

What matters in a transferable neural network model for relation classification in the biomedical domain?

Sunil Kumar Sahu and Ashish Anand

Department of Computer Science and Engineering

Indian Institute of Technology Guwahati

Assam, India

{sunil.sahu, anand.ashish}@iitg.ernet.in

Abstract

Lack of sufficient labeled data often limits the applicability of advanced machine learning algorithms to real life problems. However efficient use of *Transfer Learning* (TL) has been shown to be very useful across domains. TL utilizes valuable knowledge learned in one task (*source task*), where sufficient data is available, to the task of interest (*target task*). In biomedical and clinical domain, it is quite common that lack of sufficient training data do not allow to fully exploit machine learning models. In this work, we present two unified recurrent neural models leading to three transfer learning frameworks for relation classification tasks. We systematically investigate effectiveness of the proposed frameworks in transferring the knowledge under multiple aspects related to source and target tasks, such as, similarity or relatedness between source and target tasks, and size of training data for source task. Our empirical results show that the proposed frameworks in general improve the model performance, however these improvements do depend on aspects related to source and target tasks. This dependence then finally determine the choice of a particular TL framework.

1 Introduction

Recurrent neural network (RNN) and its variants, such as long short term memory (LSTM) network have shown to be an effective model for many natural language processing tasks (Mikolov et al., 2010; Graves, 2013; Karpathy and Fei-Fei, 2014; Zhang and Wang, 2015; Chiu and Nichols, 2015; Zhou et al., 2016). However, the requirement

of huge gold standard labeled dataset for training makes it difficult to apply them on low resource tasks such as in biomedical domain. In the biomedical domain, obtaining labeled data not only is time consuming and costly but also requires domain knowledge. Transfer learning (TL) has been used successfully in such scenario across multiple domains. The aim of transfer learning is to store the knowledge gained while training a model for a Task-A (*Source Task*), where we have sufficient gold standard labeled data, and apply it to a different Task-B (*Target Task*) where we do not have enough training data (Pan and Yang, 2010). In literature various TL frameworks have been proposed (Pan and Yang, 2010; Mou et al., 2016; Yosinski et al., 2014). With the recent surge in applications of TL using neural network based models in computer vision and image processing (Yosinski et al., 2014; Azizpour et al., 2015) as well as in NLP (Mou et al., 2016; Zoph et al., 2016; Yang et al., 2017), this work explores TL frameworks using neural model for relation classification in biomedical domain.

A very common approach to apply TL is to train learning model on source and target tasks in sequence. We refer to this approach as *sequential TL*. Further if there exists a bijection mapping between label sets of source and target tasks, then the entire model trained on source task can be transferred to the target task, otherwise only partial model can be utilized. In NLP, transferring of feature representation is the most common form of partial model transfer. Instead of performing the training in sequential manner, alternative way could be training the model on both source and target data simultaneously (Yang et al., 2017). This would be very similar to the *multi-task learning* (Collobert and Weston, 2008). This way of simultaneous training can be done in multiple ways. These options give possibilities to design several

variants of TL framework.

Apart from the options of using training data in different ways, using partial or complete model transfer, and presence or absence of bijection mapping between two label sets, other aspects such as *selection of source task, its size and relatedness or similarity with the target task* determine the selection of relevant TL model. Intuitively, it is preferred to have source task as much similar to the target task as we can obtain. For example, if target task is of binary classification of drug-drug interaction (DDI) mentioned in social media text or in doctor’s notes, then we should look for the source task of binary classification of DDI mentioned in research articles. Here, the difference lie in the nature of texts appearing in the two corpora. In the first case of doctor’s notes, text is likely to be short and precise compared to the research articles. In other words feature spaces representing data for source and target tasks differ from each other, although the two label sets are same. On the other hand, it is also possible that there does not exist any bijection between labels of source and target tasks. We can modify our previous example by making the target task as multi-class classification of DDI, to illustrate one such possible scenario.

Given that the various possibilities are arising in light of above discussion, we present two LSTM based models and corresponding three different TL frameworks in this study. Our motivation is to systematically explore various TL frameworks for the task of relation classification in biomedical domain and try to empirically analyze answers to few relevant questions. Our contribution can be summarized as follows:

- We present and evaluate three TL framework variants based on LSTM models for different relation classification tasks of biomedical and clinical text.
- We analyze effect of relatedness (implicit or explicit) between source and target tasks on the effectiveness of TL framework.
- We also explore how the size of the training data corresponding to source task affects effectiveness of TL frameworks.

2 Model Architectures

In this section we first explain a generic architecture of LSTM for relation classification task. Then

we explain three ways of using this architecture for transferring knowledge from source task to target task. We assume that relation exists between two entities, referred to as *target entities*, positions of whom within the sentence are known.

