Language Image Transformer

by

Raghavendran Ramakrishnan

A Dissertation Presented in Partial Fulfillment
of the Requirements for the Degree
Master of Science

Approved April 2020 by the
Graduate Supervisory Committee:

Sethuraman Panchanathan, Co-Chair
Hemanth Venkateswara, Co-Chair
Troy McDaniel

ARIZONA STATE UNIVERSITY

May 2020

ABSTRACT

Humans perceive the environment using multiple modalities like vision, speech (language), touch, taste, and smell. The knowledge obtained from one modality usually complements the other. Learning through several modalities helps in constructing an accurate model of the environment. Most of the current vision and language models are modality-specific and, in many cases, extensively use deep-learning based attention mechanisms for learning powerful representations. This work discusses the role of attention in associating vision and language for generating shared representation. Language Image Transformer (LIT) is proposed for learning multi-modal representations of the environment. It uses a training objective based on Contrastive Predictive Coding (CPC) to maximize the Mutual Information (MI) between the visual and linguistic representations. It learns the relationship between the modalities using the proposed cross-modal attention layers. It is trained and evaluated using captioning datasets, MS COCO, and Conceptual Captions. The results and the analysis offers a perspective on the use of Mutual Information Maximization (MIM) for generating generalizable representations across multiple modalities.

i

TABLE OF CONTENTS

LIST OF TABLES

## LIST OF FIGURES

Chapter 1

INTRODUCTION

We have achieved success in learning from a labeled data set using supervised tasks in vision and language [Devlin *et al.* (2018); He *et al.* (2015)]. But there is limited success in solving tasks in different modalities by learning a common representation. Extracting shared high-level features between modalities [Oord *et al.* (2018)] is one of the key challenges in representation learning. Several recent works have focused on generating a shared representation that can perform well in a variety of downstream tasks.

Vision and language have individually received wider attention in recent years. Several important contributions [He *et al.* (2016); Devlin *et al.* (2018)] have been made in both modalities. In many cases, development in one modality has transformed well and received adoption in another [Sun *et al.* (2019)]. Hence, there is a high degree of similarity between the model architectures in vision and language. This similarity motivates us to hypothesize that there could be benefits in developing multi-modal representation that can transform well between several modalities with a minimal overhead of specialization.

Human language has evolved specialized words for representing high-level objects, events, and actions. Hence, it is beneficial for models to exploit the language's use for developing a high-level representation of the objects in the visual context [Oord *et al.* (2018)]. The objects that are present in the visual context can effectively be represented using sentences in the language domain.

Image captioning is the ability of the model to automatically generate captions that describe the scene presented in the image. It reflects the ability of humans to

compress a large amount of visual information in descriptive language. This efficient compression technique is rooted in the ability to understand the scene and extract the relevant aspects of the scene. Hence the task forces the model to detect the important objects in the image and understand their relationship. The model learns to represent this relationship in natural language.

There have been several efforts to develop representations that transform well from vision to language in the field of image captioning [Johnson *et al.* (2015)]. The cost of acquiring and captioning the images has led to increased demand for automated extraction and annotation of the images [Sharma *et al.* (2018)]. Automatic image captioning contributes to reducing human bias in the annotations. It exploits the availability of a large amount of visual content on the internet to increase the diversity of the images in the corpus.

Several downstream tasks [Agrawal *et al.* (2015); Wang *et al.* (2017)] in vision and language have focused on creating multi-modal embeddings that can bridge the semantic gap between the two modalities. Many works (Lee et al., 2018, Yu.et.al., 2018) that have achieved success in the shared downstream tasks, have studied the latent alignment between the modalities. In case of vision and language models, it can be alignment between image regions and text. While many of these models have achieved state-of-art in their respective tasks, they have used very task-specific design approaches to accomplish it.

These factors motivated me to develop a model that can generate multi-modal representation in the shared embedding space between vision and language. I have adopted a transformer architecture and used attention mechanism across the network. I have trained the model to generate suitable captions for the presented image.

I utilize learning techniques from self-supervised learning that have achieved success in creating visual-linguistic representations. The techniques use 'proxy' training tasks to learn the deep semantic relationship between the data sources. The proxy training tasks leverage the structure within the data to generate pseudo supervised objectives. The model can learn the semantic relationship within the data by training on the objectives. The self-supervised learning techniques have shown the greatest impact in language models. BERT model uses a masking technique to learn the semantic relationship between the tokens in the input sequence. I also use a similar objective to train the model on image captioning task.

Chapter 2

RELATED WORKS

## 2.1  Language Modeling

Language models are designed to be efficient in handling sequential data such as words, tokens, and sentences. Sequence-to-sequence based architectures have achieved great success in language modeling. Recurrent Neural Networks (RNN) [Schuster and Paliwal (1997)] and Long Short-Term Memory (LSTM) [Hochreiter and Schmidhuber (1997)] models and their variants have traditionally been used in designing language models. There have been several works to address the issues with space and time complexity due to the sequential processing of input in sequence-to-sequence architectures [Schuster and Paliwal (1997)]. One of the recent works in this domain is Transformer by Vaswani *et al.* (2017). The authors managed to achieve state-of-art results in machine translation. It significantly reduced the computations by designing the transformer model to handle the sequences in parallel and leveraging the computation power offered by GPUs. Several models were introduced that use transformers as their primary blocks [Devlin *et al.* (2018); Sun *et al.* (2019)]. Our work aligns with this field of study, by adopting a transformer-based architecture. We have extended the transformer model to use visual embedding for sequence generation.

### 2.1.1  Attention in Language Models

Sutskever *et al.* (2014) introduced the seq2seq model architecture for Neural Machine Translation. The seq2seq model consists of a decoder stacked sequentially on top of an encoder. The model computes a probability $P(E|F)$ of output $E$ given an

input $F$. If the input and output are sentences belonging to different languages, then the process is Neural Machine Translation.

In a sequence-to-sequence architecture, encoder learns to generate a hidden state representation of the input sequence. The decoder network utilizes the hidden state-representation to create the output sequence. The encoder and the decoder networks learn the function map between the input and output sequences.

I will initially discuss the problems associated with the encoder-decoder architecture in the context of Neural Machine Translation. Later, i will also discuss how attention resolves many issues with seq2seq architecture.

Ideally, a seq2seq architecture that is large and is well-trained can perform machine translation accurately. But multiple issues can arise in a seq2seq model.

It is difficult for the model to capture long-range dependencies between the words in sentences of longer length. The model tries to embed the information from the sequences of varying lengths with the hidden-vector of a fixed size. It results in poor representation as the hidden-vector doesn't accurately fit the length of the input sequence. While small dimensions of hidden-vector are insufficient for lengthy sequences, oversized dimensions can be excessive. As the amount of data available for training is limited, large networks may overfit to the training sequences.

In attention, the model doesn't learn a single hidden-vector for the input sentence. Instead, the model learns a vector for every word and refers to them at every step in decoding. So the number of vectors learned is equal to the number of words in the input sequence. Attention, therefore, provides an efficient representation of the input sequence.

## 2.2 Attention in Bi-directional RNN

In the case of bi-directional RNN, the encoder first generates representation from left to right ($\vec{h}_j^f$) and right to left ($\overleftarrow{h}_j^f$). The encoder outputs the concatenated vector of both embedding $h_j^f = [\vec{h}_j^f, \overleftarrow{h}_j^f]$. By concatenating all the bi-directional vectors, we obtain a single matrix, $H_j = concat[h_1^f \ldots h^f{}_{\text{mod } F}]$. In matrix $H_j$ each column corresponds to one word in the input sentence.

The number of columns in matrix $H_j$ depends on the length of the input sentence. Hence, the size of the vector can vary based on the input sentence. It is essential to generate a single vector $c_t$ by combining the columns of the matrix $H$.

$$c_t = H^f \alpha_t \tag{2.1}$$

where $\alpha_t$ is the attention vector .The attention vector is the measure of focus required on a particular word while predicting the next word in the output sentence.

## 2.3 Types of Attention

Bahdanau *et al.* (2014) introduced the Attention mechanism for Neural Machine Translation. Attention is an elegant technique that shares characteristics with the natural perception mechanism (Rensink, 2000;Corbetta Shulman, 2002). Instead of compressing the input into a static representation, attention dynamically selects the salient features in the input as required [Xu *et al.* (2015)]. Attention can be defined as a mechanism of assigning weights to the inputs by computing their relative importance score. It is widely applied to solve tasks in various domains including language and vision [Vaswani *et al.* (2017); Devlin *et al.* (2018); Kaiser and Bengio (2016); Sun *et al.* (2019)]. Several types of attention such as Bottom-Up and Top-Down[Anderson *et al.* (2018)], Self-attention [Vaswani *et al.* (2017)] have been introduced over the

years. Self-attention relates several positions in the same input for generating a representation that captures the relationship between the elements in the same input. My model extensively uses Self-attention as part of the transformer architecture to generate a representation that captures the semantics of the input image and caption.

## 2.4   Image/Video Captioning

Image captioning models traditionally used CNN and RNN based architecture [Xu *et al.* (2015)]. CNN was used to encode the image features from the given image. The encoder was followed by an RNN, which generated text sequences or captions conditioned on the image features generated by CNN.

Humans' ability to describe a scene using diverse contexts and references arise from their capability to strongly relate to the various aspects of the scene. Understanding the semantics of the scene remains a challenging task for deep learning models. First, there is significant progress in object detection. The models are trained to capture the various objects present in a scene and label them into a set of well- defined classes. Second, there is rapid progress in language representation. The label space has expanded to include sentences that can capture the underlying context in the image. DenseCap [Johnson *et al.* (2015)] has unified the progress in the two tasks using a dense captioning task. The model was trained to predict a set of descriptions by detecting the objects present in the image.

Xu *et al.* (2015) studied the role of attention in image captioning. The authors introduced a soft deterministic attention mechanism and trained it using backpropagation. They also introduced a hard attention mechanism that was trained by maximizing an approximate variational lower bound. The Encoder model extracted the convolutional features present in the image. These features were fed to a Decoder using LSTM (Long Short Term Memory) network. The decoder used attention to

learn and attend to the appropriate positions in the input image and predict their captions. Our model uses an Encoder-Decoder based architecture

## 2.5 Bidirectional Encoder Representations from Transformer (BERT)

BERT [Devlin *et al.* (2018)] learns deep bidirectional representations from the unlabeled text by jointly conditioning on the left and right context in the input sequence. It uses a proxy task called masked language modelling to learn the left and right context in the sequence. In masked language modeling, tokens in the input sequence are randomly selected and replaced with the mask. The objective is to predict the vocabulary id of the masked token using the context. BERT achieved state-of-art in eleven NLP tasks. The use of pre-training and self-attention are two reasons for the effectiveness of BERT. The model is pre-trained on the Books Corpus (800M words)Zhu *et al.* (2015) and English Wikipedia (2,500M words). Two types of proxy tasks are used for pre-training. The first task is the masked language modelling. The second task is predicting if the two given sentences follow each other in the text corpus. After pre-training the model, it is fine-tuned for the downstream task. During fine-tuning, task-specific inputs and outputs are fed into BERT and all parameters are fine-tuned from end-to-end. At the output, the token representation from BERT is fed into the output layer for the specific downstream task. My model uses a similar transformer-based sequence-to-sequence architecture and is trained using MLM proxy task. While BERT is completely designed for language modeling, My model is adopted to visual inference.

