ANALYSIS OF DATASET, OBJECT TAG, AND OBJECT ATTRIBUTE
COMPONENTS IN NOVEL OBJECT CAPTIONING

A THESIS SUBMITTED TO
THE GRADUATE SCHOOL OF NATURAL AND APPLIED SCIENCES
OF
MIDDLE EAST TECHNICAL UNIVERSITY

BY

ENES MUVAHHİD ŞAHİN

IN PARTIAL FULFILLMENT OF THE REQUIREMENTS
FOR
THE DEGREE OF MASTER OF SCIENCE
IN
ELECTRICAL AND ELECTRONICS ENGINEERING

JULY 2022

Approval of the thesis:

**ANALYSIS OF DATASET, OBJECT TAG, AND OBJECT ATTRIBUTE COMPONENTS IN NOVEL OBJECT CAPTIONING**

submitted by **ENES MUVAHHİD ŞAHİN** in partial fulfillment of the requirements for the degree of **Master of Science in Electrical and Electronics Engineering Department, Middle East Technical University** by,

Prof. Dr. Halil Kalıpçılar
Dean, Graduate School of **Natural and Applied Sciences**       ⎯⎯⎯⎯⎯⎯

Prof. Dr. İlkay Ulusoy
Head of Department, **Electrical and Electronics Engineering**       ⎯⎯⎯⎯⎯⎯

Prof. Dr. Gözde Bozdağı Akar
Supervisor, **Electrical and Electronics Engineering, METU**       ⎯⎯⎯⎯⎯⎯

**Examining Committee Members:**

Prof. Dr. İlkay Ulusoy
Electrical and Electronics Engineering, METU       ⎯⎯⎯⎯⎯⎯

Prof. Dr. Gözde Bozdağı Akar
Electrical and Electronics Engineering, METU       ⎯⎯⎯⎯⎯⎯

Prof. Dr. Ece Güran Schmidt
Electrical and Electronics Engineering, METU       ⎯⎯⎯⎯⎯⎯

Assist. Prof. Dr. Ramazan Gökberk Cinbiş
Computer Engineering, METU       ⎯⎯⎯⎯⎯⎯

Assist. Prof. Dr. Osman Serdar Gedik
Computer Engineering, Ankara Yıldırım Beyazıt University       ⎯⎯⎯⎯⎯⎯

Date: 01.07.2022

**I hereby declare that all information in this document has been obtained and presented in accordance with academic rules and ethical conduct. I also declare that, as required by these rules and conduct, I have fully cited and referenced all material and results that are not original to this work.**

Name, Surname:    Enes Muvahhid Şahin

Signature        :

# ABSTRACT

## ANALYSIS OF DATASET, OBJECT TAG, AND OBJECT ATTRIBUTE COMPONENTS IN NOVEL OBJECT CAPTIONING

Şahin, Enes Muvahhid

M.S., Department of Electrical and Electronics Engineering

Supervisor: Prof. Dr. Gözde Bozdağı Akar

July 2022, 119 pages

Image captioning is a popular yet challenging task which lies at the intersection of Computer Vision and Natural Language Processing. A specific branch of image captioning called *Novel Object Captioning* draw attention in recent years. Different from general image captioning, Novel Object Captioning focuses on describing images with novel objects which are not seen during training. Recently, numerous image captioning approaches are proposed in order to increase quality of the generated captions for both general image captioning and Novel Object Captioning. These methods benefit from large object detection datasets for Novel Object Captioning. They also utilize specific set of object tags (class names) in the image. Even though these approaches are very successful in many aspects, they require GPU-weeks of training on several large datasets. Furthermore, captions generated by these methods may lack visual grounding and overlook details in the image. Thus, in this thesis, we analyze the dataset, object tag, and object attribute components for Novel Object Captioning. We perform Visual Vocabulary Pretraining (VIVO) [1] on small-scale [2] and large-scale [3] datasets and compare the captioning performances of a state-of-the-art method [4] in order to analyze the effect of dataset size. To analyze the effect of tag

quality on Novel Object Captioning performance, we compare the performance of captioning methods [4] trained with two different set of object tags: a large set of tags but lacking novel objects, a small set of tags with novel objects. Finally, to obtain richer captions and alleviate overlooked details in the image, we propose a novel approach in which object attributes in the image are exploited. Experimental results are demonstrated on both Novel Object Captioning and general image captioning tasks. The results show that novel object tags play a vital role for Novel Object Captioning and proposed method generates richer and more detailed captions compared to the baseline.

# ÖZ

## ÖZGÜN NESNE ALTYAZILAMA'DA VERİ KÜMESİ, NESNE ETİKETİ VE NESNE SIFATI BİLEŞENLERİNİN ANALİZİ

Şahin, Enes Muvahhid

Yüksek Lisans, Elektrik ve Elektronik Mühendisliği Bölümü

Tez Yöneticisi: Prof. Dr. Gözde Bozdağı Akar

Temmuz 2022 , 119 sayfa

İmge altyazılama Bilgisayarlı Görü ve Doğal Dil İşleme alanlarının kesişiminde yer alan hem popüler hem de zorlayıcı bir iştir. İmge altyazılamanın özel bir alt kolu olan *Özgün Nesne Altyazılama* son yıllarda ilgi görmektedir. Özgün Nesne Altyazılama eğitim sırasında görülmemiş özgün nesneler içeren imgeler için altyazı üretmeye odaklanmaktadır. Son yıllarda hem genel imge altyazılama hem de Özgün Nesne Altyazılama için üretilen altyazıların kalitesini arttırmak amacıyla çok sayıda yaklaşım önerilmiştir. Bu yaklaşımlar Özgün Nesne Altyazılama için büyük çaplı nesne tespiti veri kümelerinden faydalanmaktadır. Ayrıca, bu yöntemler imgedeki nesne etiketlerini (sınıf isimleri) kullanmaktadırlar. Bu yaklaşımlar birçok açıdan oldukça başarılı olsa da birkaç büyük veri kümesinde, haftalarca Grafik İşleme Birimi (GPU) üzerinde eğitilmektedir. Üstelik bu yöntemler tarafından üretilen altyazılar görsel temellendirme açısından zayıf kalabilmekte ve imgedeki detayları gözden kaçırabilmektedir. Bu nedenle, bu tezde, Özgün Nesne Altyazılama için veri kümesi, nesne etiketi ve nesne sıfatı bileşenlerinin analizi gerçekleştirilmiştir. Veri kümesi boyutunun imge altyazılama performansı üzerindeki etkisinin analizi için küçük çaplı [2] ve büyük

çaplı [3] veri kümelerinde Görsel Kelime Hazinesi Ön Eğitimleri (Visual Vocabulary Pretraining) [1] yapılarak en gelişmiş imge altyazılama yönteminin [4] imge altyazılama performansı karşılaştırılmıştır. Nesne etiketinin kalitesinin Özgün Nesne Altyazılama performansına etkisinin analizi için iki farklı nesne etiketi kümesi kullanılarak eğitilmiş yöntemler karşılaştırılmıştır: özgün nesne içermeyen büyük nesne etiketi kümesi, özgün nesne içeren küçük nesne etiketi kümesi. Son olarak, daha zengin altyazılama elde etmek ve imgedeki gözden kaçan detayları azaltmak amacıyla nesne sıfatlarından faydalanan özgün bir yaklaşım önerilmiştir. Deneysel sonuçlar hem Özgün Nesne Altyazılama hem de genel imge altyazılama görevlerinde gösterilmiştir. Deneysel sonuçlar özgün nesne etiketlerinin Özgün Nesne Altyazılama'da kritik bir rol oynadığını ve önerilen yaklaşımın baz alınan yaklaşıma göre daha zengin ve detaylı altyazılar ürettiğini ortaya koymuştur.

Anahtar Kelimeler: imge altyazılama, özgün nesne altyazılama, görü ve dil ön eğitimi, nesne etiketleri, nesne sıfatları

To my family

# ACKNOWLEDGMENTS

I would like to express my deep and sincere gratitude to my supervisor Prof. Dr. Gözde Bozdağı Akar for her guidance, assistance, support, and patience throughout the thesis studies.

I would like to thank my family for supporting and believing in me throughout my life.

# TABLE OF CONTENTS

# LIST OF TABLES

TABLES

# LIST OF FIGURES

FIGURES

xvi

xvii

# LIST OF ABBREVIATIONS

| | |
|---|---|
| CV | Computer Vision |
| NLP | Natural Language Processing |
| VL | Vision and Language |
| VLP | Vision and Language Pretraining |
| CNN | Convolutional Neural Network |
| RNN | Recurrent Neural Network |
| LSTM | Long Short-Term Memory |
| BUTD | Bottom-Up Top-Down |
| BERT | Bidirectional Encoder Representations from Transformers |
| BLEU | Bilingual Evaluation Understudy |
| ROUGE | Recall-Oriented Understudy for Gisting Evaluation |
| METEOR | Metric for Evaluation of Translation with Explicit ORdering |
| CIDEr | Consensus-based Image Description Evaluation |
| SPICE | Semantic Propositional Image Caption Evaluation |
| MLM | Masked Language Modelling |
| MTL | Masked Token Loss |
| CL3 | 3-Way Contrastive Loss |
| VIVO | Visual Vocabulary Pretraining |

# CHAPTER 1

# INTRODUCTION

Vision is one of the most powerful senses of human beings. It allows perceiving and discovering the world around us. On the other hand, it is also crucial to express what we perceive to the other people. Language is a way of communicating with others and a tool to describe what we perceive and see. Recently, there have been significant improvements in *Computer Vision (CV)* and *Natural Language Processing (NLP)* which allow machines to gain similar capabilities as humans such as perceiving images and videos, understanding and generating sentences and paragraphs.

Thanks to recent breakthrough in CV, machines can classify [5], [6], detect [7], [8], segment [9], [10] objects in images. Recent progress in NLP also allows machines to translate sentences from one language to the other [11], [12], perform sentiment classification [13], answer questions [13], [14].

Another research area of interest, called *Vision and Language (VL)*, has drawn attention in recent years [15]. VL lies at the intersection of CV and NLP. It borrows ideas from both of these domains. Hence, advancements in these two domains directly contribute to the VL methods. There are many sub-tasks in VL such as *image captioning, dense image captioning, image paragraph generation, video captioning, visual question answering, image-text retrieval, visual reasoning, vision and language navigation* [15], [16]. In this thesis, among these sub-tasks, we will concentrate on *image captioning*.

Image captioning models aim to generate a sentence which describes the given image. Image captioning can be utilized in many areas such as guiding and helping visually-impaired people [17], caption generation for news articles [18], generating biomedical

1

reports from medical images (X-RAY, CT) [19], or making robots navigate in an environment via vision and language [20].

An image captioning model takes an image as the input and produces a description for the input image as the output. An example image-caption pair is given in Figure 1.1. A good caption should be grammatically correct, natural sounding, rich, and

A couple of people standing in the snow flying a kite

Figure 1.1: An example image and its corresponding caption [4].

grounded on the image [21], [22], [23]. The properties of a good caption guide and construct the areas of interest for image captioning research.

Recently, several research on image captioning focus on generating captions for images which contain objects that are not present in the image captioning training dataset, [24], [1], [25], [26]. The task of describing images with unseen objects during training is called *Novel Object Captioning*.

Another group of approaches exploit datasets with large amount of image-caption

pairs ([27], [28], [29], [2]) in order to generate rich and high-quality captions, [30], [16], [4].

Yet another branch of research targets visually-grounded image captioning. There are studies demonstrating that image captioning models are inclined towards copying phrases from training dataset without paying attention to the input image, [31], [32]. Furthermore, these models might hallucinate non-existing objects in the image or overlook important details [31], [22]. An example of such a caption generated by a state-of-the-art model, VinVL [4], is given in Figure 1.2 where the generated caption hallucinates a *chair*, overlook important details such as the *fence* of the garden and the *car* in the background. Some methods try to overcome these issues and generate detailed and visually-grounded captions, [31], [32].



A garden with a chair and a plant on the ground.

Figure 1.2: An example image and a poor caption generated by VinVL [4].

Datasets and object detection components play crucial roles in the success of image captioning models [21]. State-of-the-art methods [1], [4], require large pretraining datasets and specific set of object tags (classes) for Novel Object Captioning.

In this study, we question the necessity of large datasets and specific object tags for Novel Object Captioning. We analyze the effect of pretraining datasets and object tags in Novel Object Captioning for state-of-the-art methods [1], [4].

As the main contribution in this study, we propose a novel approach for enriching image captioning results with object attributes. The proposed approach aims to generate richer and more detailed captions and improve the poor captions like the one in Figure 1.2.

Experimental results are demonstrated on both Novel Object Captioning and general image captioning tasks.

## 1.1 The Outline of the Thesis

The thesis contains 6 chapters. In Chapter 1, brief descriptions of the image captioning and Novel Object Captioning are provided. Problems with the existing methods are stated and our contributions are summarized. In Chapter 2, background information for the concepts utilized and mentioned throughout the thesis is provided. In Chapter 3, literature review for the image captioning is given. Applications areas and taxonomy of image captioning models are discussed. Datasets and evaluation metrics for image captioning are explained. In Chapter 4, the architectural and training details of the baseline method employed in this study are explained in detail. The proposed novel approach is also explained in this chapter. In Chapter 5, experimental results are provided. Analysis of dataset, object tag, and object attribute components in novel object captioning is performed. The baseline method and the proposed method are compared theoretically and practically. Finally, in Chapter 6, we summarize the work conducted in this thesis, make some conclusions and discuss possible future works as an extension to this study. Extra visual examples for the comparison of the baseline and the proposed method are provided in Appendix A as an addition to the ones in Chapter 5.

## CHAPTER 2

## BACKGROUND INFORMATION

In this chapter, background information for the core topics which are utilized and mentioned throughout the thesis is provided.

## 2.1 Image Classification and Object Detection

*Image Classification* and *Object Detection* are two hot topics of the CV. Models developed for these topics are also utilized by the baseline and the proposed image captioning methods in the thesis.

*Image Classification* is the task of assigning class labels to images. Usually, there is one dominant object in the image and the model is asked to predict the class of that dominant object. This is called *Single-Label Image Classification*. There is also an extension of *Single-Label Image Classification* called *Multi-Label Image Classification*. In *Multi-Label Image Classification*, more than one class is assigned to a given image. The model is asked to predict classes of all objects in the image. *Object Detection* takes *Multi-Label Image Classification* task a step further and aims to predict bounding boxes for the objects in the image, along with the classes of those objects. In this thesis, *Multi-Label Image Classification* and *Object Detection* models are utilized. Models used for these tasks are discussed in Section 2.1.1 and Section 2.1.2.

### 2.1.1  Multi-Label Classification

In the scope of this thesis, a multi-label classification method is used in order to extract object tags (class names of the objects in the image) which are utilized for image-text alignment as in [16]. ML-Decoder [33] is a state-of-the-art multi-label classification method. It is proposed to perform scalable multi-label classification even when number of classes ($N$) are large.

In ML-Decoder, input image $I$ is first passed through a Convolutional Neural Network (CNN) [34] and image embeddings (feature map) $M$ are obtained as in (2.1).

$$M = f(I) \tag{2.1}$$

where function $f$ models the CNN, $M$ represents the image embeddings of size $wh \times D$ where $D$ is the embedding dimension and $w, h$ are the width and height of the output feature map of the CNN, respectively.

Image embeddings $M$ and group queries (word embeddings for query classes) $E \in \mathbb{R}^{K \times D}$ for a set of fixed number of query classes are fed into a transformer (see Section 2.2 for the details about transformers) for cross-attention operation where $K$ is the number of query classes and $D$ is the embedding dimension. Outputs of the transformer for query classes are fed into a feed-forward layer. These operations are demonstrated in (2.2).

$$V = g_1(E, M)$$
$$Y = g_2(V) \tag{2.2}$$

where function $g_1$ models the transformer and cross-attention operation, $V \in \mathbb{R}^{K \times D}$ represents the output of the transformer for $K$ query classes, and $g_2$ function models the feed-forward network which produces final embeddings to be decoded, $Y \in \mathbb{R}^{K \times D}$.

Finally, outputs for $K$ classes are decoded to original $N$ classes via *group decoding*. In *group-decoding* a set of trainable linear layers are employed. Each query embedding $Y_i \in \mathbb{R}^{1 \times D}$ in $Y$ is decoded into $\frac{N}{K}$ class outputs and these outputs are concatenated to obtain outputs for $N$ classes according to (2.3).

$$Z_i = Y_i W_i \quad i = 1, 2, ..., K$$
$$Z = Concat(Z_i) \tag{2.3}$$

where $W_i \in \mathbb{R}^{D \times \frac{N}{K}}$ are the parameters of trainable linear layers specific for each group, $Z_i \in \mathbb{R}^{1 \times \frac{N}{K}}$ are the $\frac{N}{K}$ class outputs for each group, and $Z \in \mathbb{R}^{1 \times N}$ is the final multi-label classification output for $N$ classes.

Architectural diagram of ML-Decoder is provided in Figure 2.1.



Figure 2.1: Architecture of ML-Decoder (Adopted from [33]).

The baseline and proposed models in this thesis utilize object tags (class names) in an image for captioning. Since it is a state-of-the-art method and allows multi-label classification on a large dataset [3] with excessive number of classes ($\sim 600$), ML-Decoder model (with TResNet-M [34] backbone) trained on Open Images [3] is used in order to extract object tags.

### 2.1.2 Object Detection

Object detection methods are utilized by many recent image captioning methods for regional feature extraction [35], [24], [36], [16], [1], [4]. Furthermore, some group of methods make use of object tags (class names) in the image for captioning [16], [1], [4]. These methods utilize classification outputs of the object detection model as object tags.

Generally, Faster R-CNN [7] based object detectors are employed in image captioning methods. Faster R-CNN is a two-stage object detector. In the first stage, region proposals are obtained with a CNN-based *Region Proposal Network (RPN)*. RPN outputs a large set of bounding boxes (region proposals) which might contain objects. In the second stage, learned classification and regression heads refine the bounding boxes for proposals and filter out proposals with no objects. The architecture is shown in Figure 2.2.



Figure 2.2: Architecture of Faster R-CNN (Retrieved from [7]).

In the first stage, the input image is passed through a CNN backbone and a feature map is obtained. The feature map is fed into intermediate classification and regression heads in RPN. RPN produces candidate bounding boxes (region proposals) which might contain objects. Feature vectors corresponding to these proposals are extracted from the feature map using *RoIPooling* [37]. *RoIPooling* is an operation which ex-

tracts fixed-sized feature vector for each region proposal from the output feature map of the backbone CNN. It extracts the feature vectors by applying pooling on the projections of the bounding boxes on the feature map. Feature vectors for the region proposals are fed into the second stage.

In the second stage, feature vectors are passed through a set of fully-connected layers and fed into final classification and regression heads to eliminate proposals corresponding to background regions, refine and classify the bounding boxes.

In [35], regional features for detected objects are utilized during caption generation. They trained Faster R-CNN with ResNet-101 [5] backbone on Visual Genome dataset [38]. Visual Genome dataset contains ground truth attributes for objects in the image. In [35], they added extra attribute prediction head in order to help model learn richer features. However, attribute predictions are not utilized during caption generation. Faster R-CNN trained by [35] used in many captioning models as the standard approach [39], [40], [36].

Recently, VinVL [4] trained a larger Faster R-CNN model with ResNeXt-152 C4 [41] backbone on several datasets ([27], [3], [42], [38]) and achieved superior accuracy compared to its counterpart [16] which uses Faster R-CNN in [35].

## 2.2 Transformers

The invention of transformers [12] caused a paradigm shift in NLP and transformers have become the go-to method in recent approaches for various NLP tasks such as machine translation, sentiment classification, question answering ([13], [43], [14]). The ability to model relationships and interactions in long sequences and suitability for parallel processing are the main advantages of transformers [44]. Following the success of [12], transformer-based methods are utilized in many different areas including CV ([6], [45], [8]), VL ([30], [16] [1], [4], [46]).

The vanilla transformer proposed in [12] is composed of two blocks: *Encoder* and *Decoder*. Architectural details of these blocks and overall architecture of a transformer is given in Figure 2.3.

Figure 2.3: Architectural details of vanilla transformer and multi-head attention (Adopted from [12]).

Transformers are sequence-to-sequence (seq2seq) models. They accept a sequence of vectors as input and generate *transformed* output sequence. Unlike RNNs, transformers do not wait for new input at each time instant. Instead, they accept all inputs belonging to a sequence at once.

Given the input matrix $X \in \mathbb{R}^{N \times D_i}$ where $N$ is the sequence length, $D_i$ is the dimension of the input vectors, transformers first apply token embedding to these input vectors and add positional encoding as in (2.4) to obtain $X_{in} \in \mathbb{R}^{N \times D_e}$.

$$X_{in} = XW^E + P \tag{2.4}$$

where $W^E \in \mathbb{R}^{D_i \times D_e}$ is the embedding matrix, $P \in \mathbb{R}^{N \times D_e}$ is the positional encoding matrix. Since transformers operate on all elements of a sequence in parallel and at once, there is not a representation of order of elements in a sequence. To mitigate this problem, positional encoding is added to each input vector depending on position of the input in the sequence. Positional encoding can be fixed or learned. In [12], they experimented with both and observed that fixed positional encoding sampled from sinusoidal functions achieve similar results as the learned ones. Hence, they stick with the fixed positional embedding.

After $X_{in}$ is obtained, 3 different matrices are calculated according to (2.5).

$$\begin{aligned} Q &= X_{in}W^Q \\ K &= X_{in}W^K \\ V &= X_{in}W^V \end{aligned} \tag{2.5}$$

where $Q \in \mathbb{R}^{N \times D_k}$ is the query matrix, $K \in \mathbb{R}^{N \times D_k}$ is the key matrix, and $V \in \mathbb{R}^{N \times D_v}$ is the value matrix whereas $W^Q \in \mathbb{R}^{D_e \times D_k}$, $W^K \in \mathbb{R}^{D_e \times D_k}$, and $W^V \in \mathbb{R}^{D_e \times D_v}$ are learned projection matrices to obtain $Q, K, V$ matrices respectively.

$Q, K, V$ matrices contain query ($q_i \in \mathbb{R}^{D_k}$), key ($k_i \in \mathbb{R}^{D_k}$), and value ($v_i \in \mathbb{R}^{D_v}$) vectors in their rows for each element in the input sequence (for $i = 0, 1, ..., N - 1$). These matrices are used to calculate scaled dot-product attention according to (2.6).

$$Attention(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{D_k}}\right)V \tag{2.6}$$

Scaled dot-product attention calculates the relationship (attention) between query and key vectors for different inputs, and performs weighted sum of the value vectors to

obtain *transformed* representation for each query input. This process is followed by an addition and normalization (*Add & Norm*) layer. Then, outputs of this layer are fed into a linear (feed-forward) layer and a final *Add & Norm* layer. These layers are stacked $L$ times as in Figure 2.3. Output of final *Add & Norm* block in each layer is fed to the *Multi-Head Attention* block of the next layer.

Usually, in order to allow model to attend vectors from different representation subspaces, multi-headed approach is used. In this approach, all operations are performed $A$ times in parallel, with different learned matrices ($W^Q, W^K, W^V$) where $A$ is the number of heads, as shown in Figure 2.3. Outputs from different heads are merged according to (2.7).

$$
\begin{aligned}
MultiHeadAttention(Q, K, V) &= Concat(head_1, head_2, ..., head_A)W^O \\
\text{where } head_i &= Attention(QW_i^Q, KW_i^K, VW_i^V)
\end{aligned}
\tag{2.7}
$$

where $W_i^Q \in \mathbb{R}^{D_{model} \times D_k}$, $W_i^K \in \mathbb{R}^{D_{model} \times D_k}$, $W_i^V \in \mathbb{R}^{D_{model} \times D_v}$, and $W_i^O \in \mathbb{R}^{AD_v \times D_{model}}$. Here, $Q, K, V \in \mathbb{R}^{N \times D_{model}}$ come from the outputs of the previous transformer blocks.

The process is called *self-attention* when query, key, and value vectors belong to the same sequence. It is called *cross-attention* when key and value vectors are from one domain and query vectors are from another domain.

Decoder block operates in a sequential manner. Outputs at one time instant is given as input at the next time instant. In the decoder block *masked-self-attention* is performed. In *masked-self-attention*, query vectors are only allowed to attend key-value pairs from previous time instants. Attending to future time instants are prevented by setting softmax inputs for those time instants to $-\infty$ in (2.6). This is referred as *attention masking*. It is necessary to apply attention masking during training in order to preserve auto-regressive property since the model will not have tokens from future time instants during inference.