The generic architecture of the neural network for relation classification task can be described in following layers: *word level feature layer, embedding layer, sentence level feature extraction layer, fully connected and softmax layers*. We define features for all words in *word level feature layer*, which also includes some features relative to the two targeted entities. In *embedding layer* every feature gets mapped to a vector representation through a corresponding embedding matrix. Raw features are combined from entire sentence and a fixed length feature representation is obtained in the *sentence level feature extraction layer*. Although a convolution neural network (CNN) or other variants of recurrent neural network can be used in this layer, we use bidirectional LSTM because of its relatively better ability to take into account discontinuous features. *Fully connected and softmax layer* map thus obtained sentence level feature vectors to class probability. In summary, input for these models would be a sentence with the two targeted entities and output would be a probability distribution over each possible relation class between them.

2.1 BLSTM-RE

Suppose $w_1w_2\dots w_m$ is a sentence of length m . Two targeted entities e_1 and e_2 corresponds to some words (or phrases) w_i and w_j respectively. In this work we use word and its position from both targeted entities as features in word level feature layer. Position features are important for relation extraction task because they let model to know the targeted entities (Collobert et al., 2011). Output of *embedding layer* would be a sequence of vectors $x_1x_2\dots x_m$ where $x_i \in \mathbb{R}^{(d_1+d_2+d_3)}$ is concatenation of word and position vectors. d_1, d_2 and d_3 are embedding lengths of word, position from first entity and position from second entity respectively. We use bidirectional LSTM with *max pooling* in the *sentence level feature extraction layer*. This layer is responsible to get optimal fixed length feature vector from entire sentence. Basic architecture is shown in the figure 1a. We omit the mathematical equations as there is no modifications made in the standard bi-directional

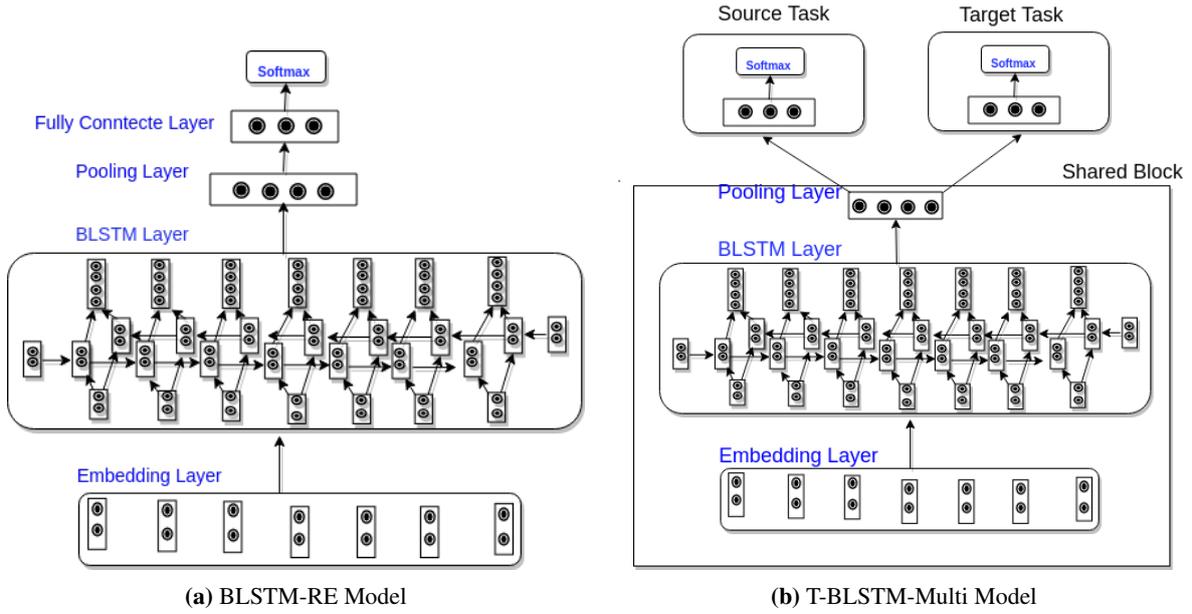


Figure 1: Proposed Model Architect: BLSTM is bidirectional long short term memory network

LSTM model (Graves, 2013).

2.2 *T-BLSTM-Mixed*

T-BLSTM-Mixed is specific way to use *BLSTM-RE* model in transfer learning framework. In this case, instances from both source and target tasks are fed into the same *BLSTM-RE* model. While training we pick one batch of data from source or target in random order with equal probability. Since training is happening simultaneously for both source and target dataset, we can say that model will learn features which is applicable for both. It is quite obvious that this model is applicable for only those cases where bijection mapping between labels of source and target tasks exists.

2.3 *T-BLSTM-Seq*

Convergence of neural network based models depends on initialization of model parameters. Several studies (Hinton et al., 2006; Bengio et al., 2007; Collobert et al., 2011) have shown that initializing parameters with other supervised or unsupervised pre-trained model’s value often improves the model convergence. In this framework of transfer we first train our model with source task’s dataset and use the learned parameters to initialize model parameters for training target task in separate model. We call this framework as *T-BLSTM-Seq*. *T-BLSTM-Seq* can be applicable for both *same label set* as well as *disparate label set* transfer. We transfer entire network parameters if there exists bijection between source and target la-

bel sets, otherwise we only share model parameters up to the second last layer of the network. The left out last layer is randomly initialized.

