# ThisIsCompetition at SemEval-2019 Task 9:
# BERT is unstable for out-of-domain samples

**Cheoneum Park**[*1] and **Juae Kim**[*2] and **Hyeon-gu Lee**[*1] and **Reinald Kim Amplayo**[*3]
**Harksoo Kim**[1] and **Jungyun Seo**[2] and **Changki Lee**[1]
**(* equal contribution)**
[1]Kangwon National University, South Korea
[2]Sogang University, South Korea
[3]University of Edinburgh, UK

## Abstract

This paper describes our system, Joint Encoders for Stable Suggestion Inference (**JESSI**), for the SemEval 2019 Task 9: Suggestion Mining from Online Reviews and Forums. JESSI is a combination of two sentence encoders: (a) one using multiple pre-trained word embeddings learned from log-bilinear regression (GloVe) and translation (CoVe) models, and (b) one on top of word encodings from a pre-trained deep bidirectional transformer (BERT). We include a domain adversarial training module when training for out-of-domain samples. Our experiments show that while BERT performs exceptionally well for in-domain samples, several runs of the model show that it is unstable for out-of-domain samples. The problem is mitigated tremendously by (1) combining BERT with a non-BERT encoder, and (2) using an RNN-based classifier on top of BERT. Our final models obtained second place with 77.78% F-Score on Subtask A (i.e. in-domain) and achieved an F-Score of 79.59% on Subtask B (i.e. out-of-domain), even without using any additional external data.

## 1 Introduction

Opinion mining (Pang and Lee, 2007) is a huge field that covers many NLP tasks ranging from sentiment analysis (Liu, 2012), aspect extraction (Mukherjee and Liu, 2012), and opinion summarization (Ku et al., 2006), among others. Despite the vast literature on opinion mining, the task on suggestion mining has given little attention. Suggestion mining (Brun and Hagège, 2013) is the task of collecting and categorizing suggestions about a certain product. This is important because while opinions indirectly give hints on how to improve a product (e.g. analyzing reviews), suggestions are direct improvement requests (e.g. tips, advice, recommendations) from people who have used the product.

To this end, Negi et al. (2019) organized a shared task specifically on suggestion mining called SemEval 2019 Task 9: Suggestion Mining from Online Reviews and Forums. The shared task is composed of two subtasks, Subtask A and B. In Subtask A, systems are tasked to predict whether a sentence of a certain domain (i.e. electronics) entails a suggestion or not given a training data of the same domain. In Subtask B, systems are tasked to do suggestion prediction of a sentence from another domain (i.e. hotels). Organizers observed four main challenges: (a) sparse occurrences of suggestions; (b) figurative expressions; (c) different domains; and (d) complex sentences. While previous attempts (Ramanand et al., 2010; Brun and Hagège, 2013; Negi and Buitelaar, 2015) made use of human-engineered features to solve this problem, the goal of the shared task is to leverage the advancements seen on neural networks, by providing a larger dataset to be used on data-intensive models to achieve better performance.

This paper describes our system **JESSI** (Joint Encoders for Stable Suggestion Inference). JESSI is built as a combination of two neural-based encoders using multiple pre-trained word embeddings, including BERT (Devlin et al., 2018), a pre-trained deep bidirectional transformer that is recently reported to perform exceptionally well across several tasks. The main intuition behind JESSI comes from our finding that although BERT gives exceptional performance gains when applied to in-domain samples, it becomes unstable when applied to out-of-domain samples, even when using a domain adversarial training (Ganin et al., 2016) module. This problem is mitigated using two tricks: (1) jointly training BERT with a CNN-based encoder, and (2) using an RNN-based encoder on top of BERT before feeding to the classifier.

JESSI is trained using only the datasets given on the shared task, without using any additional external data. Despite this, JESSI performs second on Subtask A with an F1 score of 77.78% among 33 other team submissions. It also performs well on Subtask B with an F1 score of 79.59%.