## 2.6 VideoBERT

VideoBERT [Sun *et al.* (2019)] is a joint visual-linguistic model designed to learn high-level features from YouTube videos through self-supervision from the audio. It

uses vector quantization and Automatic Speech Recognition (ASR) to generate visual-linguistic tokens. The authors achieved state-of-art in video captioning. The visual data is transformed into a discrete sequence of tokens using hierarchical vector quantization and is combined with the linguistic sentence derived from audio using ASR. The model is trained using two proxy tasks. The first task is to predict the masked token. The second task is estimating the alignment between visual and linguistic tokens. Similar to BERT, The trained model is then fine-tuned for downstream tasks. The model also learns a joint probability distribution between the two modalities.

## 2.7   Vision-and-Language BERT (ViLBERT)

ViLBERT [Lu *et al.* (2019)] introduced separate streams for vision and language processing. The model used co-attentional transformer blocks to share the parameters between the vision and language streams. The motivation is that different modalities may require distinctive pre-processing steps due to their complexity. The authors also propose that the pre-trained weights cannot accommodate for the higher number of visual tokens and may corrupt the BERT language model. The co-attentional blocks produce attention-pooled features for each modality conditioned on the other. The first task is predicting the masked portions in the image and caption. The second task is to predict if the image and text segment align with each other. The authors evaluate the model using four different downstream tasks. Unlike ViLBERT, We use unique attention mechanism called cross-modal attention. In ViLBERT, A pre-trained object detection network extracts the image features. Instead, My model functions on input pixels. I hypothesize that operating directly on input pixels helps the model to learn the semantic relationship between the pixel regions.

## 2.8 UNiversal Image Text Representation (UNITER)

UNITER [Chen *et al.* (2019)] introduced additional proxy tasks for pre-training the model for generating visual-language representation. In Masked Region Modelling (MRM), the masked regions of the image are reconstructed by conditioning on the text. The authors proposed three variants of MRM. They are Masked Region Feature Regression, Masked Region Classification using KL-divergence. While It used image embedder to encode image regions, it used a text embedder to embed text in a shared embedding space. The model used a transformer to generate a cross-modal embedding. The authors conditioned the masking on full observation of image/text instead of masking random positions in sequences. Our model conditions the generation of the masked region of captions on the image pixels. The authors feed the extracted features from an R-CNN into the model for learning the cross-modal embedding. Instead, i extract image features using attention-augmented convolution. I generate cross-modal embedding using the extracted features.

## 2.9 Stand-Alone Self-Attention in Vision Models

Convolution is extensively used in computer vision models for extracting visual features [He *et al.* (2016); Gehring *et al.* (2017)]. Many deep learning-based vision models have used convolution for their translation equivariance and weight sharing properties. One of the drawbacks with the convolution is its inability to efficiently model long-range dependencies. This is due to their poor scaling properties for large receptive fields. There have been many recent works [Ramachandran *et al.* (2019); Bello *et al.* (2019)] that have focused on augmenting convolution with attention to increase the content-based interactions. In Ramachandran *et al.* (2019), authors conclude through extensive studies that a self-attention based vision layer can act as

an effective stand-alone layer. Similar to convolution, the authors initially select a local region of pixels within a fixed spatial extent of the given pixel (*a memory block*). To limit the number of computations required by the self-attention layer, they select a fixed spatial extent around the given pixel.Consider a pixel $x_{ij}$ in position $ij$ in any channel of a given image. The attention is computed at the position $ij$ using equation (2.2).

$$y_{ij} = \sum_{a,b \in \mathcal{N}_k(i,j)} \text{softmax}_{ab}(q_{ij}^T k_{ab}) v_{ab} \qquad (2.2)$$

where query $q_{ij} = W_Q x_{ij}$, key $k_{ab} = W_K x_{ab}$ and value $v_{ab} = W_V x_{ab}$ are the linear transformation of the pixel $x_{ij}$ and its neighbouring pixels.

## 2.10   Attention Augmented Convolution

Bello *et al.* (2019) proposed a novel two-dimensional self-attention mechanism for replacing convolutions as a stand-alone layer. The authors proposed a combined framework using self-attention and convolution. While convolution allows translation invariance, self-attention contributes to an increased receptive field. Through experiments, the authors demonstrated that the combined framework applying augmented convolutional operators with self-attention delivered the best results in image classification. One limitation of self-attention is that the operation is permutation invariant. This makes it inefficient in modeling structured data like images. Hence the model uses a two-dimensional relative positional encoding. We use multiple layers of attention-augmented convolution to encode the images as attention maps. The attention maps are fed to the decoder. The output of the decoder is used for generating relevant captions for the given image.

Chapter 3

ATTENTION IN LANGUAGE MODEL - SENTIMENT ANALYSIS

## 3.1   Introduction

In this chapter, I will deviate from the current discussion on multi-modal representation and focus on the importance of attention in language models. I will consider a use case in sentiment analysis and compare the performance of traditional statistical models with the attention based models. I will demonstrate the improvement in prediction achieved using attention-based models such as BERT (Bi-directional Encoder Representation from Transformer) over traditional architectures. In this chapter, I will develop a model that can perform stock market sentiment prediction. I will consider 2 companies viz. Amazon and Apple. I will develop a model that can predict if the value of stock might increase or decrease. I utilize the sentiment from the news articles for the prediction. I will train a classifier on the sentence embeddings of the news articles and predict if the outcome will be an increase in stock price or decrease in stock price.

The possible prediction of a stock market direction may act as an unforeseen recommendation system for short-term investors and as an unforeseen financial warning methodology for long-term shareholders. The single most important factor in choosing any forecasting methods is forecasting accuracy. Research efforts in this direction to improve the accuracy of forecasting models are increasing since the last decade. The essential aim for each investor is to maximize profits on their investments. Over the years there have been many approaches to find a correlation between stock related twitter posts and a considerable surge of related stock prices following such

posts. However, a direct attempt to simply correlate both of above seems naive. Since our objective is to find significant news indicators that caused a boost in stock prices, we need a similar data source to assure ourselves that the directional stock prediction is reliable.

The financial analysts who invest in stock market are generally not aware of the stock market behavior. They constantly face the issue of trading as they cannot make a correct decision on the stocks to buy or the stocks to sell. The internet provides extensive resources to understand the current state of each company. Inferring the state from the available information is an difficult task. This will allow financial analysts to foresee the behavior of the stock that they are interested in and thus can act accordingly. The input to the model will be historical data from news sources and Apple/Amazon stock data. The prediction model will predict if the stock price of the company might rise/fall. We make a probabilistic prediction through one of the following methods.

- By evaluating the trading volume following the news announcement as an indicator of the impact of news on the stock price.

- By evaluating the diffusion rates and volumes of messages on different platforms which has the stock symbol and news links of interest present in it.

One of the very prominent sources of data is Twitter as it presents a considerably accurate platform to evaluate properties of such information diffusion and gauge their volumes. Now diffusion analysis using the three-sigma rule can be employed to investigate "viral Tweets" to create early-warning indicators that can intimate if a breakout started to emerge in its early stages. We find the URL links relevant to the breaking-news hour of Tweets and thus we can ascertain the second part of the experiment whether the information indicated by breaking Tweet volumes will contribute

to statistically considerable boost in the directional prediction accuracy for the prices of the associated stock symbols mentioned in the URL links. The characteristics of such an experimental system can be explained as follows: Recent computational advances have led to implementation of machine learning techniques for the predictive models in financial markets. In this chapter, I am using various Machine Learning models for the prediction task and compare their performance with attention based models.

## 3.2    Relevant Works

Commendable work is done by Cutler *et al.* (1989) in their publication of What moves stock? Though published in 1989, their impact analysis of news on stock price movement direction was admirable. Stock price reaction to news and no-news: drift and reversal after headlines discusses how news impact the drift and reversal in stock prices due to public news articles. Using a comprehensive database of headlines about individual companies, he examined monthly returns following public news. He compared them to stocks with similar returns, but no identifiable public news. The relation between the sentiment of news, earnings and return predictability. Tetlock *et al.* (2008) discussed how a firm's stock prices under-perform/ under-react to the public news embedded with negative words. Directional Prediction of Stock Prices using Breaking News on Twitter by Alostad and Davulcu (2015) most closely aligns with our goal to predict the stock prices using news articles. Wang *et al.* (2016) recommends use of Recurrent neural networks for small news articles. After cleaning and data processing, I extract relatively smaller volume of data containing only the information about relevant stocks.

## 3.3   Data Pre-processing

In order to achieve accurate results, the data must be well-formatted and usable by machine learning models. As we experiment with multiple models on same dataset, it should be in acceptable format for all models under consideration.

Data pre processing is done in mainly two stages.

### 3.3.1   Data Cleaning

Initially, i extracted paragraphs containing relevant stock news from news corpus for Amazon and Apple. I removed any non-English characters (other language or emojis). I used NLTK library developed by Loper and Bird (2002) to perform initial text pre-processing such as removal of stop words.

### 3.3.2   Data Labeling

I labelled each article as positive or negative class based on the increase or decrease in the stock value for the provided time period. The increase was given positive label.

## 3.4   Method

The primary goal of the project is to predict movement in the price of the stock using news articles and stock chart data. I make an assumption associating the rise/drop of the stock price to the sentiment of the article published in the same time frame. I label all the news articles that were published in the time frame as either positive/negative based on the stock price. The news articles with positive sentiment are labelled as Class '1' while news articles with negative sentiment are labelled as Class '0'.

### 3.4.1   Model

I initially perform classification using statistical Machine Learning models. I categorized news data into 2 classes: class 0 and class 1. Specifically, I used 4 supervised learning models and trained each model using Distributed Bag of words (D-BOW) and Distributed Memory (DM) respectively.

**Bi-directional Encoder Representations from Transformers (BERT)**

A BERT is a seq2seq model that uses a transformer-based architecture. It uses attention mechanism to learn relationship between words or sub-words in a sentence. The language embedding generated by BERT has shown to outperform existing models in a number of downstream tasks [Devlin *et al.* (2018)]. I use a pre-trained BERT representation for training a simple linear classifier to predict the sentiment present in the input sentence. There are two factors that contribute to BERT's performance. BERT model is pre-trained on a large corpus of text using self-supervised proxy tasks. The model uses masked-language-modelling where it is trained to predict a portion of the sentence using its context. Hence, the model learns to capture relationship between the context and the selected word effectively. Also, it helps the model to capture bi-directional context as the selected word can be present in any position in the sentence. Use of attention is another factor that contributes to BERT's performance. BERT uses multiple layers of transformer architecture. Previously, Vaswani *et al.* (2017) had shown that transformers generate powerful sequential representation.

### 3.5   Results

The tables 3.1 and 3.2 show the results obtained using different models on the Amazon and Apple data sets. I followed the steps presented in Section 3.3 to pre-

process and label the data. Initially, i applied traditional statistical models such as logistic regression, Support Vector Machine and Random Forest to study their performance in the sentiment prediction task. I tested the model with two types of sentence embeddings. I used word embedding generated using Distributed Memory and Distributed Bag-of-Words models [Le and Mikolov (2014)], two popular techniques for generating word embeddings. Each model was tested on both embeddings to understand their overall performance on the task. I used L2 based penalty for the logistical regression model. It offered moderate performance on the task with an accuracy around 59. SVM model with $'rbf'$ kernel offered relatively similar performance to logistic regression. Random Forest offered the best performance among the statistical models with the average accuracy of 60% for Distributed Memory models and 51% for Distributed Bag-of-Words model. Overall, The Distributed Memory models performed relatively well on Amazon dataset and Distributed Bag-of-Words performed well on Apple dataset. Also, the overall performance with 4 hours window was better than 24 hours window. It matches with the intuition that the sentiment in the market might fluctuate between positive and negative over a period of 24 hours. Hence, It is difficult to classify the news articles based on the difference in stock prices in a 24 hours window.

