples based on entropy weighted distances, and BADGE runs a clustering on all target samples based on the gradient embeddings. Both of them rely on the entire target data. In contrast, LAS fo-

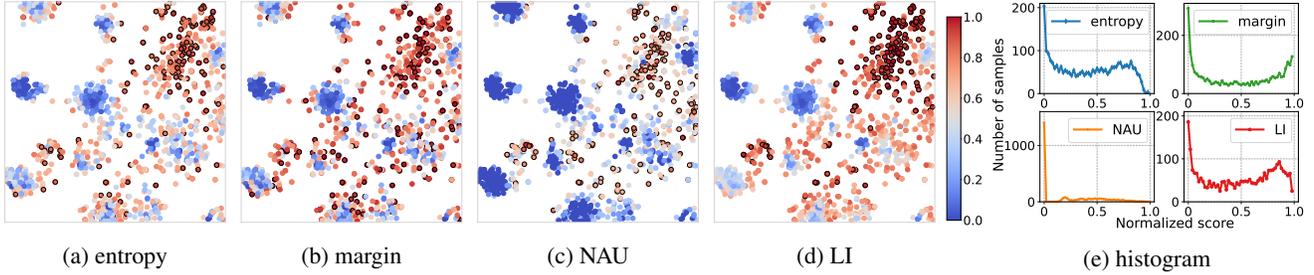


Figure 2: (a-d) t -SNE visualization of target features on Office-31 W→A. Samples are colored according to their normalized uncertainty scores, where red indicates large values and blue indicates small values. The top 10% samples with highest scores are marked with black boarders. (e) Histogram of target samples by normalized scores.

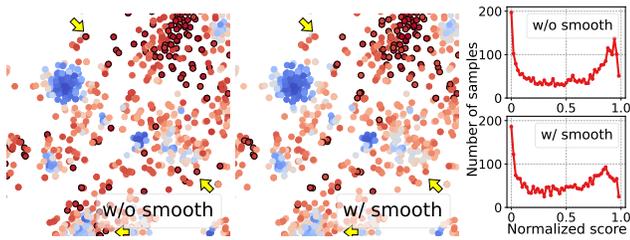


Figure 3: Zoomed in t -SNE visualization of target features and histograms on Office-31 W→A for LI. The smoothing in Eq. 2 removes some outliers (e.g., samples annotated with yellow arrows).

focuses on some local uncertain regions, and involves only about 12% target data in our experiments. When $M \geq 10$, LAS and its variants surpass CLUE and BADGE on the two tasks. In Fig. 4c when $M = 10$, LI leads to the best scores on Office-Home and Office-31. Fig. 4d studies the parameter sensitivity on K and M in LAS. Performances are close in the proximal of optimal values.

Advantages of LAS over other criteria. To better show the advantage of LAS over other active learning criteria, Fig. 5 plots the accuracies under different labeling budgets and domain adaptation methods. In the left figure, LAS consistently achieves better or comparable accuracies than other criteria with both fine-tuning and MME. The improvement is more significant when the labeling budget is small. The center figure plots the accuracy curves using 5%-budget. LAS outperforms other criteria through training process. In the right figure, no matter which DA strategy is used, LAS obtains the best scores.

To further analyze the difference between these uncertainty measures, Fig. 2 visualizes the target features on Office-31 W→A. We normalize all uncertainty scores to $[0, 1]$ and color each point accordingly. Meanwhile, the top 10% samples with highest scores are marked with black boarders. Entropy and margin do not consider the local-context. Thus they tend to include some outliers as shown in Fig. 2a and Fig. 2b. NAU is defined as $NAU = NP \times NA$, where NP is the entropy of class distribution among neigh-

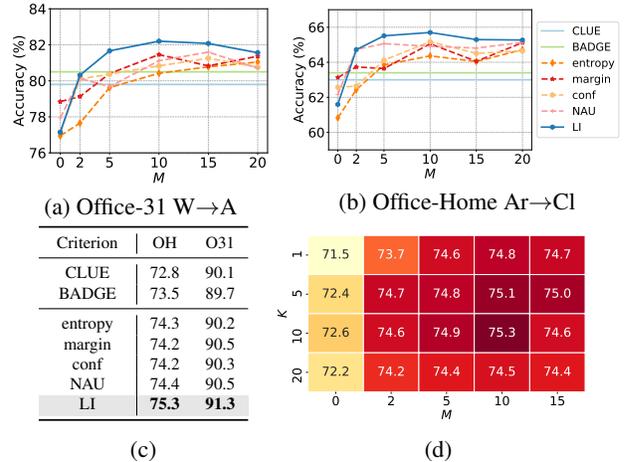


Figure 4: (a,b) Replacing LI score in LAS with different uncertainty criteria; (c) comparison results when $M = 10$; (d) varying K and M in LAS on Office-Home.

boring samples and NA is the average similarity between a sample and its neighbors. Since the number of neighbors is limited, NP has discrete values and tend to be small. Shown in Fig. 2c and Fig. 2e, the majority target data have small NAU scores. In contrast, uncertainty samples tend to have large LI scores and are more gathered. A peak around 0.8 can be observed in the histogram of LI. This also explains why a diverse sampling is important for LI. In the zoomed visualization of Fig. 3, we can see that the smoothing in LI helps as it can remove some outliers.

