

# Strong Baselines for Neural Semi-Supervised Learning under Domain Shift

Sebastian Ruder<sup>♣♣</sup> Barbara Plank<sup>♡◇</sup>

<sup>♣</sup>Insight Research Centre, National University of Ireland, Galway, Ireland

<sup>♣</sup>Aylien Ltd., Dublin, Ireland

<sup>♡</sup>Center for Language and Cognition, University of Groningen, The Netherlands

<sup>◇</sup>Department of Computer Science, IT University of Copenhagen, Denmark

sebastian@ruder.io, bplank@itu.dk

## Abstract

Novel neural models have been proposed in recent years for learning under domain shift. Most models, however, only evaluate on a single task, on proprietary datasets, or compare to weak baselines, which makes comparison of models difficult. In this paper, we re-evaluate classic general-purpose bootstrapping approaches in the context of neural networks under domain shifts vs. recent neural approaches and propose a novel *multi-task tri-training* method that reduces the time and space complexity of classic tri-training. Extensive experiments on two benchmarks are negative: while our novel method establishes a new state-of-the-art for sentiment analysis, it does not fare consistently the best. More importantly, we arrive at the somewhat surprising conclusion that classic tri-training, with some additions, outperforms the state of the art. We conclude that classic approaches constitute an important and strong baseline.

## 1 Introduction

Deep neural networks (DNNs) excel at learning from labeled data and have achieved state of the art in a wide array of supervised NLP tasks such as dependency parsing (Dozat and Manning, 2017), named entity recognition (Lample et al., 2016), and semantic role labeling (He et al., 2017).

In contrast, learning from unlabeled data, especially under domain shift, remains a challenge. This is common in many real-world applications where the distribution of the training and test data differs. Many state-of-the-art domain adaptation approaches leverage task-specific characteristics such as sentiment words (Blitzer et al., 2006; Wu and Huang, 2016) or distributional features (Schn-

abel and Schütze, 2014; Yin et al., 2015) which do not generalize to other tasks. Other approaches that are in theory more general only evaluate on proprietary datasets (Kim et al., 2017) or on a single benchmark (Zhou et al., 2016), which carries the risk of overfitting to the task. In addition, most models only compare against weak baselines and, strikingly, almost none considers evaluating against approaches from the extensive semi-supervised learning (SSL) literature (Chapelle et al., 2006).

In this work, we make the argument that such algorithms make strong baselines for any task in line with recent efforts highlighting the usefulness of classic approaches (Melis et al., 2017; Denkowski and Neubig, 2017). We re-evaluate bootstrapping algorithms in the context of DNNs. These are general-purpose semi-supervised algorithms that treat the model as a black box and can thus be used easily—with a few additions—with the current generation of NLP models. Many of these methods, though, were originally developed with in-domain performance in mind, so their effectiveness in a domain adaptation setting remains unexplored.

In particular, we re-evaluate three traditional bootstrapping methods, self-training (Yarowsky, 1995), tri-training (Zhou and Li, 2005), and tri-training with disagreement (Søgaard, 2010) for neural network-based approaches on *two* NLP tasks with different characteristics, namely, a sequence prediction and a classification task (POS tagging and sentiment analysis). We evaluate the methods across multiple domains on two well-established benchmarks, without taking any further task-specific measures, and compare to the best results published in the literature.

We make the somewhat surprising observation that classic tri-training outperforms task-agnostic state-of-the-art semi-supervised learning (Laine and Aila, 2017) and recent neural adaptation approaches (Ganin et al., 2016; Saito et al., 2017).

In addition, we propose *multi-task tri-training*, which reduces the main deficiency of tri-training, namely its time and space complexity. It establishes a new state of the art on unsupervised domain adaptation for sentiment analysis but it is outperformed by classic tri-training for POS tagging.

**Contributions** Our contributions are: a) We propose a novel multi-task tri-training method. b) We show that tri-training can serve as a strong and robust semi-supervised learning baseline for the current generation of NLP models. c) We perform an extensive evaluation of bootstrapping<sup>1</sup> algorithms compared to state-of-the-art approaches on two benchmark datasets. d) We shed light on the task and data characteristics that yield the best performance for each model.

## 2 Neural bootstrapping methods

We first introduce three classic bootstrapping methods, self-training, tri-training, and tri-training with disagreement and detail how they can be used with neural networks. For in-depth details we refer the reader to (Abney, 2007; Chapelle et al., 2006; Zhu and Goldberg, 2009). We introduce our novel multi-task tri-training method in §2.3.

### 2.1 Self-training

Self-training (Yarowsky, 1995; McClosky et al., 2006b) is one of the earliest and simplest bootstrapping approaches. In essence, it leverages the model’s own predictions on unlabeled data to obtain additional information that can be used during training. Typically the most confident predictions are taken at face value, as detailed next.

Self-training trains a model  $m$  on a labeled training set  $L$  and an unlabeled data set  $U$ . At each iteration, the model provides predictions  $m(x)$  in the form of a probability distribution over classes for all unlabeled examples  $x$  in  $U$ . If the probability assigned to the most likely class is higher than a predetermined threshold  $\tau$ ,  $x$  is added to the labeled examples with  $p(x) = \arg \max m(x)$  as pseudo-label. This instantiation is the most widely used and shown in Algorithm 1.

**Calibration** It is well-known that output probabilities in neural networks are poorly calibrated (Guo et al., 2017). Using a fixed threshold  $\tau$  is thus

<sup>1</sup>We use the term bootstrapping as used in the semi-supervised learning literature (Zhu, 2005), which should not be confused with the statistical procedure of the same name (Efron and Tibshirani, 1994).

---

### Algorithm 1 Self-training (Abney, 2007)

---

```

1: repeat
2:    $m \leftarrow \text{train\_model}(L)$ 
3:   for  $x \in U$  do
4:     if  $\max m(x) > \tau$  then
5:        $L \leftarrow L \cup \{(x, p(x))\}$ 
6: until no more predictions are confident

```

---

not the best choice. While the *absolute* confidence value is inaccurate, we can expect that the *relative* order of confidences is more robust.

For this reason, we select the top  $n$  unlabeled examples that have been predicted with the highest confidence after every epoch and add them to the labeled data. This is one of the many variants for self-training, called *throttling* (Abney, 2007). We empirically confirm that this outperforms the classic selection in our experiments.

**Online learning** In contrast to many classic algorithms, DNNs are trained online by default. We compare training setups and find that training until convergence on labeled data and then training until convergence using self-training performs best.