When we compare the performance of the classifier using representation from BERT against the statistical models, the classifier out-performs the statistical models in all categories. It can be attributed to multiple factors. One primary reason is that the features obtained from BERT representation is highly predictive of the text and the sentiment associated with it. It is particularly due to the use of attention in the transformer models such as BERT. The attention mechanism helps in providing powerful sequence representation particularly for language. It has resulted in BERT

performing well in wide range of language tasks mentioned in Devlin *et al.* (2018) including tasks related to Sentiment analysis.

| Amazon | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | DM | | | | DBOW | | | |
| Window | 4 Hours | | 24 Hours | | 4 Hours | | 24 Hours | |
| Model | Accuracy | F1 | Accuracy | F1 | Accuracy | F1 | Accuracy | F1 |
| LR | 59.06% | .58 | 58.29% | .58 | 52.33% | .52 | 53.30% | .53 |
| SVM | 58.83 % | .58 | 56.18 % | .56 | 51.85 % | .51 | 52.89 % | .50 |
| Random Forest | 61.59 % | .55 | 59.29 % | .59 | 51.83 % | .52 | 51.70 % | .52 |
| Apple | | | | | | | | |
| | DM | | | | DBOW | | | |
| Window | 4 Hours | | 24 Hours | | 4 Hours | | 24 Hours | |
| Model | Accuracy | F1 | Accuracy | F1 | Accuracy | F1 | Accuracy | F1 |
| LR | 52.33% | .52 | 53.30% | .53 | 62.91% | .63 | 62.84% | .63 |
| SVM | 51.85 % | .51 | 52.89 % | .50 | 62.89 % | .63 | 62.64 % | .63 |
| Random Forest | 51.83 % | .52 | 51.70 % | .52 | 64.03 % | .64 | 64.98 % | .65 |

Table 3.1: Results using Statistical Models

| Metric | Amazon | Apple |
|---|---|---|
| AUC | .68 | .67 |
| Accuracy | 68% | 67% |
| F1 | .70 | .68 |
| Precision | .65 | .65 |
| Recall | .76 | .72 |

Table 3.2: Results using BERT Embedding on 24 hours window

Chapter 4

VISUAL LANGUAGE REPRESENTATION

## 4.1 Formulation

Image captioning is the problem of identifying appropriate captions that capture the context presented in the given image. Given an input image $I \in \mathbb{R}^{H \times W \times C}$, the captioning model tries to find an ordered sequence vector $S : (t_1, t_2, t_3 \ldots . t_m)$ where $t_i$ represents a token in the language space. If the context of the image $I$ can be represented in a shared embedding space $H$ using a $\mathbb{R}^d$ dimensional feature vector, then the ordered sequence $S$ can be considered as the projection of $\mathbb{R}^d$ dimensional feature vector in the language space. The model uses encoder-decoder based architecture similar to many language and multi-modal models Devlin *et al.* (2018); Vaswani *et al.* (2017); Lu *et al.* (2019); Sun *et al.* (2019). The encoder $E(x)$ is the transformation function mapping the input image $I$ on to the shared embedding space $H$ and the decoder $D(x)$ transforms the encoded feature vector $E(I)$ to the language space. The output of decoder is represented as sequence vector $S$. Both encoder and decoder use attention mechanism to capture the relationship between the different input elements in the source space and also identifies the vector in the target feature space (language) that is closely aligned with the elements in the source feature space (vision). In the following sections, I will describe different components of the model in detail.

## 4.2 Attention in Language Model

Vaswani *et al.* (2017) proposed the transformer architecture which is based entirely on the attention mechanism. Attention transforms the input into 3 matrices viz. query and key-value pairs. The output is then computed as the weighted sum of the value. The weights assigned to each value is decided based on the compatibility of the query with the key. In the case of language modeling, both the input and output of the model are sequences of varying length. First, an input embedding of dimension $d_{model}$ is created using a pre-trained vector representation of words. Query $Q$, Key $K$- Value $V$ pairs are generated by a linear transformation of the input embedding vector. The transformer model uses 'Scaled Dot-Product Attention'. The queries and keys of dimension $d_k$ and the values of dimension $d_v$ are used to compute the output. In scaled dot-product attention, we compute the dot product between the query $Q$ and the key $K$. The dot product is scaled by a factor of $1/\sqrt{d_k}$ to offset the effect due to extremely small gradients that are generated by the softmax function.

$$\textbf{Attention(Q,K,V)} = \textbf{softmax}(\frac{QK^T}{\sqrt{d_k}})V \qquad (4.1)$$

## 4.3 Attention in Vision Model

Inspired by the success of attention in language models, variants of attention are developed for vision models. Bello *et al.* (2019) proposed a novel two-dimensional relative self-attention mechanism for vision models.The self-attention is combined with the traditional convolution to offer the best performance. An input image $I \in \mathbb{R}^{H \times W \times F_{in}}$ is flattened to a matrix $X \in \mathbb{R}^{HW \times F_{in}}$. The flattened matrix is treated similar to the input sequence embedding in the language model. I compute the query

$Q$, key $K$ and value $V$ by applying a linear transformation on input $X$. We later use the matrices $Q$, Key $K$ and Value $V$ to compute the output.

$$O_h = \textbf{Softmax} \left( \frac{(XW_q)(XW_k)^T}{\sqrt{d_k^h}} \right) (XW_v) \tag{4.2}$$

where $W_q, W_k \in \mathbb{R}^{F_{in} \times d_k}$ and $W_v \in \mathbb{R}^{F_{in} \times d_v}$. We use learned linear transformations $W_q, W_k$ and $W_v$ for computing the Queries $Q$ and Key $K$ - Value $V$ pairs. A linear transformation $W^O$ is applied to the output $O_h$ and then reshaped to the dimensions of the input image $\mathbb{R}^{H \times W \times d_v}$.

## 4.4   Cross-Modal Attention

The attention techniques used in the previous works [Vaswani *et al.* (2017); Bello *et al.* (2019)] have focused on applying attention to either vision or language modality. The image captioning task requires attending to both vision and language simultaneously. I have developed a novel attention mechanism to accommodate this need for attending to different modalities. In the case of the multi-modal tasks, previous works have focused on sharing the queries and key-value pairs across the layers. In case of ViLBERT [Lu *et al.* (2019)], the co-attention blocks use key-value pairs from each modality as input to the other modality's multi-headed attention block. The technique mimics the common attention mechanism in vision and language models. In my work, I focus on generating output sequence (predicted caption) from the input embedding of the image.

First, I flatten the input tensor $I \in \mathbb{R}^{H \times W \times F_{in}}$ into matrix $X \in \mathbb{R}^{HW \times F_{in}}$. I extract the input key $K$ and value $V$ matrices applying the linear transformation on the flattened input tensor. I generate relative position encoding. I sum the flattened input attention map with the relative position encoding. The computed output is fed to the cross-modal attention layer. I use the query $Q$ from the sequence to learn

22

the relevant pixels that the query needs to attend. The relative importance scores of the different image regions are obtained through a dot product between the key-value pairs from the input image and the query. The relative weights are then used to compute the weighted sum of the value vector $V$.

During training, I compute the query $Q$ from the sequence (caption) $S$ and learn the linear transformation weights $W_q^t, W_k^I$ and $W_v^I$ for the image embedding. During testing, the input query $Q$ is a vector of $< start\_token >$ of dimension $\mathbb{R}^{d_{model}}$ and the model generates the output sequence using the input image and the learned linear transformation weights.

$$O_h = \textbf{Softmax}\left(\frac{(SW_q^t)(XW_k^I)^T}{\sqrt{d_{cm}^h}}\right)(XW_v^I) \tag{4.3}$$

where $W_q^t \in \mathbb{R}^{d_{model} \times d_{cm}}$ and $W_k^I, W_v^I \in \mathbb{R}^{F_{in} \times d_{cm}}$. A linear transformation $W_O \in \mathbb{R}^{d_{cm} \times d_{model}}$ is applied on the output $O_h$ to generate the output caption $O_c \in \mathbb{R}^{n \times d_{model}}$, where $n$ is the length of the generated caption and $d_{cm}$ and $d_{model}$ are the dimensions of the cross-modal hidden embedding and the dimension of the sentence embedding. $F_{in}$ is the number of filters in the image embedding.

### 4.5  Multi-Head Attention

Similar to the Transformer [Vaswani *et al.* (2017)] architecture, I employ multi-head attention in attention blocks. The input image and caption are linearly projected $N_h$ times using different linear projections and attention applied in parallel to $N_h$ operations. The results of different attention operations are then concatenated and projected again using a linear transformation $W^O$ to obtain the final values. Multi-head attention is applied on both image and text inputs. Vaswani *et al.* (2017) claims

that multi-head attention jointly attends to information from different representation sub spaces at different positions.

$$\text{MultiHead(Q,K,V)} = \text{concat}(\text{head}_1, \ldots, \text{head}_h)W^O \tag{4.4}$$

where $\text{head}_i$ is the Attention$(QW_i^Q, KW_i^K, VW_i^V)$.

The dimensions of the linear parameters while applied to text sequences are $W_i^Q \in \mathbb{R}^{d_{model} \times d_k}$, $W_i^K \in \mathbb{R}^{d_{model} \times d_k}$, $W_i^V \in \mathbb{R}^{d_{model} \times d_v}$ and $W^O \in \mathbb{R}^{hd_v \times d_{model}}$ where $d_k = d_v = \frac{d_{model}}{h}$

The dimensions of the linear parameters while applied to images are $W_i^Q \in \mathbb{R}^{F_{in} \times d_k^h}$, $W_i^K \in \mathbb{R}^{F_{in} \times d_k^h}$, $W_i^V \in \mathbb{R}^{F_{in} \times d_v^h}$ and $W^O \in \mathbb{R}^{d_v \times d_v}$. In my model , I have used $d_k = d_v = d_{model}$ and $d_q = d_{model}/h$ while processing images.

## 4.6   Architecture

The model uses an encoder-decoder based architecture similar to other language models [Vaswani *et al.* (2017); Devlin *et al.* (2018); Sun *et al.* (2019)]. The position encoding representing the position of different pixels in the input image is computed. The position encoding is aggregated with the input image. The attention operation is position invariant and doesn't encode the position information naturally. Hence, we require position encoding to propagate position information across the layers. The encoder is a stack of $N$ identical layers. Each layer consists of $Nc$ sub layers. The decoder is also a stack of $N$ identical layers with three sub layers. I will describe each layer in detail in below sections. The figure 4.1 describes the model architecture.