2.4 *T-BLSTM-Multi*

We propose another transfer learning framework, called as *T-BLSTM-Multi*, using the same backbone of *BLSTM-RE* model. As shown in the figure 1b this model has two *fully connected* and *softmax* layers, one for source task and other is for target task. Other layers of the models are shared for the two tasks. While training, parameters of the shared block get updated with training instances from both source and target data and *fully connected* layer gets updated only with its corresponding task data. Batch of instances are picked in similar manner as *T-BLSTM-Mixed*. This way of training is also called *multi-task learning* but in that case focus is on both source and target task performance. *T-BLSTM-Multi* model is also applicable for both *disparate label set* as well as *same label set* transfer.

2.5 Training and Implementation

Pre-trained word vectors are used for initializing word embeddings and random vectors are used for other feature embedding. We use GloVe (Pennington et al., 2014) on Pubmed corpus for obtaining word vectors. Dimensions of word and position embeddings are set to 100 and 10 respectively. Adam optimization (Kingma and Ba, 2014) is used for training all models. All parameters i.e.,

word embedding, position embeddings and network parameters are being updated during training. We fixed batch size to 100 for all the experiments. In case of *T-BLSTM-Mixed* and *T-BLSTM-Multi* every time we sample one task from source and target task based on binomial distribution, we set binomial probability to half and pick one batch from that task data for training. All the remaining hyperparameters are set according to (Sahu and Anand, 2017). Entire implementation is done in Python language with *TensorFlow*¹ library.

3 Task Definitions and Used Datasets

In this section we briefly describe the tasks and corresponding datasets used in this study. Statistics of these datasets are given in Table 1.

Drug Drug Interaction Extraction (DDI): Drug-drug interaction is a state in which two or more drugs when given to a patient at the same time lead to undesired effects. Identifying DDI present in text is a kind of relation classification task, where given two drugs or pharmacological substances in a sentence we need to classify whether there is an interaction between the two or not. In *Ex.1* drug *Lithium* and *Diuretics* are interacting because they are advised to be not given together.

Ex.1: Lithium_{drug} generally should not be given with diuretics_{drug}

Drug Drug Interaction Class Extraction (DDIC): DDI can appear in text with different semantic senses, which we call as DDI class. In case of DDIC, we need to identify exact class of interactions among drugs in the sentence. For instance in example *Ex.1* type of interaction is *advise* as advice is being given. In SemEval 2013² DDI Extraction task had 4 kinds of interaction *Advise*, *Effect*, *Mechanism* and *Int*.

Adverse Drug Event Extraction (ADE): Adverse drug event is the condition in which an adverse effect happens due to consumption of a drug. In NLP, ADE extraction is the process of extracting adverse relation between a drug and a condition or disease in text. For instance in *Ex.2* for treating patient suffering from *thyrotoxicosis* disease with *methimazole* has led to adverse effect.

Ex.2: A 43 year old woman who was treated for thyrotoxicosis_{Disease} with methimazole_{Drug} developed agranulocytosis

Event Argument Extraction (EAE): In biomedical domain event is broadly described as a change on the state of a bio-molecule or bio-molecules (Pyysalo et al., 2012). Every events have their own set of arguments and EAE is the task of identifying all arguments of an event and their roles. In this task we have entities and triggers (representing an event) present in a sentence are given and the task is to find the role (relation) between all pairs of triggers and entities. For this work we don't differentiate between different types of role. This implies that if an entity is a argument of a trigger then there is positive relation between them otherwise negative. For instance in *Ex.3* (*reptin*, *regulates*), (*regulates*, *growth*) and (*growth*, *heart*) have positive relation.

Ex.3: Reptin_{Protein} regulates_{Regulation} the growth_{Growth} of the heart_{Organ}.

Clinical Relation Extraction (CRE): Clinical relation extraction is the task of identifying relation among clinical entities such as *Problem*, *Treatment* and *Test* in clinical notes or discharge summaries. In *Ex.4* *allergic* and *rash* have *problem improve problem* relation.

Ex.4: She is allergic_{Problem} to augmentin which gives her a rash_{Problem}.

Task	Corpus	Training Set	Test Set
BankDDI	Pairs	14176	3694
	Positive DDIs	3617	884
	Negative DDIs	11559	2810
MedDDI	Pairs	1319	334
	Positive DDIs	227	95
	Negative DDIs	1092	239
BankDDIC	Pairs	14176	3694
	Negative DDIs	11559	2810
	Effect	1471	298
	Mechanism	1203	278
	Advise	813	214
	Int	130	94
MedDDIC	Pairs	1319	334
	Negative DDIs	1092	239
	Effect	149	62
	Mechanism	61	24
	Advise	7	7
	Int	10	2
ADE	Pairs	8867	3802
	Positive ADEs	4177	1791
	Negative ADEs	4690	2011
EAE	Pairs	21594	11443
	Positive EAEs	4492	2202
	Negative EAEs	17102	9241
CRE	Pairs	43602	18690
	Negative CREs	36324	15995
	TeRP	2136	915
	TrAP	1832	784
	PIP	1541	660
	TrCP	368	157
	TeCP	353	150

Table 1: Statistics of all dataset

¹<https://www.tensorflow.org/>

²<https://www.cs.york.ac.uk/semeval-2013/task9/>

3.1 Source Data

BankDDI: It is manually annotated dataset for DDI extraction, collected from Drug Bank³ documents. Drug Bank contains drug information in the form of documents which has been written by medical practitioners. We collected this dataset from SemEval 2011 DDI extraction challenge (Segura Bedmar et al., 2011).

