

Dual-State Capsule Networks for Text Classification

Piyumal Demotte, Surangika Ranathunga

Department of Computer Science and Engineering, University of Moratuwa
Katubedda 10400, Sri Lanka

{piyumalanthony.16, surangika}@cse.mrt.ac.lk

Abstract

Text classification systems based on contextual embeddings are not viable options for many of the low resource languages. On the other hand, recently introduced capsule networks have shown performance in par with these text classification models. Thus, they could be considered as a viable alternative for text classification for languages that do not have pre-trained contextual embedding models. However, current capsule networks depend upon spatial patterns without considering the sequential features of the text. They are also sub-optimal in capturing the context-level information in longer sequences. This paper presents a novel Dual-State Capsule (DS-Caps) network-based technique for text classification, which is optimized to mitigate these issues. Two varieties of states, namely sentence-level and word-level, are integrated with capsule layers to capture deeper context-level information for language modeling. The dynamic routing process among capsules was also optimized using the context-level information obtained through sentence-level states. The DS-Caps networks outperform the existing capsule network architectures for multiple datasets, particularly for tasks with longer sequences of text. We also demonstrate the superiority of DS-Caps in text classification for a low resource language.

1 Introduction

Recent research has shown state-of-the-art results using BERT-like contextual embedding models (Devlin et al., 2018). However, these performances highly depend on the comprehensiveness of the pre-trained contextual embedding models. In other words, these models require a very large corpus to produce optimal results in down-stream tasks such as text classification. However, for many low resource languages, such large corpora are not available. Even if available, building contextual

embedding models for each and every language has practical concerns with respect to resource requirements. As a solution, multilingual models¹ were introduced. However, still there are many languages that are not included in these publicly available multilingual models. It has also been shown that these multilingual models give sub-optimal results compared to those trained on monolingual data (Dumitrescu et al., 2020). Thus, text classification for these languages still has to rely on techniques that do not involve BERT-like models.

Another line of research has used language-specific linguistic features such as Part of Speech (POS) as auxiliary input to neural models (Qian et al., 2016). Although these have also shown very impressive results, low resource languages that do not have such resources cannot take benefit of this line of research as well.

Among the Deep Learning models for text classification that do not rely on contextual embeddings or auxiliary linguistic features, capsule networks are in the fore-front (Zhao et al., 2018; Kim et al., 2020). For example, for the MR(2005) (Pang and Lee, 2005) results produced by Capsule Network models closely trail behind linguistically enhanced and contextual embeddings-based Deep Learning techniques (Zhao et al., 2018).

This paper further improves the state-of-the-art capsule networks and presents a novel capsule network architecture namely, Dual-State Capsule Networks (DS-Caps). In particular, this is an enhancement to the existing capsule networks model with dynamic routing, which has been employed for text classification (Zhao et al., 2018). DS-Caps is mainly inspired by the Sentence State LSTM (S-LSTM) networks (Zhang et al., 2018), which eliminates the sequential dependencies of LSTMs to enable to capturing local n-grams as well as se-

¹<https://github.com/google-research/bert/blob/master/multilingual.md>

quential features of the text simultaneously. The proposed model could be represented as a sophisticated language model that elevates the usage of n-grams, sequential patterns, and capsules vector representation for language modeling within a single architecture. DS-Caps further mitigate the limitations with the vanilla capsule network of capturing contextual-level data of the text while sequential processing.

We empirically show that our approach produces state-of-the-art performances against the existing capsule networks, for many publicly available datasets for text classification tasks. Furthermore, DS-Caps is evaluated against a text classification task for a low resource language, for which there are no pre-trained monolingual or multilingual embedding models. Even for this low resource task, DS-Caps was able to outperform other capsule network-based methods, as well as the recurrent models.

2 Related Work

Rather recent experiments suggest approaches based on combinations of sequential models like RNNs, LSTMs, and CNNs to capture both contextual information and n-gram features of text (Wang et al., 2016). Yet, these approaches still suffer from inherent limitations of sequential text processing with RNNs based approaches and convolutions over sequences of text. Zhang et al. (2018) argue that, in order to encode deep language representations within a neural language model, both n-gram features, as well as sequential information of text, should be utilized within a single language model. Their language modeling procedure includes S-LSTM, a graph RNN. S-LSTM encodes each word of a sentence as a separate state, as well as the sentence as a global state in each recurrent step. For this, they utilize message passing over graphs (Scarselli et al., 2008) in order to capture the deep neural representation of languages.

Furthermore, capsule networks that were initially employed in the image processing tasks produced state-of-the-art results with the dynamic routing procedure proposed by Sabour et al. (2017). The intention behind the capsule strategy was to represent the features of objects within the data as vector representation in order to identify the exact order or pose of the information. The dynamic routing procedure between capsules mitigates the issues of information loss of CNNs due

to max-pooling and elevates the advancement of the part-to-whole relationship between capsules for deeper capsule representation for image classification tasks.