## 4.7   Encoder

Following the architecture of Transformer [Vaswani *et al.* (2017)], each sub layer in the encoder consists of Attention Augmented Convolution [Bello *et al.* (2019)] which
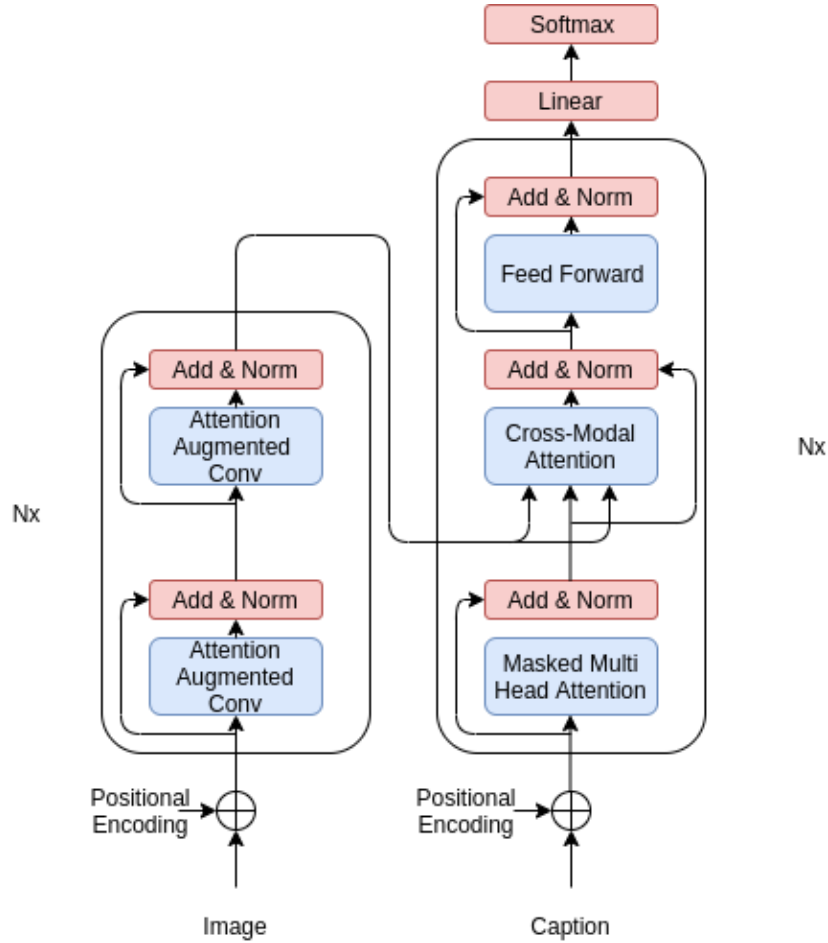
Figure 4.1: The cross modal transformer has $Nx = 6$ encoder and decoder layers. The encoder layer consists of $Nc = 2$ sub layers of Attention Augmented Convolution. The sub layers in both encoder and decoder use residual connections [He et al. (2015)] followed by Group Normalization [Wu and He (2018)]. The decoder employs masked multi-head attention which is followed by cross-modal attention and feed forward layers. During training, I feed image and relevant caption to encoder and decoder respectively. After encoding the position information, The encoder uses AA-Convolution layers to generate attention maps. After encoding the position information, the decoder uses masked multi-head attention to mask random positions in the sequence embedding. The masked sequence embedding along with the attention maps are fed to the cross-modal attention block.
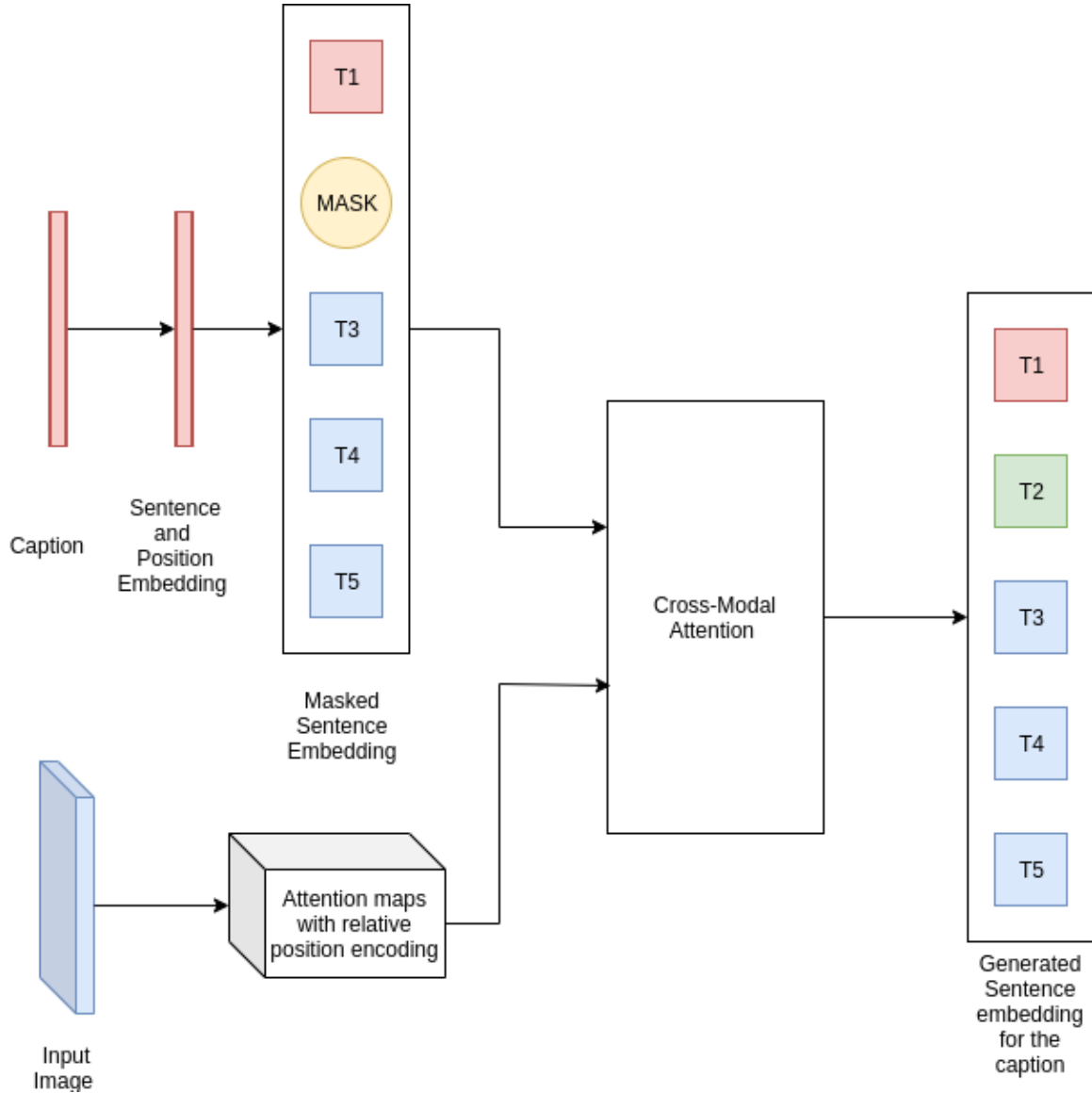
Figure 4.2: I use *Masked Language Modelling (MLM)* proxy task for training the network. First, I encode the image of dimensions $(h, w, c)$ into attention maps of same dimensions $(h, w, c)$ using Attention-Augmented Convolution layers. I generate sentence embedding of shape $n \times d_{model}$ for the caption of length $n$. The tokens in sentence embedding are randomly masked with $< mask >$ token. I use Cross-modal attention to predict the masked token. The model learns the linear transformation parameters $W_q^t$, $W_k^i$ and $W_v^i$ .

implements a multi-head self-attention followed by a position-wise fully connected feed-forward network.

I also implement residual connection [He *et al.* (2016)] around each of the sub-layers followed by Group Normalization [Wu and He (2018)]. The output of each sub-layer is GroupNorm($x$ + Sublayer($x$)), where $Sublayer(x)$ is the function implemented by the sub-layer.

### 4.7.1  Attention Augmented Convolution

Each sub layer in the encoder implements Attention Augmented Convolution [Bello *et al.* (2019)]. The architecture of the Attention Augmented Convolution is presented in Figure 4.3. Since the AA-Conv layers flattens the input image to apply attention, it results in a memory cost of $O((N_h(HW)^2)$. This prohibits us from using larger spatial dimensions. Hence, I down-sampled the input images to 124 pixels and used smaller batch sizes to meet the memory constraints.

### 4.8  Decoder

Similar to encoder, decoder is stack of $Nx$ identical layers. Each decoder layer consists of three sub layers. I use residual connections around the sub-layers ,followed by Group Normalization [Wu and He (2018)]. The attention in decoder is modified to not attend to subsequent positions ensuring that the output at any time step depends only on the known outputs and not on future positions. This is done by masking the inputs present in the future positions using $< pad >$ token. The cross-modal attention layer follows the self-attention block and accepts the inputs from the encoder and masked multi-head attention in decoder. Feed forward and normalization layers follows the cross-modal attention.

## 4.9 Position Encoding

### 4.9.1 Position Encoding of Image in Encoder

Following the technique used in AA-Convolution [Bello *et al.* (2019)] , I use relative position encoding in the encoder. Since the attention mechanism doesn't naturally encode the position information, it is essential to explicitly propagate position information through the layers. But the encoding scheme should satisfy translation equivariance to extract reliable features. The relative positional embedding independently adds relative height and width information. Consider pixels $i = (i_x, i_y)$ and $j = (j_x, j_y)$. If $q_i$ is the query vector for pixel $i$ and $k_j$ is the key vector for the pixel $j$ and $r^W_{j_x - i_x}$ and $r^H_{j_y - i_y}$ are learned embedding for relative width $j_x - i_x$ and relative height $j_y - i_y$, respectively, then the relative logit is computed as,

$$l_{i,j} = \frac{q_i^T}{\sqrt{d_k^h}}(k_j + r^W_{j_x - i_x} + r^H_{j_y - i_y}) \tag{4.5}$$

The output of the attention for each head can be computed using equation (4.6),

$$O_h = \text{Softmax}\left(\frac{QK^T + S_H^{rel} + S_W^{rel}}{\sqrt{d_k^h}}\right)V \tag{4.6}$$

where relative position matrices $S_H^{rel}, S_H^{rel} \in \mathbb{R}^{HW \times HW}$ is computed as,

$$S_H^{rel}[i,j] = q_i^T r^H_{j_y - i_y} \tag{4.7}$$

$$S_W^{rel}[i,j] = q_i^T r^W_{j_y - i_y} \tag{4.8}$$

### 4.9.2 Position Encoding of Sequences in Decoder

In case of decoder, it is essential to model the order of sequence. For this, I inject the information about the relative or absolute position of the tokens. The dimension of the position encoding is $d_{model}$ allowing them to be aggregated with the sequence

embedding. Following the encoding technique in Transformer [Vaswani *et al.* (2017)], I use sine and cosine functions of different frequencies to create position embedding. For position *pos* and dimension *i*, the encoding is given by,

$$\text{PE}_{pos,2i} = \sin\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right) \tag{4.9}$$

$$\text{PE}_{pos,2i+1} = \cos\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right) \tag{4.10}$$

## 4.10   Training

This section details the training strategy for my model. The figure 4.2 describes the training strategy.

### 4.10.1   *Microsoft COCO Dataset*

MS COCO dataset is a large dataset used for multiple vision-based tasks such as image recognition, segmentation and captioning. The dataset contains various features for images. It contains 300,000 images belonging to 80 different categories. The dataset contains 5 different human-annotated captions for each image. Several previous works [Anderson *et al.* (2018); Bello *et al.* (2019); Xu *et al.* (2015)] in image-captioning have used MS COCO dataset for training and reporting.