## 2 Related Work

**Suggestion Mining**  The task of detecting suggestions in sentences is a relatively new task, first mentioned in Ramanand et al. (2010) and formally defined in Negi and Buitelaar (2015). Early systems used manually engineered patterns (Ramanand et al., 2010) and rules (Brun and Hagège, 2013), and linguistically motivated features (Negi and Buitelaar, 2015) trained on a supervised classifier (Negi et al., 2016). Automatic mining of suggestion has also been suggested (Dong et al., 2013). Despite the recent successes of neural-based models, only few attempts were done, by using neural network classifiers such as CNNs and LSTMs (Negi et al., 2016), by using part-of-speech embeddings to induce distant supervision (Negi and Buitelaar, 2017). Since neural networks are data hungry models, a large dataset is necessary to optimize the parameters. SemEval 2019 Task 9 (Negi et al., 2019) enables training of deeper neural models by providing a much larger training dataset.

**Domain Adaptation**  In text classification, training and test data distributions can be different, and thus domain adaptation techniques are used. These include non-neural methods that map the semantics between domains by aligning the vocabulary (Basili et al., 2009; Pan et al., 2010) and generating labeled samples (Wan, 2009; Yu and Jiang, 2016). Neural methods include the use of stacked denoising autoencoders (Glorot et al., 2011), variational autoencoders (Saito et al., 2017; Ruder and Plank, 2018). Our model uses a domain adversarial training module (Ganin et al., 2016), an elegant way to effectively transfer knowledge between domains by training a separate domain classifier using an adversarial objective.

**Language Model Pretraining**  Inspired from the computer vision field, where ImageNet (Deng et al., 2009) is used to pretrain models for other tasks (Huh et al., 2016), many recent attempts in the NLP community are successful on using language modeling as a pretraining step to extract
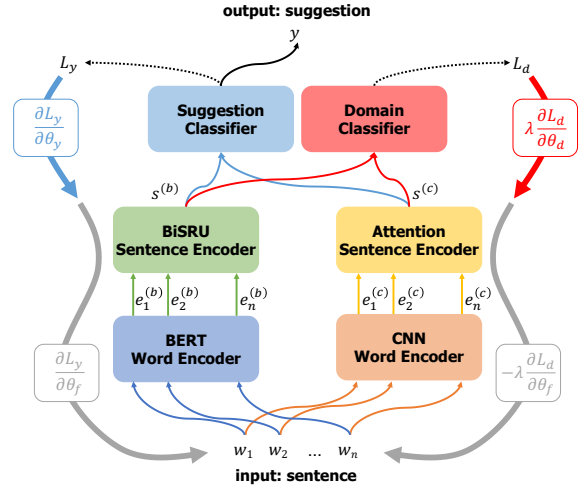


Figure 1: The overall architecture of JESSI for Subtask B. The thinner arrows correspond to the forward propagations, while the thicker arrows correspond to the backward propagations, where gradient calculations are indicated. For Subtask A, a CNN encoder is used instead of the BiSRU encoder, and the domain adversarial training module is not used.

feature representations (Peters et al., 2018), and to fine-tune NLP models (Radford et al., 2018; Devlin et al., 2018). BERT (Devlin et al., 2018) is the most recent inclusion to these models, where it uses a deep bidirectional transformer trained on masked language modeling and next sentence prediction objectives. Devlin et al. (2018) reported that BERT shows significant increase in improvements on many NLP tasks, and subsequent studies have shown that BERT is also effective on harder tasks such as open-domain question answering (Yang et al., 2019), multiple relation extraction (Wang et al., 2019), and table question answering (Hwang et al., 2019), among others. In this paper, we also use BERT as an encoder, show its problem on out-of-domain samples, and mitigate the problem using multiple tricks.

## 3 Joint Encoders for Stable Suggestion Inference

We present our model **JESSI**, which stands for Joint Encoders for Stable Suggestion Inference, shown in Figure 1. Given a sentence $x = \{w_1, w_2, ..., w_n\}$, JESSI returns a binary suggestion label $y = \{0, 1\}$. JESSI consists of four important components: (1) A BERT-based encoder that leverages general knowledge acquired from

a large pre-trained language model, (2) A CNN-based encoder that learns task-specific sentence representations, (3) an MLP classifier that predicts the label given the joint encodings, and (4) a domain adversarial training module that prevents the model to distinguish between the two domains.

