# MLP Coursework 4: Using Deep Learning Methodology For Image Captioning

G055 (s1984473, s1975313)

## Abstract

Image captioning is a challenging field of research that involves two domains: computer vision and natural language processing. Many approaches have been proposed over the last few years in order to achieve better performance when generating a description for an image. The aim of this project is to implement an image captioning model based on a combination of DenseNet and TPGN networks which has never been explored before. With only limited available resources, we have implemented the TPGN language model from scratch. Then, by using the standard MS COCO dataset, we have trained and evaluated ResNet-LSTM, DenseNet-LSTM, ResNet-TPGN, and DenseNet-TPGN models. Specifically, we have compared the new DenseNet-TPGN model against our ResNet-LSTM baseline model, to examine if the new model can perform equally well or even better on the image captioning task. Moreover, we have also validated the performance of DenseNet-based models against ResNet-based model in image captioning. A key finding of this project is that an attention mechanism in image captioning can outperform the sole usage of TPRs in TPGN.

## 1. Introduction

Image captioning has been a progressive area of research in recent years and finds its applications across a variety of domains. It is a process of generating textual descriptions of an image using both computer vision and natural language processing (NLP) techniques. The initial techniques revolved around traditional machine learning algorithms. However, in the last decade, the emergence of deep learning has significantly improved the image captioning process and has since been the state of the art technique for image captioning (Zakir Hossain et al., 2019).

Generating a well formed description for an image is challenging and requires a good understanding of the objects in the image and the relationship between those objects. The syntactic and semantic understanding of the language is also required to achieve a good caption for an image. Deep learning based techniques are very useful because they can automatically learn features from images and generate a caption for them. With increasing research interest in image captioning, a number of deep learning based techniques have been proposed.

The basic deep learning methodology in image captioning involves an encoder-decoder architecture. A Convolutional Neural Network (CNN) encoder is used to learn features in an image. Then, a Recurrent Neural Network (RNN) language model decoder is used to generate captions (Karpathy & Fei-Fei, 2017; Smolensky, 1990). The approach is popular but it lacks the ability to analyse image while generating the image descriptions. This gave rise to attention-based encoder-decoder models which consider the image spatial aspects while generating captions (Xu et al., 2015; Huang et al., 2019). A novel approach to image captioning involves an encoder-decoder model with an object detector and a slot filling mechanism (Lu et al., 2018). The main idea is to generate a sentence template with slot locations tied explicitly to specific image regions wherein the slots are filled with features of the image regions as identified by the object detector.

Image captioning, though an interesting research problem, posses some serious challenges in the image understanding domain. The extraordinary performance of deep learning based methods in general image classification and object detection tasks motivated us to explore the evolution of deep learning in image captioning which covers not only basic classification and detection but also language generation as a whole. Some of the most common encoder-decoder combinations explored before in the task of image captioning include VggNet-LSTM, ResNet-LSTM, DenseNet-LSTM, ResNet-TPGN, GoogleNet-LSTM (Staniute & Šešok, 2019). To the best of our knowledge (Staniute & Šešok, 2019; Zakir Hossain et al., 2019; Bai & An, 2018), nobody has ever tried the combination of DenseNet and Tensor Product Generation Network (TPGN) (Huang et al., 2018) for the image captioning task even though individually DenseNet (Huang et al., 2017) as an encoder and TPGN as a decoder perform better than ResNet (He et al., 2016) and Long Short-Term Memory (LSTM) (Hochreiter & Schmidhuber, 1997) respectively. This inspired us to explore the DenseNet-TPGN encoder-decoder combination and to investigate if it can give comparable or even better results in an image captioning task.

The rest of the report is structured as follows. Section 1 covers the high level aim, the research questions to be addressed and the objectives of our image captioning task. Section 2 briefly outlines the dataset used for training and evaluation of the models. We then discuss the technical methodologies of our task in Section 3. Section 4 provides details on the experiments performed followed by a critical analysis on the results obtained. In the end we discuss some related published work for certain future improvements.
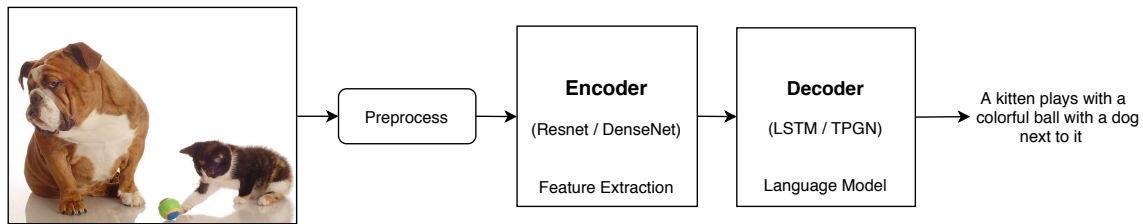
*Figure 1.* Overall pipeline of our system.

**Research Questions**

The aim of the project is to investigate the performance of a neural network architecture based on DenseNet-121 encoder and TPGN decoder combination on an image captioning task. We further aim to explore whether our DenseNet-TPGN model can give comparable or even better results on the image captioning task than the more commonly used ResNet-101 encoder and LSTM decoder combination, which serves as the baseline model for our task. We choose to study the DenseNet-TPGN encoder-decoder combination because it has never been explored before and would thus contribute to an interesting potential area of research. Based on this, the following research questions will be addressed:

1. Can TPGN perform equally well or even better as compared to our LSTM-based baseline model on MS COCO dataset?

2. Can we validate the claims of (Huang et al., 2017) that DenseNet gives a better performance than Resnet on our image captioning task?

3. How does our DenseNet-TPGN model perform as a whole on the MS COCO dataset as compared to our baseline model?

4. Which combination out of ResNet-LSTM, ResNet-TPGN, DenseNet-LSTM and DenseNet-TPGN outperforms the rest of the models, and why?

5. What are the limitations of our new model and what potential solutions can be considered for future investigation?

**Objectives**

The overall workflow of our project is to build a DenseNet-TPGN architecture, investigate its performance on the image captioning task, perform experiments to compare the different architectural combinations, and finally evaluate their performance on the MS COCO dataset using various metrics such as BLEU-4, METEOR, CIDEr and ROGUE_L scores. To specifically address the research questions, the objectives of our project are:

- Build a ResNet encoder model, and a DenseNet encoder model for image features extraction.