BankDDIC: This dataset is same as BankDDI with task of DDI class recognition (Segura Bedmar et al., 2013).

ADE: ADE extraction dataset was collected from (Gurulingappa et al., 2012b,a). The shared dataset contains manually annotated adverse drug events mentioned in a corpus of Medline abstracts.

EAE: We used MLEE⁴ corpus for event argument identification (Pyysalo et al., 2012). MLEE dataset has 20 types of events trigger and 11 entity types for relation classification.

CRE: For clinical relation classification, we collected dataset from i2b2 2010⁵ clinical information extraction challenge (Uzuner et al., 2011). We consider *TrCP*, *TrAP*, *PIP*, *TeRP* and *TeCP* classes in this case.

3.2 Target Data

MedDDI: MedDDI is manually annotated dataset for DDI extraction, collected from MedLine abstract. This dataset was also shared as part of SemEval-2013 DDI extraction challenge. There are several ways MedDDI is different from BankDDI: MedDDI was collected from MedLine abstracts whereas BankDDI is from DrugBank documents. MedLine abstracts, as being a part of research articles, contains lots of technical terms and usually sentences are longer. On the other hand DrugBank documents contains concise, relatively smaller and easily comprehensible sentences written by medical practitioners. In BankDDI and MedDDI datasets, we removed few negative instances based on same rules used in (Sahu and Anand, 2017; Zhao et al., 2016).

MedDDIC: It is same dataset with labels includes exact class of interaction. Both this dataset was used in SemEval-2013 DDI extraction task.

CRE₅: In this case we take 5% from each class of CRE’s training dataset and considered that as training set. The test set remains same.

³<https://www.drugbank.ca/>

⁴<http://nactem.ac.uk/MLEE/>

⁵<https://www.i2b2.org/NLP/Relations/>

3.3 Preprocessing

We use same preprocessing strategies for all datasets. Pre-processing steps include, all words were converted into lower case form, sentences were tokenized with geniatagger⁶, digits were replaced with *DG* symbol. Further, if any sentence have more than two entities, we create a separate instance for every pair of entities and in all sentences two targeted entities were replaced with their types and position of entity. For example, sentence in *Ex.4* will become *She is ProblemA to augmentin which gives her a ProblemB*. The complete CRE and ADE datasets were separated into training and test sets by randomly selecting 30% instances from each class as test set and remaining as training set.

4 Results and Discussion

First we discuss our experiment design to evaluate performance of the three TL frameworks under various settings. Later we analyze and discuss obtained results.

We treat the performance of bidirectional LSTM model on target task as baseline. In baseline experiments, training was done on the training set of each of the three target data and performance on the respective test sets are then reported in Table 2.

Task	Precision	Recall	F Score
MedDDI	0.561 _(0.03)	0.431 _(0.03)	0.488 _(0.02)
MedDDIC	0.684 _(0.08)	0.273 _(0.01)	0.390 _(0.03)
CRE ₅	0.529 _(0.04)	0.492 _(0.01)	0.510 _(0.009)

Table 2: Baseline Performance: Results of **BLSTM-RE** model on the three different target tasks. Numbers in Precision, Recall and F Score column indicate result corresponding to best F1 Score and subscripts are standard deviation of five runs of model

Needless to mention that baseline model do use pre-trained word embeddings, a form of unsupervised transfer learning framework. As multiple studies have already shown the superior performance of such models using pre-trained vectors than using random vectors, we do not perform any experiment related to that. Five runs with different random initialization were taken for each model and the best result in terms of F1-score along with corresponding precision and recall are shown in Tables. We experiment with different combinations of source and target tasks to analyze the ef-

