

On the Complementarity of Data Selection and Fine Tuning for Domain Adaptation

Anonymous ACL submission

Abstract

Domain adaptation of neural networks commonly relies on three training phases: pretraining, selected data training and then fine tuning. Data selection improves target domain generalization by training further on pretraining data identified by relying on a small sample of target domain data. This work examines the benefit of data selection for language modeling and machine translation. Our experiments assess the complementarity of selection with fine tuning and result in practical recommendations: (i) selected data must be similar to the fine-tuning domain but not so much as to erode the complementary effect of fine-tuning; (ii) there is a trade-off between selecting little data for fast but limited progress or much data for slow but long lasting progress; (iii) data selection can be applied early during pretraining, with performance gains comparable to long pretraining session; (iv) data selection from domain classifiers is often more effective than the popular contrastive data selection method.

1 Introduction

Machine learning models, and neural networks in particular, benefit from large training sets. However, for many application domains, the amount of training data representative of the inference conditions is limited. It is therefore common to train a model over a large amount of generic, out-of-domain data while relying on a small amount of target domain data to adapt such a model. In the recent years, a large body of work has focused on leveraging large amount of web data to train neural networks for language modeling (Peters et al., 2018; Devlin et al., 2019) or translation systems (Bañón et al., 2020; Koehn et al., 2020). Such systems are then adapted to the target distribution, typically via fine tuning (Liu et al., 2019; Raffel et al., 2020). This work studies data selection, an intermediate training phase that visits a subset of the out-of-

domain data that is deemed closer to the target domain.

Previous work has proposed conducting a data selection step after pretraining (van der Wees et al., 2017a; Wang et al., 2018; Gururangan et al., 2020; Aharoni and Goldberg, 2020), either as a final training stage or before regular fine tuning. Data selection is meant to identify a subset of the out-of-domain pretraining set which might be the most helpful to improve generalization on the target distribution. This selection is typically conducted by estimating the probability that each data point belongs to the target domain (Moore and Lewis, 2010; Axelrod et al., 2011). Recently, (Aharoni and Goldberg, 2020) introduced the use of domain classifiers for data selection.

This work examines the benefit of data selection for language modeling and machine translation. We compare different selection methods and examine their effect for short and long pretraining sessions. We also examine the benefit of selecting varying amount of training data and the impact of selection on the subsequent benefit of fine-tuning. In addition to this novel analysis, our machine translation experiments compare the benefit of selecting data with a classifier based on source language, target language or both.

The effectiveness of data selection is dependent on (i) the similarity of the pretraining data to the target domain data, (ii) the precision of the selection method to identify in-domain examples from the pretraining set, (iii) the extent to which training on the selected data is complementary to fine-tuning. This work focuses on selecting data from the pretraining set so (i) is fixed. We show that (ii) benefits from the use of domain classifiers, in particular, fine-tuned pretrained language models, outperforming the more popular contrastive methods (eg. Wang et al. (2018)) in all settings that we tested. We present the first analysis of (iii), which we refer to as the complementarity of selected data

083 to finetuning data. We show that some data selec- 133
084 tion methods can actually erode the effectiveness of 134
085 subsequent fine-tuning. In some settings, we even 135
086 report that a poor complementarity of selection and 136
087 fine tuning can result in their combination reaching 137
088 worse results than fine tuning alone. 138

089 Effective application of data selection requires 139
090 careful selection of when to switch from pretrain- 140
091 ing to selection, how much selected data to train 141
092 on and how long to train on selected data before 142
093 switching to finetuning. Much of the previous work 143
094 on data selection either evaluates small models that 144
095 converge quickly (Moore and Lewis, 2010; Axel- 145
096 rod et al., 2011) or does not describe the extent 146
097 of grid search over selection size, number of steps 147
098 of pretraining and number of steps of training on 148
099 selected data. We are the first to analyze the hy- 149
100 perparameter selection tradeoffs for data selection 150
101 on large neural models, where models may be un- 151
102 dertrained (Liu et al., 2019) and evaluating many 152
103 selection sizes may be prohibitively expensive. We 153
104 evaluate data selection on checkpoints with vari- 154
105 able numbers of pretraining steps and show that 155
106 data selection provides consistent results between 156
107 minimally and extensively pretrained models. We 157
108 also show the challenges of searching over selec- 158
109 tion sizes because smaller selection sizes always 159
110 converge more quickly but are outperformed by 160
111 larger selection sizes trained for more steps. 161

112 Our findings are the following: (i) the data se- 162
113 lection mechanism must select data that is similar, 163
114 but complementary to the fine tuning dataset (ii) 164
115 the amount of selected data introduces a trade-off 165
116 between quick but limited improvements when lim- 166
117 iting selection to the best data, and long lasting 167
118 but slow progress when selecting more data with 168
119 an overall worse quality, (iii) data selection tech- 169
120 niques are not created equal and domain classifiers 170
121 often outperform contrastive scoring, the most com- 171
122 mon data selection method, (iv) we propose three 172
123 simple variants of domain classifiers for machine 173
124 translation that can conditions the classifier on ei- 174
125 ther source, target or both. We demonstrate these 175
126 findings on language modeling and two language 176
127 pairs for neural machine translation. 177

128 2 Related Work 178

129 In Natural Language Processing (NLP), adapta- 179
130 tion methods have been applied to language mod- 180
131 eling (Moore and Lewis, 2010), machine transla- 181
132 tion (Axelrod et al., 2011; Daumé III and Jagarla-

133 mudi, 2011), dependency parsing (Finkel and Man- 134
135 ning, 2009) or sentiment analysis (Tan et al., 2009; 135
136 Glorot et al., 2011). With the growing popularity of 136
137 neural methods (Collobert et al., 2011; Bahdanau 137
138 et al., 2015; Goldberg, 2017), the adaptation of neu- 138
139 ral models via fine tuning has become wide-spread 139
140 for various NLP applications (Devlin et al., 2019; 140
141 Liu et al., 2019; Raffel et al., 2020). Data selection 141
142 is another popular technique (van der Wees et al., 142
143 2017b; Wang et al., 2018) which can be used on its 143
144 own or in combination to fine tuning. 144