- Use GloVe representation as word embedding for caption preprocessing.

- Build an attention-based LSTM decoder model, and a TPGN decoder model for caption generation.

- Compute BLEU-4, METEOR, CIDEr and ROGUE_L scores for all the models (ResNet-LSTM, ResNet-TPGN, DenseNet-LSTM, DenseNet-TPGN) and compare their performance on the MS COCO dataset.

## 2. Dataset and Task

We aim to explore our models for the image captioning task on the Microsoft Common Objects in Context (MS COCO) dataset (Lin et al., 2015) which is a standard dataset for large scale object detection, segmentation and image captioning. The dataset consists of a total of 120k images with 5 captions for each image.

For the image captioning task, the dataset is split into 113289 training images, 5000 validation images and 5000 test images. All images of the dataset are resized to 64x64 and 256x256 for uniformity. We then normalise the images by the mean and standard deviation of the ImageNet (Deng et al., 2009) images' RGB channels such that pixel values are in the range [0,1]. Captions are both the target and the inputs of the decoder as each word is used to generate the next word. The captions in the dataset are annotated with a <start> sequence token to generate the first word and an <end> sequence token to know when to stop decoding during inference. The captions are padded with <pad> tags so that all captions have the same length. Since the captions are padded, we also keep track of the lengths of each caption. This is the actual length plus two (for the <start> and <end> tokens). We use GloVe embeddings for the captions to enable semantic parsing, and extract meaning from text.

The high level overview of our image captioning task is as follows. Images are fed to an encoder model in a fixed dimension to extract features from the image which are then passed to a RNN decoder model to generate a caption for the image on the basis of these features. Figure 1 illustrates the overall pipeline for our task.

When doing inference and evaluation, a beam search (Anderson et al., 2017) of size 5 will be used to select the caption sequence that has the highest overall score from all the potential candidate captions. Then, metrics such as BLEU-4 (Papineni et al., 2002), METEOR (Banerjee & Lavie,

2005), CIDEr (Vedantam et al., 2015) and ROUGE_L (Lin, 2004) are used to evaluate our image captioning models. All these metrics output a score indicating a similarity between the candidate sentence and the reference sentences.

## 3. Methodology

This section covers the technical methodologies which are used to achieve our image captioning task.

### 3.1. Feature Extraction (Encoder)

Usually a CNN is used to extract a high level feature representation of an image due to its extraordinary performance in general image classification tasks. Each layer in a deep convolutional network architecture except the last one, automatically extracts the low level and high level features in a hierarchical order. ResNet is proposed by (He et al., 2016) in 2016 to overcome the vanishing gradient problem of very deep networks. It uses shortcut connections to perform identity mapping so that the flow of gradients can be distributed from input to output layers. DenseNet (Huang et al., 2017) differs from ResNet in a way that each layer receives a collective knowledge from the previous layer. In other words, DenseNet introduces direct connections from each layer to all the subsequent layers. The improved parameter efficiency of DenseNet makes the network easier to train which improves the performance of the network in terms of training time as well. In our research project, ResNet-101 and DenseNet-121 networks will be used for image features extraction.

### 3.2. Word Embedding

GloVe 300-dimensional word embedding method can be used to create word vectors that capture meaning in vector space (Pennington et al., 2014). The embedding method takes advantage of global count statistics instead of only local information, by aggregating the global word-to-word co-occurrence matrix from a corpus. Word vectors are trained so that their differences can predict the co-occurrence ratios. In this project, pre-trained GloVe 300-dimensional word vectors are loaded into the embedding layer of our models.

### 3.3. Attention Mechanism

Attention mechanism is used to replicate natural human behaviour in image captioning models by identifying specific regions of an image and then forming a good explanation of the relationship of objects in those regions. (Xu et al., 2015) was the first to propose an attention-based model wherein the attention mechanism was integrated in a basic encoder-decoder framework to dynamically attend salient image regions during caption generation. In this project, the soft attention method, as proposed by (Xu et al., 2015), is used to consider the weighted average across all pixels in an image. As shown in figure 2, this weighted representation of the image can be concatenated with the previously generated word at each step to generate the next word.

### 3.4. Language Models (Decoder)

We work with RNNs which act as language models in our image captioning task. RNN is a general feedforward neural network that has some internal memory. It recursively uses its current input $\mathbf{x}_t$ and its previous output $\mathbf{h}_{t-1}$ to compute its current output $\mathbf{h}_t$. For our baseline model, we choose the standard LSTM layer with attention in order to generate captions. Then, we study and implement TPGN, an architecture that is capable of doing Tensor Product Representation (TPR) computation.

**Long Short-Term Memory**

LSTMs (Hochreiter & Schmidhuber, 1997) are a special form of RNNs with the ability to remember patterns selectively for long durations of time. A LSTM unit consists of three gates which regulate the flow of information: input gate, forget gate and output gate. These gates help to identify the important bits of data, thus passing only relevant information down the long chain of sequences to make predictions. To be more precise, the input gate controls when to update the cell state. The forget gate decides what information should be kept and what should be discarded. Lastly, the output gate decides what the next hidden state should be.

**Tensor Product Generation Network**

TPGN (Huang et al., 2018) is a deep learning network that can carry out TPR computation, and capture the grammatical structures of natural language. To understand TPGN, it is useful to firstly know the concept of TPR (Smolensky, 1990).

TPR is defined as a method to represent value/variable bindings in connectionist systems. In the theory of TPR, a sentence can be embedded as a sequence of grammatical components (roles), in which each role is bound with a word (filler). This is different when compared with the common bag-of-words embedding. TPR embedding tags each word in the filler embedding space $V_F$ with a role in the role embedding space $V_R$. For instance, `Tom likes Mary` will be treated as `Tom`/SUBJ, `likes`/VERB, and `Mary`/OBJ, where SUBJ, VERB, and OBJ mean subject, verb, and object respectively. Hence, the TPR that represents `Tom likes Mary` in the $V_F \otimes V_R$ embedding space is $\mathbf{v}_{\text{TomlikesMary}} = \mathbf{t} \otimes \mathbf{r}_{\text{SUBJ}} + \mathbf{l} \otimes \mathbf{r}_{\text{VERB}} + \mathbf{m} \otimes \mathbf{r}_{\text{OBJ}}$.