4.3. Analysis on RAA/LAA

Effects of Class Balancing Rejection (CBR). To handle issues from label distribution shift, an effective way is to create a class balanced training set. We realize this through using different rejection probabilities for target samples from major or minor classes when creating the anchor set \mathcal{A} . Figure 6 visualizes the ratio of samples per class among \mathcal{A} . As can be seen, \mathcal{A} is dominated by major classes without using CBR. In contrast, the ratio of samples from minor classes increases when CBR is used. This helps to train each class

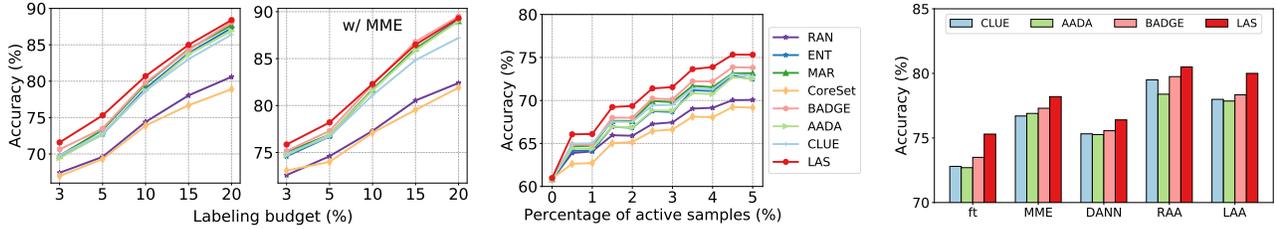


Figure 5: Analysis on Office-Home. (Left) varying labeling budget; (Center) accuracy curves using 5%-budget; (Right) combining AL criteria with different DA strategies using 5%-budget.

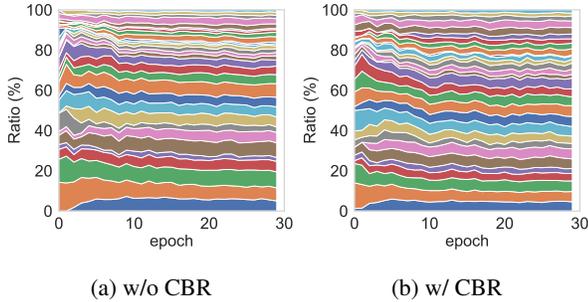


Figure 6: Effects of Class Balancing Rejection (CBR) in LAA on the ratio of samples per class among \mathcal{A} . (Averaged over six tasks on Office-Home RSUT. For better visualization, only 15 major and 15 minor classes are included.)

in a more balanced manner, as evident from +1.0% increase in per-class average accuracy in Tab. 6

When local context-aware augmentation matters? Previous results validate the effectiveness of progressive anchor set augmentation. In terms of the two strategies, random (RAA) and local context-aware (LAA), there is a trade-off in exploitation and exploration. RAA can include target samples that are distant to the queried data. However, confident samples may be misclassified due to domain distribution shift. To illustrate the situation when local context matters, we manually increase domain gap by adding noises to the latent features. Given a pretrained ResNet, we first obtain the class prototypes $\{p_c\}$ of source data. Then for each class c , we generate a random displacement $\xi_c \triangleq \sum_i r_i (p_i - p_c)$, where $r_i \sim U(-u, u)$, and apply it to the entire source data of Class- c before making predictions. Target data are unmodified. It effectively increases the domain gap, while preserving the source domain semantic structure. Fig. 7 plots the comparison results. The benefit of local context-aware augmentation is more significant when u is large where source and target domains have a large gap. The performance of RAA drops much faster due to the inclusion of false confident target samples.

Ablation studies on components of LADA. Table 6 presents ablation studies on each component of LADA. Improvements in the second row over the first row indicate that it is critical to select a diverse subset in LAS. With anchor set augmentation, the accuracies are increased by +2.6% on

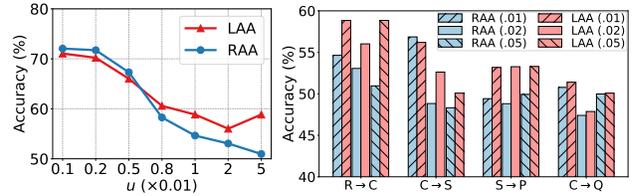


Figure 7: Increasing domain gap by adding a random displacement to source data. (left) DomainNet R→C; (right) four DomainNet tasks with different u .

Table 6: Ablation studies on Office-Home using 5%-budget and Office-Home RSUT using 10%-budget. Div-Sel: select a diverse subset in LAS; \mathcal{A} -Aug: augment \mathcal{A} with confident neighbors; CBR: class balancing rejection. (\dagger Per-class average accuracy.)

Div-Sel	\mathcal{A} -Aug	RandAug	CBR	OH	OH-RSUT	
				72.6	70.6	70.3 \dagger
✓				75.3	73.2	72.5 \dagger
✓	✓			77.9	75.5	74.2 \dagger
✓	✓		✓	79.5	75.6	74.5 \dagger
✓	✓	✓		79.6	76.2	74.3 \dagger
✓	✓	✓	✓	80.0	76.8	75.3 \dagger

Office-Home and +2.3% on Office-Home RSUT. This verifies that fine-tuning model on only queried data is ineffective. Using RandAug to create mixed images lead to further gains. To show the effectiveness of Class Balancing rejection, we report per-class average accuracies on Office-Home RSUT, where it is boosted by +1.0% in the last row.