145 Data selection is a common domain adaptation 145
146 method. It has been introduced before neural 146
147 methods were popular (Moore and Lewis, 2010; 147
148 Axelrod et al., 2011) and has later been adapted to 148
149 neural networks (Duh et al., 2013; van der Wees 149
150 et al., 2017b; Wang et al., 2018). Data selection 150
151 relies on an intermediate classifier which discrim- 151
152 inate between in-domain and out-of-domain data. 152
153 This classifier is trained relying on the small in- 153
154 domain dataset and the large out-of-domain dataset 154
155 and is then applied to the out-of-domain set to iden- 155
156 tify the examples closest to the targeted domain. 156
157 Choosing a selection model and the amount of out- 157
158 of-domain data to select have a strong impact on 158
159 the effectiveness of the selection methods (Aharoni 159
160 and Goldberg, 2020; Gururangan et al., 2020). Our 160
161 experiments explore these aspects, in addition to 161
162 the complementarity of selection with fine tuning. 162

163 Data selection can be performed in multiple 163
164 rounds, either to gradually restrict the out-of- 164
165 domain dataset to less and less data (van der Wees 165
166 et al., 2017b) or to re-evaluate out-of-domain data 166
167 as pretraining progresses (Wang et al., 2018). Data 167
168 selection can also be performed as a continuous 168
169 online process (Wang et al., 2018, 2021; Dou et al., 169
170 2020). Our work focus on single round data se- 170
171 lection, the most common setting. The benefit of 171
172 dynamic selection effectiveness has shown to be 172
173 variable (Wang et al., 2018) and its use involves 173
174 defining a complex schedule which is a research 174
175 topic in itself (Kumar et al., 2019). 175

176 Data selection for domain adaptation is also re- 176
177 lated to data selection for multitask learning. In 177
178 that case, the out-of-domain dataset is composed of 178
179 heterogeneous data from different tasks/domains 179
180 and the training algorithm favor data from some 180
181 tasks at the expense of others (Graves et al., 2017; 181
182 Wu et al., 2020; Standley et al., 2020). Contrary to 182
183 our setting, selection operates only at the task level 183
184 and the association of training examples to tasks 184

is already known. Multitask learning is an active area of research. This area has explored dynamic selection with reinforcement learning (Graves et al., 2017; Guo et al., 2019) as well as update projections to align out-of-domain gradients to in-domain gradients (Yu et al., 2020; Dery et al., 2021). Some of these ideas have later been investigated in the context of data selection for domain adaptation (Wu et al., 2018; Kumar et al., 2019; Wang et al., 2021).

3 Data Selection Methods

This section presents the selection method our experiments will focus on and introduce the trade-offs involved in choosing data selection hyperparameters.

3.1 In-Domain Data Selection

Domain adaptation has been introduced for application domains where data reflecting the inference conditions is only available in limited quantity. This setting considers that two training sets are available, a large generic out-of-domain dataset and a small specialized in-domain dataset from the targeted domain (Søgaard, 2013). Classical machine learning assumes that training and test data originate from the same distribution. At the same time, statistical modeling reaches better generalization performance with large training sets (Vapnik, 1998). Domain adaptation therefore faces a tension between using a large data set with a distribution possibly far from the test conditions and using a small training set matching the test condition.

Data selection tries to address this dilemma. It examines the out-of-domain data and identifies training examples likely to be most effective at improving the in-domain training loss. For neural methods, data selection is often used in conjunction with fine tuning in a three phases process, as shown in Algorithm 1. In a first phase, the model is pre-trained on all the out-of-domain data. In a second phase, an intermediate classifier is trained to distinguish in-domain from out-of-domain data, using both training sets. The classifier is applied to the out-of-domain set to identify examples considered close to in-domain data. The intermediate classifier is then no longer required and the main model is trained on the selected data starting from the pre-trained parameters. Finally, the main model is fine tuned, i.e. it is trained on the small in-domain training dataset starting from the parameters after the selection phase.

Algorithm 1: Data Selection & Fine Tuning for Neural Models

Input: D, T out and in domain train sets.

Output: θ trained model parameters.

Function `Select` (D, T, n):

$w \leftarrow \text{trainClassifier}(D \cup T)$

$Y \leftarrow \text{classify}(w, D)$

return $\text{argtop}_n(Y)$

Function `Main` (D, T):

$\theta_0 \leftarrow \text{initParam}()$

$\theta_{\text{pre}} \leftarrow \text{train}(\theta_0, D)$ #pretraining

$D_{\text{sel}} \leftarrow \text{select}(D, T, n)$

$\theta_{\text{sel}} \leftarrow \text{train}(\theta_{\text{pre}}, D_{\text{sel}})$

$\theta_{\text{ft}} \leftarrow \text{train}(\theta_{\text{sel}}, T)$ #fine-tuning

return θ_{ft}

Contrastive Data Selection: Commonly, classification is done by estimating the probability that a given out-of-domain example x belongs to the target domain, $P(\mathcal{T}|x)$. Such an estimation can be done by contrasting the likelihood estimated by in-domain LM, $P(\cdot|\mathcal{T})$ and an out-of-domain LM, $P(\cdot|\mathcal{D})$, i.e.

$$\log P(\mathcal{T}|x) = \log P(x|\mathcal{T}) - \log P(x|\mathcal{D}) + C \quad (1)$$

where C is a constant (log prior ratio). This method was introduced as *intelligent selection* (Moore and Lewis, 2010) and was later renamed *contrastive data selection* (CDS) (Wang et al., 2018). Initially, it relied on independent n-gram LMs for estimating $P(\cdot|\mathcal{T})$ and $P(\cdot|\mathcal{D})$, trained respectively on the (small) in-domain training set T and the (large) out-of-domain training set D (Moore and Lewis, 2010; Axelrod et al., 2011). With neural LMs, $P(\cdot|\mathcal{T})$ can be estimated by fine-tuning $P(\cdot|\mathcal{D})$ as suggested by (van der Wees et al., 2017b; Wang et al., 2018).