In fact, the tensor product $\mathbf{t} \otimes \mathbf{r}_{\text{SUBJ}}$ can simply be treated as an outer product $\mathbf{t}\mathbf{r}_{\text{SUBJ}}^\top$. This is the filler/role binding for `Tom`/SUBJ. The TPR for `Tom likes Mary` then becomes $\mathbf{v}_{\text{TomlikesMary}} = \mathbf{t}\mathbf{r}_{\text{SUBJ}}^\top + \mathbf{l}\mathbf{r}_{\text{VERB}}^\top + \mathbf{m}\mathbf{r}_{\text{OBJ}}^\top$. To uncover a word from a TPR, we can use inner product to unbind a role in the TPR and then extract the word that is associated with that role. When the role vectors $\mathbf{r}$ are orthogonal, the role vectors $\mathbf{r}$ themselves can be used as the unbinding vectors. For example, $\mathbf{v}_{\text{TomlikesMary}}\mathbf{r}_{\text{SUBJ}} = \mathbf{t}$. However, when the role vectors $\mathbf{r}$ are not orthogonal, there exist some unbinding vectors $\mathbf{u}$ that can be used for role unbinding, such as $\mathbf{v}_{\text{TomlikesMary}}\mathbf{u}_{\text{SUBJ}} = \mathbf{t}$.
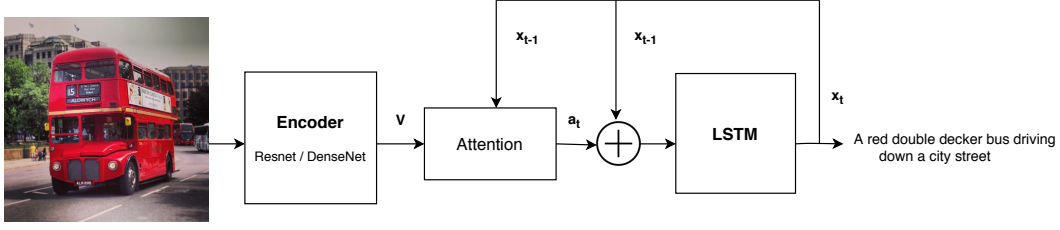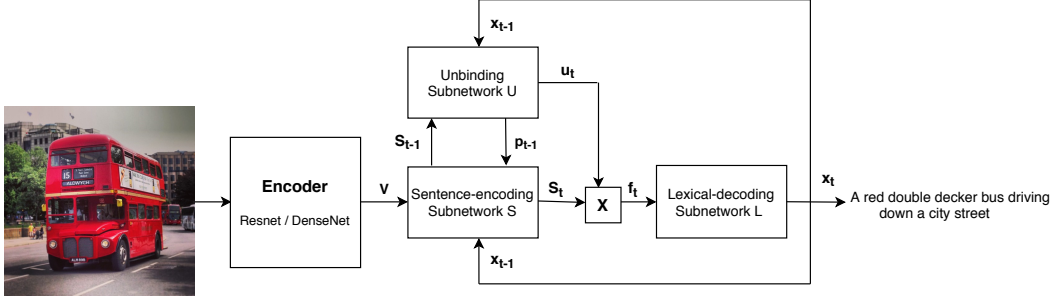
*Figure 2.* Architecture of LSTM decoder.



*Figure 3.* Architecture of TPGN decoder.

TPGN is a deep neural network that learns to represent TPRs and can be used for NLP tasks. As shown in figure 3, TPGN has three main components: a sentence-encoding subnetwork $S$, an unbinding subnetwork $U$, and a lexical-decoding subnetwork $L$. Subnetwork $S$ learns to generate an approximation to a TPR $S_t = \sum_{t=1}^{T} \mathbf{f}_t \mathbf{r}_t^\top$ that contains the filler/role binding for the subsequent caption word. Subnetwork $U$ learns to generate a role-unbinding vector $\mathbf{u}_t$. The filler word vector $\mathbf{f}_t$ is then extracted by multiplying the TPR $S_t$ with the role-unbinding vector $\mathbf{u}_t$. Lastly, the filler word vector $\mathbf{f}_t$ is fed into subnetwork $L$ to get the one-hot encoded caption word $\mathbf{x}_t$. This process is repeated until the whole caption sentence is generated.

In fact, both subnetworks $S$ and $U$ are one-layer, one-directional, and mutually-connected LSTMs. Subnetwork $S$ receives the internal hidden state from subnetwork $U$ when computing its input gate $\hat{\mathbf{i}}_{s,t}$, output gate $\hat{\mathbf{o}}_{s,t}$, forget gate $\hat{\mathbf{f}}_{s,t}$, and cell input $\hat{\mathbf{g}}_{s,t}$. Similarly, subnetwork $U$ receives the internal hidden state from subnetwork $S$ when computing its input gate $\hat{\mathbf{i}}_{u,t}$, output gate $\hat{\mathbf{o}}_{u,t}$, forget gate $\hat{\mathbf{f}}_{u,t}$, and cell input $\hat{\mathbf{g}}_{u,t}$.

