

phi-LSTM: A Phrase-based Hierarchical LSTM Model for Image Captioning

Ying Hua Tan and Chee Seng Chan

Center of Image and Signal Processing,
 Faculty of Computer Science & Information Technology,
 University of Malaya, Kuala Lumpur, Malaysia
 tanyinghua@siswa.um.edu.my; cs.chan@um.edu.my

Abstract. *A picture is worth a thousand words.* Not until recently, however, we noticed some success stories in understanding of visual scenes: a model that is able to detect/name objects, describe their attributes, and recognize their relationships/interactions. In this paper, we propose a phrase-based hierarchical Long Short-Term Memory (phi-LSTM) model to generate image description. The proposed model encodes sentence as a sequence of combination of phrases and words, instead of a sequence of words alone as in those conventional solutions. The two levels of this model are dedicated to i) learn to generate image relevant noun phrases, and ii) produce appropriate image description from the phrases and other words in the corpus. Adopting a convolutional neural network to learn image features and the LSTM to learn the word sequence in a sentence, the proposed model has shown better or competitive results in comparison to the state-of-the-art models on Flickr8k and Flickr30k datasets.

1 Introduction

Automatic caption/description generation from images is a challenging problem that requires a combination of visual information and linguistic as illustrated in Fig. 1. In other words, it requires not only complete image understanding, but also sophisticated natural language generation [1–4]. This is what makes it such an interesting task that has been embraced by both the computer vision and natural language processing communities.

One of the most common models applied for automatic caption generation is a neural network model that composes of two sub-networks [5–10], where a convolutional neural network (CNN) [11] is used to obtain feature representation

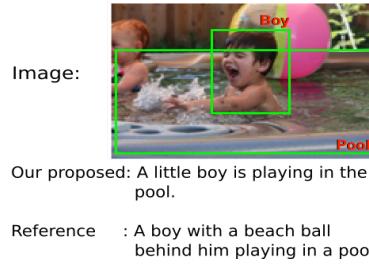


Fig. 1: Complete visual scene understanding is a holy grail in computer vision.

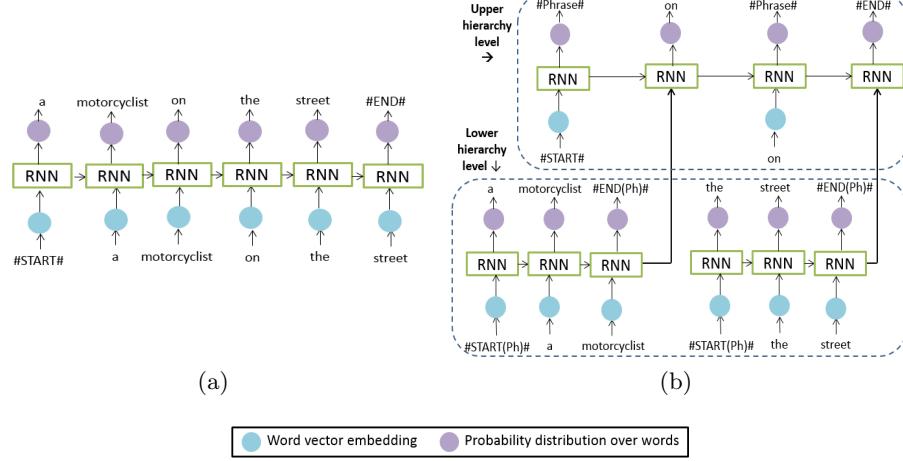


Fig. 2: Model comparison: (a) Conventional RNN language model, and (b) our proposed phrase-based model.

of an image; while a recurrent neural network (RNN)¹ is applied to encode and generate its caption description. In particular, Long Short-Term Memory (LSTM) model [12] has emerged as the most popular architecture among RNN, as it has the ability to capture long-term dependency and preserve sequence. Although sequential model is appropriate for processing sentential data, it does not capture any other syntactic structure of language at all. Nevertheless, it is undeniable that sentence structure is one of the prominent characteristics of language, and Victor Yngve - an influential contributor in linguistic theory stated in 1960 that “*language structure involving, in some form or other, a phrase-structure hierarchy, or immediate constituent organization*” [13]. Moreover, Tai et al. [14] proved that a tree-structured LSTM model that incorporates syntactic interpretation of sentence structure, can learn the semantic relatedness between sentences better than a pure sequential LSTM alone. This gives rise to question of whether is it a good idea to disregard other syntax of language in the task of generating image description.

In this paper, we would like to investigate the capability of a phrase-based language model in generating image caption as compared to the sequential language model such as [6]. To this end, we design a novel phrase-based hierarchical LSTM model, namely **phi-LSTM** to encode image description in three stages - chunking of training caption, image-relevant phrases composition as a vector representation and finally, sentence encoding with image, words and phrases. As opposed to those conventional RNN language models which process sentence as a sequence of words, our proposed method takes noun phrase as a unit in the

¹ RNN is a popular choice due to its capability to process arbitrary length sequences like language where words sequence governing its semantic is order-sensitive.

sentence, and thus processes the sentential data as a sequence of combination of both words and phrases together. Fig. 2 illustrates the difference between the conventional RNN language model and our proposal with an example. Both phrases and sentences in our proposed model are learned with two different sets of LSTM parameters, each models the probability distribution of word conditions on previous context and image. Such design is motivated by the observation that some words are more prone to appear in phrase, while other words are more likely to be used to link phrases. In order to train the proposed model, a new perplexity based cost function is defined. Experimental results using two publicly available datasets (Flickr8k [15] and Flickr30k [16]), and a comparison to the state-of-the-art results [5–7, 9, 33] have shown the efficacy of our proposed method.

2 Related Works

The image description generation task is generally inspired by two lines of research, which are (i) the learning of cross-modality transition or representation between image and language, and (ii) the description generation approaches.

