

Active Learning for Domain Adaptation: An Energy-based Approach

Binhui Xie¹ Longhui Yuan¹ Shuang Li^{1,*} Chi Harold Liu¹ Xinjing Cheng^{2,3} Guoren Wang¹

¹Beijing Institute of Technology ²Tsinghua University ³Inceptio Technology
{binhuixie, longhuiyuan, shuangli, chiliu, wanggr}@bit.edu.cn, cnorbot@gmail.com

Abstract

Unsupervised domain adaptation has recently emerged as an effective paradigm for generalizing deep neural networks to new target domains. However, there is still enormous potential to be tapped to reach the fully supervised performance. In this paper, we present a novel active learning strategy to assist knowledge transfer in the target domain, dubbed active domain adaptation. We start from an observation that energy-based models exhibit *free energy biases* when training (source) and test (target) data come from different distributions. Inspired by this inherent mechanism, we empirically reveal that a simple yet efficient energy-based sampling strategy sheds light on selecting the most valuable target samples than existing approaches requiring particular architectures or computation of the distances. Our algorithm, *Energy-based Active Domain Adaptation (EADA)*, queries groups of target data that incorporate both domain characteristic and instance uncertainty into every selection round. Meanwhile, by aligning the free energy of target data compact around the source domain via a regularization term, domain gap can be implicitly diminished. Through extensive experiments, we show that EADA surpasses state-of-the-art methods on well-known challenging benchmarks with substantial improvements, making it a useful option in the open world. Code is available at <https://github.com/BIT-DA/EADA>.

Introduction

In recent years, we have witnessed great strides in diverse machine learning problems with the success of deep neural networks (Krizhevsky, Sutskever, and Hinton 2012a). At the moment, however, these leaps in performance come only when massive labeled data are available. This limits their usage in many practical applications, such as autonomous driving with abundant unlabeled data (Yogamani et al. 2019) and medical diagnosis with high labeling cost (Ronneberger, Fischer, and Brox 2015). Moreover, even labeling all available data is not an excellent solution, as it's impossible to fully capture the way the world looks in a single dataset, let alone the fact that the test data rarely matches the data seen during training. Recognizing the challenges, extensive studies have

been explored in domain adaptation (DA), which transfer the knowledge from a label-rich source domain to an unlabeled target domain (Pan and Yang 2010; Ganin and Lempitsky 2015; Tzeng et al. 2015, 2017; Long et al. 2019, 2018; Bousmalis et al. 2017; Saito et al. 2018; Li et al. 2021a,c). The performance of DA, in spite of great success, often falls far behind that of supervised learning. In practice, it may be feasible to obtain extra annotations for a small set of the target domain. But to be effective, it is critical to identify samples with high information only via active learning (Prince 2004; Hanneke 2014; Bickel, Brückner, and Scheffer 2009).

While previous active learning studies drastically lower human annotation costs, they are impractical when test data are collected from out-of-distribution. How can we design an efficient and practical sampling strategy for domain adaptation? For one thing, it is essential to determine which target samples will, once labeled, boost the accuracy and generalization considerably. For another, it remains the boundary to explore how to effectively utilize limited labeled data from the target domain to perform adaptation. Aware of this need, researchers have developed an array of active domain adaptation (Active DA) methods (Chattopadhyay et al. 2013; Rai et al. 2010; Su et al. 2020; Fu et al. 2021; Prabhu et al. 2021; Chan and Ng 2007). Prior works mainly focus on assessing how private each target data is according to the output of a domain discriminator or calculating its distance to the cluster centroids. However, these additional procedures either select target samples that are originally well aligned with the source domain or increase the computational overhead, which limits their capability. Therefore, a simple yet efficient solution is urgently desired.

In this paper, we advocate the use of energy-based models (EBMs) (LeCun et al. 2006) to help realize the potential of active learning under domain shift. For any given x (e.g., an image), an EBM approach gives the lowest energy to the correct answer y (e.g., a label). Grathwohl et al. (2020) and Liu et al. (2020) have demonstrated that energy-based training improves calibration and better distinguishes in- and out-of-distribution samples than the standard discriminative classifier. At this point, we begin with investigating the distributions of free energy on source and target domains using diverse methods and make several observa-

*Corresponding author.

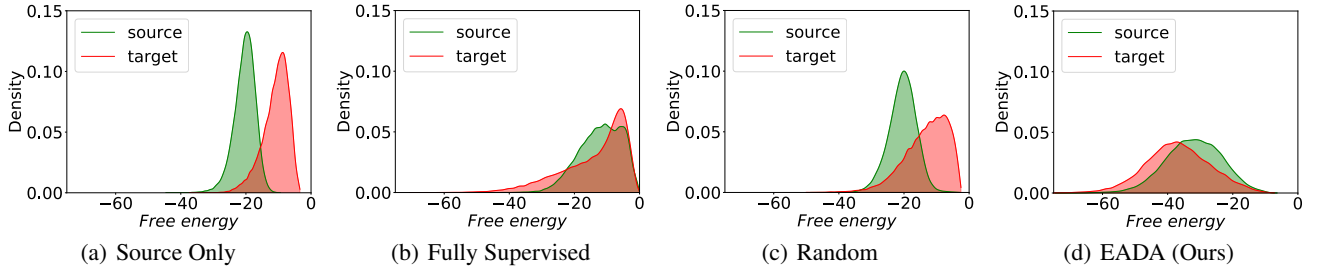


Figure 1: (a & b) Free energy distribution biases between source and target domains on VisDA-2017 from “Source Only”. We then contrast the distributions from (c) a native baseline of random selection and (d) our energy-based selection strategy. EADA exhibits better-aligned distribution than “Random” and is similar to “Full Supervised” i.e., all source and target data are labeled.

tions from Fig. 1. First, a model trained only on labeled source data will cause the free energy distribution of the supervised source data to be lower than that of the unlabeled target data, that is, *free energy biases* between the two domains (Fig. 1(a)). Then, an interesting finding is that these two distributions tend to be consistent using “Full Supervised” (Fig. 1(b)). Next, the biases eliminate slightly when a few unlabeled target data are randomly annotated in the training process (Fig. 1(c)). Lastly, using our algorithm to identify limited target instances for labeling, surprisingly, it well matches two distributions as same as the situation of “Full Supervised” (Fig. 1(d)). We conjecture that there exist redundant or trivial data in the target domain itself that do little to help the learning objective.