Decoder block employs *cross-attention* as well. Key and value vectors obtained from the outputs of the encoder block, and query vectors obtained from the previous decoder outputs are used to apply *cross-attention* using (2.6).

12

### 2.2.1 BERT: Bidirectional Encoder Representations from Transformers

After the success of vanilla transformer [12], many variants of it were proposed ([47], [43], [14]) and trained for different tasks. *Bidirectional Encoder Representations from Transformers (BERT)* [13] was one of them. It became quite popular in a short amount of time because it allowed utilization of a single bidirectional architecture on different tasks by exploiting large datasets via unsupervised pretraining [13]. BERT is also the base model for the methods analyzed ([1], [4]) and used in this thesis.

BERT is a transformer architecture which has only the encoder block. BERT has two models with the following parameters:

- $\text{BERT}_{\text{base}} : L = 12 \ A = 12 \ H = 768$

- $\text{BERT}_{\text{large}} : L = 24 \ A = 16 \ H = 1024$

where $L$ is the number of consecutive transformer blocks, $A$ is the number of attention heads, and H is the hidden embedding dimension of the transformer (dimensions for embedding, query, key, and value are same in this case).

$\text{BERT}_{\text{base}}$ is used in [16], [1], [4], and proposed method in this thesis.

BERT has special tokenizer which converts sentences to sequence of tokens (words or subwords). These tokens are the inputs of the architecture. They are first passed through embedding layer, then positional encoding is added. Different from vanilla transformer [12], they used learned positional encoding instead of a fixed one. Furthermore, BERT also adds segment embedding vector after positional encoding. Segment embeddings represent different parts of input data. These parts could be different sentences in a paragraph, a question and its corresponding answer, or a caption for an image and corresponding regional features as in [1], [4].

BERT has 3 special tokens:

- $[CLS]$: Classification Token

- $[SEP]$: Separator Token

- $[MASK]$: Mask Token

Each BERT input starts with $[CLS]$ token. Inputs belonging to different segments are separated with $[SEP]$ token. Example tokenization of two sentences (two segments) is given in 2.4.



Figure 2.4: BERT tokenization and embeddings (Retrieved from [13]).

BERT utilizes a pretraining objective called *Masked Language Modelling (MLM)*. In MLM, network is given a corpus (sequence of words). Random words in the sequence are replaced with the special $[MASK]$ token. After passing inputs through Multi-Head Self-Attention, output corresponding to $[MASK]$ token is given to a classifier head. Network tries to guess what the masked token was. The network is forced to deduce the masked word from its surrounding words with the help of attention mechanism [13]. Usually, models are pretrained on large corpora with MLM objective and finetuned on downstream tasks.

Methods utilizing BERT for sentence generation (seq2seq objective) performs autoregressive decoding [30], [16], [1], [4]. They give $[MASK]$ token as input for the word to be predicted along with previously predicted words and classify transformer output corresponding to the $[MASK]$ token in order to predict the current word. Predicted word is fed to the model in the next time instant and the next input position is masked. The network predicts the masked word at the output.

## 2.3 Assignment Problem

Assignment problem is a fundamental problem which appears in many areas where elements belonging to two sets should be matched with each other so that overall

matching cost is minimum. For example, consider there are 3 workers and 3 jobs which will be assigned to the workers. Assuming cost of an assignment is the time it takes for a worker to finish the given job, an assignment algorithm seeks to find the optimum assignment between workers and jobs so that 3 jobs are completed in minimum amount of time. There are different algorithms which solve assignment problem such as [48], [49], [50]. Among these, Hungarian Algorithm [48] is one of the most widely used assignment algorithm. Given an $N \times N$ cost matrix $C$ where rows would be assigned to columns, Hungarian Algorithm finds the optimal assignment in polynomial time. Following example demonstrates the steps of Hungarian Algorithm for a $3 \times 3$ cost matrix $C$:

$$C = \begin{bmatrix} 100 & 125 & 150 \\ 150 & 135 & 175 \\ 120 & 140 & 200 \end{bmatrix}$$

**Step 1:** Subtract the smallest entry of each row from elements of the corresponding row: *Smallest entries for first, second, and third row are 100, 135, and 120, respectively.*

$$C_1 = \begin{bmatrix} 0 & 25 & 50 \\ 15 & 0 & 40 \\ 0 & 20 & 80 \end{bmatrix}$$

**Step 2:** Subtract the smallest entry of each column from elements of the corresponding column: *Smallest entries for first, second, and third columns are 0, 0, and 40, respectively.*

$$C_2 = \begin{bmatrix} 0 & 25 & 10 \\ 15 & 0 & 0 \\ 0 & 20 & 40 \end{bmatrix}$$

**Step 3:** Cover all the zeros with minimum number of lines (vertical or horizontal): *For $C_2$, to cover all zeros, minimum of 2 lines is necessary (a line through first column and a line through second row).*

**Step 4:** If minimum number of lines to cover all zeros is equal to number of rows or columns ($N$), go to **Step 7**, else go to **Step 5**: *Since 2 is smaller than $N = 3$, we proceed with **Step 5**.*

**Step 5:** Find the smallest element which is not crossed by a line drawn in **Step 3**.

15

Subtract this element from all elements of uncovered rows. Add this element to all elements of covered columns. *Smallest element is 10 in this case. Only the second row is covered. Only the first column is covered. Subtract 10 from first and third rows to obtain $C_3$. Add 10 to the first column to obtain $C_4$.*

$$C_3 = \begin{bmatrix} -10 & 15 & 0 \\ 15 & 0 & 0 \\ -10 & 10 & 30 \end{bmatrix} \quad C_4 = \begin{bmatrix} 0 & 15 & 0 \\ 25 & 0 & 0 \\ 0 & 10 & 30 \end{bmatrix}$$

**Step 6:** Return to **Step 3**. *Now **Step 3** will result in minimum of 3 lines to cover all zeros in $C_4$. Hence, we proceed to **Step 7**.*

**Step 7:** Optimal assignment is obtained by selecting 0 entries so that each column and row has only one 0 entry selected: *Selected entries are shown with *.*

$$C_4 = \begin{bmatrix} 0 & 15 & 0^* \\ 25 & 0^* & 0 \\ 0^* & 10 & 30 \end{bmatrix} \quad C = \begin{bmatrix} 100 & 125 & 150^* \\ 150 & 135^* & 175 \\ 120^* & 140 & 200 \end{bmatrix}$$

Hence, worker-1 is assigned to job-3, worker-2 is assigned to job-2, and worker-3 is assigned to job-1. Cost of this optimum assignment is $150 + 135 + 120 = 405$.

For non-square cost matrices dummy zero rows or columns are added to make the matrix square. Afterwards, the same procedure described above is applied but during **Step 7**, entries belonging to dummy rows or columns are not selected for assignment.

The intuition behind Hungarian Algorithm is that the optimum assignment does not change if a number is subtracted from or added to all elements in a row or in a column. Hence, the optimum assignment can be found by reducing the cost matrix with addition/subtraction operations on rows/columns such that 0 entries representing the one-to-one assignment are obtained.

# CHAPTER 3

## LITERATURE REVIEW

In this chapter, literature review for image captioning is provided. Recent applications of image captioning are discussed. A taxonomy of the image captioning models and their details are provided. Datasets utilized by these models are explained. Evaluation metrics for the image captioning task are also provided.

## 3.1   Applications of Image Captioning

Image captioning is a challenging and broad topic. Different methods tackle image captioning from different angles. Even though all methods aim at generating visually-grounded, rich, natural sounding, grammatically correct, and fluent captions, they might specialize in various domains.

Some of the recent methods [30], [16], [4] focus on exploiting datasets with large amount of image-caption pairs [27], [28], [29], [2]. These methods generally follow a two-stage training pipeline where in the first stage, model is trained on large image-text paired datasets with general Vision and Language (VL) objectives. This step is called *Vision and Language Pretraining (VLP)*. In the second stage, models are finetuned on task-specific datasets for different tasks (image captioning, visual question answering). Since the models are trained on numerous large datasets, these approaches require significant computational resources for training.

Another group of methods focus on generating captions for images which contain objects which are not present in the image captioning training dataset, [24], [25], [1], [26]. This task is often referred as *Novel Object Captioning (NOC)*. Generally, NOC

methods make use of external object detection algorithms ([7], [51]) to introduce novel objects in the generated caption.

Another branch of research targets grounded image captioning. There are studies showing current image captioning models do not fully refer to the input image while generating captions, instead they are prone to biases in the training dataset and language model [31], [32]. Recent methods try to overcome this issue by introducing additional information to the captioning model or explicitly modelling visual grounding for the predicted words [31], [32].

VizWiz [17], [52] is a dataset and a challenge which aims to guide visually impaired people with image captioning and visual question answering. Hence, captions in this dataset is more specific and detailed compared to general datasets. Models developed for this challenge usually include an Optical Character Recognition (OCR) component as well [53].

Several methods focus on generating stylized captions [54], [55], [56], [57]. These methods try to generate positive, negative, emotional, humorous, romantic, or factual captions for the same image.

Some methods are proposed for captioning images in news articles [18], [58]. These methods usually exploit the article itself in order to generate more detailed and informative captions.

Another area of interest is image paragraph captioning. Image paragraph captioning models describe the details in the image by generating a paragraph, possibly with more than one sentence [59], [60], [61].

## 3.2   Taxonomy of Image Captioning Methods

Since image captioning lies at the intersection of Computer Vision (CV) and Natural Language Processing (NLP), it borrows ideas from both of these domains. Consequently, in general, image captioning models contain two blocks: a *Visual Encoder* block and a *Language Model* block [21]. Visual Encoder block processes the input image and generates encoded representation of it, called *context vector*. Language

Model block generates caption describing the image based on the context vector. Visual Encoder block is referred as *encoder* and Language Model block is referred as *decoder*. The process of caption generation is illustrated in Figure 3.1. The feedback path from Language Model to Visual Encoder is optional and it exists in models which utilize *attention* mechanism, as will be discussed in Section 3.2.1.



Figure 3.1: General image captioning block diagram.

Image captioning models differ in methods they utilize in one or both of these two blocks. Similar to [21], image captioning models can be categorized considering the approaches they use in these two blocks. Overall taxonomy of image captioning models are given in Figure 3.2.



Figure 3.2: General image captioning taxonomy [21].

A literature review of image captioning models are provided in the following sections according to the taxonomy shown in 3.2.

### 3.2.1 Visual Encoder

Visual encoding is very important in image captioning because encoded representation (context vector) is what the language model sees while generating words. It is challenging to obtain a good visual encoding since a good encoding should distill the information in the image while preserving fine-grained details [21].

As large scale datasets such as ImageNet [62], COCO [27] are collected and computational power of computers increased, Convolutional Nerual Networks (CNNs) [63] were reborn and became the go-to method for image understanding and representation tasks [64], [65], [5], [7], [9]. Hence, CNNs are also at the heart of the Visual Encoder block for most of the captioning methods.

As presented in Figure 3.2, approaches for Visual Encoder can be examined under 3 groups: *Non-Attentive* approaches, *Additive Attention* based approaches, *Self-Attention (Transformer)* based approaches.

Earliest and pioneering deep learning based image captioning methods directly utilized output of a CNN as the visual encoding [66], [67]. In such a scenario, CNN acts as a global feature extractor and generates a context vector for the input image. Language Model is conditioned on this context vector while generating words for the caption. Whole output feature map of the CNN is directly used while generating context vector without any attention, hence the name Non-Attentive. This process is illustrated in Figure 3.3 - a. Among Non-Attentive approaches, Show and Tell [66] was one of the first and prominent methods which utilized deep learning for image caption generation. They use Inception [65] network for the Visual Encoder and feed the context vector to LSTM-based (Long Short-Term Memory) [68] Language Model. After the success of Show and Tell, other methods [69], [70], [71] also used Non-Attentive Visual Encoders with improved CNN architectures [72], [5]. Non-Attentive approach is simple and it encodes the input image as a whole, globally. On the other hand, since the context vector is calculated for once at the beginning, all words in the language model are conditioned on the same vector during decoding. Hence, this approach may overlook fine-grained details in the image while generating words [21].

Additive Attention based approaches aim to overcome aforementioned disadvantages of Non-Attentive based approaches due to context vector being fixed and global. Additive Attention based approaches obtain weights for a set of feature vectors in which weights depend on the state of the language model. Afterwards, using these weights, weighted averaging of the feature vectors are performed in order to obtain context vector for the next step of decoding. The idea of Additive Attention originates from a machine translation method [73] where a new context vector is obtained by applying Additive Attention over hidden states of encoder Recurrent Neural Network (RNN) at each time instant of the decoding stage.

Given an image captioning model with CNN-based Visual Encoder and RNN-based Language Model blocks, let $V = v_1, v_2, ..., v_L$ be the set of feature vectors for the input image. Then, the context vector $c_t$ for the time instant $t$ can be calculated as in (3.1).

$$e_{ti} = f_{att}(v_i, h_{t-1})$$
$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{k=1}^{L} \exp(e_{tk})} \qquad (3.1)$$
$$c_t = \sum_{v_i \in V} \alpha_{ti} v_i$$

where $f_{att}$ is the attention function which calculates the effect ($e_{ti}$) of the feature vector $v_i$ to the state vector of the RNN at the previous time instant, $h_{t-1}$. $\alpha_{ti}$ is the attention weight for the feature vector $v_i$ at time instant $t$ which is calculated by applying softmax to $e_{ti}$ for $i = 1, 2, ..., L$.

Inspired from [73], Show, Attend, and Tell [74] was the first image captioning method which utilized Additive Attention based Visual Encoder instead of global and Non-Attentive one. They apply Additive Attention over grid features of the CNN output as in Figure 3.3-b. After the success of Show, Attend, and Tell, other methods also utilized variants of Additive Attention over grid features in their Visual Encoder block ([75], [76], [77]).

In [35], another pioneering method, BUTD (Bottom-Up Top-Down), for image captioning was proposed. In BUTD, two types of attention is utilized which they call bottom-up attention and top-down attention. Top-down attention is the aforementioned attention mechanism which is based on the relationship between feature vec-

Figure 3.3: Different visual encoding approaches (a: Non-Attentive, b: Additive Attention over Grid Features, c: Additive Attention over Regional Features, d: Self-Attention Followed by Additive Attention, Dashed red line indicates feedback from Language Model block which is used for attention calculations).

tors and state of the language model, as in (3.1). Bottom-up attention is used to decide which regions of the image contain valuable (salient) information. For this purpose, an object detection method, Faster R-CNN [7], is utilized to generate region proposals and regional feature vectors for objects in the image. Then, the set of feature vectors, $V = \{v_1, v_2, ..., v_L\}$, are used in Additive Attention instead of grid features as in (3.1). This process is visualized in Figure 3.3-c. Faster R-CNN detector used in this approach is trained on Visual Genome [38] dataset with an extra attribute prediction head and loss in order to enhance learned feature representations [35]. After the success of BUTD, it became standard approach to apply Additive Attention over regional feature vectors. Many other approaches utilized regional features coming from an object detector in their Visual Encoder block with Additive Attention [39], [40], [36].

The invention of transformers [12] caused a paradigm shift in NLP and transformers have become the go-to method in recent approaches [13], [43], [14]. Transformers introduce and make use of *self-attention* and *cross-attention* concepts [12]. Details

22

of these concepts and transformers are provided in Section 2.2. Following the success of transformers, self-attention based methods are utilized in many different areas including Computer Vision ([6], [8]), Vision and Language ([1], [4], [46]).

Since self-attention is a good way of representing inter-connections in a set of vectors and performing *message passing* between them [78], image captioning methods also started utilizing self-attention mechanism. One of the first methods with Self-Attention (Transformer) based Visual Encoder block was proposed in [31]. They employed self-attention block in order to transform regional features to a space where interactions (relationships) between objects are also modelled. Other methods also enhanced regional feature vectors with self-attention block [79], [80]. In [81] and [82] they used modified-self-attention mechanisms such that geometric relationships between objects are also taken into consideration for attention weights. In [83], they extended self-attention mechanism with learned memory vectors which model a priori information between regional features. In general, these methods exploit relationships between regional feature vectors via self-attention. Then, enhanced feature vectors are fed into an additive attention block to obtain final context vector. This process is illustrated in Figure 3.3-d.

Recently transformers are also used as a feature extractor similar to CNNs [6], [45]. These approaches divide the image into patches and perform self-attention on these patches. These architectures are called Vision Transformers [6]. Inspired from these approaches, some image captioning methods employed Vision Transformers in their Visual Encoder blocks [84], [85].

Another branch of work utilizing self-attention focuses on unified architecture for Visual Encoder and Language Model blocks. This kind of approach was first introduced in [30], called Unified-VLP (Unified Vision and Language Pretraining). In Unified-VLP, they use a single transformer network for both encoding and decoding steps. It is also unified in the sense that, a single architecture is pretrained on large image-text paired datasets and then finetuned on task-specific datasets for downstream tasks such as image captioning or visual question answering. The architecture of BERT (Bidirectional Encoder Representations from Transformers) [13] is directly used for captioning. In this approach, regional image features are extracted as in [35]. Re-

gional image feature vectors and word embeddings are fed into transformer together. During pretraining, network is trained with Masked Language Modelling (MLM) objective similar to BERT as explained in Section 2.2.1. During fintuning for the image captioning task, the network is trained with sequence-to-sequence (seq2seq) objective. OSCAR [16] further improved Unified-VLP with the integration of object tags which helps alignment between image features and language features. VinVL [4] was proposed as an improvement over OSCAR [16]. In VinVL, the captioning architecture is the same as OSCAR but the object detection model is improved by training a larger Faster R-CNN [7] model on a collection of datasets ([27], [3], [42], [38]) instead of just one ([38]). This helps extracting richer regional features and results in better captioning performance [21].

VIVO [1] is another Unified-VLP based method which was proposed for Novel Object Captioning. In VIVO, pretraining is performed on an object detection dataset [3] and finetuning is performed on image captioning dataset [27].

VinVL achieves state-of-the-art results on many VL tasks [4]. Approaches in VIVO and VinVL are taken as baseline in this thesis. Analyses and improvement strategies are employed on this baseline. Thus, these methods are discussed in more detail in Chapter 4.

### 3.2.2 Language Model

Language Model is a crucial component of image captioning methods since it generates final output of the overall system. A good Language Model should be able to generate fluent and grammatically correct captions which are also grounded on the context vector generated by Visual Encoder.

In NLP, the language model is defined as a model which predicts probability of occurrence of a sequence of words in a sentence [21]. In image captioning, Language Model models the probability of occurrence of the caption given the context vector. This probability is defined in (3.2).

$$P(y_N, y_{N-1}, ..., y_1 \mid \boldsymbol{X}) = \prod_{i=1}^{N} P(y_i \mid y_1, y_2, ..., y_{i-1}, \boldsymbol{X}) \qquad (3.2)$$

24

where $X$ is the context vector, $y_i$ is the word in the caption at the $i$-th time instant, and $N$ is the length of the caption.

Usually, during caption generation (decoding) step, Language Model is given a special token (vector) called *start token* as the first word. Afterwards, model predicts most probable next word. In an auto-regressive manner, new words are predicted by conditioning on previously predicted words. This process usually stops when the Language Model predicts another special token called *end-of-sequence* or when the maximum sequence length is reached [21].

As presented in Figure 3.2, there are 3 main approaches in Language Model block of the image captioning models: *LSTM-based* methods, *CNN-based* methods, *Self-Attention (Transformer)* based methods.

Caption generation is basically a sequence generation task. RNNs are frequently used in sequence generation tasks in NLP domain since they are well-suited for sequential and time-series data [11], [73]. Hence, it is also natural to use RNNs, specifically LSTMs [68], for caption generation task. Earliest and pioneering approaches in image captioning utilized LSTM-based approaches in their Language Model block.

In [66], [74], [67], context vector output of the Visual Encoder is fed into LSTM. For the first time step, special start token (*<START>*) is given as input and decoding starts. Afterwards, LSTM generates words in an auto-regressive manner as in (3.2). This process is illustrated in Figure 3.4.

There are other approaches with more complex Language Model blocks which also utilize LSTMs. For example, in [69], [35], [24], they used 2-layer LSTM where first layer generates a hidden state vector which is utilized in attention block and fed as input to the second layer. Second layer acts as the language model similar to the ones in [66], [74], [67]. Since they have high representation power, 2-layer LSTMs are heavily utilized in Language Model block until Transformers [12] supersede them [21].

In order to overcome complex structure and sequential nature of RNNs, [86] used a CNN-based approach in the Language Model block inspired from [87]. Context vector and word embeddings are fed into a CNN. The CNN operates on all inputs in

Figure 3.4: LSTM-based language model.

parallel with convolutional kernels and generates output words. In order to preserve sequential nature of the language and prevent predicting words based on the future time instants, they utilized masking mechanism during training similar to [12]. Even though parallel training was an improvement compared to LSTMs, this approach did not gain too much popularity due to inferior performance and rise of transformer-based approaches [21].

In [12], transformer architecture is proposed for machine translation task and later it became more and more popular in the NLP domain [13], [43], [14]. Transformers allow parallel processing of the input data and modelling dependencies between them, even for long sequences due to self-attention [12]. Thanks to these benefits, image captioning models also started utilizing transformers in their Language Model block.

Image captioning models with Transformer-based Language Model blocks became popular in recent years. Some methods such as [81], [82], [85] directly replaced RNNs with transformer decoder blocks. These models take the context vector from Visual Encoder block, apply cross-attention between encoder output and decoder states, and apply self-attention between word embeddings in an auto-regressive manner to perform caption generation.

As discussed in Section 3.2.1, there are also methods such as Unified-VLP [30], OS-

CAR [16], VIVO [1], VinVL [4] which utilize unified transformer architecture for Visual Encoder and Language Model blocks.

## 3.3 Datasets

Datasets play a crucial role for image captioning models since image captioning is a data-driven task [21]. As computational capabilities of machines increase, models get bigger, larger datasets are collected and utilized.

Summary of the common datasets used by image captioning models are summarized in Table 3.1.

One of the first systematically collected dataset for image captioning was Pascal Sentences [88] dataset. For this dataset, 1000 images are randomly sampled from Pascal VOC2008 [92] dataset and they are human-annotated with 5 captions per image. Currently, this dataset is not used since it is relatively simple and larger datasets exist.

Flickr8k [89] and Flickr30k [2] datasets are composed of images obtained from Flickr [93]. These datasets contain around 8,000 and 30,000 images respectively, and they have 5 ground-truth captions per image. These two datasets are still used by recent image captioning methods.

COCO [27] is the most widely used dataset in the image captioning literature since it is a large and diverse dataset [21]. It contains around 80,000 images for training, 40,000 images for validation, and 40,000 images for testing. Each image in the training and validation sets have 5 ground-truth human-annotated captions. These annotations are publicly available. A small portion (5,000 images) of the test set contains 40 ground-truth captions per image and the rest contains 5 ground-truth captions per image. Ground-truth annotations for the test set are not publicly available, they are used in the online challenge [94]. Hence, in [67] they have come up with different splits of training and validation sets such that 5,000 images are used for validation, 5,000 images are used for testing and the rest (113,288) are used for training. These splits are called *Karpathy splits* and accepted as standard by the image captioning literature [21].

Table 3.1: Summary of common datasets utilized in image captioning.

| Dataset | Task | Number of Images | Details |
|---|---|---|---|
| Pascal Sentences [88] | Image Captioning | 1,000 | – 5 ground-truth captions per image. |
| Flickr8k [89] | Image Captioning | 8,000 | – 5 ground-truth captions per image. |
| Flickr30k [2] | Image Captioning | 30,000 | – 5 ground-truth captions per image. |
| COCO [27] | Image Captioning | Train: 114,000 Validation: 5000 Test: 5,000 | – 5 ground-truth captions per image. |
| SBU [29] | Image Captioning | 1,000,000 | – 1 ground-truth caption per image. |
| CC [90] | Image Captioning | 12,000,000 | – Contains web images and their alt-text as captions. |
| *nocaps* [91] | Novel Object Captioning | Validation: 4,500 Test: 10,600 | – Images are from Open Images [3] object detection dataset. – 10 ground-truth captions per image. – Contains 3 subsets based on the novelty of objects in the image: *in-domain, near-domain, out-of-domain.* |
| Open Images [3] | Object Detection | 1,700,000 | – Contains around 600 novel object classes. – Utilized while training novel object captioning models for *nocaps* challenge. |
| Visual Genome [38] | Scene Graph Generation | 109,000 | – Contains ground-truth object classes, boxes, attributes, and relationships. – Utilized while training regional feature extractors in image captioning models. |

SBU [29] is a dataset with around 1 million images from Flickr with each having 1 description. This dataset is mostly used during the pretraining stage of Vision and Language Pretraining (VLP) based methods.