5. Conclusion

In this paper, we advocate to utilize the local context of queried data for active domain adaptation. We propose a local context-aware active selection method based on the local inconsistency of model predictions. It consistently selects more informative samples than previous criteria. Then we propose a progressive anchor set augmentation module to mitigate issues from small labeling budgets. It utilizes queried data and their expanded neighbors to refine the model. Extensive experiments validate that our full method, named LADA (Local context-aware Active Domain Adaptation), surpasses state-of-the-art ADA solutions.

A. Dataset Details

Office-31 [28] contains 31 classes of 4,110 office environment related images. It has three domains: Amazon (A), DSLR (D) and Webcam (W). **Office-Home** is a similar dataset, containing 15,500 office images from 65 classes, split in four domains: Product (Pr), Clip Art (Cl), Artistic (Ar) and Real-World (Rw). **Office-Home RSUT** [36] is a subset of Office-Home created with the protocol of Reverse-unbalanced Source and Unbalanced Target to have a large label distribution shift. The major classes in the ‘RS’ fold become minor classes in the ‘UT’ fold, while the minor classes in the ‘RS’ fold become major classes in the ‘UT’ fold. **VisDA** [21] is a large-scale Synthetic-to-Real dataset of 12 objects. The training set contains 152,397 synthetic 2D renderings of 3D models and the validation set contains 55,388 real images. We use the training set as the source domain and the validation set as the target domain. **DomainNet** [20] consists of about 0.6 million images from 345 classes, distributed in six domains. Following [22, 8], we use five domains: Real (R), Clipart (C), Painting (P), Sketch (S), and Quickdraw (Q) for experiments.

B. Implementation Details

We implement all experiments with PyTorch 1.8. Results are run on servers with NVIDIA A5000/A6000 GPU. Following previous ADA works [5, 42, 41], we use ResNet-50 [7] pretrained on ImageNet [27] as the backbone network, a bottleneck layer (`Linear->BatchNorm1d`), and a classification head of one single `Linear` layer. The bottleneck feature dimension is 256. Training images are first resized to 256×256 , and then randomly cropped to 224×224 . Test images use center cropping instead. We adopt Adadelta optimizer with learning rate of 0.1 and a batch size of 32. On Office-Home and Office-31, we first train the models on only source data for 10 epochs, and then train on both source and target data with active domain adaptation for 30 epochs. At the epoch of 10, 12, 14, 16, 18, $B/5$ target data are selected for querying labels, where B is the labeling budget. On VisDA, we conduct source-only training for 1 epoch and ADA for 10 epochs. On DomainNet, we conduct source-only training for 10 epochs and ADA for another 10 epochs. Mean accuracies of 3 repeated experiments are reported.

C. Additional Results and Analyses

Running time. Table A.1 reports running time in seconds with one A6000 GPU, including active sampling (AL) time averaged over 5 rounds (10%-budget) and model update (DA) time averaged over all training epochs. LAS consumes much less time than CLUE and slightly more than other AL methods. RAA/LAA is comparable to MCC and faster than CDAC.

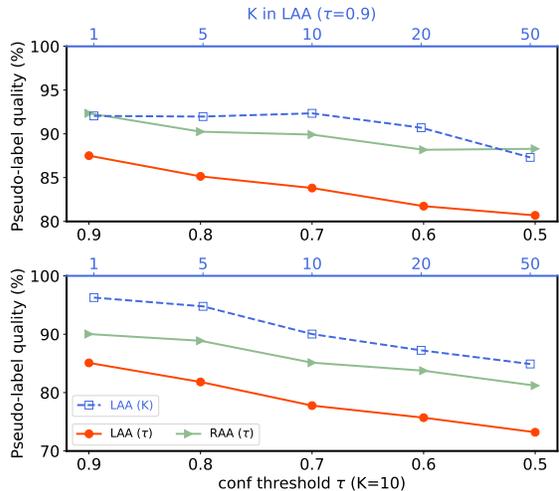


Figure A.1: Pseudo-label quality averaged over Office-Home $Cl \rightarrow \{Ar, Pr, Rw\}$ (upper) and Office-Home RSUT $C \rightarrow \{P, R\}$ (lower) using 10%-budget.

Pseudo-label quality of LAA/RAA. We conduct experiments with different confidence thresholds τ and report the percentage of target samples in the anchor set with correct pseudo-labels. Shown in Fig. A.1, as τ reduces, the pseudo-label quality decreases. LAA has better pseudo-label quality than RAA. We also report results by fixing $\tau = 0.9$ and varying neighborhood size K in LAA (dashed lines). Overall, $\tau = 0.9$ used in the paper leads to a decent pseudo-label quality.