Intuitively, we provide both mathematical insights and empirical evidence that an energy-based active learning scheme is desirable for active domain adaptation. A central theme of this work is that we design an approach, Energy-based Active Domain Adaptation (EADA), which adequately ensures samples that are representative of the entire target domain to be selected by considering both domain characteristic and instance uncertainty. More precisely, as mentioned above, the free energies of most labeled source data are lower than that of unlabeled target data. Thus, we can treat the intrinsic free energy of an unlabeled target sample as a surrogate metric to reflect the domain characteristic. Naturally, the target samples with higher free energy are more dissimilar to source data, and thus be typical for target distribution. In addition, we assess the value of minimum energy versus second-minimum energy (MvSM) for each unlabeled target data to quantify its uncertainty under the current model. To this end, given the labeling budget in each round, we first maintain a candidate set from unlabeled target data with higher free energies and then select samples with significant MvSM values from candidates. Furthermore, free energy can also serve as a regularization signal in the form of an alignment loss to implicitly diminish domain shift, which is complementary to our active strategy.

In summary, our work makes the following contributions:

- We provide a new perspective to select a highly informative subset of unlabeled target data under domain shift via exploiting *free energy biases* between the two domains.
- We complement empirical results with theoretical inves-

tigations in the method section and establish an intuitive sufficient condition when it would help.

- Though simple, EADA attains excellent results with quite limited labeling expenses. Extensive experiments and in-depth analysis demonstrate its effectiveness.

Related Work

Active learning (AL) has been studied for decades in both theory and practice (Settles 2009; Dasgupta 2011; Bachman, Sordoni, and Trischler 2017; Gal, Islam, and Ghahramani 2017). A case in point is to search informative data for labeling in order to learn a satisfactory model at a low annotation cost. Most popular algorithms formulate and solve it by uncertainty sampling. They select samples about which the current model is uncertain (Schohn and Cohn 2000; Joshi, Porikli, and Papanikolopoulos 2009; Wang and Shang 2014). Another line of work turns to representative sampling (Sener and Savarese 2018; Sinha, Ebrahimi, and Darrell 2019; Gissin and Shalev-Shwartz 2019), which picks a set of typical samples via clustering or core-set selection.

Recently, several studies have leveraged a hybrid of the above active sampling objectives to achieve promising results, such as Ash et al. (2020). However, these conventional AL methods cannot deal with the domain shift issues for domain adaptation, whereas our method aims to overcome this challenge by leveraging a simple energy-based strategy.

Unsupervised domain adaptation (UDA) studies the problem of transferring knowledge gained from an abundant labeled source domain to a target domain where labeled data are scarce (Ganin and Lempitsky 2015; Long et al. 2019, 2017; Xu et al. 2019; Li et al. 2018, 2020, 2021b; Hoffman et al. 2018; Zou, Yang, and Wu 2021; Gong et al. 2012). A series of works minimizes the domain discrepancy at the uppermost layer of deep neural networks using maximum mean discrepancy (Gretton et al. 2007) or adversarial training (Goodfellow et al. 2014). Recently, some methods allow a few target data labeled, e.g., semi-supervised DA (Saito et al. 2019) and few-shot DA (Teshima, Sato, and Sugiyama 2020). Though impressive, they randomly select a few data to annotate, neglecting which target samples should be labeled given a fixed labeling budget. Consequently, some selected samples are originally well predicted by the current model. In contrast, our work differentiates itself by allowing

the model to acquire labels for valuable target samples via an oracle. As such, it would have the best potential performance gain compared with randomly picking labels.

Active domain adaptation (Active DA). The seminal work (Rai et al. 2010) has demonstrated the synergy between active learning and domain adaptation, which facilitates AL in a domain of interest with the aid of the knowledge from a related domain. Recently, Su et al. (2020) and Fu et al. (2021) incorporate Active DA with advanced tools, such as adversarial training, both of which identify domainness via a learned domain discriminator. However, it may give identically high scores to most target data, thus not adequately ensuring that selected samples are representative of the entire target distribution. A parallel line of work instead proposes to select active samples via clustering. For example, Prabhu et al. (2021) cluster deep embeddings of target data weighted by the uncertainty and select nearest neighbors to the inferred cluster centroids for labeling. However, clustering-based strategies have some drawbacks in nature. First, they encounter a computational burden and could hardly be applied on large data sets. Second, the clustering is sensitive to noise and easy to collapse.

Originating from energy-based models, our method adapts the concept of energy to identify limited target samples that are most unique to the target distribution and meanwhile complementary to labeled source data. It yields a new sampling protocol that accounts for domain characteristic and instance uncertainty together. Also, it has no extra parameters that need to be optimized and learning is efficient.

Method

In active domain adaptation (Active DA), we have access to a labeled source domain $\mathcal{S} = \{(x_s, y_s)\}$ and an unlabeled target domain $\mathcal{T} = \{x_t\}$ from different distributions. Following the standard Active DA setting (Fu et al. 2021; Prabhu et al. 2021), B active samples are selected in the target domain for annotation, which are much smaller than the amount of \mathcal{T} . Therefore, the entire target domain consists of a labeled pool \mathcal{T}_l and an unlabeled pool \mathcal{T}_u , i.e., $\mathcal{T} = \mathcal{T}_l \cup \mathcal{T}_u$. The goal is to learn a neural network with parameter θ that brings good generalization on the target. In this work, we introduce an energy-based strategy to select the most valuable target samples to assist the knowledge transfer.

Energy-based Models Revisit

The essence of machine learning is to encode dependencies between variables. Let us consider an energy-based model (EBM) with two sets of variables x (a high-dimensional variable) and y (a discrete variable). Training this model consists in finding an energy function i.e., $E(x, y)$ that gives the lowest energy to correct answer and higher energy to all other (incorrect) answers¹. Precisely, the model must produce the value y^* for which $E(x, y)$ is the smallest:

$$y^* = \arg \min_{y \in \mathcal{Y}} E(x, y). \quad (1)$$

Generally, the size of set \mathcal{Y} is small for classification, hence the inference procedure can simply compute $E(x, y)$ for all possible values of $y \in \mathcal{Y}$ and pick the smallest.

With the energy function, the joint probability of input x and label y can be estimated through the Gibbs distribution:

$$p(x, y) = \exp(-E(x, y)) / Z, \quad (2)$$

where $Z = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} \exp(-E(x, y))$ is called the partition function that marginalizes over x and y . It should be noted that the above transformation of energy into probability is only possible if Z converges. By marginalizing out y , we obtain the probability density for x as well,

$$p(x) = \sum_{y \in \mathcal{Y}} p(x, y) = \sum_{y \in \mathcal{Y}} \exp(-E(x, y)) / Z. \quad (3)$$

Intuitively, in Active DA, to select the most representative target samples, one can directly estimate the probability of occurrence for each target sample from Eq. (3) and then those samples with lower probabilities should be selected.