I used MSCOCO 2017 captions dataset [Lin *et al.* (2014)] in this work. I applied karpathy splits [Karpathy and Fei-Fei (2017)] similar to previous implementations on image captioning [Anderson *et al.* (2018)]. The split contains 113,287 training images with five captions each, and 5K images for validation and testing each. I perform few text pre-processing steps on the caption sentences. The captions are converted to lower case. The model vocabulary is of size 10,000 words.
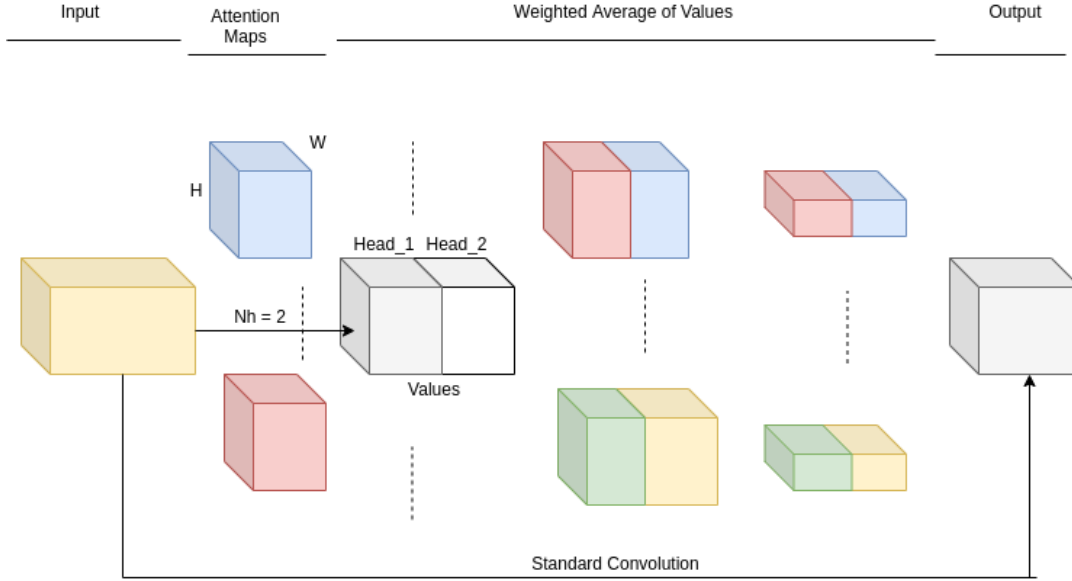
Figure 4.3: Following the technique proposed in Attention-Augmented Convolution [Bello *et al.* (2019)], I perform multi-head attention on the input image of dimensions $(h, w, c)$ and generate $N_h$ attention maps using the queries and keys of the input image. I use the attention maps to compute the $N_h$ weighted averages of values $V$. I concatenate the results of the multi-head attention and reshape them to match the dimensions of the original image. I compute standard convolution on the input image in parallel and concatenate the results of the attention and convolution.

### 4.10.2   Conceptual Captions

Sharma *et al.* (2018) used Automatic Image Captioning technique to compile Conceptual Captions dataset. The dataset contains an order of magnitude more images than the MS COCO dataset. It contains a wide variety of images and caption styles collected from a large number of web pages. It contains only JPEG formatted images with dimensions greater than 400 pixels. Each image contains an alt-text description extracted from the online page. Also, each image contains a conceptual caption. It contains 3.3M images and description pairs. The validation and test

```
A brown and white dog on a skateboard next to group of people.
People watch from the sidelines as a bulldog rides on a skateboard
A dog runs alongside a skateboard with one paw on.
a close up of a dog riding a skate board
a dog is pushing a skateboard while kids walk.
```



Figure 4.4: Sample 1 for category 'dog' in MS COCO. The dataset contains five different captions for each image. The dataset has a wide variety of samples. It contains 300,000 images belonging to 80 different categories

set contain 28 K and 22.5 K images respectively. Unlike the COCO dataset, the Conceptual Caption dataset contains images harvested from the internet and hence represent a wider style of images and captions. The conceptual captions have unique word ratios covering various POS tags. The dataset contains a wide variety including natural images, professional images, product images, and drawings.

```
A man on a skateboard and his dog running down the street.
A boy speeds along a road with his greyhound.
A young man riding a skateboard down a road next to a sheep.
A man is crouched down next to a running dog while skateboarding.
Man on skateboard with dog running next to him on roadway.
```



Figure 4.5: Sample 2 for category 'dog' in MS COCO

### 4.10.3  Hardware and Training Schedule

I trained the model on one machine using 2 NVIDIA V100 32 GPUs. I used mini batches of size 4 and trained it for 10 epochs. Each epoch comprised of 30 K steps. Each epoch took around 10 hours to complete.

### 4.10.4  Optimizer

Following Transformer [Vaswani *et al.* (2017)], I used an Adam Optimizer Kingma and Ba (2015) with $\beta_1$ =0.9 and $\beta_2$ =0.98 and $\epsilon = 10^{-9}$.

**Alt-text**: A Pakistani worker helps to clear the debris from the Taj Mahal Hotel November 7, 2005 in Balakot, Pakistan.

**Conceptual Captions**: a worker helps to clear the debris.

**Alt-text**: Musician Justin Timberlake performs at the 2017 Pilgrimage Music & Cultural Festival on September 23, 2017 in Franklin, Tennessee.

**Conceptual Captions**: pop artist performs at the festival in a city.

Figure 4.6: Sharma *et al.* (2018): The dataset contains 3.3 M images and alt-text extracted from the web pages. It also contains pre-processed Conceptual Caption for each image.

I used the following strategy to increase the learning rate linearly for first *warm up_steps*.

$$lrate = d_{model}^{-0.5} \cdot min(step^{0.5}, step \cdot warmup^{-1.5}) \tag{4.11}$$

### 4.10.5  Masked Language Modelling

Similar to training strategy in BERT [Devlin *et al.* (2018)], I use Masked Language Modelling (MLM) task for training the model. In this task, I replace random positions in the embedded vector of the sequence with the $< MASK >$ token. I train the model to predict the token in the masked position using the context offered by the other tokens in the sequence. This training strategy optimizes the below joint probability,

$$\log P(x|\theta) = \frac{1}{Z(\theta)} \sum_{I=1}^{N} \log \phi_i(x|\theta), \tag{4.12}$$

where $\phi_i$ is the potential function for the $i^{th}$ input element, $\theta$ is the parameter and $z(\theta)$ is the partition function. The potential function is given by,

$$\log \phi_i(x|\theta) = x_i^T f_i(x_{\setminus i}|\theta)_i, \tag{4.13}$$

where $f_i(x_{\setminus i}|\theta)$ is the output of the model for the $i^{th}$ position, where $x_{\setminus i} = \{x_1, \ldots x_{i-1}, [\text{MASK}], x_{i+1}, \ldots, x_N\}$. For a random sentence $x$ and random position $i$ chosen for masking, the loss is computed as ,

$$\text{L}_{MLM}(\theta) = -\text{E}_{x \sim D, I \sim \{1, \ldots, N\}} \log \phi_i(x) \tag{4.14}$$

## 4.11    Results

I trained the model using the configuration presented in Section[4]. I used $Nx{=}6$ layers of encoders and decoders. Encoder contained $Nc{=}2$ layers of the Attention Augmented convolution network. I trained the model on the MS COCO dataset for 10 epochs. I present the results of the image captioning task in figures 4.7, 4.8 and 4.9. In many cases, the model is able to capture the context and understand the relationship between the different objects present in the images. The generated captions are coherent and are grammatically correct with proper punctuation. The model is, therefore, able to learn semantics and the context of the captions.

### 4.11.1    Inductive Bias

The generated captions are sometimes affected by the inductive bias due to human-annotated captions. Since many images contain humans, the generated captions usually begin with phrases such as 'A man/ A woman / A person'. The issue is visible in the third image in figure 4.8. The image contains a dog sitting in front of

a laptop. But the model incorrectly identifies the dog as woman. I hypothesize that training the model with different categories would help in resolving the issue.

### 4.11.2   Attention to Context

The inability to understand the central context of the image is another issue reflected in some of the images. Sometimes the model doesn't capture the central theme of the image, instead describes the events in the background. The second image in figure 4.5 is an example of this issue. The actual caption describes the person in a blue shirt as the person occupies major portion of the frame. The model instead describes the group of people in the background. I suspect the issue is due to the limited amount of training data which is insufficient for the model to understand the central aspect in the given image. Additionally, the relevance of caption generated can be subjective. It is a challenge to estimate the relevance of a caption generated to the associated image.

### 4.12   Conclusion

I proposed a model that can generate cross-modal representation of a visual and linguistic representations in the shared visual-linguistic space. I use a transformer-based architecture where I replace the sequence encoder with an image encoder. The image encoder uses several layers of Attention Augmented convolution, which employs self-attention to generate attention maps. I introduce a novel cross-modal attention layer to perform attention across vision and language modalities. I use a masked language modeling technique to train the model. I present the samples from the image captioning task to demonstrate its effectiveness. In this work, I have focused on applying cross-modal attention to the vision-to-language task.We can further extend the architecture and the cross-modal attention technique to accommodate the lan-

| Image | Actual Caption | Generated Caption |
|---|---|---|
|  | A guy skateboarding on the bars of stairs | A person sitting on a skateboard at a skate park |
|  | A young man standing in front of a fence holding a skateboard | A man is standing on the side of a bench |
|  | A dog sitting in front of a laptop on top of a bed | A woman is sitting on a couch with a laptop |
|  | Two women in glasses using Nintendo wiI remote controllers | A person playing a video game in a room |

Figure 4.7: Visualizations of image captions from MS COCO. Qualitative samples from the results. For each example, I show the input image (left), human-annotated caption (center), model-generated caption (right).

| Image | Actual Caption | Generated Caption |
|---|---|---|
|  | fish and broccolI plated neatly on a dish sitting on a table | a close up of a plate of food |
|  | a man in a blue shirt playing with a white ball | a couple of people on a bench near a tree . |
|  | A referee and tennis player speaking while fans watch in the stands . | A man in a a blue short and white shirt and a woman . |
|  | A woman sits on a motorcycle with a sidecar . | A man in a a kitchen with a stove and a |

Figure 4.8: Visualizations of image captions from MS COCO. Qualitative samples from the results. For each example, I show the input image (left), human-annotated caption (center), model-generated caption (right).
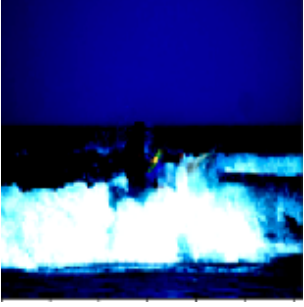
| Image | Actual Caption | Generated Caption |
|---|---|---|
|  | A man on a surfboard riding a wave . | A man with a surfboard on the beach . |
|  | A skier stands on a small ledge in the snow | A man riding a skateboard ramp . |
|  | A man and women looking at a small child . | A man A woman with a |
|  | A white plate topped with a sandwich and sliced veggies | A plate of pizza on a white plate |

Figure 4.9: Visualizations of image captions from MS COCO. Qualitative samples from the results. For each example, I show the input image (left), human-annotated caption (center), model-generated caption (right).

guage to vision tasks. I hypothesize that such an architecture would be capable of generalizing to many vision and language tasks . Such a model would further improve state-of-art in several downstream tasks in vision and language.

