# UDALM: Unsupervised Domain Adaptation through Language Modeling

**Constantinos Karouzos[1], Georgios Paraskevopoulos[1,4], Alexandros Potamianos[1,2,3]**

[1] School of ECE, National Technical University of Athens, Athens, Greece
[2] Signal Analysis and Interpretation Laboratory (SAIL), USC, Los Angeles, CA, USA
[3] Behavioral Signal Technologies, Los Angeles, CA, USA
[4] Institute for Language and Speech Processing, Athena Research Center, Athens, Greece
ckarouzos@gmail.com, geopar@central.ntua.gr, potam@central.ntua.gr

## Abstract

In this work we explore Unsupervised Domain Adaptation (UDA) of pretrained language models for downstream tasks. We introduce UDALM, a fine-tuning procedure, using a mixed classification and Masked Language Model loss, that can adapt to the target domain distribution in a robust and sample efficient manner. Our experiments show that performance of models trained with the mixed loss scales with the amount of available target data and the mixed loss can be effectively used as a stopping criterion during UDA training. Furthermore, we discuss the relationship between A-distance and the target error and explore some limitations of the Domain Adversarial Training approach. Our method is evaluated on twelve domain pairs of the Amazon Reviews Sentiment dataset, yielding 91.74% accuracy, which is an 1.11% absolute improvement over the state-of-the-art.

## 1 Introduction

Deep architectures have achieved state-of-the-art results in a variety of machine learning tasks. However, real world deployments of machine learning systems often operate under domain shift, which leads to performance degradation. This introduces the need for adaptation techniques, where a model is trained with data from a specific domain, and then can be optimized for use in new settings. Efficient techniques for model re-usability can lead to faster and cheaper development of machine learning applications and facilitate their wider adoption. Especially techniques for Unsupervised Domain Adaptation (UDA) can have high real world impact, because they do not rely on expensive and time-consuming annotation processes to collect labeled data for domain-specific supervised training, further streamlining the process.

UDA approaches in the literature can be grouped in three major categories, namely pseudo-labeling techniques (e.g. Yarowsky, 1995; Zhou and Li, 2005), domain adversarial training (e.g. Ganin et al., 2016) and pivot-based approaches (e.g. Blitzer et al., 2006; Pan et al., 2010). Pseudo-labeling approaches use a model trained on the source labeled data to produce pseudo-labels for unlabeled target data and then train a model for the target domain in a supervised manner. Domain adversarial training aims to learn a domain-independent mapping for input samples by adding an adversarial cost during model training, that minimizes the distance between the source and target domain distributions. Pivot-based approaches aim to select domain-invariant features (pivots) and use them as a basis for cross-domain mapping. This work does not fall under any of these categories, rather we aim to optimize the fine-tuning procedure of pretrained language models (LMs) for learning under domain-shift.

Transfer learning from language models pretrained in massive corpora (Howard and Ruder, 2018; Devlin et al., 2019; Yang et al., 2019; Liu et al., 2019; Brown et al., 2020) has yielded significant improvements across a wide variety of NLP tasks, even when small amounts of data are used for fine-tuning. Fine-tuning a pretrained model is a straightforward framework for adaptation to target tasks and new domains, when labeled data are available. However, optimizing the fine-tuning process in UDA scenarios, where only labeled out-of-domain and unlabeled in-domain data are available is challenging.

In this work, we propose UDALM, a fine-tuning method for BERT (Devlin et al., 2019) in order to address the UDA problem. Our method is based on simultaneously learning the task from labeled data in the source distribution, while adapting to the language in the target distribution using multi-task learning. The key idea of our method is that by simultaneously minimizing a task-specific loss on the source data and a language modeling loss on the target data during fine-tuning, the model

will be able to adapt to the language of the target domain, while learning the supervised task from the available labeled data.

Our key contributions are: (a) We introduce UDALM, a novel, simple and robust unsupervised domain adaptation procedure for downstream BERT models based on multitask learning, (b) we achieve state-of-the-art results for the Amazon reviews benchmark dataset, surpassing more complicated approaches and (c) we explore how A-distance and the target error are related and conclude with some remarks on domain adversarial training, based on theoretical concepts and our empirical observations. Our code and models are publicly available[1].

## 2 Related Work

Traditionally, UDA has been performed using pseudo-labeling approaches. Pseudo-labeling techniques are semi-supervised algorithms that either use the same model (self-training) (Yarowsky, 1995; McClosky et al., 2006; Abney, 2007) or multiple ensembles of models (tri-training) (Zhou and Li, 2005; Søgaard, 2010) in order to produce pseudo-labels for the target unlabeled data. Saito et al. (2017) proposed an asymmetric tri-training approach. Ruder and Plank (2018) introduced a multi-task tri-training method. Rotman and Reichart (2019) and Lim et al. (2020) study pseudo-labeling with contextualized word representations. Ye et al. (2020) combine self-training with XLM-R (Conneau et al., 2020) to reduce the produced label noise and propose CFd, class aware feature self-distillation.

Another line of UDA research includes pivot-based methods, focusing on extracting cross-domain features. Structural Correspondence Learning (SCL) (Blitzer et al., 2006) and Spectral Feature Alignment (Pan et al., 2010) aim to find domain-invariant features (pivots) to learn a mapping between two domain distributions. Ziser and Reichart (2017, 2018, 2019) combine SCL with neural network architectures and language modeling. Miller (2019) propose to jointly learn the task and pivots. Li et al. (2018b) learn pivots with hierarchical attention networks. Pivot-based methods have also been used in conjunction with BERT (Ben-David et al., 2020).

Domain adversarial training is a dominant approach for UDA (Ramponi and Plank, 2020), in-

spired by the theory for learning from different domains introduced in Ben-David et al. (2007, 2010). Ganin et al. (2016); Ganin and Lempitsky (2015) propose to learn a task while not being able to distinguish if samples come from the source or the target distribution, through use of an adversarial cost. This approach has been adopted for a diverse set of problems, e.g. sentiment analysis, tweet classification and universal dependency parsing (Li et al., 2018a; Alam et al., 2018; Sato et al., 2017). Du et al. (2020) pose domain adversarial training in the context of BERT models. Zhao et al. (2018) propose multi-source domain adversarial networks. Guo et al. (2018) propose a mixture-of-experts approach for multi-source UDA. Guo et al. (2020) explore distance measures as additional losses and use them to construct dynamic multi-armed bandit controller for the source domains. Shen et al. (2018) learn domain invariant features via Wasserstein distance. Bousmalis et al. (2016) introduce domain seperation networks with private and shared encoders.

