# ATTENTION-BASED LSTM FOR PSYCHOLOGICAL STRESS DETECTION FROM SPOKEN LANGUAGE USING DISTANT SUPERVISION

*Genta Indra Winata, Onno Pepijn Kampman, Pascale Fung*

Human Language Technology Center
Department of Electronic and Computer Engineering
Hong Kong University of Science and Technology, Clear Water Bay, Hong Kong
{giwinata, opkampman}@connect.ust.hk, pascale@ece.ust.hk

## ABSTRACT

We propose a Long Short-Term Memory (LSTM) with attention mechanism to classify psychological stress from self-conducted interview transcriptions. We apply distant supervision by automatically labeling tweets based on their hashtag content, which complements and expands the size of our corpus. This additional data is used to initialize the model parameters, and which it is fine-tuned using the interview data. This improves the model's robustness, especially by expanding the vocabulary size. The bidirectional LSTM model with attention is found to be the best model in terms of accuracy (74.1%) and f-score (74.3%). Furthermore, we show that distant supervision fine-tuning enhances the model's performance by 1.6% accuracy and 2.1% f-score. The attention mechanism helps the model to select informative words.

***Index Terms***— Psychological Stress Detection, LSTM, Natural Language Processing, Distant Supervision, Attention Mechanism

## 1. INTRODUCTION

Psychological stress has a serious effect on mental health and is often a precursor for more severe conditions. Although stress is a natural stimulant, persistent increased levels yield adverse effects, such as heart attacks [1], hypertension [2], and addiction [3]. Prolonged stress is also linked to mental health issues like anxiety [4] and depression [5]. Its prevalence has been increasing in the past decade [6] and affects the way people speak and their choice of spoken language. Emotional support is known to alleviate stress, yet less than 50% of the stressed population receives enough support from friends, family and professionals [6]. Linguistic studies have shown that language choice contains pointers to levels of stress and mental health [7]. The potential of text data from social media and Twitter for predicting major depression occurrence has also been demonstrated [8, 9].

Research on sentence-level stress detection has been mostly focused on written text collected from social media such as micro-blogs [10, 11]. The authors of these two works used a framework that combines linguistic, visual and social attributes in classifying stress categories. Lin et al. explored tweets to find stressors and stressful events, by building a stressors and stress subject dictionary. They collected thousands of written Weibo tweets and manually categorized them into 10 groups [11]. On word-level stress detection, a simple bidirectional RNN can achieve good results on Russian speech transcriptions [12].

In this work, we propose to build attention-based Long Short-Term Memory (LSTM) models fed with word embeddings for detecting psychological stress on sentence-level from interview transcriptions. It takes a long dependency across words in an utterance. Then, the attention mechanism weighs the importance of every word and chooses what to retrieve from the memory. It outputs the weighted combination of all words to the network for predicting the stress level. We apply distant supervision by adding unlabeled tweets from Twitter to our training set. This technique refers to extracting noisy signals from text as label [13]. In our case, we manually pick hashtags that indicate either a stressed or unstressed state of mind of the author, and use them to scrape stressed (positive labels) and unstressed (negative labels) tweets. We need to include more data during training because our interview corpus is relatively small and covers a limited number of topics, mostly related to academia. The major contribution of this paper is to show that unlabeled data collected from Twitter can improve the classification performance on our interview transcriptions corpus, and that applying an attention mechanism helps the model to effectively choose important words.

## 2. METHODOLOGY

Our objective is to determine whether someone is stressed or not, given an utterance as input. We explore several different models. For the LSTM and bidirectional LSTM (BLSTM)
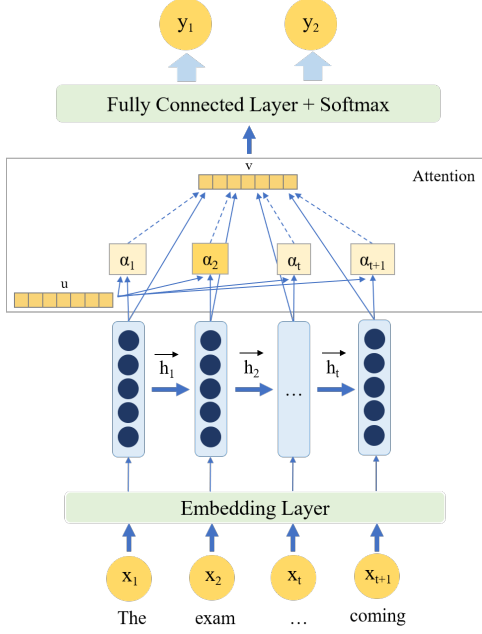
**Fig. 1**. Attention-based LSTM architecture



**Fig. 2**. Attention-based BLSTM architecture

models, we use a trainable embedding layer whose vectors should eventually form stressed and unstressed term clusters. LSTMs can capture temporal dynamics of words in a sentence.