**BERT-based Encoder** Fine-tuning a pre-trained BERT (Devlin et al., 2018) classifier then using the separately produced classification encoding [CLS] has shown to produce significant improvements. Differently, JESSI uses a pre-trained BERT as a word encoder, that is instead of using [CLS], we use the word encodings $e_1^{(b)}, e_2^{(b)}, ..., e_n^{(b)}$ produced by BERT. BERT is still fine-tuned during training.

We append a sentence encoder on top of BERT, that returns a sentence representation $s^{(b)}$, which is different per subtask. For Subtask A, we use a CNN encoder with max pooling (Kim, 2014) to create the sentence embedding. For Subtask B, we use a bidirectional simple recurrent units (Lei et al., 2018, BiSRU), a type of RNN that is highly parallelizable, as the sentence encoder.

**CNN-based Encoder** To make the final classifier more task-specific, we use a CNN-based encoder that is trained from scratch. Specifically, we employ a concatenation of both pre-trained GloVe (Pennington et al., 2014) and CoVe (McCann et al., 2017) word embeddings as input $w_i, 1 \leq i \leq n$. Then, we do convolution operations $\text{Conv}(w_i, h_j)$ using multiple filter sizes $h_j$ to a window of $h_j$ words. We use different paddings for different filter sizes such that the number of output for each convolution operation is $n$. Finally, we concatenate the outputs to obtain the word encodings, i.e. $e_i^{(c)} = \oplus_j(\text{Conv}(w_i, h_j))$, where $\oplus$ is the sequence concatenate operation.

We pool the word encodings using attention mechanism to create a sentence representation $s^{(c)}$. That is, we calculate attention weights using a latent variable $v$ that measures the importance of the words $e_i^{(c)}$, i.e., $a_i = \text{softmax}(v^\top f(e_i^{(c)}))$, where $f(\cdot)$ is a nonlinear function. We then use $a_i$ to weight-sum the words into one encoding, i.e., $s^{(c)} = \sum_i a_i e_i^{(c)}$.

**Suggestion Classifier** Finally, we use a multi-layer perceptron (MLP) as our classifier, using a concatenation of outputs from both the BERT- and CNN-based encoders, i.e., $p(y) =$ MLP$_y([s^{(b)}; s^{(c)}])$. Training is done by minimizing the cross entropy loss, i.e., $\mathbb{L} = -\log p(y)$.

**Domain Adversarial Training** For Subtask B, the model needs to be able to classify out-of-domain samples. Using the model as is decreases performance significantly because of cross-domain differences. To this end, we use a domain adversarial training module (Ganin et al., 2016) to prevent the classifier on distinguishing differences between domains. Specifically, we create another MLP classifier that classifies the *domain* of the text using the concatenated sentence encoding with *reverse gradient function* GradRev$(\cdot)$, i.e., $p(d) =$ MLP$_d(\text{GradRev}([s^{(b)}; s^{(c)}]))$. The reverse gradient function is a function that performs equivalently with the identity function when propagating forward, but reverses the sign of the gradient when propagating backward.

Through this, we eliminate the possible ability of the classifier to distinguish the domains of the text. We train the domain classifier using the available trial datasets for each domain. We also use a cross entropy loss as the objective of this classifier. Overall, the objective of JESSI is to minimize the following loss: $\mathbb{L} = -\log p(y) - \lambda \log p(d)$, where $\lambda$ is set increasingly after each epoch, following Ganin et al. (2016).

## 4 Experimental Setup

**Dataset and Preprocessing** We use the dataset provided in the shared task: a training dataset from the electronics domain, and labeled trial and unlabeled test datasets from both the electronics and hotels domain. Table 1 summarizes the dataset statistics and shows the distribution differences between two domains. During training, we use the labeled training dataset to train the suggestions classifier, and trial datasets, without the suggestion labels, to train the domains classifier. For preprocessing, we lowercased and tokenized using the Stanford CoreNLP toolkit[1] (Manning et al., 2014).