Classic self-training has shown mixed success. In parsing it proved successful only with small datasets (Reichart and Rappoport, 2007) or when a generative component is used together with a reranker in high-data conditions (McClosky et al., 2006b; Suzuki and Isozaki, 2008). Some success was achieved with careful task-specific data selection (Petrov and McDonald, 2012), while others report limited success on a variety of NLP tasks (Plank, 2011; Van Asch and Daelemans, 2016; van der Goot et al., 2017). Its main downside is that the model is not able to correct its own mistakes and errors are amplified, an effect that is increased under domain shift.

### 2.2 Tri-training

Tri-training (Zhou and Li, 2005) is a classic method that reduces the bias of predictions on unlabeled data by utilizing the agreement of three independently trained models. Tri-training (cf. Algorithm 2) first trains three models  $m_1$ ,  $m_2$ , and  $m_3$  on bootstrap samples of the labeled data  $L$ . An unlabeled data point is added to the training set of a model  $m_i$  if the other two models  $m_j$  and  $m_k$  agree on its label. Training stops when the classifiers do not change anymore.

Tri-training *with disagreement* (Søgaard, 2010)

---

**Algorithm 2** Tri-training (Zhou and Li, 2005)

---

```
1: for  $i \in \{1..3\}$  do
2:    $S_i \leftarrow \text{bootstrap\_sample}(L)$ 
3:    $m_i \leftarrow \text{train\_model}(S_i)$ 
4: repeat
5:   for  $i \in \{1..3\}$  do
6:      $L_i \leftarrow \emptyset$ 
7:     for  $x \in U$  do
8:       if  $p_j(x) = p_k(x) (j, k \neq i)$  then
9:          $L_i \leftarrow L_i \cup \{(x, p_j(x))\}$ 
10:         $m_i \leftarrow \text{train\_model}(L \cup L_i)$ 
10: until none of  $m_i$  changes
11: apply majority vote over  $m_i$ 
```

---

is based on the intuition that a model should only be strengthened in its weak points and that the labeled data should not be skewed by easy data points. In order to achieve this, it adds a simple modification to the original algorithm (altering line 8 in Algorithm 2), requiring that for an unlabeled data point on which  $m_j$  and  $m_k$  agree, the other model  $m_i$  disagrees on the prediction. Tri-training with disagreement is more data-efficient than tri-training and has achieved competitive results on part-of-speech tagging (Søgaard, 2010).

**Sampling unlabeled data** Both tri-training and tri-training with disagreement can be very expensive in their original formulation as they require to produce predictions for each of the three models on all unlabeled data samples, which can be in the millions in realistic applications. We thus propose to sample a number of unlabeled examples at every epoch. For all traditional bootstrapping approaches we sample 10k candidate instances in each epoch. For the neural approaches we use a linearly growing candidate sampling scheme proposed by (Saito et al., 2017), increasing the candidate pool size as the models become more accurate.

**Confidence thresholding** Similar to self-training, we can introduce an additional requirement that pseudo-labeled examples are only added if the probability of the prediction of at least one model is higher than some threshold  $\tau$ . We did not find this to outperform prediction without threshold for traditional tri-training, but thresholding proved essential for our method (§2.3).

The most important condition for tri-training and tri-training with disagreement is that the models are diverse. Typically, bootstrap samples are used

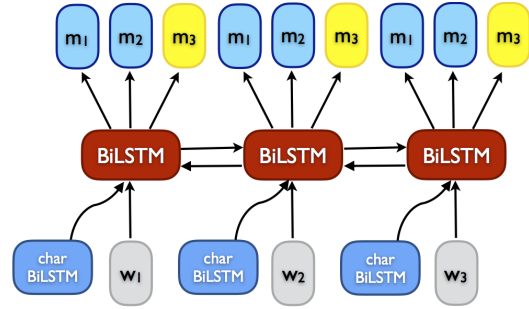


Figure 1: Multi-task tri-training (MT-Tri).

to create this diversity (Zhou and Li, 2005; Søgaard, 2010). However, training separate models on bootstrap samples of a potentially large amount of training data is expensive and takes a lot of time. This drawback motivates our approach.

### 2.3 Multi-task tri-training

In order to reduce both the time and space complexity of tri-training, we propose Multi-task Tri-training (MT-Tri). MT-Tri leverages insights from multi-task learning (MTL) (Caruana, 1993) to share knowledge across models and accelerate training. Rather than storing and training each model separately, we propose to share the parameters of the models and train them jointly using MTL.<sup>2</sup> All models thus collaborate on learning a joint representation, which improves convergence.

The output softmax layers are model-specific and are only updated for the input of the respective model. We show the model in Figure 1 (as instantiated for POS tagging). As the models leverage a joint representation, we need to ensure that the features used for prediction in the softmax layers of the different models are as diverse as possible, so that the models can still learn from each other’s predictions. In contrast, if the parameters in all output softmax layers were the same, the method would degenerate to self-training.

To guarantee diversity, we introduce an orthogonality constraint (Bousmalis et al., 2016) as an additional loss term, which we define as follows:

$$\mathcal{L}_{orth} = \|W_{m_1}^\top W_{m_2}\|_F^2 \quad (1)$$

where  $\|\cdot\|_F^2$  is the squared Frobenius norm and  $W_{m_1}$  and  $W_{m_2}$  are the softmax output parameters

<sup>2</sup>Note: we use the term multi-task learning here albeit all tasks are of the same kind, similar to work on multi-lingual modeling treating each language (but same label space) as separate task e.g., (Fang and Cohn, 2017). It is interesting to point out that our model is further doing implicit multi-view learning by way of the orthogonality constraint.

of the two source and pseudo-labeled output layers  $m_1$  and  $m_2$ , respectively. The orthogonality constraint encourages the models not to rely on the same features for prediction. As enforcing pairwise orthogonality between three matrices is not possible, we only enforce orthogonality between the softmax output layers of  $m_1$  and  $m_2$ ,<sup>3</sup> while  $m_3$  is gradually trained to be more target-specific. We parameterize  $\mathcal{L}_{orth}$  by  $\gamma=0.01$  following (Liu et al., 2017). We do not further tune  $\gamma$ .