Unfortunately, one cannot compute or even reliably estimate Z . Therefore, we turn to *free energy* i.e., $\mathcal{F}(x)$, a function hidden in EBMs that serves as the ‘‘rationality’’ of the occurrence of the variable x . Mathematically, the probability density for x can also be expressed as

$$p(x) = \frac{\exp(-\mathcal{F}(x))}{\sum_{x \in \mathcal{X}} \exp(-\mathcal{F}(x))}. \quad (4)$$

This formulation indicates that $\mathcal{F}(x)$ could be substituted for $p(x)$ to select the target samples that have lower probabilities. By connecting Eq. (3) and Eq. (4), we have

$$\mathcal{F}(x) = -\log \sum_{y \in \mathcal{Y}} \exp(-E(x, y)). \quad (5)$$

Energy-based Active Domain Adaptation

We now wish to take advantage of a new perspective of the energy-based model (EBM) to gain the benefits of active domain adaptation, where the biases of free energy between source- and target-domain data allow effective selection and adaptation. In the following, we first describe how to train an EBM with several loss functions. We then describe using an energy-based sampling strategy to identify the most informative unlabeled target data to annotate. At last, we provide an intuitive sufficient condition when it helps.

Training process Given a set of labeled source samples $\mathcal{S} = \{(x_s, y_s)\}$, we want to train a well-behaved EBM that gives the lowest energy to the correct answer and higher energy to all other (incorrect) answers. To this end, we utilize a commonly used loss in EBMs, i.e., the negative log-likelihood loss that comes from probabilistic modeling to train a model for classification, and it can be formulated as

$$\mathcal{L}_{nll}(x, y; \theta) = E(x, y; \theta) + \frac{1}{\tau} \log \sum_{c \in \mathcal{Y}} \exp(-\tau E(x, c; \theta)), \quad (6)$$

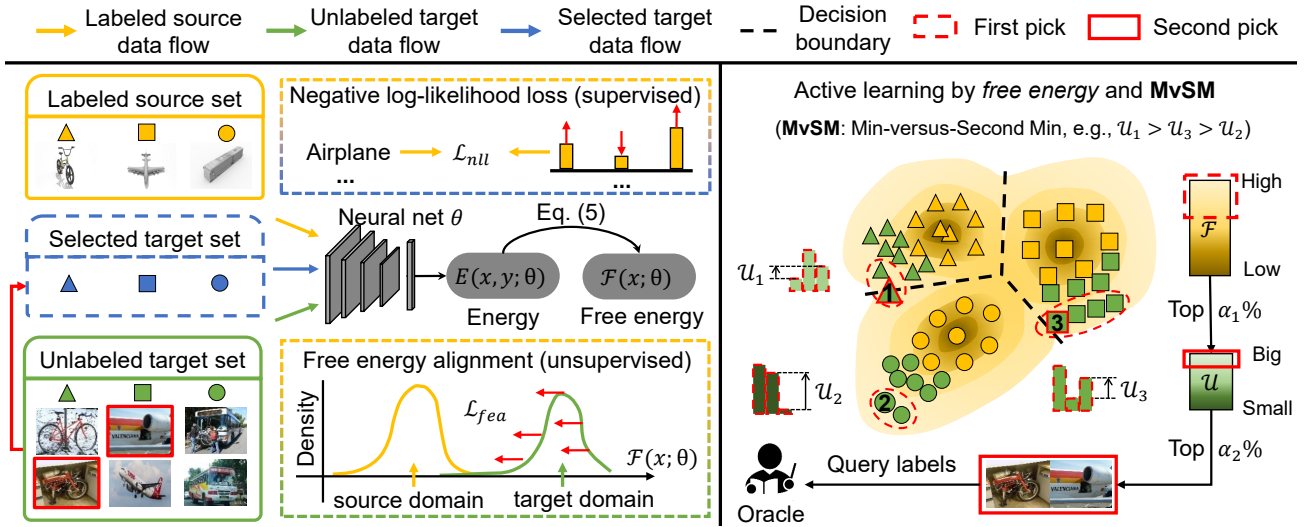
where τ ($\tau > 0$) is the reverse temperature and a low value corresponds to smooth partition of energy over the space \mathcal{Y} . For simplicity, we fix $\tau=1$, and then we have

$$\mathcal{L}_{nll}(x, y; \theta) = E(x, y; \theta) - \mathcal{F}(x; \theta). \quad (7)$$

The second term in Eq. (7) will cause the energies of all answers to be pulled up. The energy of the correct answer is also pulled up, but not as hard as it is pushed down by the first term. An analysis of gradient is presented in Appendix².

¹See (LeCun et al. 2006) for a comprehensive tutorial.

²Appendix can be found at <https://arxiv.org/abs/2112.01406>.



However, we observe that the values of free energy on target samples are considerably higher than those on source ones, called *free energy biases*. Naturally, one can treat it as a surrogate to reflect the domain divergence. By designing a simple regularization term, these biases can be reduced, which to some extent aligns the distribution across domains. And the free energy alignment loss \mathcal{L}_{fea} is defined as:

$$\mathcal{L}_{fea}(x; \theta) = \max(0, \mathcal{F}(x; \theta) - \Delta), \quad (8)$$

where $\Delta = \mathbb{E}_{x \sim \mathcal{S}} \mathcal{F}(x; \theta)$ is the average value of the free energy over source data. During training, Δ is estimated via exponential moving average: $\Delta_t = \lambda \Delta_{t-1} + (1 - \lambda) \Delta'_t$, where Δ_t is the estimation of average value in all t mini-batches and Δ'_t is the average value in t^{th} mini-batch and λ is a weight sampled from the uniform distribution, $\lambda \sim U(0, 1)$. Additionally, we experimentally found that such way is comparable with calculating average value over the whole source domain data while improving the efficiency.

Overall, the full learning objective is given by:

$$\min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{S} \cup \mathcal{T}_l} \mathcal{L}_{nll}(x, y; \theta) + \gamma \mathbb{E}_{x \sim \mathcal{T}_u} \mathcal{L}_{fea}(x; \theta), \quad (9)$$

where γ is a loss weight hyperparameter.

Selection process The goal in Active DA is to identify more valuable target samples that, once labeled and used for training, improve the model's accuracy and generalization performance significantly. In practice, we suggest a two-step sampling strategy to adequately ensure such samples by incorporating domain characteristic and instance uncertainty. To be clear, we summarize the training and selection processes based on the above discussion as Algorithm 1.