Conceptual Captions [28] is a dataset which is also being used during the pretraining stage of VLP based methods. It contains more than 3 million images obtained from web. Alt-text of these images are filtered and cleaned in order to obtain an image-text paired dataset. Recently, a larger version of this dataset with 12 million images is published [90]. YFCC100M [95] is also a similar dataset containing web images and auto-generated descriptions.

A specific branch of image captioning is called *Novel Object Captioning*. Novel Object Captioning models aim to generate captions for images which contain novel objects which are not seen in the captioning (image-text paired) dataset. *nocaps* challenge [91], [96] is specifically designed for Novel Object Captioning. In this challenge, two datasets are allowed to be used: Open Images [3] and COCO [27]. Open Images is a large object detection dataset with 600 different classes. In *nocaps*, it is expected from models to learn novel objects from Open Images object detection dataset and transfer that knowledge to image captioning via training on COCO dataset which contains only 80 objects. A subset (4,500 images) of Open Images validation set and a subset (10,600 images) of Open Images test set are labelled for captioning with 10 captions per image. These two sets are used as validation and test sets of the *nocaps* challenge. It is not allowed to use image-text paired dataset other than COCO. Validation and test sets of the *nocaps* dataset are divided into three subsets for evaluation purposes: *in-domain*, *near-domain*, and *out-of-domain*. Images in *in-domain* subset contain objects from COCO classes and do not contain novel objects from Open Images dataset. Images in *near-domain* subset contain objects from both COCO classes and novel classes from Open Images dataset. Images in *out-of-domain* subset only contain novel objects from Open Images dataset. Training and evaluation setups and subsets are illustrated in Figure 3.5.

Example images and corresponding ground truth captions from most popular image captioning datasets are given in Figure 3.6

There are also other datasets for specific branches of image captioning. For example,

Figure 3.5: *nocaps* task setup [91].



Figure 3.6: Examples from popular image captioning datasets.

VizWiz [97] is a dataset and challenge which aims to help visually impaired people. Hence, captions in this dataset is more specific and detailed compared to general

datasets. In [59], a new dataset is proposed for image paragraph captioning where a model is expected to describe details in the image by generating a paragraph, possibly with more than one sentence. FlickrStyle10K [55] is a dataset which contains stylized captions such as humorous or romantic.

Visual Genome [38] is a dataset which is indirectly utilized by most of the image captioning methods. Visual Genome is a Scene Graph Generation dataset where not only objects in the scene, but also their attributes and relationships with each other are also labelled. Hence, an image in this dataset contains objects, their bounding boxes, attributes, and relationships with other objects. An example image from Visual Genome and its scene graph is given in Figure 3.7.



Figure 3.7: An example image from Visual Genome dataset and its scene graph. Pink elements represent *objects*, green elements represent *attributes*, blue elements represent *relationships* in the scene graph.

Since it is a large and informative dataset, Visual Genome is utilized by [35] in order to train an object detection model for regional feature extraction. VinVL [4] further improved this approach by training a larger model on a collection of datasets (COCO [27], Open Images [3], Objects365 [42], Visual Genome [38]). These models train the object detector on Visual Genome dataset with an extra attribute prediction head but they do not use attribute predictions during training or inference of the captioning models. Aim in this approach is helping model extract richer regional features [35].

## 3.4 Evaluation Metrics

As in all scientific works, evaluation metrics are just as important as the proposed method itself. Evaluation metrics are critical to assess the quality of a proposed method and compare it with the existing approaches. In image captioning models, model output and ground truth targets are sentences. Hence, it is natural to use metrics which aim to compare similarity between a candidate sentence and a set of reference sentences. For image captioning, a good caption should also be grounded to input image along with being natural sounding and grammatically correct [21].

Standard evaluation metrics for image captioning compare similarity of a proposed caption to set of ground truth reference sentences [21]. [98] further classifies these metrics into 2: The ones that are borrowed from NLP (translation, summarization), the ones that are developed exclusively for image captioning. BLEU [99], METEOR [100], ROUGE [101] belong to the first category whereas CIDEr [102] and SPICE [103] are specifically proposed for image captioning task.

Most of these metrics (except SPICE) utilize a concept called $n$-grams [98]. An $n$-gram is a sequence of $n$ words which are consecutive in a sentence. For example, the set of *bigrams* in a sentence contains each possible consecutive two words in the sentence.

### 3.4.1 BLEU: Bilingual Evaluation Understudy

BLEU (Bilingual Evaluation Understudy) [99] is a metric proposed for machine translation task. Its aim is measuring similarity of a candidate translation to a set of reference translations. It calculates $n$-gram precision between candidate sentence and reference sentences. It can be formulated as in (3.3) for different $n$ values as $n$-grams.

$$\text{BLEU}_n(x, Y) = BP \times \frac{\sum_{w_n \in x} \min\left(c_x(w_n), \max_{i=1,...,k} c_{Y_i}(w_n)\right)}{\sum_{w_n \in x} c_x(w_n)} \tag{3.3}$$

where $x$ is the candidate sentence, $Y = \{Y_1, Y_2, ..., Y_k\}$ is the set of reference sentences, $w_n$ is $n$-gram, $c_a(w_n)$ is the number of occurrences of $n$-gram $w_n$ in sentence $a$. $BP$ is called Brevity Penalty and it is used to penalize short sentences. It is defined

as in (3.4).

$$BP = \begin{cases} 1 & \text{if } c > r \\ e^{(1-\frac{r}{c})} & \text{if } c \leq r \end{cases} \tag{3.4}$$

where $c$ is the length of the candidate sentence and $r$ is the length of the reference sentence whose length is closest to the that of candidate sentence.

Generally cumulative BLEU-4 score is reported as the overall BLEU score and it is the geometric mean of $\text{BLEU}_n$ scores for $n = 1, ..., 4$

$$\text{BLEU-4} = \left( \prod_{j=1}^{4} \text{BLEU}_n \right)^{\frac{1}{4}} \tag{3.5}$$

### 3.4.2 ROUGE: Recall-Oriented Understudy for Gisting Evaluation

ROUGE (Recall-Oriented Understudy for Gisting Evaluation) [101] is a metric introduced to measure the quality of a summary. It measures similarity between a predicted summary and a set of ground truth summaries by considering $n$-gram recall. It can be formulated as in (3.6) for different $n$ values as $n$-grams.

$$\text{ROUGE}_n(x, Y) = \frac{\sum_{Y_i \in Y} \sum_{w_n \in Y_i} c_x(w_n)}{\sum_{Y_i \in Y} \sum_{w_n \in Y_i} c_{Y_i}(w_n)} \tag{3.6}$$

where $x$ is the candidate sentence, $Y = \{Y_1, Y_2, ..., Y_k\}$ is the set of reference sentences, $w_n$ is $n$-gram, $c_a(w_n)$ is the number of occurrences of $n$-gram $w_n$ in sentence $a$.

In [101], different variants of ROUGE metric utilizing Longest Common Subsequence (LCS) [104], weighted LCS, skip-bigrams are proposed.

As opposed to BLEU [99] which focuses on $n$-gram precision, ROUGE focuses on $n$-gram recall. Hence, it favors long sentences [102].

### 3.4.3 METEOR: Metric for Evaluation of Translation with Explicit ORdering

METEOR (Metric for Evaluation of Translation with Explicit ORdering) [100] is also a metric proposed for machine translation task. It takes both precision and recall into

account as opposed to BLEU [99] or ROUGE [101] which focuses on just precision or recall, respectively.

It can be formulated as in (3.7).

$$\text{METEOR} = F_{mean} \times (1 - penalty) \tag{3.7}$$

where $F_{mean}$ is the harmonic mean of unigram precision and recall with more focus on recall as defined in (3.8).

$$F_{mean} = \frac{10PR}{R + 9P} \tag{3.8}$$

where

$$\begin{aligned} P &= \frac{m}{w_t} \\ R &= \frac{m}{w_r} \end{aligned} \tag{3.9}$$

with $m$ being number of unigrams in the candidate sentence which also appear in the reference sentence, $w_t$ being number of unigrams in the candidate sentence, $w_r$ being number of unigrams in the reference sentence.

$penalty$ is a term to penalize large number of *chunks*. A *chunk* is defined as a set of unigrams adjacent in both candidate and reference sentence. For a perfect translation where candidate sentence is the same as the reference sentence, number of *chunks* is 1. $penalty$ term is calculated according to (3.10).

$$penalty = 0.5 \times \left(\frac{c}{m}\right)^3 \tag{3.10}$$

where $c$ is the number of chunks and $m$ is the number of matched unigrams (unigrams that exist in both candidate and reference).

METEOR also allows matching synonyms and words with identical stem. Hence, as opposed to BLEU [99] and ROUGE [101], METEOR captures semantic similarities between candidate and reference sentences to some extent [98].

### 3.4.4 CIDEr: Consensus-based Image Description Evaluation

CIDEr (Consensus-based Image Description Evaluation) [102] is specifically proposed for image captioning task. Its aim is measuring consensus between a proposed

caption and a set of human-annotated ground truth captions for an image. To measure such consensus, an $n$-gram based approach is proposed. In this approach, $n$-grams which are less informative are down-weighted. To decide informativeness (weight) of $n$-grams in the candidate sentence and ground truth sentences, *Term Frequency-Inverse Document Frequency (tf-idf)* based approach is used. *tf-idf* weight for an $n$-gram is computed as in (3.11).

$$g_k(s_{ij}) = \frac{h_k(s_{ij})}{\sum_{\omega_l \in \Omega} h_l(s_{ij})} \log \left( \frac{|I_s|}{\sum_{I_p \in I} \min \left( 1, \sum_q h_k(s_{pq}) \right)} \right) \tag{3.11}$$

where $h_k(s_{ij})$ is count of $n$-gram $\omega_k$ in reference sentence $s_{ij}$, $\Omega$ is the vocabulary of all $n$-grams, $I_s$ is the set of all images in the dataset, $|I_s|$ is number of images in the dataset. The first term in (3.11) is *term frequency* and it measures how frequent $n$-gram $\omega_k$ occurs in the sentence $s_{ij}$. If it is frequent, it may be informative for current sentence, hence, larger weight. The second term is *inverse document frequency* and it measures discriminativeness of $\omega_k$. Number of images in the dataset is divided by number of images whose reference sentences contain $\omega_k$. A higher *idf* means $\omega_k$ is rare and discriminative, hence, larger weight. CIDEr score for $n$-grams of length $n$ is defined using cosine similarity of *tf-idf* weighting vectors as in (3.12).

$$\text{CIDEr}_n(x, Y) = \frac{1}{|Y|} \sum_{Y_i \in Y} \frac{g^n(x) \cdot g^n(Y_i)}{\|g^n(x)\| \|g^n(Y_i)\|} \tag{3.12}$$

where $x$ is the candidate sentence, $Y = \{Y_1, Y_2, ..., Y_k\}$ is the set of reference sentences, $g^n(a)$ is vector formed by $tf - idf$ weights for all $n$-grams in a sentence $a$ according to equation (3.11). Overall CIDEr score is the mean of $\text{CIDEr}_n$ scores for $n$-grams of length from 1 to 4 ((3.13)).

$$\text{CIDEr} = \frac{1}{4} \sum_{n=1}^{4} \text{CIDEr}_n(x, Y) \tag{3.13}$$

Even though this metric is proposed for image captioning, it does not utilize image and solely based on sentences and the language itself. Hence, *tf-idf* weighting may over-weight/under-weight unnecessary/important $n$-grams for image description and it may be a suboptimal metric for image captioning evaluation [98].

### 3.4.5 SPICE: Semantic Propositional Image Caption Evaluation

SPICE (Semantic Propositional Image Caption Evaluation) [103] is a recently proposed metric for evaluation of image captioning methods. It is claimed that $n$-gram overlap based methods do not take semantic information into account and capturing human judgements requires considering semantic relationships between candidate sentence and reference sentences. Hence, a *scene graph* based image captioning evaluation metric is proposed. A *scene graph* is a graph structure which models *objects*, *attributes* of the objects, and *relationships* between objects in an image [105].

As the first step of evaluation process, a scene graph is constructed via semantic parsing from set of reference ground truth sentences. Semantic parsing of a caption $c$ can be formulated as in (3.14).

$$G(c) = \langle O(c), E(c), K(c) \rangle \tag{3.14}$$

where $O(c)$ is the objects in caption $c$, $E(c)$ defines relationships between two objects in $O(c)$, $K(c)$ defines attributes of objects in $O(c)$.

For a set of reference captions, one common scene graph is constructed. An example for such a construction is given in Figure 3.8.

Once scene graphs for reference sentences and candidate caption are constructed, a function $T$ is used to extract the set of tuples from the scene graphs.

$$T(G(c)) \triangleq O(c) \cup E(c) \cup K(c) \tag{3.15}$$

A tuple is either *(object)*, *(object, attribute)*, or *(object, relationship, object)*. Following tuple extraction, precision $(P)$, recall $(R)$, and SPICE metric are defined as in (3.16) for a candidate caption $x$ and a set of reference captions $Y$.

$$P(x, Y) = \frac{|T(G(x)) \otimes T(G(Y))|}{|T(G(x))|}$$

$$R(x, Y) = \frac{|T(G(x)) \otimes T(G(Y))|}{|T(G(Y))|} \tag{3.16}$$

$$\text{SPICE}(x, Y) = F_1(x, Y) = \frac{2 \cdot P(x, Y) \cdot R(x, Y)}{P(x, Y) + R(x, Y)}$$

where $\otimes$ is binary matching operator which returns matched tuples between two sets of tuples in two scene graphs. While performing matching, similar to [100], syn-

Figure 3.8: Set of reference captions for the given image (left) and corresponding scene graph (right). Pink nodes represent *objects*, green nodes represent *attributes*, cyan nodes represent *relationships* (Retrieved from [103]).

onyms, words with same lemmatized forms, or words belonging to same WordNet ID (synset ID) [106] are considered to be matched.

While semantic relationship between candidate caption and reference captions are exploited, grammatical correctness and fluency of the captions are not measured [103], [21]. Another point is that the performance of this metric is negatively affected if parsing function $G(c)$ produces incorrect scene graphs.

# CHAPTER 4

# PROPOSED METHODOLOGY

In this chapter, methodology followed for the work conducted within the scope of this thesis is discussed. The baseline method VinVL [4] and the proposed novel approach for enriching image captioning results with object attributes are explained. Different training strategies and caption generation process are discussed.

Notations are explained wherever they are used. Additionally, frequently used notations are described in Table 4.1 for convenience.

Table 4.1: Table of notations.

| Symbol | Description |
|---|---|
| $I$ | : Input image |
| $Z$ | : Feature map for input image $I$ |
| $N$ | : Number of objects in input image $I$ |
| $M$ | : Number of masked tag and/or caption tokens |
| $L$ | : Number of object attributes in input image $I$ |
| $Q$ | : Number of masked attribute tokens |
| $T$ | : Length of the ground-truth caption |
| $\alpha$ | : Optimum assignment for object tags |
| $\beta$ | : Optimum assignment for object attributes |
| $C = \{c_0, c_1, \ldots, c_{N-1}\}$ | : Set of object tags |
| $B = \{b_0, b_1, \ldots, b_{N-1}\}$ | : Set of bounding boxes |
| $F = \{f_0, f_1, \ldots, f_{N-1}\}$ | : Set of regional feature vectors |
| $A = \{a_0, a_1, \ldots, a_{L-1}\}$ | : Set of object attributes |
| $S = \{s_0, s_1, \ldots, s_{T-1}\}$ | : Set of tokens in ground-truth caption |
| $Y$ | : Set of auto-regressively obtained input tokens during decoding |
| $H = \{h_0, h_1, \ldots, h_{T+N-1}\}$ | : Set of all maskable tokens for VLP (Union of $S$ and $C$) |
| $\mathcal{D} = \{d_0, d_1, \ldots, d_{M-1}\}$ | : Set of masked tag and/or caption tokens at the transformer input |
| $E = \{e_0, e_1, \ldots, e_{Q-1}\}$ | : Set of masked object attributes at the transformer input |
| $\overline{X}$ | : Set of transformer outputs corresponding to the input set $X$. $X$ can be one of $(C, B, F, S, H, \mathcal{D}, Y, A, E)$ above. |

## 4.1 Baseline Method

Transformer-based [12] architectures gained popularity in NLP and CV domains in the recent years. As discussed in Section 3.2, transformer architecture came into use in image captioning methods as well. Recently, Unified-VLP [30] approach achieved remarkable results by using a unified transformer architecture for Visual Encoder and Language Model blocks. OSCAR [16] further improved Unified-VLP by integrating object tags (class names) in the image to the same transformer architecture. VIVO [1] used the same architecture as OSCAR, and targeted Novel Object Captioning with a novel pretraining methodology on an object detection dataset [3]. Finally, VinVL [4] surpassed the results of OSCAR on general image captioning and the results of VIVO on Novel Object Captioning with the same architecture and training strategy as them by just improving regional feature extraction model. Hence, VinVL is taken as the baseline method. Its architectural details, training strategies, and caption generation process are explained in the next subsections.

### 4.1.1 Architecture

As stated in Section 3.2.1, many methods utilize regional features during caption generation. VinVL [4] also extracts regional features and tags for the objects in the input image. VinVL utilize Faster R-CNN [7] with ResNeXt-152 C4 [41] backbone trained on a collection of datasets (COCO [27], Open Images [3], Objects365 [42], Visual Genome [38]) for regional feature and object tag extraction.

Given an image with size $w \times h$ where $w$ is the width, $h$ is the height of the image, a feature map $Z \in \mathbb{R}^{H_Z \times W_Z \times D_Z}$ is extracted by the backbone CNN where $H_Z, W_Z, D_Z$ are the height, width, and number of channels for the feature map. Feature map $Z$ is fed into classification and regression heads for the bounding box prediction.

The object detection network predicts two output sets for tags (class names) and bounding boxes of the objects: $C = \{c_0, c_1, \ldots, c_{N-1}\}$ and $B = \{b_0, b_1, \ldots, b_{N-1}\}$ where $b_i = \{(x_i^{tl}, y_i^{tl}), (x_i^{br}, y_i^{br})\}$ represents the top-left and bottom-right coordinates of the predicted bounding box and $c_i$ represents the tag, for the $i$-th object in the image, and $N$ is the number of detected objects in the image. Regional feature vector $f_i$

Figure 4.1: Feature and object tag extraction for VinVL [4].

for the $i$-th object is calculated as in (4.1).

$$z_i = RoIPooling(Z, b_i)$$

$$f_i = Concat\left(z_i, \frac{x_i^{tl}}{w}, \frac{y_i^{tl}}{h}, \frac{x_i^{br}}{w}, \frac{y_i^{br}}{h}, \frac{x_i^{br} - x_i^{tl}}{w}, \frac{y_i^{br} - y_i^{tl}}{h}\right) \tag{4.1}$$

where *RoIPooling* [37] operation generates fixed-sized intermediate feature vector $z_i \in \mathbb{R}^{D_z}$ from projection of $b_i$ on feature map $Z$ (as shown in the middle in Figure 4.1) by applying max pooling. *Concat* operation combines intermediate feature vector with positional information of the bounding box (coordinates, width, and height) by applying concatenation in order to obtain final regional feature vector $f_i \in \mathbb{R}^{D_z+4}$. $D_z$ is the size of intermediate feature vectors, $z_i$. It is equal to 2048 in VinVL. Hence, the dimension of regional feature vectors ($f_i$) is 2054.

The process of extracting regional features and object tags is illustrated in Figure 4.1. After extraction, object tags $C = \{c_0, c_1, \dots, c_{N-1}\}$ and regional feature vectors $F = \{f_0, f_1, \dots, f_{N-1}\}$ for the detected objects in the image are fed into the transformer block in the next step.

Transformer block acts as a unified architecture for visual encoding and language modelling. It takes inputs from both visual and textual modalities and performs self-attention. For the transformer block, uncased version of the $\text{BERT}_{\text{base}}$ model [13]

41

is used. Hence, the model does not differentiate between lower-case and upper-case letters. Tokenizer, positional embedding, and segment embedding of the BERT are directly utilized. Moreover, pretrained weights of the $\text{BERT}_{\text{base}}$ trained on large English corpora with Masked Language Modelling (MLM) objective are used, and the model is further trained for captioning. Details of transformers and BERT are provided in Section 2.2 and Section 2.2.1, respectively.

Let $S = \{s_0, s_1, \ldots, s_{T-1}\}$ be the ground truth caption of length $T$ for an image and $s_i$ be the $i$-th (tokenized) word in $S$. Transformer block takes ground truth caption $S$, outputs of the object detector $F$ (regional features), and object tags $C$ as input and applies self-attention according to (4.2).

$$\overline{S}, \overline{C}, \overline{F} = \text{Transformer}(S, C, F) \tag{4.2}$$

where sets $\overline{S}, \overline{C}, \overline{F}$ are transformed representations of sets $S, C$, and $F$, respectively. Transformer block projects all inputs ($s_i \in S, c_i \in C, f_i \in F$) to a common embedding dimension $D_E$. Transformer block exploits the relationships between caption, tag, and feature inputs by applying self-attention and outputs the transformed representations of the inputs. This operation is illustrated in in Figure 4.2 for the example provided in Figure 4.1.
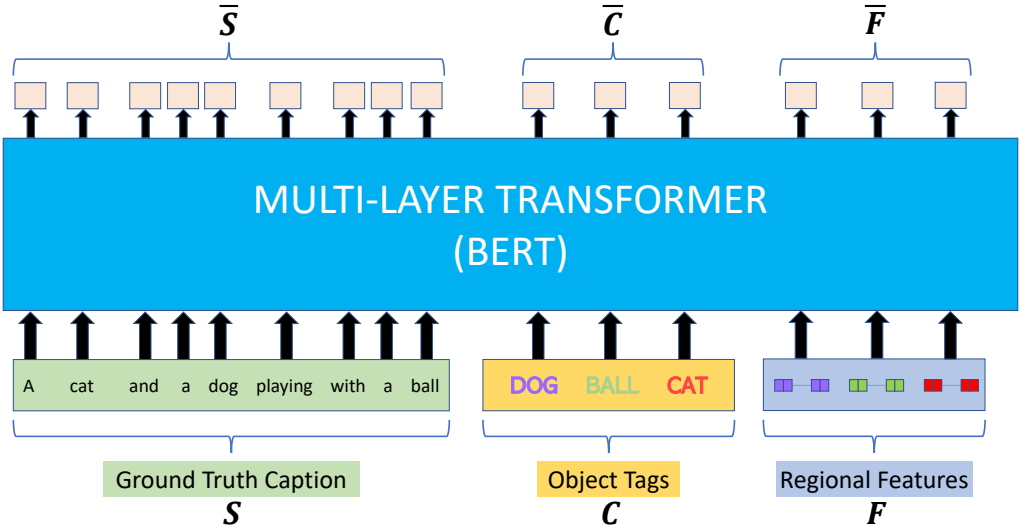


Figure 4.2: Transformer (Self-Attention) block for VinVL [4] for the example provided in Figure 4.1.

### 4.1.2 Training Strategies

The feature and object tag extraction method and self-attention method explained in Section 4.1.1 are used with following training strategies with different objective functions and inputs:

- Vision and Language Pretraining (VLP) [16], [4]

- Visual Vocabulary Pretraining for Novel Object Captioning (VIVO) [1], [4]

- Finetuning for Image Captioning [16], [4]

Details for these strategies are explained in the subsequent sections.

### 4.1.2.1 Vision and Language Pretraining (VLP)

Vision and Language Pretraining (VLP) refers to training a model on image-text paired datasets [30]. Pretrained models are further finetuned for downstream tasks (image captioning, visual question answering) on task-specific datasets. Usually, pretraining is performed on several large image-text paired datasets. Aim in VLP is modelling and establishing the relationship between visual and textual modalities. Datasets used in VLP might belong to different tasks. In VinVL, they use datasets belonging to 3 different tasks. For each of these tasks, the transformer block takes 3 input sets, $S, C,$ and $F$. These tasks and corresponding inputs are provided in Table 4.2.

Combination of all image-text paired datasets provided in Table 4.2 are used during VLP. $S, C,$ and $F$ inputs are fed into the transformer along with special $[CLS]$ and $[SEP]$ tokens in BERT. Input sequence starts with $[CLS]$ token followed by sets $S, C,$ and $F$ where each input set is separated with $[SEP]$ token, as illustrated in Figure 4.3.