The fine tuning strategy is particularly efficient when one performs data selection to adapt a language model. In that case, there is no need for an intermediate model. The pretrained language model to adapt is itself fine-tuned in a few steps on T and is itself used to score the out-of-domain set.

Classifier Selection: Discriminative classification (DC), introduced by Aharoni and Goldberg (2020); Jacovi et al. (2021), trains a binary classifier to distinguish T and D examples. This classifier is either trained from scratch or fine tuned from a pre-

trained model (Devlin et al., 2019; Liu et al., 2019). Aharoni and Goldberg (2020) train the domain classifier, which they refer to as “Domain-Finetune”, only on the source (English) side of the parallel corpus. We propose two alternative domain classifiers, that instead condition the classifier on either the target language or both source and target concatenated. To finetune language models on the target language data, we use BERT models that are pretrained on German (deepset.ai), Russian (Kurato and Arkhipov, 2019) and multilingual BERT (Devlin et al., 2018).

The motivation for these alternative classifiers are two fold: (1) noisy web crawled translation datasets often have incorrect translations (or even languages) which could be missed by the domain classifier if only conditioning on the English source data, (2) the multilingual domain classifier is able to model the interaction between the source and target and is more analogous to the *bilingual cross-entropy difference* proposed by Axelrod et al. (2011)

Compared to CDS, DC trains a different model which adds training overhead. On the other hand, a distinct intermediate model offers more flexibility. The classifier might be pretrained on a different task (e.g. masked LM to select translation data) and its capacity can be selected independently from the hyperparameter of the model to be adapted. Both aspects are important since intermediate models can easily overfit given the small size of the target domain set T .

Nearest Neighbor Selection: A lesser used methods is sentence embedding nearest neighbors (Gururangan et al., 2020; Aharoni and Goldberg, 2020). Embedding nearest neighbors relies on a pretrained model (Devlin et al., 2019; Liu et al., 2019) to represent sentences as vectors and then measure a domain-score by comparing the distance between a candidate sentence vector $v(x)$ and the average in-domain sentence vector $\frac{1}{|T|} \sum_{x \in T} x$.

In our experiments, we evaluate both contrastive data selection, the most common method by far, and selection with discriminative classifiers as it has been shown more effective in subsequent work (Aharoni and Goldberg, 2020). Previous work and our preliminary experiments indicated that nearest neighbor selection was not competitive with other baselines so we do not include it in our analysis.

3.2 Hyperparameter Trade-offs

Data selection for domain adaptation requires selecting several hyperparameters: the *number of pretraining steps*, i.e. when to transition from training on the full out-of-domain set to the selected subset; the *number of selection steps*, i.e. how long to train the model on the selected data; the *fraction of selected data*, i.e. the size of the selected subset.

These parameters are important as they impact the computational cost of training and the target domain generalization performance. To examine these trade-offs, the difference between pretraining and fine-tuning is important. Pretraining on a large dataset starts with an initial strong generalization improvement, followed by a long session where the rate of generalization improvement is still positive but ever diminishing. Fine tuning gives a strong generalization improvement in a few steps before overfitting quickly. The fraction of selected data allows trading off between these two extremes: a large fraction of selected data results in a large training set with a distribution close to the out-of-domain distribution while a small fraction results in small training set with a distribution close to the in-domain distribution. This means that settings with large fractions can perform more steps with generalization improvement albeit at a slower pace compared to lower fraction settings. Thus the number of selection steps and the selected fraction parameter interact. Our experiments investigate this interaction.

We characterize the effects of overfitting of the intermediate selection classifier, which uniquely affects data selection in conjunction with finetuning. The intermediate classifier is trained on the small target domain set T . As any machine learning model, it is biased toward its training set and the data it selects can reflect this bias. The selected out-of-domain examples might resemble the examples of T more than other in-domain examples unseen during training. This bias transferred to the selected data is itself inherited by the model trained on the selected data. This indirect overfitting is crucial for later fine tuning: we report that, in some cases, the selected data is too similar to T . There, the complementary value of selection and fine tuning vanishes as data selection fails to identify data providing updates complementary to those provided later by fine tuning on T .

4 Experiments

We evaluate domain adaptation with data selection on two tasks, language modeling (LM) and machine translation (MT). For both tasks, we have a large out-of-domain dataset and a small number of examples from the target domain. Both sets of data fulfil two functions each. The out-of-domain data is used to pretrain the model and all the selected data come from the out-of-domain set. The small target domain set is used to train the intermediate model that scores examples for data selection and, critically, this same set is used for finetuning the final model. For evaluation, we also have a validation set and test set from the target domain. The validation set is used to select hyperparameters and early stopping points and the test set is only used for the final model evaluation.

For language modeling, we use the 4.5 million sentences from the One Billion Word corpus (Chelba et al., 2013) as the out-of-domain set and 5k sentences from the Yelp corpus as the target domain. This dataset was used for domain adaptation by (Oren et al., 2019) and we use their filtered and preprocessed version of the data, including the 1k Yelp validation set and 10k Yelp test set. We train 2 language models; a 2-layer LSTM recurrent network (Zaremba et al., 2014) and a base-size transformer (Vaswani et al., 2017).

Our machine translation experiments focus on English-to-German and English-to-Russian. For the out-of-domain set, we use 4.5 million English-to-German pairs and 5.2 million English-to-Russian pairs taken from filtered Paracrawl (Esplà et al., 2019). Paracrawl is composed of translations crawled from the web. Even though we use the filtered version of the dataset, Paracrawl is still noisy including examples of entirely mismatched sentences and occasionally incorrect languages. As in domain data, we rely on news data from the News Commentary Dataset (Tiedemann, 2012), which are high quality translations from the news domain. Our in-domain set is limited to 6k sentence pairs. We use an additional 3k for validation and 10k as the test set. As a neural MT model, we train a base transformer (Vaswani et al., 2017). Code to reproduce our experiments is available¹. Models are implemented with Flax (Heek et al., 2020).

We finetune on the small in-domain set by grid searching for a learning rate and using the validation set for early stopping.