Comparison with other criteria on Office-31 Figure A.2 presents analyses on Office-31 similar to that on Office-Home in the Fig. 2 of the main paper. In the left figure, our LAS outperforms other active learning criteria for labeling budgets ranging from 3% to 20%. When the labeling budget is small (*e.g.*, 3% or 5%), LAS boosts the accuracy by a large margin. Since in ADA the situation with a small labeling budget is more important, it shows the effectiveness of LAS. In the center figure with 5%-budget, the curve of LAS lies above others after 1% samples are selected. In the right figure, when combined with five different domain adaptation strategies, LAS consistently achieves the highest accuracies than three other active learning criteria that are previous arts.

Remarks on SSDA. Semi-supervised DA (SSDA) is closely related to Active DA. In both task, a few labeled target data and many unlabeled target data are available. Yet there are some differences. In SSDA, all labeled target data are provided for once at the beginning of training and fixed afterwards. While in ADA, labeled target data are actively selected. The active querying and model update interleave for several rounds during the training of ADA.

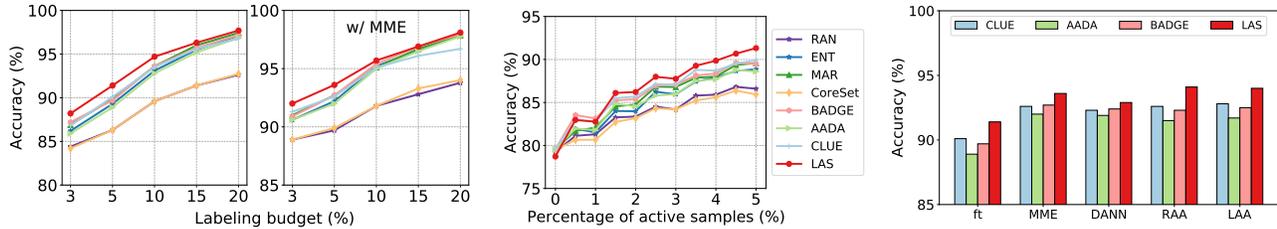


Figure A.2: Analysis on Office-31. (Left) varying labeling budget; (Center) accuracy curves with 5%-budget; (Right) combining AL criteria with different DA strategies.

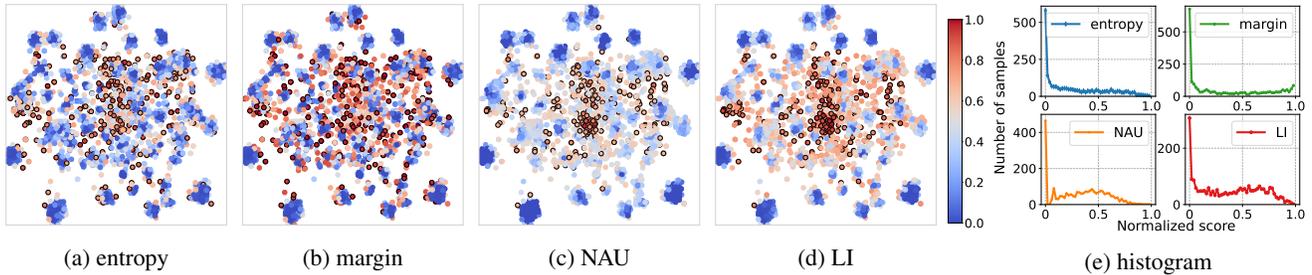


Figure A.3: (a-d) *t*-SNE visualization of target features on Office-Home $Rw \rightarrow Ar$. Samples are colored according to their normalized uncertainty scores, where red indicates large values and blue indicates small values. The top 10% samples with highest scores are marked with black boarders. (e) Histogram of target samples by normalized scores.

Table A.1: Running time on **Office-Home** $Ar \rightarrow Rw$ and **VisDA** in seconds.

	AL					DA				
	BADGE	AADA	CLUE	MHPL	LAS	ft	MCC	CDAC	RAA	LAA
$Ar \rightarrow Rw$	17.6 \pm 0.5	18.4 \pm 0.8	24.0 \pm 0.9	17.8 \pm 0.4	18.2 \pm 0.7	30.2 \pm 8.7	39.6 \pm 8.6	67.4 \pm 16.5	63.0 \pm 18.2	62.0 \pm 10.8
VisDA	80.4 \pm 2.4	57.0 \pm 1.7	733.2 \pm 56.6	70.0 \pm 2.6	105.6 \pm 17.9	1136.6 \pm 45.3	1623.4 \pm 18.4	2662.0 \pm 17.2	1449.8 \pm 15.8	1624.2 \pm 245.2

Existing SSDA methods can be directly applied in ADA. In the paper, we have compared different active query methods when using MME [29] and CDAC [15] as the model adaptation methods. Additional results on Office-Home with CDAC is provided in Tab. A.2. Generally, the proposed LAS can select more informative samples than other active selection criteria.

Nevertheless, it may be sub-optimal to simply combine active query with existing SSDA methods. A unified ADA solution that considers both active query and model adaptation would be better effective. Tables A.3, A.4 present comparison results with two state-of-the-art SSDA methods, ECACL [16] and CDAC [15]. The comparison results are taken from [8]. It should be noted that although ADA can select more informative labeled samples, the performance is also affected by the way to utilize unlabeled data. From the tables, ECACL and CDAC surpass three early ADA methods. The state-of-the-art LAMDA method selects target data to approximate the entire target distribution, and addresses the issue of label distribution mismatch

between source and target domains. It obtains better performances than SSDA arts. Our proposed LADA (LAS w/ LAA) selects locally-representative samples, and progressively expand the labeled data with confident samples in a class-balanced manner. LADA outperforms LAMDA by +2.3% on Office-Home and +1.0% on Office-Home RSUT. When replacing LAA with CDAC in LADA, the performance drops, as we show in the paper.