⁶<http://www.nactem.ac.uk/GENIA/tagger/>

Type	Model	Precision	Recall	F Score	Δ
Similar	<i>T-BLSTM-Mixed</i> _(BankDDI⇒MedDDI)	0.656 _(0.02)	0.705 _(0.03)	0.680 _(0.03)	39.34%
	<i>T-BLSTM-Seq</i> _(BankDDI⇒MedDDI)	0.678 _(0.02)	0.621 _(0.03)	0.648 _(0.03)	32.78%
	<i>T-BLSTM-Multi</i> _(BankDDI⇒MedDDI)	0.701 _(0.05)	0.568 _(0.05)	0.627 _(0.02)	28.48%
	<i>T-BLSTM-Mixed</i> _(BankDDIC⇒MedDDIC)	0.631 _(0.04)	0.505 _(0.02)	0.561 _(0.01)	43.84%
	<i>T-BLSTM-Seq</i> _(BankDDIC⇒MedDDIC)	0.600 _(0.03)	0.463 _(0.02)	0.550 _(0.01)	41.02%
	<i>T-BLSTM-Multi</i> _(BankDDIC⇒MedDDIC)	0.579 _(0.01)	0.421 _(0.006)	0.487 _(0.006)	24.87%
Dissimilar	<i>T-BLSTM-Mixed</i> _(ADE⇒MedDDI)	0.494 _(0.02)	0.515 _(0.03)	0.505 _(0.02)	3.48%
	<i>T-BLSTM-Seq</i> _(ADE⇒MedDDI)	0.595 _(0.02)	0.294 _(0.03)	0.394 _(0.02)	-19.26%
	<i>T-BLSTM-Multi</i> _(ADE⇒MedDDI)	0.533 _(0.02)	0.505 _(0.02)	0.518 _(0.01)	6.14%
	<i>T-BLSTM-Mixed</i> _(EAE⇒MedDDI)	0.540 _(0.03)	0.557 _(0.05)	0.549 _(0.02)	12.5%
	<i>T-BLSTM-Seq</i> _(EAE⇒MedDDI)	0.544 _(0.03)	0.515 _(0.04)	0.529 _(0.02)	8.40%
	<i>T-BLSTM-Multi</i> _(EAE⇒MedDDI)	0.538 _(0.02)	0.589 _(0.05)	0.562 _(0.02)	15.16%

Table 3: Results of TL frameworks in case of **same label set transfer**. Here ($X \Rightarrow Y$) indicates transferring from X dataset to Y dataset and $Type$ indicates nature of source and target dataset. Numbers in Precision, Recall and F Score column indicate result corresponding to best F1 Score and subscripts are standard deviation of five runs of model. Δ is relative percentage improvement over baseline (without TL) method

fect of similarity of feature spaces corresponding to source and target tasks as well as their label sets on the choice of TL frameworks.

4.1 Performance on Same Label Set Transfer

Let’s first look at the relative improvement of various TL models over the baseline results on DDI and DDIC tasks. Table 3 shows the performance of all TL models on these two tasks under various settings. $Type$ in the Table 3 indicates semantic relatedness of source and target tasks. For example, data in both *BankDDI* and *MedDDI* indicate drug-drug interaction, and hence are of the same semantic type. But *EAE* gives existence of trigger argument relation which is not of the same semantic type as drug-drug interaction, although both tasks fall into binary classification. As the results indicate, *T-BLSTM-Mixed* model gave the best performance (in terms of F1-score) for the *similar* type tasks, whereas *T-BLSTM-Multi* gave the worst. However all the TL models gave significant improvement over the baseline results. *T-BLSTM-Mixed* obtained approximately 40% relative improvement over the baseline for the DDI task and approximately 44% for the DDIC task. On the other hand, *T-BLSTM-Multi* gave best performance for the *dissimilar* type tasks and *T-BLSTM-Seq* gave the worst. In fact, *T-BLSTM-Seq* gave relatively poor performance than baseline in one case.

4.2 Performance on Disparate Label Set Transfer

Next we examine performance of relevant TL models when there does not exist a bijection be-

tween the two label sets corresponding to source and target tasks (Table 4). As *T-BLSTM-Mixed*, by design, require the existence of bijection between the two label sets, we exclude this model for this case. Among the rest two, *T-BLSTM-Multi* always led to significantly improved performance compared to the respective baseline results. On the other hand, performance of the *T-BLSTM-Seq* model is not so consistent specially when *ADE* was used as source data.

4.3 Analyzing Similarity between Source and Target Tasks

Earlier we observe that similarity between source and target tasks affects the relative performance of each TL framework. Let us try to analyze this observation. *T-BLSTM-Mixed* and *T-BLSTM-Seq* models transfer full knowledge or in other words, both model share the complete model between source and target tasks. This allow the last layers to see more examples and to be adaptive to samples from both source and target data. On the other hand, *T-BLSTM-Multi* only share the partial model upto the second last layers. In this case, last layer for source and target tasks are trained separately and is not being shared between the two. Thus the last layers are being specific to the respective tasks. When there is similarity between source and target tasks, as well as there exists a bijection, target tasks gets benefited by sharing full model. In such scenario co-training seems better suited as *T-BLSTM-Mixed* was found to be the best among three frameworks. On the other hand, *T-BLSTM-Multi* fail to exploit the full knowledge present in