Step one: we observe that biases of free energy distribution between source and target domains exhibit. Thus, we can utilize this intrinsic free energy of an unlabeled target sample as a surrogate metric to reflect the domain characteristic. Certainly, the target samples with higher free energy

Algorithm 1: EADA algorithm

- 1: **Input:** Labeled source data \mathcal{S} , unlabeled target data \mathcal{T}_u and labeled target set $\mathcal{T}_l = \emptyset$, maximum epoch M , selection rounds R , selection ratios α_1, α_2
- 2: **Output:** Final model parameters θ_M
- 3: **for** $m = 1$ to M **do**
- 4: Update model θ_m via Eq. (9)
- 5: **if** m in R **then**
- 6: $\forall x \in \mathcal{T}_u$, compute free energy $\mathcal{F}(x)$ (Eq. (5)) to serve as measure of domain characteristic
- 7: $\mathcal{T}_l^r \leftarrow$ select $\alpha_1\%$ of \mathcal{F} with the highest values
- 8: $\forall x \in \mathcal{T}_l^r$, compute MvSM $\mathcal{U}(x)$ (Eq. (10)) to serve as measure of instance uncertainty
- 9: $\mathcal{T}_l^r \leftarrow$ select $\alpha_2\%$ of \mathcal{U} with the highest values as active samples for annotating, getting $\mathcal{T}_l = \mathcal{T}_l \cup \mathcal{T}_l^r$
- 10: **end if**
- 11: **end for**

are unique to the target distribution and meanwhile complementary to the labeled source data.

Step two: to measure instance uncertainty, existing methods rely primarily on the entropy score (Su et al. 2020; Prabhu et al. 2021). In contrast, we consider the difference between the energy values of the two answers with the lowest estimated energy value as a measure of uncertainty. Since it is a comparison of the minimum answer and the second minimum answer, we refer to it as the Min-versus-Second-Min (MvSM) strategy and it can be formulated as

$$\mathcal{U}(x) = E(x, y^*; \theta) - E(x, y'; \theta), \quad (10)$$

where $y^* = \arg \min_{y \in \mathcal{Y}} E(x, y; \theta)$ is the lowest energy output and $y' = \arg \min_{y \in \mathcal{Y} \setminus \{y^*\}} E(x, y; \theta)$ is the second-lowest energy output. Such a measure is a more direct way of estimating confusion about class membership from a classi-

fication standpoint. Using the MvSM measure, the instances around the decision boundaries in Fig. 2 during the selection procedure will be selected to query an oracle.

Theoretical Analysis

This section contains our preliminary study of why *free energy biases* exhibit between two different domains. For an energy-based model, we prove that positive gradient inner product between the negative log-likelihood loss function and *free energy* leads to a lower value of *free energy* on labeled source samples during the training process. Limited by space, all the proofs are left for the Appendix.

Before stating our main theoretical result, we first illustrate the general intuition with a toy problem. Considering a simple energy-based model on the classification task, where the network is a one layer linear network parameterized by $\mathbf{W} = (\omega_1 \ \cdots \ \omega_C)^\top \in \mathbb{R}^{C \times N}$, $x \in \mathbb{R}^N$ denotes a source sample, $y \in \{1, \dots, C\}$ denotes the label, we have

$$E(x, j; \mathbf{W}) = \omega_j^\top x, \quad j = 1, \dots, C,$$

$$\mathcal{F}(x; \mathbf{W}) = -\log \sum_{c=1}^C \exp(-\omega_c^\top x), \quad (11)$$

$$\mathcal{L}_{nll}(x, y; \mathbf{W}) = E(x, y; \mathbf{W}) - \mathcal{F}(x; \mathbf{W}).$$

Now we update the the weight matrix \mathbf{W} by one step of gradient descent on \mathcal{L}_{nll} as follows:

$$\mathbf{W}' = \mathbf{W} - \eta \nabla \mathcal{L}_{nll}(x, y; \mathbf{W}), \quad (12)$$

where η is the learning rate and \mathbf{W}' is the updated matrix.

Then we have two lemmas to show that the inner product between the gradients of negative log-likelihood loss function and free energy is positive, and the value of the free energy of a labeled source sample is descending with a step of gradient descent on negative log-likelihood loss function.

Lemma 1. Assume that a toy model correctly predict a labeled source sample (x, y) , we have

$$\langle \nabla \mathcal{L}_{nll}(x, y; \mathbf{W}), \nabla \mathcal{F}(x; \mathbf{W}) \rangle > 0, \quad (13)$$

where $\langle \cdot, \cdot \rangle$ denotes the inner product of gradients.

Lemma 2. Assume that a toy model correctly predict a labeled source sample (x, y) with learning rate $\eta > 0$ we have

$$\mathcal{F}(x; \mathbf{W}) > \mathcal{F}(x; \mathbf{W}'), \quad (14)$$

To summarize, if the positive gradient inner product between the negative log-likelihood loss function and *free energy*, free energy biases are exhibited. Our main theoretical results extend this to general deep neural networks.

Theorem 1. Let $\mathcal{L}_{nll}(x, y; \theta)$ denote the negative log-likelihood loss on source domain (x, y) with parameters of deep network θ and $\mathcal{F}(x; \theta)$ denote the free energy of x . Assume that $\forall (x, y)$, $\mathcal{L}_{nll}(x, y; \theta)$ is differentiable, β -smooth in θ and $\forall \theta$, $\|\nabla \mathcal{L}_{nll}(x, y; \theta)\| < G$, $\|\nabla \mathcal{F}(x; \theta)\| < G$. With learning rate $\eta \in (0, \frac{2\varepsilon}{\beta G^2})$, and for every (x, y) such that

$$\langle \nabla \mathcal{L}_{nll}(x, y; \theta), \nabla \mathcal{F}(x; \theta) \rangle > \varepsilon, \quad (15)$$

where $\varepsilon > 0$, we have

$$\mathcal{F}(x; \theta) > \mathcal{F}(x; \theta'), \quad (16)$$

where $\theta' = \theta - \eta \nabla \mathcal{L}_{nll}(x, y; \theta)$ i.e., supervised training with one step of gradient descent, and $\langle \cdot, \cdot \rangle$ denotes the inner product of gradients.

Experiments

We evaluate EADA against state-of-the-art approaches on various scenarios including a toy problem, three popular image classification datasets: **VisDA-2017** (Peng et al. 2017), **Office-Home** (Venkateswara et al. 2017) and **Office-31** (Saenko et al. 2010), as well as a challenging semantic segmentation task, i.e., **GTAV** (Richter et al. 2016) to **Cityscapes** (Cordts et al. 2016). All methods are implemented based on PyTorch, employing ResNet (He et al. 2016) models pre-trained on ImageNet (Krizhevsky, Sutskever, and Hinton 2012b). We follow the standard protocols for Active DA as (Su et al. 2020; Fu et al. 2021). Meanwhile, the various compared active learning, active domain adaptation and domain adaptation algorithms are **Source Only** (ResNet), **Random** (randomly select target samples to label), **BvSB** (Joshi, Porikli, and Papanikolopoulos 2009), **Entropy** (Wang and Shang 2014), **CoreSet** (Sener and Savarese 2018), **WAAL** (Shui et al. 2020), **BADGE** (Ash et al. 2020), **AADA** (Su et al. 2020), **DBAL** (de Mathelin, Mougeot, and Vayatis 2021), **TQS** (Fu et al. 2021), **CLUE** (Prabhu et al. 2021), **AdaptSegNet** (Tsai et al. 2018), and **PLCA** (Kang et al. 2020). Notably, we carry out experiments with five different random seeds and report the average accuracy. More details are presented in Appendix.