### 2.1. Long Short-Term Memory (LSTM)

First, we build a unidirectional LSTM [14] taking word embedding as input. We denote $V$ as the number of unique words in our corpus and $k$ as the dimension of the word embedding vectors. Each word is a one-hot vector $\boldsymbol{x} \in \mathbb{R}^{|V|}$ and performs a multiplication with the embedding layer $A \in \mathbb{R}^{|V| \cdot k}$, where $k = 100$. The resulting vector is $\boldsymbol{b} \in \mathbb{R}^k$

$$\boldsymbol{b} = A^T \boldsymbol{x} \tag{1}$$

The model architecture is shown in Figure 1. The LSTM consists of one recurrent layer that propagates the embedding vector $b_t$ for the word at time $t$ (i.e. a column of $\boldsymbol{b}$) through the LSTM network to find hidden state $h_t$

$$\overrightarrow{h_t} = LSTM(b_t), t \in [1, T] \tag{2}$$

All hidden states are fed into a subsequent attention layer [15]. We added this layer because not all words contribute equally to the stress classifier. The word importance vector $u_t$ is calculated with Equation 3. The normalized word weight $\alpha_t$ is obtained through a softmax function (Equation 4). The aggregate of all the information in the sentence $\boldsymbol{v}$ is the weighted sum of each $h_t$ with $\alpha_t$ as corresponding weights.
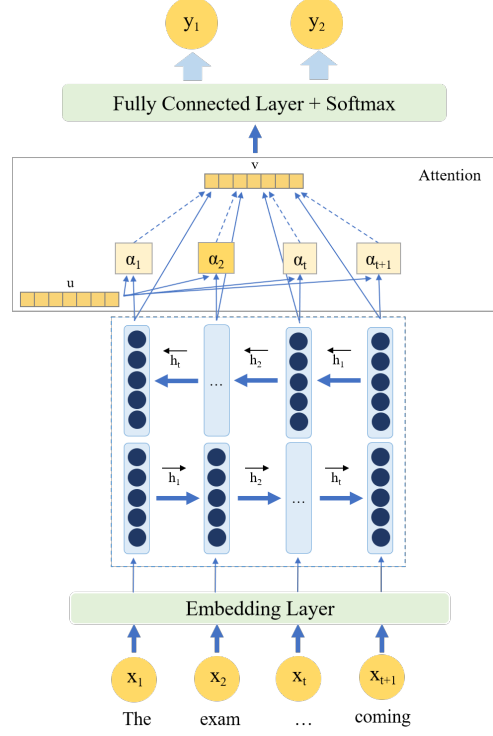
$$u_t = tanh(W \overrightarrow{h_t} + b) \tag{3}$$

$$\alpha_t = \frac{exp(u_t^T u)}{\sum_t exp(u_t^T u)} \tag{4}$$

$$\boldsymbol{v} = \sum_t \alpha_t \overrightarrow{h_t} \tag{5}$$

This vector $\mathbf{v}$ is then fed to a fully connected layer with softmax activation to perform the final classification. The prediction is a vector $\boldsymbol{y} \in \mathbb{R}^2$ with the probabilities of being unstressed and stressed. We choose the highest probability by using $argmax$ as the model's prediction.

### 2.2. Bidirectional Long Short-Term Memory (BLSTM)

We train a BLSTM model in identical fashion to process the word sequence in both forward ($\overrightarrow{h_t}$) and backward direction ($\overleftarrow{h_t}$). This recurrent neural network uses two LSTMs, one for each direction. The architecture is shown in Figure 2.

### 2.3. Support Vector Machine (SVM)

As a baseline, we build an SVM [16] with a Radial Basis Function (RBF) kernel. We extract `word2vec` [17] vector embeddings for each word in a given sentence. The embeddings have dimensionality of $k = 300$ and were pre-trained

**Table 1**. Twitter hashtags.

| # stressed |
| --- |
| amstressed, busylife, collegestress, distress, distressed, familystress, feelingbusy, feelingfrustrated, feelingoverwhelmed, feelingstress, feelstress, feelstressed, frustrated, frustrating, frustration, iamstressed, ifrustated, imstressed, overwhelm, overwhelmed, overwhelming, panic, sostress, sostressed, sostressful, stress, stressed, stressedlife, stressedout, stresses, stressful, stressfulllife, stressingout, stresslife, stressor, stressors, stresss, stressss, stresssss, verystressed, workstress |
| **# unstressed** |
| blessed, comfort, feelingrelax, feelingrelaxed, grateful, iamblessed, iamgrateful, iamrelaxed, imblessed, imgrateful, nostress, peaceful, relax, relaxed, relaxing |