More formally, let us illustrate the model by taking the sequence prediction task (Figure 1) as illustration. Given an utterance with labels  $y_1, \dots, y_n$ , our Multi-task Tri-training loss consists of three task-specific ( $m_1, m_2, m_3$ ) tagging loss functions (where  $\vec{h}$  is the uppermost Bi-LSTM encoding):

$$\mathcal{L}(\theta) = - \sum_i \sum_{1, \dots, n} \log P_{m_i}(y|\vec{h}) + \gamma \mathcal{L}_{orth} \quad (2)$$

In contrast to classic tri-training, we can train the multi-task model with its three model-specific outputs jointly and *without* bootstrap sampling on the labeled source domain data until convergence, as the orthogonality constraint enforces different representations between models  $m_1$  and  $m_2$ . From this point, we can leverage the pair-wise agreement of two output layers to add pseudo-labeled examples as training data to the third model. We train the third output layer  $m_3$  only on pseudo-labeled target instances in order to make tri-training more robust to a domain shift. For the final prediction, majority voting of all three output layers is used, which resulted in the best instantiation, together with confidence thresholding ( $\tau = 0.9$ , except for high-resource POS where  $\tau = 0.8$  performed slightly better). We also experimented with using a domain-adversarial loss (Ganin et al., 2016) on the jointly learned representation, but found this not to help. The full pseudo-code is given in Algorithm 3.

**Computational complexity** The motivation for MT-Tri was to reduce the space and time complexity of tri-training. We thus give an estimate of its efficiency gains. MT-Tri is  $\sim 3\times$  more space-efficient than regular tri-training; tri-training stores one set of parameters for each of the three models, while MT-Tri only stores one set of parameters (we use three output layers, but these make up a comparatively small part of the total parameter budget). In terms of time efficiency, tri-training first

<sup>3</sup>We also tried enforcing orthogonality on a hidden layer rather than the output layer, but this did not help.

---

### Algorithm 3 Multi-task Tri-training

---

```

1:  $m \leftarrow \text{train\_model}(L)$ 
2: repeat
3:   for  $i \in \{1..3\}$  do
4:      $L_i \leftarrow \emptyset$ 
5:     for  $x \in U$  do
6:       if  $p_j(x) = p_k(x) (j, k \neq i)$  then
7:          $L_i \leftarrow L_i \cup \{(x, p_j(x))\}$ 
8:       if  $i = 3$  then  $m_i = \text{train\_model}(L_i)$ 
9:       else  $m_i \leftarrow \text{train\_model}(L \cup L_i)$ 
10: until end condition is met
11: apply majority vote over  $m_i$ 

```

---

requires to train each of the models from scratch. The actual tri-training takes about the same time as training from scratch and requires a separate forward pass for each model, effectively training three independent models simultaneously. In contrast, MT-Tri only necessitates one forward pass as well as the evaluation of the two additional output layers (which takes a negligible amount of time) and requires about as many epochs as tri-training until convergence (see Table 3, second column) while adding fewer unlabeled examples per epoch (see Section 3.4). In our experiments, MT-Tri trained about  $5\text{-}6\times$  faster than traditional tri-training.

MT-Tri can be seen as a self-ensembling technique, where different variations of a model are used to create a stronger ensemble prediction. Recent approaches in this line are *snapshot ensembling* (Huang et al., 2017) that ensembles models converged to different minima during a training run, *asymmetric tri-training* (Saito et al., 2017) (ASYM) that leverages agreement on two models as information for the third, and *temporal ensembling* (Laine and Aila, 2017), which ensembles predictions of a model at different epochs. We tried to compare to temporal ensembling in our experiments, but were not able to obtain consistent results.<sup>4</sup> We compare to the closest most recent method, asymmetric tri-training (Saito et al., 2017). It differs from ours in two aspects: a) ASYM leverages only pseudo-labels from data points on which  $m_1$  and  $m_2$  agree, and b) it uses only one task ( $m_3$ ) as final predictor. In essence, our formulation of MT-Tri is closer to the original tri-training formulation (agreements on two provide pseudo-labels to the third) thereby incorporating more diversity.

<sup>4</sup>We suspect that the sparse features in NLP and the domain shift might be detrimental to its unsupervised consistency loss.



	Domain	# labeled	# unlabeled
POS tagging	Answers	3,489	27,274
	Emails	4,900	1,194,173
	Newsgroups	2,391	1,000,000
	Reviews	3,813	1,965,350
	Weblogs	2,031	524,834
	WSJ	30,060	100,000
Sentiment	Book	2,000	4,465
	DVD	2,000	3,586
	Electronics	2,000	5,681
	Kitchen	2,000	5,945

Table 1: Number of labeled and unlabeled sentences for each domain in the SANCL 2012 dataset (Petrov and McDonald, 2012) for POS tagging (above) and the Amazon Reviews dataset (Blitzer et al., 2006) for sentiment analysis (below).

### 3 Experiments

In order to ascertain which methods are robust across different domains, we evaluate on two widely used unsupervised domain adaptation datasets for two tasks, a sequence labeling and a classification task, cf. Table 1 for data statistics.

#### 3.1 POS tagging

For POS tagging we use the SANCL 2012 shared task dataset (Petrov and McDonald, 2012) and compare to the top results in both low and high-data conditions (Schnabel and Schütze, 2014; Yin et al., 2015). Both are strong baselines, as the FLORS tagger has been developed for this challenging dataset and it is based on contextual distributional features (excluding the word’s identity), and hand-crafted suffix and shape features (including some language-specific morphological features). We want to gauge to what extent we can adopt a nowadays fairly standard (but more lexicalized) general neural tagger.

Our POS tagging model is a state-of-the-art Bi-LSTM tagger (Plank et al., 2016) with word and 100-dim character embeddings. Word embeddings are initialized with the 100-dim Glove embeddings (Pennington et al., 2014). The BiLSTM has one hidden layer with 100 dimensions. The base POS model is trained on WSJ with early stopping on the WSJ development set, using patience 2, Gaussian noise with  $\sigma = 0.2$  and word dropout with  $p = 0.25$  (Kiperwasser and Goldberg, 2016).

Regarding data, the source domain is the Ontonotes 4.0 release of the Penn treebank Wall Street Journal (WSJ) annotated for 48 fine-grained POS tags. This amounts to 30,060 labeled sen-

tences. We use 100,000 WSJ sentences from 1988 as unlabeled data, following Schnabel and Schütze (2014).<sup>5</sup> As target data, we use the five SANCL domains (answers, emails, newsgroups, reviews, weblogs). We restrict the amount of unlabeled data for each SANCL domain to the first 100k sentences, and do not do any pre-processing. We consider the development set of ANSWERS as our only target dev set to set hyperparameters. This may result in suboptimal per-domain settings but better resembles an unsupervised adaptation scenario.

