# Sentylic at IEST 2018: Gated Recurrent Neural Network and Capsule Network Based Approach for Implicit Emotion Detection

**Prabod Rathnayaka, Supun Abeysinghe, Chamod Samarajeewa**
**Isura Manchanayake, Malaka Walpola**
Department of Computer Science and Engineering
University of Moratuwa, Sri Lanka
{prabod.14,supun.14,chamod.14,isura.14,malaka}@cse.mrt.ac.lk

## Abstract

In this paper, we present the system we have used for the Implicit WASSA 2018 Implicit Emotion Shared Task. The task is to predict the emotion of a tweet of which the explicit mentions of emotion terms have been removed. The idea is to come up with a model which has the ability to implicitly identify the emotion expressed given the context words. We have used a Gated Recurrent Neural Network (GRU) and a Capsule Network based model for the task. Pre-trained word embeddings have been utilized to incorporate contextual knowledge about words into the model. GRU layer learns latent representations using the input word embeddings. Subsequent Capsule Network layer learns high-level features from that hidden representation. The proposed model managed to achieve a macro-F1 score of 0.692.

## 1 Introduction

Emotion is a complex aspect of the human behavior which makes the humanity distinguishable from other biological behaviors of creatures. Emotions are typically originated as a response to a situation. Since the emergence of social media, people often express opinions as responses to daily encounters by posting on these platforms. These microblogs contain emotions related to the topics the author have discussed. Thus emotion detection is useful to understand more specific sentiments held by the author towards the discussed topics. Hence this is a challenge with a significant business value.

Emotion analysis can be considered as an extension of sentiment analysis. Even though there has been a notable amount of research in sentiment analysis in the literature, research on emotion analysis has not gained much attention. The related work suggests this task can be handled using emojis or hashtags present in the text i.e. distance supervision techniques (Felbo et al., 2017). However, these features can be unreliable due to noise, thus affect the accuracy of the results.

Although explicit words related to emotions (happy, sad, etc.) in a document directly affect the emotion detection task, other linguistic features play a major role as well. Implicit Emotion Recognition Shared Task introduced in Klinger et al. (2018) aims at developing models which can classify a text into one of the emotions; *Anger, Fear, Sadness, Joy, Surprise, Disgust* without having access to an explicit mention of an emotion word. Participants were given a tweet from which any of the above emotion terms or one of their synonyms is removed. The task is to predict the emotion that the excluded word expresses.

E.g.:

*It's [#TARGETWORD#] when you feel like you are invisible to others.*

The [#TARGETWORD#] in the given example corresponds to sadness ("sad").

In this paper, we propose an approach based on Gated Recurrent Units (GRU) (Cho et al., 2014) followed by Capsule Networks (Sabour et al., 2017) to tackle the challenge. This model managed to achieve a macro-F1 score of **0.692** and ranked **5th** in WASSA 2018 implicit emotion detection task.

## 2 Methodology

We have used a sentence classification model which is based on bidirectional GRUs and Capsule networks. First, the raw tweets are preprocessed, then mapped into a continuous vector space using an embedding layer. Afterward, we used a Bidirectional Gated Recurrent Unit (Bi-GRU) (Cho et al., 2014) layer to encode sentences into a fixed length representation. The fixed length represen-

tation is then fed into a Capsule Network (Sabour et al., 2017) where it will learn the features and emotional context of the sentences. Finally, the Capsule network is followed by a fully connected dense layer with softmax activation for the classification.

## 2.1 Preprocessing

Microblogs typically contain informal language usages such as short terms, emojis, misspellings, and hashtags. Hence, preprocessing steps should be employed in order to clean these informal and noisy text data. Moreover, efficient preprocessing plays a vital role in achieving a good performance. Ekphrasis tool (Baziotis et al., 2017) is used for initial preprocessing of the tweets. Tweet tokenizing, word normalization, spell correcting and word segmentation for hashtags are done as preprocessing steps.