2.1 Multimodal Representation and Transition

To model the relationship between image and language, some works associate both modalities by embedding their representations into a common space [17–20]. First, they obtain the image features using a visual model like CNN [18, 19], as well as the representation of sentence with a language model such as recursive neural network [19]. Then, both of them are embedded into a common multimodal space and the whole model is learned with ranking objective for image and sentence retrieval task. This framework was also tested at object level by Karpathy et al. [20] and proved to yield better results for the image and sentence bi-directional retrieval task. Besides that, there are works that learn the probability density over multimodal inputs using various statistical approaches. These include Deep Boltzmann Machines [21], topic models [22], log-bilinear neural language model [8, 23] and recurrent neural networks [5–7] etc. Such approaches fuse different input modalities together to obtain a unified representation of the inputs. It is notable to mention that there are also some works which do not explicitly learn the multimodal representation between image and language, but transit between modalities with retrieval approach. For example, Kuznetsova et al. [24] retrieve images similar to the query image from their database, and extract useful language segments (such as phrases) from the descriptions of the retrieved images.

2.2 Description Generation

On the other hand, caption generation approaches can generally be grouped into three categories as below:

Template-based. These approaches generate sentence from a fixed template [25–29]. For example, Farhadi et al. [25] infer a single triplet of object, action and scene from an image and convert it into a sentence with fixed template. Kulkarni et al. [26] use complex graph of detections to infer elements in sentence with conditional random field (CRF), but the generation of sentences is still based on the template. Mitchell et al. [28] and Gupta et al. [29] use a more powerful language parsing model to produce image description. In overall, all these approaches generate description which is syntactically correct, but rigid and not flexible.

Composition Method. These approaches extract components related to the images and stitch them up to form a sentence [24, 30, 31]. Description generated in such manner is broader and more expressive compared to the template-based approach, but is more computationally expensive at test time due to its non-parametric nature.

Neural Network. These approaches produce description by modeling the conditional probability of a word given multimodal inputs. For instance, Kiros et al. [8, 23] developed multimodal log-bilinear neural language model for sentence generation based on context and image feature. However, it has a fixed window context. The other popular model is recurrent neural network [5–7, 9, 32], due to its ability to process arbitrary length of sequential inputs such as sequence of words. This model is usually connected with a deep CNN that generates image features. The variants on how this sub-network is connected to the RNN have been investigated by different researchers. For instance, the multimodal recurrent neural network proposed by Mao et al. [5] introduces a multimodal layer at each time step of the RNN, before the softmax prediction of words. Vinyals et al. [6] treat the sentence generation task as a machine translation problem from image to English, and thus image feature is employed in the first step of the sequence trained with their LSTM RNN model.

2.3 Relation to Our Work

Automatic image caption generated via template-based [25–29] and composition methods [24, 30, 31] are typically two-stage approaches, where relevant elements such as objects (noun phrases) and relations (verb and prepositional phrases) are generated first before a full descriptive sentence is formed with the phrases. With the capability of LSTM model in processing long sequence of words, neural network based method that uses a two-stage approach deem unnecessary. However, we are still interested to find out how sequential model with phrase as a unit of sequence performs. The closest work related to ours is the one proposed by Lebret et al. [33]. They obtain phrase representation with simple word vector addition and learn its relevancy with image by training with negative samples. Sentence is then generated as a sequence of phrases, predicted using a statistical framework conditioned on previous phrases and its chunking tags. While their

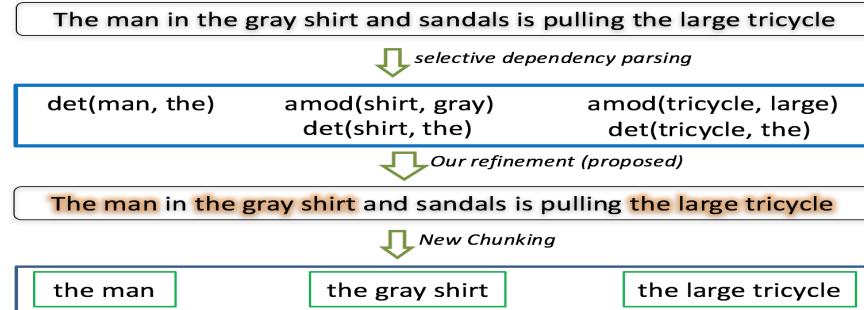


Fig. 3: Phrase chunking from dependency parse.

aim was to design a phrase-based model that is simpler than RNN, we intend to compare RNN phrase-based model with its sequential counterpart. Hence, our proposed model generates phrases and recomposes them into sentence with two sub-networks of LSTM, which are linked to form a hierarchical structure as shown in Fig. 2(b).

3 Our Proposed phi-LSTM Model

This section details how the proposed method encodes image description in three stages - i) chunking of image description, ii) encode words and phrases into distributed representations, and finally iii) encodes sentence with the phi-LSTM model.

3.1 Phrase Chunking

A quick overview on the structure of image descriptions reveals that, key elements which made up the majority of captions are usually noun phrases that describe the content of the image, which can be either objects or scene. These elements are linked with verb and prepositional phrases. Thus, noun phrase essentially covers over half of the corpus in a language model trained to generate image description. And so, in this paper, our idea is to partition the learning of noun phrase and sentence structure so that they can be processed more evenly, compared to extracting all phrases without considering their part of speech tag.

To identify noun phrases from a training sentence, we adopt the dependency parse with refinement using Stanford CoreNLP tool [34], which provides good semantic representation over a sentence by providing structural relationships between words. Though it does not chunk sentence directly as in constituency parse and other chunking tools, the pattern of noun phrase extracted is more flexible as we can select desirable structural relations. The relations we selected are:

- determiner relation (*det*),

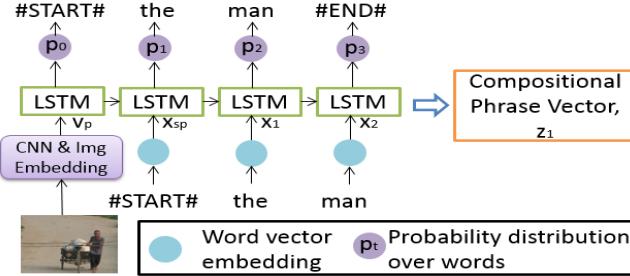


Fig. 4: Composition of phrase vector representation in phi-LSTM.

- numeric modifier (*nummod*),
- adjectival modifier (*amod*),
- adverbial modifier (*advmod*), but is selected only when the meaning of adjective term is modified, e.g. “*dimly lit room*”,
- compound (*compound*),
- nominal modifier for possessive alteration (*nmod:of* & *nmod:poss*).