Unsupervised pretraining on domain-specific corpora can be an effective adaptation process. For example BioBERT (Lee et al., 2020) and SciBERT (Beltagy et al., 2019) are specialized BERT variants, where pretraining is extended on large amounts of biomedical and scientific corpora respectively. Sun et al. (2019) propose continuing the pretraining of BERT with target domain data and multitask learning using relevant tasks for BERT fine-tuning. Xu et al. (2019) introduce a review reading comprehension task and a post-training approach for BERT with an auxiliary loss on a question-answering task. Continuing pretraining on multiple phases, from general to domain specific (DAPT) and task specific data (TAPT), further improves performance of pretrained language models, as shown by Gururangan et al. (2020). Han and Eisenstein (2019) propose AdaptaBERT, which includes a second phase of unsupervised pretraining, in order to use BERT in a unsupervised domain adaptation context.

Recent works have highlighted the merits of using Language Modeling as an auxiliary task during fine-tuning. Chronopoulou et al. (2019) use an auxiliary LM loss to avoid catastrophic forgetting in transfer learning and Jia et al. (2019) adopt this approach for cross-domain named-entity recognition. We draw inspiration from these approaches and utilize auxiliary Language Modeling for UDA.
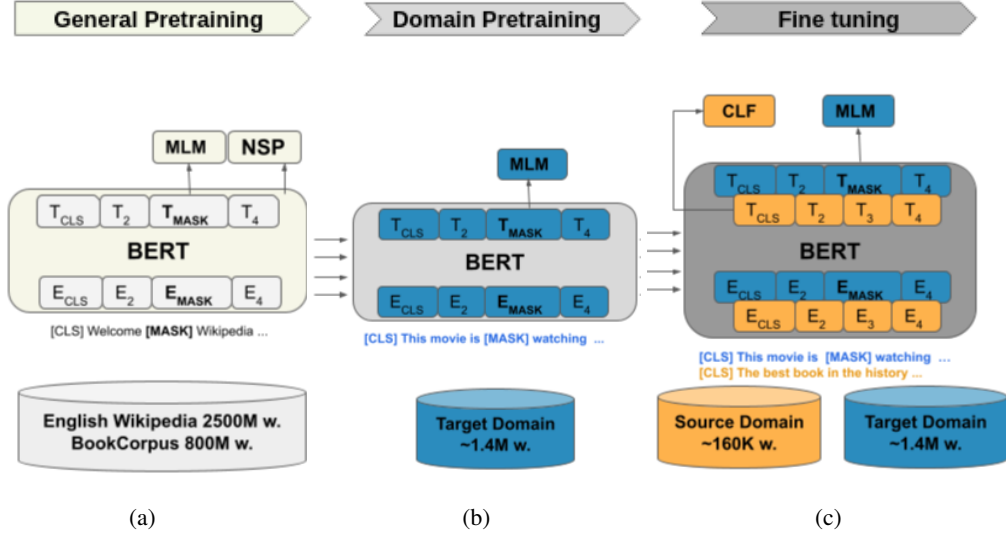
Figure 1: (a) BERT (Devlin et al., 2019) is pretrained on English Wikipedia and BookCorpus with the Masked Language Modeling (MLM) and the Next Sentence Prediction (NSP) tasks. (b) We continue the pretraining of BERT on unlabeled target domain data using the MLM task. (c) We train a task classifier with source domain labeled data, while we keep the MLM objective on unlabeled target domain data.

## 3 Problem Definition

Let $X$ be the input space and $Y$ the set of labels. For binary classification tasks $Y = \{0, 1\}$. In domain adaptation there are two different distributions over $X \times Y$, called the source domain $D_S$ and the target domain $D_T$. In the unsupervised setting labels are provided for samples drawn from $D_S$, while samples drawn from $D_T$ are unlabeled. The goal is to train a model that performs well on samples drawn from the target distribution $D_T$. This is summarized in Eq. 1.

$$S = (x_i, y_i)_{i=1}^{n} \sim (D_S)^n$$
$$T = (x_i)_{i=n+1}^{n+m} \sim (D_T^X)^m \qquad (1)$$

where $D_T^X$ is the marginal distribution of $D_T$ over $X$, $n$ is the number of samples from the source domain and $m$ is the number of samples from the target domain.

## 4 Proposed Method

Fig. 1 gives an overview of the proposed Unsupervised Domain Adaptation through Language Modeling (UDALM). Starting from a model that is pretrained in general corpora (Fig. 1a), we keep pretraining it on target domain data using the masked language modeling task (Fig. 1b). On the final fine-tuning step (Fig. 1c) we update the model weights using both a classification loss on the labeled source data and Masked Language Modeling loss on the unlabeled target data.

In Fig. 1a we see the BERT general pretraining phase. BERT (Devlin et al., 2019) is based on the Transformer architecture (Vaswani et al., 2017). During BERT pretraining, input tokens are randomly selected to be masked. BERT is trained using the Masked Language Modeling (MLM) objective, which consists of predicting the most probable tokens for the masked positions. Additionally it uses a Next Sentence Prediction (NSP) loss, which classifies whether the pair of input sentences are continuous or not. If a labeled dataset is available, a pretrained BERT model can be fine-tuned for the downstream task in a supervised manner with the addition of an output layer.

In Fig. 1b we initialize a model using the weights of a generally pretrained BERT and continue pretraining on an unsupervised set of in-domain data, in order to adapt to the target domain. This step does not require use of supervised data, since we use the MLM objective.

For the final fine-tuning step, shown in Fig. 1c we perform supervised fine-tuning on the source data, while we keep the MLM objective on the target data as an auxiliary task. Following standard practice, we use the `[CLS]` token representation for classification. The classifier consists of a single feed-forward layer.