#### 3.2 Sentiment analysis

For sentiment analysis, we evaluate on the Amazon reviews dataset (Blitzer et al., 2006). Reviews with 1 to 3 stars are ranked as negative, while reviews with 4 or 5 stars are ranked as positive. The dataset consists of four domains, yielding 12 adaptation scenarios. We use the same pre-processing and architecture as used in (Ganin et al., 2016; Saito et al., 2017): 5,000-dimensional tf-idf weighted unigram and bigram features as input; 2k labeled source samples and 2k unlabeled target samples for training, 200 labeled target samples for validation, and between 3k-6k samples for testing. The model is an MLP with one hidden layer with 50 dimensions, sigmoid activations, and a softmax output. We compare against the Variational Fair Autoencoder (VFAE) (Louizos et al., 2015) model and domain-adversarial neural networks (DANN) (Ganin et al., 2016).

#### 3.3 Baselines

Besides comparing to the top results published on both datasets, we include the following baselines:

- a) the task model trained on the source domain;
- b) self-training (Self);
- c) tri-training (Tri);
- d) tri-training with disagreement (Tri-D); and
- e) asymmetric tri-training (Saito et al., 2017).

Our proposed model is multi-task tri-training (MT-Tri). We implement our models in DyNet (Neubig et al., 2017). Reporting single evaluation scores might result in biased results (Reimers and Gurevych, 2017). Throughout the paper, we report mean accuracy and standard deviation over five runs for POS tagging and over ten runs for

<sup>5</sup>Note that our unlabeled data might slightly differ from theirs. We took the first 100k sentences from the 1988 WSJ dataset from the BLLIP 1987-89 WSJ Corpus Release 1.

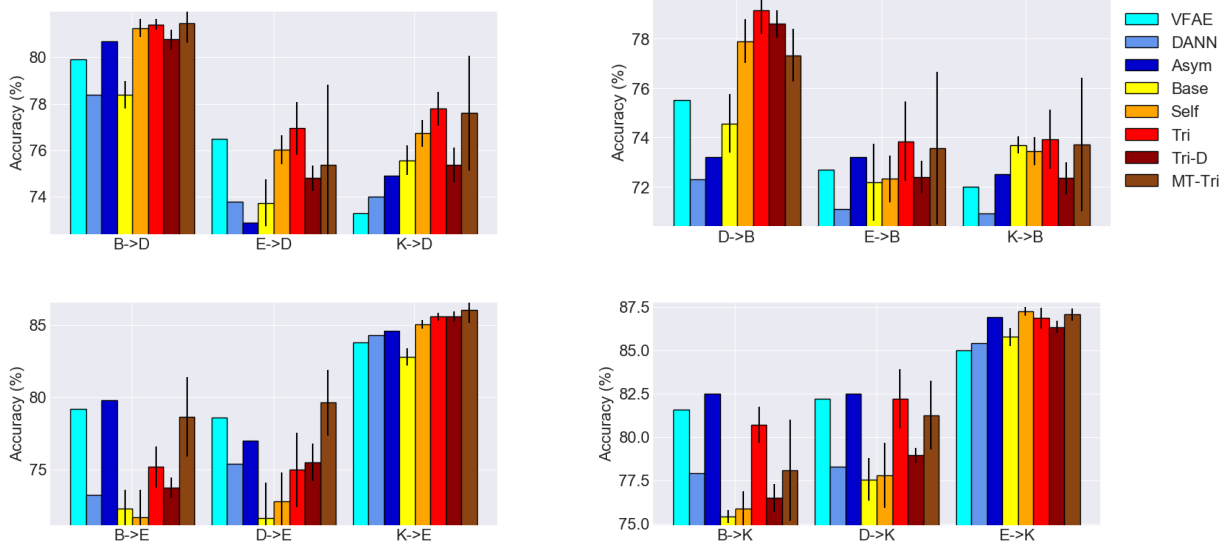


Figure 2: Average results for unsupervised domain adaptation on the Amazon dataset. Domains: B (Book), D (DVD), E (Electronics), K (Kitchen). Results for VFAE, DANN, and Asym are from Saito et al. (2017).

sentiment analysis. Significance is computed using bootstrap test. The code for all experiments is released at: <https://github.com/bplank/semi-supervised-baselines>.

### 3.4 Results

**Sentiment analysis** We show results for sentiment analysis for all 12 domain adaptation scenarios in Figure 2. For clarity, we also show the accuracy scores averaged across each target domain as well as a global macro average in Table 2.

Model	D	B	E	K	Avg
VFAE*	76.57	73.40	80.53	82.93	78.36
DANN*	75.40	71.43	77.67	80.53	76.26
Asym*	76.17	72.97	80.47	<b>83.97</b>	78.39
Src	75.91	73.47	75.61	79.58	76.14
Self	78.00	74.55	76.54	80.30	77.35
Tri	<b>78.72</b>	<b>75.64</b>	78.60	83.26	79.05
Tri-D	76.99	74.44	78.30	80.59	77.58
MT-Tri	78.14	74.86	<b>81.45</b>	82.14	<b>79.15</b>

Table 2: Average accuracy scores for each SA target domain. \*: result from Saito et al. (2017).

Self-training achieves surprisingly good results but is not able to compete with tri-training. Tri-training with disagreement is only slightly better than self-training, showing that the disagreement component might not be useful when there is a strong domain shift. Tri-training achieves the best

average results on two target domains and clearly outperforms the state of the art on average.

MT-Tri finally outperforms the state of the art on 3/4 domains, and even slightly traditional tri-training, resulting in the overall best method. This improvement is mainly due to the B->E and D->E scenarios, on which tri-training struggles. These domain pairs are among those with the highest  $\mathcal{A}$ -distance (Blitzer et al., 2007), which highlights that tri-training has difficulty dealing with a strong shift in domain. Our method is able to mitigate this deficiency by training one of the three output layers only on pseudo-labeled target domain examples.

In addition, MT-Tri is more efficient as it adds a smaller number of pseudo-labeled examples than tri-training at every epoch. For sentiment analysis, tri-training adds around 1800-1950/2000 unlabeled examples at every epoch, while MT-Tri only adds around 100-300 in early epochs. This shows that the orthogonality constraint is useful for inducing diversity. In addition, adding fewer examples poses a smaller risk of swamping the learned representations with useless signals and is more akin to fine-tuning, the standard method for supervised domain adaptation (Howard and Ruder, 2018).