Note that the dependency parse only extracts triplet made up of a governor word and a dependent word linked with a relation. So, in order to form phrase chunk with the dependency parse, we made some refinements as illustrated in Fig. 3. The triplets of selected relations in a sentence are first located, and those consecutive words (as highlighted in the figure, e.g. “the”, “man”) are grouped as a single phrase, while the standalone word (e.g. “in”) will remain as a unit in the sentence.

3.2 Compositional Vector Representation of Phrase

This section describes how compositional vector representation of a phrase is computed, given an image.

Image Representation. A 16-layer VggNet [35] pre-trained on ImageNet [36] classification task is applied to learn image feature in this work. Let $\mathbf{I} \in \mathbb{R}^D$ be an image feature, it is embedded into a K -dimensional vector, \mathbf{v}_p with image embedding matrix, $\mathbf{W}_{ip} \in \mathbb{R}^{K \times D}$ and bias $\mathbf{b}_{ip} \in \mathbb{R}^K$.

$$\mathbf{v}_p = \mathbf{W}_{ip}\mathbf{I} + \mathbf{b}_{ip}. \quad (1)$$

Word Embedding. Given a dictionary \mathcal{W} with a total of V vocabulary, where word $w \in \mathcal{W}$ denotes word in the dictionary, a word embedding matrix $\mathbf{W}_e \in \mathbb{R}^{K \times V}$ is defined to encode each word into a K -dimensional vector representation, \mathbf{x} . Hence, an image description with words $w_1 \dots w_M$ will correspond to vectors $\mathbf{x}_1 \dots \mathbf{x}_M$ accordingly.

Composition of Phrase Vector Representation. For each phrase extracted from the sentence, a LSTM-based RNN model similar to [6] is used to encode its sequence as shown in Fig. 4. Similar to [6], we treat the sequential modeling from image to phrasal description as a machine translation task, where the embedded image vector is inputted to the RNN on the first time step, followed by a start token $\mathbf{x}_{\text{sp}} \in \mathbb{R}^K$ indicating the translation process. It is trained to predict the next word at each time step by outputting $\mathbf{p}_{t_p+1} \in \mathbb{R}^{K \times V}$, which is modeled as the probability distribution over all words in the corpus. The last word of the phrase will predict an end token. So, given a phrase P which is made up by L words, the input \mathbf{x}_{t_p} at each time step are:

$$\mathbf{x}_{t_p} = \begin{cases} \mathbf{v}_P, & \text{if } t_p = -1 \\ \mathbf{x}_{\text{sp}}, & \text{if } t_p = 0 \\ \mathbf{W}_e w_{t_p}, & \text{for } t_p = 1 \dots L \end{cases} \quad (2)$$

For a LSTM unit at time step t_p , let $\mathbf{i}_{t_p}, \mathbf{f}_{t_p}, \mathbf{o}_{t_p}, \mathbf{c}_{t_p}$ and \mathbf{h}_{t_p} denote the input gate, forget gate, output gate, memory cell and hidden state at the time step respectively. Thus, the LSTM transition equations are:

$$\mathbf{i}_{t_p} = \sigma(\mathbf{W}_i \mathbf{x}_{t_p} + \mathbf{U}_i \mathbf{h}_{t_p-1}), \quad (3)$$

$$\mathbf{f}_{t_p} = \sigma(\mathbf{W}_f \mathbf{x}_{t_p} + \mathbf{U}_f \mathbf{h}_{t_p-1}), \quad (4)$$

$$\mathbf{o}_{t_p} = \sigma(\mathbf{W}_o \mathbf{x}_{t_p} + \mathbf{U}_o \mathbf{h}_{t_p-1}), \quad (5)$$

$$\mathbf{u}_{t_p} = \tanh(\mathbf{W}_u \mathbf{x}_{t_p} + \mathbf{U}_u \mathbf{h}_{t_p-1}), \quad (6)$$

$$\mathbf{c}_{t_p} = \mathbf{i}_{t_p} \odot \mathbf{u}_{t_p} + \mathbf{f}_{t_p} \odot \mathbf{c}_{t_p-1}, \quad (7)$$

$$\mathbf{h}_{t_p} = \mathbf{o}_{t_p} \odot \tanh(\mathbf{c}_{t_p}), \quad (8)$$

$$\mathbf{p}_{t_p+1} = \text{softmax}(\mathbf{h}_{t_p}). \quad (9)$$

Here, σ denotes a logistic sigmoid function while \odot denotes elementwise multiplication. The LSTM parameters $\{\mathbf{W}_i, \mathbf{W}_f, \mathbf{W}_o, \mathbf{W}_u, \mathbf{U}_i, \mathbf{U}_f, \mathbf{U}_o, \mathbf{U}_u\}$ are all matrices with dimension of $\mathbb{R}^{K \times K}$. Intuitively, each gating unit controls the extent of information updated, forgotten and forward-propagated while the memory cell holds the unit internal memory regarding the information processed up to current time step. The hidden state is therefore a gated, partial view of the memory cell of the unit. At each time step, the probability distribution of words outputted is equivalent to the conditional probability of word given the previous words and image, $P(w_t | w_{1:t-1}, I)$. On the other hand, the hidden state at the last time step L is used as the compositional vector representation of the phrase, $\mathbf{z} \in \mathbb{R}^K$, where $\mathbf{z} = \mathbf{h}_L$.

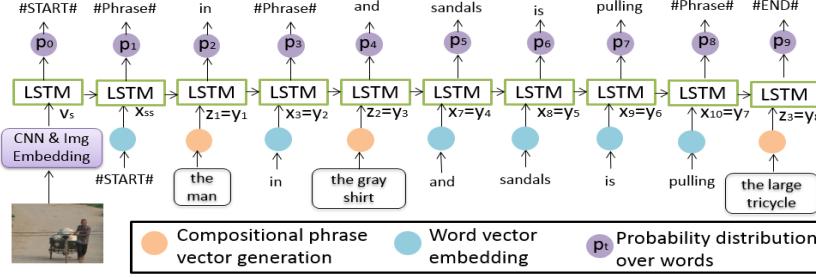


Fig. 5: Sentence encoding with phi-LSTM.