VinVL is trained with two objectives during VLP: *Masked Token Loss ($\mathcal{L}_{MTL}$)* and *3-way Contrastive Loss ($\mathcal{L}_{CL3}$)*. Total loss is the sum of these two losses as in (4.3).

$$\mathcal{L}_{total} = \mathcal{L}_{MTL} + \mathcal{L}_{CL3} \tag{4.3}$$

Table 4.2: Different tasks, corresponding datasets and inputs to the transformer block for VLP.

| Task | Datasets | Inputs to Transformer | | |
| --- | --- | --- | --- | --- |
| | | $S$ | $C$ | $F$ |
| Image Captioning | COCO [27], CC [28], SBU [29], Flickr30k [2] | Ground Truth Caption | Machine-Generated Object Tags | Regional Features |
| Visual Question Answering | GQA [107], VQA [108], VG [38] | Question | Answer | Regional Features |
| Image Tagging | OI [3] | Machine-Generated Caption | Ground Truth Object Tags | Regional Features |



Figure 4.3: Input sequence to transformer for VinVL [4] for the example provided in Figure 4.1.

*Masked Token Loss ($\mathcal{L}_{MTL}$)* is first proposed in OSCAR [16]. The idea is similar to MLM in BERT [13]. Words belonging to sets $S$ and $C$ are randomly replaced with the special $[MASK]$ token and the sequence is fed into the transformer. The objective is predicting the masked word by classifying the corresponding output of the transformer for the $[MASK]$ token. The goal of $\mathcal{L}_{MTL}$ is forcing the network deduce the masked word from its surrounding words and regional features. This helps constructing the alignment between textual and visual modalities [16].

Let $H = S \cup C$ be the union of sets $S$ and $C$. A word *cat* is masked with 15% probability where $h_i \in H$ is the token corresponding to the word *cat*. Following procedure is applied if token $h_i$ is to be masked:

- Replace the $h_i$ with $[MASK]$ token 80% of the time.

- Replace the $h_i$ with a random word in the vocabulary 10% of the time.

- Keep $h_i$ as it is 10% of the time.

In all three cases, $h_i$ is treated as a masked token and the network should classify transformer output corresponding to $h_i$ as *cat*. This idea is inspired from BERT and it helps generalization and learning contextual representation for each input token [13].

Cross-Entropy loss is used in $\mathcal{L}_{MTL}$. Let $\mathcal{D}$ be the set of masked words in $H$ with size $|\mathcal{D}| = M$. $\mathcal{L}_{MTL}$ is calculated according to (4.4) for a single training example (image-text pair).

$$\mathcal{L}_{MTL} = -\frac{1}{M} \sum_{h_i \in \mathcal{D}} \sum_{j}^{V} y_i^j \log p(\bar{h}_i^j) \tag{4.4}$$

where $V$ is the size of the vocabulary (number of output classes, equals to 30522 in all experiments in this thesis, following BERT [13]), $p(\bar{h}_i^j)$ is the predicted output probability for output class $j$ for the masked token $h_i$, and $y_i^j$ is the binary ground truth label for class $j$ for the masked word $h_i$. VLP with $\mathcal{L}_{MTL}$ is illustrated in Figure 4.4.

*3-way Contrastive Loss ($\mathcal{L}_{CL3}$)* also aims to align textual and visual modalities. In $\mathcal{L}_{CL3}$ objective, a single word in $S$ or $C$ is randomly replaced (*polluted*) with another word with probability 0.5 and the sequence is fed into the transformer. The objective for the network is predicting whether the input sequence is polluted, if it is, which set ($S$ or $C$) is polluted. It is a 3-way classification task with 3 classes: *unpolluted*, *S-polluted*, and *C-polluted*. The classification is performed at the transformer output corresponding to $[CLS]$ input token since it is viewed as an encoded representation for the whole input $(S, C, F)$.

Cross-Entropy loss is used in $\mathcal{L}_{CL3}$. $\mathcal{L}_{CL3}$ is calculated according to (4.5) for a single training example (image-text pair).

$$\mathcal{L}_{CL3} = -\sum_{j=0}^{2} y_{cls}^j \log p(\bar{h}_{cls}^j) \tag{4.5}$$

where $p(\bar{h}_{cls}^j)$ is the predicted output probability for output class $j$ for the $[CLS]$ token $h_{cls}$, and $y_{cls}^j$ is the binary ground truth label for class $j$ for the input sequence. There

Figure 4.4: Illustration of training with Masked Token Loss ($\mathcal{L}_{MTL}$) for VinVL [4] for the example provided in Figure 4.1.

are 3 possible output classes: *unpolluted* ($y_{cls}^0 = 1$), *S-polluted* ($y_{cls}^1 = 1$), and *C-polluted* ($y_{cls}^2 = 1$). VLP with $\mathcal{L}_{CL3}$ is illustrated in Figure 4.5 for an *S-polluted* input sequence.

As seen in Figure 4.4 and Figure 4.5, full (square) attention mask is used for VLP with $\mathcal{L}_{MTL}$ and $\mathcal{L}_{CL3}$. So, all tokens are allowed to attend all other tokens. Since aim is modelling the relationship between visual and textual modalities, masked tokens can attend to tokens for the future words in the sequence as well.

The model pretrained with VLP is finetuned on an image-caption paired dataset for image captioning as will be described in Section 4.1.2.3.

Figure 4.5: Illustration of training with 3-way Contrastive Loss ($\mathcal{L}_{CL3}$) for VinVL [4] for the example provided in Figure 4.1. Input sequence is $S$-polluted with the replacement word *SKY*.

#### 4.1.2.2 Visual Vocabulary Pretraining (VIVO)

Visual Vocabulary Pretraining (VIVO) [1] is a method proposed for Novel Object Captioning and *nocaps* challenge [91]. The same model architecture described in Section 4.1.1 is used. In VIVO, Novel Object Captioning problem is broken into two stages. In the first stage, the model is pretrained on a large object detection dataset, Open Images [3]. In the second stage, the pretrained network is finetuned for image captioning. The idea behind this methodology is that object detection datasets are abundant compared to image-caption paired datasets. An image captioning model should be able to learn objects from an object detection dataset and mention these objects in generated captions. Hence, in the pretraining stage, the model is forced to learn a *visual vocabulary* for the novel objects which models a joint embedding space.

47

The model maps tags for similar objects and their visual features to vectors which are closer to each other in this embedding space. In the second stage, the model is expected to learn caption generation while also exploiting knowledge from pretraining in order to generate captions for images with novel objects. Decoupling visual vocabulary learning and caption generation learning allows utilizing large amount of image-tag paired data and generalization to novel objects which are unseen in the ground truth captions of the image captioning datasets.

In VIVO pretraining, the model is fed with two sets of inputs $C$ and $F$, where $C$ is the set of object tags and $F$ is the set of regional features for the objects in the image. Object tags are extracted with an object detection model. In VIVO, an object detection model trained on Open Images dataset is used. A multi-label classification network [33] trained on Open Images dataset can also be used as a tag extractor, as we did in this thesis. The regional features are extracted as described in Section 4.1.2.1. An example image from the object detection dataset and corresponding regional features and tags are shown in Figure 4.6.
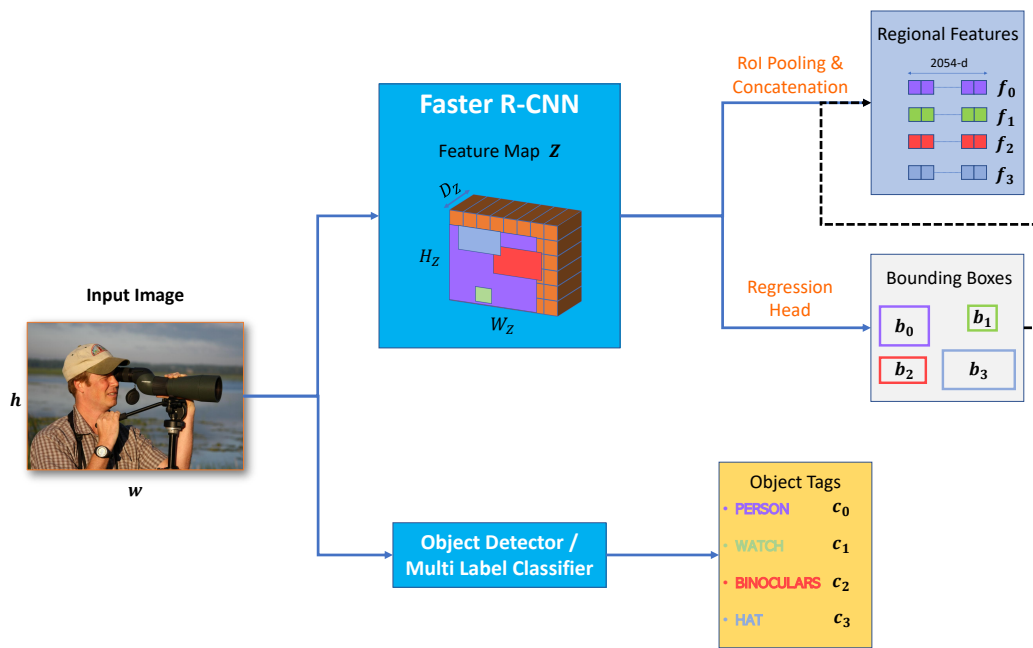


Figure 4.6: Feature and object tag extraction for VIVO.

Similar to BERT [13] and Unified-VLP [30], masking is performed for object tags. The idea is similar to Masked Token Loss ($\mathcal{L}_{MTL}$) described in Section 4.1.2.1. The training objective is predicting masked tags at the transformer output by attending to other object tags and regional features.

Let $\mathcal{D} = \{d_0, d_1, \ldots, d_{M-1}\}$ be the set of tokens for masked object tags in $C \supseteq \mathcal{D}$ and $\overline{\mathcal{D}} = \{\bar{d}_0, \bar{d}_1, \ldots, \bar{d}_{M-1}\}$ be the transformer outputs corresponding to elements in $\mathcal{D}$ where $M$ is the number of masked tags. A token $c_i \in C$ is masked with 15% probability where $c_i$ is the token corresponding to the masked tag, for example *binoculars*. Following procedure is applied if token $c_i$ is to be masked:

- Replace the $c_i$ with $[MASK]$ token 80% of the time.

- Replace the $c_i$ with a random word in the vocabulary 10% of the time.

- Keep $c_i$ as it is 10% of the time.

In all three cases, $c_i$ is treated as a masked token and the network should classify transformer output corresponding to $c_i$ as *binoculars*.

An illustration of VIVO pretraining is given in Figure 4.7.

As seen in Figure 4.7, there may be multiple masked tags for an image. The model should predict masked tags at the output of transformer. Different from captions, object tags do not have a notion of order. In fact, object tags for an image could be given to the network in any order. This nature of object tags create ambiguity during prediction of masked words. The model can predict both *person* and *binoculars* at either of transformer outputs, $\bar{d}_0$ and $\bar{d}_1$, without enforcing original order of input tags. To resolve this ambiguity, network outputs should be assigned to masked tags, and each output should be responsible (contribute to the loss) from its assigned tag.

A one-to-one assignment, $\alpha$, between network inputs and outputs is necessary such that each output $\bar{d}_i \in \overline{\mathcal{D}}$ should be responsible from a unique input tag $d_j \in \mathcal{D}$. This assignment should be the most optimum one in terms of the overall training loss. To obtain the optimal assignment $\alpha$, *Hungarian Assignment Algorithm* [48] is used. Details of the Hungarian Assignment Algorithm is provided in Section 2.3.

Figure 4.7: Illustration of VIVO pretraining.

Let $t_j$ be the vocabulary (token) id for the masked input $d_j$. The square cost matrix $J \in \mathbb{R}^{M \times M}$ for this assignment is constructed such that

$$J_{ij} = 1 - p(\bar{d}_i^{t_j}) \tag{4.6}$$

where $p(\bar{d}_i^{t_j})$ is the probability of transformer output $\bar{d}_i$ being the word with vocabulary id $t_j$. $\alpha$ assigns an input index $i$ to an output index $j$ as in (4.7).

$$\alpha(i) = j \quad i, j \in [0, M - 1] \tag{4.7}$$

Hungarian Assignment Algorithm finds $\alpha$ such that cost function in (4.8) is minimized globally for all masked words.

$$J_\alpha = \sum_{i=0}^{M-1} 1 - p(\bar{d}_i^{\alpha(i)}) \tag{4.8}$$

After the assignment is obtained, $\mathcal{L}_{MTL}$ objective is applied similar to VLP. In $\mathcal{L}_{MTL}$, Cross-Entropy loss is calculated using the assigned pairs according to (4.9).

$$\mathcal{L}_{MTL} = -\frac{1}{M} \sum_{d_i \in \mathcal{D}} \sum_{j}^{V} y_i^j \log p(\bar{d}_i^j) \tag{4.9}$$

where $V$ is the size of the vocabulary (number of output classes), $p(\bar{d}_i^j)$ is the predicted output probability for output class $j$ for the masked token $d_i$. $y_i^j$ is the binary ground truth label for class $j$ for the masked word $d_i$ and it is related to assignment $\alpha$ as in (4.10).

$$y_i^j = \begin{cases} 1, & \text{if } \alpha(i) = j \\ 0, & \text{otherwise} \end{cases} \tag{4.10}$$

As illustrated in 4.7, full attention mask is used during VIVO pretraining, similar to VLP. All object tags and regional features can attend to each other since the objective is learning the *visual vocabulary*.

The model pretrained with VIVO is finetuned on COCO [27] captioning dataset for image captioning as will be described in Section 4.1.2.3.

### 4.1.2.3 Finetuning for Image Captioning

After the model is trained with VLP (Section 4.1.2.1) or VIVO (Section 4.1.2.2), it is finetuned on an image-caption paired dataset, COCO [27], for caption generation.

Similar to VLP and VIVO, two stage pipeline is used in finetuning.

In the first stage, object tags and regional features for the images in the captioning dataset are extracted. First stage is illustrated in Figure 4.8.

During VLP for VinVL, object tags and regional features are extracted using the same Faster R-CNN [7] network. On the other hand, for VIVO pretraining, VinVL utilize external tag prediction network trained on Open Images [3] dataset for detection of novel objects. Same methodology is followed during finetuning as well. Hence, Faster R-CNN and Object Detector models in Figure 4.8 are the same for general image captioning whereas for Novel Object Captioning, external object detection model is utilized.
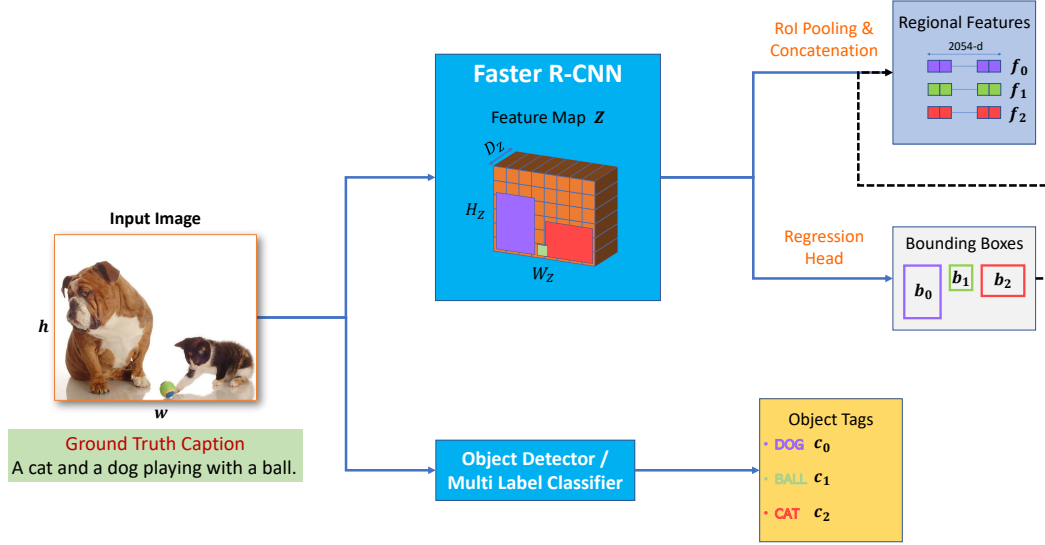
Figure 4.8: Illustration of object tag and regional feature extraction for image captioning finetuning.

In the second stage, regional features, $F$, object tags, $C$, and ground truth caption, $S$, are fed into the transformer network for image captioning finetuning. Similar to BERT [13], VLP, and VIVO, finetuning is performed with $\mathcal{L}_{MTL}$ objective. Different from VLP and VIVO, masking is performed on captions only. Object tags are not masked.

Let $\mathcal{D} = \{d_0, d_1, \ldots, d_{M-1}\}$ be the set of tokens for masked words in $S \supseteq D$ and $\overline{\mathcal{D}} = \{\bar{d}_0, \bar{d}_1, \ldots, \bar{d}_{M-1}\}$ be the transformer outputs corresponding to elements in $\mathcal{D}$ where $M$ is the number of masked words in $S$. A token $s_i \in S$ is masked with 15% probability where $s_i$ is the token corresponding to the masked word in the ground truth caption, for example *playing*. Following procedure is applied if token $s_i$ is to be masked:

- Replace the $s_i$ with $[MASK]$ token 80% of the time.

- Replace the $s_i$ with a random word in the vocabulary 10% of the time.

- Keep $s_i$ as it is 10% of the time.

In all three cases, $s_i$ is treated as a masked token and the network should classify transformer output corresponding to $s_i$ as *playing*.

The network is trained with $\mathcal{L}_{MTL}$ objective. Cross-Entropy loss is used in $\mathcal{L}_{MTL}$. $\mathcal{L}_{MTL}$ is calculated according to (4.11).

$$\mathcal{L}_{MTL} = -\frac{1}{M} \sum_{d_i \in \mathcal{D}} \sum_{j}^{V} y_i^j \log p(\bar{d}_i^j) \tag{4.11}$$

where $V$ is the size of the vocabulary (number of output classes), $p(\bar{d}_i^j)$ is the predicted output probability for output class $j$ for the masked token $d_i$. $y_i^j$ is the binary ground truth label for class $j$ for the masked word $d_i$.

Finetuning for image captioning is illustrated in Figure 4.9 for the example in Figure 4.8.
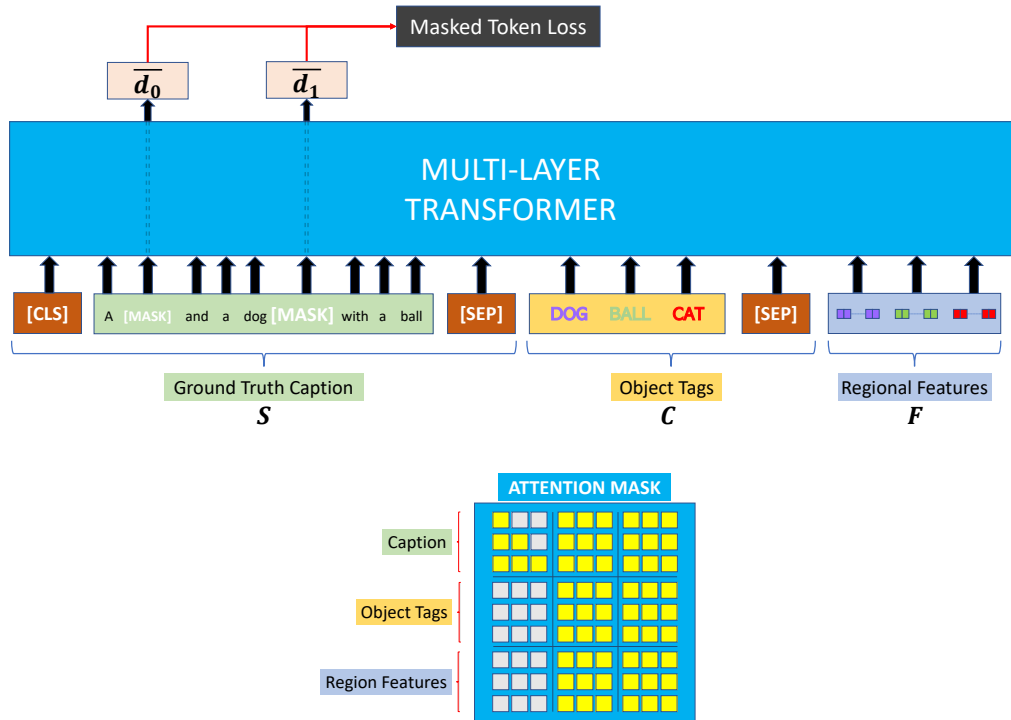


Figure 4.9: Illustration of finetuning for image captioning for the example in Figure 4.8.

Different from VLP and VIVO, unidirectional attention mask is employed for tokens belonging to the ground truth caption, $S$, as shown in Figure 4.9. Hence, a token $s_i \in S$ can attend to tokens $s_j \in S$ for $j < i$ in the ground truth caption, $S$, all object

53

tags, $C$, and all regional features, $F$. Furthermore, tokens belonging to object tags and regional features cannot attend to tokens of the ground truth caption. Otherwise, this would violate unidirectional attention for tokens in $S$ due to transitive property of attention.

Unidirectional attention mask is necessary for sequence generation tasks. Because the model should adapt to the sequential nature of the language. Decoding (inference) is performed in an auto-regressive manner as described in Section 2.2.1 and Section 4.1.3. During caption generation, tokens for the future time instants are not available. Hence, the model can only attend to previous time instants, object tags, and regional features. It is forced to predict next word based on previous words, object tags, and regional features. A full attention mask would be cheating during training and it would create discrepancy between training and inference. Thus, unidirectional attention mask is employed during finetuning for image captioning as well.

So, the network acts as a language model. It models the probability of a caption given previous words, object tags, and regional features as in (4.12).

$$P(y_{N-1}, y_{N-2}, ..., y_0 \mid \boldsymbol{C}, \boldsymbol{F}) = \prod_{i=0}^{N-1} P(y_i \mid y_0, y_1, ..., y_{i-1}, \boldsymbol{C}, \boldsymbol{F}) \qquad (4.12)$$

where $\boldsymbol{C}$ is the set of object tags, $\boldsymbol{F}$ is the set of regional features, $y_i$ is the word in the caption at the $i$-th time instant, and $N$ is the length of the caption, $S$.

### 4.1.3 Caption Generation (Decoding)

After the model is finetuned on image captioning dataset, it can be used for captioning inference. The process of generating captions for an input image is referred as *decoding*.

Object tags and regional features are extracted as in Figure 4.8.

The structure of inputs to the transformer is similar to the one in finetuning (Section 4.1.2.3. There are 3 sets of inputs: Auto-regressive input tokens - $Y$, object tags - $C$, regional features - $F$. Object tags and regional features are the same as the ones in finetuning. Different from finetuning, there are no ground truth captions. Hence,

initially the set for input tokens contain $T$ $[MASK]$ tokens where $T$ is the maximum number of words for the generated caption.

The auto-regressive decoding procedure is provided in Algorithm 1. Forward propagation is applied using the given inputs. The transformer output $\bar{d}_t$ at time instant $t$ (corresponding to the $t$-th masked input token) is fed into the classifier and the network predicts the word $\bar{y}_t$. In the next time instant $t+1$, the $t$-th masked word is replaced with previously predicted word $\bar{y}_t$. After forward pass, the network makes prediction for $\bar{y}_{t+1}$. This process is followed in an auto-regressive manner until maximum caption length is achieved or the network predicts a special end-of-sentence token, $[EOS]$, at the output.

---

**Algorithm 1** Auto-regressive decoding process.

---

$Y = \{[CLS], [MASK], [MASK], \ldots, [MASK], [SEP]\}, \quad |Y| = T + 2$

$\overline{Y} = \{\}$                           $\triangleright$ Set of Words for the Generated Caption

$t = 0$

**while** $|\overline{Y}| \neq T$ **and** $\bar{y}_t \neq [EOS]$ **do**        $\triangleright$ $[EOS]$: End of Sentence Token

    $\bar{y}_t \leftarrow \text{Model}(Y, C, F)$

    $\overline{Y} \leftarrow \overline{Y} \cup \{\bar{y}_t\}$

    $y_t \leftarrow \bar{y}_t$

    $t \leftarrow t + 1$

**end while**

---

An example decoding snapshot for the example in Figure 4.8 at time instant $t = 5$ is given in Figure 4.10.