### 2.1.1 Tweet Tokenizing

Tokenizing is the first and the most important step of preprocessing. Ability to correctly tokenize a tweet directly impacts the quality of a system. Since there is a large variety of vocabulary and expressions present in short texts such as Twitter, it is a challenging task to correctly tokenize a given tweet. Twitter markup, emoticons, emojis, dates, times, currencies, acronyms, censored words (e.g. s**t), words with emphasis (e.g. *very*) are recognized during tokenizing and treated as a separate token.

### 2.1.2 Word Normalization

Upon tokenizing, set of transformations including converting to lowercase and transforming URLs, usernames, emails, phone numbers, dates, times, hashtags to a predefined set of tags (e.g @user1 → <user>) are applied. This method helps to reduce the vocabulary size and generalize the tweet.

### 2.1.3 Spell Correcting and Word Segmentation

As the last step in preprocessing, we apply spell correcting and word segmentation to hashtags. (e.g. #makeitrain → make it rain)

## 2.2 Model

An overview of the model is shown in figure 1 and each segment of the model is described in the following sub sections.
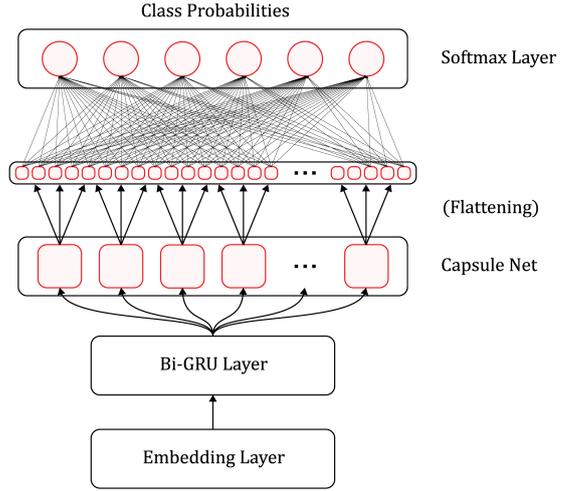


Figure 1: Overall model architecture

### 2.2.1 Word Embedding Layer

Word embedding layer is the first layer of the model. Each token will be mapped into a continuous vector space using a set of pretrained word embeddings. We used 300 dimensional, pretrained, Word2Vec embeddings introduced in Mikolov et al. (2013). Given an input tweet, $S = [s_1, s_2, ., s_i, .., s_n]$ where $s_i$ is the token at position $i$, the embedding matrix $W_e$, the output of the embedding layer $X$ is,

$$X = SW_e \qquad (1)$$

## 2.3 Bidirectional GRU Layer

The word embedding layer is followed by a bidirectional GRU (Cho et al., 2014) layer. There is a forward GRU ($\overrightarrow{h_t}$) and a backward GRU ($\overleftarrow{h_t}$)) and the latent representation output by the two GRUs is concatenated to get the final output ($\overrightarrow{h_t}, \overleftarrow{h_t}$) of the layer. Following set of equations follows the standard notation used in Cho et al. (2014).

$$r_t = \sigma(W_{ir}x_t + b_{ir} + W_{hr}h_{(t-1)} + b_{hr}) \qquad (2)$$

$$z_t = \sigma(W_{iz}x_t + b_{iz} + W_{hz}h_{(t-1)} + b_{hz}) \qquad (3)$$

$$n_t = \tanh(W_{in}x_t + b_{in} + r_t(W_{hn}h_{(t-1)} + b_{hn})) \qquad (4)$$

$$h_t = (1 - z_t)n_t + z_t h_{(t-1)} \qquad (5)$$

## 2.4 Capsule Layer

Features encoded by the bidirectional GRU layer is then passed to a Capsule Network (Sabour et al.,

2017). Capsule Network consists of a set of capsules where each capsule corresponds to a high level feature. Each capsule outputs a vector, of which the magnitude represents the probability of the corresponding feature existence. Following set of equations follows the standard notation used in Sabour et al. (2017).