3.3 Encoding of Image Description

Once the compositional vector of phrases are obtained, they are linked with the remaining words in the sentence using another LSTM-based RNN model as shown in Fig. 5. Another start token $\mathbf{x}_{ss} \in \mathbb{R}^K$ and image representation $\mathbf{v}_s \in \mathbb{R}^K$ are introduced, where

$$\mathbf{v}_s = \mathbf{W}_{is}\mathbf{I} + \mathbf{b}_{is}, \quad (10)$$

with $\mathbf{W}_{is} \in \mathbb{R}^{K \times D}$ and bias $\mathbf{b}_{is} \in \mathbb{R}^K$ as embedding parameters. Hence, the input units of the LSTM in this level will be the image representation \mathbf{v}_s , start token \mathbf{x}_{ss} , followed by either compositional vector of phrase \mathbf{z} or word vector \mathbf{x} in accordance to the sequence of its description.

For simplicity purpose, the arranged input sequence will be referred as \mathbf{y} . Therefore, given the example in Fig. 4-5, the LSTM input sequence of the sentence will be $\{\mathbf{v}_s, \mathbf{x}_{ss}, \mathbf{y}_1 \dots \mathbf{y}_N\}$ where $N = 8$, and it is equivalent to sequence $\{\mathbf{v}_s, \mathbf{x}_{ss}, \mathbf{z}_1, \mathbf{x}_3, \mathbf{z}_2, \mathbf{x}_7, \mathbf{x}_8, \mathbf{x}_9, \mathbf{x}_{10}, \mathbf{z}_3\}$, as in Fig. 5. Note that a phrase token is added to the vocabulary, so that the model can predict it as an output when the next input is a noun phrase.

The encoding of the sentence is similar to the phrase vector composition. Eq. ??-?? are applied here using \mathbf{y}_{t_s} as input instead of \mathbf{x}_{t_p} , where t_p and t_s represent time step in phrase and sentence respectively. A new set of model parameters with same dimensional size is used in this hierarchical level.

4 Training the phi-LSTM Model

The proposed phi-LSTM model is trained with log-likelihood objective function computed from the perplexity² of sentence conditioned on its corresponding image in the training set. Given an image \mathbf{I} and its description \mathbf{S} , let R be the number of phrases of the sentence, P_i correspond to the number of LSTM blocks processed to get the compositional vector of phrase i , Q is the length of composite sequence of sentence \mathbf{S} , while \mathbf{p}_{tp} and \mathbf{p}_{ts} are the probability output

² Perplexity is a standard approach to evaluate language model.

of LSTM block at time step $t_p - 1$ and $t_s - 1$ for phrase and sentence level respectively. The perplexity of sentence \mathbf{S} given its image \mathbf{I} is

$$\log_2 \mathcal{PPL}(\mathbf{S}|\mathbf{I}) = -\frac{1}{N} \left[\sum_{t_s=-1}^Q \log_2 \mathbf{p}_{ts} + \sum_{i=1}^R \left[\sum_{t_p=-1}^{P_i} \log_2 \mathbf{p}_{tp} \right] \right], \quad (11)$$

where

$$N = Q + \sum_{i=1}^R P_i. \quad (12)$$

Hence, with M number of training samples, the cost function of our model is:

$$\mathcal{C}(\theta) = -\frac{1}{L} \sum_{j=1}^M [N_j \log_2 \mathcal{PPL}(\mathbf{S}_j|\mathbf{I}_j)] + \lambda_\theta \cdot \|\theta\|_2^2, \quad (13)$$

where

$$L = M \times \sum_{j=1}^M N_j. \quad (14)$$

It is the average log-likelihood of word given their previous context and the image described, summed with a regularization term, $\lambda_\theta \cdot \|\theta\|_2^2$, average over the number of training samples. Here, θ is the parameters of the model.

This objective however, does not discern on the appropriateness of different inputs at each time step. So, given multiple possible inputs, it is unable to distinguish which phrase is the most probable input at that particular time step during the decoding stage. That is, when a phrase token is inferred as the next input, all possible phrases will be inputted in the next time step. The candidate sequences are then ranked according to their perplexity up to this time step, where only those with high probability are kept. Unfortunately, this is problematic because subject in an image usually has much lower perplexity as compared to object and scene. Thus, such algorithm will end up generating description made up of only variants of subject noun phrases.

To overcome this limitation, we introduce a phrase selection objective during the training stage. At all time steps when an input is a phrase, H number of randomly selected phrases that are different from the ground truth input is feed into the phi-LSTM model as shown in Fig. 6. The model will then produce two outputs, which are the next word prediction solely based on the actual input, and a classifier output that distinguishes the actual one from the rest. Though the number of inputs at these time steps increases, the memory cell and hidden state that is carried to the next time step keep only information of the actual input. The cost function for phrase selection objective of a sentence is

$$\mathcal{C}_{PS} = \sum_{t_s \in \mathcal{P}} \sum_{k=1}^{H+1} \kappa_{t_s k} \sigma(1 - y_{t_s k} h_{t_s k} \mathbf{W}_{PS}). \quad (15)$$

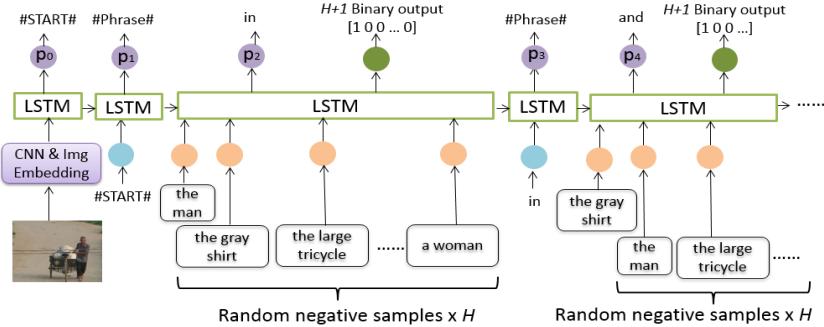


Fig. 6: Upper hierarchy of phi-LSTM with phrase selection objective.