During this procedure the model learns the task through the classification objective using the labeled source domain samples, and simultaneously

it adapts to the target domain data through the MLM objective. The model is trained on the source domain labeled data for the classification task and target domain unlabeled data for the masked language modeling task. We mask only the target domain data. During training we interleave source and target data and feed them to the BERT encoder. Features extracted from the source data are then used for classification, while target features are used for Masked Language Modeling.

The mixed loss used for the fine-tuning step, is the sum of the classification loss $L_{CLF}$ and the auxiliary MLM loss $L_{MLM}$. $L_{CLF}$ is a cross-entropy loss, calculated on labeled examples from source domain, while $L_{MLM}$ is used to predict masked tokens for unlabeled examples from target domain. We train the model over mixed batches, that include both source and target data, used for the respective tasks. The mixed loss is presented in Eq. 2:

$$L(\mathbf{s}, \mathbf{t}) = \lambda L_{CLF}(\mathbf{s}) + (1 - \lambda)L_{MLM}(\mathbf{t}) \quad (2)$$

We process $n$ labeled source samples $\mathbf{s} \sim D_S$ and $m$ unlabeled target samples $\mathbf{t} \sim D_T$ on a batch. The weighting factor $\lambda$ is selected as the ratio of labeled source data over the sum of labeled source and unlabeled target data, as stated in Eq. 3:

$$\lambda = \frac{n}{n + m} \quad (3)$$

## 5 Experiments

### 5.1 Dataset

We evaluate UDALM on the Amazon reviews multi-domain sentiment dataset (Blitzer et al., 2007), a standard benchmark dataset for domain adaptation. Reviews with one or two stars are labeled as negative, while reviews with four or five stars are labeled as positive. The dataset contains reviews on four product domains: *Books* (B), *DVDs* (D), *Electronics* (E) and *Kitchen appliances* (K), yielding 12 adaptation scenarios of source-target domain pairs. Balanced sets of 2000 labeled reviews are available for each domain. We use 20000 (randomly selected) unlabeled reviews for (B), (D) and (E). For (K) 17805 unlabeled reviews are available. For each of the 12 adaptation scenarios we use 20% of both labeled source and unlabeled target data for validation, while labeled target data are used for testing exclusively and are not seen during training or validation.

### 5.2 Implementation Details

We use $BERT_{BASE}$ (uncased) as the Language Model on which we apply domain pretraining. The $BERT_{BASE}$ original English model is a 12-layer, 768-hidden, 12-heads, 110M parameter transformer architecture, trained on the BookCorpus with 800M words and a version of the English Wikipedia with 2500M words. We convert source and target sentences to WordPieces (Wu et al., 2016). For target sentences we randomly mask 15% of WordPiece tokens, as in (Devlin et al., 2019). If a token in a specific position is selected to be masked 80% of the time is replaced with a `[MASK]` token, 10% of the time with a random token and 10% of the time remains unchanged.

The maximum sequence length is set to 512 by truncation of inputs. During domain pretraining we train with batch size of 8 for 3 epochs (2 hours on two GTX-1080Ti cards). During the final fine-tuning step of UDALM we train with batch size 36, consisting of $n = 1$ source sub-batch of 4 samples and $m = 8$ target sub-batches of 4 samples each. We update parameters after every 5 accumulated sub-batches. We train for 10 epochs with early stopping on the mixed loss in Eq. 2. For all experiments we use AdamW optimizer (Loshchilov and Hutter, 2018) with learning rate $10^{-5}$. Each adaptation scenario requires one hour on one GTX-1080Ti. For the domain adversarial experiments we set $\lambda_d = 0.01$ in Eq. 4 [2] and train for 10 epochs. Models are developed with PyTorch (Paszke et al., 2019) and HuggingFace Transformers (Wolf et al., 2019).

### 5.3 Baselines - Compared methods

We select three state-of-the-art methods for comparison. Each of the selected methods represents a different line of UDA research, namely domain adversarial training **BERT-DAAT** (Du et al., 2020), self-training XLM-R based **p+CFd** (Ye et al., 2020) and pivot-based **R-PERL** (Ben-David et al., 2020). We report results for the following settings with BERT models:

**Source only (SO)**: We fine-tune BERT on source domain labeled data, without using target data.

**Domain Pretraining (DPT)**: We use the target domain unlabeled data in order to continue pretraining of BERT with MLM loss (as in Fig. 1b) and then

---

[2]We also manually experimented with $\lambda_d = 1$ and $lambda_d = 0.1$, and a sigmoid schedule for $\lambda_d$. We report best results.