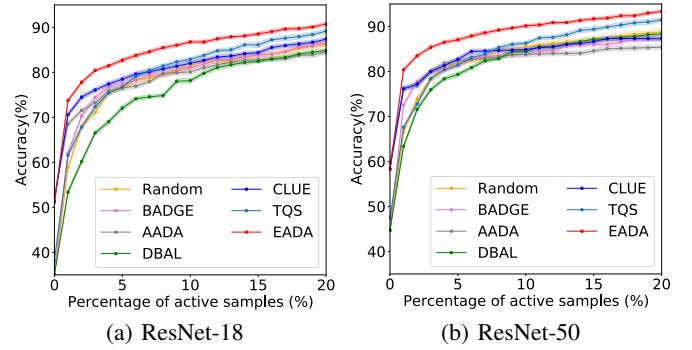


Figure 3: Comparison results of varying the percentage of labeled target samples on **VisDA-2017** with ResNet-18/50.

Main Results

VisDA-2017. The experimental results of different methods with 5% labeling budget on VisDA-2017 are shown in the first column in Table 1, proving that EADA is superior to all the baselines. Randomly selecting samples achieves better performance than ResNet, which implies that active learning is a promising and complementary solution for DA.

In addition, to further validate the effectiveness of EADA, we vary the target labeling budget from 0% to 20% with different backbones ResNet-18/50 and report the performance after each round in Fig. 3. We can observe that EADA consistently outperforms alternative methods across rounds. For instance, with shallower ResNet-18, we improve upon the state-of-the-art method, i.e., TQS by 2-6% over rounds, and obtain comparable results against other methods using deeper ResNet-50 at some rounds. This demonstrates

Method	VisDA-2017	Office-Home												
		Ar→Cl	Ar→Pr	Ar→Rw	Cl→Ar	Cl→Pr	Cl→Rw	Pr→Ar	Pr→Cl	Pr→Rw	Rw→Ar	Rw→Cl	Rw→Pr	Mean
Source Only	44.7 ± 0.1	42.1	66.3	73.3	50.7	59.0	62.6	51.9	37.9	71.2	65.2	42.6	76.6	58.3
Random	78.1 ± 0.6	52.5	74.3	77.4	56.3	69.7	68.9	57.7	50.9	75.8	70.0	54.6	81.3	65.8
BvSB	81.3 ± 0.4	56.3	78.6	79.3	58.1	74.0	70.9	59.5	52.6	77.2	71.2	56.4	84.5	68.2
Entropy	82.7 ± 0.3	58.0	78.4	79.1	60.5	73.0	72.6	60.4	54.2	77.9	71.3	58.0	83.6	68.9
CoreSet	81.9 ± 0.3	51.8	72.6	75.9	58.3	68.5	70.1	58.8	48.8	75.2	69.0	52.7	80.0	65.1
WAAL	83.9 ± 0.4	55.7	77.1	79.3	61.1	74.7	72.6	60.1	52.1	78.1	70.1	56.6	82.5	68.3
BADGE	84.3 ± 0.3	58.2	79.7	79.9	61.5	74.6	72.9	61.5	56.0	78.3	71.4	60.9	84.2	69.9
AADA	80.8 ± 0.4	56.6	78.1	79.0	58.5	73.7	71.0	60.1	53.1	77.0	70.6	57.0	84.5	68.3
DBAL	82.6 ± 0.3	58.7	77.3	79.2	61.7	73.8	73.3	62.6	54.5	78.1	72.4	59.9	84.3	69.6
TQS	83.1 ± 0.4	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
CLUE	85.2 ± 0.4	58.0	79.3	80.9	68.8	77.5	76.7	66.3	57.9	81.4	75.6	60.8	86.3	72.5
EADA	88.3 ± 0.1	63.6	84.4	83.5	70.7	83.7	80.5	73.0	63.5	85.2	78.4	65.4	88.6	76.7

Table 1: Comparison results on **VisDA-2017** and **Office-Home** with **5%** target samples as the labeling budget.

Method	A→D	A→W	D→A	D→W	W→A	W→D	Mean
Source Only	81.5	75.0	63.1	95.2	65.7	99.4	80.0
Random	87.1	84.1	75.5	98.1	75.8	99.6	86.7
BvSB	89.8	87.9	78.2	99.0	78.6	100.0	88.9
Entropy	91.0	89.2	76.1	99.7	77.7	100.0	88.9
CoreSet	82.5	81.1	70.3	96.5	72.4	99.6	83.7
WAAL	88.4	89.6	76.4	100.0	76.0	100.0	88.4
BADGE	90.8	89.1	79.8	99.6	79.6	100.0	89.8
AADA	89.2	87.3	78.2	99.5	78.7	100.0	88.8
DBAL	88.2	88.9	75.2	99.4	77.0	100.0	88.1
TQS	92.8	92.2	80.6	100.0	80.4	100.0	91.1
CLUE	92.0	87.3	79.0	99.2	79.6	99.8	89.5
EADA	97.7	96.6	82.1	100.0	82.8	100.0	93.2

Table 2: Comparison results on **Office-31** with **5%** target samples as the labeling budget.

that EADA can indeed select more representative and informative target data using our novel energy-based criterion. Additional comparison results with standard active learning methods are shown in Appendix.

Office-Home & Office-31. The results on Office-Home and Office-31 are reported in Table 1 & 2, respectively, which show the best performance across all tasks. It can be observed that most Active DA methods generally outperform the traditional active learning methods since the latter does not take the domain shift into account. EADA performs much better than all the baselines with a large margin, especially for very hard scenarios e.g., Ar→Cl, Pr→Cl, D→A and W→A, which emphasizes the benefit of jointly capturing domain characteristic and instance uncertainty for active sampling in combination with free energy alignment.

GTAV → Cityscapes. While prior works restrict their task to image classification, it is important to also study Active DA in the context of related tasks. Now we focus on task of semantic segmentation adapting from GTAV to Cityscapes and we use the same setting as (Tsai et al. 2018; Kang et al. 2020), which adopts DeepLab-v2 (Chen et al. 2018) with ResNet-101 as backbone network. We select 5% target images to query for pixel-level labels of the whole image. The experimental results are shown in Fig. 4(a). There is a large performance gap between the UDA methods and the full supervised version, such as a popular adversarial approach, AdaptSegNet, lags behind 25.2% mIoU. Surprisingly, our EADA brings a significant boost and shows performance comparable to that of fully supervised at the final round.