where \mathcal{P} is the set of all time steps where the input is phrase, $h_{t_s k}$ is the hidden state output at time step t_s from input k , and $y_{t_s k}$ is its label which is +1 for the actual input and -1 for the false inputs. $\mathbf{W}_{ps} \in \mathbb{R}^{K \times 1}$ is trainable parameters for the classifier while $\kappa_{t_s k}$ scales and normalizes the objective based on the number of actual and false inputs at each time step. The overall objective function is then

$$\mathcal{C}_F(\theta) = -\frac{1}{L} \sum_{j=1}^M [N_j \log_2 \mathcal{P} \mathcal{P} \mathcal{L}(\mathbf{S}_j | \mathbf{I}_j) + \mathcal{C}_{PSj}] + \lambda_\theta \cdot \|\theta\|_2^2 . \quad (16)$$

This cost function is minimized and backpropagated with RMSprop optimizer [37] and trained in a minibatch of 100 image-sentence pair per iteration. We cross-validate the learning rate and weight decay depending on dataset, and dropout regularization [38] is employed over the LSTM parameters during training to avoid overfitting.

5 Image Caption Generation

Generation of textual description using the phi-LSTM model given an image is similar to other statistical language models, except that the image relevant phrases are generated first in the lower hierarchical level of the proposed model. Here, embedded image feature of the given image followed by the start token of phrase are inputted into the model, acting as the initial context required for phrase generation. Then, the probability distribution of the next word over the vocabulary is obtained at each time step given the previous contexts, and the word with the maximum probability is picked and fed into the model again to predict the subsequent word. This process is repeated until the end token for phrase is inferred. As we usually need multiple phrases to generate a sentence, beam search scheme is applied and the top K phrases generated are kept as the candidates to form the sentence. To generate a description from the phrases, the upper hierarchical level of the phi-LSTM model is applied in a similar fashion. When a phrase token is inferred, K phrases generated earlier are used as the

inputs for the next time step. Keeping only those phrases which generate positive result with the phrase selection objective, inference on the next word given the previous context and the selected phrases is performed again. This process iterates until the end token is inferred by the model.

Some constraints are added here, which are i) each predicted phrase may only appear once in a sentence, ii) maximum number of unit (word or phrase) that made up a sentence is limited to 20, iii) maximum number of words forming a phrase is limited to 10, and iv) generated phrases with perplexity higher than threshold T are discarded.

6 Experiment

6.1 Datasets

The proposed phi-LSTM model is tested on two benchmark datasets - Flickr8k [15] and Flickr30k [16], and compared to the state-of-the-art methods [5–7,9,33]. These datasets consist of 8000 and 31000 images respectively, each annotated with five ground truth descriptions from crowd sourcing. For both datasets, 1000 images are selected for validation and another 1000 images are selected for testing; while the rest are used for training. All sentences are converted to lower case, with frequently occurring punctuations removed and word that occurs less than 5 times (Flickr8k) or 8 times (Flickr30k) in the training data discarded. The punctuations are removed so that the image descriptions are consistent with the data shared by Karpathy and L. Fei-Fei [7].

6.2 Results Evaluated with Automatic Metric

Sentence generated using the phi-LSTM model is evaluated with automatic metric known as the bilingual evaluation understudy (BLEU) [39]. It computes the n-gram co-occurrence statistic between the generated description and multiple reference sentences by measuring the n-gram precision quality. It is the most commonly used metric in this literature.

Table 1 shows the performance of our proposed model in comparison to the current state-of-the-art methods. NIC [6] which is used as our baseline is a reimplementation, and thus its BLEU score reported here is slightly different from the original work. Our proposed model performs better or comparable to the state-of-the-art methods on both Flickr8k and Flickr30k datasets. In particular, we outperform our baseline on both datasets, as well as PbIC [33] - a work that is very similar to us on Flickr30k dataset by at least 5-10%.

As mentioned in Section 5, we generate K phrases from each image and discard those with perplexity higher than a threshold value T , when generating the

³ The BLEU score reported here is computed on our implementation of NIC [6], and the bracketed value is the reported score by the author.

⁴ The BLEU score reported here is cited from [7], and the bracketed value is the reported score by the author.

Table 1: BLEU score of generated sentence on (a) Flickr8k and (b) Flickr30k dataset.

(a)					(b)				
Models	Flickr8k				Flickr30k				
	B-1	B-2	B-3	B-4	B-1	B-2	B-3	B-4	
DeepVs [7]	57.9	38.3	24.5	16.0	57.3	36.9	24.0	15.7	
NIC [6] ³	60.2(63)	40.4	25.9	16.5	60	41	28	19	
phi-LSTM	63.6	43.6	27.6	16.6	66.3(66)	42.3	27.7	18.3	

Models	Flickr30k			
	B-1	B-2	B-3	B-4
DeepVS [7]	57.3	36.9	24.0	15.7
mRNN [5]	60	41	28	19
NIC [6] ⁴	66.3(66)	42.3	27.7	18.3
LRCNN [9]	58.7	39.1	25.1	16.5
PbIC [33]	59	35	20	12
phi-LSTM	66.6	45.8	28.2	17.0

Table 2: Vocab size, word occurrence and average caption length in train data, test data, and generated description on Flickr8k dataset.

	Train Data		Test Data			Gen. Caption	
	Number of sentence	30000	5000	1000	1000	NIC [6]	phi-LSTM
Actual	30000	5000	1000	1000			
Trained							
Size of vocab	7371	2538	3147	1919	1507	1187	128
Number of words	324481	316423	54335	52683	11139	10806	8275
Avg. caption length	10.8	10.5	10.9	10.5	11.1	10.8	8.3
							6.8

image caption. In order to understand how these two parameters affect our generated sentence, we use different K and T to generate the image caption with our proposed model trained on the Flickr30k dataset. Changes of the BLEU score against T and K are plotted in Fig. 7. It is shown that K does not have a significant effect on the BLEU score, when T is set to below 5.5. On the other hand, unigram and bi-gram BLEU scores improve with lower perplexity threshold, in contrast to tri-gram and 4-gram BLEU scores that reach an optimum value when $T=5.2$. This is because the initial (few) generated phrases with the lowest perplexity are usually different variations of phrase describing the same entity, such as ‘*a man*’ and ‘*a person*’. Sentence made with only such phrases has higher chance to match with the reference descriptions, but it would hardly get a match on tri-gram and 4-gram. In order to avoid generating caption made from only repetition of similar phrases, we select T and K which yield the highest 4-gram BLEU score, which are $T=6.5$ and $K=6$ on Flickr8k dataset, and $T=5.2$ and $K=5$ on Flickr30k dataset. A few examples are shown in Fig. 8.