The following equations are used by subnetwork $S$ to generate the TPR $S_t$:

$$\hat{\mathbf{f}}_{s,t} = \sigma(W_{s,f}\mathbf{p}_{t-1} - D_{s,f}W_e\mathbf{x}_{t-1} + U_{s,f}\hat{S}_{t-1}) \quad (1)$$

$$\hat{\mathbf{i}}_{s,t} = \sigma(W_{s,i}\mathbf{p}_{t-1} - D_{s,i}W_e\mathbf{x}_{t-1} + U_{s,i}\hat{S}_{t-1}) \quad (2)$$

$$\hat{\mathbf{o}}_{s,t} = \sigma(W_{s,o}\mathbf{p}_{t-1} - D_{s,o}W_e\mathbf{x}_{t-1} + U_{s,o}\hat{S}_{t-1}) \quad (3)$$

$$\hat{\mathbf{g}}_{s,t} = \tanh(W_{s,c}\mathbf{p}_{t-1} - D_{s,c}W_e\mathbf{x}_{t-1} + U_{s,c}\hat{S}_{t-1}) \quad (4)$$

$$\mathbf{c}_{s,t} = \hat{\mathbf{f}}_{s,t} \odot \mathbf{c}_{s,t-1} + \hat{\mathbf{i}}_{s,t} \odot \hat{\mathbf{g}}_{s,t} \quad (5)$$

$$\hat{S}_t = \hat{\mathbf{o}}_{s,t} \odot \tanh(\mathbf{c}_{s,t}) \quad (6)$$

$$S_t = \begin{bmatrix} \hat{S}_t & 0 & \cdots & 0 \\ 0 & \hat{S}_t & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \hat{S}_t \end{bmatrix} \quad (7)$$

where $\hat{\mathbf{f}}_{s,t}, \hat{\mathbf{i}}_{s,t}, \hat{\mathbf{o}}_{s,t}, \hat{\mathbf{g}}_{s,t}, \mathbf{c}_{s,t}, \hat{S}_t \in \mathbb{R}^{d \times d}$; $\mathbf{p}_t \in \mathbb{R}^d$; $W_{s,f}, W_{s,i}, W_{s,o}, W_{s,c}, D_{s,f}, D_{s,i}, D_{s,o}, D_{s,c} \in \mathbb{R}^{d \times d \times d}$; $U_{s,f}, U_{s,i}, U_{s,o}, U_{s,c} \in \mathbb{R}^{d \times d \times d \times d}$; $S_t \in \mathbb{R}^{d^2 \times d^2}$; $\mathbf{x}_t \in \mathbb{R}^v$ is the one-hot encoded caption word vector; $W_e \in \mathbb{R}^{d \times v}$ is the word embedding matrix; $\sigma(\cdot)$ is the logistic sigmoid function. In this paper, $d$ means the decoder output dimension, $e$ means the encoder output dimension, and $v$ means the caption vocabulary size. Besides, for clarity, all biases are not shown in the equations.

On the other hand, the following equations are used by subnetwork $U$ to generate the role unbinding vector $\mathbf{u}_t$:

$$\hat{\mathbf{f}}_{u,t} = \sigma(\hat{S}_{t-1}\mathbf{w}_{u,f} - D_{u,f}W_e\mathbf{x}_{t-1} + U_{u,f}\mathbf{p}_{t-1}) \quad (8)$$

$$\hat{\mathbf{i}}_{u,t} = \sigma(\hat{S}_{t-1}\mathbf{w}_{u,i} - D_{u,i}W_e\mathbf{x}_{t-1} + U_{u,i}\mathbf{p}_{t-1}) \quad (9)$$

$$\hat{\mathbf{o}}_{u,t} = \sigma(\hat{S}_{t-1}\mathbf{w}_{u,o} - D_{u,o}W_e\mathbf{x}_{t-1} + U_{u,o}\mathbf{p}_{t-1}) \quad (10)$$

$$\hat{\mathbf{g}}_{u,t} = \tanh(\hat{S}_{t-1}\mathbf{w}_{u,c} - D_{u,c}W_e\mathbf{x}_{t-1} + U_{u,c}\mathbf{p}_{t-1}) \quad (11)$$

$$\mathbf{c}_{u,t} = \hat{\mathbf{f}}_{u,t} \odot \mathbf{c}_{u,t-1} + \hat{\mathbf{i}}_{u,t} \odot \hat{\mathbf{g}}_{u,t} \quad (12)$$

$$\mathbf{p}_t = \hat{\mathbf{o}}_{u,t} \odot \tanh(\mathbf{c}_{u,t}) \quad (13)$$

$$\mathbf{u}_t = \tanh(W_u\mathbf{p}_t) \quad (14)$$

where $\hat{\mathbf{f}}_{u,t}, \hat{\mathbf{i}}_{u,t}, \hat{\mathbf{o}}_{u,t}, \hat{\mathbf{g}}_{u,t}, \mathbf{c}_{u,t}, \mathbf{w}_{u,f}, \mathbf{w}_{u,i}, \mathbf{w}_{u,o}, \mathbf{w}_{u,c} \in \mathbb{R}^d$;

$D_{u,f}, D_{u,i}, D_{u,o}, D_{u,c}, U_{u,f}, U_{u,i}, U_{u,o}, U_{u,c} \in \mathbb{R}^{d \times d}$; $W_u \in \mathbb{R}^{d^2 \times d}$; $\mathbf{u}_t \in \mathbb{R}^{d^2}$.

To initialise both subnetworks $S$ and $U$, the initial hidden state $\hat{S}_0$ is set to $\hat{S}_0 = C_s(\mathbf{v} - \bar{\mathbf{v}})$, where $\mathbf{v} \in \mathbb{R}^e$ is the encoded image features, $\bar{\mathbf{v}}$ is the mean of all the encoded image features, and $C_s \in \mathbb{R}^{d \times d \times e}$. Besides, the initial hidden state $\mathbf{p}_0 = \mathbf{0}$, which is a zero vector. Apart from that, the first one-hot encoded caption word $\mathbf{x}_0$ is initialised using the <start> tag. At each time step $t$, the filler word vector $\mathbf{f}_t \in \mathbb{R}^{d^2}$ can be extracted from TPR $S_t$ by using the role unbinding vector $\mathbf{u}_t$. Lastly, a softmax function is used by the lexical subnetwork $L$ to decode the one-hot encoded caption word vector $\mathbf{x}_t$ from the filler word vector $\mathbf{f}_t$. $W_x \in \mathbb{R}^{v \times d^2}$ is the word de-embedding matrix.

$$\mathbf{f}_t = S_t \mathbf{u}_t \tag{15}$$

$$\mathbf{x}_t = \text{softmax}(W_x \mathbf{f}_t) \tag{16}$$

## 4. Experiments

This section covers the experiments, the quantitative results obtained and the critical analysis of the results which validates the efficiency of our models for image caption generation. The experiments are performed to achieve the two main goals of this project: Firstly, to investigate the impacts of DenseNet-TPGN model on the image captioning task as compared to our baseline model; Secondly, to compare the different architectural combinations and explore which model performs the best. We discuss in detail the implementation settings and the training parameters used for model tuning for all of our models, followed by the results and analysis.