Inspired by the impressive performance of the capsule network architectures for image classification, Wang et al. (2018) applied the same for sentiment classifications with the combination of RNNs, which produced state-of-the-art results at that time. Zhao et al. (2018) conducted an empirical experiment of capsule networks with dynamic routing to validate utilization of capsule networks for text classification. The implementation of capsule-A and capsule-B, considering different n-gram variations produced the optimal performances for text classification tasks. With even deeper analysis, Kim et al. (2020) produced an approach based on static routing between capsules for text classification tasks. This procedure alleviates the limitations provided by the variations of text with background noise for capsules with dynamic routing. However, capsules that are solely built upon convolutions do not illustrate the ability to capture the context of the text when the length of sequences increases.

3 Model Architecture

DS-Caps further improves the language representation capability of capsule networks (Zhao et al., 2018) with the use of sentence-level states and word-level states, an idea borrowed from the S-LSTM architecture (Zhang et al., 2018). As illustrated in Figure 1, for the purpose of capturing n-grams features with sequential dependencies of the text, sentence-level states, and word-level states were employed on top of feature embeddings.

In particular, vanilla capsule models leverage pre-trained word-embeddings as the input features. The idea behind the word-level states is to capture both n-gram and sequential features of the text simultaneously and incorporate with vector representation of capsules to enhance language modeling. Therefore instead of using word-embeddings as feature inputs for vanilla capsule, feature maps based on concatenated word-level states were utilized.

Furthermore, the sentence-level states which carry global context-level information were fed into the dynamic routing process between capsules. This improves the context-sensitivity of capsules compared to vanilla capsule networks considering the variability of background information of text.

3.1 Sentence-level and Word-level States

For the purpose of extracting n-gram features while preserving sequential dependencies of text, **(1.) global sentence-level states were utilized for context-level feature extraction, and (2.) word-level states were employed for each word to extract dependencies among words in the sequences of text as n-gram features.**

Generally, an encoded hidden state of a sentence using sentence-level states (Ψ_1^t and Ψ_2^t) and word-level states (h_i^t) at a given time step t could be represented as Eq. 1.

$$H^t = \langle h_1^t, h_2^t, h_3^t, \dots, h_n^t, \Psi_1^t, \Psi_2^t \rangle. \quad (1)$$

Here the state of the neural representation of encoded hidden state H^t of S-LSTM layer consists of a sub-state h_i^t for each word w_i . Since two dynamic routing procedures exist between three capsule layers in a vanilla capsule network, two sentence-level sub-states Ψ_1^t and Ψ_2^t are introduced to utilize towards the optimizations of dynamic routing procedure.

The suggested approach for the neural encoding of sentences utilizes recurrent information exchange between word-level states and sentence-level states to incrementally achieve a rich neural representation in each time step. The initial states are set as $H^0 = h_i^0 = \Psi_1^0 = \Psi_2^0 = h^0$ in the form of model parameters. The state transition of a word state h_i^t could be illustrated with the following equations.

$$\begin{aligned} \phi_i^t &= [h_{i-1}^{t-1}, h_i^{t-1}, h_{i+1}^{t-1}] \\ \tilde{\Psi} &= \text{avg}(\Psi_1^{t-1}, \Psi_2^{t-1}) \\ \hat{i}_i^t &= \sigma(W_i \phi_i^t + U_i x_i + V_i \tilde{\Psi} + b_i) \\ \hat{l}_i^t &= \sigma(W_l \phi_i^t + U_l x_i + V_l \tilde{\Psi} + b_l) \\ \hat{r}_i^t &= \sigma(W_r \phi_i^t + U_r x_i + V_r \tilde{\Psi} + b_r) \\ \hat{f}_i^t &= \sigma(W_f \phi_i^t + U_f x_i + V_f \tilde{\Psi} + b_f) \\ \hat{s}_{i1}^t &= \sigma(W_{s1} \phi_i^t + U_{s1} x_i + V_{s1} \tilde{\Psi} + b_{s1}) \\ \hat{s}_{i2}^t &= \sigma(W_{s2} \phi_i^t + U_{s2} x_i + V_{s2} \tilde{\Psi} + b_{s2}) \\ o_i^t &= \sigma(W_o \phi_i^t + U_o x_i + V_o \tilde{\Psi} + b_o) \\ u_i^t &= \tanh(W_u \phi_i^t + U_u x_i + V_u \tilde{\Psi} + b_u) \\ i_i^t, l_i^t, r_i^t, f_i^t, s_i^t &= \text{softmax}(\hat{i}_i^t, \hat{l}_i^t, \hat{r}_i^t, \hat{f}_i^t, \hat{s}_{i1}^t, \hat{s}_{i2}^t) \\ c_i^t &= l_i^t \odot c_{i-1}^{t-1} + f_i^t \odot c_i^{t-1} + r_i^t \odot c_{i+1}^{t-1} + s_{i1}^t \odot c_{\Psi_1}^{t-1} \\ &\quad + s_{i2}^t \odot c_{\Psi_2}^{t-1} + i_i^t \odot u_i^t \\ h_i^t &= o_i^t \odot \tanh(c_i^t) \end{aligned}$$