Prediction vector $\hat{u}_{j|i}$ is calculated by multiplying the output $h_i$ from the GRU layer with a weight matrix.

$$\hat{u}_{j|i} = W_{ij}h_i \qquad (6)$$

Total input to a capsule $s_j$ is a weighted sum over all the prediction vectors $\hat{u}_{j|i}$.

$$s_j = \sum_i c_{ij}\hat{u}_{j|i} \qquad (7)$$

$c_{ij}$ represents the coupling coefficients found through the iterative dynamic routing.

A non-linear *"Squash"* function is used to scale the vectors such that the magnitude is mapped to a value between 0 and 1.

$$v_j = \frac{\|s_j\|^2}{1 + \|s_j\|^2} \frac{s_j}{\|s_j\|} \qquad (8)$$

Dynamic Routing process introduced by Sabour et al. (2017) is used as the routing mechanism between capsules.

## 2.5 Classification Layer

The flattened output from the capsule layer (Say $C$) is fed to a dense layer which has a softmax activation. It outputs a vector of size 6 (number of classes). The values in the vector components are probabilities for the presence of each of the six emotions. The emotion with the highest probability is selected as the output.

$$Y = W_{dense}C \qquad (9)$$

For all $y_i \in Y$, $f_i$ is calculated as follows.

$$f_i = \frac{e^{-y_i}}{\sum_{y_j \in Y} e^{-y_j}} \qquad (10)$$

Then the class with highest $f_i$ is taken as the output.

$$output = \arg\max_i f_i \qquad (11)$$

## 2.6 Regularization

Gaussian noise is added to both the embedding layer and the softmax classification layer for the purpose of making the model more robust to overfitting. Further, dropout is applied to the Capsule network output and a spatial dropout is applied to the embedding Layer to reduce overfitting.

## 3 Experiments and Results

### 3.1 Experimental setup

#### 3.1.1 Training

We used Adam optimizer (Kingma and Ba, 2014) for optimizing our network with a batch size of 512. Gradient Clipping (Pascanu et al., 2013) was employed to address the exploding gradient problem where all the gradients were clipped at 1. Keras (Chollet et al., 2015) was used to develop the model and experiments were done using both Tensorflow (Abadi et al., 2015) and Theano (Theano Development Team, 2016) backends. Google Colaboratory[1] was used as the runtime environment for training the model.

#### 3.1.2 Hyper-Parameters

We have employed Word2Vec (Mikolov et al., 2013) embeddings of 300 dimensions for the embedding layer. The GRU layer consists of 128 cells for both directions. We have used 16 capsules each with an output size of 32 and 5 routing iterations. Spatial dropout of 0.3 is applied to the embeddings and dropout of 0.25 is applied to the Capsule network. Gaussian noise of 0.1 is added to both the embedding layer and the Capsule network.

### 3.2 Results

We ranked 5th among 30 contestants in the competition. We achieved a macro-F1 score of 0.692 which is 0.155 improvement compared to the baseline model (Maximum Entropy Model using bag of words (BoW) and bigrams). The top-ranked model has a 0.031 improvement compared to our model. Table 1 shows the macro-F1 scores of the top 10 competitors and the baseline model.