Source⇒Target	<i>T-BLSTM-Seq</i>			<i>T-BLSTM-Multi</i>		
	Precision	Recall	F Score	Precision	Recall	F Score
BankDDI⇒MedDDIC	0.50 _(0.08)	0.378 _(0.03)	0.431 _(0.008)	0.603 _(0.05)	0.368 _(0.03)	0.457 _(0.03)
CRE⇒MedDDIC	0.448 _(0.05)	0.368 _(0.03)	0.404 _(0.02)	0.468 _(0.02)	0.389 _(0.02)	0.425 _(0.01)
EAE⇒MedDDIC	0.596 _(0.04)	0.326 _(0.03)	0.421 _(0.02)	0.488 _(0.03)	0.452 _(0.06)	0.469 _(0.04)
ADE⇒MedDDIC	0.512 _(0.06)	0.221 _(0.01)	0.308 _(0.01)	0.447 _(0.04)	0.400 _(0.03)	0.422 _(0.02)
BankDDI⇒CRE ₅	0.555 _(0.02)	0.485 _(0.01)	0.518 _(0.01)	0.546 _(0.01)	0.508 _(0.02)	0.526 _(0.006)
BankDDIC⇒CRE ₅	0.564 _(0.04)	0.447 _(0.02)	0.498 _(0.01)	0.523 _(0.01)	0.557 _(0.02)	0.539 _(0.007)
EAE⇒CRE ₅	0.543 _(0.02)	0.533 _(0.02)	0.538 _(0.006)	0.587 _(0.01)	0.548 _(0.01)	0.567 _(0.01)
ADE⇒CRE ₅	0.516 _(0.003)	0.503 _(0.01)	0.509 _(0.006)	0.598 _(0.03)	0.483 _(0.02)	0.535 _(0.007)
BankDDIC⇒MedDDI	0.605 _(0.04)	0.452 _(0.03)	0.518 _(0.01)	0.623 _(0.04)	0.557 _(0.03)	0.588 _(0.02)
CRE⇒MedDDI	0.569 _(0.11)	0.473 _(0.17)	0.517 _(0.03)	0.631 _(0.05)	0.505 _(0.02)	0.561 _(0.02)

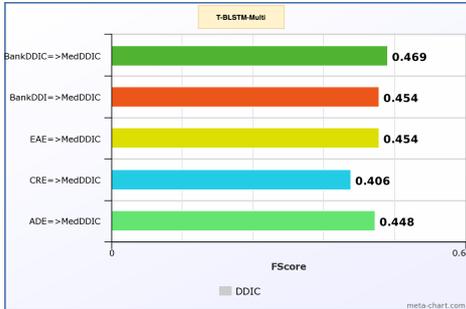
Table 4: Results of *T-BLSTM-Seq* and *T-BLSTM-Multi* on **disparate label set transfer** task. Here ($X \Rightarrow Y$) indicates transferring from X dataset to Y dataset. Numbers in Precision, Recall and F Score column indicate result corresponding to best F1 Score and subscripts are standard deviation of 5 runs of model.

the training data of source task. But this becomes advantageous for the *T-BLSTM-Multi* framework in case of absence of bijection between source and target label sets. The last layer, in *T-BLSTM-Multi*, takes the shared knowledge and tune it into the specific target task. This observation also fits well with the observations made earlier in (Yosinski et al., 2014) that the initial layers are relatively generic and become more specific as we go towards the last layer.

performance difference. Hence to take this effect out from the consideration, all source data was made of the same size (8867) as the minimum among all source training data. During random selection, proportion of instances from all classes were maintained as in the original set. We have shown only the results obtained by *T-BLSTM-Multi* in Figure 2 but similar results are obtained for other models as well. We observe the performance obtained from using different source data but of same size match with performance obtained with the same set of source data but of different sizes. This indicates that context and label mapping played more crucial role than size of selected source data.



(a) T-BLSTM-Multi Model (DDI Extraction)



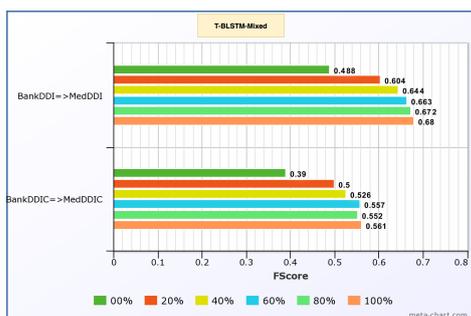
(b) T-BLSTM-Multi : Model (DDIC Extraction)

Figure 2: Performance of proposed models with different source task on same size data

Variability in training set size of different source data could have influenced the observed perfor-

4.4 Analyzing Size of Source Task Dataset

One of the important arguments generally given for the use of transfer learning is insufficient dataset for the target task hinders the performance of learning algorithms. Performance can be enhanced by utilizing information available in relatively higher amount of source data. In this section we investigate the effect of different size of source data on the performance improvement of the *T-BLSTM-Mixed* and *T-BLSTM-Multi* models. Figure 3 shows the results on both similar and dissimilar tasks. In both scenario, even having 20% of source data significantly improves the performance. However, there is a consistent increasing trend in improvement observed for the similar tasks, whereas performance was found to be little fluctuating for the dissimilar tasks. The fluctuation could be due to too much of source data may confuse the model.