Here, ϕ_i^t is the representation of concatenated hidden vectors of the context window. $\tilde{\Psi}$ represents the averaged state of two sentence-level states Ψ_1 and Ψ_2 . Following the S-LSTM network (Zhang et al., 2018) six gates were applied to control information flow from cell state of previous time step, left cell state of previous time step, right cell state of previous time step, two global sentence states of previous time step, and input state of a given word to cell state of any given time step. These gates are denoted respectively as $f_i^t, l_i^t, r_i^t, s_{i1}^t, s_{i2}^t$ and i_i^t . o_i^t represents the gate that controls the information flow from the current cell state c_i^t to the current hidden state h_i^t . W_x, U_x, V_x and b_x are model parameters where $x \in \{i, o, l, r, f, s1, s2, u\}$.

The following equations illustrate the state transitions for global sentence-level states Ψ_1 and Ψ_2 .

$$\begin{aligned} \tilde{h} &= \text{avg}(h_1^{t-1}, h_2^{t-1}, \dots, h_n^{t-1}) \\ \hat{f}_{\Psi_1}^t &= \sigma(W_{\Psi_1} \Psi_1^{t-1} + U_{\Psi_1} \tilde{h} + b_{\Psi_1}) \\ \hat{f}_{\Psi_2}^t &= \sigma(W_{\Psi_2} \Psi_2^{t-1} + U_{\Psi_2} \tilde{h} + b_{\Psi_2}) \\ \hat{f}_i^t &= \sigma(W_f \tilde{\Psi} + U_{f_i} h_i^{t-1} + b_{f_i}) \\ o_{\Psi_1}^t &= \sigma(W_{o1} \Psi_1^{t-1} + U_{o1} \tilde{h} + b_{o1}) \\ o_{\Psi_2}^t &= \sigma(W_{o2} \Psi_2^{t-1} + U_{o2} \tilde{h} + b_{o2}) \\ f_1^t, f_2^t, \dots, f_n^t, f_{\Psi_1}^t, f_{\Psi_2}^t &= \text{softmax}(\hat{f}_1^t, \hat{f}_2^t, \dots, \hat{f}_n^t, \\ &\quad \hat{f}_{\Psi_1}^t, \hat{f}_{\Psi_2}^t) \\ c_{\Psi_1}^t &= f_{\Psi_1}^t \odot c_{\Psi_1}^{t-1} + \sum f_{i1}^t \odot c_i^{t-1} \\ c_{\Psi_2}^t &= f_{\Psi_2}^t \odot c_{\Psi_2}^{t-1} + \sum f_i^t \odot c_i^{t-1} \\ \Psi_1^t &= o_1^t \odot \tanh(c_{\Psi_1}^t) \\ \Psi_2^t &= o_2^t \odot \tanh(c_{\Psi_2}^t) \end{aligned}$$

Here, $f_1^t, f_2^t, \dots, f_n^t$ and $f_{\Psi_1}^t, f_{\Psi_2}^t$ are gates controlling information respectively from $c_1^{t-1}, c_2^{t-1}, \dots, c_n^{t-1}$ and $c_{\Psi_1}^{t-1}, c_{\Psi_2}^{t-1}$ to $c_{\Psi_1}^t$ and $c_{\Psi_2}^t$. $o_{\Psi_1}^t$ and $o_{\Psi_2}^t$ are gates that manipulate information flow from $c_{\Psi_1}^t$ and $c_{\Psi_2}^t$ to Ψ_1 and Ψ_2 . W_x, U_x and b_x are model parameters where $x \in \{\Psi_1, \Psi_2, f, o1, o2\}$.

3.2 Capsule Networks

As displayed in Figure 1, As a major enhancement to vanilla capsule network consisting of 4 layers namely n-gram convolutional layer, primary capsules, convolutional capsules and text capsules, concatenated word-level states, $\langle h_1^t, h_2^t, h_3^t, \dots, h_n^t \rangle$ were fed as input features. Furthermore, two global sentence-level states (Ψ_1

and Ψ_2) that are consolidated with context-level-information, are integrated with the dynamic routing process between primary capsules and convolutional capsules, and between convolutional capsules and text capsules (Zhao et al., 2018; Kim et al., 2020). We empirically evaluate the usage of the number of sentence-level states for optimal performance in Section 5.3.