## 4 Analysis

### 4.1 Investigated Approaches

Recurrent Neural Networks (RNN) (Socher et al., 2013) based model achieves state-of-the-art in

---

[1] https://colab.research.google.com/

| Team | Macro-F1 |
|---|---|
| **Amobee** | **0.714** |
| IIIDYT | 0.710 |
| NTUA-SLP | 0.703 |
| UBC-NLP | 0.693 |
| **Sentylic** | **0.692** |
| HUMIR | 0.686 |
| nlp | 0.685 |
| DataSEARCH | 0.680 |
| YNU1510 | 0.676 |
| EmotiKLUE | 0.671 |
| **Baseline** | **0.599** |

Table 1: Competition results of top 10 competitors and the maximum entropy baseline classifier.

| Model | Macro-F1 |
|---|---|
| GRU + Hierarchical Attention | 0.671 |
| GRU + CNN | 0.657 |
| **GRU + Capsnet** | **0.692** |

Table 2: Performance analysis of the best models in each investigated approaches.

| Model | Macro-F1 |
|---|---|
| **GRU (1 layer) + Capsnet** | **0.692** |
| LSTM (1 layer) + Capsnet | 0.687 |
| GRU (2 layers) + Capsnet | 0.678 |

Table 3: Performance analysis of different variants of the proposed system

sentence classification tasks. RNNs have the capability to capture sequential features present in sentences. Further, when they are incorporated with attention mechanisms the accuracy of the models increases notably (Yang et al., 2016; Tang et al., 2015). Hence, we have first implemented a model which uses a bidirectional GRU (Cho et al., 2014) layer to learn latent representations followed by a hierarchical attention mechanism. Attention mechanisms have the ability to capture important keywords in sentences and give a higher weight to those words. This is one of the prominent approaches that typically results in a good performance in regular text classification tasks. Table 2 shows that this approach yielded a reasonable accuracy, yet it was not the best performing approach.

Another approach is to use a Convolution Neural Network (CNN) (Kim, 2014) layer on top of RNNs instead of attention mechanisms. Intuition is that the CNN layers will act as a different attention mechanism and captures high-level features from the features learned by the below layers. Hence, the second approach we investigated was using CNNs instead of the attention mechanism. As the table 2 shows, this approach resulted in a slight drop in performance compared to the previous approach.

Our next approach was to use Capsule networks (Sabour et al., 2017) instead of Convolution Neural Networks (CNN). Capsule networks have shown promising results in the field of computer vision. Sabour et al. (2017) argues that it is essential to preserve the hierarchical translational and rotational features of the identified high-level fea-

tures in order to perform image classification and object detection in the field of computer vision. However, traditional CNNs with max-pooling layers tend to lose this spatial information related to identified features. Sabour et al. (2017) introduces capsule networks to tackle these issues identified in traditional CNNs. Nonetheless, the usability of Capsule networks has not researched much in the Natural Language Processing (NLP) community. Along the same lines, we can intuitively argue that CNN based models with pooling layers will cause loss of information in text related classification tasks as well. Hence, we have investigated the usability of capsule networks for improving the performance of text classification models. The use of Capsule networks instead of CNNs has improved the performance of the model slightly and assisted in gaining the best performing model.

### 4.2 Model Architecture Variants

We have tried several variants of the proposed model. Table 3 shows the performance of each of those variants. We have tried approaches using Long Short Term Memory networks (LSTM) (Hochreiter and Schmidhuber, 1997) which is one of the other prominent types of RNNs. However, the results showed a minor drop. Another variant is to use two layers of GRUs instead of using a single layer. Even this approach made the performance of the model slightly lesser. A potential reason for this could be model over-fitting. Using a single GRU layer followed by the Capsnet gave the best performance.

| Emotion | Precision | Recall | F1-Score |
|---------|-----------|--------|----------|
| Anger | 0.631 | 0.614 | 0.622 |
| Disgust | 0.689 | 0.687 | 0.688 |
| Fear | 0.728 | 0.731 | 0.730 |
| Joy | 0.803 | 0.774 | 0.788 |
| Sad | 0.682 | 0.655 | 0.668 |
| Surprise | 0.625 | 0.689 | 0.656 |
| **Micro Avg.** | 0.694 | 0.694 | 0.694 |
| **Macro Avg.** | 0.693 | 0.692 | 0.692 |

Table 4: Precision, recall and F1-score of each class in test set using our proposed model.