Chapter 5

LANGUAGE IMAGE TRANSFORMER

## 5.1 Introduction

One of the fundamental characteristics of human intelligence is the ability to associate information from different modalities, such as vision and language. Understanding the context in various modalities can be crucial for learning in the real-world. The complexity of the real-world environment is a challenge for representing in a single modality. While an image is particularly useful for understanding finer details about a scene, a sentence description can capture high-level qualitative aspects about the context. Each modality is useful for optimally learning a different aspect of the scene. Learning to associate different modalities would help in generating a better model of the real world.

With the recent progress in self-supervised learning, there is an increased focus on developing techniques that can generalize across domains and tasks. I have achieved great success in domain-specific tasks associated with vision and language. Many popular techniques in both domains use some kind of pre-training proxy task to learn the association between the modalities. These techniques have also proved to be successful across a wide range of downstream tasks in both domains. The similarities between the techniques in vision and language motivate us to hypothesize that there could be a common learning approach that can be developed for learning simultaneously from vision and language. A natural starting point towards this goal would be to explore the techniques that have achieved success in both vision and language. Instead of developing separate proxy tasks for learning from vision and language, I

explore the possibility of using a common proxy task for learning simultaneously from both modalities.

Contrastive Predictive Coding is a unsupervised learning technique that has shown significant results in image recognition. It learns to model the global structure in the data by predicting spatio-temporal variation in the data. It has recently also been applied to vision-language tasks [Sun *et al.* (2019)]. With the ability to model the high-level distribution and adaptability across domains makes it a great candidate for the use in cross-domain applications. Attention techniques has also received wide adaption in both vision and language applications.

In this chapter, I propose a transformer-based architecture for learning multi-modal representation by maximizing the Mutual Information (MI) between the vision and language representations. I also propose a novel attention layer that can accept representation from two modalities and generate a shared representation.

## 5.2 Approach

### 5.2.1 Mutual Information Maximization

There are several recent successes in applying the InfoMax Principle to maximizing Mutual Information. The tasks involve learning a representation that maximizes the Mutual Information between the input signal and the encoded output. The MI between two random variables $X$ and $Y$ can be defined as the amount of information that can be learnt about $Y$ by observing $X$. MI is formally defined as Kullback-Leibler (KL) divergence measure between the joint probability distribution $p(x, y)$ and the product of its marginals $p(x)p(y)$. $p(x, y)$ is the joint probability distribution between the random variables $X$ and $Y$, while $p(x)$ and $p(y)$ are the marginal distributions.

$$I(X;Y) = D_{KL}(p(x,y)||p(x)p(y)) = \mathbb{E}_{p(x,y)} \left[ \log \frac{p(x,y)}{p(x)p(y)} \right]$$

It is a challenging task to estimate MI in higher-dimensional space. Hence, recent efforts have focused on maximizing a tractable lower bound on MI [Poole *et al.* (2019)]. Oord *et al.* (2018) introduced Contrastive Predictive Coding to maximize the lower bound in MI. Sun *et al.* (2019) adapted BERT to learn multi-modal representation that maximizes the lower bound.

One of the common approaches is to maximize MI between two related encoded inputs in relatively lower dimensional space than the original input. For example, in unsupervised image representation learning, two overlapping views of image $X$ are constructed and encoded using different encoders. Finally, to learn a representation of the image, the mutual information between the two encoded views are maximized.

While learning a shared representation space between modalities such as image and text, the related inputs from individual modalities are handled similar to views in the previous example. In this case , the views describe a shared context in different modalities. In the task of image captioning, the image and its related caption are two different views describing the same underlying setting. Although present in different modalities, the representations share many characteristics describing the shared context. In order to learn a common representation, I maximize the MI between the related inputs in the individual modalities. I learn encoders in the two modalities and then maximize the MI between the output representations of the two encoders. Consider an input image $X_m^{(i)}$ and a related caption $X_t^{(i)}$ where $i \in \{1 \ldots n\}$ from the $n$ paired samples. The caption $X_t^{(i)}$ is a sentence describing the context presented in the image $X_m^{(i)}$. Our objective is to learn shared representation for the image-caption pair $(X_m^{(i)}; X_t^{(i)})$. In order to achieve this objective, the model initially learn image representation $g_m(X_m^{(i)})$ using a generic image encoder $g_m(.)$. Similarly, the
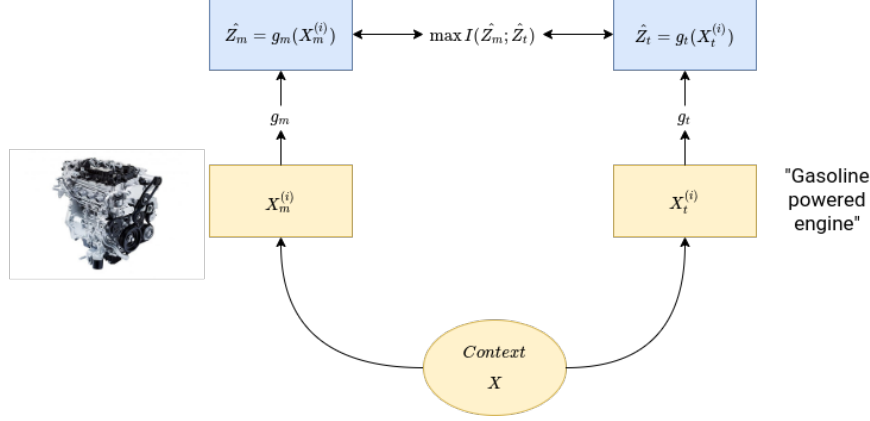
Figure 5.1: We can describe hidden context 'Engine' $(X)$ with an image $(X_m^{(i)})$ displaying its different components or a caption $(X_t^{(i)})$ describing its characteristic. Both modalities provide different perspectives of the context. We try to learn from both modalities using a shared representation. Initially, image encoder $g_m(.)$ generates a representation of the image $(\hat{Z}_m)$, while text encoder $g_t(.)$ generates a sentence representation of the caption $(\hat{Z}_t)$. Finally, we learn a shared representation that maximizes the mutual information $I(\hat{Z}_m; \hat{Z}_t)$ between the modalities.

model also learns a text representation $g_t(X_t^{(i)})$ using a sentence encoder $g_t(.)$. Once the model learns individual representations in the two modalities, it learns a shared representation that maximizes the MI between the two representations $g_m(X_m^{(i)})$ and $g_t(X_t^{(i)})$.

$$\max_{g_m \in \mathbb{G}_m, g_t \in \mathbb{G}_t} I_{\text{EST}}\left(g_m(X_m^{(i)}); g_t(X_t^{(i)})\right) \text{ with } i \in \{1 \ldots n\} \tag{5.1}$$

where $I_{\text{EST}}$ is a estimator of true MI $I(X; Y)$ between the encoded representations $g_m(X_m^{(i)})$ and $g_t(X_t^{(i)})$.

As described in Micheal et. al. 2019, there are multiple advantages in choosing the objective in 5.1 . First, it provides the flexibility to choose the dimension of representations $g_m(X_m^{(i)})$ and $g_t(X_t^{(i)})$. We can choose a representation space that is computationally feasible and is best suitable for the downstream task in hand.

Second, it provides the flexibility to choose the modelling architecture for the image encoder $g_m(.)$ and language encoder $g_t(.)$ that can complement each other is capturing the overall context.

### 5.2.2   Contrastive Predictive Coding (CPC)

Oord *et al.* (2018) proposed CPC, a self-supervised learning technique to extract useful high level representations in speech, images and text. The CPC model uses powerful autoregressive model to predict the representation of future observations using those of the past observation. By predicting the future observations, the model learns to discard noisy low-level features and instead captures shared high-level global features essential for predicting future observations. In other words, the objective maximizes the mutual information between the encoded representations of the input signals. The technique uses a contrastive loss, where the model classifies future observations amongst the set of unrelated negative samples.

This work uses Contrastive Predictive Coding (CPC) in the image stream to capture high level image features in the shared representation. Unlike the original implementation, the model is not trained on a uni-directional prediction task. Instead, it learns the features by predicting randomly masked input representations amongst the set of negative samples. I hypothesize that the modified objective helps the model to capture the overall context using the neighbouring patches in all direction.

The model uses an encoder $g_m$ to map the sequence of input observation $x_m$ to a sequence of low-dimensional latent representations $z_m = g_m(x_m)$. A set of $n$ representations are randomly chosen and masked from the sequence. An auto-regressive model $g_{ar}(.)$ then summarizes the partially-masked input representation as a context vector $c_m = g_{ar}(z_{m \setminus n})$. A linear transformation $\hat{z}_n = W_k c_m$ is later applied to the

context vector $c_m$. The weights $W_k$ are learnt by predicting the masked representation $\hat{z}_n$ among the set of randomly chosen negative samples $\{z_l\}$.

When applied to images, The input image is initially segmented into a set of overlapping patches $x_{i,j}$ which are encoded using a network $f_\theta$ into a embedding vector $z_{i,j} = f_\theta(x_{i,j})$. A random set of input representation $\{z_n\}_{n<(i \times j)}$ are masked in the embedding vector of size $i \times j$. The partially-masked vector $z_{(i,j)\backslash n}$ is summarized using a transformer model $g_{ar}$ into a context vector $c_{i,j}$. The transformer model discussed in 5.3.3 is specially adapted to handle cross-modal inputs from different modalities. A linear transformation is then applied to the output representation $c_{i,j}$ from the transformer layer to generate predicted feature vector $\hat{z}_n = W_k c_{i,j}$.

Similar to the original implementation, the model uses contrastive loss to evaluate the prediction task. The goal is to recognize the targets $z_n$ amongst the randomly sampled feature vectors $\{z_l\}$ from the dataset. I apply softmax to obtain a probability distribution over the sampled feature vectors.The model evaluates the output prediction using cross-entropy loss.

$$L_{CPC} = -\sum_{i,j} \log P(z_{i,j}|\hat{z}_{i,j}, \{z_l\}) = -\sum_{i,j} \log \frac{\exp(\hat{z}_{i,j}^T z_{i,j})}{\hat{z}_{i,j}^T z_{i,j} + \sum_l \exp \hat{z}_{i,j}^T z_l}$$

### 5.2.3 Attention

Vaswani et. al. initially proposed the Transformer architecture for the Neural Machine Translation (NMT). The architecture was based entirely on the attention mechanism. Attention transforms the input signal into 3 matrices viz. query $Q$ and key-value $(K-V)$ pairs. The matrices are obtained from the linear transformations of the input signal. The query $Q$ and the key $K$ vectors define the similarity between the different components in the input sequences considered. The similarity/compatiblity score between the components is obtained using the dot product between relevant

Query and Key vectors. The scores are used for assigning the relative weights to the value vector. The output of the attention block is computed as the weighted sum of the value vector. The weight assigned to each value vector is decided based on the compatibility of the query with the key. In the case of language modeling, both the input and output of the model are sequences of varying length. First, an input embedding of dimension $d_{model}$ is created using a pre-trained vector representation of words. Query $Q$, Key $K$- Value $V$ pairs are generated by a linear transformation of the input embedding vector. The transformer model uses scaled dot product attention. The queries and keys of dimension $dk$ and the values of dimension $dv$ are used to compute the output. In scaled dot-product attention, We compute the dot product between the query Q and the key K. The dot product is scaled by a factor of $\frac{1}{\sqrt{d_k}}$ to offset the effect due to extremely small gradients that are generated by the softmax function.

$$\text{Attention(Q,K,V)} = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V \tag{5.2}$$

The model uses a variant of the attention block, as a autoregressive function to summarize the context vector. The original implementation of the transformer uses attention for learning sequence representation of language sentences. Several recent efforts have focused on adapting attention-based models such as BERT for sequence generation tasks in various domains. Sun *et al.* (2019) focused on adapting BERT to learn shared representation for videos. Similarly, models such as ViLBERT [Lu *et al.* (2019)] , UNITER [Chen *et al.* (2019)] adapted BERT for performing a wide variety of Image-Language tasks such as Visual Question Answering [Antol *et al.* (2015)] and Image Retrieval. Our model falls under a similar category. The attention block is used to align the image and language embedding together. In other words,

it generates a shared representation that maximizes the MI between the modality-specific representations.