3.2.1 N-gram Convolutional Layer

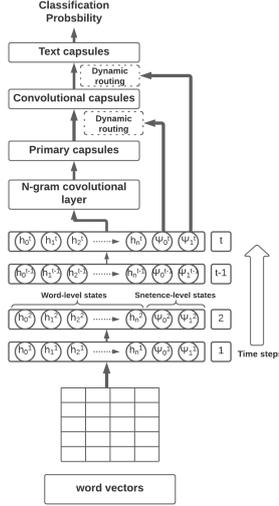


Figure 1: Dual-State Capsule Networks

For the purpose of extracting n-gram features from different positions of the concatenated word-level states (h_w^t), a standard convolutional layer was applied. Let $H \in \mathbb{R}^{d \times l}$ denotes the concatenated word-level states with l number of d dimensional word-level states. The convolution operation consists of filters of $C \in \mathbb{R}^{d \times f}$, which are applied to a context window of size f on top of the concatenated word states $H^{i:i+f}$, where i is the starting index of the context window. This produces new features governed by the Eq. 2.

$$x_i = \sigma(H^{i:i+f} \circ C + b_0) \quad (2)$$

Here, $x_i \in \mathbb{R}$ represents the generated feature, and \circ denotes the element-wise multiplication between context window and filter while σ and b_0 are the non-linear activation function (ReLU or tanh) and bias term respectively. This whole convolution process generates $M_i \in \mathbb{R}^{d-f+1}$ as feature column consists of features obtained through Eq. 2. With N number of filters, ultimately N number of feature columns were generated according to the following

Eq. 3.

$$M = [M_1, M_2, \dots, M_N] \in \mathbb{R}^{(d-f+1) \times N} \quad (3)$$

3.2.2 Primary Capsules

To convert the scalar output of the convolutional layer to the vector representation of capsules, a capsule layer was applied with $v_i \in \mathbb{R}^u$ as the instantiated parameters of the capsules with dimension u . The capsules are generated based on the matrix multiplication that includes matrix filters of $F \in \mathbb{R}^{N \times u}$. $(d-f+1)$ number of capsules were generated with the matrix multiplication process governed by the following Eq. 4.

$$v_i = \text{squash}(F \otimes M_i + b_1) \quad (4)$$

Here, v_i consists of a map of capsules that includes $(d-f+1)$ number of capsules. \otimes and b_1 denote the matrix multiplication operation and bias term, respectively. The squash function refers to the non-linearity displayed in Eq 10. With D number of filters, D number of maps were created, which consist of d dimensional $(d-f+1)$ number of capsules. Therefore, the primary capsule layer could be represented as $V \in \mathbb{R}^{(d-f+1) \times D \times u}$.

3.2.3 Convolutional Capsules

Inspired by the convolutional capsules introduced by Zhao et al. (2018), a convolutional capsule layer was applied immediately after the primary capsule layer, where a local region from the primary capsule layer is connected to the capsules of the convolutional capsule layer. This procedure improves the model performance due to the nature of the text, where certain objects of text within the local regions of text could get activated when identifying the most probable features towards the final outcome of the classification process.

Let the primary capsules in the region $l \times D$ to be mapped to the convolutional capsules using the weight matrix $W^c \in \mathbb{R}^{E \times u \times u}$, where $l \times D$ number of capsules are connected to each convolutional capsule to incrementally learn child-to-parent relationship. Here E is the number of parent capsules in the convolutional layer. Given each child capsule, the parent convolutional capsules are generated according to the following Eq. 5.

$$\hat{u}_{j|i} = W_j^c u_i + \hat{b}_{j|i} \quad (5)$$

Here $\hat{b}_{j|i}$ is the bias term for the mapping. The child capsules in the region $l \times D$ are denoted as u_i ,

which are mapped to the parent convolutional capsules. W_j^c is the j^{th} weight matrix for the learning by agreement procedure, which maps u_i to the $\hat{u}_{j|i}$. Thus, $(d - f - l + 2) \times D$ total number of u dimensional capsules were generated as convolutional capsules.

3.2.4 Text Capsules

The final capsule layer was designed to contain a number of capsules based on the number of classes in the text classification tasks. The capsules in the layer below were transformed based on the matrix multiplication process. This generates text capsules while learning the child-to-parent relationship process through the routing procedure.

The convolutional capsules were flattened to a list of capsules and were transformed into text capsules by a transformation matrix $W^d \in \mathbb{R}^{Z \times d \times d}$ considering the routing procedure to learn child-to-parent relationship. The final capsules were produced as $x_j \in \mathbb{R}^d$. Activation of the capsule denotes the class probability of a given text category as $a_j \in \mathbb{R}$ for each class. Here, Z denotes the number of text capsules.