¹Hidden for anonymity

4.1 Selection Methods

Contrastive Data Selection The base pretrained (PT) model is fine-tuned (FT) on the small target domain dataset. This model acts as the “intermediate” model in this setting. Each example in the out-of-domain dataset is scored by the difference between the log likelihoods of the fine-tuned model and the pretrained model. The full dataset can be ranked by this score and a threshold is selected to train on a uniform distribution of only the top examples.

Discriminative Classifier The target domain dataset is used as positive examples and random samples from the out-of-domain dataset are used as negative examples to train a discriminative domain classifier. The classifier can be a new model trained from random weights, the base model with a binary classification head or a pretrained model from another task (such as a generic masked language model). Unlike CDS, the base model is not necessarily reused. The input features to the classifier may either be representations learned from the pretrained base model, other embeddings or the raw text data. In the case of machine translation, the classifier can be trained on the source, target or both.

In our transformer experiments, we evaluate CDS and two classifiers, (i) a logistic regression model on bytepair encodings (Sennrich et al., 2016) and (ii) a fine-tuned BERT classifier (deepset.ai; Kuratov and Arkipov, 2019; Devlin et al., 2018). We use four settings for the BERT classifier, training on the source, target, mean of the former two, and concatenated language pairs, using the respective language specific pretrained BERT. For the concatenated case, we use a multilingual BERT.

	En-De		En-Ru	
	logPPL	BLEU	logPPL	BLEU
PT	1.666	<i>23.71</i>	1.815	<i>23.20</i>
+FT	1.612	<i>26.89</i>	1.708	<i>24.92</i>
PT + CDS	1.626	<i>26.77</i>	1.757	<i>24.08</i>
+FT	1.608	<i>27.27</i>	1.707	<i>25.08</i>
PT + DC (LogReg)	1.624	<i>26.22</i>	1.762	<i>23.43</i>
+FT	1.575	<i>27.54</i>	1.666	<i>25.35</i>
PT + DC (BERT)	1.599	<i>26.33</i>	1.752	<i>23.66</i>
+FT	1.550	27.78	1.645	25.52

Table 1: Data selection for machine translation of English to German and English to Russian. BLEU in italics next to log-perplexity (log PPL). For both datasets, models were trained to 200K steps of pretraining and 15k steps of data selection.

	En-De		En-Ru		LM
	lgPPL	BLEU	lgPPL	BLEU	lgPPL
PT	1.00	1.00	1.00	1.00	1.00
+FT	1.00	1.00	1.00	0.992	1.00
CDS	1.00	1.00	1.00	1.00	1.00
+FT	1.00	0.998	1.00	0.975	1.00
DC-LR	1.00	1.00	1.00	1.00	1.00
+FT	0.951	0.890	0.840	0.742	0.998
DC-BERT	1.00	1.00	1.00	1.00	1.00
+FT	-	-	-	-	-

Table 2: Paired bootstrap comparison: each value reports the fraction of samples with worse mean performance than PT + DC-BERT + FT for 1k samples of 10k sentences sampled from a 10k sample test set.

4.2 Training on Selected Data

Machine Translation Table 1 reports the log-perplexity and BLEU scores on two language pairs for each of the selection methods described above. Data selection always outperforms the baseline without selection, with the BERT domain classifier producing the best log-probability and BLEU on both datasets. The effectiveness of DC compared to CDS is a surprising result given the popularity of CDS. We fix the number of training steps on the selected data to 15K and pretrain the baseline model for an additional 15k steps so there is the same number of pretraining + finetuning steps for all settings. We search the optimal selection size for this cutoff of training steps, which we found to be 1 million for En-Ru and 500k for En-De. We report results before and after finetuning to highlight the variation in effectiveness of finetuning after the alternative selection methods. This is particularly noticeable for En-Ru where CDS outperforms the logistic regression classifier before finetuning but is worse after finetuning. In all settings, finetuning is more effective after data selection with a discriminative classifier rather than with CDS. Section 4.3 provides insight as to why this is the case.

Table 2 reports the paired bootstrap resampling (Koehn, 2004) where the PT + DC (BERT) + FT model is compared to the baseline models, in terms of loss and BLEU, corresponding to Table 1. Each value is computed from the 10,000 example test set. We draw 1,000 bootstrap samples of 10,000 points each, with replacement. This test shows that the classifier method of data selection outperforms CDS with over 99% statistical significance on log-perplexity.

Figure 1 shows the log-probabilities at different checkpoints ranging from 50k to 1 million steps

of training. The relative benefit of FT and DC+FT over PT is diminishing as training progresses. However, there are consistent benefits from data selection, so longer pretraining on large models is not sufficient to replace data selection. Even pretraining up to 1m steps and finetuning (log ppl = 1.530) does not reach the loss from DC + FT at 400k (log ppl = 1.519). The relative improvement between methods is surprisingly constant across pretraining steps with a slight decline in the complementary benefit of combining fine tuning with selection. This means that comparing the adaptation methods early in the pretraining process is indicative of their relative loss at a later stage.

Further evaluation of performance at different checkpoints throughout pretraining can be found in the Appendix.

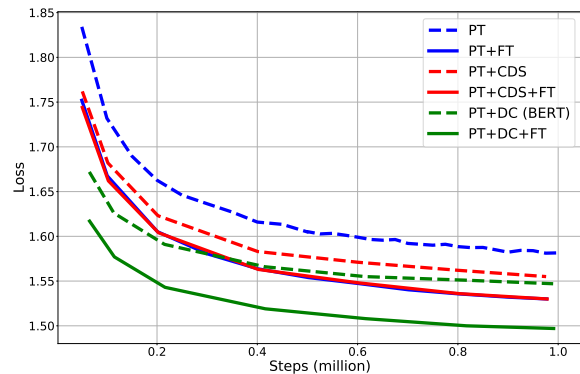


Figure 1: The validation loss curves for pretraining, data selection and finetuning (MT En-De). The pretraining loss (PT) is a single training run, whereas all the other points are checkpoints from the base run that were trained on selected data and/or finetuned.