We observe an asymmetry in the results between some of the domain pairs, e.g. B->D and D->B. We hypothesize that the asymmetry may be due to properties of the data and that the domains are relatively far apart e.g., in terms of  $\mathcal{A}$ -distance. In fact, asymmetry in these domains is already reflected

Model	$ep$	Target domains					Avg	WSJ	$\mu_{pseudo}$
		Answers	Emails	Newsgroups	Reviews	Weblogs			
Src (+glove)		87.63 $\pm$ .37	86.49 $\pm$ .35	<b>88.60</b> $\pm$ .22	90.12 $\pm$ .32	92.85 $\pm$ .17	89.14 $\pm$ .28	95.49 $\pm$ .09	—
Self	(5)	87.64 $\pm$ .18	86.58 $\pm$ .30	88.42 $\pm$ .24	90.03 $\pm$ .11	92.80 $\pm$ .19	89.09 $\pm$ .20	95.36 $\pm$ .07	.5k
Tri	(4)	88.42 $\pm$ .16	87.46 $\pm$ .20	87.97 $\pm$ .09	90.72 $\pm$ .14	93.40 $\pm$ .15	89.56 $\pm$ .16	95.94 $\pm$ .07	20.5k
Tri-D	(7)	<b>88.50</b> $\pm$ .04	<b>87.63</b> $\pm$ .15	88.12 $\pm$ .05	<b>90.76</b> $\pm$ .10	<b>93.51</b> $\pm$ .06	<b>89.70</b> $\pm$ .08	<b>95.99</b> $\pm$ .03	7.7K
Asym	(3)	87.81 $\pm$ .19	86.97 $\pm$ .17	87.74 $\pm$ .24	90.16 $\pm$ .17	92.73 $\pm$ .16	89.08 $\pm$ .19	95.55 $\pm$ .12	1.5k
MT-Tri	(4)	87.92 $\pm$ .18	87.20 $\pm$ .23	87.73 $\pm$ .37	90.27 $\pm$ .10	92.96 $\pm$ .07	89.21 $\pm$ .19	95.50 $\pm$ .06	7.6k
FLORS		89.71	88.46	89.82	92.10	94.20	90.86	95.80	—

Table 3: Accuracy scores on dev set of target domain for POS tagging for 10% labeled data. Avg: average over the 5 SANCL domains. Hyperparameter  $ep$  (epochs) is tuned on Answers dev.  $\mu_{pseudo}$ : average amount of added pseudo-labeled data. FLORS: results for Batch (u:big) from (Yin et al., 2015) (see §3).

Model	Target domains dev sets					Avg on targets	WSJ
	Answers	Emails	Newsgroups	Reviews	Weblogs		
TnT*	88.55	88.14	88.66	90.40	93.33	89.82	95.75
Stanford*	88.92	88.68	89.11	91.43	94.15	90.46	96.83
Src	88.84 $\pm$ .15	88.24 $\pm$ .12	89.45 $\pm$ .23	91.24 $\pm$ .03	93.92 $\pm$ .17	90.34 $\pm$ .14	96.69 $\pm$ .08
Tri	89.34 $\pm$ .18	88.83 $\pm$ .07	89.32 $\pm$ .21	91.62 $\pm$ .06	94.40 $\pm$ .06	90.70 $\pm$ .12	96.84 $\pm$ .04
Tri-D	89.35 $\pm$ .16	88.66 $\pm$ .09	89.29 $\pm$ .12	91.58 $\pm$ .05	94.32 $\pm$ .05	90.62 $\pm$ .09	96.85 $\pm$ .06
Src (+glove)	89.35 $\pm$ .16	88.55 $\pm$ .14	<b>90.12</b> $\pm$ .31	91.48 $\pm$ .15	94.48 $\pm$ .07	90.80 $\pm$ .17	96.90 $\pm$ .04
Tri	<b>90.00</b> $\pm$ .03	<b>89.06</b> $\pm$ .16	90.04 $\pm$ .25	<b>91.98</b> $\pm$ .11	<b>94.74</b> $\pm$ .06	<b>91.16</b> $\pm$ .12	<b>96.99</b> $\pm$ .02
Tri-D	89.80 $\pm$ .19	88.85 $\pm$ .10	90.03 $\pm$ .22	<b>91.98</b> $\pm$ .09	94.70 $\pm$ .05	91.01 $\pm$ .13	96.95 $\pm$ .05
Asym	89.51 $\pm$ .15	88.47 $\pm$ .19	89.26 $\pm$ .16	91.60 $\pm$ .20	94.28 $\pm$ .15	90.62 $\pm$ .17	96.56 $\pm$ .01
MT-Tri	89.45 $\pm$ .05	88.65 $\pm$ .04	89.40 $\pm$ .22	91.63 $\pm$ .23	94.41 $\pm$ .05	90.71 $\pm$ .12	97.37 $\pm$ .07
FLORS*	90.30	89.44	90.86	92.95	94.71	91.66	96.59
Model	Target domains test sets					Avg on targets	WSJ
	Answers	Emails	Newsgroups	Reviews	Weblogs		
TnT*	89.36	87.38	90.85	89.67	91.37	89.73	96.57
Stanford*	89.74	87.77	91.25	90.30	92.32	90.28	97.43
Src (+glove)	90.43 $\pm$ .13	87.95 $\pm$ .18	91.83 $\pm$ .20	90.04 $\pm$ .11	92.44 $\pm$ .14	90.54 $\pm$ .15	<b>97.50</b> $\pm$ .03
Tri	<b>91.21</b> $\pm$ .06	<b>88.30</b> $\pm$ .19	<b>92.18</b> $\pm$ .19	<b>90.06</b> $\pm$ .10	<b>92.85</b> $\pm$ .02	<b>90.92</b> $\pm$ .11	97.45 $\pm$ .03
Asym	90.62 $\pm$ .26	87.71 $\pm$ .07	91.40 $\pm$ .05	89.89 $\pm$ .22	92.37 $\pm$ .27	90.39 $\pm$ .17	97.19 $\pm$ .03
MT-Tri	90.53 $\pm$ .15	87.90 $\pm$ .07	91.45 $\pm$ .19	89.77 $\pm$ .26	92.35 $\pm$ .09	90.40 $\pm$ .15	97.37 $\pm$ .07
FLORS*	91.17	88.67	92.41	92.25	93.14	91.53	97.11

Table 4: Accuracy for POS tagging on the dev and test sets of the SANCL domains, models trained on full source data setup. Values for methods with \* are from (Schnabel and Schütze, 2014).

in the results of Blitzer et al. (2007) and is corroborated in the results for asymmetric tri-training (Saito et al., 2017) and our method.