Similar to finetuning stage, unidirectional attention mask is utilized since the model has no information about the future words in the caption. It is expected for the model to predict current masked word by attending to previously predicted words, object tags, and regional features.

There are different decoding techniques utilized during caption generation. The most straight forward one is picking the word with highest probability at the model output. This is referred as *greedy decoding*. Another widely-used approach is called *beam search* [35]. In beam search, at each time instant, top-$n$ most probable sentences
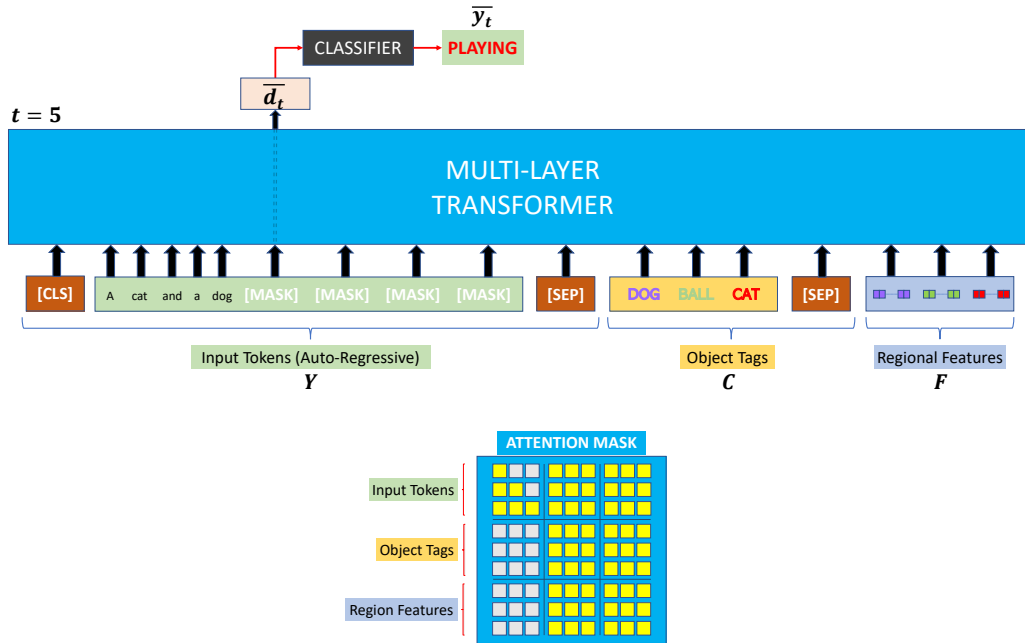
Figure 4.10: Illustration of auto-regressive decoding (caption generation).

are selected and the model is fed with each of these sentences separately in the next time instant where $n$ is called *beam width*. Probability of a sentence is calculated by multiplying individual classification probabilities of words in the sentence. Since the beam search considers probability of a sentence as a whole, it generally produces better captions compared to the greedy decoding which only focuses on the independent word probabilities. We use beam width of 5 in our experiments.

## 4.2  Proposed Method

Even though recent image captioning methods usually produce high-quality captions, they may overlook details in the image. There are studies which claim that current image captioning models do not fully refer to the input image while generating captions. Instead, they are prone to biases in the training dataset and language model [31], [32]. In [31], it is shown that methods disregard attributes of objects in the image. Furthermore, during testing, they copy frequent phrases in the training dataset instead of fully referring to the input image. Hence, they might ignore details in the image which de-

creases the descriptive capability of the caption. An example poor caption generated by the baseline method is given in Figure 4.11 where the model hallucinates a *chair* and fails to mention about the *car* in the background. Furthermore, the caption lacks properties of the objects in the scene.



A garden with a chair and a plant on the ground.

Figure 4.11: An example image and a poor caption generated by VinVL [4].

Recent methods try to overcome this issue by introducing additional information (object attributes, relationships) to the captioning model or explicitly modelling visual grounding for the predicted words [31], [32]. In conjunction with these methods, we propose to exploit object attributes in order to produce richer and more descriptive captions. Different from existing approaches, we apply this methodology to a state-of-the-art *self-attention*-based baseline, VinVL [4] which is described in Section 4.1.

The architecture for the proposed approach is the same as VinVL, which is described

in Section 4.1.1. We employ $\text{BERT}_{\text{base}}$ transformer model similar to VinVL.

The training objectives are also similar to that of VinVL, which are described in Section 4.1.2. Due to computational limitations, we demonstrate our results following Novel Object Captioning setup since training datasets are restricted in *nocaps* [91] challenge as opposed to general image captioning in which several large datasets are utilized. However, proposed approach is also applicable to general image captioning. We provide experimental results on *nocaps* challenge and COCO [27] test set in Chapter 5.

Similar to VinVL, we use pretrained weights of the $\text{BERT}_{\text{base}}$ trained on large English corpora with Masked Language Modelling (MLM) objective. We employ two-stage training pipeline where in the first stage VIVO pretraining (Section 4.1.2.2) is applied on Open Images [3] dataset and in the second stage image captioning finetuning on COCO [27] dataset is performed. We exploit object attributes in both stages.

As the first step in both pretraining and finetuning, we extract a set of object attributes along with object tags and regional features for a given input image. For novel object tag extraction we use a state-of-the-art multi-label classification model, ML-Decoder [33], trained on Open Images [3] dataset. Details of this model is discussed in Section 2.1.1. For object tag and regional feature extraction we use the Faster R-CNN model in VinVL. This model is finetuned on Visual Genome [38] with extra attribute prediction head. Number of output attribute classes is 524 [4]. The model predicts a set of attributes $M_i = \{m_{i0}, m_{i1}, \ldots, m_{iG_i-1}\}$ for each detected object $o_i$ in the image where $G_i$ is the number of predicted attributes for the object $o_i$. An example image with detected objects and their attributes is given in Figure 4.12.

Given an input image $I$, we obtain the set of overall object attributes $A = \{a_0, a_1, \ldots, a_{L-1}\}$ where $L$ is the number of attributes for the input image $I$. We only add maximum of 1 attribute from $M_i$ to the overall set $A$. For an object $o_i$ in $I$, the attribute with the highest confidence in $M_i$ which is also not already available in the overall set $A$ is added to $A$. Hence, the attributes in $A$ are unique. This process is described in Algorithm 2.

Object tag extraction, regional feature extraction, and object attribute extraction are
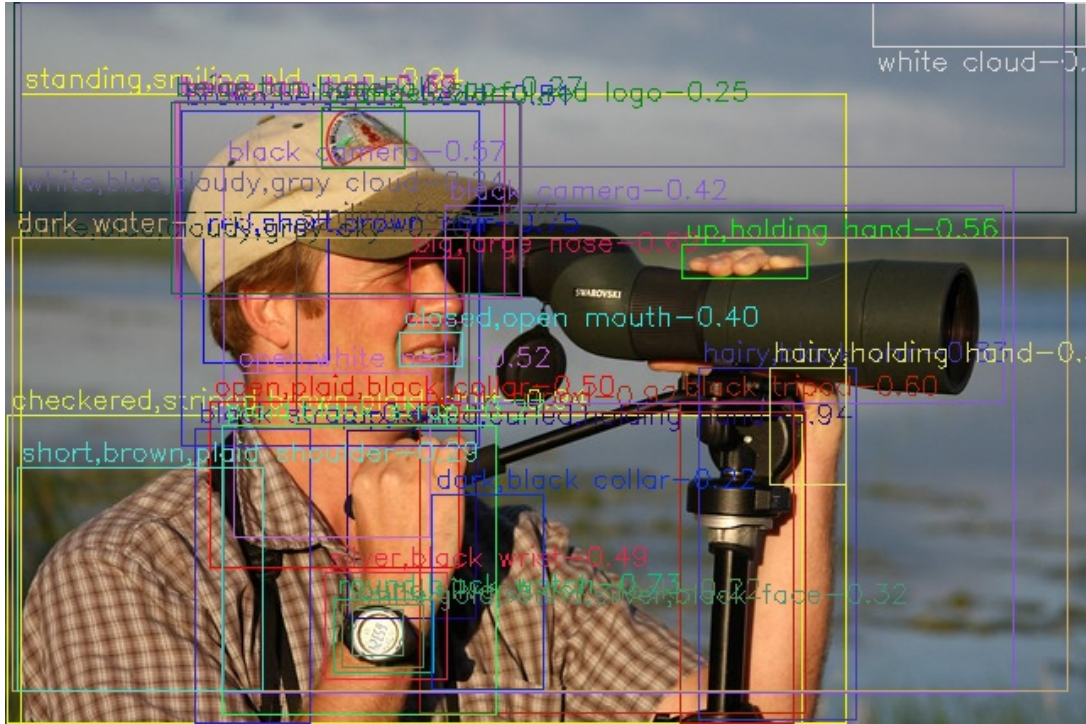
Figure 4.12: An example object-attribute detection result.

---

**Algorithm 2** Attribute selection process.

---

$M = \{M_0, M_1, \ldots, M_{N-1}\}$            $\triangleright$ N: Number of detected objects

$M_i = \{m_{i0}, m_{i1}, \ldots, m_{iG_i-1}\}$      $\triangleright$ $G_i$: Number of attributes for the object $o_i$

$M_i^C = \{m_{i0}^C, m_{i1}^C, \ldots, m_{iG_i-1}^C\}$      $\triangleright$ $m_{ij}^C$: Confidence score for attribute $m_{ij}$

$A = \{\}$                      $\triangleright$ Set of Object Attributes

$k = 0$

**for** $i \leftarrow 0$ **to** $N - 1$ **do**

     $j_{max} \leftarrow \arg \max M_i^C$

     $m_{i_{max}} \leftarrow m_{ij_{max}}$

     **if** $m_{i_{max}} \notin A$ **then**

         $a_k \leftarrow m_{i_{max}}$

         $A \leftarrow A \cup \{a_k\}$

         $k \leftarrow k + 1$

     **end if**

**end for**
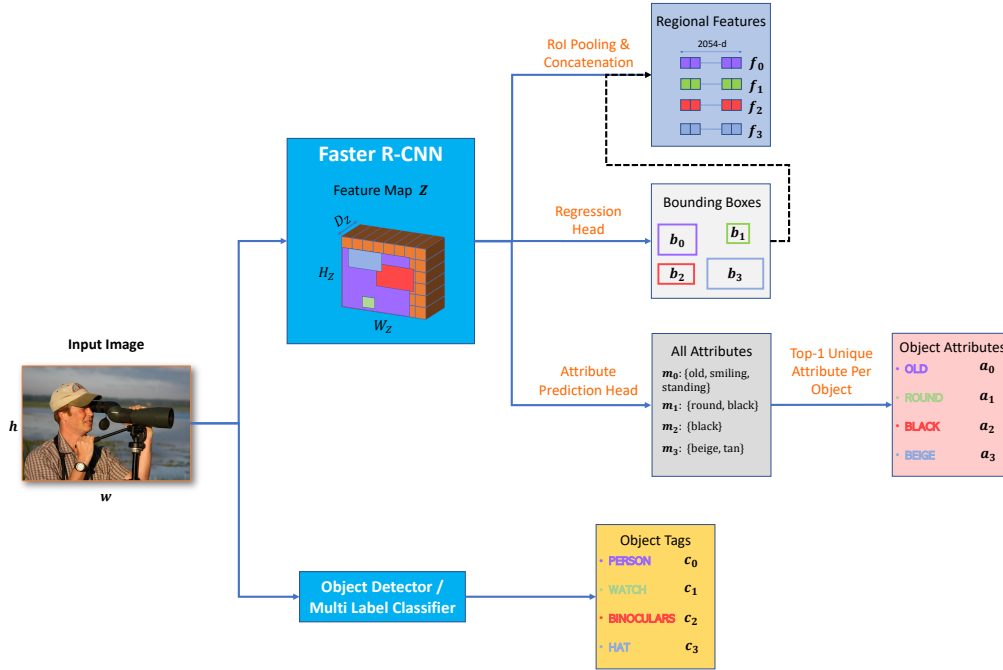
---

illustrated in Figure 4.13.



Figure 4.13: Feature, tag, and attribute extraction for the proposed method for an image to be used in pretraining.

After obtaining input sets, $C, A,$ and $F$ for object tags, object attributes, and regional features, we apply VIVO pretraining. Object attributes are fed into the transformer following object tags with an additional $[SEP]$ token. Different from the baseline (Section 4.1.2.2), masking is performed for object attributes in addition to object tags.

Let $\mathcal{D} = \{d_0, d_1, \ldots, d_{M-1}\}$ be the set of tokens for masked object tags in $C \supseteq \mathcal{D}$ and $\overline{\mathcal{D}} = \{\bar{d}_0, \bar{d}_1, \ldots, \bar{d}_{M-1}\}$ be the transformer outputs corresponding to elements in $\mathcal{D}$. Let $E = \{e_0, e_1, \ldots, e_{Q-1}\}$ be the set of tokens for masked object attributes in $A \supseteq E$ and $\overline{E} = \{\bar{e}_0, \bar{e}_1, \ldots, \bar{e}_{Q-1}\}$ be the transformer outputs corresponding to elements in $E$ where $M$, and $Q$ are the number of masked tags and the number of masked attributes, respectively.

A token $c_i \in C$ is masked with 15% probability where $c_i$ is the token corresponding to the masked tag, for example *binoculars*. Following procedure is applied if token $c_i$ is to be masked:

- Replace the $c_i$ with $[MASK]$ token 80% of the time.

- Replace the $c_i$ with a random word in the vocabulary 10% of the time.

- Keep $c_i$ as it is 10% of the time.

In all three cases, $c_i$ is treated as a masked token and the network should classify transformer output corresponding to $c_i$ as *binoculars*.

Similarly, a token $a_i \in A$ is masked with 15% probability where $a_i$ is the token corresponding to the masked attribute, for example *round*. Following procedure is applied if token $a_i$ is to be masked:

- Replace the $a_i$ with $[MASK]$ token 80% of the time.

- Replace the $a_i$ with a random word in the vocabulary 10% of the time.

- Keep $a_i$ as it is 10% of the time.

In all three cases, $a_i$ is treated as a masked token and the network should classify transformer output corresponding to $a_i$ as *round*.

Illustration of VIVO-like pretraining apporach with object attributes is given in Figure 4.14.

Similar to Section 4.1.2.2, a one-to-one assignment is necessary between masked input tokens (tags or attributes) and network outputs for masked tokens. The reason for this necessity is that tokens for object tags and attributes do not have a notion of order. In fact, these tokens could be given to the network in any order. This nature of object tags and object attributes create ambiguity during prediction of masked words. For example, considering Figure 4.14, the model can predict both masked attributes *round* and *black* at either of transformer outputs, $\bar{e}_0$ and $\bar{e}_1$, without enforcing original order of input attributes. The same ambiguity exists for masked object tags as well. To resolve this ambiguity, network outputs should be assigned to masked tokens, and each output should be responsible (contribute to the loss) from its assigned token.

Since object tags and object attributes belong to different parts of speech, they are treated separately. Two optimal one-to-one assignments are obtained using Hungarian
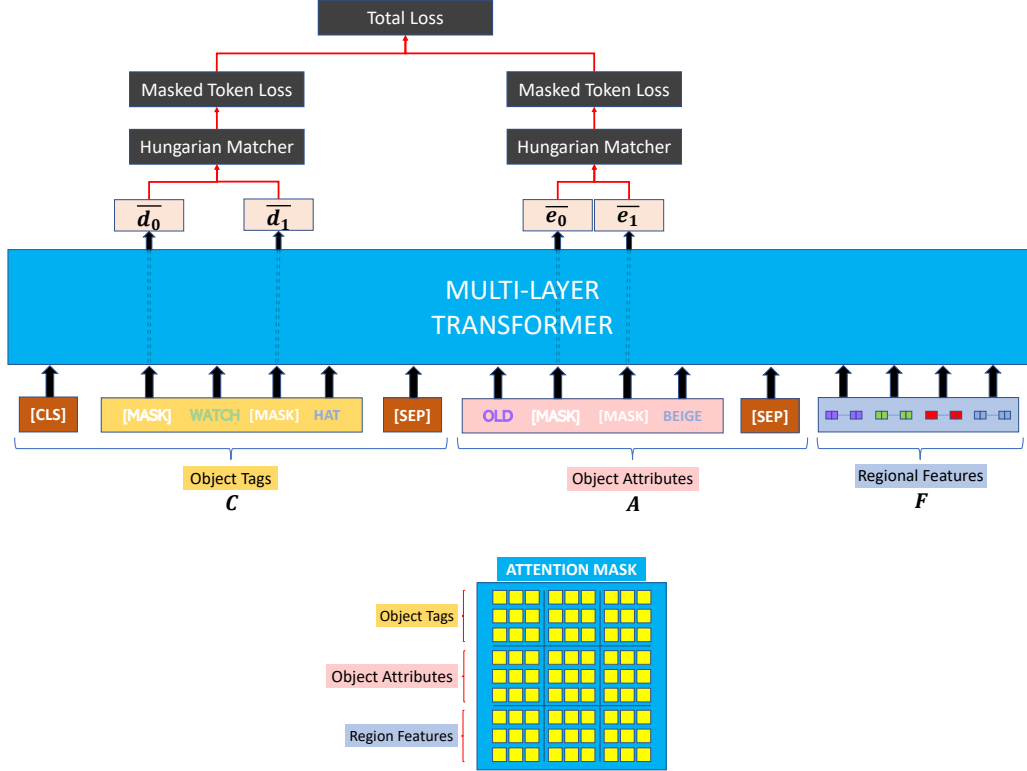
Figure 4.14: Proposed VIVO-like pretraining approach with object attributes.

Assignment Algorithm [48]: $\alpha$ and $\beta$. $\alpha$ assigns masked object tags to transformer outputs for these tags and $\beta$ assigns masked object attributes to transformer outputs for these attributes. To obtain $\alpha$ and $\beta$, Hungarian Assignment Algorithm is solved for two square cost matrices $J^\alpha \in \mathbb{R}^{M \times M}$ and $J^\beta \in \mathbb{R}^{Q \times Q}$, respectively for $\alpha$ and $\beta$. These cost matrices are constructed according to (4.13)

$$
\begin{aligned}
J_{ij}^\alpha &= 1 - p(\bar{d}_i^{t_j}) \\
J_{mn}^\beta &= 1 - p(\bar{e}_m^{t_n})
\end{aligned}
\tag{4.13}
$$

where $p(\bar{d}_i^{t_j})$ and $p(\bar{e}_m^{t_n})$ are the probabilities of transformer outputs $\bar{d}_i$ and $\bar{e}_m$ being the words with vocabulary id $t_j$ and $t_n$, for masked tag and attribute, respectively. $\alpha$ and $\beta$ assigns input indexes to output indexes as in (4.14).

$$
\begin{aligned}
\alpha(i) &= j & i,j \in [0, M-1] \\
\beta(m) &= n & m,n \in [0, Q-1]
\end{aligned}
\tag{4.14}
$$

Hungarian Assignment Algorithm finds $\alpha$ and $\beta$ such that cost functions in (4.15) are

minimized.

$$J_\alpha = \sum_{i=0}^{M-1} 1 - p(\bar{d}_i^{\alpha(i)})$$

$$J_\beta = \sum_{m=0}^{Q-1} 1 - p(\bar{e}_m^{\alpha(m)}) \tag{4.15}$$

After the assignments are obtained, the network is trained with Masked Token Loss ($\mathcal{L}_{MTL}$). $\mathcal{L}_{MTL}$ for object tags and object attributes are calculated separately according to (4.16).

$$\mathcal{L}_{MTL}^{tag} = -\frac{1}{M} \sum_{d_i \in \mathcal{D}} \sum_{j}^{V} y_i^j \log p(\bar{d}_i^j)$$

$$\mathcal{L}_{MTL}^{attribute} = -\frac{1}{Q} \sum_{e_m \in E} \sum_{n}^{V} y_m^n \log p(\bar{e}_m^n) \tag{4.16}$$

where $V$ is the size of the vocabulary (number of output classes), $p(\bar{d}_i^j)$ is the predicted output probability for output class $j$ for the masked tag token $d_i$, and $p(\bar{e}_m^n)$ is the predicted output probability for output class $n$ for the masked attribute token $e_m$. $y_i^j$ is the binary ground truth label for class $j$ for the masked word $d_i$ and $y_m^n$ is the binary ground truth label for class $n$ for the masked word $e_m$. $y_i^j$ and $y_m^n$ are calculated using the assignments $\alpha$ and $\beta$ as in (4.17).

$$y_i^j = \begin{cases} 1, & \text{if } \alpha(i) = j \\ 0, & \text{otherwise} \end{cases} \qquad y_m^n = \begin{cases} 1, & \text{if } \beta(m) = n \\ 0, & \text{otherwise} \end{cases} \tag{4.17}$$

Total $\mathcal{L}_{MTL}$ used during pretraining is the sum of individual losses for tags and attributes as in (4.18) where $\lambda$ is used to weight object tag and object attribute losses.

$$\mathcal{L}_{MTL}^{total} = (1 - \lambda)\mathcal{L}_{MTL}^{tag} + \lambda\mathcal{L}_{MTL}^{attribute} \tag{4.18}$$

As illustrated in Figure 4.14, full attention mask is used during pretraining. All object tags, object attributes, and regional features can attend to each other since the objective is learning the *visual vocabulary* while also exploiting the attributes. It is expected for the model to establish the relationship between objects, their attributes, and regional features.

The pretrained model with $\mathcal{L}_{MTL}$ objective for tags and attributes is finetuned for image captioning on COCO [27] dataset, similar to VinVL baseline which is described in Section 4.1.2.3.

As the first step, object tags, object attributes, and regional features are extracted in the same way as in pretraining (Figure 4.13) for the images in captioning dataset. This process is illustrated in Figure 4.15.
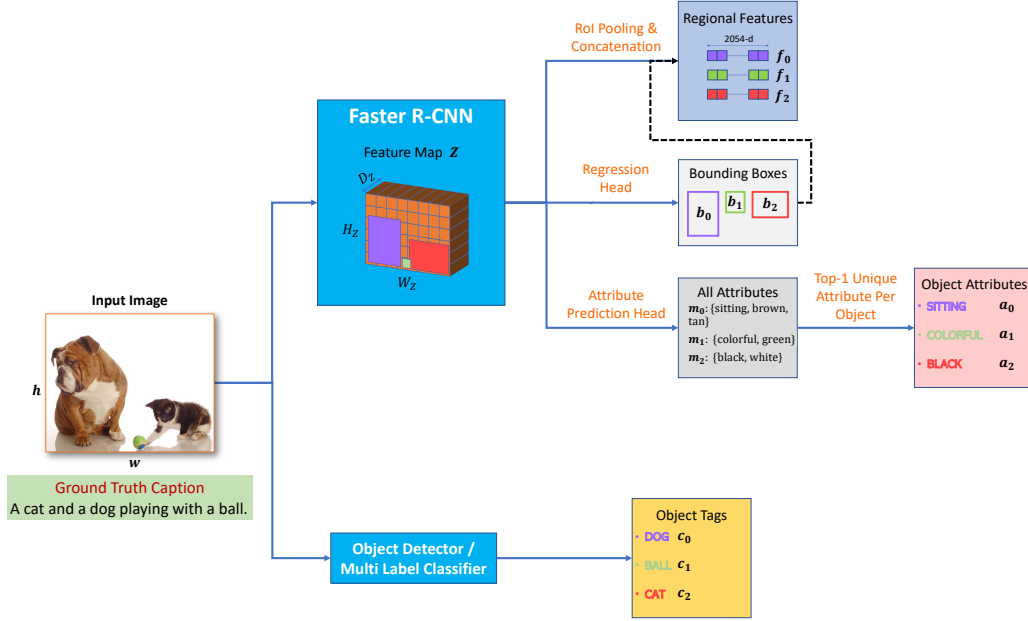


Figure 4.15: Feature, tag, and attribute extraction for the proposed method for an image in captioning dataset.

In the second step, ground truth caption, $S$, object tags, $C$, object attributes, $A$, and regional features, $F$, are fed into the transformer for image captioning finetuning. $\mathcal{L}_{MTL}$ objective is used during finetuning as well. Different from pretraining, object tags and attributes are not masked. Masking is only performed for tokens in ground truth caption.