**Implementation** We use the pre-trained BERT models[2] provided by the original authors to initialize the parameters of BERT. We use BERT-large

---

[1] https://stanfordnlp.github.io/CoreNLP/
[2] https://github.com/google-research/bert

| Subtask | A | B |
|---|---|---|
| Domain | Electronics | Hotels |
| #Training | 8,230 | 0 |
| #Trial | 592 | 808 |
| #Test | 833 | 824 |
| #Vocabulary | 10,897 | 3,570 |
| Ave. Tokens | 19.0 | 16.8 |

Table 1: Dataset Statistics

for Subtask A and BERT-base for Subtask B[3]. For our CNNs, we use three filters with sizes $\{3, 5, 7\}$, each with 200 dimensions. For the BiSRU, we use hidden states with 150 dimensions and stack with two layers. The MLP classifier contains two hidden layers with 300 dimensions.

We use dropout (Srivastava et al., 2014) on all nonlinear connections with a dropout rate of 0.5. We also use an $l_2$ constraint of 3. During training, we use mini-batch size of 32. Training is done via stochastic gradient descent over shuffled mini-batches with the Adadelta (Zeiler, 2012) update rule. We perform early stopping using the trial sets. Moreover, since the training set is relatively small, multiple runs lead to different results. To handle this, we perform an ensembling method as follows. We first run 10-fold validation over the training data, resulting into ten different models. We then pick the top three models with the highest performances, and pick the class with the most model predictions.

## 5 Experiments

In this section, we show our results and experiments. We denote JESSI-A as our model for Subtask A (i.e., BERT→CNN+CNN→ATT), and JESSI-B as our model for Subtask B (i.e., BERT→BiSRU+CNN→ATT+DomAdv). The performance of the models is measured and compared using the F1-score.

**Ablation Studies** We present in Table 2 ablations on our models. Specifically, we compare JESSI-A with the same model, but without the CNN-based encoder, without the BERT-based encoder, and with the CNN sentence encoder of the BERT-based encoder replaced with the BiSRU variant. We also compare JESSI-B with the same

| Model | F-Score |
|---|---|
| JESSI-A | 88.78 |
| + BERT→BiSRU | 86.01 |
| – CNN→ATT | 85.14 |
| – BERT→CNN | 83.89 |

(a) Subtask A

| Model | F-Score |
|---|---|
| JESSI-B | 87.31 |
| – CNN→ATT | 84.01 |
| – BERT→BiSRU | 81.13 |
| + BERT→CNN | 70.21 |
| – DomAdv | 47.48 |

(b) Subtask B

Table 2: Ablation results for both subtasks using the provided trial sets. The + denotes a *replacement* of the BERT-based encoder, while the – denotes a *removal* of a specific component.

model, but without the CNN-based encoder, without the BERT-based encoder, without the domain adversarial training module, and with the BiSRU sentence encoder of the BERT-based encoder replaced with the CNN variant. The ablation studies show several observations. First, jointly combining both BERT- and CNN-based encoders help improve the performance on both subtasks. Second, the more effective sentence encoder for the BERT-based encoder (i.e., CNN versus BiSRU) differs for each subtask; the CNN variant is better for Subtask A, while the BiSRU variant is better for Subtask B. Finally, the domain adversarial training module is very crucial in achieving a significant increase in performance.

**Out-of-Domain Performance** During our experiments, we noticed that BERT is unstable when predicting out-of-domain samples, even when using the domain adversarial training module. We show in Table 3 the summary statistics of the F-Scores of 10 runs of the following models: (a) vanilla BERT that uses the [CLS] classification encoding, (b-c) our BERT-based encoders BERT→CNN and BERT→BiSRU that use BERT as a word encoder and use an additional CNN/BiSRU as a sentence encoder, (d) JESSI-B that uses BERT→BiSRU and CNN→ATT as joint encoders, and (e) CNN→ATT that does not employ BERT in any way. The results show that while CNN→ATT performs similarly on different runs, BERT performs very unsta-

| Model | min | max | mean | std |
|---|---|---|---|---|
| BERT | 0.00 | 70.59 | 22.52 | 31.0 |
| BERT→CNN | 0.00 | 74.62 | 28.23 | 34.1 |
| BERT→BɪSRU | 54.00 | 88.83 | 74.86 | 8.8 |
| JESSI-B | 69.28 | 89.21 | 82.41 | 5.6 |
| CNN→Aᴛᴛ | 68.19 | 77.06 | 72.50 | 2.5 |