We note a weakness of this dataset is high variance. Existing approaches only report the mean, which makes an objective comparison difficult. For this reason, we believe it is essential to evaluate proposed approaches also on other tasks.

**POS tagging** Results for tagging in the low-data regime (10% of WSJ) are given in Table 3.

Self-training does not work for the sequence prediction task. We report only the best instantia-

tion (throttling with  $n=800$ ). Our results contribute to negative findings regarding self-training (Plank, 2011; Van Asch and Daelemans, 2016).

In the low-data setup, tri-training *with disagreement* works best, reaching an overall average accuracy of 89.70, closely followed by classic tri-training, and significantly outperforming the baseline on 4/5 domains. The exception is newsgroups, a difficult domain with high OOV rate where none of the approaches beats the baseline (see §3.4). Our proposed MT-Tri is better than asymmetric tri-training, but falls below classic tri-training. It beats

	Ans	Email	Newsg	Rev	Web1
% unk tag	0.25	0.80	0.31	0.06	0.0
% OOV	8.53	10.56	10.34	6.84	8.45
% UWT	2.91	3.47	2.43	2.21	1.46
Accuracy on OOV tokens					
Src	54.26	57.48	<b>61.80</b>	59.26	<b>80.37</b>
Tri	<b>55.53</b>	<b>59.11</b>	61.36	<b>61.16</b>	79.32
Asym	52.86	56.78	56.58	59.59	76.84
MT-Tri	52.88	57.22	57.28	58.99	77.77
Accuracy on unknown word-tag (UWT) tokens					
Src	<b>17.68</b>	<b>11.14</b>	<b>17.88</b>	<b>17.31</b>	<b>24.79</b>
Tri	16.88	10.04	17.58	16.35	23.65
Asym	17.16	10.43	17.84	16.92	22.74
MT-Tri	16.43	11.08	17.29	16.72	23.13
FLORS*	17.19	15.13	21.97	21.06	21.65

Table 5: Accuracy scores on dev sets for OOV and unknown word-tag (UWT) tokens.

the baseline significantly on only 2/5 domains (answers and emails). The FLORS tagger (Yin et al., 2015) fares better. Its contextual distributional features are particularly helpful on unknown word-tag combinations (see § 3.4), which is a limitation of the lexicalized generic bi-LSTM tagger.

For the high-data setup (Table 4) results are similar. Disagreement, however, is only favorable in the low-data setups; the effect of avoiding easy points no longer holds in the full data setup. Classic tri-training is the best method. In particular, traditional tri-training is complementary to word embedding initialization, pushing the non-pre-trained baseline to the level of SRC with Glove initialization. Tri-training pushes performance even further and results in the best model, significantly outperforming the baseline again in 4/5 cases, and reaching FLORS performance on weblogs. Multi-task tri-training is often slightly more effective than asymmetric tri-training (Saito et al., 2017); however, improvements for both are not robust across domains, sometimes performance even drops. The model likely is too simplistic for such a high-data POS setup, and exploring shared-private models might prove more fruitful (Liu et al., 2017). On the test sets, tri-training performs consistently the best.

**POS analysis** We analyze POS tagging accuracy with respect to word frequency<sup>6</sup> and unseen word-tag combinations (UWT) on the dev sets. Table 5 (top rows) provides percentage of un-

<sup>6</sup>The binned log frequency was calculated with base 2 (bin 0 are OOVs, bin 1 are singletons and rare words etc).

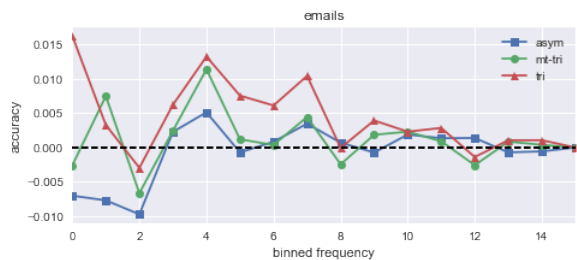


Figure 3: POS accuracy per binned log frequency.

known tags, OOVs and unknown word-tag (UWT) rate. The SANCL dataset is overall very challenging: OOV rates are high (6.8-11% compared to 2.3% in WSJ), so is the unknown word-tag (UWT) rate (answers and emails contain 2.91% and 3.47% UWT compared to 0.61% on WSJ) and almost all target domains even contain unknown tags (Schnabel and Schütze, 2014) (unknown tags: ADD,GW,NFP,XX), except for weblogs. Email is the domain with the highest OOV rate and highest unknown-tag-for-known-words rate. We plot accuracy with respect to word frequency on email in Figure 3, analyzing how the three methods fare in comparison to the baseline on this difficult domain.

Regarding OOVs, the results in Table 5 (second part) show that classic tri-training outperforms the source model (trained on only source data) on 3/5 domains in terms of OOV accuracy, except on two domains with high OOV rate (newsgroups and weblogs). In general, we note that tri-training works best on OOVs and on low-frequency tokens, which is also shown in Figure 3 (leftmost bins). Both other methods fall typically below the baseline in terms of OOV accuracy, but MT-Tri still outperforms Asym in 4/5 cases. Table 5 (last part) also shows that no bootstrapping method works well on unknown word-tag combinations. UWT tokens are very difficult to predict correctly using an unsupervised approach; the less lexicalized and more context-driven approach taken by FLORS is clearly superior for these cases, resulting in higher UWT accuracies for 4/5 domains.

## 4 Related work

**Learning under Domain Shift** There is a large body of work on domain adaptation. Studies on unsupervised domain adaptation include early work on *bootstrapping* (Steedman et al., 2003; McClosky et al., 2006a), *shared feature representations* (Blitzer et al., 2006, 2007) and *instance weighting* (Jiang and Zhai, 2007). Recent ap-



proaches include *adversarial learning* (Ganin et al., 2016) and *fine-tuning* (Sennrich et al., 2016). There is almost no work on bootstrapping approaches for recent neural NLP, in particular under domain shift. Tri-training is less studied, and only recently re-emerged in the vision community (Saito et al., 2017), albeit is not compared to classic tri-training.