Uncertainty measures in LAS. Figure A.3 visualizes the target features on Office-Home $Rw \rightarrow Ar$. Similar to the plots in Fig. 3 of the paper, entropy and margin include some outliers in the top 10% samples (see circles with black boarders in the bottom part of Figs. A.3a, A.3b). For NAU in Fig. A.3c, target data have small normalized scores. Target data with high normalized LI scores form several small clusters in Fig. A.3d. From the histogram in Fig. A.3e, a maximal can be observed around 0.6-0.7 for LI. These phenomena are similar to Office-31 $W \rightarrow A$ in the paper.

Training with a joint labeled set. In the implementation of some early ADA works [5, 42], the queried labeled data are

Table A.2: Accuracies (%) on **Office-Home** with 5% labeled target samples. ([†]Training mini-batches are sampled from a joint labeled set.)

AL method	DA method	Ar→Cl	Ar→Pr	Ar→Rw	Cl→Ar	Cl→Pr	Cl→Rw	Pr→Ar	Pr→Cl	Pr→Rw	Rw→Ar	Rw→Cl	Rw→Pr	Avg.
RAN		60.5	78.4	79.0	60.2	74.4	72.7	61.5	56.2	77.6	69.9	59.6	82.8	69.4
ENT		62.8	81.5	82.7	64.1	78.1	75.7	63.4	57.5	81.2	72.8	62.5	87.3	72.5
MAR		64.0	81.8	82.7	64.1	79.0	74.9	64.6	59.9	80.7	73.2	64.6	87.8	73.1
CoreSet		58.5	77.1	79.3	60.8	72.4	71.8	60.9	54.9	77.3	70.9	58.8	81.3	68.7
BADGE	ft w/ CE loss [†]	65.5	83.6	82.1	63.1	79.8	75.3	64.9	61.0	80.8	73.1	65.1	87.1	73.5
AADA		61.8	82.0	82.1	62.3	77.7	76.0	63.1	59.4	<u>81.8</u>	72.9	62.4	87.2	72.4
CLUE		65.3	81.8	81.7	62.6	78.5	74.8	63.9	61.4	79.9	72.9	63.1	87.6	72.8
CONF		63.4	81.9	82.9	63.8	78.2	75.8	64.2	60.2	81.6	73.3	63.2	87.4	73.0
MHPL		65.6	82.1	82.9	<u>65.3</u>	79.1	74.6	64.7	61.4	81.6	73.3	63.7	88.1	73.5
LAS		<u>67.2</u>	<u>84.3</u>	<u>83.1</u>	65.1	<u>80.9</u>	<u>77.0</u>	<u>65.3</u>	<u>62.5</u>	81.4	<u>73.8</u>	<u>66.7</u>	<u>89.0</u>	<u>74.7</u>
RAN		61.6	78.8	80.1	67.7	80.2	77.6	68.7	61.9	79.7	74.1	63.0	85.2	73.2
ENT		62.9	81.9	83.4	69.0	82.0	80.0	70.3	63.3	84.2	75.6	67.7	87.1	75.6
MAR		65.6	83.8	83.3	69.0	83.7	81.0	70.2	65.7	84.6	75.9	67.0	88.1	76.5
CoreSet		58.9	77.7	79.6	67.1	77.9	77.2	67.4	58.6	81.6	73.6	63.4	83.4	72.2
BADGE	CDAC	63.3	80.4	81.3	69.6	83.0	78.8	70.4	62.7	83.6	76.1	67.4	88.0	75.4
AADA		61.8	81.8	82.8	69.6	83.2	80.4	70.7	63.5	84.3	76.2	66.2	87.2	75.6
CLUE		65.2	83.6	82.3	68.8	84.4	79.8	69.7	64.9	83.6	75.2	68.0	87.5	76.1
CONF		62.6	82.9	<u>83.8</u>	70.6	83.6	79.7	70.0	64.6	84.3	76.4	66.8	88.1	76.1
MHPL		65.5	82.4	<u>82.7</u>	70.8	84.1	<u>81.7</u>	70.5	66.2	84.3	76.9	68.7	87.8	76.8
LAS		<u>67.4</u>	<u>85.4</u>	83.1	<u>71.0</u>	<u>85.0</u>	<u>81.7</u>	<u>72.1</u>	<u>67.8</u>	<u>85.1</u>	<u>77.4</u>	<u>70.4</u>	<u>89.5</u>	<u>78.0</u>
LAS	RAA	71.2	88.1	85.3	73.2	87.8	83.8	72.6	72.2	86.6	79.2	74.4	91.7	80.5
LAS	LAA	71.2	87.4	84.6	72.1	87.0	83.6	71.5	71.6	85.3	79.3	75.5	90.4	80.0

Table A.3: Comparison with Semi-Supervised Domain Adaptation methods on **Office-Home** using 10%-budget.