|  | R-PERL | DAAT | p+CFd | SO BERT | DAT BERT | DPT BERT | UDALM |
|---|---|---|---|---|---|---|---|
| $B \to D$ | 87.8 | 90.9 | 87.7 | $89.51 \pm 0.76$ | $87.31 \pm 2.14$ | $90.49 \pm 0.38$ | $\mathbf{90.97 \pm 0.22}$ |
| $B \to E$ | 87.2 | 88.9 | 91.3 | $90.51 \pm 0.51$ | $86.91 \pm 2.71$ | $90.38 \pm 1.59$ | $\mathbf{91.69 \pm 0.31}$ |
| $B \to K$ | 90.2 | 88.0 | 92.5 | $91.75 \pm 0.28$ | $90.59 \pm 1.17$ | $92.66 \pm 0.43$ | $\mathbf{93.21 \pm 0.22}$ |
| $D \to B$ | 85.6 | 89.7 | $\mathbf{91.5}$ | $90.26 \pm 0.64$ | $86.30 \pm 3.10$ | $91.02 \pm 0.75$ | $91.00 \pm 0.42$ |
| $D \to E$ | 89.3 | 90.1 | 91.6 | $88.71 \pm 1.48$ | $87.85 \pm 1.24$ | $91.03 \pm 0.82$ | $\mathbf{92.30 \pm 0.47}$ |
| $D \to K$ | 90.4 | 88.8 | 92.5 | $91.22 \pm 0.69$ | $89.95 \pm 1.53$ | $92.30 \pm 0.42$ | $\mathbf{93.66 \pm 0.37}$ |
| $E \to B$ | 90.2 | 89.6 | 88.7 | $87.96 \pm 0.89$ | $85.65 \pm 1.91$ | $88.52 \pm 0.55$ | $\mathbf{90.61 \pm 0.30}$ |
| $E \to D$ | 84.8 | $\mathbf{89.3}$ | 88.2 | $87.37 \pm 0.64$ | $83.99 \pm 1.31$ | $87.85 \pm 0.47$ | $88.83 \pm 0.61$ |
| $E \to K$ | 91.2 | 91.7 | 93.6 | $93.30 \pm 0.50$ | $92.45 \pm 1.35$ | $94.39 \pm 0.72$ | $\mathbf{94.43 \pm 0.24}$ |
| $K \to B$ | 83.0 | $\mathbf{90.8}$ | 89.8 | $88.15 \pm 0.64$ | $85.07 \pm 1.03$ | $88.83 \pm 0.81$ | $90.29 \pm 0.51$ |
| $K \to D$ | 85.6 | $\mathbf{90.5}$ | 87.8 | $87.23 \pm 0.49$ | $84.11 \pm 0.62$ | $88.52 \pm 0.69$ | $89.54 \pm 0.59$ |
| $K \to E$ | 91.2 | 93.2 | 92.6 | $93.23 \pm 0.34$ | $92.07 \pm 0.24$ | $93.42 \pm 0.40$ | $\mathbf{94.34 \pm 0.26}$ |
| Average | 87.50 | 90.12 | 90.63 | $89.93 \pm 0.65$ | $87.68 \pm 1.53$ | $90.78 \pm 0.67$ | $\mathbf{91.74 \pm 0.38}$ |

Table 1: Accuracy of unsupervised domain adaptation on twelve domain pairs of Amazon Reviews Multi Domain Sentiment Dataset.

fine-tune the resulting model on source domain labeled data.

**Domain Adversarial (DAT)**: Domain Adversarial Training with BERT. Starting from the domain pre-trained BERT (as in Fig. 1b), we then fine-tune the model with domain adversarial training as in Ganin et al. (2016). For a BERT model with parameters $\theta$, with $L_{CLF}$ being a cross-entropy loss for supervised task prediction, $L_{ADV}$ being a cross-entropy loss for domain prediction and $\lambda_d$ being a weighting factor, domain adversarial training consists of the minimization criterion described in Eq. 4.

$$\min_{\theta} L_{CLF}(\theta; D_S) - \lambda_d L_{ADV}(\theta; D_S, D_T) \quad (4)$$

**UDALM**: The proposed method, where we fine-tune the model created in the domain pretraining step using the mixed loss in Eq. 2.

## 6 Experimental Results

### 6.1 Comparison to state-of-the-art

We present results for all 12 domain adaptation settings in Table 1. Results for SO BERT, DAT BERT, DPT BERT and UDALM are averaged over five runs and we include standard deviations The last line of Table 1 contains the macro-averaged accuracy and deviations over all domain pairs. UDALM surpasses all other techniques, yielding an absolute improvement of $1.81\%$ over the SO BERT baseline. For fair comparison, we compare only with methods based on pretrained models, mostly BERT. We observe that BERT fine-tuned only with the source domain labeled data, without any knowledge of the target domain, is a competitive baseline. This source-only model even surpasses state-of-the-art methods developed for UDA, e.g. R-PERL (Ben-David et al., 2020).

We reproduce the domain adversarial training procedure and present results in the DAT BERT column of Table 1. Adversarial training proved to be unstable in our experiments, even after careful tuning of the adversarial loss weighting factor $\lambda_d$. This is evidenced by the high standard deviations in the DAT BERT experiments. We observe that adversarial training does not manage to outperform the source-only baseline.[3]

Domain pretraining increases the average accuracy with an absolute improvement of $0.85\%$ over the source-only baseline. Continuing MLM pre-training on the target domain data leads to better model adaptation, and therefore improved performance, on the target domain. This is consistent with previous works on supervised (Gururangan et al., 2020; Xu et al., 2019; Sun et al., 2019) and unsupervised settings (Han and Eisenstein, 2019; Du et al., 2020).

UDALM yields an additional $0.96\%$ absolute improvement of average accuracy over domain pre-training. Keeping the MLM loss during fine-tuning therefore, leads to better adaptation and acts as a regularizer that prevents the model from overfitting on the source domain. We also observe smaller standard deviations when using UDALM, which indicates that including the MLM loss during fine-tuning can result to more robust training.

### 6.2 Sample efficiency

UDALM surpasses in terms of macro-average accuracy all other approaches for unsupervised domain adaptation on the Amazon reviews multi-domain sentiment dataset. Specifically, our method improves on the state-of-the-art pseudo-labeling

---

[3]Note that we did not have to perform extensive tuning for the other methods, including UDALM.
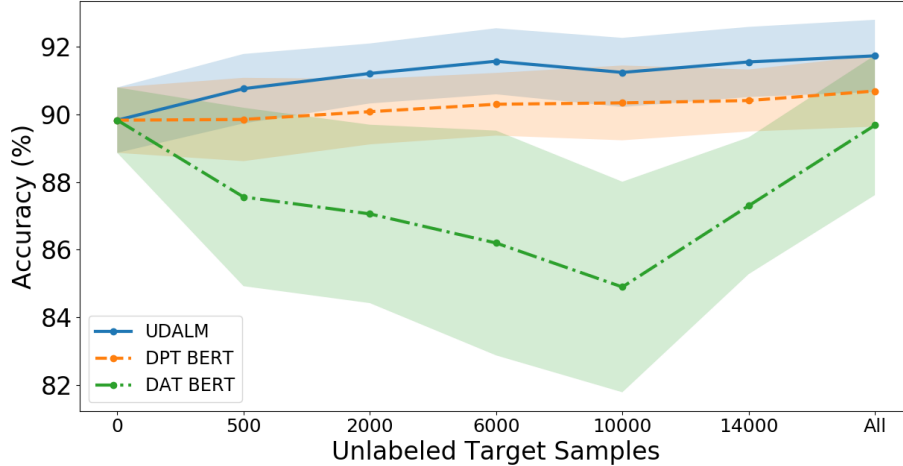
Figure 2: Average accuracy for different amount of target domain unlabeled samples of: (1) DPT BERT (2) DAT BERT and (3) UDALM.