3.2.5 Dynamic Routing Algorithm

The routing by agreement procedure under the dynamic routing algorithm incrementally ensures that appropriate child capsules are sent to parent capsules by iteratively creating a non-linear mapping between capsules (Sabour et al., 2017). Under this experiment, two dynamic routing procedure exist in the proposed architecture, between primary capsules and convolutional capsules, and between convolutional capsules and text capsules.

The dynamic routing procedure is initiated by initializing the log prior probabilities, which are the core components for the routing by agreement. As the proposed enhancement, **instead of initializing the log prior probabilities to values in statistical distribution or zeros, the log prior probabilities are assigned with values obtained through reshape operation on, global sentence-level states.** This strategy carries context-level information of the text to validate the routing procedure which incrementally improves child-to-parent relationship understanding background knowledge. The log prior probabilities are initialized as in Eq. 6.

$$b_{ij} \leftarrow \text{reshape}(\Psi) \quad (6)$$

Here b_{ij} represents the log prior probability between the child capsule i and the parent capsule j ,

while the reshape operation includes a matrix multiplication that maps the dimensions of sentence-level state Ψ to the dimensions of b_{ij} (Both Ψ_1 and Ψ_2 sentence-level states were considered under this operation.) The log-le prior probabilities between each child capsule i and each parent capsule j were normalized using the standard *softmax* function during the iterative routing procedure, which sums up coupling coefficients of each child capsule i to all parent capsules j in the layer above, to 1 as the Eq. 7.

$$c_{ij} \leftarrow \frac{\exp(b_{ij})}{\sum_k \exp(b_{ik})} \quad (7)$$

Here, c_{ij} represents the coupling coefficient that is used when each child capsule i in the layer below maps to each parent capsule j in the layer above. These log prior probabilities b_{ij} could be iteratively learned at the same time as all other weights, where the log prior probabilities only depend on the location of the child and parent capsules but not on the given sequence of the input text. Then the log prior probabilities between the child and parent capsules are updated in an iterative manner, considering the agreement measurement that repeatedly measures the similarity among predicted parent vectors by child capsules $\hat{u}_{j|i}$ (predicted parent capsule j given child capsule i) and current capsule s_j . The similarity measurement is a simple scalar product between the predicted and current parent capsules.

$$s_j = \sum_i c_{ij} * \hat{u}_{j|i} \quad (8)$$

$$b_{ij} \leftarrow b_{ij} + \hat{u}_{j|i} * s_j \quad (9)$$

The current parent capsule is calculated based on the coupling coefficients between the child and parent capsules c_{ij} , and the predicted parent capsules by the child capsules denoted as $\hat{u}_{j|i}$ according to the Eq. 8. The log prior probabilities b_{ij} are iteratively updated according to the Eq. 9 considering the measurement of the agreement between current and predicted parent capsules.

After the iterative routing procedure, the output capsules are normalized using the non-linear squash function, which transforms the length of each vector to represent the probability of the existence of the entity present within a capsule. The squash function converts the length of long vectors near to 1 while shrinking the length of short vectors to 0, which represents the probability as the

potential of an entity represented by each capsule. The non-linear squash function applied for parent capsules s_j that outputs v_j is defined as Eq. 10.

$$v_j = \frac{\|s_j\|^2}{1 + \|s_j\|^2} \frac{s_j}{\|s_j\|^2} \quad (10)$$

Here $\|s_j\|$ denotes the standard norm for s_j , and as illustrated in **Algorithm 1**, the non-linear squash function is applied outside the iterative mapping procedure. This results in eliminating the degradation of log prior probabilities in each routing iteration.

Algorithm 1: Dynamic Routing using Sentence State

```

1 procedure ROUTING( $\hat{u}_{j|i}, \Psi, r, l$ )
2 forall capsules  $i$  in layer  $l$  and  $j$  in
   layer  $l + 1$  do
3    $b_{ij} \leftarrow \text{reshape}(\Psi)$ 
4 for  $r$  iterations do
5   forall capsules  $i$  in layer  $l$  do
6      $c_i \leftarrow \text{softmax}(b_i)$ 
7   forall capsules  $j$  in layer  $l + 1$  do
8      $s_j \leftarrow \sum_i c_{ij} * \hat{u}_{ij}$ 
9   forall capsules  $i$  in layer  $l$  and  $j$  in
   layer  $l + 1$  do
10     $b_{ij} \leftarrow b_{ij} + u_j * s_j$ 
11 forall capsules  $j$  in layer  $l+1$  do
12    $v_j \leftarrow \text{squash}(s_j)$ 
13 return  $v_j$ 

```

3.2.6 Loss Function

For the text classification task, for each text capsule, we used a separate margin loss (Sabour et al., 2017) function to identify where a given text category is present within a given capsule. For text capsule s , the margin loss L_s is given by;

$$L_s = T_s \max(0, m^+ - \|v_s\|)^2 + \lambda(1 - T_s) \max(0, \|v_s\| - m^-)^2 \quad (11)$$

Here $T_s = 1$ if the text category exists within the text capsule, otherwise it is set to 0. m^+ and m^- are set as 0.9 and 0.1 accordingly. The down-weighting coefficient λ is set to 0.25 with the optimal performance.