6.3 Comparison of phi-LSTM with Its Sequence Model Counterpart

To compare the differences between a phrase-based hierarchical model and a pure sequence model in generating image caption, the phi-LSTM model and NIC [6] are both implemented using the same training strategy and parameter tuning. We are interested to know how well the corpus is trained by both models. Using

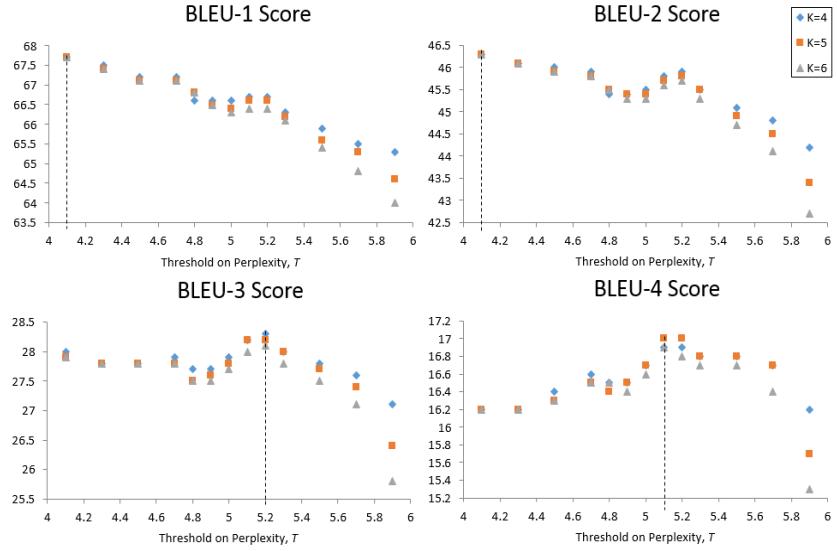


Fig. 7: Effect of perplexity threshold T and maximum number of phrases used for generating sentence, K on BLEU score.

the Flickr8k dataset, we computed the corpus information of i) the training data, ii) the reference sentences in the test data and iii) the generated captions as tabulated in Table 2. We remove words that occur less than 5 times in the training data, and it results in 4833 words being removed. However, this reduction in term of word count is only 2.48%. Furthermore, even though the model is evaluated in comparison to all reference sentences in the test data, there are actually 1228 words within the references that are not in our training corpus. Thus, it is impossible for the model to predict those words, and this is a limitation on scoring with references in all language models. For a better comparison with the 1000 generated captions, we also compute another reference corpus based on the first sentence of each test image. From Table 2, it can be seen that even though there are at least 1187 possible words to be inferred with images in the test set, the generated descriptions are made up from only 128 and 154 words in NIC [6] and phi-LSTM model, respectively. These numbers show that the actual number of words learned by these two models are barely 10%, suggesting more research is necessary to improve the learning efficiency in this field. Nevertheless, it shows that introducing the phrase-based structure in sequential model still improves the diversity of caption generated.

To get further insight on how the word occurrence in the training corpus affects the word prediction when generating caption, we record the top five, most trained words that are missing from the corpus of generated captions, and the top five, least trained words that are predicted by both models when generating description, as shown in Table 3. We consider only those words that appear in

Table 3: Top 5 (a) least trained word found, and (b) most trained word missing, from generated caption in the Flickr8k dataset.

(a)		(b)					
NIC [6]		phi-LSTM		NIC [6]		phi-LSTM	
Word	Occurrence	Word	Occurrence	Word	Occurrence	Word	Occurrence
<i>obstacle</i>	93	<i>overlooking</i>	81	<i>to</i>	2306	<i>while</i>	1443
<i>surfer</i>	127	<i>obstacle</i>	93	<i>his</i>	1711	<i>green</i>	931
<i>bird</i>	148	<i>climber</i>	96	<i>while</i>	1443	<i>by</i>	904
<i>woods</i>	155	<i>course</i>	106	<i>three</i>	1052	<i>one</i>	876
<i>snowboarder</i>	166	<i>surfer</i>	127	<i>small</i>	940	<i>another</i>	713

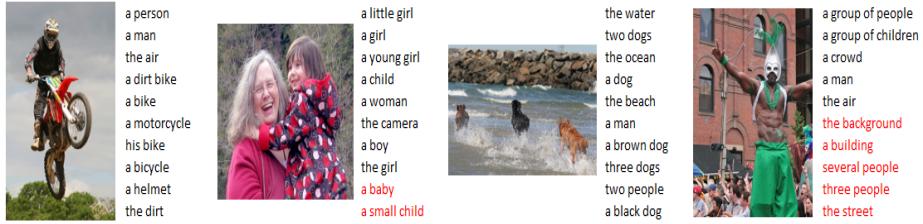


Fig. 8: Example of phrases generated from image by lower hierarchical level of phi-LSTM model. Red fonts indicate that the log probability of that phrases is below threshold.

the reference sentences to ensure that these words are related to the images in the test data. It appears that the phrase-based model is able to infer more words which are less trained, compared to the sequence model. Among the top five words that are not predicted, even though they have high occurrence in the training corpus, it can be seen that those words are either not very observable in the images, or are more probable to be described with other alternative. For example, *the* is a more probable alternative of *another*.

A few examples of the image description generated with our proposed model and NIC model [6] are shown in Fig. 9. It can be seen that both models are comparable qualitatively. An interesting example is shown in the first image where our model mis-recognizes the statue as a person, but is able to infer the total number of “persons” within the image. The incorrect recognition stems from insufficient training data on the word *statue* in the Flickr8k dataset, as it only occurs for 48 times, which is about 0.015% in the training corpus.

7 Additional Results

In order to further demonstrate the capability of our proposed model - the phi-LSTM, additional results from the test set of both Flickr8k and Flickr30k datasets are shown in Fig. 10 and Fig. 11, respectively. The results are selected such that images with very similar content are not repeatedly reported.