**Table 2**. Interview Corpus statistics.

| utterances | tokens | speakers | vocab size |
| --- | --- | --- | --- |
| 2,272 | 36,538 | 38 | 3,127 |

**Table 3**. Twitter Corpus statistics.

| tweets | tokens | stressed | unstressed | vocab size |
| --- | --- | --- | --- | --- |
| 367,312 | 5,439,427 | 59,768 | 307,544 | 135,463 |

on Google News data (around 100 billion words with around 3 million unique words). Since the utterances have a variable number of words, we compute the average sentence vector for each embedding feature and train on the remaining feature space. We average the sum of all vectors to get a new word embedding vector, such that

$$b_j = \frac{\sum_{i=<N>} a_{i,j}}{N} \qquad (6)$$

where $a_{i,j}$ is the word embedding vector of word $i$ in sentence $j$, and $b_j$ is the sentence vector. Thus, for the SVM, the input is represented as an input matrix consisting of $N$ utterance vectors.

## 3. EXPERIMENTS

### 3.1. Corpora

For our experiments, we used two different corpora: an interview corpus, the Natural Stress Emotion corpus [18], and the Stress Twitter Corpus. The former corpus contains 25 students (13 females) answering the same set of 12 interview questions. The questions were designed to be progressively stress provoking. Additionally, we expanded the dataset by conducting 13 more interviews (3 with females) with identical setup. All answers were binary labeled for stress by three judges, from which we took the majority vote. It has four hours of recordings in total with 36,538 word tokens (see Figure 2). For the text-based models described here, we only used the English transcriptions. Because this corpus is small,

**Table 4**. Model performance.

| method | accu. | prec. | recall | f-score |
| --- | --- | --- | --- | --- |
| SVM | 68.7 | 72.0 | 61.2 | 66.2 |
| LSTM | 70.0 | 70.3 | 68.1 | 69.2 |
| LSTM w/ attention | **73.8** | **74.7** | **71.9** | **73.2** |
| BLSTM | 72.2 | 74.5 | 67.5 | 70.8 |
| BLSTM w/ attention | 72.5 | 73.1 | 71.2 | 72.2 |

**Table 5**. Fine-tuning performance.

| method | accu. | prec. | recall | f-score |
| --- | --- | --- | --- | --- |
| LSTM | 73.4 | 73.6 | 73.1 | 73.4 |
| LSTM w/ attention | 73.8 | 74.4 | 72.5 | 73.4 |
| BLSTM | 73.8 | **74.7** | 71.9 | 73.2 |
| BLSTM w/ attention | **74.1** | 73.6 | **75.0** | **74.3** |

we collected more data from Twitter and selected tweets with a set of filtering heuristics based on [19]. We only kept tweets with the hashtag at the end, and having less than four hashtags in total. We filtered out tweets containing URLs and images, and applied distant supervision [13] to label the unlabeled tweets. That is, we manually chose hashtags that indicate either stressed or unstressed state of mind of the author (See Figure 1), and used these to automatically label our scraped tweets. Not all text is created equally and it's important here to be aware of the differences between spoken language and written language on social media.

### 3.2. Setup

For the LSTM and BLSTM experiments, the recurrent layer consists of 64 units. In order to regularize the model, a dropout layer [20] with probability of 0.2 is inserted between the recurrent and attention layers. We use batch gradient descent using Adam [21] as optimizer, with batches of 128 samples. We also run both the LSTM and BLSTM without attention mechanism for comparison.

We take 160 random samples from each class from the interview corpus as our test set, as we want to evaluate on spoken language. The remainder is used as training set. All sentences are padded to 35 words. In order to balance the distribution of our training set, we oversample the minority class (stressed) within the training set. We validate the model to find the best setting. For two iterations, the model is trained only with twitter data and afterwards, the model subsequently fine-tuned with interview data. Twitter data is inherently different from spoken transcripts, and both are noisy (absence of correct grammar) in their own way. Since the Twitter corpus is imbalanced, we random sample and take 49,000 tweets from each class every iteration.
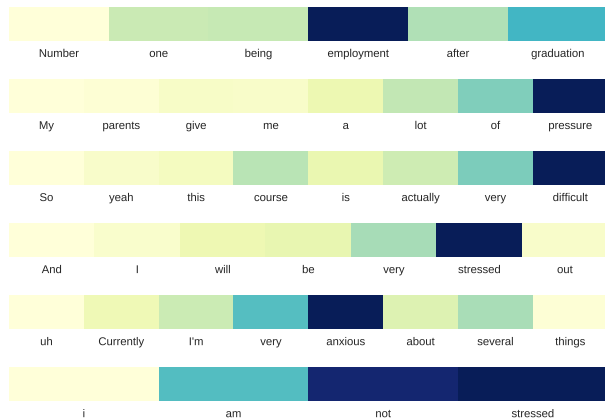
**Fig. 3**. Heatmap of attention layer weights for stressed utterances.