4 Experiments

4.1 Data Sets

Experiments were conducted on six benchmark datasets covering multiple classification tasks. These include movie review classification (MR(2004) (Pang and Lee, 2004), MR(2005), IMDB), news article classification (Reuters10, MPQA (Wiebe et al., 2005)), and question categorization (TREC-QA (Li and Roth, 2002)). The details of each dataset are shown in table 1. For measuring the performance of the model against low resource language processing, a publicly available Sinhala dataset (Senevirathne et al.) was utilized. This includes 15059 news comments annotated with 4 sentiment categories with an average comment length of 10.

4.2 Implementation

For the experimental analysis, we utilized the 300-dimensional GloVe² word vectors, which consisted of 840 billion words. Adam optimizer was used for the optimization process with exponential learning rate decay. The models were trained on Google Colab with Tensorflow as the implementation tool. The optimal hyperparameters for models in each data set are indicated in Table 2.

For each experiment, the learning rate was set to $1e - 3$, and the learning rate decay was set to 0.95. The hidden word-level state dimension was set to 300, while sentence-level hidden state size was set to 600 dimensions. Max sentence length was chosen as the sentence length, considering the variations of the datasets. The n-gram convolutional layers were instantiated to extract 3-grams from each context window with 32 filters. Each capsule in the primary capsule layer was instantiated with 8-dimensional vectors, while each convolutional and text capsule was instantiated with 16-dimensional vectors. The length of each capsule denotes the existence of an entity within a capsule, which is further utilized with text capsules to identify the text categories within a given sequence of text.

4.3 Baseline models

In this experiment, we empirically evaluated our models against several baseline models such as LSTM, Bi-LSTM, CNN with randomly initialized vectors, CNN with non-trainable embeddings (CNN-static), CNN with trainable embed-

²<https://nlp.stanford.edu/projects/glove/>

Data	c	Train	Dev	Test	l_{avg}	$ V $
MR(2004)	2	1620	180	200	779	40693
MR(2005)	2	8635	960	1067	22	18764
Reuters10	10	6472	720	2787	168	28482
TREC-QA	6	4843	539	500	9	8689
MPQA	2	8587	955	1067	3	6246
IMDB	2	22500	2500	25000	231	112540

Table 1: Properties of datasets. c: Number of text categories, Train, Dev, Test: Size of training, development, and test sets (respectively), l_{avg} : Average sentence length, $|V|$: Size of the vocabulary

Data	Batch	Routing	Steps	CW	Epochs
MR(2004)	4	3	4	1	50
MR(2005)	8	3	8	2	20
Reuters10	8	2	4	2	20
TREC-QA	4	3	4	1	50
MPQA	4	4	4	1	20
IMDB	8	3	8	2	10

Table 2: Optimal hyperparameters for each dataset. Batch: Batch size, Routing: Number of iterations in the dynamic routing procedure, Steps: Number of recurrent time steps used to encode the sequence of text utilizing word and sentence states, CW: Number of word-states considered in both left and right context for encoding word-states in a particular time step.

dings (CNN-non-static) (Kim, 2014) and S-LSTM (Zhang et al., 2018) as primary baseline models. Furthermore, we evaluated the DS-Caps approach against several capsule networks based techniques namely, Capsule-A, Capsule-B (Zhao et al., 2018), and capsule networks with static routing (Kim et al., 2020). Some evaluation records for datasets under certain baseline techniques are not displayed in Table 3 as the results are not reported in the literature.

5 Evaluation

5.1 Evaluation on Benchmark Data Sets

Accuracy was chosen as the evaluation metric for our experiments, following related research in the same domain (Zhao et al., 2018). The experiment results are summarized in the Table 3, against the six benchmark datasets.

Our DS-Caps network was able to achieve the best result for four out of six benchmark results including MR(2004), MR(2005), Reuters and IMDB datasets, outperforming the previous studies on capsule networks with static routing for MR(2004) and MR(2005) datasets, Capsule-B for Reuters10 dataset, and CNN-non-static for IMDB dataset. In particular, the DS-Caps network extensively

and consistently defeats sequential neural networks such as LSTM, BiLSTM, Tree LSTM and S-LSTM networks, and spatial models such as CNN-rand, CNN-static, CNN-non-static on all six datasets. The observation was expected due to the language representation of capsules, where the vector representation could be greatly associated with the exact order or pose of sequences of text. The dynamic routing procedure introduced under the capsule networks further eliminated the loss of information due to the pooling strategy used under CNNs. This observation could be justified as the ability of the S-LSTM layer to encode the neural representation of text more efficiently, which subsequently enriches the dynamic routing procedure providing context-level information through global sentence-level information.