**Neural network ensembling** Related work on self-ensembling approaches includes snapshot ensembling (Huang et al., 2017) or temporal ensembling (Laine and Aila, 2017). In general, the line between “explicit” and “implicit” ensembling (Huang et al., 2017), like dropout (Srivastava et al., 2014) or temporal ensembling (Saito et al., 2017), is more fuzzy. As we noted earlier our multi-task learning setup can be seen as a form of self-ensembling.

**Multi-task learning in NLP** Neural networks are particularly well-suited for MTL allowing for parameter sharing (Caruana, 1993). Recent NLP conferences witnessed a “tsunami” of deep learning papers (Manning, 2015), followed by what we call a multi-task learning “wave”: MTL has been successfully applied to a wide range of NLP tasks (Cohn and Specia, 2013; Cheng et al., 2015; Luong et al., 2015; Plank et al., 2016; Fang and Cohn, 2016; Søgaard and Goldberg, 2016; Ruder et al., 2017; Augenstein et al., 2018). Related to it is the pioneering work on adversarial learning (DANN) (Ganin et al., 2016). For sentiment analysis we found tri-training and our MT-Tri model to outperform DANN. Our MT-Tri model lends itself well to shared-private models such as those proposed recently (Liu et al., 2017; Kim et al., 2017), which extend upon (Ganin et al., 2016) by having separate source and target-specific encoders.

## 5 Conclusions

We re-evaluate a range of traditional general-purpose bootstrapping algorithms in the context of neural network approaches to semi-supervised learning under domain shift. For the two examined NLP tasks classic tri-training works the best and even outperforms a recent state-of-the-art method. The drawback of tri-training is its time and space complexity. We therefore propose a more efficient multi-task tri-training model, which outperforms both traditional tri-training and recent alternatives in the case of sentiment analysis. For POS tagging, classic tri-training is superior, performing especially well on OOVs and low frequency to-

kens, which suggests it is less affected by error propagation. Overall we emphasize the importance of comparing neural approaches to strong baselines and reporting results across several runs.

## Acknowledgments

We thank the anonymous reviewers for their valuable feedback. Sebastian is supported by Irish Research Council Grant Number EBPPG/2014/30 and Science Foundation Ireland Grant Number SFI/12/RC/2289. Barbara is supported by NVIDIA corporation and thanks the Computing Center of the University of Groningen for HPC support.

## References

- Steven Abney. 2007. *Semisupervised learning for computational linguistics*. CRC Press.
- Isabelle Augenstein, Sebastian Ruder, and Anders Søgaard. 2018. **Multi-task Learning of Pairwise Sequence Classification Tasks Over Disparate Label Spaces**. In *Proceedings of NAACL-HLT 2018*.
- John Blitzer, Mark Dredze, and Fernando Pereira. 2007. **Biographies, bollywood, boom-boxes and blenders: Domain adaptation for sentiment classification**. *Annual Meeting-Association for Computational Linguistics*, 45(1):440.
- John Blitzer, Ryan McDonald, and Fernando Pereira. 2006. **Domain Adaptation with Structural Correspondence Learning**. *Proceedings of the 2006 Conference on Empirical Methods in Natural Language Processing (EMNLP '06)*, pages 120–128.
- Konstantinos Bousmalis, George Trigeorgis, Nathan Silberman, Dilip Krishnan, and Dumitru Erhan. 2016. **Domain Separation Networks**. *NIPS*.
- Rich Caruana. 1993. **Multitask learning: A knowledge-based source of inductive bias**. In *Proceedings of the Tenth International Conference on Machine Learning*.
- Olivier Chapelle, Bernhard Schölkopf, and Alexander Zien. 2006. *Semi-Supervised Learning*, volume 1. MIT press.
- Hao Cheng, Hao Fang, and Mari Ostendorf. 2015. **Open-domain name error detection using a multi-task rnn**. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 737–746. Association for Computational Linguistics.
- Trevor Cohn and Lucia Specia. 2013. **Modelling annotator bias with multi-task gaussian processes: An application to machine translation quality estimation**. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1:*

- Long Papers*), pages 32–42, Sofia, Bulgaria. Association for Computational Linguistics.
- Michael Denkowski and Graham Neubig. 2017. Stronger baselines for trustable results in neural machine translation. *arXiv preprint arXiv:1706.09733*.
- Timothy Dozat and Christopher D. Manning. 2017. [Deep Biaffine Attention for Neural Dependency Parsing](#). In *Proceedings of ICLR 2017*.
- Bradley Efron and Robert J Tibshirani. 1994. *An introduction to the bootstrap*. CRC press.
- Meng Fang and Trevor Cohn. 2016. Learning when to trust distant supervision: An application to low-resource pos tagging using cross-lingual projection. In *Proceedings of CoNLL-16*.
- Meng Fang and Trevor Cohn. 2017. [Model transfer for tagging low-resource languages using a bilingual dictionary](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 587–593. Association for Computational Linguistics.
- Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, Francois Laviolette, Mario Marchand, and Victor Lempitsky. 2016. [Domain-Adversarial Training of Neural Networks](#). *Journal of Machine Learning Research*, 17:1–35.
- Rob van der Goot, Barbara Plank, and Malvina Nissim. 2017. To normalize, or not to normalize: The impact of normalization on part-of-speech tagging. In *Proceedings of the 3rd Workshop on Noisy User-generated Text*, pages 31–39, Copenhagen, Denmark. Association for Computational Linguistics.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. 2017. [On Calibration of Modern Neural Networks](#). *Proceedings of ICML 2017*.
- Luheng He, Kenton Lee, Mike Lewis, and Luke Zettlemoyer. 2017. Deep semantic role labeling: What works and what’s next. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 473–483, Vancouver, Canada. Association for Computational Linguistics.
- Jeremy Howard and Sebastian Ruder. 2018. Universal Language Model Fine-tuning for Text Classification. In *Proceedings of ACL 2018*.
- Gao Huang, Yixuan Li, Geoff Pleiss, Zhuang Liu, John E. Hopcroft, and Kilian Q. Weinberger. 2017. Snapshot Ensembles: Train 1, get M for free. In *Proceedings of ICLR 2017*.
- Jing Jiang and ChengXiang Zhai. 2007. Instance weighting for domain adaptation in nlp. In *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics*, pages 264–271. Association for Computational Linguistics.
- Young-Bum Kim, Karl Stratos, and Dongchan Kim. 2017. Adversarial adaptation of synthetic or stale data. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1297–1307, Vancouver, Canada. Association for Computational Linguistics.
- Eliyahu Kiperwasser and Yoav Goldberg. 2016. Simple and accurate dependency parsing using bidirectional lstm feature representations. *Transactions of the Association for Computational Linguistics*, 4:313–327.
- Samuli Laine and Timo Aila. 2017. Temporal Ensembling for Semi-Supervised Learning. In *Proceedings of ICLR 2017*.
- Guillaume Lample, Miguel Ballesteros, Sandeep Subramanian, Kazuya Kawakami, and Chris Dyer. 2016. [Neural Architectures for Named Entity Recognition](#). In *NAACL-HLT 2016*.
- Pengfei Liu, Xipeng Qiu, and Xuanjing Huang. 2017. Adversarial multi-task learning for text classification. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1–10, Vancouver, Canada. Association for Computational Linguistics.
- Christos Louizos, Kevin Swersky, Yujia Li, Max Welling, and Richard Zemel. 2015. The variational fair autoencoder. *arXiv preprint arXiv:1511.00830*.
- Minh-Thang Luong, Quoc V Le, Ilya Sutskever, Oriol Vinyals, and Lukasz Kaiser. 2015. Multi-task sequence to sequence learning. *arXiv preprint arXiv:1511.06114*.
- Christopher D Manning. 2015. Computational linguistics and deep learning. *Computational Linguistics*, 41(4):701–707.
- David McClosky, Eugene Charniak, and Mark Johnson. 2006a. Effective self-training for parsing. In *Proceedings of the Human Language Technology Conference of the NAACL, Main Conference*, New York City, USA. Association for Computational Linguistics.
- David McClosky, Eugene Charniak, and Mark Johnson. 2006b. [Reranking and Self-Training for Parser Adaptation](#). *International Conference on Computational Linguistics (COLING) and Annual Meeting of the Association for Computational Linguistics (ACL)*, (July):337–344.
- Gábor Melis, Chris Dyer, and Phil Blunsom. 2017. [On the State of the Art of Evaluation in Neural Language Models](#). In *arXiv preprint arXiv:1707.05589*.
- Graham Neubig, Chris Dyer, Yoav Goldberg, Austin Matthews, Waleed Ammar, Antonios Anastasopoulos, Miguel Ballesteros, David Chiang, Daniel Clothiaux, Trevor Cohn, et al. 2017. Dynet: The dynamic neural network toolkit. *arXiv preprint arXiv:1701.03980*.

- Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. Glove: Global vectors for word representation. In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543.
- Slav Petrov and Ryan McDonald. 2012. Overview of the 2012 shared task on parsing the web. *Notes of the First Workshop on Syntactic Analysis of Non-Canonical Language (SANCL)*, 59.
- Barbara Plank. 2011. *Domain adaptation for parsing*. University Library Groningen.
- Barbara Plank, Anders Søgaard, and Yoav Goldberg. 2016. [Multilingual Part-of-Speech Tagging with Bidirectional Long Short-Term Memory Models and Auxiliary Loss](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*.
- Roi Reichart and Ari Rappoport. 2007. Self-training for enhancement and domain adaptation of statistical parsers trained on small datasets. In *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics*, pages 616–623.
- Nils Reimers and Iryna Gurevych. 2017. Reporting score distributions makes a difference: Performance study of lstm-networks for sequence tagging. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 338–348, Copenhagen, Denmark. Association for Computational Linguistics.
- Sebastian Ruder, Joachim Bingel, Isabelle Augenstein, and Anders Søgaard. 2017. [Learning what to share between loosely related tasks](#). *arXiv preprint arXiv:1705.08142*.
- Kuniaki Saito, Yoshitaka Ushiku, and Tatsuya Harada. 2017. [Asymmetric Tri-training for Unsupervised Domain Adaptation](#). In *ICML 2017*.
- Tobias Schnabel and Hinrich Schütze. 2014. FLORS: Fast and Simple Domain Adaptation for Part-of-Speech Tagging. *TACL*, 2:15–26.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving neural machine translation models with monolingual data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 86–96, Berlin, Germany. Association for Computational Linguistics.
- Anders Søgaard. 2010. Simple semi-supervised training of part-of-speech taggers. In *Proceedings of the ACL 2010 Conference Short Papers*, pages 205–208.
- Anders Søgaard and Yoav Goldberg. 2016. Deep multi-task learning with low level tasks supervised at lower layers. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*, pages 231–235, Berlin, Germany. Association for Computational Linguistics.
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. [Dropout: A Simple Way to Prevent Neural Networks from Overfitting](#). *Journal of Machine Learning Research*, 15:1929–1958.
- Mark Steedman, Rebecca Hwa, Stephen Clark, Miles Osborne, Anoop Sarkar, Julia Hockenmaier, Paul Ruhlen, Steven Baker, and Jeremiah Crim. 2003. Example selection for bootstrapping statistical parsers. In *Proceedings of the 2003 Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics*.
- Jun Suzuki and Hideki Isozaki. 2008. Semi-supervised sequential labeling and segmentation using gigaword scale unlabeled data. pages 665–673.
- Vincent Van Asch and Walter Daelemans. 2016. Predicting the effectiveness of self-training: Application to sentiment classification. *arXiv preprint arXiv:1601.03288*.
- Fangzhao Wu and Yongfeng Huang. 2016. Sentiment Domain Adaptation with Multiple Sources. *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (ACL 2016)*, pages 301–310.
- David Yarowsky. 1995. Unsupervised Word Sense Disambiguation Rivaling Supervised Methods. In *Proceedings of the 33rd annual meeting on Association for Computational Linguistics*.
- Wenpeng Yin, Tobias Schnabel, and Hinrich Schütze. 2015. Online Updating of Word Representations for Part-of-Speech Tagging. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, September, pages 1329–1334.
- Guangyou Zhou, Zhiwen Xie, Jimmy Xiangji Huang, and Tingting He. 2016. Bi-transferring deep neural networks for domain adaptation. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 322–332, Berlin, Germany. Association for Computational Linguistics.
- Zhi-Hua Zhou and Ming Li. 2005. [Tri-Training: Exploiting Unlabeled Data Using Three Classifiers](#). *IEEE Trans.Data Eng.*, 17(11):1529–1541.
- Xiaojin Zhu. 2005. Semi-Supervised Learning Literature Survey. Technical Report 1530, Computer Sciences, University of Wisconsin-Madison.
- Xiaojin Zhu and Andrew B Goldberg. 2009. Introduction to semi-supervised learning. *Synthesis lectures on artificial intelligence and machine learning*, 3(1):1–130.