Table 3: Summary statistics of the F-Scores of 10 runs of different models on the trial set of Subtask B when doing a 10-fold validation over the available training data. All models include the domain adversarial training module (+DᴏᴍAᴅᴠ), which is omitted for brevity.

bly, achieving varying F-Scores as low as zero and as high as 70.59, with a standard deviation of 31. Appending a CNN-based sentence encoder (i.e., BERT→CNN) increases the performance, but worsens the stability of the model. Appending an RNN-based sentence encoder (i.e., BERT→BɪSRU) both increases the performance and improves the model stability. Finally, combining a separate CNN-based encoder (i.e., JESSI-B) improves the performance and stability further.

**Test Set Results** Table 4 presents how JESSI compared to the top performing models during the competition proper. Overall, JESSI-A ranks second out of 33 official submissions with an F-Score of 77.78%. Although we were not able to submit JESSI-B during the submission phase, JESSI-B achieves an F-Score of 79.59% on the official test set. This performance is similar to the performance of the model that obtained sixth place in the competition. We emphasize that JESSI does not use any labeled and external data for Subtask B, and thus is just exposed to the hotels domain using the available *unlabeled* trial dataset, containing 808 data instances. We expect the model to perform better when additional data from the hotels domain.

**Performance by Length** We compare the performance of models on data with varying lengths to further investigate the increase in performance of JESSI over other models. More specifically, for each range of sentence length (e.g., from 10 to 20), we look at the accuracy of JESSI-A, BERT→BɪSRU, and BERT→CNN on Subtask A, and the accuracy of JESSI-B, BERT→BɪSRU, and BERT→CNN, all with domain adversarial training module, on Subtask B. Figure 2 shows the
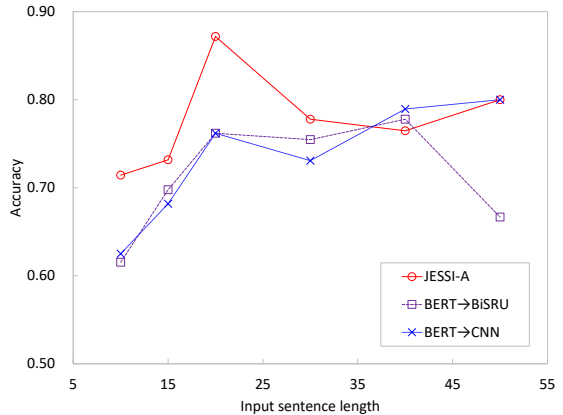
| Rank | Model | F-Score |
|---|---|---|
| 1 | OleNet | 78.12 |
| 2 | JESSI-A | 77.78 |
| 3 | m_y | 77.61 |

(a) Subtask A

| Rank | Model | F-Score |
|---|---|---|
| 1 | NTUA-ISLab | 85.80 |
| 2 | OleNet | 85.79 |
| 3 | NL-FIIT | 83.13 |
| * | JESSI-B | 79.59 |
| 11 | CNN→Aᴛᴛ+DᴏᴍAᴅᴠ | 64.86 |

(b) Subtask B

Table 4: F-Scores of JESSI and top three models for each subtask. Due to time constraints, we were not able to submit JESSI-B during the competition. For clarity, we also show our final official submission (CNN→Aᴛᴛ+DᴏᴍAᴅᴠ).



(a) Subtask A



(b) Subtask B

Figure 2: Accuracy over various input sentence length on the test set.

plots of the experiments on both subtasks. On both experiments, JESSI outperforms the other models

when the sentence length is short, suggesting that the increase in performance of JESSI can be attributed to its performance in short sentences. This is more evident in Subtask B, where the difference of accuracy between JESSI and the next best model is approximately 20%. We can also see a consistent increase in performance of JESSI over other models on Subtask B, which shows the robustness of JESSI when predicting out-of-domain samples.