Furthermore, the results obtained for datasets with higher sequence lengths (MR(2004) and IMDB) illustrate greater performance compared to the baseline models. This provides shreds of evidence for the capability of the DS-Caps network to classify rather longer sequences of text utilizing sequential features and n-gram features of the text, collaborating with the vector representations of capsules.

5.2 Evaluation on Low Resource Languages

As mentioned in Section 1, to demonstrate that DS-Caps gives state-of-the-art results on low resource languages, we selected Sinhala. Sinhala is a low resource language, and the largest reported dataset is the CommonCrawl dataset, which is just above 100M (Lakmal et al., 2020). Furthermore, Sinhala does not have a pre-trained BERT mode, nor is it included in multiBERT. Using auxiliary features such as POS tags is also not an option for this language, as the best-reported results for POS tagging are not optimal (Fernando and Ranathunga, 2018). The proposed DS-caps network was evaluated against a Sinhala multi-class sentiment analysis task (Senevirathne et al.). As the results displayed in Table 4, DS-caps outperformed all previously tested deep learning approaches for Sinhala sentiment analysis. The hyper-parameters used for the DS-Caps networks were the same as the parameter used by Senevirathne et al. for Capsule-A and Capsule-B and fastText embeddings were used as primary features for sentiment analysis task.

DS-caps produced promising results in the context of this low resourced language, suggesting the possibility to be effective for any language as a uni-

	MR(2004)	MR(2005)	Reuters	TREC-QA	MPQA	IMDB
CNN-rand (Kim, 2014)	-	76.1	-	91.2	-	83.4
CNN-static (Kim, 2014)	-	81.0	-	92.8	-	89.6
CNN-non-static (Kim et al., 2020)	88.0	81.5	87.4	92.7	89.9	90.4
LSTM (Zhao et al., 2018)	-	75.9	-	86.8	-	-
Bi-LSTM (Zhao et al., 2018)	-	79.3	-	89.6	-	-
Tree-LSTM (Zhao et al., 2018)	-	80.7	-	91.8	-	-
S-LSTM (Zhang et al., 2018)	82.4	-	-	-	-	87.2
Capsule-A (Kim et al., 2020)	85.0	79.4	87.7	90.8	88.3	89.3
Capsule-B (Kim et al., 2020)	89.5	79.0	88.0	89.8	88.4	89.3
Capsule-static routing (Kim et al., 2020)	89.6	81.0	87.5	94.8	90.6	89.7
Dual-state capsule networks (DS-Caps)	90.5	82.1	88.6	92.8	89.2	90.6

Table 3: Text classification accuracies for benchmark datasets.

Model	Accuracy	Precision	Recall	F1-score
RNN	58.98	42.93	54.98	42.30
LSTM	62.88	70.95	51.93	54.50
GRU	62.78	60.93	62.78	54.83
Bi-LSTM	63.81	61.17	63.81	57.71
Stacked-BiLSTM	63.13	69.71	63.18	59.42
HAHNN	61.16	71.08	48.54	59.25
Capsule-A	61.89	56.12	61.89	53.55
Capsule-B	63.23	59.84	63.23	59.11
DS-Caps	64.03	61.68	64.03	61.33

Table 4: 10-fold cross-validated, weighted evaluation metrics for performance on Sinhala multi-class sentiment analysis (Senevirathne et al.)

No. of Sentence states	Loss function	Context window	Time steps	Accuracy
0	Margin	2	5	80.6
0	Cross entropy	2	7	79.4
0	Spread	2	10	79.9
0	Margin + L2	3	4	80.3
1	Margin	1	3	80.5
1	Cross entropy	2	3	79.1
1	Spread	2	4	80.4
1	Margin + L2	2	4	80.1
2	Margin	1	7	82.1
2	Cross entropy	2	7	79.9
2	Spread	2	3	79.7
2	Margin + L2	2	4	80.9

Table 5: Performance of DS-Caps networks varying no. of sentence-level states against lost function for MR(2005) dataset.

versal approach. Even though BERT models (Devlin et al., 2018) indicated superior performance in text classification tasks, they cannot be readily applied for low resourced languages due to lack of pre-trained models and high computation bottleneck. Also, the lack of linguistic resources for low resource languages, makes it impossible to use highly resource-intensive methodologies. Therefore DS-caps could be used as a proper substitution in place of text classification models that are based on computationally intensive BERT models.