### 4.3 Analysis on Predictions

Table 4 shows the performance of the proposed model for each class. As evident from the results, anger shows a significantly lower F1-score. Other emotions show similar results whereas joy stands out with a notably higher F1-score. Anger has been misclassified as sad in several examples.

e.g.- *Girls will get [#TARGETWORD#] that her man cheated with an ugly girl more than the fact he actually cheated.*

In the above example, it is unclear whether the emotion is anger or sadness. Such ambiguity of anger has affected the reduction of F1-score values. There were few other similar cases where it is challenging even for humans to clearly discriminate emotions due to nuance nature of emotions expressed.

## 5 Conclusion

WASSA 2018 Implicit Emotion Shared Task (Klinger et al., 2018) introduces a task to predict the emotion of a tweet of which the explicit mentions of emotion terms have been removed. We have experimented with several deep learning based approaches to tackle this task. We have used pre-trained Word2Vec embeddings. All the approaches we tried utilize an initial GRU layer which learns latent representations from the input word embeddings. Different alternative methods have been investigated for the subsequent layer. These methods include attention layer, CNN layer, and Capsnet layer. Model with the Capsnet layer achieved the best results among the experimented alternatives. Potential future work includes investigating the possibility of using Capsule networks for other tasks in Natural Language Processing, especially where CNNs are involved. Another line of future work could be to follow the ap-

proach mentioned in Felbo et al. (2017) and apply transfer learning on the model trained using this semi-automatically annotated dataset to test on human annotated datasets such as Mohammad et al. (2018).

## References

Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. 2015. TensorFlow: Large-scale machine learning on heterogeneous systems. Software available from tensorflow.org.

Christos Baziotis, Nikos Pelekis, and Christos Doulkeridis. 2017. Datastories at semeval-2017 task 4: Deep lstm with attention for message-level and topic-based sentiment analysis. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 747–754, Vancouver, Canada. Association for Computational Linguistics.

Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*.

François Chollet et al. 2015. Keras. https://keras.io.

Bjarke Felbo, Alan Mislove, Anders Søgaard, Iyad Rahwan, and Sune Lehmann. 2017. Using millions of emoji occurrences to learn any-domain representations for detecting sentiment, emotion and sarcasm. *arXiv preprint arXiv:1708.00524*.

Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.

Yoon Kim. 2014. Convolutional neural networks for sentence classification. *arXiv preprint arXiv:1408.5882*.

Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

Roman Klinger, Orphée de Clercq, Saif M. Mohammad, and Alexandra Balahur. 2018. Iest: Wassa-2018 implicit emotions shared task. In *Proceedings*

*of the 9th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis*, Brussels, Belgium. Association for Computational Linguistics.

Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.

Saif M Mohammad, Felipe Bravo-Marquez, Mohammad Salameh, and Svetlana Kiritchenko. 2018. Semeval-2018 task 1: Affect in tweets. In *Proceedings of International Workshop on Semantic Evaluation (SemEval-2018), New Orleans, LA, USA*.

Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. 2013. On the difficulty of training recurrent neural networks. In *International Conference on Machine Learning*, pages 1310–1318.

Sara Sabour, Nicholas Frosst, and Geoffrey E Hinton. 2017. Dynamic routing between capsules. In *Advances in Neural Information Processing Systems*, pages 3856–3866.

Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pages 1631–1642.

Duyu Tang, Bing Qin, and Ting Liu. 2015. Document modeling with gated recurrent neural network for sentiment classification. In *Proceedings of the 2015 conference on empirical methods in natural language processing*, pages 1422–1432.

Theano Development Team. 2016. Theano: A Python framework for fast computation of mathematical expressions. *arXiv e-prints*, abs/1605.02688.

Zichao Yang, Diyi Yang, Chris Dyer, Xiaodong He, Alex Smola, and Eduard Hovy. 2016. Hierarchical attention networks for document classification. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1480–1489.