Let $\mathcal{D} = \{d_0, d_1, \ldots, d_{M-1}\}$ be the set of tokens for masked words in $S \supseteq D$ and $\overline{\mathcal{D}} = \{\bar{d}_0, \bar{d}_1, \ldots, \bar{d}_{M-1}\}$ be the transformer outputs corresponding to elements in $\mathcal{D}$ where $M$ is the number of masked words in $S$. A token $s_i \in S$ is masked with 15% probability where $s_i$ is the token corresponding to the masked word in the ground

truth caption, for example *playing*. Following procedure is applied if token $s_i$ is to be masked:

- Replace the $s_i$ with $[MASK]$ token 80% of the time.

- Replace the $s_i$ with a random word in the vocabulary 10% of the time.

- Keep $s_i$ as it is 10% of the time.

In all three cases, $s_i$ is treated as a masked token and the network should classify transformer output corresponding to $s_i$ as *playing*.

The network is trained with $\mathcal{L}_{MTL}$ objective. Cross-Entropy loss is used in $\mathcal{L}_{MTL}$. $\mathcal{L}_{MTL}$ is calculated according to (4.19).

$$\mathcal{L}_{MTL} = -\frac{1}{M} \sum_{d_i \in \mathcal{D}} \sum_{j}^{V} y_i^j \log p(\bar{d}_i^j) \tag{4.19}$$

where $V$ is the size of the vocabulary (number of output classes), $p(\bar{d}_i^j)$ is the predicted output probability for output class $j$ for the masked token $d_i$. $y_i^j$ is the binary ground truth label for class $j$ for the masked word $d_i$.

Finetuning with attributes for image captioning is illustrated in Figure 4.16 for the example in Figure 4.15.

As shown in Figure 4.16, object attributes are integrated as another set of input but the training strategy is the same as VinVL baseline. Unidirectional attention mask is employed for tokens belonging to the ground truth caption, $S$, as shown in Figure 4.16. Hence, a token $s_i \in S$ can attend to tokens $s_j \in S$ for $j < i$ in the ground truth caption, $S$, all object tags, $C$, all object attributes, $A$, and all regional features, $F$. Furthermore, tokens belonging to object tags, object attributes, and regional features cannot attend to tokens of the ground truth caption. Otherwise, this would violate unidirectional attention for tokens in $S$ due to transitive property of attention.

Due to unidirectional attention mask, the model is forced to predict next word based on previous words, object tags, object attributes, and regional features. Object attributes are exploited as an additional source of information for the input image.
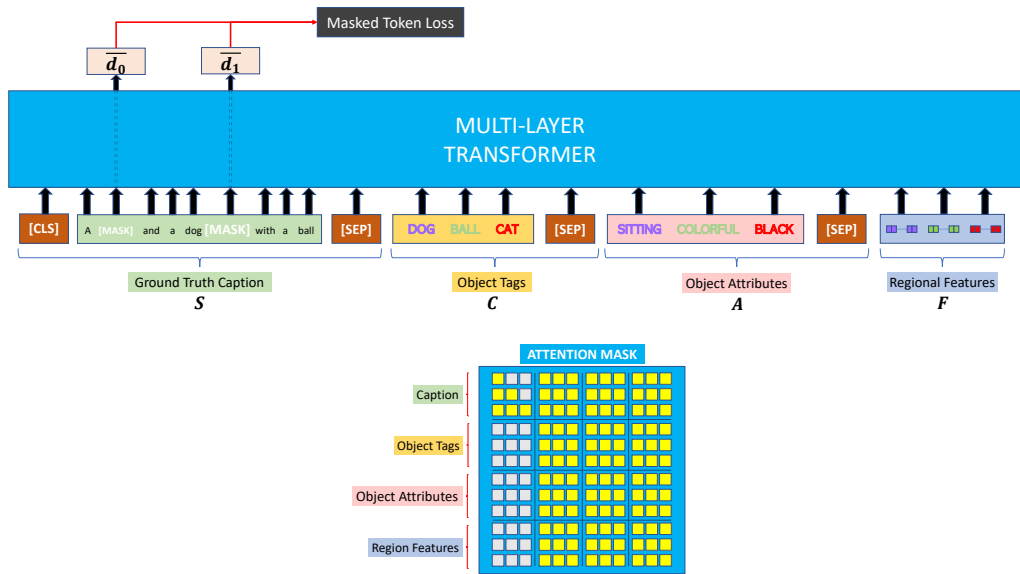
Figure 4.16: Illustration of finetuning for image captioning with attributes for the example in Figure 4.15.

After finetuning on image captioning dataset is performed, model can be used for caption generation (decoding). For decoding, object tags, object attributes, and regional features are extracted as in Figure 4.15.

The decoding step is the same as VinVL baseline described in Section 4.1.3 with additional set of inputs, object attributes.

There are 4 sets of inputs: Auto-regressive input tokens - $Y$, object tags - $C$, object attributes - $A$, and regional features - $F$.

The auto-regressive decoding procedure is provided in Algorithm 3. Forward propagation is applied using the given 4 sets of inputs. The transformer output $\bar{d}_t$ at time instant $t$ (corresponding to the $t$-th masked input token) is fed into the classifier and the network predicts the word $\bar{y}_t$. In the next time instant $t+1$, the $t$-th masked word is replaced with previously predicted word $\bar{y}_t$. After forward pass, the network makes prediction for $\bar{y}_{t+1}$. This process is followed in an auto-regressive manner until maximum caption length is achieved or the network predicts a special end-of-sentence token, $[EOS]$, at the output.

**Algorithm 3** Auto-regressive decoding process.

$Y = \{[CLS], [MASK], [MASK], \ldots, [MASK], [SEP]\}, \quad |Y| = T + 2$

$\overline{Y} = \{\}$                  ▷ Set of Words for the Generated Caption

$t = 0$

**while** $|\overline{Y}| \neq T$ **and** $\bar{y}_t \neq [EOS]$ **do**        ▷ $[EOS]$: End of Sentence Token

    $\bar{y}_t \leftarrow \text{Model}(Y, C, F)$

    $\overline{Y} \leftarrow \overline{Y} \cup \{\bar{y}_t\}$

    $y_t \leftarrow \bar{y}_t$

    $t \leftarrow t + 1$

**end while**

An example decoding snapshot for the example in Figure 4.15 at time instant $t = 7$ is given in Figure 4.17 (Note the extra attribute information, *black* and *sitting*, in the generated caption. The model is expected to generate more detailed captions with attributes).
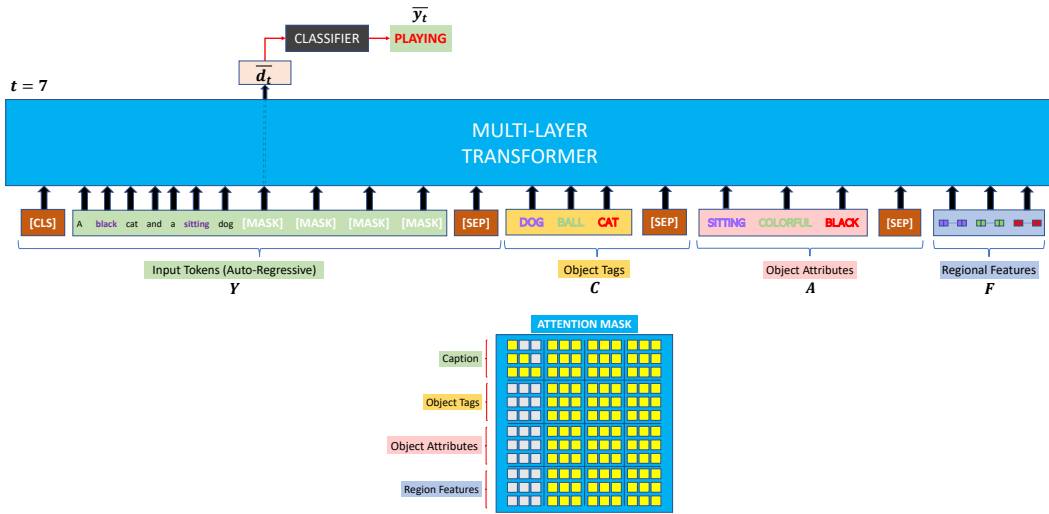


Figure 4.17: Illustration of auto-regressive decoding (caption generation) with attributes.

Similar to the finetuning stage, unidirectional attention mask is utilized since the model has no information about the future words in the caption. It is expected for the model to predict current masked word by attending to previously predicted words, object tags, object attributes, and regional features.

In the end, aim in proposed approach is exploiting additional object attribute information available in the image, during VIVO pretraining, image captioning finetuning, and caption generation, in order to produce richer and more informative captions.

# CHAPTER 5

## EXPERIMENTAL RESULTS

In this chapter experimental results will be discussed. Analysis of dataset and object tag components in Novel Object Captioning will be provided. Effect of object attribute component in Novel Object Captioning will be examined using the proposed novel approach which is explained in Section 4.2.

All models are first pretrained on a dataset (Open Images [3] or Flickr30k [2]) with VIVO [1] pretraining as described in Section 4.1.2.2. Then, the models are finetuned for image captioning on COCO [27] dataset as described in Section 4.1.2.3. Models are evaluated on *nocaps* [91] validation set and COCO test set (*Karpathy's* split [67]), for Novel Object Captioning and general image captioning, respectively. Details of datasets are provided in Section 3.3.

For Novel Object Captioning, an external object detection model trained on Open Images dataset is utilized in VinVL [4]. The model and its details are not shared. Hence, for the tag prediction in VinVL baseline models, we utilize a state-of-the-art multi-label classification model, ML-Decoder [33], trained on Open Images dataset. Its details are explained in Section 2.1.1.

It is relatively easy to train models for *nocaps* challenge since datasets are limited in number and size compared to VLP-based general image captioning approaches which utilize several large datasets as explained in Section 4.1.2.1. Hence, *nocaps* challenge is aimed for the purpose of demonstrating work conducted in this thesis but the ideas are applicable to other domains as well, given enough computational resources. Even though the models are trained for Novel Object Captioning, results on COCO [27] general image captioning dataset are also provided. Furthermore, proposed novel

approach is directly applicable to VLP-based general image captioning models as well.

## 5.1 Implementation Details

The models are pretrained and finetuned on a PC with NVIDIA GeForce RTX 3090 GPU with 24 GB memory and Intel Core i7-9700K CPU. All models are initialized with the weights of the $\text{BERT}_{\text{base}}$ [13] trained on large English corpora with Masked Language Modelling (MLM) objective. Batch size is used as 64. Following VIVO [1] and VinVL [4], learning rate of $5 \times 10^{-5}$ and $1 \times 10^{-5}$ are used during pretraining and finetuning, respectively. Following VIVO [1] and VinVL [4], we used maximum sequence lengths as in (5.1) and (5.2) for the input sets mentioned in Section 4.2, for pretraining and finetuning, respectively. The sets of inputs $S$, $C$, $A$, and $F$ do not necessarily have to have the same sizes. Some objects or attributes might be missed by the tag prediction, attribute prediction, or regional feature extraction networks. Hence, there is not a one-to-one correspondence between different sets of inputs.

$$
\begin{aligned}
|F| &= 50 \\
|C| &= 35 \\
|A| &= 35
\end{aligned}
\tag{5.1}
$$

$$
\begin{aligned}
|F| &= 50 \\
|S| &= 40 \\
|C| &= 30 \\
|A| &= 15
\end{aligned}
\tag{5.2}
$$

Maximum sentence length during decoding is 20, as in VIVO [1].

## 5.2 Effect of Pretraining Datasets

In *nocaps* [91] challenge, two datasets are provided for training: Open Images [3] and COCO [27]. Open Images is an object detection dataset. Training set of Open Images dataset is used in order to allow models to learn the concept of novel object classes

since the dataset contains objects from around 600 novel classes. COCO dataset is used for captioning training.

Since Open Images dataset is very large ( 1.7 million images) it takes several GPU weeks to just pretrain the network, let alone finetuning for captioning. We question the necessity of such large dataset for Novel Object Captioning. For this purpose, we trained two VinVL models [4]. The baseline model is pretrained on Open Images dataset with VIVO pretraining and finetuned on COCO dataset for captioning. This is the standard apporach used in VIVO [1] and VinVL [4]. The second model is pretrained on a much smaller dataset, Flickr30k [2] ( 30k images), with VIVO pretraining and finetuned on COCO dataset for captioning.

During the pretraining on Flickr30k with VIVO objective, we still used an external tag prediction network trained on Open Images in order to allow the network to learn *visual vocabulary* for the novel objects. Hence, despite the model is pretrained on a smaller dataset, it is still exposed to novel classes of Open Images dataset.

Evaluation results of two models on COCO and *nocaps* datasets are given in Table 5.1 and Table 5.2, respectively. As seen in Table 5.1, two models perform very

Table 5.1: Evaluation results on COCO Karpathy test split [67] for VinVL [4] pretrained on Open Images [3] and Flickr30k [2] datasets and finetuned on COCO [27] dataset.

| | | Model | |
| | | Pre-Trained on | Pretrained on |
| Dataset | Metric | Open Images (Baseline) | Flickr30k |
| | BLEU-4 | **35.27** | 34.8 |
| | METEOR | 28.39 | **28.47** |
| COCO | ROUGE-L | **57.14** | 56.71 |
| | CIDEr | **119.41** | 118.16 |
| | SPICE | 21.62 | **21.73** |

similar on COCO test set. COCO dataset have less variation in terms of objects in the

Table 5.2: Evaluation results on *nocaps* validation set for VinVL [4] pretrained on Open Images [3] and Flickr30k [2] datasets and finetuned on COCO [27] dataset.

| Dataset | Subset | Metric | Model | |
|---|---|---|---|---|
| | | | Pre-Trained on Open Images (Baseline) | Pretrained on Flickr30k |
| *nocaps* | in-domain | BLEU-4 | 27.43 | **27.45** |
| | | METEOR | 28.50 | **28.59** |
| | | ROUGE-L | 56.08 | **56.73** |
| | | CIDEr | 85.48 | **88.92** |
| | | SPICE | 12.84 | **13.16** |
| | near-domain | BLEU-4 | 21.89 | **23.65** |
| | | METEOR | 25.94 | **26.72** |
| | | ROUGE-L | 54.03 | **54.87** |
| | | CIDEr | 77.22 | **81.31** |
| | | SPICE | 12.3 | **12.54** |
| | out-of-domain | BLEU-4 | 11.71 | **13.38** |
| | | METEOR | 21.62 | **22.68** |
| | | ROUGE-L | 48.17 | **49.31** |
| | | CIDEr | 59.53 | **65.51** |
| | | SPICE | 10.14 | **10.77** |
| | entire | BLEU-4 | 20.62 | **22.11** |
| | | METEOR | 25.43 | **26.17** |
| | | ROUGE-L | 53.13 | **54.01** |
| | | CIDEr | 74.82 | **79.2** |
| | | SPICE | 11.96 | **12.29** |

image. Hence, the effect of learning visual vocabulary for a large set of novel objects is not critical on COCO dataset. Furthermore, the models are finetuned on COCO datasets. This allows models to fit their parameters according to image captioning task on COCO dataset. So, the effect of VIVO pretraining dataset on general image

captioning is insignificant. This outcome is in parallel with VIVO [1] where they claim that the effect of VIVO pretraining is diminishing on COCO dataset.

The model pretrained on Flickr30k with VIVO pretraining objective surpasses the baseline model on *nocaps* validation set, as illustrated in Table 5.2. The results suggest that for a better Novel Object Captioning performance a large pretraining dataset is not mandatory. As will be discussed in Section 5.3, it is more crucial for the network to see novel object tags during pretraining and finetuning.

Regarding the performance gap between two models, during the pretraining stage, the model trained on Flickr30k converged rapidly after 30 epochs whereas for the baseline model trained on Open Images it took around 200 epochs. This is due to the fact that Open Images dataset being significantly larger ( 55 times) than Flickr30k dataset. The pretraining stage on Open Images took approximately 1 month. We hypothesize that with more careful training strategies and hyperparameter optimizations, the model could fit better to such large dataset. However, due to computational constraints, we could not perform further experiments for this hypothesis.

## 5.3 Effect of Different Set of Object Tags

For the *nocaps* challenge, models utilize external object detection networks which are trained on Open Images dataset which contains around 600 novel classes. For regional feature and object tag extraction in general image captioning, VinVL utilize an object detection network, Faster R-CNN [7], as explained in Section 4.1.1. This network is trained on a collection of datasets ([27], [3], [42], [38]) and it predicts 1594 output classes for objects in the image. For *nocaps* challenge, it is only used as regional feature extractor. An object detector trained on Open Images is used as object tag extractor for the *nocaps* challenge. 313 classes of Open Images dataset also exist in 1594 output classes of the Faster R-CNN object detector. This leaves around 200 novel classes which are not covered by the 1594 classes.

In order to analyze the effect of quantity and novelty of object tags on Novel Object Captioning and general image captioning, we train two models with two different object tag extractors. The baseline model utilize ML-Decoder [33] trained on Open

Images with 523 novel object classes as the tag predictor. The second model utilize Faster R-CNN object detection model with 1594 output classes as the tag predictor. Both captioning models are pretrained on Open Images dataset with VIVO pretraining objective and finetuned on COCO dataset for image captioning, following the restrictions of *nocaps* challenge. Even though the second model lacks around 200 novel classes compared to the baseline, output vocabulary size is the same for both models (30522 words). Hence, in theory, both models could generate the same set of words during caption generation.

Evaluation results of two models on COCO and *nocaps* datasets are given in Table 5.3 and Table 5.4, respectively. As seen in Table 5.3, two models perform pretty

Table 5.3: Evaluation results on COCO Karpathy test split [67] for VinVL [4] trained with 523 OI (Open Images) classes as object tags and 1594 VG (Visual Genome) classes as object tags.

| | | Model | |
| --- | --- | --- | --- |
| | | OI Classes (523) | VG Classes (1594) |
| Dataset | Metric | As Tags (Baseline) | As Tags |
| | BLEU-4 | **35.27** | 34.93 |
| | METEOR | 28.39 | **28.47** |
| COCO | ROUGE-L | **57.14** | 57.13 |
| | CIDEr | 119.41 | **119.47** |
| | SPICE | 21.62 | **21.85** |

similar on COCO test set. The reason for this outcome is that COCO dataset does not contain novel objects. 80 object classes in the COCO dataset also exist in 523 Open Images classes and 1594 Visual Genome classes. Hence, both models cover all COCO classes. Furthermore, the models are finetuned on COCO dataset. This allows models to fit their parameters according to image captioning task on COCO dataset. Hence, the effect of novelty or abundance of object tags on general image captioning is insignificant. This outcome is in parallel with VIVO [1] where they claim that the effect of VIVO pretraining is diminishing on COCO dataset.

Table 5.4: Evaluation results on *nocaps* validation set for VinVL [4] trained with 523 OI (Open Images) classes as object tags and 1594 VG (Visual Genome) classes as object tags.

| | | | Model | |
|---|---|---|---|---|
| | | | OI Classes (523) | VG Classes (1594) |
| Dataset | Subset | Metric | As Tags (Baseline) | As Tags |
| | | BLEU-4 | **27.43** | 27.03 |
| | | METEOR | **28.5** | 27.79 |
| | in-domain | ROUGE-L | **56.08** | 55.93 |
| | | CIDEr | **85.48** | **85.48** |
| | | SPICE | **12.84** | 12.81 |
| | | BLEU-4 | 21.89 | **22.08** |
| | | METEOR | **25.94** | 25.55 |
| | near-domain | ROUGE-L | **54.03** | 53.75 |
| | | CIDEr | **77.22** | 73.67 |
| | | SPICE | **12.30** | 12.08 |
| *nocaps* | | BLEU-4 | **11.71** | 10.82 |
| | | METEOR | **21.62** | 20.83 |
| | out-of-domain | ROUGE-L | **48.17** | 47.30 |
| | | CIDEr | **59.53** | 54.18 |
| | | SPICE | **10.14** | 9.81 |
| | | BLEU-4 | **20.62** | 20.51 |
| | | METEOR | **25.43** | 24.91 |
| | entire | ROUGE-L | **53.13** | 52.76 |
| | | CIDEr | **74.82** | 71.41 |
| | | SPICE | **11.96** | 11.77 |

In Table 5.4, it is shown that baseline model utilizing Open Images classes as object tags surpasses the model utilizing Visual Genome classes as object tags on *nocaps* validation set. The difference between models is small for *in-domain* subset.

This is expected since the images in *in-domain* subset contains only COCO classes. The baseline model outperforms the other model when images contain novel objects (*near-domain* and *out-of-domain* subsets). Even though the second model is trained with much larger set of object tags, it fails to generate captions for novel objects. This result indicates that seeing a novel object class during training increases the accuracy of generated caption for an image which contains the novel object. During training, the model should be exposed to the set of object classes which exist in the test images. Hence, novelty of object tags is crucial for Novel Object Captioning. Increasing the novelty and number of object classes seen during pretraining and finetuning might lead even better results. An expanded set of object tags with both novel classes of Open Images and diverse classes of Visual Genome could yield richer and more grounded captions for both Novel Object Captioning and general image captioning. This may be investigated as a future work.

## 5.4  Effect of Object Attributes

Results of Section 5.3 demonstrate that it is crucial for the network to be exposed to novel object tags if we would like to generate captions which mention novel objects in the image. Following this outcome, as explained in Section 4.2, we propose a novel approach which aims to enrich generated captions with object attributes. We expose the network to object attributes during pretraining and finetuning in order to generate captions with attributes.

Two methods are compared in this experiment. The baseline model is pretrained on Open Images with VIVO objective and finetuned on COCO for captioning. The proposed model is pretrained on Open Images with VIVO-like pretraining and finetuned on COCO for captioning while also exploiting object attributes in both stages, as explained in Section 4.2. We used $\lambda = 0.5$ in (4.18). Hence, during pretraining, both tags and attributes have equal effect on overall loss. In this experiment, we investigate the effect of proposed method in generating richer captions in addition to the theoretical captioning performance.

Evaluation results of baseline method and proposed method on COCO and *nocaps*

datasets are given in Table 5.5 and Table 5.6, respectively.    As seen in Table 5.5 and

Table 5.5: Evaluation results on COCO Karpathy test split [67] for VinVL [4] baseline and proposed method which utilize object attributes.

| | | Model | |
| Dataset | Metric | Baseline (without Attributes) | Proposed Method (with Attributes) |
| --- | --- | --- | --- |
| | BLEU-4 | **35.27** | 32.96 |
| | METEOR | **28.39** | 27.08 |
| COCO | ROUGE-L | **57.14** | 55.09 |
| | CIDEr | **119.41** | 110.62 |
| | SPICE | **21.62** | 20.70 |

Table 5.6, the baseline method surpasses the proposed method on both COCO and *nocaps* datasets. We think there are a couple of reasons for this result. Firstly, during evaluation, predicted captions are compared with a set of reference ground truth captions. These ground truth captions are human-annotated. Hence, if ground truth captions do not contain object attributes, the model which mentions object attributes in its predicted caption might get lower evaluation score. Secondly, integrating attribute information during pretraining and finetuning might have affected the model negatively such that generated captions might be wrong or grammatically incorrect. In order to understand the reasoning behind the theoretical results and observe the effect of proposed method in generated captions, we analyze the visual outputs of baseline method and proposed method. In visual examples, predicted object attributes by the proposed method are shown in bold. These attributes are one of the 524 output classes of the attribute prediction network.

An example image from COCO test set, generated captions by the baseline and proposed methods, and ground truth captions are given in Table 5.7.    In Table 5.7, the prediction of the baseline method is not fully correct since it hallucinates a *chair*. Furthermore, it does not mention about the *car* in the background and does not give any details or attributes about the objects in the image. On the other hand, our proposed

Table 5.6: Evaluation results on *nocaps* validation set for VinVL [4] baseline and proposed method which utilize object attributes.

| Dataset | Subset | Metric | Model | |
| --- | --- | --- | --- | --- |
| | | | Baseline (without Attributes) | Proposed Method (with Attributes) |
| *nocaps* | in-domain | BLEU-4 | **27.43** | 24.54 |
| | | METEOR | **28.5** | 26.75 |
| | | ROUGE-L | **56.08** | 54.02 |
| | | CIDEr | **85.48** | 79.29 |
| | | SPICE | **12.84** | 12.44 |
| | near-domain | BLEU-4 | **21.89** | 20.13 |
| | | METEOR | **25.94** | 24.94 |
| | | ROUGE-L | **54.03** | 52.63 |
| | | CIDEr | **77.22** | 72.1 |
| | | SPICE | **12.3** | 11.74 |
| | out-of-domain | BLEU-4 | **11.71** | 11.31 |
| | | METEOR | **21.62** | 20.86 |
| | | ROUGE-L | **48.17** | 47.63 |
| | | CIDEr | **59.53** | 55.07 |
| | | SPICE | **10.14** | 9.81 |
| | entire | BLEU-4 | **20.62** | 18.97 |
| | | METEOR | **25.43** | 24.37 |
| | | ROUGE-L | **53.13** | 51.82 |
| | | CIDEr | **74.82** | 69.68 |
| | | SPICE | **11.96** | 11.48 |

approach which exploit object attributes successfully describe the scene. It mentions about the *car* in the background, it gives details about the objects in the image by mentioning about their attributes such as *small*, *wooden*, and *red*. This was what we aimed for in our proposed approach. This example reveals a few important points as

Table 5.7: An example image and generated caption by the baseline method, generated caption by the proposed method and ground truth reference sentences for the image.



| Baseline | – A garden with a chair and a plant on the ground. |
| --- | --- |
| **Proposed Method** | – A **small** garden with a **wooden** fence and a **red** car behind it. |
| **Ground Truth** | – A simple, fenced in garden containing several kinds of plants.<br>– A plot of dirt with a fence around it with flower pots next to it.<br>– A few plants in a garden near a fence.<br>– A small garden that has vegetables blooming in it.<br>– A car some dirt a garden grass bushes and trees. |

well.