### 4.1. Implementation and Training

In this project, the PyTorch framework is used to implement the deep learning models. We start by training a baseline ResNet-LSTM model using 64x64 resized MS COCO images. Then, to investigate our research questions, we develop and train DenseNet-LSTM, ResNet-TPGN, and DenseNet-TPGN models on the both 64x64 and 256x256 resolution images.

**Baseline**

Our ResNet-LSTM baseline model, inspired by (Xu et al., 2015), is trained using 64x64 resized MS COCO images. A pre-trained ResNet-101 model from PyTorch is used. Since we do not need to do image classification, we strip the last two layers of the pre-trained ResNet-101, which consist of an average pooling layer and a fully-connected linear layer.

A soft attention mechanism with an attention dimension of 512 is implemented according to (Xu et al., 2015). Then, a pre-trained GloVe embedding layer (Pennington et al., 2014) with a dimension of 300 is loaded. For the LSTM decoder, we set the output dimension to 512. We initialise all biases to 0 and randomly initialise all the weights from a uniform distribution of -0.1 to 0.1.

Due to GPU memory limitation, a minitbatch size of 28 is used. To improve the performance, we fine-tune both the ResNet-101 encoder and the GloVe embedding layer. However, only convolutional blocks 2 to 4 in the pre-trained ResNet-101 are fine-tuned. This is because we think that the first convolutional block of the pre-trained model might contain some fundamental learning that are vital for image processing such as recognising lines and edges.

To calculate the loss for each word in a minibatch of data, the PyTorch cross entropy loss function is used. Besides, an additional loss component $\lambda \sum_i^L (1 - \sum_t^T \alpha_{i,t})^2$ that can act as a doubly stochastic regularization is used for the attention mechanism, as recommended by (Xu et al., 2015). Here, $\alpha$ means the attention weight. This is to make $\sum_t^T \alpha_{i,t} \approx 1$, encouraging the models to pay equal attention in all parts of an image over the course of generating the whole caption sentence. In addition, we do not compute losses over the padded regions, which are indicated by the <pad> tokens in the generated caption.

Adam optimization (Kingma & Ba, 2015) is used with the hyperparameter settings of: encoder learning rate = $1 \times 10^{-4}$, decoder learning rate = $4 \times 10^{-4}$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, and $\epsilon = 1 \times 10^{-8}$. A very small encoder learning rate is used here because we think that the pre-trained encoder model is already quite optimized and we do not want to update the model parameters too quickly. Moreover, a dropout (Srivastava et al., 2014) rate of 0.5 is used to regularise the training of the models. To avoid exploding gradient problem, we also clip the parameter gradients element-wise for each minibatch of data to be within the range of -5 to 5 before updating the parameters. To aid and speed up the decoder learning process, we have used a technique called teacher forcing (Williams & Zipser, 1989), in which the ground-truth caption word is supplied as input to the decoder, regardless of the previous generated word.

Lastly, as recommended by (Xu et al., 2015), the commonly reported BLEU-4 score calculated from the validation dataset is used to do early stopping and model selection. Within 20 epochs of training, our best baseline model candidate achieves a validation BLEU-4 score of 18.75. In fact, due to teacher forcing, this BLEU-4 score might not reflect the real performance of the baseline model. However, we think that it is still reasonable to be used as an indicator of the best validation model because teacher forcing has been consistently applied in every epoch.

**DenseNet-LSTM**

When training our DenseNet-LSTM model, a pre-trained DenseNet-121 model from PyTorch is used. Again, we strip the last fully-connected linear layer of the pre-trained DenseNet-121 model. Similarly to our baseline model, an attention dimension of 512, a Glove embedding dimension of 300, and a LSTM decoder output dimension of 512 are used. During training, we fine-tune both the pre-trained DenseNet-121 model and the GloVe embedding layer, but only dense blocks 2 to 4 in the pre-trained DenseNet-121 model are fine-tuned. The same baseline training hyper-

parameters, optimization algorithm, and model selection approach are used here for the same reasons mentioned in baseline model training.

### ResNet-TPGN and DenseNet-TPGN

ResNet-TPGN and DenseNet-TPGN are the next two models that are developed and trained. By using the PyTorch `torch.nn.LSTMCell()` source codes, we implement the TPGN decoder from scratch. To train our TPGN decoder, we have used the TPGN end-to-end training algorithm from (Huang et al., 2018):

---

**Algorithm 1** TPGN Training (Huang et al., 2018)

---

**input** Image feature vector $\mathbf{v}^i$ and its corresponding caption $X^i = [\mathbf{x}_1^i, \cdots, \mathbf{x}_T^i]$ for $i = 1, \cdots, N$. Here $N$ is the total number of samples.

1: Initialise $\hat{S}_0$ by $C_s(\mathbf{v} - \bar{\mathbf{v}})$;
2: Initialise $\mathbf{x}_0$ as the one-hot encoded caption vector corresponding to the <start> tag.
3: Initialise $\mathbf{p}_0$ as the zero vector.
4: Randomly initialise $C_s$ and all weights and biases in (1) – (16);
5: **for** $n$ from 1 to $N$ **do**
6:   **for** $t$ from 1 to $T$ **do**
7:     Calculate (1) – (6) to obtain $\hat{S}_t$;
8:     Calculate (7) to obtain $S_t$;
9:     Calculate (8) – (13) to obtain $\mathbf{p}_t$;
10:     Calculate (14) to obtain $\mathbf{u}_t$;
11:     Calculate (15) to obtain $\mathbf{f}_t$;
12:     Calculate (16) to obtain $\mathbf{x}_t$;
13:     Use back-propagation to update $C_s$ and all weights and biases in (1) – (16).
14:   **end for**
15: **end for**
**output**

---

Again, a GloVe embedding layer with dimension 300 is used. By referring to (Huang et al., 2018), we use the same dimension $d = 25$ for our TPGN decoder. In this case, the output filler word vector $\mathbf{f}_t$ has a dimension of $d^2 = 625$. Due to the TPGN decoder architecture, attention mechanism is not used here. Hence, the additional loss component for attention mechanism is not used. Lastly, similar to the training of our baseline and DenseNet-LSTM models, we use the same fine-tuning strategy, training hyperparameters, optimization algorithm, and model selection approach.