(a) T-BLSTM-Mixed : transfer on similar task



(b) T-BLSTM-Multi : transfer on dissimilar task

Figure 3: Performance of proposed models with different training set size corresponding to source task

4.5 Comparison with State-of-art Results

At the end we compare our results with the state-of-the-art results obtained on target tasks of *DDI* and *DDIC*. Table 5 shows best results of the SemEval 2010 DDI extraction challenge (Segura Bedmar et al., 2013) as well as the results obtained by BLSTM-RE and *T-BLSTM-Mixed* models on dissimilar tasks. We can observe that although BLSTM-RE can not outperform the best results of the challenge but under the TL framework, *T-BLSTM-Multi* even using dissimilar tasks improved the state-of-the-art results.

Models	DDIC	DDI
FBK-Irst(Chowdhury and Lavelli, 2013)	0.398	0.530
SCAI(Bobic et al.)	0.420	0.47
WBI(Thomas et al., 2013)	0.365	0.503
UTurku(Björne et al., 2013)	0.286	0.479
UMAD(Rastegar-Mojarad et al., 2013)	0.312	0.479
BLSTM-RE	0.390	0.488
<i>T-BLSTM-Multi</i> _{EAE=>⊕}	0.469	0.562
<i>T-BLSTM-Multi</i> _{CRE=>⊕}	0.425	0.561
<i>T-BLSTM-Multi</i> _{ADE=>⊕}	0.422	0.518

Table 5: Performance comparison of existing methods for DDI and DDIC task. Here values indicate **F1 Score** of both task. ⊕ is MedDDI or MedDDIC

5 Related Work

Recurrent neural network and its variants have been successfully applied to many semantic relation classification tasks. Authors in (Zhang and Wang, 2015; Zhang et al., 2015) have used recurrent neural network with max pooling operation while Zhou et al. (2016) used RNN with attentive pooling for the same relation classification task. Li et al. (2016) partitioned the sentence with targeted entities and a separate RNN was trained to obtain division specific features and use them for classification. Li et al. (2016) used recurrent neural network for semantic relations as well as Bio Event argument relation extraction tasks. Although none of these works have used transfer learning but have used RNN models for various relation classification tasks.

The proposed TL frameworks are closely related to the works of (Yang et al., 2017; Mou et al., 2016; Collobert and Weston, 2008). (Yang et al., 2017) have introduced variety of TL frameworks using gated recurrent neural network (GRU). They have evaluated the proposed frameworks on different sequence labeling tasks, such as *PoS tagging* and *chunking*. Mou et al. (2016), similar to the study by Yosinski et al. (2014) for image processing tasks, evaluated CNN and RNN based TL frameworks for sentence classification and sentence pair modeling tasks. Collobert and Weston (2008) have used window based neural network and convolution neural networks for several sequence labeling tasks in the multi-task learning framework. Zoph et al. (2016) have explored transfer learning for neural machine translation tasks. They have shown significant improvement in many low resource language translation tasks. Their model repurpose the learned model, trained on high resource language translation dataset (source task), for target task.

6 Conclusions

In this work we present various transfer learning frameworks based on LSTM models for relation classification task in biomedical domain. We observe that in general transfer learning do help in improving the performance. However, similarity of source tasks with the target task as well as size of corresponding source data affects the performance and hence plays important role in selection of appropriate TL framework.