Task	Method	Ar→Cl	Ar→Pr	Ar→Rw	Cl→Ar	Cl→Pr	Cl→Rw	Pr→Ar	Pr→Cl	Pr→Rw	Rw→Ar	Rw→Cl	Rw→Pr	Avg.
SSDA	ECACL	72.2	86.7	82.8	70.5	85.0	82.6	70.9	71.5	82.9	76.0	74.0	88.9	78.7
	CDAC	69.5	83.2	80.2	66.9	82.4	78.7	66.1	70.6	80.9	72.3	70.5	87.2	75.7
ADA	TQS	64.3	84.8	83.5	66.1	81.0	76.7	66.5	61.4	82.0	73.7	65.9	88.5	74.5
	CLUE	62.1	80.6	73.9	55.2	76.4	75.4	53.9	62.1	80.7	67.5	63.0	88.1	69.9
	S ³ VAADA	67.8	83.9	82.9	67.0	81.4	79.5	65.8	65.9	82.4	74.8	68.6	87.9	75.7
	LAMDA	74.8	88.5	86.9	73.8	88.2	83.3	74.6	75.5	86.9	80.8	77.8	91.7	81.9
	LADA	77.2	91.9	88.1	76.9	91.1	86.8	76.6	78.1	88.3	82.0	79.0	93.8	84.2

Table A.4: Comparison with SSDA methods on **Office-Home RSUT** using 10%-budget.

Task	Method	C→P	C→R	P→C	P→R	R→C	R→P	Avg.
SSDA	ECACL	78.6	68.6	59.5	77.1	61.9	82.0	71.3
	CDAC	73.0	58.7	55.8	73.3	50.3	77.3	64.7
ADA	TQS	69.4	65.7	53.0	76.3	53.1	81.1	66.4
	CLUE	69.7	65.9	57.1	73.4	59.5	82.7	68.1
	S ³ VAADA	73.0	63.0	50.7	69.6	52.6	78.3	64.5
	LAMDA	81.2	75.7	64.1	81.6	65.1	87.2	75.8
	LADA	83.2	77.2	63.8	83.0	65.4	88.1	76.8

added to the source labeled data, and training mini-batches are sampled from this joint labeled set. The objective is

$$\mathcal{L} = \mathbb{E}_{(x,y) \sim \mathcal{D}_s \cup \mathcal{D}_{t1}} \ell_{ce}(h(x), y) \quad (\text{A.6})$$

where ℓ_{ce} is the cross entropy loss. Differently, recent works [41, 8] and ours adopt

$$\mathcal{L} = \mathbb{E}_{(x,y) \sim \mathcal{D}_s} \ell_{ce}(h(x), y) + \mathbb{E}_{(x,y) \sim \mathcal{D}_{t1}} \ell_{ce}(h(x), y) \quad (\text{A.7})$$

Comparing Eq. A.6 with Eq. A.7, the advantage of training with a joint labeled set is that it only needs to back-propagate through one batch of data, thus reducing the

memory and computation usage. The disadvantage is that the labeled data set is dominated by the source data. When there is a large domain gap (e.g. when label distribution shift exists), the performance may be hurt.

Nevertheless, to better demonstrate the effectiveness of LAS, Table A.2 lists the results using fine-tuning with a joint labeled set. Accuracies are slightly lower than their counterparts in Table 1 of the main paper. LAS still achieves the best scores among all AL methods.

Visualization of standard deviations. Figure A.4 plots the standard deviations on two Office-31 tasks over 3 repeated experiments. Performances are relatively stable to different random initializations. Of all active selection methods, our proposed LAS obtains the highest average accuracies.

Visualization of LAS. To visualize how LAS selects target samples, we present *t*-SNE plots of target features on Office-Home Pr→Ar and Ar→Pr in Fig. A.5 and Fig. A.6, respectively. We choose the 10th epoch, where 10% target data are selected as candidates based on LI-scores, of which 1% target data from cluster centroids are selected for querying labels. Candidate, selected and remaining target

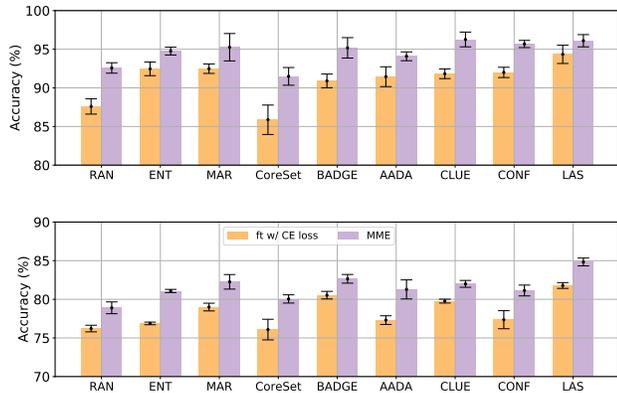


Figure A.4: Visualization of standard deviations on Office-31 A→W (upper) and W→A (lower) using 5%-budget.

samples are marked with squares, stars and points, respectively. The top 20 candidates and queried images are also displayed under the t -SNE plots. As can be seen, the candidates (*i.e.*, samples with large LI-scores) generally lie in the regions where model predictions are inconsistent. It is also difficult to distinguish their semantic labels visually, especially for Pr→Ar, indicating that these images are hard cases. There are some highly similar images in the candidates. After the second step of diverse selection, images selected for querying labels become much more diverse.