Firstly, ground truth captions describe the image from different perspectives. Some describe the *vegetables* in it, some describe the *potted plants* next to the *garden*, and only the last ground truth caption mentions about the *car* in the background. Moreover, the last ground truth caption is grammatically ambiguous since it is not

a proper sentence. We have observed similar issues with ground truth captions of the other images from COCO dataset. These kind of ground truth captions which are grammatically problematic and overlooking the details in the image might result in poor evaluation results for the models. More importantly, such erroneous ground truth captions in the training set would affect performance of the models and result in degraded performance.

Secondly, we observe that some of the ground truth captions mention about object attributes such as *simple, small, fenced, dirt*. However, some attributes for the objects are missing. For example, *wooden*, and *red*. To develop models which generate richer captions, it might be better to enrich ground truth captions as a first step.

In Table 5.8, Table 5.9, and Table 5.10 example images from COCO test set, generated captions by the baseline and proposed methods, ground truth captions, and evaluation scores for the predictions of models are provided. In these examples it is shown that the proposed method successfully describes the scenes with additional attribute information (*white, pink, metal, stainless steel,* and *sitting*). However, in all of these examples, output of the proposed method gets lower scores compared to the baseline. The reason for this result is that ground-truth captions in these examples describe the scene from different perspectives and they do not contain object attributes. Hence, the proposed method is penalized for integrating attributes even though the predicted sentences describe the scenes perfectly. Regardless of the theoretical results, the proposed pretraining and finetuning approach which exploits attributes seems to be helping generating richer captions by mentioning attributes of objects in the captions. These points support our claim that assessment of the proposed method based solely on the theoretical evaluation metrics might be misleading. Since ground truth captions are generated by humans, they are open-ended. Hence, captions with missing details and attributes yield poor evaluation results. So, it is important to take visual examples and outputs into consideration as well.

An example image from *nocaps* validation set, generated captions by the baseline and proposed methods, and ground truth captions are given in Table 5.11. In *nocaps* dataset, each image has 10 ground truth captions as opposed to COCO dataset which has 5 ground turth captions per image. In Table 5.11, we can see that the baseline

Table 5.8: An example of poor theoretical evaluation due to mediocre ground-truth captions.



| Baseline | – A close up of a street sign with a sky background. |
|---|---|
| **Proposed Method** | – A couple of **white** street signs **sitting** on top of a **metal** pole. |
| **Ground Truth** | – A street sign that reads " Greta Garbo Strafze ".<br>– A street sign that reads Greta Garbo Strafe.<br>– Two intersecting white street signs at an intersection.<br>– The street sign is reading Greta Garbo on the side of the pole.<br>– A street sign that is pointing in different ways. |
| **Score** | **Baseline**: CIDEr: **61.27**    SPICE: **16.6**<br>**Proposed**: CIDEr: 38.04    SPICE: 14.29 |

method does not give details about the objects in the image. Moreover, it hallucinated a *dishwasher* which does not exist in the image. On the other hand, the proposed method correctly described the scene while also mentioning about the object attributes such as *metal* and *wooden*. It is also clear that ground truth captions in *nocaps* dataset is more detailed. However there might still be missing attributes. For example, the *metal* attribute for the object *sink* is not mentioned in any on the ground truth captions while it is mentioned in the caption generated by the proposed method. This might result in poor evaluation result for the proposed method. The reason for performance of the all models being worse on *nocaps* dataset compared to COCO dataset might be the difference in detail levels of the ground truth captions. The models are trained on

Table 5.9: An example of poor theoretical evaluation due to mediocre ground-truth captions.



| Baseline | – A bathroom with a toilet, sink, and soap dispenser. |
|---|---|
| **Proposed Method** | – A **white** toilet **sitting** next to a **stainless steel** sink. |
| **Ground Truth** | – We are looking at a small airline toilet. |
| | – A compact sized bathroom with toilet and sink inside. |
| | – A toilet and sink in the bathroom of an airplane. |
| | – A bathroom toilet that has the seat down. |
| | – This is a lavatory of an airplane. |
| **Score** | **Baseline**: <u>CIDEr</u>: **63.0**    <u>SPICE</u>: **35.71** |
| | **Proposed**: <u>CIDEr</u>: 43.64    <u>SPICE</u>: 14.81 |

COCO dataset which has broad ground truth captions with less details and they are evaluated on *nocaps* dataset which has very detailed ground truth captions.

Another example image from *nocaps* validation set, generated captions by the baseline and proposed methods, and ground truth captions are given in Table 5.12. In Table 5.12, we observe that both the baseline method and the proposed method generate good captions. The baseline method also mentioned about attributes such as *yellow, black,* and *sitting*. During experiments, we observed that the baseline method can predict simple attributes such as colors and shapes from time to time. The reason for this might be that such simple and common attributes also exist in the ground truth captions of the COCO dataset. However, we observe that the baseline model chose the word *black* to describe the bird while the bird is *blue*. Proposed method successfully described the bird as *yellow* and *blue* due to explicit integration of object

Table 5.10: An example of poor theoretical evaluation due to mediocre ground-truth captions.



| | |
|---|---|
| **Baseline** | – A vase filled with flowers sitting on top of a table. |
| **Proposed Method** | – A vase filled with **white** and **pink** flowers on top of a table. |
| **Ground Truth** | – Creative centerpiece floral arrangement at an outdoor event. <br> – A wedding centerpiece made of flowers and various other plants. <br> – A vase of flowers sitting on an outdoor table. <br> – A vase filled with flowers on top of a table. <br> – A floral arrangement inside a clear cylinder shaped vase. |
| **Score** | **Baseline**: CIDEr: **207.57**     SPICE: **33.3** <br> **Proposed**: CIDEr: 138.99     SPICE: 33.3 |

attributes during training and decoding. Furthermore, the proposed method used the word *perched* to describe the bird instead of using the word *sitting*. *perched* is one of the 524 output classes of the attribute prediction network.

These results illustrate the effectiveness of attribute component in Novel Object Captioning and general image captioning in order to generate richer and more visually

Table 5.11: An example image and generated caption by the baseline method, generated caption by the proposed method and ground truth reference sentences for the image.



| Baseline | – A kitchen with a sink, cabinets, and a dishwasher. |
| --- | --- |
| **Proposed Method** | – A kitchen with a **metal** sink and **wooden** cabinets. |
| **Ground Truth** | – Wood cabinetry surrounding a sink with a high curved faucet and kitchen knives on the left side of the counter.<br>– The interior of a kitchen showing the sink and surrounding counter top.<br>– A kitchen with a cabinet, sink and counter top.<br>– A clean counter top with a cool cabinet furniture downward.<br>– A kitchen counter and double sink with a towel laying on the counter.<br>– Kitchen sink with granite countertop and wooden drawers.<br>– A sink with knives, plates, and a towel on top of the counter.<br>– A grey sink is placed in the middle of a countertop.<br>– A kitchen sink with many cabinets under the sink.<br>– Knives, dish soap, a plate, an orange container and a towel lie on a kitchen top of the sink. |

grounded captions.

Even though the proposed method generates richer and more detailed captions compared to the baseline, we have observed that there are some drawbacks of the proposed

Table 5.12: An example image and generated caption by the baseline method, generated caption by the proposed method and ground truth reference sentences for the image.



| Baseline | – A yellow and black bird sitting on a tree branch. |
|---|---|
| **Proposed Method** | – A **yellow** and **blue** bird **perched** on a tree branch. |
| **Ground Truth** | – The yellow and blue bird is peeking on the plant.<br>– A blue and yellow bird is on branch next to green leaves.<br>– A blue and yellow bird perched on a twig.<br>– A blue and yellow bird is sitting in a tree with green leaves.<br>– A yellow and blue bird on a branch in a tree with a lot of leaves.<br>– A blue and yellow bird sitting on a branch with some leaves.<br>– A blue and gold bird is sitting on a branch surrounded by leaves.<br>– A blue and yellow bird is perched on a tree branch.<br>– A yellow and blue bird sitting in a wooden branch.<br>– A blue and yellow bird on top of a twig. |

method for some images. An example visual result demonstrating such a case for an image from COCO test set is given in Table 5.13. The example in Table 5.13 shows that proposed method associates *bunch* and *sliced* attributes for the *orange* object correctly. Description of *oranges* is more detailed compared to the baseline. However, we observe that the proposed method fails to mention about the *plate* and the *bananas* in the background. The baseline method mentions about these objects correctly. However, it fails to predict any attributes for the objects in the image. The

Table 5.13: An example image and generated caption by the baseline method, generated caption by the proposed method and ground truth reference sentences for the image.
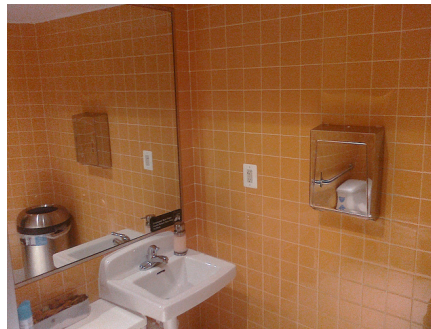


| Baseline | – A close up of oranges and bananas on a plate. |
|---|---|
| **Proposed Method** | – A **bunch** of **sliced** oranges sitting on a table. |
| **Ground Truth** | – A plate is piled high of orange slices while a bunch of bananas sits next to it. |
| | – A plate full of sliced oranges next to a bunch of bananas. |
| | – A white plate filled with slice oranges next to a pile of bananas. |
| | – A close up of orange slices and bananas. |
| | – A close-up picture of cut oranges on a plate, with blurred bananas in the background. |

failure of the proposed model to describe the *plate* and the *bananas* may be because of the proposed model being focused on attributes too much. During pretraining, we set $\lambda = 0.5$ in (4.18). Thus, both the masked object tags and masked object attributes have the same effect on the total loss. Decreasing $\lambda$ to down weight the contribution of the masked attributes on total loss might help model to focus on objects during decoding. On the other hand, such an approach may lead to less richer captions in terms of object attributes. During finetuning and decoding we set the maximum sequence length for the object attributes as half of the maximum sequence length for the object tags. Maximum sequence length for the attributes can be reduced further to decrease the effect of object attributes in generated captions. However, it should be noted that not integrating object attributes at all results in captions without details as seen in the

prediction of the baseline model.

An example visual result where proposed method generates a rich caption but fails to mention all objects in the image is given in Table 5.14. In Table 5.14, it is shown

Table 5.14: An example image and generated caption by the baseline method, generated caption by the proposed method and ground truth reference sentences for the image.



| **Baseline** | – A bath room with a sink a mirror and a towel dispenser. |
|---|---|
| **Proposed Method** | – A bathroom with **orange tile** walls and a **white** sink. |
| **Ground Truth** | – The bath room is clean with brown tile, white sink and large mirror. <br> – A public bathroom area with orange tile walls. <br> – A sink sits in a public bathroom with bright orange tile. <br> – A public restroom with focus on the sink and towel dispenser. <br> – A bathroom has gold tile and a silver box on the wall. |

that proposed method describes the objects in the scene with several attributes such as *white, orange,* and *tile*. However, it fails to mention about *towel dispenser* and *mirror* whereas the baseline method mentions about these objects. On the other hand, the baseline method does not predict any attributes and it fails to mention about the walls unlike the proposed method.

These results illustrate that it is important to have a balance between object attributes and object tags during training and inference in order to fully describe an image with

attribute details.

More visual examples comparing the baseline method and the proposed method are given in Appendix A.

# CHAPTER 6

# SUMMARY, CONCLUSIONS AND FUTURE WORK

## 6.1 Summary and Conclusion

In this study, we have analyzed the effect of dataset, object tag, and object attribute components in Novel Object Captioning. We have provided results and analysis for general image captioning as well. During experiments, we have performed both theoretical evaluation and visual analysis of different methods.

Recently, transformer-based networks replaced the RNN-based ones in sequence understanding and generation tasks [12], [13], [14]. Compared to RNNs, transformers provide better modelling of relationships and interactions in long sequences, and they are more suitable for parallel processing [44]. Due to these advantages, transformers are also utilized by many image captioning methods in recent years [30], [16], [1] [4], [83]. Hence, in this work, we have taken a state-of-the-art transformer-based image captioning model, VinVL [4], as baseline.

VinVL requires a large pretraining dataset and a specific set of object tags for Novel Object Captioning (*nocaps* [91]). We investigate the effect of dataset and object tag components and question the necessity of them. Furthermore, image captioning models generally lack visual grounding and copy phrases from the training dataset [31], [32]. To alleviate this problem, we have proposed a novel approach which exploits additional object attribute information to generate richer and visually grounded captions.

As the first experiment, we have pretrained the VinVL model on a much smaller dataset, Flickr30k [2], while keeping all other components fixed. As the second ex-

periment, we have pretrained and finetuned the VinVL model using a much larger set of object tags (1594 classes) which lacks around 200 novel objects available in the baseline that contains a total of 523 novel object classes. Experimental results demonstrate that it is very crucial for the model to be exposed to the novel object tags during training in order to be successful on Novel Object Captioning dataset. On the contrary, as long as the model sees novel object tags during training, similar Novel Object Captioning performance could be obtained with a much smaller pretraining dataset. Changing the pretraining dataset or object tags do not seem to be affecting the captioning performance on general image captioning significantly because COCO [27] dataset is simpler compared to the *nocaps* dataset and the models are finetuned directly on the COCO dataset.

These results demonstrated that in order to mention about some novel objects in the generated caption, it is essential for the network to see those keywords during training. Hence, in order to enrich generated captions, we have proposed to expose the network to the object attributes during training and decoding. We have integrated object attributes during pretraining with a Masked Token Loss (MTL) training objective which allows the network to build the relationships (attentions) between object tags, object attributes, and regional features. Experimental results show that the proposed method successfully utilizes attributes for the objects in the image which results in richer and visually grounded captions. Even though we performed model training according to the *nocaps* challenge due to limited computational resources, the proposed method is directly applicable to the general image captioning tasks which utilize several large image-tag paired datasets as well.

## 6.2 Future Work

Visual analysis on the baseline and proposed methods revealed that the proposed method generates grammatically and semantically correct captions in general. The generated captions are also richer in terms of object attributes compared to the baseline method. However the proposed method falls behind the baseline method when it comes to theoretical evaluation. To figure out the reasons behind this result we have analyzed datasets, ground truth captions, and predictions of models. The anal-

ysis on datasets and ground truth captions shows that different people describe an image from very different perspectives resulting in ground truth captions with high variance. Furthermore, some ground truth captions are very broad and do not mention about important details in the image such as attributes of objects. This results in poor evaluation performance for the proposed method. Hence, as a future work, we think that it is crucial to improve the quality of the ground truth captions in order to boost the performances of captioning models and perform a fair evaluation. Furthermore, a captioning dataset with detailed ground truth captions for the training set would help models to produce richer and detailed captions during inference.

In addition, the proposed method for integrating object attributes can also be applied to the object relationships as well by employing a visual relationship detection network. That might lead to even better and more detailed captions.

Lastly, object tags, object attributes, and regional features are fed to the network directly without association between them. Explicitly aligning these three sets of inputs as *(tag, attribute, feature)* triplets might allow the network to build up the relationships between them during training and predict words which are more grounded on tags and attributes during decoding.

# REFERENCES

[1] X. Hu, X. Yin, K. Lin, *et al.*, "VIVO: Visual Vocabulary Pre-Training for Novel Object Captioning", *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 2, pp. 1575–1583, 2 May 18, 2021, ISSN: 2374-3468.

[2] P. Young, A. Lai, M. Hodosh, and J. Hockenmaier, "From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions", *Transactions of the Association for Computational Linguistics*, vol. 2, pp. 67–78, 2014. DOI: `10.1162/tacl_a_00166`.

[3] A. Kuznetsova, H. Rom, N. Alldrin, *et al.*, "The Open Images Dataset V4", *International Journal of Computer Vision*, 2020. DOI: `10.1007/s11263-020-01316-z`.

[4] P. Zhang, X. Li, X. Hu, *et al.*, "VinVL: Revisiting Visual Representations in Vision-Language Models", presented at the Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2021, pp. 5579–5588.

[5] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition", in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2016, pp. 770–778. DOI: `10.1109/CVPR.2016.90`.

[6] A. Dosovitskiy, L. Beyer, A. Kolesnikov, *et al.*, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", *ICLR*, 2021.

[7] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", in *Advances in Neural Information Processing Systems*, vol. 28, Curran Associates, Inc., 2015.

[8] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-End Object Detection with Transformers", in *Computer Vision – ECCV 2020*, A. Vedaldi, H. Bischof, T. Brox, and J.-M. Frahm, Eds., ser. Lec-

ture Notes in Computer Science, Cham: Springer International Publishing, 2020, pp. 213–229, ISBN: 978-3-030-58452-8. DOI: 10.1007/978-3-030-58452-8_13.

[9] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional Networks for Biomedical Image Segmentation", *MICCAI*, 2015. DOI: 10.1007/978-3-319-24574-4_28.

[10] K. He, G. Gkioxari, P. Dollar, and R. Girshick, "Mask R-CNN", presented at the Proceedings of the IEEE International Conference on Computer Vision, 2017, pp. 2961–2969.

[11] K. Cho, B. van Merriënboer, C. Gulcehre, *et al.*, "Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation", in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Doha, Qatar: Association for Computational Linguistics, Oct. 2014, pp. 1724–1734. DOI: 10.3115/v1/D14-1179.

[12] A. Vaswani, N. Shazeer, N. Parmar, *et al.*, "Attention is All you Need", in *Advances in Neural Information Processing Systems*, vol. 30, Curran Associates, Inc., 2017.

[13] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding", in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, Minneapolis, Minnesota: Association for Computational Linguistics, Jun. 2019, pp. 4171–4186. DOI: 10.18653/v1/N19-1423.

[14] T. Brown, B. Mann, N. Ryder, *et al.*, "Language Models are Few-Shot Learners", in *Advances in Neural Information Processing Systems*, vol. 33, Curran Associates, Inc., 2020, pp. 1877–1901.

[15] A. Mogadala, M. Kalimuthu, and D. Klakow, "Trends in Integration of Vision and Language Research: A Survey of Tasks, Datasets, and Methods", *J. Artif. Intell. Res.*, 2021. DOI: 10.1613/jair.1.11688.

[16] X. Li, X. Yin, C. Li, *et al.*, "Oscar: Object-Semantics Aligned Pre-training for Vision-Language Tasks", in *ECCV*, 2020. DOI: `10.1007/978-3-030-58577-8_8`.

[17] D. Gurari, Y. Zhao, M. Zhang, and N. Bhattacharya, "Captioning Images Taken by People Who Are Blind", in *ECCV*, 2020. DOI: `10.1007/978-3-030-58520-4_25`.

[18] A. F. Biten, L. Gómez, M. Rusiñol, and D. Karatzas, "Good News, Everyone! Context Driven Entity-Aware Captioning for News Images", *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. DOI: `10.1109/CVPR.2019.01275`.

[19] J. Pavlopoulos, V. Kougia, and I. Androutsopoulos, "A Survey on Biomedical Image Captioning", in *Proceedings of the Second Workshop on Shortcomings in Vision and Language*, Minneapolis, Minnesota: Association for Computational Linguistics, Jun. 2019, pp. 26–36. DOI: `10.18653/v1/W19-1803`.

[20] A. Majumdar, A. Shrivastava, S. Lee, P. Anderson, D. Parikh, and D. Batra, "Improving Vision-and-Language Navigation with Image-Text Pairs from the Web", in *Computer Vision – ECCV 2020*, A. Vedaldi, H. Bischof, T. Brox, and J.-M. Frahm, Eds., ser. Lecture Notes in Computer Science, Cham: Springer International Publishing, 2020, pp. 259–274, ISBN: 978-3-030-58539-6. DOI: `10.1007/978-3-030-58539-6_16`.

[21] M. Stefanini, M. Cornia, L. Baraldi, S. Cascianelli, G. Fiameni, and R. Cucchiara, "From Show to Tell: A Survey on Deep Learning-based Image Captioning", *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–1, 2022, ISSN: 1939-3539. DOI: `10.1109/TPAMI.2022.3148210`.

[22] A. Rohrbach, L. A. Hendricks, K. Burns, T. Darrell, and K. Saenko, "Object Hallucination in Image Captioning", in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, Brussels, Belgium: Association for Computational Linguistics, Oct. 2018, pp. 4035–4045. DOI: `10.18653/v1/D18-1437`.

[23] Y. Zhou, M. Wang, D. Liu, Z. Hu, and H. Zhang, "More Grounded Image Captioning by Distilling Image-Text Matching Model", presented at the Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020, pp. 4777–4786.

[24] J. Lu, J. Yang, D. Batra, and D. Parikh, "Neural Baby Talk", in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, Jun. 2018, pp. 7219–7228. DOI: `10.1109/CVPR.2018.00754`.

[25] B. Demirel, R. G. Cinbis, and N. Ikizler-Cinbis, "Image Captioning with Unseen Objects", in *BMVC*, 2019.

[26] B. Demirel and R. G. Cinbis, *Detection and Captioning with Unseen Object Classes*, Aug. 13, 2021. DOI: `10.48550/arXiv.2108.06165`. arXiv: `2108.06165 [cs]`.

[27] X. Chen, H. Fang, T.-Y. Lin, *et al.*, "Microsoft COCO Captions: Data Collection and Evaluation Server", *ArXiv*, 2015.

[28] P. Sharma, N. Ding, S. Goodman, and R. Soricut, "Conceptual Captions: A Cleaned, Hypernymed, Image Alt-text Dataset For Automatic Image Captioning", in *ACL*, 2018. DOI: `10.18653/v1/P18-1238`.

[29] V. Ordonez, G. Kulkarni, and T. Berg, "Im2Text: Describing Images Using 1 Million Captioned Photographs", in *Advances in Neural Information Processing Systems*, vol. 24, Curran Associates, Inc., 2011.

[30] L. Zhou, H. Palangi, L. Zhang, H. Hu, J. J. Corso, and J. Gao, "Unified Vision-Language Pre-Training for Image Captioning and VQA", *AAAI*, 2020. DOI: `10.1609/AAAI.V34I07.7005`.

[31] X. Yang, H. Zhang, and J. Cai, "Learning to Collocate Neural Modules for Image Captioning", *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019. DOI: `10.1109/ICCV.2019.00435`.

[32] C.-Y. Ma, Y. Kalantidis, G. AlRegib, P. Vajda, M. Rohrbach, and Z. Kira, "Learning to Generate Grounded Visual Captions Without Localization Supervision", in *ECCV*, 2020. DOI: `10.1007/978-3-030-58523-5_21`.

[33] T. Ridnik, G. Sharir, A. Ben-Cohen, E. Ben-Baruch, and A. Noy, *ML-Decoder: Scalable and Versatile Classification Head*, Dec. 31, 2021. DOI: `10.48550/arXiv.2111.12933`. arXiv: `2111.12933 [cs]`.

[34] T. Ridnik, H. Lawen, A. Noy, E. Ben, B. G. Sharir, and I. Friedman, "TRes-Net: High Performance GPU-Dedicated Architecture", in *2021 IEEE Winter Conference on Applications of Computer Vision (WACV)*, Jan. 2021, pp. 1399–1408. DOI: `10.1109/WACV48630.2021.00144`.

[35] P. Anderson, X. He, C. Buehler, *et al.*, "Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering", presented at the Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018, pp. 6077–6086.