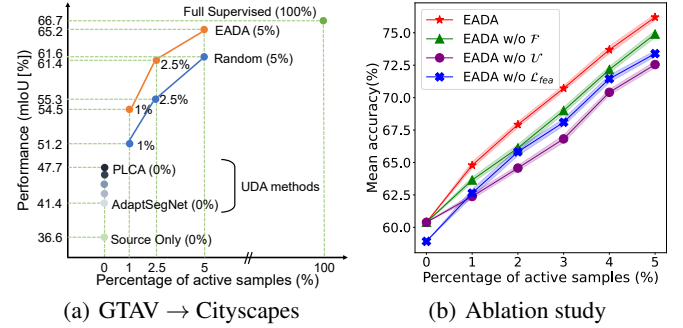


Figure 4: (a) Experimental results on GTAV→Cityscapes. (b) Mean accuracy of EADA and its variants on Office-Home.

Insight Analysis

Ablation study. To investigate the efficacy of key components of the proposed EADA, we conduct a thorough ablation study with the following variants on all 12 tasks of Office-Home: (i) EADA w/o \mathcal{F} : removing the free energy sampling from selection process; (ii) EADA w/o \mathcal{U} : removing the instance uncertainty sampling from the selection process; (iii) EADA w/o \mathcal{L}_{fea} : removing \mathcal{L}_{fea} from Eq. (9).

The results are shown in Fig. 4(b), it is clear that the full method outperforms other variants and achieves large improvements. We also observe that EADA surpasses EADA (w/o \mathcal{F} and w/o \mathcal{U}), manifesting that domain characteristic sampling and instance uncertain sampling are both necessary to select representative and informative data. Further, the consistent and notable increases from EADA w/o \mathcal{L}_{fea} to EADA justify our decision to use a regularization term to align free energy distributions between both domains, which is beneficial to reducing the domain shift implicitly.

Toy example. To better explain why the energy-based label acquisition strategy works well and what kind of sample is more representative and informative, we perform a toy example, a binary classification task with domain shift. As shown in Fig. 5, from the left to the right: Source Only, Random, BADGE, and our EADA are shown one by one and the target errors are 52.0%, 8.5%, 4.2%, 1.0%, respectively.

From the experimental results, we can make several insightful findings. (i) *Free energy biases*: the values of free

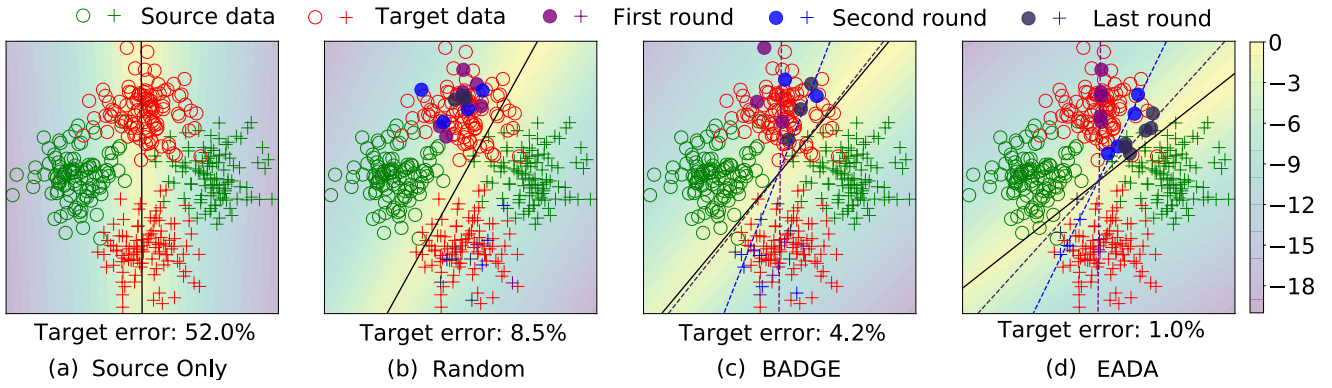


Figure 5: (Best viewed in color.) Illustrative comparison of sampling strategies on a toy example. Red points and green points denote unlabeled target data and labeled source data, respectively. Source data are drawn from two different Gaussian distributions denoted as circle class and plus class and target data are generated by rotating source data directly. We train a single layer fully-connected network and perform 3 rounds of active learning with a per-round budget is 2% of target samples. In (c) & (d), we draw the decision boundary before each selection round with a dash line, and the final decision boundary with a solid line.

α_1 / α_2 (%)	10 / 10	25 / 4	50 / 2	75 / 1.3	100 / 1	1 / 100
Office-31	91.9	92.6	93.2	91.5	90.8	90.4
Office-Home	76.2	76.3	76.7	76.0	74.7	72.6
VisDA-2017	87.2	87.6	88.3	87.1	86.6	85.7

Table 3: Effect of selection ratios.

energy on target samples are considerably higher than those on source samples in Fig. 5(a). Motivated by this, we design a free energy sampling as a surrogate measure to describe domain characteristic. (ii) *Redundant/Trivial selection*: in Fig. 5(b), we can observe that a large portion of samples selected by “Source Only” resides in an area where the target data density is high, leading to many redundant instances. BADGE (a state-of-the-art active learning method) runs a clustering scheme on “gradient embedding” to incorporate both uncertainty and diversity, which slightly mitigates the dilemma of redundancy. However, when we deeply study the relationship between decision boundary and the selected samples in each round, we find that BADGE still selects a few well-aligned samples and the selected samples are not the most uncertain samples of the current classifier. (iii) *Free energy versus decision boundary*: the final decision boundary is the area with the highest free energy. Accordingly, we explore a MvSM metric to precisely quantify the uncertainty of a target sample under the current model. The results in Fig. 5(d) validate the effectiveness of our method. In short, we define that a target sample with the highest free energy and located around the decision boundary serves as the most valuable, both representative and informative, sample.

Effect of selection ratios. In Table 3, we show the accuracy on three image classification benchmarks with varying α_1 (α_2). Our EADA can achieve consistent performance within a wide range. It is worth noting that excluding any step (α_1 or $\alpha_2 = 100$) will lead to a performance drop. We leave it as future work to explore other more complex combinations like self-adaptive α_1 and weighted calculation.