References

- [1] Jordan T Ash, Chicheng Zhang, Akshay Krishnamurthy, John Langford, and Alekh Agarwal. Deep batch active learning by diverse, uncertain gradient lower bounds. In *ICLR*, 2020. 2, 5
- [2] Maria-Florina Balcan, Andrei Broder, and Tong Zhang. Margin based active learning. In *COLT*, pages 35–50, 2007. 2
- [3] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data augmentation with a reduced search space. In *CVPRW*, pages 702–703, 2020. 5
- [4] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *ICLR*, 2021. 2
- [5] Bo Fu, Zhangjie Cao, Jianmin Wang, and Mingsheng Long. Transferable query selection for active domain adaptation. In *CVPR*, pages 7272–7281, 2021. 1, 3, 4, 5, 9, 10
- [6] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *JMLR*, 17(1):2096–2030, 2016. 1, 2
- [7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *ICCV*, pages 770–778, 2016. 5, 9
- [8] Sehyun Hwang, Sohyun Lee, Sungyeon Kim, Jungseul Ok, and Suha Kwak. Combating label distribution shift for active domain adaptation. In *ECCV*, 2022. 3, 5, 6, 9, 10, 11
- [9] Ying Jin, Ximei Wang, Mingsheng Long, and Jianmin Wang. Minimum class confusion for versatile domain adaptation. In *ECCV*, pages 464–480, 2020. 6
- [10] Ajay J Joshi, Fatih Porikli, and Nikolaos P Papanikolopoulos. Scalable active learning for multiclass image classification. *TPAMI*, 34(11):2259–2273, 2012. 1, 5
- [11] Kwanyoung Kim, Dongwon Park, Kwang In Kim, and Se Young Chun. Task-aware variational adversarial active learning. In *CVPR*, pages 8166–8175, 2021. 2
- [12] Andreas Kirsch, Joost Van Amersfoort, and Yarin Gal. Batchbald: Efficient and diverse batch acquisition for deep bayesian active learning. *NeurIPS*, 32, 2019. 2
- [13] David D Lewis and Jason Catlett. Heterogeneous uncertainty sampling for supervised learning. In *Machine Learning Proceedings 1994*, pages 148–156. Elsevier, 1994. 2
- [14] Bo Li, Yezhen Wang, Tong Che, Shanghang Zhang, Sicheng Zhao, Pengfei Xu, Wei Zhou, Yoshua Bengio, and Kurt Keutzer. Rethinking distributional matching based domain adaptation. *arXiv preprint arXiv:2006.13352*, 2020. 2
- [15] Jichang Li, Guanbin Li, Yemin Shi, and Yizhou Yu. Cross-domain adaptive clustering for semi-supervised domain adaptation. In *CVPR*, pages 2505–2514, 2021. 1, 2, 6, 10
- [16] Kai Li, Chang Liu, Handong Zhao, Yulun Zhang, and Yun Fu. Ecacl: A holistic framework for semi-supervised domain adaptation. In *ICCV*, pages 8578–8587, 2021. 10
- [17] Jian Liang, Dapeng Hu, and Jiashi Feng. Do we really need to access the source data? source hypothesis transfer for unsupervised domain adaptation. In *ICML*, pages 6028–6039, 2020. 2
- [18] Hong Liu, Jianmin Wang, and Mingsheng Long. Cycle self-training for domain adaptation. *NeurIPS*, 34, 2021. 1, 2
- [19] Zimo Liu, Jingya Wang, Shaogang Gong, Huchuan Lu, and Dacheng Tao. Deep reinforcement active learning for human-in-the-loop person re-identification. In *ICCV*, pages 6122–6131, 2019. 2
- [20] Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment matching for multi-source domain adaptation. In *ICCV*, pages 1406–1415, 2019. 9
- [21] Xingchao Peng, Ben Usman, Neela Kaushik, Judy Hoffman, Dequan Wang, and Kate Saenko. Visda: The visual domain adaptation challenge. *arXiv preprint arXiv:1710.06924*, 2017. 5, 9
- [22] Viraj Prabhu, Arjun Chandrasekaran, Kate Saenko, and Judy Hoffman. Active domain adaptation via clustering uncertainty-weighted embeddings. In *ICCV*, pages 8505–8514, 2021. 1, 3, 4, 5, 9
- [23] Viraj Prabhu, Shivam Khare, Deeksha Kartik, and Judy Hoffman. Sentry: Selective entropy optimization via committee consistency for unsupervised domain adaptation. In *ICCV*, pages 8558–8567, 2021. 2



Figure A.6: Visualization of LAS sampling on Office-Home Ar→Pr. The first row presents t -SNE plots. Squares denote candidate target samples based on LI-scores; stars denote selected target samples for querying labels; and points denote the rest target samples. Each marker is colored according to its (left) ground-truth label and (right) pseudo label from the current model. The last two rows plot top 20 candidate samples and queried samples with largest LI-scores, respectively.