## 6 Conclusion

We presented JESSI (Joint Encoders for Stable Suggestion Inference), our system for the SemEval 2019 Task 9: Suggestion Mining from Online Reviews and Forums. JESSI builds upon jointly combined encoders, borrowing pre-trained knowledge from a language model BERT and a translation model CoVe. We found that BERT alone performs bad and unstably when tested on out-of-domain samples. We mitigate the problem by appending an RNN-based sentence encoder above BERT, and jointly combining a CNN-based encoder. Results from the shared task show that JESSI performs competitively among participating models, obtaining second place on Subtask A with an F-Score of 77.78%. It also performs well on Subtask B, with an F-Score of 79.59%, even without using any additional external data.

## Acknowledgement

## References

Roberto Basili, Diego De Cao, Danilo Croce, Bonaventura Coppola, and Alessandro Moschitti. 2009. Cross-language frame semantics transfer in bilingual corpora. In *Computational Linguistics and Intelligent Text Processing, 10th International Conference, CICLing 2009, Mexico City, Mexico, March 1-7, 2009. Proceedings*, pages 332–345.

Caroline Brun and Caroline Hagège. 2013. Suggestion mining: Detecting suggestions for improvement in users' comments. *Research in Computing Science*, 70:199–209.

Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Fei-Fei Li. 2009. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2009), 20-25 June 2009, Miami, Florida, USA*, pages 248–255.

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.

Li Dong, Furu Wei, Yajuan Duan, Xiaohua Liu, Ming Zhou, and Ke Xu. 2013. The automated acquisition of suggestions from tweets. In *Proceedings of the Twenty-Seventh AAAI Conference on Artificial Intelligence, July 14-18, 2013, Bellevue, Washington, USA*.

Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor S. Lempitsky. 2016. Domain-adversarial training of neural networks. *Journal of Machine Learning Research*, 17:59:1–59:35.

Xavier Glorot, Antoine Bordes, and Yoshua Bengio. 2011. Domain adaptation for large-scale sentiment classification: A deep learning approach. In *Proceedings of the 28th International Conference on Machine Learning, ICML 2011, Bellevue, Washington, USA, June 28 - July 2, 2011*, pages 513–520.

Mi-Young Huh, Pulkit Agrawal, and Alexei A. Efros. 2016. What makes imagenet good for transfer learning? *CoRR*, abs/1608.08614.

Wonseok Hwang, Jinyeung Yim, Seunghyun Park, and Minjoon Seo. 2019. A comprehensive exploration on wikisql with table-aware word contextualization. *arXiv preprint arXiv:1902.01069*.

Yoon Kim. 2014. Convolutional neural networks for sentence classification. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL*, pages 1746–1751.

Lun-Wei Ku, Yu-Ting Liang, and Hsin-Hsi Chen. 2006. Opinion extraction, summarization and tracking in news and blog corpora. In *Computational Approaches to Analyzing Weblogs, Papers from the 2006 AAAI Spring Symposium, Technical Report SS-06-03, Stanford, California, USA, March 27-29, 2006*, pages 100–107.

Tao Lei, Yu Zhang, Sida I. Wang, Hui Dai, and Yoav Artzi. 2018. Simple recurrent units for highly parallelizable recurrence. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, pages 4470–4481.

Bing Liu. 2012. *Sentiment Analysis and Opinion Mining*. Synthesis Lectures on Human Language Technologies. Morgan & Claypool Publishers.

Christopher D. Manning, Mihai Surdeanu, John Bauer, Jenny Rose Finkel, Steven Bethard, and David McClosky. 2014. The stanford corenlp natural language processing toolkit. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics, ACL 2014, June 22-27, 2014, Baltimore, MD, USA, System Demonstrations*, pages 55–60.

Bryan McCann, James Bradbury, Caiming Xiong, and Richard Socher. 2017. Learned in translation: Contextualized word vectors. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, pages 6297–6308.