One of the general approaches developed for vision-language representation involves defining separate streams for Image and Language. The model uses a shared cross modal transformer to share the parameters between the image and language streams.

## 5.3 Architecture

### 5.3.1 Image Encoder

Following the architecture presented in Contrastive Predictive Coding (CPC) [Oord *et al.* (2018)], the model uses a ResNet-18 architecture [He *et al.* (2016)] for generating the image embedding. The model reshapes images to 256×256 and flips some images horizontally. It segments each image into $7 \times 7$ grid. Each segment in the grid is of shape $64 \times 64$ with 32 pixels overlap with its neighbors. The model extracts the embeddings for each of the segmented patches from the penultimate layer of the ResNet. Finally, the model flattens each image embedding to a sequence of length 49 with the embedding length of 512. The embedding is used as input for the image stream to the cross-modal transformer.

### 5.3.2 Text Encoder

The model uses learned word embedding of 512 dimensions from Glove [Pennington *et al.* (2014)] to represent the input caption associate with the input image. It applies masked self-attention to the embedding vector, similar to the one used in the decoder layer of the transformer. Each position in the layer output can attend to all the previous and current positions in the input embedding but not the future
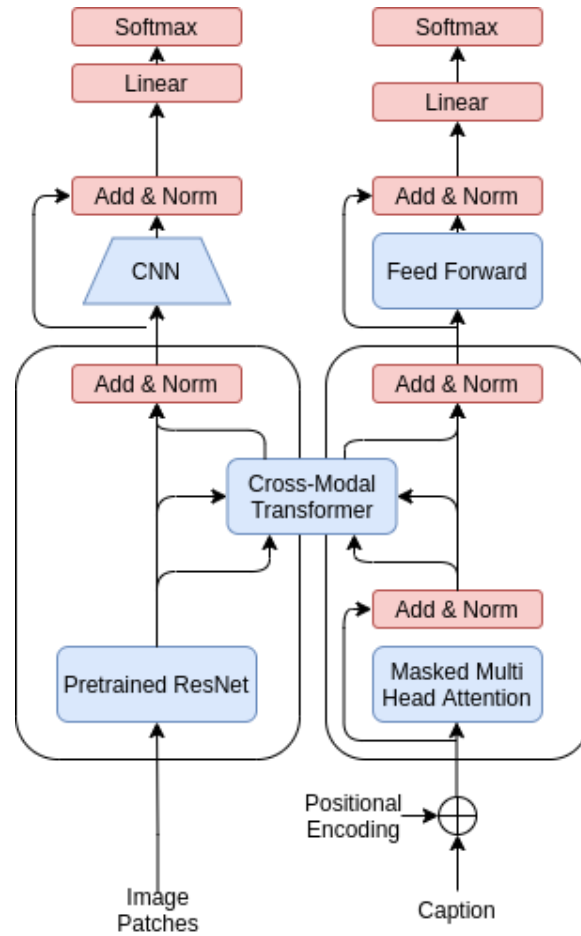
Figure 5.2: Language Image Transformer consists of two streams, vision and language. The vision stream consists of a pre-trained ResNet as encoder. It generates visual representations for the input image. The language stream consists of a position encoder followed by a masked multi-head attention layer. It has similar architecture to the decoder in transformer. The masked multi-head attention layer predicts the word in the input sentence using the previous observed words. Finally, a Cross Modal Transformer is used to generate shared representations using the image and language representation from vision and language streams. The decoders uses the shared representations to predict the masked image embedding and word embedding simultaneously.

observations. It prevents the leftward information flow to preserve the autoregressive property. I mask out (set values to -$\infty$) the invalid positions in the input of the softmax.

### 5.3.3 Cross Modal Attention

The model uses a Cross-Modal Transformer layer to generate multi-modal embedding using the inputs from the image and language streams. The cross-modal transformer acts as an auto-regressive model to predict the future states using the observed states. ViLBERT [Lu *et al.* (2019)] uses the co-attentional transformer layer to learn a shared representation of intermediate visual and linguistic representations. The co-attentional transformer consists of two sub-modules. One module performs language-conditioned image attention in the visual stream, and the other performs image-conditioned language attention in the language stream.

The model passes keys and values from one modality to the other modality attention block to achieve it. Consider $Q_v$, $k_v$ and $V_v$ are the query, key and value matrices corresponding to image stream. The three matrices are linear transformations of the image representation $H_v$. Similarly, $Q_w$, $k_w$ and $V_w$ are the query, key and value matrices corresponding to the language stream. They are linear transformations of the language representation $H_w$.

In ViLBERT, the multi-head attention module in the image stream performs $attention(Q_v, K_w, V_w)$ and the language stream performs $attention(Q_W, K_v, V_v)$. I replace the two modules with a single two-stage operation. The model applies two linear transformations on the intermediate representation of image $H_v$ to obtain matrices $S_v, V_v$. Similarly, apply linear transformations on language representation $H_w$ to obtain the matrices $S_w, V_w$. It completes the first stage. Next, it computes the scaled dot product between the matrices $S_v$ and $S_w$. The resulting matrix signifies the

dependency between the inputs of the two modalities. It calculates two sets of scores by applying softmax across the rows and columns of the resulting matrix. The first set of scores obtained by $softmax(\frac{S_v S_w^T}{d_{model}})$ is scaled by value matrix $V_w$. It results in the output for the image stream. The other set of scores obtained by $softmax(\frac{S_w S_v^T}{d_{model}})$ is scaled by value matrix $V_v$ resulting in the output for the language stream.

## 5.4   Results

Papineni *et al.* (2002) introduced BLEU score as a metric for Machine Translation. The basic intuition behind BLEU score is that, closer the machine translated sentence is to the professional human translation, the better it is. It uses weighted average of variable length phrase matches against the reference translation. The BLEU score uses a modified $n$-gram precision score. The $n$-gram counts of the candidate sentences and their corresponding maximum overlap with the reference sentences are computed. The candidate counts are clipped by their corresponding maximum values and averaged across the total number of candidate $n$-grams. The BLEU score computes the geometric mean and adds a brevity penalty to discourage short sentences. BLEU4, which is a popular metric, uses 1-grams up to 4-grams . BLUE scores are usually computed at corpus level as the correlation of human judgement with the individual sentences might be low. In this work, I have used BLUE1 till BLUE4 for evaluating individual sentences generated by the model as part of image captioning.

Although BLEU score seems to be a good candidate for evaluating translations, it is not ideal for evaluating generated captions [Kulkarni *et al.* (2013)]. This is primarily because there can be vast variations in the captions that can be generated for a single image. As a result, there can be multiple instances where a semantically correct caption generated by model can be assigned a low BLEU score. Kulkarni

*et al.* (2013) demonstrated that the generated BLUE score in many cases were not in correlation with the human judgement.

ROGUE [Lin (2004)] is a similar metric used commonly for evaluating the generated sentences. Recall Oriented Understudy of Gisting Evaluation (ROGUE) follows a similar approach as BLEU but uses recall instead of precision for computing the scores. Metric for Evaluation of Translation with Explicit ORdering (METEOR) [Banerjee and Lavie (2005)] computes the clipped F-score based on the overlap with the set of references. It computes the similarity between the words using exact matches, stemming and semantic similarity.

In order to better align with the human judgements, I have also presented the CIDEr (Consensus-based Image Description Evaluation) score [Vedantam *et al.* (2015)]. CIDEr score uses a consensus protocol to measure the similarity of candidate sentences to a majority of how most people describe the image (reference sentences). CIDEr score initially computes the Term Frequency Inverse Document Frequency (TF-IDF) of the $n$-grams across the corpus of all captions. CIDEr score is then computed as the average cosine similarity between the candidate and the reference sentences, accounting for both precision as well as recall.

The figures 5.3, 5.4, 5.5 and 5.6 show qualitative examples from the generated captions of the model. The examples highlight that the model is able to understand the context and generate captions that is relevant to the image. The results can serve as an empirical evidence to infer that learning from both image and language embedding simultaneously helps in improving the quality of captions. The tables 5.1 and 5.2 show the experimental results of the image captioning obtained by applying Language Image Transformer (LIT) on the MS COCO and Conceptual Captions dataset respectively. The LIT model performs better than similar attention based models. The performance closely resembles the Soft and Hard Attention models. I

| Model | BLEU-1 | BLEU-2 | BLEU-3 | BLEU-4 | METEOR | ROGUE | CIDEr |
|-------|--------|--------|--------|--------|--------|-------|-------|
| LIT | .699 | .536 | .403 | .301 | .270 | .515 | .865 |
| BRNN | .642 | .451 | .304 | .203 | - | - | - |
| NIC | .666 | .461 | .329 | .246 | - | - | - |
| Soft | .707 | .492 | .344 | .243 | .239 | - | - |
| Hard | .718 | .504 | .357 | .250 | .234 | - | - |

Table 5.1: The table highlights the scores on the comparison metrics used to evaluate the Language Image Transformer (LIT) Model against similar models that use attention. BRNN was proposed by Karpathy and Li (2015), Google NIC from Vinyals *et al.* (2014) and soft/hard attentions proposed by Xu *et al.* (2015). All the models are evaluated using MS COCO.

| Model | BLEU-1 | BLEU-2 | BLEU-3 | BLEU-4 | METEOR | ROGUE | CIDEr |
|-------|--------|--------|--------|--------|--------|-------|-------|
| LIT | .191 | .95 | .059 | .040 | .077 | .177 | .327 |
| 5en_CE | - | - | - | - | - | .278 | 1.04 |
| Ultra | - | - | - | - | - | .26 | .98 |

Table 5.2: The table highlights the scores of Language Image Transformer (LIT) Model on Conceptual Captions dataset with top entries in the leader board of the Conceptual Captions .

hypothesize the reason behind the improvement in performance can be attributed to the use of unified model and a unified training objective.

### 5.4.1   Sharing of Parameters

Most of the other multi-modal models such as ViLBERT and VLBERT handle vision and language tasks independently and use separate streams for training using image and language modalities. With a unified architecture, there is better sharing of parameters between the vision and language streams. It is achieved in LIT using novel Cross Modal Attention layer that generates the shared Vision and Language representations simultaneously.