**Fig. 4**. Heatmap of attention layer weights for unstressed utterances.

## 3.3. Results

Relevant evaluation results on the test set are shown in Tables 4 and 5. The performance of the BLSTM with attention outperforms the other classifiers in terms of accuracy and f-score. The fine-tuning process helps the model to classify sentences related to stress better, but significantly increases the recall. Models with attention mechanism are slightly better.

## 3.4. Discussion

To visualize the attention mechanism, we extract the attention weights from the best trained model and evaluate several stressed and unstressed utterances (see Figures 3 and 4). The figures show the contribution of each word in the classification task. Darker colors represent stronger word contributions to the classification task. Interestingly, the added attention layer captures key terms related to stress. For instance, in the first stressed example, words such as "employment" and "graduation" weigh heavier than others. These words have stronger relation with stress. Furthermore, we can see from other stressed samples that words such as "employment", "pressure", "difficult", "stressed", and "anxious" have similar weights. Conversely, words such as "my", "Number", "I", and "And" are least considered for the classification because they are not related to stress.

Distant supervision seems to slightly improve the performance, especially the recall (i.e. false negatives are turned to true positive). This is likely due to the interview corpus covering a limited domain (mainly academic issues), and all interviewees answering the same questions. The domain of the Twitter corpus is more general because it includes other domains and is approximately 40 times larger than the interview corpus, thus adding it makes the model more robust. However, the model did not learn more complicated grammatical
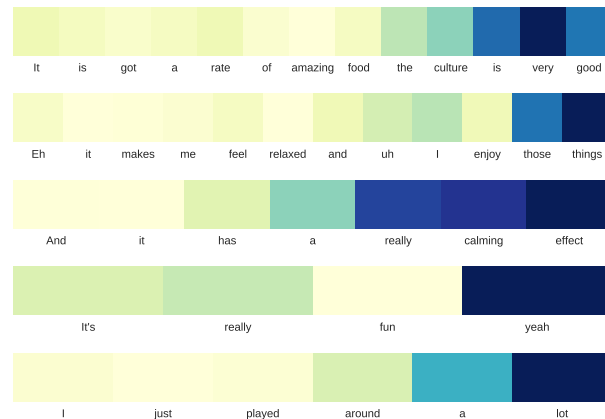
structures. For example, from Figure 3 we can see that "I am not stressed" is classified as stressed. The observation that the model does not learn the semantic meaning of negation, could be caused by a lack of data. Also, we believe tweets are not a great source for models to learn proper grammatical structures.

Our models learn statistical characteristics of language choice under stressed psyche. That is, they learn which word combinations and sequences are expressed more often when someone is stressed. It is showed more explicitly in Figure 3. Obviously, a person can talk about stressful topics, yet still remain calm, and vice versa. Although there is an obvious signal, not all stress information is encoded in language choice. This inherently limits our model. The interview corpus only contains transcribed spoken language. A more complete stress detection framework would also include context and prosodic features of spoken utterances.

## 4. CONCLUSION

We have presented methods for classifying interviewee stress level from interview transcriptions. The best performance was found for our bidirectional LSTM model, which outperformed the other models in terms of accuracy and f-score. The two-phase training method with the out-of-domain stress tweets dataset improves the learning performance. Future work includes multi-modal learning using linguistic and acoustic features. We are also interested in gathering more grammatically correct sentences for transfer learning purposes, so the model may learn how to deal with negation (among others). Furthermore, we will incorporate the model described here into our virtual therapist platform [22], where it is fed with Automatic Speech Recognition output. This makes the system aware of user stress, to which it responds with appropriate stress management advice and exercises.

# 5. REFERENCES

[1] John S Yudkin, Meena Kumari, Steve E Humphries, and Vidya Mohamed-Ali, "Inflammation, obesity, stress and coronary heart disease: is interleukin-6 the link?," *Atherosclerosis*, vol. 148, no. 2, pp. 209–214, 2000.