	AL Strategy	Query Complexity	Query Time (Ar→Cl, VisDA-2017)
cluster	CoreSet	$\mathcal{O}(DN^2)$	(0.1s, 1.3m)
	BADGE	$\mathcal{O}(CDN^2)$	(4.7s, 3.5m)
	DBAL	$\mathcal{O}(DN(M+N))$	(0.4s, 5.3m)
	CLUE	$\mathcal{O}(DN(N+TB))$	(0.5s, 2.9m)
rank	Entropy	$\mathcal{O}(N \log N)$	(0.04s, 2.1s)
	AADA	$\mathcal{O}(N \log N)$	(0.03s, 2.2s)
	TQS	$\mathcal{O}(N \log N)$	(0.04s, 1.7s)
	EADA	$\mathcal{O}(N \log N)$	(0.02s, 0.9s)

Table 4: Comparison results on query complexity and query time. C, M, N denote number of classes, source instances and target instances respectively. D denotes feature dimension, B is labeling budget, T denotes clustering rounds.

Time complexity. Table 4 lists the query complexity and query time for EADA and comparable baseline methods. BADGE and CLUE achieve better mean accuracy (see Table 1 and Table 2) but are slower due to a clustering step. Our EADA obtains the best accuracy and is significantly more efficient than the competitive baselines as well.

Conclusion

In this paper, we present Energy-based Active Domain Adaptation (EADA), an algorithm to tackle performance limitations of domain adaptation at minimal label cost. We propose a novel energy-based sampling strategy into domain adaptation, for the selection of limited target samples that are representative and informative. On top of that, we further explore a regularization term to implicitly diminish the domain gap. In addition, theoretical results about when and why EADA is expected to work are elaborated. Through our experiments, we demonstrate its effectiveness in various transfer scenarios. More generally, our work is but a small step toward alleviating the intensive workload of annotation. This offers encouraging evidence that there remains value to be explored to go beyond the fully supervised method.

Acknowledgments

This work was supported in part by the National Natural Science Foundation of China (No. 61902028) and the National Research and Development Program of China (No. 2019YQ1700).

References

- Ash, J. T.; Zhang, C.; Krishnamurthy, A.; Langford, J.; and Agarwal, A. 2020. Deep Batch Active Learning by Diverse, Uncertain Gradient Lower Bounds. In *ICLR*.
- Bachman, P.; Sordoni, A.; and Trischler, A. 2017. Learning Algorithms for Active Learning. In Precup, D.; and Teh, Y. W., eds., *ICML*, 301–310.
- Bickel, S.; Brückner, M.; and Scheffer, T. 2009. Discriminative Learning Under Covariate Shift. *J. Mach. Learn. Res.*, 10: 2137–2155.
- Bousmalis, K.; Silberman, N.; Dohan, D.; Erhan, D.; and Krishnan, D. 2017. Unsupervised Pixel-Level Domain Adaptation with Generative Adversarial Networks. In *CVPR*, 95–104.
- Chan, Y. S.; and Ng, H. T. 2007. Domain Adaptation with Active Learning for Word Sense Disambiguation. In *ACL*.
- Chattopadhyay, R.; Fan, W.; Davidson, I.; Panchanathan, S.; and Ye, J. 2013. Joint Transfer and Batch-mode Active Learning. In *ICML*, 253–261.
- Chen, L.; Papandreou, G.; Kokkinos, I.; Murphy, K.; and Yuille, A. L. 2018. DeepLab: Semantic Image Segmentation with Deep Convolutional Nets, Atrous Convolution, and Fully Connected CRFs. *IEEE Trans. Pattern Anal. Mach. Intell.*, 40(4): 834–848.
- Cordts, M.; Omran, M.; Ramos, S.; Rehfeld, T.; Enzweiler, M.; Benenson, R.; Franke, U.; Roth, S.; and Schiele, B. 2016. The Cityscapes dataset for semantic urban scene understanding. In *CVPR*, 3213–3223.
- Dasgupta, S. 2011. Two faces of active learning. *Theor. Comput. Sci.*, 412(19): 1767–1781.
- de Mathelin, A.; Mougeot, M.; and Vayatis, N. 2021. Discrepancy-Based Active Learning for Domain Adaptation. *CoRR*, abs/2103.03757.
- Fu, B.; Cao, Z.; Wang, J.; and Long, M. 2021. Transferable Query Selection for Active Domain Adaptation. In *CVPR*, 7272–7281.
- Gal, Y.; Islam, R.; and Ghahramani, Z. 2017. Deep Bayesian Active Learning with Image Data. In *ICML*, 1183–1192.
- Ganin, Y.; and Lempitsky, V. 2015. Unsupervised Domain Adaptation by Backpropagation. In *ICML*, 1180–1189.
- Gissin, D.; and Shalev-Shwartz, S. 2019. Discriminative Active Learning. *CoRR*, abs/1907.06347.
- Gong, B.; Shi, Y.; Sha, F.; and Grauman, K. 2012. Geodesic flow kernel for unsupervised domain adaptation. In *CVPR*, 2066–2073.
- Goodfellow, I.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A.; and Bengio, Y. 2014. Generative adversarial nets. In *NeurIPS*, 2672–2680.
- Grathwohl, W.; Wang, K.; Jacobsen, J.; Duvenaud, D.; Norouzi, M.; and Swersky, K. 2020. Your classifier is secretly an energy based model and you should treat it like one. In *ICLR*.
- Gretton, A.; Borgwardt, K. M.; Rasch, M.; Schölkopf, B.; and Smola, A. J. 2007. A kernel method for the two-sample problem. In *NeurIPS*, 513–520.
- Hanneke, S. 2014. Theory of disagreement-based active learning. *Found. Trends Mach. Learn.*, 7(2-3): 131–309.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *CVPR*, 770–778.
- Hoffman, J.; Tzeng, E.; Park, T.; Zhu, J.; Isola, P.; Saenko, K.; Efros, A. A.; and Darrell, T. 2018. CyCADA: Cycle-Consistent Adversarial Domain Adaptation. In *ICML*, 1994–2003.
- Joshi, A. J.; Porikli, F.; and Papanikolopoulos, N. 2009. Multi-class active learning for image classification. In *CVPR*, 2372–2379.
- Kang, G.; Wei, Y.; Yang, Y.; Zhuang, Y.; and Hauptmann, A. G. 2020. Pixel-Level Cycle Association: A New Perspective for Domain Adaptive Semantic Segmentation. In *NeurIPS*.
- Krizhevsky, A.; Sutskever, I.; and Hinton, G. E. 2012a. ImageNet Classification with Deep Convolutional Neural Networks. In *NeurIPS*, 1097–1105.
- Krizhevsky, A.; Sutskever, I.; and Hinton, G. E. 2012b. ImageNet Classification with Deep Convolutional Neural Networks. In *NeurIPS*, 1097–1105.
- LeCun, Y.; Chopra, S.; Hadsell, R.; Ranzato, M.; and Huang, F. 2006. A tutorial on energy-based learning. *Predicting structured data*, 1(0).
- Li, S.; Liu, C. H.; Lin, Q.; Wen, Q.; Su, L.; Huang, G.; and Ding, Z. 2021a. Deep Residual Correction Network for Partial Domain Adaptation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 43(7): 2329–2344.
- Li, S.; Liu, H. C.; Lin, Q.; Xie, B.; Ding, Z.; Huang, G.; and Tang, J. 2020. Domain Conditioned Adaptation Network. In *AAAI*, 11386–11393.
- Li, S.; Lv, F.; Xie, B.; Liu, C. H.; Liang, J.; and Qin, C. 2021b. Bi-Classifer Determinacy Maximization for Unsupervised Domain Adaptation. In *AAAI*, 8455–8464.
- Li, S.; Song, S.; Huang, G.; Ding, Z.; and Wu, C. 2018. Domain Invariant and Class Discriminative Feature Learning for Visual Domain Adaptation. *IEEE Trans. Image Process.*, 27(9): 4260–4273.
- Li, S.; Xie, M.; Gong, K.; Liu, C. H.; Wang, Y.; and Li, W. 2021c. Transferable Semantic Augmentation for Domain Adaptation. In *CVPR*, 11516–11525.
- Liu, W.; Wang, X.; Owens, J. D.; and Li, Y. 2020. Energy-based Out-of-distribution Detection. In *NeurIPS*.
- Long, M.; Cao, Y.; Cao, Z.; Wang, J.; and Jordan, M. I. 2019. Transferable Representation Learning with Deep Adaptation Networks. *IEEE Trans. Pattern Anal. Mach. Intell.*, 41(12): 3071–3085.