References

- Hossein Azizpour, Ali Sharif Razavian, Josephine Sullivan, Atsuto Maki, and Stefan Carlsson. 2015. From generic to specific deep representations for visual recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*. pages 36–45.
- Yoshua Bengio, Pascal Lamblin, Dan Popovici, Hugo Larochelle, et al. 2007. Greedy layer-wise training of deep networks. *Advances in neural information processing systems* 19:153.
- Jari Björne, Suwisa Kaewphan, and Tapio Salakoski. 2013. Uturku: drug named entity recognition and drug-drug interaction extraction using svm classification and domain knowledge. In *Proceedings of the 7th International Workshop on Semantic Evaluation (SemEval 2013)*. pages 651–659.
- Tamara Bobic, Juliane Fluck, and Martin Hofmann-Apitius. 2013. Scai: Extracting drug-drug interactions using a rich feature vector .
- Jason P. C. Chiu and Eric Nichols. 2015. [Named entity recognition with bidirectional lstm-cnns](#). *CoRR* abs/1511.08308. <http://arxiv.org/abs/1511.08308>.
- Md Faisal Mahbub Chowdhury and Alberto Lavelli. 2013. Fbk-irst: a multi-phase kernel based approach for drug-drug interaction detection and classification that exploits linguistic information. *Atlanta, Georgia, USA* 351:53.
- Ronan Collobert and Jason Weston. 2008. A unified architecture for natural language processing: Deep neural networks with multitask learning. In *Proceedings of the 25th international conference on Machine learning*. ACM, pages 160–167.
- Ronan Collobert, Jason Weston, Léon Bottou, Michael Karlen, Koray Kavukcuoglu, and Pavel Kuksa. 2011. [Natural language processing \(almost\) from scratch](#). *J. Mach. Learn. Res.* 12:2493–2537. <http://dl.acm.org/citation.cfm?id=1953048.2078186>.
- Alex Graves. 2013. Generating sequences with recurrent neural networks. *arXiv preprint arXiv:1308.0850* .
- Harsha Gurulingappa, Abdul Mateen-Rajpu, and Luca Toldo. 2012a. Extraction of potential adverse drug events from medical case reports. *Journal of biomedical semantics* 3(1):1.
- Harsha Gurulingappa, Abdul Mateen Rajput, Angus Roberts, Juliane Fluck, Martin Hofmann-Apitius, and Luca Toldo. 2012b. Development of a benchmark corpus to support the automatic extraction of drug-related adverse effects from medical case reports. *Journal of biomedical informatics* 45(5):885–892.
- Geoffrey E Hinton, Simon Osindero, and Yee-Whye Teh. 2006. A fast learning algorithm for deep belief nets. *Neural computation* 18(7):1527–1554.
- Andrej Karpathy and Li Fei-Fei. 2014. Deep visual-semantic alignments for generating image descriptions. *arXiv preprint arXiv:1412.2306* .
- Diederik Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* .
- Fei Li, Meishan Zhang, Guohong Fu, Tao Qian, and Donghong Ji. 2016. [A bi-lstm-rnn model for relation classification using low-cost sequence features](#). *CoRR* abs/1608.07720. <http://arxiv.org/abs/1608.07720>.
- Tomas Mikolov, Martin Karafiát, Lukás Burget, Jan Cernocký, and Sanjeev Khudanpur. 2010. Recurrent neural network based language model. In *INTERSPEECH 2010, 11th Annual Conference of the International Speech Communication Association, Makuhari, Chiba, Japan, September 26-30, 2010*. pages 1045–1048.
- Lili Mou, Zhao Meng, Rui Yan, Ge Li, Yan Xu, Lu Zhang, and Zhi Jin. 2016. How transferable are neural networks in nlp applications? *arXiv preprint arXiv:1603.06111* .
- Sinno Jialin Pan and Qiang Yang. 2010. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering* 22(10):1345–1359.
- Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of EMNLP*.
- Sampo Pyysalo, Tomoko Ohta, Makoto Miwa, Han-Cheol Cho, Jun’ichi Tsujii, and Sophia Ananiadou. 2012. Event extraction across multiple levels of biological organization. *Bioinformatics* 28(18):i575–i581.
- Majid Rastegar-Mojarad, Richard D Boyce, and Rashmi Prasad. 2013. Uwm-triads: classifying drug-drug interactions with two-stage svm and post-processing. In *Proceedings of the 7th International Workshop on Semantic Evaluation*. pages 667–674.
- Sunil Kumar Sahu and Ashish Anand. 2017. [Drug-drug interaction extraction from biomedical text using long short term memory network](#). *CoRR* abs/1701.08303. <http://arxiv.org/abs/1701.08303>.
- Isabel Segura Bedmar, Paloma Martínez, and María Herrero Zazo. 2013. Semeval-2013 task 9: Extraction of drug-drug interactions from biomedical texts (ddiextraction 2013). Association for Computational Linguistics.
- Isabel Segura Bedmar, Paloma Martinez, and Daniel Sánchez Cisneros. 2011. The 1st ddiextraction-2011 challenge task: Extraction of drug-drug interactions from biomedical texts .
- Philippe Thomas, Mariana Neves, Tim Rocktäschel, and Ulf Leser. 2013. Wbi-ddi: drug-drug interaction extraction using majority voting. In *Second Joint*

*Conference on Lexical and Computational Semantics (*SEM)*. volume 2, pages 628–635.

Özlem Uzuner, Brett R South, Shuying Shen, and Scott L DuVall. 2011. 2010 i2b2/va challenge on concepts, assertions, and relations in clinical text. *Journal of the American Medical Informatics Association* 18(5):552–556.

Zhilin Yang, Ruslan Salakhutdinov, and William W. Cohen. 2017. Transfer learning for sequence tagging with hierarchical recurrent networks. ICLR-2017.

Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. 2014. How transferable are features in deep neural networks? In *Advances in neural information processing systems*. pages 3320–3328.

Dongxu Zhang and Dong Wang. 2015. Relation classification via recurrent neural network. *arXiv preprint arXiv:1508.01006* .

Shu Zhang, Dequan Zheng, Xinchun Hu, and Ming Yang. 2015. Bidirectional long short-term memory networks for relation classification. In *PACLIC*.

Zehuan Zhao, Zhihao Yang, Ling Luo, Hongfei Lin, and Jian Wang. 2016. Drug drug interaction extraction from biomedical literature using syntax convolutional neural network. *Bioinformatics* page btw486.

Peng Zhou, Wei Shi, Jun Tian, Zhenyu Qi, Bingchen Li, Hongwei Hao, and Bo Xu. 2016. Attention-based bidirectional long short-term memory networks for relation classification. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. Association for Computational Linguistics, Berlin, Germany, pages 207–212. <http://anthology.aclweb.org/P16-2034>.

Barret Zoph, Deniz Yuret, Jonathan May, and Kevin Knight. 2016. Transfer learning for low-resource neural machine translation. *arXiv preprint arXiv:1604.02201* .