(p+CFd Ye et al., 2020), domain adversarial (DAAT Du et al., 2020) and pivot-based (R-PERL Ben-David et al., 2020) approaches by $1.11\%$, $1.62\%$ and $4.24\%$ respectively.

We further investigate the impact of using different amount of target domain unlabeled data on model performance, to study the sample efficiency of UDALM. We experiment with settings of 500, 2000, 6000, 10000 and 14000 samples, by randomly limiting the number of unlabeled target domain data. For each setting we conduct three experiments with BERT models: (1) DPT, (2) DAT and (3) UDALM. When no target data are available, all methods are equivalent to a source only fine-tuned BERT. Again, we do not tune the hyper-parameters for DPT or UDALM. Fig. 2 shows the average accuracy on the twelve adaptation scenarios of the studied dataset. We see that UDALM produces robust performance improvement when we limit the amount of target data, indicating that it can be used in low-resource settings. However, training BERT in a domain adversarial manner shows instabilities. This is further discussed in Section 7.

### 6.3 On the stopping criteria for UDA training

A common problem when performing UDA is the lack of target labeled data that can be used for hyperparameter validation. For example, Ruder and Plank (2018) use a small set of labeled target data for validation, putting the problem in a semi-supervised setting. When training under a domain shift, optimization of model performance on the source data may not result to optimal performance for the target data.

To illustrate this, we examine if the minimization of the mixed loss can be used as a stopping criterion for UDA training. We compare five stopping criteria: (1) fixed training for 1 epoch, (2) fixed training for 3 epochs, (3) fixed training for 10 epochs, (4) stop when the minimum classification loss is reached for the source data and (5) stop when the minimum mixed loss ( Eq. 2) is reached. For (4) and (5) we train for 10 epochs with patience 3. We report average accuracy of the five stopping criteria over the twelve adaptation scenarios of Amazon Reviews dataset on Table 2. Training for a fixed number of 10 epochs and stopping when the minimum mixed loss perform best, yielding comparable accuracies of $91.75\%$ and $91.73\%$ respectively. Note that stopping when the minimum source loss stops the fine-tuning process too soon and does not allow the model to learn the target domain effectively. Overall, we observe that the mixed loss can be effectively used for early stopping, regularizing the model and alleviating the need for extensive search for the optimal number of training steps. This is an indication that the mixed loss could be used for model validation.

| Stopping Criterion | Epochs | Av. Acc. |
|---|---|---|
| Fixed | 1 | 90.98 |
| Fixed | 3 | 91.65 |
| Fixed | 10 | **91.75** |
| Min source loss | 10, patience 3 | 91.30 |
| Min mixed loss | 10, patience 3 | **91.74** |

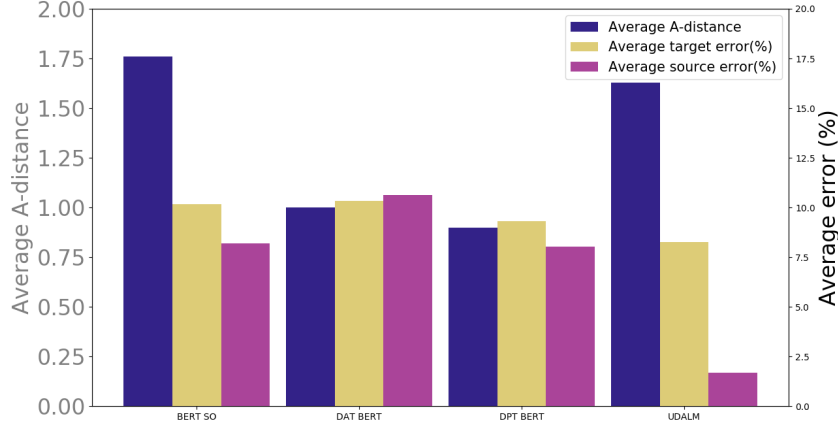Table 2: Comparison of average accuracy for various validation settings.

Figure 3: Comparison of average A-distance, average source error and average target error rate of different methods over all source - target pairs of the Amazon reviews dataset.

## 7 Discussion

### 7.1 Background Theory

Ben-David et al. (2007, 2010) provide a theory of learning from different domains. A key outcome of this work is the following theorem:

**Theorem** (Ben-David et al., 2007, 2010) Let $H$ be the hypothesis space and let $D_S, D_T$ be the two domains and $\epsilon_S, \epsilon_T$ be the corresponding error functions. Then for any $h \in H$:

$$\epsilon_T(h) \leq \epsilon_S(h) + \frac{1}{2}d_{H\Delta H}(D_S, D_T) + C \quad (5)$$

where $d_{H\Delta H}(D_S, D_T)$ is the $H\Delta H$-divergence (Kifer et al., 2004) between two domains, that is a measure of distance between domains that can be estimated from finite samples.

Eq. 5 defines an upper bound for the expected error $\epsilon_T(h)$ of a hypothesis $h$ on the target domain as the sum of three terms, namely the expected error on the source domain $\epsilon_S(h)$, the divergence between the source and target domain distributions $\frac{1}{2}d_{H\Delta H}(D_S, D_T)$ and the error of the ideal joint hypothesis $C$. When such an hypothesis exists, the term is considered relatively small and in practice ignored. The first term, bounds the expected error on the target domain by the expected error in the source domain and is expected to be small, due to supervised learning on the source domain. The second term, gives a notion of distance between the source and target domain extracted features. Intuitively this equation states: "if there exists a hypothesis $h$ that has small error on the source data and the source feature space is close to the target feature space, then this hypothesis will have low error

on the target data". Domain Adversarial Training aims to learn features that simultaneously result to low source error and low distance between target and source feature spaces based on the combined loss in Eq. 4.