Arjun Mukherjee and Bing Liu. 2012. Aspect extraction through semi-supervised modeling. In *The 50th Annual Meeting of the Association for Computational Linguistics, Proceedings of the Conference, July 8-14, 2012, Jeju Island, Korea - Volume 1: Long Papers*, pages 339–348.

Sapna Negi, Kartik Asooja, Shubham Mehrotra, and Paul Buitelaar. 2016. A study of suggestions in opinionated texts and their automatic detection. In *Proceedings of the Fifth Joint Conference on Lexical and Computational Semantics, *SEM@ACL 2016, Berlin, Germany, 11-12 August 2016*.

Sapna Negi and Paul Buitelaar. 2015. Towards the extraction of customer-to-customer suggestions from reviews. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, EMNLP 2015, Lisbon, Portugal, September 17-21, 2015*, pages 2159–2167.

Sapna Negi and Paul Buitelaar. 2017. Inducing distant supervision in suggestion mining through part-of-speech embeddings. *CoRR*, abs/1709.07403.

Sapna Negi, Tobias Daudert, and Paul Buitelaar. 2019. Semeval-2019 task 9: Suggestion mining from online reviews and forums. In *Proceedings of the 13th International Workshop on Semantic Evaluation (SemEval-2019)*.

Sinno Jialin Pan, Xiaochuan Ni, Jian-Tao Sun, Qiang Yang, and Zheng Chen. 2010. Cross-domain sentiment classification via spectral feature alignment. In *Proceedings of the 19th International Conference on World Wide Web, WWW 2010, Raleigh, North Carolina, USA, April 26-30, 2010*, pages 751–760.

Bo Pang and Lillian Lee. 2007. Opinion mining and sentiment analysis. *Foundations and Trends in Information Retrieval*, 2(1-2):1–135.

Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL*, pages 1532–1543.

Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2018, New Orleans, Louisiana, USA, June 1-6, 2018, Volume 1 (Long Papers)*, pages 2227–2237.

Alec Radford, Karthik Narasimhan, Time Salimans, and Ilya Sutskever. 2018. Improving language understanding with unsupervised learning. Technical report, Technical report, OpenAI.

J. Ramanand, Krishna Bhavsar, and Niranjan Pedanekar. 2010. Wishful thinking - finding suggestions and 'buy' wishes from product reviews. In *Proceedings of the NAACL HLT 2010 Workshop on Computational Approaches to Analysis and Generation of Emotion in Text*, pages 54–61. Association for Computational Linguistics.

Sebastian Ruder and Barbara Plank. 2018. Strong baselines for neural semi-supervised learning under domain shift. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 1: Long Papers*, pages 1044–1054.

Kuniaki Saito, Yoshitaka Ushiku, and Tatsuya Harada. 2017. Asymmetric tri-training for unsupervised domain adaptation. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, pages 2988–2997.

Nitish Srivastava, Geoffrey E. Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: a simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(1):1929–1958.

Xiaojun Wan. 2009. Co-training for cross-lingual sentiment classification. In *ACL 2009, Proceedings of the 47th Annual Meeting of the Association for Computational Linguistics and the 4th International Joint Conference on Natural Language Processing of the AFNLP, 2-7 August 2009, Singapore*, pages 235–243.

Haoyu Wang, Ming Tan, Mo Yu, Shiyu Chang, Dakuo Wang, Kun Xu, Xiaoxiao Guo, and Saloni Potdar. 2019. Extracting multiple-relations in one-pass with pre-trained transformers. *arXiv preprint arXiv:1902.01030*.

Wei Yang, Yuqing Xie, Aileen Lin, Xingyu Li, Luchen Tan, Kun Xiong, Ming Li, and Jimmy Lin. 2019. End-to-end open-domain question answering with bertserini. *arXiv preprint arXiv:1902.01718*.

Jianfei Yu and Jing Jiang. 2016. Learning sentence embeddings with auxiliary tasks for cross-domain sentiment classification. In *Proceedings of the 2016*

*Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016*, pages 236–246.

Matthew D. Zeiler. 2012. ADADELTA: an adaptive learning rate method. *CoRR*, abs/1212.5701.