Image:			
NIC:	A group of people are standing in front of a building.	A man is doing a trick on a skateboard.	Two dogs play in the grass.
phi-LSTM:	Three people are standing in front of three men.	A skateboarder does a trick on a ramp.	Three dogs play in a grassy field.
Reference:	A group of tourists stand around as a lady puts her hand near the mouth of a statue.	A skateboarder in the air at a big outdoor ramp.	The three dogs ran in the ya

Fig. 9: Examples of caption generated with phi-LSTM model, in comparison to NIC [6].

Fig. 10 shows the outputs from the Flickr8k dataset. In the first row, it can be seen that our proposed model is able to distinguish different actions performed by the same subject (i.e. dog), from “*playing in the field*” to “*racing*” to “*jumping to catch a toy*”. In the second row, we demonstrate the capability of the proposed phi-LSTM model in identifying three different sports with very similar appearance in action. In particular, our model managed to detect and recognize a bicycle in the third image, even though the size of the bicycle is very small. Beside that, we also show that our proposed model is able to determine the number of subject(s) to certain extent. For example, it can identify “*two dogs*” and “*a group of women*”.

Fig. 11 presents the outputs from the Flickr30k. Images in first row show three running actions performed by a dog and a horse in different scenes, in which the captions generated by our proposed model have correctly described them. Then, all images in the second row and the first image in the third row once again demonstrate the capability of the phi-LSTM in identifying subjects, number of subjects and scene correctly. The last two images in the third row show that our proposed model is capable of recognizing a bike, regardless the object is displayed in a partial view or a complete view. Lastly, all images in the final row display different subjects in the water, and our proposed method is able to describe each of the subjects correctly (i.e. girl, man and surfer).

Also, note that these results show that the captions generated from our proposed model are in free form, instead of fixed template like subject-verb-object or subject-action-scene. Some descriptions may describe the scene while others may not, and verb is also an optional in the description generated. The only recurring element is the subject, which is essential in the task of image description.

Fig. 12 shows some examples of our proposed method that have some errors in the generated captions, such as the number of subjects, actions, negligence of simultaneous action performed by subject and more specific object etc. However, it is still able to generate description that is somewhat related to the image. From

our investigation, in any case, there are hardly any generated captions that infer a totally unrelated subject in the test set.

8 Conclusion

In this paper, we present the phi-LSTM model, which is a neural network model trained to generate reasonable description on image. The model consists of a CNN sub-network connected to a two-hierarchical level RNN, in which the lower level encodes noun phrases relevant to the image; while the upper level learns the sequence of words describing the image, with phrases encoded in the lower level as a unit. A phrase selection objective is coupled when encoding the sentence. It is designed to aid the generation of caption from relevant phrases. This design preserves syntax of sentence better, by treating it as a sequence of phrases and words instead of a sequence of words alone. Such adaptation also splits the content to be learned by the model into two, which are stored in two sets of parameters. Thus, it can generate sentence which is more accurate and with more diverse corpus, as compared to a pure sequence model.

Acknowledgement. This research is supported by the Fundamental Research Grant Scheme (FRGS) MoHE Grant from the Ministry of Higher Education Malaysia. We also would like to thank NVIDIA for the GPU donation.



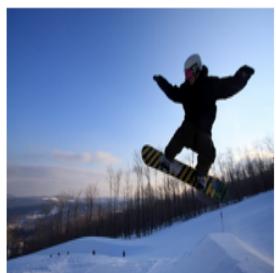
Two dogs play in a grassy field.



A dog in a race.



A small dog jumps to catch a toy.



A snowboarder in the air.



A skateboarder does a trick on a skateboard.



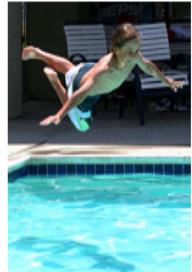
A person does a trick on a bicycle.



A person in a helmet is riding a dirt bike.



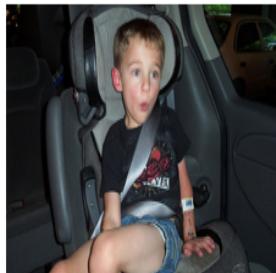
A surfer in a wave.



A young boy jumps into a swimming pool.



A group of women in the camera.



A little boy in a car.



A child in a swing.

Fig. 10: Flickr8k dataset: Sample image captioning results. It can be seen that our proposed model is able to distinguish different actions performed by the same subject.



A brown dog is running through the grass.



A dog is running through the snow.



A cowboy is riding a horse.



A group of people in the snow.



A woman in the snow.



A woman in the street.



A group of people in a field.



A person is riding a dirt bike.



A man is riding a bike.



A girl in the water.



A man in the water.



A surfer in the water.

Fig. 11: Flickr30k dataset: Sample image captioning results. It can be seen that our proposed model is able to distinguish different actions performed by the same subject.



A man in a boat in the water.



A child in a slide.



A man is playing a guitar.



A dog jumps over an obstacle course.



A little girl in a red jacket is standing in the snow.



A woman is holding a young boy.

Fig. 12: Examples of image captions generated with minor errors.

References

1. Sadeghi, M.A., Farhadi, A.: Recognition using visual phrases. In: CVPR. (2011) 1745–1752
2. Gupta, A., Verma, Y., Jawahar, C.: Choosing linguistics over vision to describe images. In: AAAI. (2012) 606–612
3. Bernardi, R., Cakici, R., Elliott, D., Erdem, A., Erdem, E., Ikizler-Cinbis, N., Keller, F., Muscat, A., Plank, B.: Automatic description generation from images: A survey of models, datasets, and evaluation measures. *Journal of Artificial Intelligence Research* **55** (2016) 409–442
4. Rasiwasia, N., Costa Pereira, J., Coviello, E., Doyle, G., Lanckriet, G.R., Levy, R., Vasconcelos, N.: A new approach to cross-modal multimedia retrieval. In: Proceedings of the ACM international conference on Multimedia. (2010) 251–260
5. Mao, J., Xu, W., Yang, Y., Wang, J., Huang, Z., Yuille, A.: Deep captioning with multimodal recurrent neural networks (m-rnn). arXiv preprint arXiv:1412.6632 (2014)
6. Vinyals, O., Toshev, A., Bengio, S., Erhan, D.: Show and tell: A neural image caption generator. In: CVPR. (2015) 3156–3164
7. Karpathy, A., Fei-Fei, L.: Deep visual-semantic alignments for generating image descriptions. In: CVPR. (2015) 3128–3137
8. Kiros, R., Salakhutdinov, R., Zemel, R.: Unifying visual-semantic embeddings with multimodal neural language models. arXiv preprint arXiv:1411.2539 (2014)
9. Donahue, J., Anne Hendricks, L., Guadarrama, S., Rohrbach, M., Venugopalan, S., Saenko, K., Darrell, T.: Long-term recurrent convolutional networks for visual recognition and description. In: CVPR. (2015) 2625–2634

10. Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A., Salakhudinov, R., Zemel, R., Bengio, Y.: Show, attend and tell: Neural image caption generation with visual attention. In: ICML. (2015) 2048–2057
11. Krizhevsky, A., Sutskever, I., Hinton, G.: Imagenet classification with deep convolutional neural networks. In: NIPS. (2012) 1097–1105
12. Hochreiter, S., Schmidhuber, J.: Long short-term memory. *Neural Computation* **9** (1997) 1735–1780
13. Yngve, V.: A model and an hypothesis for language structure. *Proceedings of the American Philosophical Society* **104** (1960) 444–466
14. Tai, K.S., Socher, R., Manning, C.: Improved semantic representations from tree-structured long short-term memory networks. arXiv preprint arXiv:1503.00075 (2015)
15. Rashtchian, C., Young, P., Hodosh, M., Hockenmaier, J.: Collecting image annotations using amazon’s mechanical turk. In: Proceedings of the NAACL HLT 2010 Workshop on Creating Speech and Language Data with Amazon’s Mechanical Turk. (2010) 139–147
16. Young, P., Lai, A., Hodosh, M., Hockenmaier, J.: From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics* **2** (2014) 67–78
17. Hodosh, M., Young, P., Hockenmaier, J.: Framing image description as a ranking task: Data, models and evaluation metrics. *Journal of Artificial Intelligence Research* (2013) 853–899
18. Frome, A., Corrado, G.S., Shlens, J., Bengio, S., Dean, J., Mikolov, T., et al.: Devise: A deep visual-semantic embedding model. In: NIPS. (2013) 2121–2129
19. Socher, R., Karpathy, A., Le, Q.V., Manning, C.D., Ng, A.: Grounded compositional semantics for finding and describing images with sentences. *Transactions of the Association for Computational Linguistics* **2** (2014) 207–218
20. Karpathy, A., Joulin, A., Fei-Fei, L.: Deep fragment embeddings for bidirectional image sentence mapping. In: NIPS. (2014) 1889–1897
21. Srivastava, N., Salakhutdinov, R.: Multimodal learning with deep boltzmann machines. In: NIPS. (2012) 2222–2230
22. Jia, Y., Salzmann, M., Darrell, T.: Learning cross-modality similarity for multinomial data. In: ICCV. (2011) 2407–2414
23. Kiros, R., Salakhutdinov, R., Zemel, R.: Multimodal neural language models. In: ICML. (2014) 595–603
24. Kuznetsova, P., Ordonez, V., Berg, T., Choi, Y.: Treetalk: Composition and compression of trees for image descriptions. *Transactions of the Association for Computational Linguistics* **2** (2014) 351–362
25. Farhadi, A., Hejrati, M., Sadeghi, M.A., Young, P., Rashtchian, C., Hockenmaier, J., Forsyth, D.: Every picture tells a story: Generating sentences from images. In: ECCV. (2010) 15–29
26. Kulkarni, G., Premraj, V., Ordonez, V., Dhar, S., Li, S., Choi, Y., Berg, A., Berg, T.: Babytalk: Understanding and generating simple image descriptions. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **35** (2013) 2891–2903
27. Yang, Y., Teo, C.L., Daumé III, H., Aloimonos, Y.: Corpus-guided sentence generation of natural images. In: Proceedings of the Conference on Empirical Methods in Natural Language Processing. (2011) 444–454
28. Mitchell, M., Han, X., Dodge, J., Mensch, A., Goyal, A., Berg, A., Yamaguchi, K., Berg, T., Stratos, K., Daumé III, H.: Midge: Generating image descriptions from computer vision detections. In: Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics. (2012) 747–756

29. Gupta, A., Mannem, P.: From image annotation to image description. In: ICONIP. (2012) 196–204
30. Li, S., Kulkarni, G., Berg, T., Berg, A., Choi, Y.: Composing simple image descriptions using web-scale n-grams. In: Proceedings of the Fifteenth Conference on Computational Natural Language Learning. (2011) 220–228
31. Kuznetsova, P., Ordonez, V., Berg, A., Berg, T., Choi, Y.: Collective generation of natural image descriptions. In: ACL. (2012) 359–368
32. Chen, X., Zitnick, L.: Learning a recurrent visual representation for image caption generation. arXiv preprint arXiv:1411.5654 (2014)
33. Lebret, R., Pinheiro, P., Collobert, R.: Phrase-based image captioning. arXiv preprint arXiv:1502.03671 (2015)
34. Manning, C., Surdeanu, M., Bauer, J., Finkel, J., Bethard, S., McClosky, D.: The Stanford CoreNLP natural language processing toolkit. In: Proceedings of 52nd Annual Meeting of the Association for Computational Linguistics: System Demonstrations. (2014) 55–60
35. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556 (2014)
36. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: CVPR. (2009) 248–255
37. Hinton, G., Srivastava, N., Swersky, K.: Lecture 6a overview of mini-batch gradient descent. Coursera Lecture slides <https://class.coursera.org/neuralnets-2012-001/lecture>, [Online] (2012)
38. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.: Dropout: A simple way to prevent neural networks from overfitting. The Journal of Machine Learning Research **15** (2014) 1929–1958
39. Papineni, K., Roukos, S., Ward, T., Zhu, W.J.: Bleu: a method for automatic evaluation of machine translation. In: Proceedings of the 40th annual meeting on association for computational linguistics. (2002) 311–318