Domain Classifier Variants Table 3 reports the log-perplexities and BLEU scores for the four variants of the BERT domain classifier. Aharoni and Goldberg (2020) propose the Source DC method. We propose also exploring target-language-conditioned domain classifiers, and in fact, find that the Target DC selection method outperforms Source DC on En-DE. Concatenation DC does not yield the best results despite having access to the most data (ie. both source and target). This may be because of the pretraining mismatch, in that Multilingual BERT was not trained on pairs of segments from different languages. We also take evaluate using the mean score of the source and target models as a simple alternative to the multilingual BERT approach. Future work may explore alterna-

518 tive methods for fusing source and target language
 519 representations for training a domain classifier.

	En-De		En-Ru	
	log PPL	BLEU	log PPL	BLEU
Target DC	1.550	27.78	1.653	25.21
Source DC	1.557	27.52	1.645	25.52
Concat DC	1.560	27.68	1.657	25.20
Mean DC	1.555	27.71	1.647	25.29

Table 3: Different types of BERT classifiers, target uses the target language (De/Ru), the source is English and *Concat* concatenates source and target and trains classifier on multilingual BERT. Mean takes the mean scores from source and target classifiers. All models are evaluated at 200k pretraining steps, similar to Table 1.

	LSTM	Transformer
PT	4.978	4.582
+FT	4.284	4.145
PT + CDS	4.548	4.392
+FT	4.183	4.151
PT + DC (LogReg)	4.644	4.456
+FT	4.183	4.108
PT + DC (LM Hidden)	4.603	-
+FT	4.179	-
PT + DC (BERT)	-	4.385
+FT	-	4.069

Table 4: Language modeling results (log-perplexity) across selection methods for an LSTM and a base-transformer. The LSTM was pretrained for 115k steps and the transformer was trained for 20k steps.

520 **Language Modeling** For language modeling we
 521 evaluate on both a modestly sized LSTM and a
 522 base-size transformer. For the LSTM domain clas-
 523 sifier, we reuse the pretrained language model as
 524 the feature representation for a simple linear do-
 525 main classifier (LM Hidden), as a smaller domain
 526 classifier seems appropriate given the smaller lan-
 527 guage model. We see similar results for the two
 528 models despite the large differences in number of
 529 parameters, training steps and proximity to conver-
 530 gence. The LM results in Table 4 show that fine
 531 tuning (PT+FT) and data selection (CDS, DC) are
 532 improving the pretrained model on target domain
 533 validation data. The benefit of FT alone is generally
 534 greater than selection alone but both approaches are
 535 complementary with the best result obtained with
 536 combined approaches (CDS+FT, DC+FT). When
 537 comparing methods we observe that DC is worse
 538 than CDS on its own but it is equivalent or bet-
 539 ter in combination with fine tuning (DC+FT vs
 540 CDS+FT). This indicates that the methods differ
 541 in their complementarity with FT and evaluating

542 selection approaches before fine tuning is not suffi-
 543 cient.

4.3 Overfitting and Complementarity

544 Our work compares two common data selection
 545 techniques, contrastive data selection (CDS) and
 546 a discriminative domain classifier (DC). As dis-
 547 cussed in the previous section, we found the combi-
 548 nation of DC+FT to be the most effective combina-
 549 tion both for our LM and MT settings. One reason
 550 of this success is the complementarity of DC with
 551 FT. CDS did not benefit as much from subsequent
 552 fine tuning as DC selection.

553 In Figure 2 (left), we show the learning curves
 554 for both CDS and DC (BERT) with the same se-
 555 lection size of 1m for MT with 200k steps of pre-
 556 training. The red dotted curve show that the CDS
 557 model reaches excellent performance on the target-
 558 domain training set, but fail to perform as well on
 559 the target-domain validation set. This means that
 560 the MT model trained on CDS selected data suffers
 561 more from overfitting than the MT model trained
 562 on DC selected data. This is particularly surpris-
 563 ing given the large selection size of nearly 1/4th of
 564 pretraining data. The data selected by CDS is too
 565 specific to the target-domain training set. This bias
 566 also certainly explains the worse complementarity
 567 of CDS and FT, i.e. if CDS selects a training set
 568 whose effect is similar to the target-domain training
 569 set T , the updates from T at fine-tuning are less
 570 beneficial.

571 Lastly, we examine important pitfalls to avoid
 572 when comparing selection methods and validat-
 573 ing their parameters. Figure 2 (middle) shows
 574 that when considering selection sets of different
 575 sizes, training curves converges at different rates.
 576 Small selected subsets progress at the fastest rate
 577 but reaches their best generalization quickly, and
 578 subsequently overfit, while large subsets progress
 579 at a slower rate but their best generalization later.
 580 This means that short diagnostics to pick the sub-
 581 set size will under estimate the value of large sub-
 582 sets. This is problematic for efficiently defining
 583 curriculum with data selection (Kumar et al., 2019).
 584 Similarly, the generalization loss of model which
 585 went through a data selection phase but prior to fine
 586 tuning is also misleading to predict its loss after
 587 fine tuning as illustrated in Figure 2 (right).
 588