### Training with Higher Resolution Images

To further improve the models' performance, we also train our ResNet-LSTM, DenseNet-LSTM, ResNet-TPGN, and DenseNet-TPGN models on the higher resolution 256x256 resized MS COCO images. All the training strategy and hyperparameters stay the same, except a smaller minibatch size of 20 is used due to the limitation of GPU memory.

### 4.2. Results and Analysis

We have trained our models on both 64x64 and 256x256 resolution images. The training accuracy, training loss, and validation BLEU-4 score during the training of all the models are plotted and shown in figure 4.

Firstly, it is easy to notice that TPGN models always have higher training accuracy and lower training loss than LSTM models. These significant gaps of training performance can be explained by the way we calculate the training losses. For the attention-based LSTM models, an additional loss that can act as a doubly stochastic regularization is used. This leads to an increase in both the training and validation losses, and subsequently a decrease in both the training and validation accuracies.

To select the best validation model for each encoder-decoder combinations, validation BLEU-4 score is used. As shown in figure 4(c) and 4(f), all the models converge pretty quickly, some within 10 epochs, and others within 20 epochs. We think that this is because of the usage of teacher forcing, in which it makes the models to learn faster. In general, most TPGN models have lower validation BLEU-4 scores than LSTM models. It seems that lower training losses do not help TPGN models to achieve a better generalization performance due to overfitting.

Next, all the best validation models are evaluated using the test dataset. The test scores are shown in table 1. It is observed that for models that are trained with 64x64 resolution images, DenseNet-LSTM model performs the best. Similarly, for models that are trained with 256x256 resolution images, DenseNet-LSTM model also performs the best. This might be explained by the attention mechanism that is used in LSTM models. In fact, (Huang et al., 2018) only compare their TPGN image captioning model with non-attention-based LSTM models. In our experiments, it is observed that the usage of TPRs to encode extra grammatical information does not supersede LSTM models when an attention mechanism is present.

When comparing between ResNet and DenseNet models, DenseNet-LSTM combination always beats ResNet-LSTM combination in all test scores. Similarly, DenseNet-TPGN combination also beats ResNet-TPGN combination in all test scores most of the time. The only exception is when being trained with 256x256 resolution images, ResNet-TPGN model has a slightly better BLEU-4 and CIDEr test scores than DenseNet-TPGN model. Hence, we are able to validate the claims of densenet paper as DenseNet performs better than ResNet in most cases.

Although DenseNet does not show a significant advantage to ResNet when being trained on the higher resolution images, it has improved parameter efficiency which makes it much easier and faster to train. In fact, this is reflected in the training time required per epoch, as shown in Table 1. Regarding training time, all TPGN models spend a substantial 7.5 to 10.5 hours per epoch. We think that this inefficiency in training might be explained by the imple-
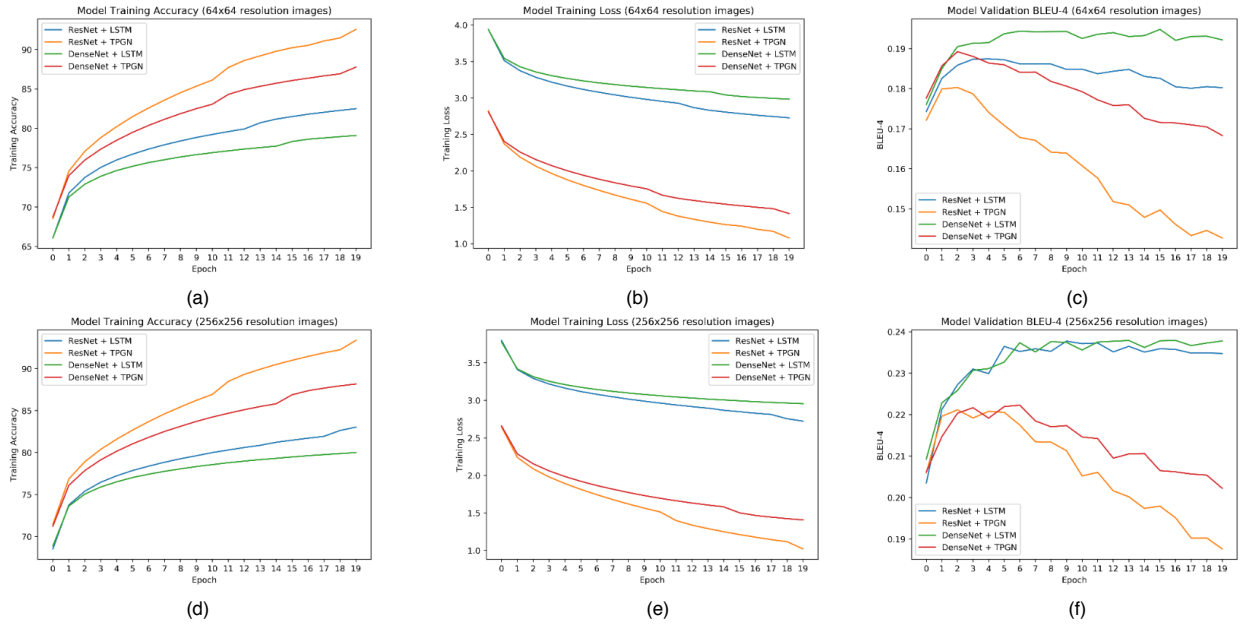
Figure 4. models training for 20 epochs with 64x64 resolution images: (a) Training accuracy, (b) Training loss, (c) Validation BLEU-4 for early stopping and best validation model selection; models training for 20 epochs with 256x256 resolution images: (d) Training accuracy, (e) Training loss, (f) Validation BLEU-4 for early stopping and best validation model selection.