- [24] Piyush Rai, Avishek Saha, Hal Daumé III, and Suresh Venkatasubramanian. Domain adaptation meets active learning. In *Proceedings of the NAACL HLT 2010 Workshop on Active Learning for Natural Language Processing*, pages 27–32, 2010. [3](#)
- [25] Harsh Rangwani, Arihant Jain, Sumukh K Aithal, and R Venkatesh Babu. S3vaada: Submodular subset selection for virtual adversarial active domain adaptation. In *ICCV*, pages 7516–7525, 2021. [3](#), [5](#)
- [26] Pengzhen Ren, Yun Xiao, Xiaojun Chang, Po-Yao Huang, Zhihui Li, Brij B Gupta, Xiaojiang Chen, and Xin Wang. A survey of deep active learning. *ACM Computing Surveys (CSUR)*, 54(9):1–40, 2021. [2](#)
- [27] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet large scale visual recognition challenge. *IJCV*, 115(3):211–252, 2015. [9](#)
- [28] Kate Saenko, Brian Kulis, Mario Fritz, and Trevor Darrell. Adapting visual category models to new domains. In *ECCV*, pages 213–226, 2010. [5](#), [9](#)
- [29] Kuniaki Saito, Donghyun Kim, Stan Sclaroff, Trevor Darrell, and Kate Saenko. Semi-supervised domain adaptation via minimax entropy. In *ICCV*, pages 8050–8058, 2019. [1](#), [2](#), [6](#), [10](#)
- [30] Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set approach. In *ICLR*, 2018. [1](#), [2](#), [5](#)
- [31] Samarth Sinha, Sayna Ebrahimi, and Trevor Darrell. Variational adversarial active learning. In *ICCV*, pages 5972–5981, 2019. [2](#)
- [32] Jong-Chyi Su, Yi-Hsuan Tsai, Kihyuk Sohn, Buyu Liu, Subhransu Maji, and Manmohan Chandraker. Active adversarial domain adaptation. In *WACV*, pages 739–748, 2020. [1](#), [3](#), [5](#)
- [33] Tao Sun, Cheng Lu, and Haibin Ling. Prior knowledge guided unsupervised domain adaptation. In *ECCV*, pages 639–655. Springer, 2022. [2](#)
- [34] Tao Sun, Cheng Lu, and Haibin Ling. Domain adaptation with adversarial training on penultimate activations. In *AAAI*, pages 9935–9943, 2023. [2](#)
- [35] Tao Sun, Cheng Lu, Tianshuo Zhang, and Haibin Ling. Safe self-refinement for transformer-based domain adaptation. In *CVPR*, 2022. [1](#), [2](#)
- [36] Shuhan Tan, Xingchao Peng, and Kate Saenko. Class-imbalanced domain adaptation: an empirical odyssey. In *ECCV Workshops*, pages 585–602, 2020. [5](#), [9](#)
- [37] Hui Tang, Ke Chen, and Kui Jia. Unsupervised domain adaptation via structurally regularized deep clustering. In *CVPR*, pages 8725–8735, 2020. [2](#)
- [38] Eric Tzeng, Judy Hoffman, Ning Zhang, Kate Saenko, and Trevor Darrell. Deep domain confusion: Maximizing for domain invariance. *arXiv preprint arXiv:1412.3474*, 2014. [2](#)
- [39] Hemant Venkateswara, Jose Eusebio, Shayok Chakraborty, and Sethuraman Panchanathan. Deep hashing network for unsupervised domain adaptation. In *CVPR*, pages 5018–5027, 2017. [5](#)
- [40] Fan Wang, Zhongyi Han, Zhiyan Zhang, Rundong He, and Yilong Yin. Mhpl: Minimum happy points learning for active source free domain adaptation. In *CVPR*, pages 20008–20018, 2023. [3](#), [6](#)
- [41] Binhui Xie, Longhui Yuan, Shuang Li, Chi Harold Liu, Xinjing Cheng, and Guoren Wang. Active learning for domain adaptation: An energy-based approach. In *AAAI*, 2021. [1](#), [3](#), [4](#), [5](#), [9](#), [11](#)
- [42] Ming Xie, Yuxi Li, Yabiao Wang, Zekun Luo, Zhenye Gan, Zhongyi Sun, Mingmin Chi, Chengjie Wang, and Pei Wang. Learning distinctive margin toward active domain adaptation. In *CVPR*, 2022. [1](#), [3](#), [4](#), [5](#), [9](#), [10](#)
- [43] Zhao Xu, Kai Yu, Volker Tresp, Xiaowei Xu, and Jizhi Wang. Representative sampling for text classification using support vector machines. In *ECIR*, pages 393–407, 2003. [2](#)
- [44] Luyu Yang, Yan Wang, Mingfei Gao, Abhinav Shrivastava, Kilian Q Weinberger, Wei-Lun Chao, and Ser-Nam Lim. Deep co-training with task decomposition for semi-supervised domain adaptation. In *ICCV*, pages 8906–8916, 2021. [2](#)
- [45] Shiqi Yang, Joost van de Weijer, Luis Herranz, and Shangling Jui. Exploiting the intrinsic neighborhood structure for source-free domain adaptation. *NeurIPS*, 34, 2021. [2](#)
- [46] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115, 2021. [1](#)