[2] Karen A Matthews, Charles R Katholi, Heather Mc-Creath, Mary A Whooley, David R Williams, Sha Zhu, and Jerry H Markovitz, "Blood pressure reactivity to psychological stress predicts hypertension in the cardia study," *Circulation*, vol. 110, no. 1, pp. 74–78, 2004.

[3] Natalie Slopen, Emily Z Kontos, Carol D Ryff, John Z Ayanian, Michelle A Albert, and David R Williams, "Psychosocial stress and cigarette smoking persistence, cessation, and relapse over 9–10 years: a prospective study of middle-aged adults in the united states," *Cancer Causes & Control*, vol. 24, no. 10, pp. 1849–1863, 2013.

[4] Carlo Faravelli and Stefano Pallanti, "Recent life events and panic disorder," *Am J Psychiatry*, vol. 146, no. 5, pp. 622–6, 1989.

[5] Naomi Breslau, Lonni Schultz, and Edward Peterson, "Sex differences in depression: a role for preexisting anxiety," *Psychiatry research*, vol. 58, no. 1, pp. 1–12, 1995.

[6] APA, "Stress in america, coping with change 2017," *American Psychology Association*, 2017.

[7] Bridianne O'Dea, Mark E Larsen, Philip J Batterham, Alison L Calear, and Helen Christensen, "A linguistic analysis of suicide-related twitter posts," *Crisis*, 2017.

[8] Munmun De Choudhury, Michael Gamon, Scott Counts, and Eric Horvitz, "Predicting depression via social media.," in *ICWSM*, 2013, p. 2.

[9] Patricia A Cavazos-Rehg, Melissa J Krauss, Shaina Sowles, Sarah Connolly, Carlos Rosas, Meghana Bharadwaj, and Laura J Bierut, "A content analysis of depression-related tweets," *Computers in human behavior*, vol. 54, pp. 351–357, 2016.

[10] Huijie Lin, Jia Jia, Quan Guo, Yuanyuan Xue, Jie Huang, Lianhong Cai, and Ling Feng, "Psychological stress detection from cross-media microblog data using deep sparse neural network," in *Multimedia and Expo (ICME), 2014 IEEE International Conference on*. IEEE, 2014, pp. 1–6.

[11] Huijie Lin, Jia Jia, Liqiang Nie, Guangyao Shen, and Tat-Seng Chua, "What does social media say about your stress?.," in *IJCAI*, 2016, pp. 3775–3781.

[12] Maria Ponomareva, Kirill Milintsevich, Ekaterina Chernyak, and Anatoly Starostin, "Automated word stress detection in russian," in *Proceedings of the First Workshop on Subword and Character Level Models in NLP*, 2017, pp. 31–35.

[13] Micol Marchetti-Bowick and Nathanael Chambers, "Learning for microblogs with distant supervision: Political forecasting with twitter," in *Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics*. Association for Computational Linguistics, 2012, pp. 603–612.

[14] Sepp Hochreiter and Jürgen Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.

[15] Zichao Yang, Diyi Yang, Chris Dyer, Xiaodong He, Alexander J Smola, and Eduard H Hovy, "Hierarchical attention networks for document classification.," in *HLT-NAACL*, 2016, pp. 1480–1489.

[16] Corinna Cortes and Vladimir Vapnik, "Support-vector networks," *Machine learning*, vol. 20, no. 3, pp. 273–297, 1995.

[17] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean, "Distributed representations of words and phrases and their compositionality," in *Advances in neural information processing systems*, 2013, pp. 3111–3119.

[18] Xin Zuo, L Lin, and Pascale Fung, "A multilingual database of natural stress emotion," in *Proceedings of the Eight International Conference on Language Resources and Evaluation (LREC12)*, 2012, pp. 1174–1178.

[19] Wenbo Wang, Lu Chen, Krishnaprasad Thirunarayan, and Amit P Sheth, "Harnessing twitter" big data" for automatic emotion identification," in *Privacy, Security, Risk and Trust (PASSAT), 2012 International Conference on and 2012 International Confernece on Social Computing (SocialCom)*. IEEE, 2012, pp. 587–592.

[20] Nitish Srivastava, Geoffrey E Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov, "Dropout: a simple way to prevent neural networks from overfitting.," *Journal of machine learning research*, vol. 15, no. 1, pp. 1929–1958, 2014.

[21] Diederik P Kingma and Jimmy Lei Ba, "Adam: Amethod for stochastic optimization," .

[22] Genta Indra Winata, Onno Kampman, Yang Yang, Anik Dey, and Pascale Fung, "Nora the empathetic psychologist," *Proc. Interspeech 2017*, pp. 3437–3438, 2017.