| Model | Resolution | Train Time (Hour/Epoch) | BLEU-4 | METEOR | ROUGE_L | CIDEr |
|---|---|---|---|---|---|---|
| ResNet-LSTM (Baseline) | 64x64 | 3.5 | 23.2 | 20.3 | 46.3 | 65.8 |
| DenseNet-LSTM | 64x64 | 2.0 | **24.4** | **20.9** | **47.2** | **69.1** |
| ResNet-TPGN | 64x64 | 10.0 | 20.9 | 19.1 | 44.6 | 57.1 |
| DenseNet-TPGN | 64x64 | 7.5 | 22.5 | 20.1 | 45.6 | 63.4 |
| ResNet-LSTM | 256x256 | 5.5 | 32.4 | 25.7 | 53.8 | 100.8 |
| DenseNet-LSTM | 256x256 | 3.5 | **33.1** | **25.8** | **54.1** | **101.7** |
| ResNet-TPGN | 256x256 | 10.5 | 29.9 | 24.4 | 51.9 | 92.2 |
| DenseNet-TPGN | 256x256 | 10.5 | 29.7 | 24.6 | 51.9 | 92.1 |

Table 1. Test scores achieved by all models.

mentation of our codes. In fact, more advanced PyTorch features such as TorchScript can be used to optimise our TPGN model for GPU training. We did not realise this when we implemented our TPGN model, and this can be a future improvement to our research project.

In order to present our results more visually and to depict our models' qualitative performance, a sample image from the test dataset is used. Figure 5 illustrates the sample image with its five ground-truth captions and the captions generated by the different models. For this particular sample image, while most of the models are able to understand the scene and generate similar captions, TPGN specifically is able to pick up phrases like 'red jacket', 'down a hill' and 'snowy surface'.

Some very interesting results are obtained from the experiments performed while addressing the main goals of the project. These results validate the effectiveness of our models and better results can be achieved with future improvements as discussed in the next section.

## 5. Related Work

Image captioning though a recently emerged field of research, has been gaining popularity since the last few years. In the last 5 years, a large number of papers have been published on image captioning, each focusing on new and improved techniques. In this paper, we have covered two popular approaches of image captioning: basic encoder-decoder framework for all our models and additional attention mechanism incorporated in our baseline ResNet-LSTM model. There are several other techniques which can be adopted in our project and which would further help in improving the results obtained so far.

### 5.1. Attentive Tensor Product Learning

(Huang et al., 2019) introduces a novel architecture called Attentive Tensor Product Learning (ATPL) which exploits TPR and attention mechanism to learn a vector representation that can capture the grammatical structures of the

| Original Captions |
| --- |
| A woman posing for the camera standing on skies |
| a woman standing on skiis while posing for the camera |
| A woman in a red jacket skiing down a slope |
| A young woman is skiing down the mountain slope |
| a person on skis makes her way through the snow |

| Model | Generated Captions |
| --- | --- |
| ResNet-LSTM with GloVe, 64x64 | \<start\> a woman riding skis down a snow covered slope \<end\> |
| Resnet-TPGN with GloVe, 64x64 | \<start\> a woman standing on skis in the snow \<end\> |
| DenseNet-LSTM with Glove, 64x64 | \<start\> a woman riding skis down a snow covered slope \<end\> |
| DenseNet-TPGN with GLove, 64x64 | \<start\> a person riding skis on a snowy surface \<end\> |
| ResNet-LSTM with GloVe, 256x256 | \<start\> a woman riding skis down a snow covered slope \<end\> |
| Resnet-TPGN with GloVe, 256x256 | \<start\> a woman riding skis down a snow covered slope \<end\> |
| DenseNet-LSTM with Glove, 256x256 | \<start\> a woman riding skis down a snow covered slope \<end\> |
| DenseNet-TPGN with GLove, 256x256 | \<start\> a woman in a red jacket skiing down a hill \<end\> |

*Figure 5.* An example image with its ground truth captions and the captions that are generated by the models.

natural language descriptions. An attention module over an input vector $\mathbf{v}$ in an ATPL architecture is defined as $Attn(\mathbf{v}) = \sigma(W\mathbf{v} + \mathbf{b})$, where $\sigma$ is the sigmoid function and $W$ and $\mathbf{b}$ are two sets of parameters. ATPL employs two attention modules to place weights on each dimension of the vector $\mathbf{v}$, so that it can be used to compute the TPR of the sentence $S$ and the unbinding vector $\mathbf{u}$. This approach would be a significant improvement over our general TPR-based TPGN model due to its use of attention modules to compute TPR.

### 5.2. Image Captioning Using Object Detector and Slot Filling Mechanism

(Lu et al., 2018) introduces a new approach to image captioning called Neural Baby Talk, which is based on an encoder-decoder framework with an object detector and a slot filling mechanism. The main idea is to generate a sentence template with slot locations tied explicitly to specific image regions wherein the slots are filled with features of the image regions as identified by the object detector. This helps to generate natural language captions explicitly grounded in entities identified by the object detectors. Moreover, the captions generated using this technique are syntactically and semantically more accurate and much more related to the image than captions generated using trivial attentive encoder-decoder frameworks. For a given image $I$ and its corresponding caption $\mathbf{c}$, the candidate grounding regions are obtained through a Faster-RCNN network (Ren et al., 2015) and the caption template is generated using a RNN. A pointer network is used to generate the slot for visual words and an attention-based LSTM layer which takes convolutional maps as input acts as the decoder for caption generation in the Neural Baby Talk model.

## 6. Conclusions

With the aim of investigating the performance of DenseNet-TPGN encoder-decoder model for image captioning, we implemented TPGN decoder from scratch by only referring to (Huang et al., 2018) because there is no complete reference codes available online.

According to our experiment results, we find that the DenseNet-LSTM performs the best among all the encoder-decoder models (Research question 4). Although our DenseNet-TPGN model that is trained with 256x256 resolution images performs better than our ResNet-LSTM baseline model, its test scores are lower than the ResNet-LSTM model that is also trained with 256x256 resolution images. In fact, all TPGN-based models lose in comparison to its corresponding LSTM-based model. This shows that while TPGN is better than non-attention-based LSTM (Huang et al., 2018), the usage of TPRs to encode extra grammatical information does not supersede LSTM-based models when an attention mechanism is present (Research questions 1, 3). Moreover, we validated the claims of (Huang et al., 2017) in which for most cases, DenseNet-based models perform better than ResNet-based models. The training process also shows that DenseNet models are easier to train and require a shorter training time (Research question 2).