### 7.2 A-distance only provides an upper bound for target error

According to Ben-David et al. (2007) the $H\Delta H$-divergence can be approximated by proxy A-distance, that is defined by Eq. 6 given the domain classification error $\epsilon_D$.

$$d_A = 2(1 - 2\epsilon_D) \quad (6)$$

We calculate an approximation of the distance between domains. Following prior work (Ganin et al., 2016; Saito et al., 2017) we create an SVM domain classifier. We feed the SVM with BERT's $[CLS]$ token representations, measure the domain classification error, and compute A-distance as in Eq. 6. We train the domain classifier on 2000 samples from each source and target domains. Fig. 3 shows the A-distance along with the source and the target error, averaged over the twelve available domain pairs using representations obtained from four methods, namely BERT SO, DAT BERT, DPT BERT and UDALM. DAT BERT minimizes the distance between domains. DPT BERT also reduces the A-distance, to similar levels with DAT, without using an explicit loss to minimize A-distance. To our surprise we found that, although it achieves the lowest error rate, UDALM does not significantly reduce the proxy A-distance compared to the source-only baseline. Additionally, we observe that the source error is correlated to model performance on the target task, i.e. models with lower source error

have also lower target error. UDALM specifically, achieves high accuracy on the source task and is able to transfer the task knowledge across domains, while DAT is able to bring domain representations closer, but at the cost of achieving weaker performance on the task at hand.

Overall, we do not observe a correlation between the resulting A-distance and model performance on target domain. Therefore, lower distance between domains, achieved intentionally or not, is not a necessary condition for good performance on the target domain[4], and our efforts could be better spent towards synergistic learning of the supervised source task and the target domain distribution.

### 7.3 Limitations of Domain Adversarial Training

Domain adversarial training (Ganin et al., 2016) faces some critical limitations that make the method difficult to be reproduced due to high hyper-parameter sensitivity and instability during training.

Such limitations have been highlighted by other authors in the UDA literature. For example, according to Shen et al. (2018) when a domain classifier can perfectly distinguish target from source representations, there will be a gradient vanishing problem. Shah et al. (2018) state that domain adversarial training is unstable and needs careful hyper-parameter tuning for their experiments. Wang et al. (2020) report results over three multi-domain NLP datasets, where domain adversarial training in conjunction with BERT under-performs. Ruder and Plank (2018) found that the domain adversarial loss did not help for their experiments on the Amazon reviews dataset.

In our experiments we note that domain-adversarial training results to worse performance than naive source only training. Furthermore, we experienced the need for extensive tuning of the $\lambda_d$ parameter from Eq. 4 every time the experimental setting changed (e.g. when testing for different amounts of available target data as in Section 6.2). This motivated us to further investigate the behavior of BERT fine-tuned with the adversarial cost. For visual inspection, we perform T-SNE (Maaten and Hinton, 2008) on representations extracted

from BERT, under four UDA setings in Fig. 4. In Fig. 4a we observe features extracted using BERT with Domain Adversarial Training and we compare it with features from SO BERT (Fig. 4b), DPT BERT (Fig. 4c) and UDALM (Fig. 4d). We observe that domain adversarial training manages to group tightly target and source samples, especially in the case of positive samples. Nevertheless, in the process, DAT introduces significant distortion in the semantic space, which is reflected in model performance[5].

We can attribute this behavior to two factors. First, The formulation of the adversarial loss in Eq. (4) can lead to trivial solutions. In order to maximize the $L_{ADV}$ term of Eq. (4), the model can just flip all domain labels, namely just predict that source samples belong to the target domain and vice-versa. In this case the model can still discriminate between domains and domain-independent representations are not encouraged. We empirically observed this behavior in our experiments with DAT, and only extensive hyper-parameter tuning could alleviate this issue. Additionally, Eq. (4) aims to minimize the upper bound of the target error $\epsilon_T(h)$ in Eq. (5). While this is desirable, reduction of the upper bound does not necessarily result in reduction of the bounded term in all scenarios. Furthermore, optimizing the $L_{ADV}(\theta; D_S, D_T)$ term can lead to increasing $L_{CLF}(\theta; D_S)$, and therefore one must find a balance between the two adversarial terms, again through careful hyper-parameter tuning. These issues could potentially be alleviated by including regularization terms that discourage trivial solutions and improve robustness. Therefore, given the lack of guarantees for good performance and the practical considerations, further investigation should be conducted regarding the robustness and reproducibility of DAT for UDA.

## 8 Conclusions and Future Work

Unsupervised domain adaptation of pretrained language models is a challenging problem with direct real world applications. In this work we propose UDALM, a robust, plug and play training strategy, which is able to improve performance in the target domain, achieving state-of-the-art results across 12 adaptation settings in the multi-domain Ama-

---

[4]Shu et al. (2018) state that feature distribution matching is a weak constraint when high-capacity feature extractors are used. Intuitively, a high-capacity feature extractor can perform arbitrary transformations to the input features in order to match the distributions.

[5]Note, we include this visualization for a single source-domain pair as an example. We performed multiple runs of T-SNE over all 12 source-domain pairs and this behavior appeared consistently.
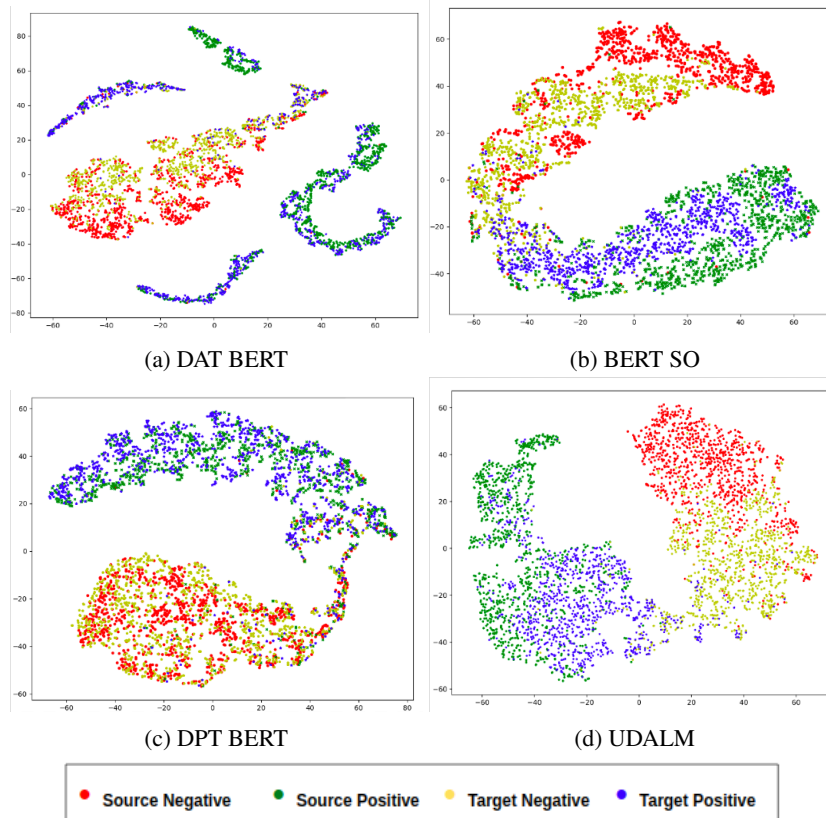
Figure 4: $2D$ representations of BERT $[CLS]$ features using t-SNE for the $D \to K$ task. The goal is to maximize separation between target positive (blue) and target negative (yellow) samples.