5.3 Ablation Study

For the purpose of analyzing the performance of the model varying different components, an ablation study was conducted on MR(2005) dataset and reported in Table 5. The performance was analyzed with respect to loss functions proposed by (Zhao et al., 2018) for capsule networks and further adding L2 regularization to the optimal loss function. Furthermore, the effect of the number of sentence-level states was analyzed to find out the optimal number of sentence-level states to be integrated with the dynamic routing process. Under the experiments, the batch size was kept as 8 and, the number of epochs and dynamic routing iterations were fixed as 20 and 3 respectively. Margin loss with two sentence-level states produced the best performance for MR(2005) dataset.

This result was expected due to the efficiency of margin loss for vanilla capsule network (Zhao et al., 2018). Further two distinct sentence-level states facilitate separate learning procedures for dynamic routing algorithms to learn log prior probabilities which elevate part-to-whole relationship considering low-level and high-level-capsules.

6 Conclusion and Future Work

In this paper, we proposed a novel Dual-State Capsule (DS-Caps) network architecture, which incrementally improves the language representation considering local and global information of the text. As further enhancements, attention mechanisms could be integrated with DS-Caps. It is also worthwhile to integrate language resources such lexicons and contextual embeddings such as BERT (Devlin et al., 2018) with DS-Caps since most of the other deep learning approaches have produced greater performances with such strategies (Saha et al., 2020).

References

- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Stefan Daniel Dumitrescu, Andrei-Marius Avram, and Sampo Pyysalo. 2020. The birth of romanian bert. *arXiv preprint arXiv:2009.08712*.
- Sandareka Fernando and Surangika Ranathunga. 2018. Evaluation of different classifiers for sinhala pos tagging. In *2018 Moratuwa Engineering Research Conference (MERCCon)*, pages 96–101. IEEE.
- Jaeyoung Kim, Sion Jang, Eunjeong Park, and Sungchul Choi. 2020. Text classification using capsules. *Neurocomputing*, 376:214–221.
- Yoon Kim. 2014. Convolutional neural networks for sentence classification. *arXiv preprint arXiv:1408.5882*.
- Dimuthu Lakmal, Surangika Ranathunga, Saman Peramuna, and Indu Herath. 2020. Word embedding evaluation for sinhala. In *Proceedings of The 12th Language Resources and Evaluation Conference*, pages 1874–1881.
- Xin Li and Dan Roth. 2002. Learning question classifiers. In *COLING 2002: The 19th International Conference on Computational Linguistics*.
- Bo Pang and Lillian Lee. 2004. A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts. *arXiv preprint cs/0409058*.
- Bo Pang and Lillian Lee. 2005. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. *arXiv preprint cs/0506075*.
- Qiao Qian, Minlie Huang, Jinhao Lei, and Xiaoyan Zhu. 2016. Linguistically regularized lstms for sentiment classification. *arXiv preprint arXiv:1611.03949*.
- Sara Sabour, Nicholas Frosst, and Geoffrey E Hinton. 2017. Dynamic routing between capsules. In *Advances in neural information processing systems*, pages 3856–3866.
- Tulika Saha, Srivatsa Ramesh Jayashree, Sriparna Saha, and Pushpak Bhattacharyya. 2020. Bert-caps: A transformer-based capsule network for tweet act classification. *IEEE Transactions on Computational Social Systems*.
- Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. 2008. The graph neural network model. *IEEE Transactions on Neural Networks*, 20(1):61–80.
- Lahiru Senevirathne, Piyumal Demotte, Binod Karunanayake, Udyogi Munasinghe, and Surangika Ranathunga. Sentiment analysis for sinhala using deep learning techniques. *arXiv preprint arXiv:2011.07280v1*.
- Xingyou Wang, Weijie Jiang, and Zhiyong Luo. 2016. Combination of convolutional and recurrent neural network for sentiment analysis of short texts. In *Proceedings of COLING 2016, the 26th international conference on computational linguistics: Technical papers*, pages 2428–2437.
- Yequan Wang, Aixin Sun, Jialong Han, Ying Liu, and Xiaoyan Zhu. 2018. Sentiment analysis by capsules. In *Proceedings of the 2018 world wide web conference*, pages 1165–1174.
- Janyce Wiebe, Theresa Wilson, and Claire Cardie. 2005. Annotating expressions of opinions and emotions in language. *Language resources and evaluation*, 39(2-3):165–210.
- Yue Zhang, Qi Liu, and Linfeng Song. 2018. Sentence-state LSTM for text representation. *arXiv preprint arXiv:1805.02474*.
- Wei Zhao, Jianbo Ye, Min Yang, Zeyang Lei, Suofei Zhang, and Zhou Zhao. 2018. Investigating capsule networks with dynamic routing for text classification. *arXiv preprint arXiv:1804.00538*.