In future, we would like to further investigate the performance of an attention-based DenseNet-TPGN encoder-decoder model. In addition, the usage of object detector and slot filling mechanism (Lu et al., 2018) together with TPGN also holds an interesting research direction. Lastly, we wish to improve and optimise the implementation of our TPGN codes by using more advanced PyTorch development practice such as TorchScript.

# References

Anderson, Peter, Fernando, Basura, Johnson, Mark, and Gould, Stephen. Guided open vocabulary image captioning with constrained beam search. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pp. 936–945, Copenhagen, Denmark, September 2017. Association for Computational Linguistics. doi: 10.18653/v1/D17-1098. URL https://www.aclweb.org/anthology/D17-1098.

Bai, Shuang and An, Shan. A survey on automatic image caption generation. *Neurocomputing*, 2018. ISSN 18728286. doi: 10.1016/j.neucom.2018.05.080.

Banerjee, Satanjeev and Lavie, Alon. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pp. 65–72, Ann Arbor, Michigan, June 2005. Association for Computational Linguistics. URL https://www.aclweb.org/anthology/W05-0909.

Deng, J., Dong, W., Socher, R., Li, L., Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 248–255, 2009.

He, Kaiming, Zhang, Xiangyu, Ren, Shaoqing, and Sun, Jian. Deep residual learning for image recognition. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, volume 2016-December, pp. 770–778, 2016. ISBN 9781467388504. doi: 10.1109/CVPR.2016.90.

Hochreiter, Sepp and Schmidhuber, Jürgen. Long Short-Term Memory. *Neural Computation*, 1997. ISSN 08997667. doi: 10.1162/neco.1997.9.8.1735.

Huang, Gao, Liu, Zhuang, Van Der Maaten, Laurens, and Weinberger, Kilian Q. Densely connected convolutional networks. In *Proceedings - 30th IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017*, volume 2017-January, pp. 2261–2269, 2017. ISBN 9781538604571. doi: 10.1109/CVPR.2017.243.

Huang, Qiuyuan, Smolensky, Paul, He, Xiaodong, Deng, Li, and Wu, Dapeng. Tensor product generation networks for deep NLP modeling. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pp. 1263–1273, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-1114. URL https://www.aclweb.org/anthology/N18-1114.

Huang, Qiuyuan, Deng, Li, Wu, Dapeng Oliver, Liu, Chang, and He, Xiaodong. Attentive tensor product learning. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019,*

*The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019*, pp. 1344–1351. AAAI Press, 2019. doi: 10.1609/aaai.v33i01.33011344. URL https://doi.org/10.1609/aaai.v33i01.33011344.

Karpathy, Andrej and Fei-Fei, Li. Deep Visual-Semantic Alignments for Generating Image Descriptions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017. ISSN 01628828. doi: 10.1109/TPAMI. 2016.2598339.

Kingma, Diederik P. and Ba, Jimmy Lei. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015 - Conference Track Proceedings*, 2015.

Lin, Chin-Yew. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pp. 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics. URL https://www.aclweb.org/anthology/W04-1013.

Lin, Tsung-Yi, Maire, Michael, Belongie, Serge, Bourdev, Lubomir, Girshick, Ross, Hays, James, Perona, Pietro, Ramanan, Deva, Zitnick, C. Lawrence, and Dollár, Piotr. Microsoft COCO: Common Objects in Context. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2015. ISSN 10636919. doi: 10.1109/CVPR.2014.471.

Lu, Jiasen, Yang, Jianwei, Batra, Dhruv, and Parikh, Devi. Neural Baby Talk. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2018. ISBN 9781538664209. doi: 10.1109/CVPR.2018.00754.

Papineni, Kishore, Roukos, Salim, Ward, Todd, and Zhu, Wei-Jing. Bleu: A method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*, ACL '02, pp. 311–318, USA, 2002. Association for Computational Linguistics. doi: 10.3115/1073083.1073135. URL https://doi.org/10.3115/1073083.1073135.

Pennington, Jeffrey, Socher, Richard, and Manning, Christopher D. GloVe: Global vectors for word representation. In *EMNLP 2014 - 2014 Conference on Empirical Methods in Natural Language Processing, Proceedings of the Conference*, 2014. ISBN 9781937284961. doi: 10.3115/v1/d14-1162.

Ren, Shaoqing, He, Kaiming, Girshick, Ross, and Sun, Jian. Faster r-cnn: Towards real-time object detection with region proposal networks, 2015.

Smolensky, Paul. Tensor product variable binding and the representation of symbolic structures in connectionist systems. *Artificial Intelligence*, 1990. ISSN 00043702. doi: 10.1016/0004-3702(90)90007-M.

Srivastava, Nitish, Hinton, Geoffrey, Krizhevsky, Alex, Sutskever, Ilya, and Salakhutdinov, Ruslan. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 2014. ISSN 15337928.

Staniute, Raimonda and Šešok, Dmitrij. A systematic literature review on image captioning, 2019. ISSN 20763417.

Vedantam, R., Zitnick, C. L., and Parikh, D. Cider: Consensus-based image description evaluation. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4566–4575, 2015.

Williams, Ronald J. and Zipser, David. A Learning Algorithm for Continually Running Fully Recurrent Neural Networks. *Neural Computation*, 1989. ISSN 0899-7667. doi: 10.1162/neco.1989.1.2.270.

Xu, Kelvin, Ba, Jimmy Lei, Kiros, Ryan, Cho, Kyunghyun, Courville, Aaron, Salakhutdinov, Ruslan, Zemel, Richard S., and Bengio, Yoshua. Show, attend and tell: Neural image caption generation with visual attention. In *32nd International Conference on Machine Learning, ICML 2015*, 2015. ISBN 9781510810587.

Zakir Hossain, M. D., Sohel, Ferdous, Shiratuddin, Mohd Fairuz, and Laga, Hamid. A comprehensive survey of deep learning for image captioning, 2019. ISSN 15577341.