zon reviews dataset. Our method produces robust results with little hyper-parameter tuning and the proposed mixed-loss can be used for model validation, allowing for fast model development. Furthermore, UDALM scales with the amount of available unsupervised data from the target domain, allowing for adaptation in low-resource settings. In our analysis, we discuss the relationship between the A-distance and the target error. We observe that low A-distance may not suggest low target error for high capacity models. Additionally, we examine limitations of Domain Adversarial Training and highlight that the adversarial cost may lead to distortion of the feature space and negatively impact performance.

In the future we plan to apply UDALM to other tasks under domain-shift, such as sequence classification, question answering and part-of-speech tagging. Furthermore, we plan to extend our method for temporal and style adaptation, by adding more relevant auxiliary tasks that model language shift over time and over different platforms. Finally, we want to investigate the effectiveness of the proposed fine-tuning approach in supervised scenarios.

## References

Steven Abney. 2007. *Semisupervised learning for computational linguistics*. CRC Press.

Firoj Alam, Shafiq Joty, and Muhammad Imran. 2018. Domain adaptation with adversarial training and

graph embeddings. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1077–1087, Melbourne, Australia. Association for Computational Linguistics.

Iz Beltagy, Kyle Lo, and Arman Cohan. 2019. SciBERT: A pretrained language model for scientific text. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3615–3620, Hong Kong, China. Association for Computational Linguistics.

Eyal Ben-David, Carmel Rabinovitz, and Roi Reichart. 2020. Perl: Pivot-based domain adaptation for pre-trained deep contextualized embedding models. *Transactions of the Association for Computational Linguistics*, 8:504–521.

Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman Vaughan. 2010. A theory of learning from different domains. *Machine learning*, 79(1-2):151–175.

Shai Ben-David, John Blitzer, Koby Crammer, and Fernando Pereira. 2007. Analysis of representations for domain adaptation. In *Advances in neural information processing systems*, pages 137–144.

John Blitzer, Mark Dredze, and Fernando Pereira. 2007. Biographies, Bollywood, boom-boxes and blenders: Domain adaptation for sentiment classification. In *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics*, pages 440–447, Prague, Czech Republic. Association for Computational Linguistics.

John Blitzer, Ryan McDonald, and Fernando Pereira. 2006. Domain adaptation with structural correspondence learning. In *Proceedings of the 2006 Conference on Empirical Methods in Natural Language Processing*, pages 120–128, Sydney, Australia. Association for Computational Linguistics.

Konstantinos Bousmalis, George Trigeorgis, Nathan Silberman, Dilip Krishnan, and Dumitru Erhan. 2016. Domain separation networks. In *Advances in neural information processing systems*, pages 343–351.

Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*.

Alexandra Chronopoulou, Christos Baziotis, and Alexandros Potamianos. 2019. An embarrassingly simple approach for transfer learning from pre-trained language models. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2089–2095, Minneapolis, Minnesota. Association for Computational Linguistics.

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Chunning Du, Haifeng Sun, Jingyu Wang, Qi Qi, and Jianxin Liao. 2020. Adversarial and domain-aware BERT for cross-domain sentiment analysis. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4019–4028, Online. Association for Computational Linguistics.

Yaroslav Ganin and Victor Lempitsky. 2015. Unsupervised domain adaptation by backpropagation. volume 37 of *Proceedings of Machine Learning Research*, pages 1180–1189, Lille, France. PMLR.

Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. 2016. Domain-adversarial training of neural networks. *The Journal of Machine Learning Research*, 17(1):2096–2030.

Han Guo, Ramakanth Pasunuru, and Mohit Bansal. 2020. Multi-source domain adaptation for text classification via distancenet-bandits. In *AAAI*, pages 7830–7838.

Jiang Guo, Darsh Shah, and Regina Barzilay. 2018. Multi-source domain adaptation with mixture of experts. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4694–4703, Brussels, Belgium. Association for Computational Linguistics.

Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith. 2020. Don't stop pretraining: Adapt language models to domains and tasks. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8342–8360, Online. Association for Computational Linguistics.

Xiaochuang Han and Jacob Eisenstein. 2019. Unsupervised domain adaptation of contextualized embeddings for sequence labeling. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4238–4248, Hong Kong, China. Association for Computational Linguistics.

Jeremy Howard and Sebastian Ruder. 2018. Universal language model fine-tuning for text classification. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 328–339, Melbourne, Australia. Association for Computational Linguistics.

Chen Jia, Xiaobo Liang, and Yue Zhang. 2019. Cross-domain NER using cross-domain language modeling. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2464–2474, Florence, Italy. Association for Computational Linguistics.

Daniel Kifer, Shai Ben-David, and Johannes Gehrke. 2004. Detecting change in data streams. In *VLDB*, volume 4, pages 180–191. Toronto, Canada.

J Lee, W Yoon, S Kim, D Kim, S Kim, CH So, and J Kang. 2020. Biobert: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics (Oxford, England)*, 36(4):1234.