4.4 Effectiveness of Data Selection

589 The purpose of the intermediate data selection
 590 model is to rank all the out-of-domain data from
 591

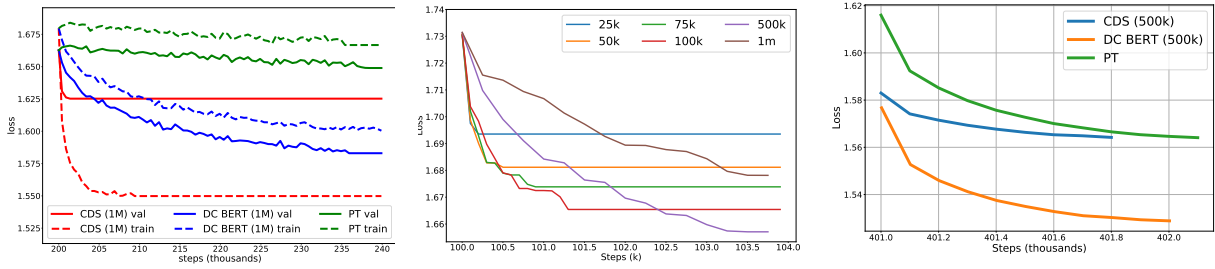


Figure 2: Effects of overfitting and complementarity: **Left:** Validation and training loss on the target domain during training on selected data (MT En-De). The dotted line falling below the solid line indicates the model is overfitting to the small target domain dataset despite never seeing this data in training. **Middle:** Loss curves for 6 different data selection sizes for DC (BERT) at the 100k checkpoint (MT En-De). Larger sizes improve loss more slowly but can be trained for longer to eventually outperform the smaller sets. For readability, we display the best checkpoint up to each step. **Right:** Validation loss on MT En-De during finetuning. Both data selection methods start at a loss that is better than pretraining but CDS does not benefit much from finetuning, reaching a loss similar to finetuning without data selection. Classifier selection has large a improvement from finetuning.

most to least similar with respect to the in-domain data. We evaluate and report the performance of CDS and DC for both LM and MT tasks. The data selection model is never used explicitly as a binary classifier but rather as a scorer. However, as a proxy for the quality of scoring, we evaluate the binary classification accuracy on an unseen set of in-domain and out-of-domain data. We also report the average quantile of the in-domain validation data which simulates where in the ranking true in-domain examples would appear. We split the out-of-domain data into 100 equal bins and take the average of the bin index that each in-domain example would fall into by its data selection score.

Table 5 shows good performance of CDS and DC for language modeling but clear underperformance of CDS as a binary classifier in the MT setting. Also, it is noteworthy that logistic regression on byte-pair unigrams outperforms CDS and approaches the performance of BERT while having many fewer parameters and a much lower training cost.

	Classifier	Accuracy	Avg Quant.
LM	CDS	91.65%	3.6
	MLP	89.02%	4.9
MT (En-De)	CDS	66.94%	26.0
	LogReg	87.52%	3.9
	BERT	93.51%	2.0

Table 5: Binary classification accuracy of domain classifier and average quantile of in-domain data when binned with ranked out-of-domain data.

5 Conclusions

This work explores data selection, a popular method for domain adaption for neural language modeling and neural machine translation. Data selection typically divides a training run into three phases: pretraining on out-of-domain data, training on out-of-domain data selected to resemble target domain data and fine tuning on target domain data. We compare the most common selection methods, contrastive data selection and discriminative model classifier and measure their complementarity with fine tuning.

Our experiments motivate several practical recommendations for the practitioner: (i) pretraining followed by data selection and fine tuning can reach a given generalization loss several time faster in terms of total training steps than pretraining with fine tuning; (ii) a data selection method should not be evaluated before fine tuning since not all methods/parameters bring the same complementary value compared to fine tuning; (iii) data selection should care about overfitting to the in-domain training set, since this type of overfitting results in selected data very similar to the fine tuning set and impacts the complementarity of data selection and fine tuning; (iv) longer pretraining runs are always beneficial to later adaptation stages for fine-tuning, data selection and their combination but pretraining has diminishing return; (v) despite the popularity of contrastive data selection, discriminative domain classifiers consistently outperformed this method in our experiments.

References

Roei Aharoni and Yoav Goldberg. 2020. [Unsupervised domain clusters in pretrained language models](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7747–7763, Online. Association for Computational Linguistics.

Amittai Axelrod, Xiaodong He, and Jianfeng Gao. 2011. [Domain adaptation via pseudo in-domain data selection](#). In *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing*, pages 355–362, Edinburgh, Scotland, UK. Association for Computational Linguistics.

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. [Neural machine translation by jointly learning to align and translate](#). In *ICLR*.

Marta Bañón, Pinzhen Chen, Barry Haddow, Kenneth Heafield, Hieu Hoang, Miquel Esplà-Gomis, Mikel L. Forcada, Amir Kamran, Faheem Kirefu, Philipp Koehn, Sergio Ortiz Rojas, Leopoldo Pla Sempere, Gema Ramírez-Sánchez, Elsa Sarriás, Marek Strelec, Brian Thompson, William Waites, Dion Wiggins, and Jaume Zaragoza. 2020. [ParaCrawl: Web-scale acquisition of parallel corpora](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4555–4567, Online. Association for Computational Linguistics.

Ciprian Chelba, Tomas Mikolov, Mike Schuster, Qi Ge, Thorsten Brants, Phillipp Koehn, and Tony Robinson. 2013. [One billion word benchmark for measuring progress in statistical language modeling](#). Technical report, Google.

Ronan Collobert, Jason Weston, Léon Bottou, Michael Karlen, Koray Kavukcuoglu, and Pavel Kuksa. 2011. [Natural language processing \(almost\) from scratch](#). *Journal of Machine Learning Research*, 12(76):2493–2537.

Hal Daumé III and Jagadeesh Jagarlamudi. 2011. [Domain adaptation for machine translation by mining unseen words](#). In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 407–412, Portland, Oregon, USA. Association for Computational Linguistics.

deepset.ai. [Open sourcing german bert](#). <https://deepset.ai/german-bert>.

Lucio Dery, Yann Dauphin, and David Grangier. 2021. [Auxiliary task update decomposition: The good, the bad and the neutral](#). In *International Conference on Learning Representation (ICLR)*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. [BERT: pre-training of deep bidirectional transformers for language understanding](#). *CoRR*, abs/1810.04805.

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.

Zi-Yi Dou, Antonios Anastasopoulos, and Graham Neubig. 2020. [Dynamic data selection and weighting for iterative back-translation](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5894–5904, Online. Association for Computational Linguistics.

Kevin Duh, Graham Neubig, Katsuhito Sudoh, and Hajime Tsukada. 2013. [Adaptation data selection using neural language models: Experiments in machine translation](#). In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 678–683, Sofia, Bulgaria. Association for Computational Linguistics.