- Long, M.; Cao, Z.; Wang, J.; and Jordan, M. I. 2018. Conditional adversarial domain adaptation. In *NeurIPS*, 1647–1657.
- Long, M.; Zhu, H.; Wang, J.; and Jordan, M. I. 2017. Deep Transfer Learning with Joint Adaptation Networks. In *ICML*, 2208–2217.
- Pan, S. J.; and Yang, Q. 2010. A survey on transfer learning. *IEEE Trans. Knowl. Data Eng.*, 22(10): 1345–1359.
- Peng, X.; Usman, B.; Kaushik, N.; Hoffman, J.; Wang, D.; and Saenko, K. 2017. VisDA: The Visual Domain Adaptation Challenge. *CoRR*, abs/1710.06924.
- Prabhu, V.; Chandrasekaran, A.; Saenko, K.; and Hoffman, J. 2021. Active Domain Adaptation via Clustering Uncertainty-weighted Embeddings. In *ICCV*, 8505–8514.
- Prince, M. 2004. Does active learning work? A review of the research. *Journal of engineering education*, 93(3): 223–231.
- Rai, P.; Saha, A.; Daumé III, H.; and Venkatasubramanian, S. 2010. Domain adaptation meets active learning. In *Proceedings of the NAACL HLT 2010 Workshop on Active Learning for Natural Language Processing*, 27–32.
- Richter, S. R.; Vineet, V.; Roth, S.; and Koltun, V. 2016. Playing for data: Ground truth from computer games. In *ECCV*, 102–118.
- Ronneberger, O.; Fischer, P.; and Brox, T. 2015. U-Net: Convolutional Networks for Biomedical Image Segmentation. In *MICCAI*, 234–241.
- Saenko, K.; Kulis, B.; Fritz, M.; and Darrell, T. 2010. Adapting visual category models to new domains. In *ECCV*, 213–226.
- Saito, K.; Kim, D.; Sclaroff, S.; Darrell, T.; and Saenko, K. 2019. Semi-Supervised Domain Adaptation via Minimax Entropy. In *ICCV*, 8050–8058.
- Saito, K.; Watanabe, K.; Ushiku, Y.; and Harada, T. 2018. Maximum classifier discrepancy for unsupervised domain adaptation. In *CVPR*, 3723–3732.
- Schohn, G.; and Cohn, D. 2000. Less is More: Active Learning with Support Vector Machines. In *ICML*, 839–846.
- Sener, O.; and Savarese, S. 2018. Active Learning for Convolutional Neural Networks: A Core-Set Approach. In *ICLR*.
- Settles, B. 2009. Active learning literature survey.
- Shui, C.; Zhou, F.; Gagné, C.; and Wang, B. 2020. Deep Active Learning: Unified and Principled Method for Query and Training. In *AISTATS*, volume 108, 1308–1318.
- Sinha, S.; Ebrahimi, S.; and Darrell, T. 2019. Variational Adversarial Active Learning. In *ICCV*, 5971–5980.
- Su, J.; Tsai, Y.; Sohn, K.; Liu, B.; Maji, S.; and Chandraker, M. 2020. Active Adversarial Domain Adaptation. In *WACV*, 728–737.
- Teshima, T.; Sato, I.; and Sugiyama, M. 2020. Few-shot Domain Adaptation by Causal Mechanism Transfer. In *ICML*, 9458–9469.
- Tsai, Y.; Hung, W.; Schuster, S.; Sohn, K.; Yang, M.; and Chandraker, M. 2018. Learning to Adapt Structured Output Space for Semantic Segmentation. In *CVPR*, 7472–7481.
- Tzeng, E.; Hoffman, J.; Darrell, T.; and Saenko, K. 2015. Simultaneous Deep Transfer Across Domains and Tasks. In *ICCV*, 4068–4076.
- Tzeng, E.; Hoffman, J.; Saenko, K.; and Darrell, T. 2017. Adversarial discriminative domain adaptation. In *CVPR*, 7167–7176.
- Venkateswara, H.; Eusebio, J.; Chakraborty, S.; and Panchanathan, S. 2017. Deep hashing network for unsupervised domain adaptation. In *CVPR*, 5018–5027.
- Wang, D.; and Shang, Y. 2014. A new active labeling method for deep learning. In *IJCNN*, 112–119.
- Xu, R.; Li, G.; Yang, J.; and Lin, L. 2019. Larger Norm More Transferable: An Adaptive Feature Norm Approach for Unsupervised Domain Adaptation. In *ICCV*, 1426–1435.
- Yogamani, S. K.; Witt, C.; Rashed, H.; Nayak, S.; Mansoor, S.; Varley, P.; Perrotton, X.; O’Dea, D.; Pérez, P.; Hughes, C.; Horgan, J.; Sistu, G.; Chennupati, S.; Uricár, M.; Milz, S.; Simon, M.; and Amende, K. 2019. WoodScape: A Multi-Task, Multi-Camera Fisheye Dataset for Autonomous Driving. In *ICCV*, 9307–9317.
- Zou, H.; Yang, J.; and Wu, X. 2021. Unsupervised Energy-based Adversarial Domain Adaptation for Cross-domain Text Classification. In *ACL*, 1208–1218.