[36] L. Wang, Z. Bai, Y. Zhang, and H. Lu, "Show, Recall, and Tell: Image Captioning with Recall Mechanism", *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 07, pp. 12 176–12 183, 07 Apr. 3, 2020, ISSN: 2374-3468. DOI: `10.1609/aaai.v34i07.6898`.

[37] R. Girshick, "Fast R-CNN", in *2015 IEEE International Conference on Computer Vision (ICCV)*, Dec. 2015, pp. 1440–1448. DOI: `10.1109/ICCV.2015.169`.

[38] R. Krishna, Y. Zhu, O. Groth, *et al.*, "Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations", *International Journal of Computer Vision*, vol. 123, no. 1, pp. 32–73, May 1, 2017, ISSN: 0920-5691. DOI: `10.1007/s11263-016-0981-7`.

[39] Y. Qin, J. Du, Y. Zhang, and H. Lu, "Look Back and Predict Forward in Image Captioning", in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2019, pp. 8359–8367. DOI: `10.1109/CVPR.2019.00856`.

[40] L. Ke, W. Pei, R. Li, X. Shen, and Y.-W. Tai, "Reflective Decoding Network for Image Captioning", in *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, Oct. 2019, pp. 8887–8896. DOI: `10.1109/ICCV.2019.00898`.

[41] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, "Aggregated Residual Transformations for Deep Neural Networks", in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 5987–5995. DOI: `10.1109/cvpr.2017.634`.

[42] S. Shao, Z. Li, T. Zhang, *et al.*, "Objects365: A Large-Scale, High-Quality Dataset for Object Detection", in *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, Oct. 2019, pp. 8429–8438. DOI: `10.1109/ICCV.2019.00852`.

[43] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, "Language Models are Unsupervised Multitask Learners", *undefined*, 2019.

[44] Q. Wen, T. Zhou, C. Zhang, *et al.*, *Transformers in Time Series: A Survey*, Mar. 7, 2022. DOI: `10.48550/arXiv.2202.07125`. arXiv: `2202.07125 [cs, eess, stat]`.

[45] H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, and H. Jegou, "Training data-efficient image transformers & distillation through attention", in *Proceedings of the 38th International Conference on Machine Learning*, PMLR, Jul. 1, 2021, pp. 10 347–10 357.

[46] J. Lu, D. Batra, D. Parikh, and S. Lee, "ViLBERT: Pretraining Task-Agnostic Visiolinguistic Representations for Vision-and-Language Tasks", in *Advances in Neural Information Processing Systems*, vol. 32, Curran Associates, Inc., 2019.

[47] A. Radford and K. Narasimhan, "Improving Language Understanding by Generative Pre-Training", *undefined*, 2018.

[48] H. W. Kuhn, "The Hungarian method for the assignment problem", *Naval Research Logistics Quarterly*, vol. 2, no. 1-2, pp. 83–97, 1955, ISSN: 1931-9193. DOI: `10.1002/nav.3800020109`.

[49] R. Jonker and A. Volgenant, "A shortest augmenting path algorithm for dense and sparse linear assignment problems", *Computing*, vol. 38, no. 4, pp. 325–340, Dec. 1, 1987, ISSN: 1436-5057. DOI: `10.1007/BF02278710`.

[50] D. P. Bertsekas, "The auction algorithm: A distributed relaxation method for the assignment problem", *Annals of Operations Research*, vol. 14, no. 1, pp. 105–123, Dec. 1, 1988, ISSN: 1572-9338. DOI: `10.1007/BF02186476`.

[51] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You Only Look Once: Unified, Real-Time Object Detection", in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2016, pp. 779–788. DOI: `10.1109/CVPR.2016.91`.

[52] `https://vizwiz.org/`.

[53] X. Gong, H. Zhu, Y. Wang, *et al.*, "Multiple Transformer Mining for VizWiz Image Caption", p. 2,

[54] A. Mathews, L. Xie, and X. He, "SentiCap: generating image descriptions with sentiments", in *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*, ser. AAAI'16, Phoenix, Arizona: AAAI Press, Feb. 12, 2016, pp. 3574–3580.

[55] C. Gan, Z. Gan, X. He, J. Gao, and L. Deng, "StyleNet: Generating Attractive Visual Captions with Styles", in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jul. 2017, pp. 955–964. DOI: `10.1109/CVPR.2017.108`.

[56] T. Chen, Z. Zhang, Q. You, *et al.*, ""Factual" or "Emotional": Stylized Image Captioning with Adaptive Learning and Attention", in *ECCV*, 2018. DOI: `10.1007/978-3-030-01249-6_32`.

[57] L. Guo, J. Liu, P. Yao, J. Li, and H. Lu, "MSCap: Multi-Style Image Captioning With Unpaired Stylized Text", in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2019, pp. 4199–4208. DOI: `10.1109/CVPR.2019.00433`.

[58] A. Tran, A. Mathews, and L. Xie, "Transform and Tell: Entity-Aware News Image Captioning", in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2020, pp. 13 032–13 042. DOI: `10.1109/CVPR42600.2020.01305`.

[59] J. Krause, J. Johnson, R. Krishna, and L. Fei-Fei, "A Hierarchical Approach for Generating Descriptive Image Paragraphs", in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jul. 2017, pp. 3337–3345. DOI: `10.1109/CVPR.2017.356`.

[60] L. Melas-Kyriazi, A. Rush, and G. Han, "Training for Diversity in Image Paragraph Captioning", in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, Brussels, Belgium: Association for Computational Linguistics, Oct. 2018, pp. 757–761. DOI: `10.18653/v1/D18-1084`.

[61] Z.-J. Zha, D. Liu, H. Zhang, Y. Zhang, and F. Wu, "Context-Aware Visual Policy Network for Fine-Grained Image Captioning", *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 2, pp. 710–722, 2022, ISSN: 1939-3539. DOI: `10.1109/TPAMI.2019.2909864`.

[62] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A large-scale hierarchical image database", in *2009 IEEE Conference on Computer Vision and Pattern Recognition*, Jun. 2009, pp. 248–255. DOI: `10.1109/CVPR.2009.5206848`.

[63] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition", *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, Nov. 1998, ISSN: 1558-2256. DOI: `10.1109/5.726791`.

[64] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet Classification with Deep Convolutional Neural Networks", in *Advances in Neural Information Processing Systems*, vol. 25, Curran Associates, Inc., 2012.

[65] S. Ioffe and C. Szegedy, "Batch normalization: accelerating deep network training by reducing internal covariate shift", in *Proceedings of the 32nd International Conference on International Conference on Machine Learning - Volume 37*, ser. ICML'15, Lille, France: JMLR.org, Jul. 6, 2015, pp. 448–456.

[66] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan, "Show and tell: A neural image caption generator", in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2015, pp. 3156–3164. DOI: `10.1109/CVPR.2015.7298935`.

[67] A. Karpathy and L. Fei-Fei, "Deep visual-semantic alignments for generating image descriptions", in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2015, pp. 3128–3137. DOI: `10.1109/CVPR.2015.7298932`.

[68] S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory", *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, Nov. 1997, ISSN: 0899-7667. DOI: `10.1162/neco.1997.9.8.1735`.

[69] J. Donahue, L. A. Hendricks, M. Rohrbach, *et al.*, "Long-Term Recurrent Convolutional Networks for Visual Recognition and Description", *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 4, pp. 677–691, Apr. 2017, ISSN: 1939-3539. DOI: `10.1109/TPAMI.2016.2599174`.

[70] J. Mao, W. Xu, Y. Yang, J. Wang, and A. Yuille, "Deep Captioning with Multimodal Recurrent Neural Networks (m-RNN)", *ICLR*, 2015.

[71] S. J. Rennie, E. Marcheret, Y. Mroueh, J. Ross, and V. Goel, "Self-Critical Sequence Training for Image Captioning", in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jul. 2017, pp. 1179–1195. DOI: `10.1109/CVPR.2017.131`.

[72] K. Simonyan and A. Zisserman, "Very Deep Convolutional Networks for Large-Scale Image Recognition", *ICLR*, 2015.

[73] D. Bahdanau, K. Cho, and Y. Bengio, "Neural Machine Translation by Jointly Learning to Align and Translate", *ICLR*, 2015.

[74] K. Xu, J. Ba, R. Kiros, *et al.*, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", in *Proceedings of the 32nd International Conference on Machine Learning*, PMLR, Jun. 1, 2015, pp. 2048–2057.

[75] J. Lu, C. Xiong, D. Parikh, and R. Socher, "Knowing When to Look: Adaptive Attention via a Visual Sentinel for Image Captioning", in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jul. 2017, pp. 3242–3250. DOI: `10.1109/CVPR.2017.345`.

[76] X. Chen, L. Ma, W. Jiang, J. Yao, and W. Liu, "Regularizing RNNs for Caption Generation by Reconstructing the Past with the Present", in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, Jun. 2018, pp. 7995–8003. DOI: 10.1109/CVPR.2018.00834.

[77] L. Chen, H. Zhang, J. Xiao, *et al.*, "SCA-CNN: Spatial and Channel-Wise Attention in Convolutional Networks for Image Captioning", *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. DOI: 10.1109/CVPR.2017.667.

[78] T. Lin, Y. Wang, X. Liu, and X. Qiu, "A Survey of Transformers", *ArXiv*, 2021.

[79] G. Li, L. Zhu, P. Liu, and Y. Yang, "Entangled Transformer for Image Captioning", in *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, Oct. 2019, pp. 8927–8936. DOI: 10.1109/ICCV.2019.00902.

[80] L. Huang, W. Wang, J. Chen, and X.-Y. Wei, "Attention on Attention for Image Captioning", in *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, Oct. 2019, pp. 4633–4642. DOI: 10.1109/ICCV.2019.00473.

[81] S. Herdade, A. Kappeler, K. Boakye, and J. Soares, "Image Captioning: Transforming Objects into Words", in *NeurIPS*, 2019.

[82] L. Guo, J. Liu, X. Zhu, P. Yao, S. Lu, and H. Lu, "Normalized and Geometry-Aware Self-Attention Network for Image Captioning", in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2020, pp. 10 324–10 333. DOI: 10.1109/CVPR42600.2020.01034.

[83] M. Cornia, M. Stefanini, L. Baraldi, and R. Cucchiara, "Meshed-Memory Transformer for Image Captioning", presented at the Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020, pp. 10 578–10 587.

[84] W. Liu, S. Chen, L. Guo, X. Zhu, and J. Liu, "CPTR: Full Transformer Network for Image Captioning", *ArXiv*, 2021.

[85] Z. Wang, J. Yu, A. W. Yu, Z. Dai, Y. Tsvetkov, and Y. Cao, "SimVLM: Simple Visual Language Model Pretraining with Weak Supervision", presented at the International Conference on Learning Representations, May 11, 2022.

[86] J. Aneja, A. Deshpande, and A. G. Schwing, "Convolutional Image Captioning", in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, Jun. 2018, pp. 5561–5570. DOI: 10.1109/CVPR.2018.00583.

[87] J. Gehring, M. Auli, D. Grangier, D. Yarats, and Y. N. Dauphin, "Convolutional Sequence to Sequence Learning", in *Proceedings of the 34th International Conference on Machine Learning*, PMLR, Jul. 17, 2017, pp. 1243–1252.

[88] C. Rashtchian, P. Young, M. Hodosh, and J. Hockenmaier, "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*, ser. CSLDAMT '10, USA: Association for Computational Linguistics, Jun. 6, 2010, pp. 139–147.

[89] M. Hodosh, P. Young, and J. Hockenmaier, "Framing Image Description as a Ranking Task: Data, Models and Evaluation Metrics", *Journal of Artificial Intelligence Research*, vol. 47, pp. 853–899, Aug. 30, 2013, ISSN: 1076-9757. DOI: 10.1613/jair.3994.

[90] S. Changpinyo, P. Sharma, N. Ding, and R. Soricut, "Conceptual 12M: Pushing Web-Scale Image-Text Pre-Training To Recognize Long-Tail Visual Concepts", presented at the Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2021, pp. 3558–3568.

[91] H. Agrawal, K. Desai, Y. Wang, *et al.*, "Nocaps: novel object captioning at scale", in *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, Oct. 2019, pp. 8947–8956. DOI: 10.1109/ICCV.2019.00904.

[92] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman, "The Pascal Visual Object Classes (VOC) Challenge", *International Journal of Computer Vision*, vol. 88, no. 2, pp. 303–338, Jun. 1, 2010, ISSN: 1573-1405. DOI: 10.1007/s11263-009-0275-4.

[93] https://flickr.com/.

[94] https://cocodataset.org/#captions-leaderboard.

[95] B. Thomee, D. A. Shamma, G. Friedland, *et al.*, "YFCC100M: the new data in multimedia research", *Communications of the ACM*, vol. 59, no. 2, pp. 64–73, Jan. 25, 2016, ISSN: 0001-0782. DOI: 10.1145/2812802.

[96] https://nocaps.org/.

[97] D. Gurari, Q. Li, A. Stangl, *et al.*, "VizWiz Grand Challenge: Answering Visual Questions from Blind People", *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018. DOI: 10.1109/CVPR.2018.00380.

[98] M. Kilickaya, A. Erdem, N. Ikizler-Cinbis, and E. Erdem, "Re-evaluating Automatic Metrics for Image Captioning", in *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, Valencia, Spain: Association for Computational Linguistics, Apr. 2017, pp. 199–209.

[99] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, "BLEU: a method for automatic evaluation of machine translation", in *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*, ser. ACL '02, USA: Association for Computational Linguistics, Jul. 6, 2002, pp. 311–318. DOI: 10.3115/1073083.1073135.

[100] S. Banerjee and A. Lavie, "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*, Ann Arbor, Michigan: Association for Computational Linguistics, Jun. 2005, pp. 65–72.

[101] C.-Y. Lin, "ROUGE: A Package for Automatic Evaluation of Summaries", in *Text Summarization Branches Out*, Barcelona, Spain: Association for Computational Linguistics, Jul. 2004, pp. 74–81.

[102] R. Vedantam, C. L. Zitnick, and D. Parikh, "CIDEr: Consensus-based image description evaluation", in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2015, pp. 4566–4575. DOI: 10.1109/CVPR.2015.7299087.

[103] P. Anderson, B. Fernando, M. Johnson, and S. Gould, "SPICE: Semantic Propositional Image Caption Evaluation", in *ECCV*, 2016. DOI: `10.1007/978-3-319-46454-1_24`.

[104] L. Bergroth, H. Hakonen, and T. Raita, "A survey of longest common subsequence algorithms", in *Proceedings Seventh International Symposium on String Processing and Information Retrieval. SPIRE 2000*, Sep. 2000, pp. 39–48. DOI: `10.1109/SPIRE.2000.878178`.

[105] J. Johnson, R. Krishna, M. Stark, *et al.*, "Image retrieval using scene graphs", in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2015, pp. 3668–3678. DOI: `10.1109/CVPR.2015.7298990`.

[106] G. A. Miller, "WordNet: a lexical database for English", *Communications of the ACM*, vol. 38, no. 11, pp. 39–41, Nov. 1, 1995, ISSN: 0001-0782. DOI: `10.1145/219717.219748`.

[107] D. A. Hudson and C. D. Manning, "GQA: A New Dataset for Real-World Visual Reasoning and Compositional Question Answering", *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. DOI: `10.1109/CVPR.2019.00686`.

[108] Y. Goyal, T. Khot, D. Summers-Stay, D. Batra, and D. Parikh, "Making the V in VQA Matter: Elevating the Role of Image Understanding in Visual Question Answering", *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. DOI: `10.1109/CVPR.2017.670`.

# APPENDIX A

# VISUAL EXAMPLES

In this part, visual examples for the comparison of the baseline method and the proposed method are given as an extension to the Section 5.4. In Section A.1, good examples where the proposed method successfully describes the scene with attribute-wise rich sentences are provided. In Section A.2, some problematic examples where the proposed method generates imperfect captions (hallucinated objects, missing details) are provided.

## A.1   Successful Examples

The proposed method generates very good captions with details and additional attribute information for the examples in this section. It generally describes the scene better than the baseline method.

| Baseline | – A bowl of food with a spoon in it. |
|---|---|
| **Proposed Method** | – A **white** bowl **filled** with **chopped** veggies and a **wooden** spoon. |
| **Ground Truth** | – A pan filled with veggies and a block of butter. |
| | – Butter or cheese on top of some vegetables in a pan. |
| | – Spoon in a bowl of chopped vegetables with butter. |
| | – A pan with butter and other ingredients for a dish. |
| | – Pot with variety of chopped vegetables, big chunk of butter and spoon. |



| Baseline | – An old truck sitting in the middle of a forest. |
|---|---|
| **Proposed Method** | – An **old rusted** truck **sitting** in the middle of **yellow** flowers. |
| **Ground Truth** | – A old truck with a busted window in the tall bushes. |
| | – A rusty old truck sitting in an overgrown field. |
| | – A rusted out truck parked next to some yellow flowers. |
| | – Old pick-up with flowers growing in front of it. |
| | – A truck is shown decaying among flowers without a window. |

| Baseline | – A couple of zebra standing next to each other on a field. |
|---|---|
| **Proposed Method** | – A group of zebra standing next to each other on a **dry grass** field. |
| **Ground Truth** | – Zebra grazing on dry grass in a fried up field.<br>– A mother zebra and her young on a grassy plain.<br>– A zebra and a young zebra in a field<br>– A zebra with its offspring grazing in the desert.<br>– Two Zebras in a brown field grazing and eating. |



| Baseline | – A group of buses that are sitting in the street. |
|---|---|
| **Proposed Method** | – A group of **double decker** buses **parked** in front of a building. |
| **Ground Truth** | – A tall building with four double decker buses driving along a parking lot.<br>– A large long bus on a city street.<br>– Two white double decker buses on a street.<br>– Double-decker buses sit at the curb in front of an old building.<br>– Double-Decker buses line up at a bus stop. |

| Baseline | – A white toilet sitting next to a sink in a bathroom. |
|---|---|
| **Proposed Method** | – A bathroom with a **black** and **white checkered** floor. |
| **Ground Truth** | – A kitchen sink sitting next to a toilet. |
| | – A black and white, checkerboard bathroom with a red towel hanging up. |
| | – A red towel hanging in a black and white bathroom. |
| | – The bathroom is tiled black and white floors. |
| | – A bath room with a toilet a sink and a mirror. |

| | |
|---|---|
| **Baseline** | – A white door is open outside of a building. |
| **Proposed Method** | – An **old white** building on the side of a road. |
| **Ground Truth** | – A shabby garage stands next to a brown-red residential building. |
| | – An run down garage, next to a building. |
| | – An old white garage is next to a brown building in front of trees. |
| | – A white garage in disrepair next to a house. |
| | – A garage that is not in a good shape next to a building. |
| | – The large building appears to be a garage. |
| | – A tiny one car garage made of dented sheet metal in the middle of a driveway. |
| | – A garage contains a dented metallic side wall. |
| | – An old white aluminum windowless garage, in disrepair, next to a building. |
| | – A one car garage with galvanized steel outside walls. |

| Baseline | – A statue of a statue of a man's head. |
|---|---|
| Proposed Method | – A close up of a statue of a **shiny metal silver** head. |
| Ground Truth | – A bronze sculpture of a bust that is shiny. |
| | – With face distorted, the bust in bronze stares ahead. |
| | – A bronze sculpture with a strange expression stands on a brick surface. |
| | – The sculpture of a man made with cooper has the mouth open. |
| | – A bust of a bronzed man is shown with a brick background. |
| | – A statue of a man's face is outside near a brick way. |
| | – A bronze statue of a man who is making a face and has heavy eyebrows. |
| | – A bronze statue of a human face that is not good looking. |
| | – A statue of a man has a bald spot. |
| | – A bronze sculpture of a bearded japanese samurai. |

| | |
|---|---|
| **Baseline** | – A man walking a dog on a leash. |
| **Proposed Method** | – A man in **camouflage** with a dog in front of him. |
| **Ground Truth** | – A man dressed in camouflage next to a dog on a leash. |
| | – A military man holding the leash to a service dog. |
| | – A man in camoflage is holding the leash of his military dog. |
| | – A soldier standing with a leashed dog at his side. |
| | – A soldier standing with a dog on a leash. |
| | – A man in camo holding onto a dog on the road. |
| | – A man in a military camo uniform holding the leash of a dog. |
| | – A man in fatigues stands with his dog at his side. |
| | – A person, holding a leash, standing next to his dog. |
| | – A man in a uniform is sitting with a German Shepard dog on pavement. |

| **Baseline** | – A man in a baseball uniform standing in the woods. |
|---|---|
| **Proposed Method** | – A man in **camouflage** gear walking through a forest. |
| **Ground Truth** | – A man standing in uniform and sunglasses near some plants.<br>– A Russian soldier dressed in camouflage gear and goggles.<br>– The person is standing in the forest with sunglasses.<br>– A female person, wearing camouflage gear and sunglasses is standing, surrounded by brown and green foliage.<br>– A person in a camouflage outfit and sunglasses stands in some tall brush with his right hand on his chest.<br>– A man in camouflage and black sunglasses standing in between some plants.<br>– Man in military clothing and sunglasses walking through an area with bushes.<br>– A military man in the woods with hat and gloves on.<br>– A person with military clothing, helmet and sunglasses looks forward.<br>– A uniformed soldier standing in the bush. |

## A.2 Problematic Examples

Examples in this section contain problematic captions generated by the proposed method. The baseline method also generates poor captions for some of the examples. Generated captions might be problematic in the sense that some objects are hallucinated, some details (objects, attributes) in the image are missed.



| Baseline | – A wooden table topped with lots of different types of fruit. |
|---|---|
| **Proposed Method** | – A **wooden** table topped with lots of **ripe** bananas. |
| **Ground Truth** | – Bananas sitting on top of apples, pears, and other fruits. |
| | – A large plate topped with lots of fresh fruit. |
| | – A platter full of bananas, apples, oranges, peaches, pears and limes. |
| | – There is a large platter of fruit on the table. |
| | – A tray heaped with fruit including bananas apples oranges and limes. |

| | |
|---|---|
| **Baseline** | – A bike parked next to a bench near a river. |
| **Proposed Method** | – A **red** bicycle **parked** next to a **metal** pole. |
| **Ground Truth** | – A bike leaning on a metal fence next to some flowing water. |
| | – A bicycle parked next to a flooded river. |
| | – A bicycle chained to a rail near a flooded area. |
| | – A bicycle leans against a fence near some flood waters. |
| | – A bicycle parked on a rail next to some flooding water. |

| Baseline | – A close up of a watch on the side of a street. |
|---|---|
| **Proposed Method** | – A **red** and **white digital** clock **mounted** to the side of a road. |
| **Ground Truth** | – A digital stop watch with a red, silver, and black face.<br>– A runner's watch displays thier performance, and times on this man's wrist.<br>– A late watch with some letters on it.<br>– A sports watch used for many different timing needs.<br>– The person has a freckled arm and is wearing a fitness watch.<br>– A Forerunner305 watch that is on a person's wrist.<br>– An arm displays a sports watch with several different functions.<br>– A persons arm with a garmin brand watch on it.<br>– A person has a watch on their wrist.<br>– A watch rests on the arm of a person. |

| **Baseline** | – A man is walking through a train station at night. |
|---|---|
| **Proposed Method** | – A **blurry** image of a person with a skateboard. |
| **Ground Truth** | – A person stands in a bowling alley filled with different neon lights.<br>– A man in a black and blue shirt bowling.<br>– A person is standing by a bowling alley.<br>– A small boy standing in front of bowling lanes, each containing different colored pins.<br>– A young boy looking down bowling lane at the lights.<br>– A boy standing at the edge of the bowling lane looking ahead.<br>– A person bowls in a neon bowling alley.<br>– In a spray of neon lights, a young boy prepares to throw the bowling ball down the lane.<br>– A person bowling at a bowling alley.<br>– Someone stands at the end of the bowling alley. |

| Baseline | – Two men are sitting on top of a green train. |
|---|---|
| **Proposed Method** | – Three men **sitting** on top of a **military** helicopter. |
| **Ground Truth** | – Two men sitting on top of an army tank. |
| | – Two men are riding on a large tank in the woods. |
| | – Two men sit atop a tank that is painted with camouflage. |
| | – Two older men sitting on a camouflage tank. |
| | – Two men are sitting on a tank that is camo. |
| | – A small tank carrying two people through the woods. |
| | – Two men ride on top of a tank. |
| | – Two men are riding a tank vehicle outside. |
| | – Two men sitting on a military tank in the woods. |
| | – Two men in a military style vehicle with a small cannon it it. |