Miquel Esplà, Mikel Forcada, Gema Ramírez-Sánchez, and Hieu Hoang. 2019. [ParaCrawl: Web-scale parallel corpora for the languages of the EU](#). In *Proceedings of Machine Translation Summit XVII Volume 2: Translator, Project and User Tracks*, pages 118–119, Dublin, Ireland. European Association for Machine Translation.

Jenny Rose Finkel and Christopher D. Manning. 2009. [Hierarchical Bayesian domain adaptation](#). In *Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 602–610, Boulder, Colorado. Association for Computational Linguistics.

Xavier Glorot, Antoine Bordes, and Yoshua Bengio. 2011. [Domain adaptation for large-scale sentiment classification: A deep learning approach](#). In *ICML*.

Yoav Goldberg. 2017. [Neural network methods for natural language processing](#). *Synthesis lectures on human language technologies*, 10(1):1–309.

Alex Graves, Marc G. Bellemare, Jacob Menick, Rémi Munos, and Koray Kavukcuoglu. 2017. [Automated curriculum learning for neural networks](#). In *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 1311–1320. PMLR.

Han Guo, Ramakanth Pasunuru, and Mohit Bansal. 2019. [AutoSeM: Automatic task selection and mixing in multi-task learning](#). 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 3520–3531, Minneapolis, Minnesota. Association for Computational Linguistics.

759	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 <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 8342–8360, Online. Association for Computational Linguistics.	814
760		815
761		816
762		817
763		818
764		819
765		820
766		821
767	Jonathan Heek, Anselm Levskaya, Avital Oliver, Marvin Ritter, Bertrand Rondepierre, Andreas Steiner, and Marc van Zee. 2020. Flax: A neural network library and ecosystem for JAX .	822
768		823
769		824
770		825
771	Alon Jacovi, Gang Niu, Yoav Goldberg, and Masashi Sugiyama. 2021. Scalable evaluation and improvement of document set expansion via neural positive-unlabeled learning . In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 581–592, Online. Association for Computational Linguistics.	826
772		827
773		828
774		829
775		830
776		831
777		832
778		833
779	Philipp Koehn. 2004. Statistical significance tests for machine translation evaluation . In <i>Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing</i> , pages 388–395, Barcelona, Spain. Association for Computational Linguistics.	834
780		835
781		836
782		837
783		838
784		839
785	Philipp Koehn, Vishrav Chaudhary, Ahmed El-Kishky, Naman Goyal, Peng-Jen Chen, and Francisco Guzmán. 2020. Findings of the WMT 2020 shared task on parallel corpus filtering and alignment . In <i>Proceedings of the Fifth Conference on Machine Translation</i> , pages 726–742, Online. Association for Computational Linguistics.	840
786		841
787		842
788		843
789		844
790		845
791		846
792	Gaurav Kumar, George Foster, Colin Cherry, and Maxim Krikun. 2019. Reinforcement learning based curriculum optimization for neural machine translation . In <i>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)</i> , pages 2054–2061, Minneapolis, Minnesota. Association for Computational Linguistics.	847
793		848
794		849
795		850
796		851
797		852
798		853
799		854
800		855
801	Yuri Kuratov and Mikhail Arkhipov. 2019. Adaptation of deep bidirectional multilingual transformers for russian language . <i>arXiv preprint arXiv:1905.07213</i> .	856
802		857
803		858
804	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 .	859
805		860
806		861
807		862
808		863
809	Robert C. Moore and William Lewis. 2010. Intelligent selection of language model training data . In <i>Proceedings of the ACL 2010 Conference Short Papers</i> , pages 220–224, Uppsala, Sweden. Association for Computational Linguistics.	864
810		865
811		866
812		867
813		868
	Yonatan Oren, Shiori Sagawa, Tatsunori Hashimoto, and Percy Liang. 2019. Distributionally robust language modeling . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 4227–4237, Hong Kong, China. Association for Computational Linguistics.	869
		870
	Matthew Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 2227–2237, New Orleans, Louisiana. Association for Computational Linguistics.	871
		872
	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer . <i>Journal of Machine Learning Research</i> , 21(140):1–67.	873
		874
	Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units . In <i>Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.	875
		876
	Anders Søgaard. 2013. Semi-supervised learning and domain adaptation in natural language processing . <i>Synthesis Lectures on Human Language Technologies</i> , 6(2):1–103.	877
		878
	Trevor Standley, Amir Roshan Zamir, Dawn Chen, Leonidas J. Guibas, Jitendra Malik, and Silvio Savarese. 2020. Which tasks should be learned together in multi-task learning? In <i>Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event</i> , volume 119 of <i>Proceedings of Machine Learning Research</i> , pages 9120–9132. PMLR.	879
		880
	Songbo Tan, Xueqi Cheng, Yuefen Wang, and Hongbo Xu. 2009. Adapting naive bayes to domain adaptation for sentiment analysis . In <i>European Conference on Information Retrieval</i> , pages 337–349. Springer.	881
		882
	Jörg Tiedemann. 2012. Parallel data, tools and interfaces in opus . In <i>Proceedings of the Eight International Conference on Language Resources and Evaluation (LREC'12)</i> , Istanbul, Turkey. European Language Resources Association (ELRA).	883
		884
	Marlies van der Wees, Arianna Bisazza, and Christof Monz. 2017a. Dynamic data selection for neural machine translation . In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing</i> , pages 1400–1410.	885
		886
		887
		888
		889

870 Marlies van der Wees, Arianna Bisazza, and Christof
 871 Monz. 2017b. [Dynamic data selection for neural](#)
 872 [machine translation](#). In *Proceedings of the 2017*
 873 *Conference on Empirical Methods in Natural Lan-*
 874 *guage Processing*, pages 1400–1410, Copenhagen,
 875 Denmark. Association for Computational Linguistics.
 876

877 V.N. Vapnik. 1998. *Statistical Learning Theory*. A
 878 Wiley-Interscience publication. Wiley.

879 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob
 880 Uszkoreit, Llion Jones, Aidan N. Gomez, undefin-
 881 dukasz Kaiser, and Illia Polosukhin. 2017. Attention
 882 is all you need. In *Proceedings of the 31st Interna-*
 883 *tional Conference on Neural Information Processing*
 884 *Systems, NIPS’17*, page 6000–6010, Red Hook, NY,
 885 USA. Curran Associates Inc.