Yitong Li, Timothy Baldwin, and Trevor Cohn. 2018a. What's in a domain? learning domain-robust text representations using adversarial training. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 474–479, New Orleans, Louisiana. Association for Computational Linguistics.

Zheng Li, Ying Wei, Yu Zhang, and Qiang Yang. 2018b. Hierarchical attention transfer network for cross-domain sentiment classification. In *Thirty-Second AAAI Conference on Artificial Intelligence*.

KyungTae Lim, Jay Yoon Lee, Jaime Carbonell, and Thierry Poibeau. 2020. Semi-supervised learning on meta structure: Multi-task tagging and parsing in low-resource scenarios. 34:8344–8351.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.

Ilya Loshchilov and Frank Hutter. 2018. Decoupled weight decay regularization. In *International Conference on Learning Representations*.

Laurens van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov):2579–2605.

David McClosky, Eugene Charniak, and Mark Johnson. 2006. Reranking and self-training for parser adaptation. In *Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the Association for Computational Linguistics*, pages 337–344, Sydney, Australia. Association for Computational Linguistics.

Timothy Miller. 2019. Simplified neural unsupervised domain adaptation. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 414–419, Minneapolis, Minnesota. Association for Computational Linguistics.

Sinno Jialin Pan, Xiaochuan Ni, Jian-Tao Sun, Qiang Yang, and Zheng Chen. 2010. Cross-domain sentiment classification via spectral feature alignment. In *Proceedings of the 19th international conference on World wide web*, pages 751–760.

Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, volume 32, pages 8026–8037. Curran Associates, Inc.

Alan Ramponi and Barbara Plank. 2020. Neural unsupervised domain adaptation in nlp—a survey. *arXiv preprint arXiv:2005.14672*.

Guy Rotman and Roi Reichart. 2019. Deep contextualized self-training for low resource dependency parsing. *Transactions of the Association for Computational Linguistics*, 7:695–713.

Sebastian Ruder and Barbara Plank. 2018. Strong baselines for neural semi-supervised learning under domain shift. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1044–1054, Melbourne, Australia. Association for Computational Linguistics.

Kuniaki Saito, Yoshitaka Ushiku, and Tatsuya Harada. 2017. Asymmetric tri-training for unsupervised domain adaptation. volume 70 of *Proceedings of Machine Learning Research*, pages 2988–2997, International Convention Centre, Sydney, Australia. PMLR.

Motoki Sato, Hitoshi Manabe, Hiroshi Noji, and Yuji Matsumoto. 2017. Adversarial training for cross-domain Universal Dependency parsing. In *Proceedings of the CoNLL 2017 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*, pages 71–79, Vancouver, Canada. Association for Computational Linguistics.

Darsh Shah, Tao Lei, Alessandro Moschitti, Salvatore Romeo, and Preslav Nakov. 2018. Adversarial domain adaptation for duplicate question detection. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1056–1063, Brussels, Belgium. Association for Computational Linguistics.

Jian Shen, Yanru Qu, Weinan Zhang, and Yong Yu. 2018. Wasserstein distance guided representation learning for domain adaptation. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*, pages 4058–4065. AAAI Press.

Rui Shu, Hung Bui, Hirokazu Narui, and Stefano Ermon. 2018. A dirt-t approach to unsupervised domain adaptation. In *International Conference on Learning Representations*.

Anders Søgaard. 2010. Simple semi-supervised training of part-of-speech taggers. In *Proceedings of the ACL 2010 Conference Short Papers*, pages 205–208, Uppsala, Sweden. Association for Computational Linguistics.

Chi Sun, Xipeng Qiu, Yige Xu, and Xuanjing Huang. 2019. How to fine-tune bert for text classification? In *China National Conference on Chinese Computational Linguistics*, pages 194–206. Springer.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30, pages 5998–6008. Curran Associates, Inc.

Chengyu Wang, Minghui Qiu, Jun Huang, and Xiaofeng He. 2020. Meta fine-tuning neural language models for multi-domain text mining. *arXiv preprint arXiv:2003.13003*.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2019. Huggingface's transformers: State-of-the-art natural language processing. *ArXiv*, pages arXiv–1910.

Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. 2016. Google's neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*.

Hu Xu, Bing Liu, Lei Shu, and Philip Yu. 2019. BERT post-training for review reading comprehension and aspect-based sentiment analysis. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2324–2335, Minneapolis, Minnesota. Association for Computational Linguistics.

Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. In *Advances in neural information processing systems*, pages 5753–5763.

David Yarowsky. 1995. Unsupervised word sense disambiguation rivaling supervised methods. In *33rd Annual Meeting of the Association for Computational Linguistics*, pages 189–196, Cambridge, Massachusetts, USA. Association for Computational Linguistics.

Hai Ye, Qingyu Tan, Ruidan He, Juntao Li, Hwee Tou Ng, and Lidong Bing. 2020. Feature Adaptation of Pre-Trained Language Models across Languages and Domains for Text Classification. *arXiv:2009.11538 [cs]*. ArXiv: 2009.11538.

Han Zhao, Shanghang Zhang, Guanhang Wu, José M. F. Moura, Joao P Costeira, and Geoffrey J Gordon. 2018. Adversarial multiple source domain adaptation. In *Advances in Neural Information Processing Systems*, volume 31, pages 8559–8570. Curran Associates, Inc.

Zhi-Hua Zhou and Ming Li. 2005. Tri-training: Exploiting unlabeled data using three classifiers. *IEEE Transactions on knowledge and Data Engineering*, 17(11):1529–1541.

Yftah Ziser and Roi Reichart. 2017. Neural structural correspondence learning for domain adaptation. In *Proceedings of the 21st Conference on Computational Natural Language Learning (CoNLL 2017)*, pages 400–410, Vancouver, Canada. Association for Computational Linguistics.

Yftah Ziser and Roi Reichart. 2018. Pivot based language modeling for improved neural domain adaptation. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1241–1251, New Orleans, Louisiana. Association for Computational Linguistics.

Yftah Ziser and Roi Reichart. 2019. Task refinement learning for improved accuracy and stability of unsupervised domain adaptation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5895–5906, Florence, Italy. Association for Computational Linguistics.