886 Wei Wang, Taro Watanabe, Macduff Hughes, Tetsuji
 887 Nakagawa, and Ciprian Chelba. 2018. [Denoising](#)
 888 [neural machine translation training with trusted](#)
 889 [data and online data selection](#). In *Proceedings of*
 890 *the Third Conference on Machine Translation: Re-*
 891 *search Papers*, pages 133–143, Brussels, Belgium.
 892 Association for Computational Linguistics.

893 Xinyi Wang, Ankur Bapna, Melvin Johnson, and Orhan
 894 Firat. 2021. [Gradient-guided loss masking for neu-](#)
 895 [ral machine translation](#).

896 Jiawei Wu, Lei Li, and William Yang Wang. 2018. [Re-](#)
 897 [inforced co-training](#). In *Proceedings of the 2018*
 898 *Conference of the North American Chapter of the*
 899 *Association for Computational Linguistics: Human*
 900 *Language Technologies, Volume 1 (Long Papers)*,
 901 pages 1252–1262, New Orleans, Louisiana. Associ-
 902 ation for Computational Linguistics.

903 Sen Wu, Hongyang R. Zhang, and Christopher Ré.
 904 2020. [Understanding and improving information](#)
 905 [transfer in multi-task learning](#). In *8th International*
 906 *Conference on Learning Representations, ICLR*
 907 *2020, Addis Ababa, Ethiopia, April 26-30, 2020*.
 908 OpenReview.net.

909 Tianhe Yu, Saurabh Kumar, Abhishek Gupta, Sergey
 910 Levine, Karol Hausman, and Chelsea Finn. 2020.
 911 [Gradient surgery for multi-task learning](#). In *Ad-*
 912 *vances in Neural Information Processing Systems*,
 913 volume 33, pages 5824–5836. Curran Associates,
 914 Inc.

915 Wojciech Zaremba, Ilya Sutskever, and Oriol Vinyals.
 916 2014. Recurrent neural network regularization.
 917 *arXiv preprint arXiv:1409.2329*.

918 A Appendix

919 A.1 Training Steps

920 Figure 3 shows the acceleration of training as a
 921 function of pretraining + finetuning (PT+FT) steps
 922 needed to reach an equivalent loss for translation.

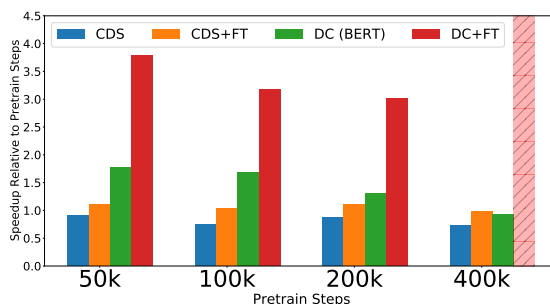


Figure 3: Data selection (MT En-De) as an acceleration method. This table shows the speedup of reaching a given loss at each checkpoint relative to how many steps of pretraining and finetuning are required to reach the same loss. Values lower than 1 indicate that the loss can be reached in fewer steps without data selection. The final bar for DC is shaded to indicate extrapolation and is off the y-axis because the loss is lower than any loss reachable in 1 million steps with pretraining and finetuning.

This figure highlights the effectiveness of pretraining since the performance obtained by data selection for early checkpoints can be matched by simply pretraining longer. Furthermore, DC+FT at 400k pretraining steps cannot be matched, even when pretraining for up to 1m steps. This figure shows that a practitioner with a given generalization requirement can consider data selection early since the target domain generalization gain for early checkpoints might avoid a long pretraining run.

At 50k steps, data selection accelerates training by a factor of about 3.5x, meaning the same performance can be reached with an additional 150k steps of pretraining. However, for later checkpoints, the marginal benefits of pretraining decreases while the improvements from data selection are steady making data selection a clear choice for later checkpoints. In particular for well trained smaller models, such as the LSTM we evaluate for language modeling, the performance after data selection may actually be unreachable just through pretraining either due to the noisiness of the training data that might be filtered from data selection or due to the limited model capacity.

947 A.2 Complementary Finetuning vs 948 Overfitting

949 Figure 4 measures the correlation between the relative
 950 difference between the train and valid best
 951 in-domain loss prior to fine tuning (selection over-
 952 fitting rate) and the relative difference between the
 953 valid loss before and after fine tuning (fine tuning

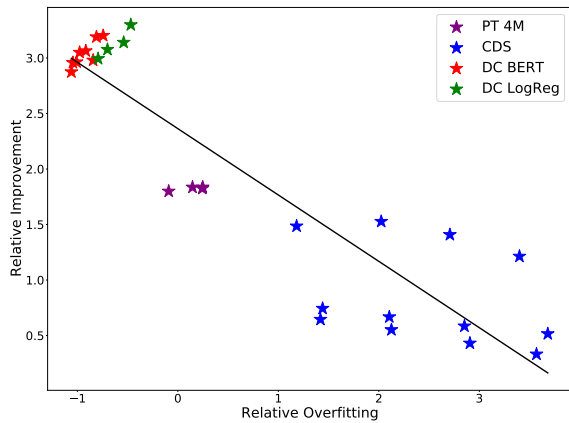


Figure 4: Impact of selection overfitting (MT En-De). When data selection overfits to the in domain set, the improvements from finetuning are lower. The x-axis is the overfitting relative difference and the y-axis is the relative improvement from finetuning. Pearson Correlation Coefficient : -0.91

954 rate). There is a strong anti-correlation between
 955 these factors, showing that overfitting at the selec-
 956 tion stage indeed impacts negatively the impact of
 957 FT. We include points on this graph selecting the
 958 top 4m examples, effectively filtering out the bot-
 959 tom 500k, which has a slight overfitting effect, to in-
 960 clude more points with an intermediate overfitting-
 961 to-improvement tradeoff.