### 5.4.2   Unified learning objective

Unlike other multi-model representation models, LIT also uses a unified and self-supervised objective for learning shared representations. It drives the model to generate a shared representations that is predictive of individual modalities. The Visual stream uses a Contrastive Predictive Coding loss ($L_{CPC}$). The Language stream uses a standard Kullback-Leibler Divergence ($L_{KLD}$) loss between the predicted word and actual word. I discuss the model and its architecture in detail in sections 5.2 and 5.3.

The overall objective uses individual penalty parameters associated with vision and language objectives for controlling the influence of the individual vision and language embeddings on the generated shared embeddings.

$$L_{overall} = \alpha * L_{CPC} + \beta * L_{KLD} \tag{5.3}$$

In equation (5.3) , I determine the values of $\alpha$ and $\beta$ based on the downstream task. For Image captioning, I set the $\alpha$ to 0.01 and $\beta$ to 0.5.

| Image | Actual Caption | Generated Caption |
|---|---|---|
|  | a number of birds flying over a body of water | A flock of seagulls flying over the ocean. |
|  | A laptop and computer monitor with the same screen. | A laptop computer sitting on top of a desk. |
|  | Two giraffes walking around in the grass and dirt. | A giraffe standing in the middle of a dirt road. |
|  | A stainless steel kitchen sink on a black granite countertop | A kitchen with a stove, sink, and a stove. |

Figure 5.3: Visualizations for image caption prediction (MS COCO). Qualitative samples from the results. For each example, I show the input image (left), human-annotated caption (center), model-generated caption (right).

54

| Image | Actual Caption | Generated Caption |
|-------|---------------|-------------------|
|  | A man playing a game of tennis on a tennis court. | A man holding a tennis racket on a clay court. |
|  | The school bus is reflected in the rear view mirror | A large green and yellow train traveling down a track. |
|  | A train engine approaches a switch in a train yard. | A train traveling down train tracks next to a forest. |
|  | A person riding on the back of a horse walking across a field | A man riding a horse across a field. |

Figure 5.4: Visualizations of image captions from MS COCO. Qualitative samples from the results. For each example, I show the input image (left), human-annotated caption (center), model-generated caption (right).

| Image | Actual Caption | Generated Caption |
|---|---|---|
|  | person, on stage, playing a musical instrument, guitar, night and indoor | person, on stage, playing a musical instrument and indoor |
|  | the view i loved when riding the ferry | the view from the window of the house. |
|  | think outside of your box | the kitchen is a great example of a classic wood cabinets and stainless steel appliances. |
|  | the latest men 's designs from the label | a model walks the runway at the fashion show during fashion week. |

Figure 5.5: Visualizations of image captions from Conceptual Captions. Qualitative samples from the results. For each example, I show the input image (left), human-annotated caption (center), model-generated caption (right).

| Image | Actual Caption | Generated Caption |
|---|---|---|
|  | buses line up in in front in the early 1980s. | the main street is a city in the world 's most beautiful cities |
|  | players and staff show their respect by holding a minutes silence for football player who passed away today before the training session | players warm up during a training session ahead of the match. |
|  | actor attends the premiere of person | actor arrives at the premiere of thriller film held |
|  | small dock and boat at the lake | the sun rises over the horizon, the sky |

Figure 5.6: Visualizations of image captions from Conceptual Captions.Qualitative samples from the results. For each example, I show the input image (left), human-annotated caption (center), model-generated caption (right).

## 5.5 Conclusion

Attention mechanism has become an essential component in vision and language models. It is primarily due to the efficiency of attention models in generating powerful representations. It's relevance to both vision and language tasks, makes it an ideal candidate for associating the two modalities.

I proposed Cross Modal Transformer (CMT) and Language Image Transformer (LIT) models for learning shared embedding between vision and language. The LIT model learns shared representation by maximizing the mutual information between the modalities. The model uses self-supervised objective based on Contrastive Predictive Coding (CPC), for learning the shared representation. The objective drives the model to learn a representation that is predictive of both modalities simultaneously. I trained the model on popular captioning datasets, MS COCO and Conceptual-Captions. I evaluated the model in generating meaningful and relevant captions.

Multi-modal representation models can highly benefit from the improvements in self-supervised learning techniques and larger vision-language datasets. Enhancing the cross modal attention mechanism to handle a diverse number of tasks is an interesting area to explore in future.

# REFERENCES

Agrawal, A., J. Lu, S. Antol, M. Mitchell, C. L. Zitnick, D. Batra and D. Parikh, "VQA: Visual Question Answering www.visualqa.org", Tech. rep., URL `www.visualqa.org` (2015).

Alostad, H. and H. Davulcu, "Directional prediction of stock prices using breaking news on twitter", in "2015 IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology (WI-IAT)", vol. 1, pp. 523–530 (IEEE, 2015).

Anderson, P., X. He, C. Buehler, D. Teney, M. Johnson, S. Gould and L. Zhang, "Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering", in "Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition", pp. 6077–6086 (IEEE Computer Society, 2018).

Antol, S., A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. Lawrence Zitnick and D. Parikh, "Vqa: Visual question answering", in "Proceedings of the IEEE international conference on computer vision", pp. 2425–2433 (2015).

Bahdanau, D., K. Cho and Y. Bengio, "NEURAL MACHINE TRANSLATION BY JOINTLY LEARNING TO ALIGN AND TRANSLATE", Tech. rep. (2014).

Banerjee, S. 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", pp. 65–72 (2005).

Bello, I., B. Zoph, A. Vaswani, J. Shlens Quoc and V. Le Google Brain, "Attention Augmented Convolutional Networks", Tech. rep. (2019).

Chen, Y.-C., L. Li, L. Yu, A. E. Kholy, F. Ahmed, Z. Gan, Y. Cheng and J. Liu, "UNITER: Learning UNiversal Image-TExt Representations", URL `http://arxiv.org/abs/1909.11740` (2019).

Cutler, D. M., J. M. Poterba and L. H. Summers, "What moves stock prices?", Journal of Portfolio Management **15**, 2, 4–12 (1989).

Devlin, J., M.-W. Chang, K. Lee and K. Toutanova, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding", URL `http://arxiv.org/abs/1810.04805` (2018).

Gehring, J., M. Auli, D. Grangier, D. Yarats and Y. N. Dauphin, "Convolutional Sequence to Sequence Learning", Tech. rep. (2017).

He, K., X. Zhang, S. Ren and J. Sun, "Deep Residual Learning for Image Recognition", URL `http://arxiv.org/abs/1512.03385` (2015).

He, K., X. Zhang, S. Ren and J. Sun, "Deep residual learning for image recognition", in "Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition", vol. 2016-December, pp. 770–778 (IEEE Computer Society, 2016).

Hochreiter, S. and J. Schmidhuber, "Long Short-Term Memory", Neural Computation **9**, 8, 1735–1780 (1997).

Johnson, J., A. Karpathy and L. Fei-Fei, "DenseCap: Fully Convolutional Localization Networks for Dense Captioning", URL `http://arxiv.org/abs/1511.07571` (2015).

Kaiser, and S. Bengio, "Can active memory replace attention?", in "Advances in Neural Information Processing Systems", pp. 3781–3789 (Neural information processing systems foundation, 2016).

Karpathy, A. and L. Fei-Fei, "Deep Visual-Semantic Alignments for Generating Image Descriptions", Tech. rep. (2017).

Karpathy, A. and F.-F. Li, "Deep visual-semantic alignments for generating image descriptions.", in "CVPR", pp. 3128–3137 (IEEE Computer Society, 2015), URL `http://dblp.uni-trier.de/db/conf/cvpr/cvpr2015.html#KarpathyL15`.

Kingma, D. P. and J. L. Ba, "Adam: A method for stochastic optimization", in "3rd International Conference on Learning Representations, ICLR 2015 - Conference Track Proceedings", (International Conference on Learning Representations, ICLR, 2015).

Kulkarni, G., V. Premraj, V. Ordonez, S. Dhar, S. Li, Y. Choi, A. C. Berg and T. L. Berg, "Babytalk: Understanding and generating simple image descriptions", IEEE Transactions on Pattern Analysis and Machine Intelligence **35**, 12, 2891–2903 (2013).

Le, Q. V. and T. Mikolov, "Distributed representations of sentences and documents.", in "ICML", vol. 32 of *JMLR Workshop and Conference Proceedings*, pp. 1188–1196 (JMLR.org, 2014), URL `http://dblp.uni-trier.de/db/conf/icml/icml2014.html#LeM14`.

Lin, C.-Y., "ROUGE: A package for automatic evaluation of summaries", in "Text Summarization Branches Out", pp. 74–81 (Association for Computational Linguistics, Barcelona, Spain, 2004), URL `https://www.aclweb.org/anthology/W04-1013`.

Lin, T.-Y., M. Maire, S. Belongie, L. Bourdev, R. Girshick, J. Hays, P. Perona, D. Ramanan, C. L. Zitnick and P. Dollár, "Microsoft COCO: Common Objects in Context", URL `http://arxiv.org/abs/1405.0312` (2014).

Loper, E. and S. Bird, "Nltk: The natural language toolkit", in "In Proceedings of the ACL Workshop on Effective Tools and Methodologies for Teaching Natural Language Processing and Computational Linguistics. Philadelphia: Association for Computational Linguistics", (2002).

Lu, J., D. Batra, D. Parikh and S. Lee, "ViLBERT: Pretraining Task-Agnostic Visiolinguistic Representations for Vision-and-Language Tasks", URL `http://arxiv.org/abs/1908.02265` (2019).

Oord, A. v. d., Y. Li and O. Vinyals, "Representation Learning with Contrastive Predictive Coding", URL `http://arxiv.org/abs/1807.03748` (2018).

Papineni, K., 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", pp. 311–318 (Association for Computational Linguistics, 2002).

Pennington, J., R. Socher and C. D. Manning, "Glove: Global vectors for word representation.", in "EMNLP", vol. 14, pp. 1532–1543 (2014).

Poole, B., S. Ozair, A. v. d. Oord, A. A. Alemi and G. Tucker, "On variational bounds of mutual information", arXiv preprint arXiv:1905.06922 (2019).

Ramachandran, P., N. Parmar, A. Vaswani, I. Bello, A. Levskaya and J. Shlens, "Stand-Alone Self-Attention in Vision Models", URL `http://arxiv.org/abs/1906.05909` (2019).

Schuster, M. and K. K. Paliwal, "Bidirectional recurrent neural networks", IEEE Transactions on Signal Processing **45**, 11, 2673–2681 (1997).

Sharma, P., N. Ding, S. Goodman and R. Soricut, "Conceptual Captions: A Cleaned, Hypernymed, Image Alt-text Dataset For Automatic Image Captioning", URL `https://ai.google/research/pubs/pub47380/` (2018).

Sun, C., A. Myers, C. Vondrick, K. Murphy and C. Schmid, "VideoBERT: A Joint Model for Video and Language Representation Learning", URL `http://arxiv.org/abs/1904.01766` (2019).

Sutskever, I., O. Vinyals and Q. V. Le, "Sequence to sequence learning with neural networks", in "Advances in neural information processing systems", pp. 3104–3112 (2014), URL `https://papers.nips.cc/paper/5346-sequence-to-sequence-learning-with-neural-networks.pdf`.

Tetlock, P. C., M. Saar-Tsechansky and S. Macskassy, "More than words: Quantifying language to measure firms' fundamentals", The Journal of Finance **63**, 3, 1437–1467 (2008).

Vaswani, A., N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Kaiser and I. Polosukhin, "Attention is all you need", in "Advances in Neural Information Processing Systems", vol. 2017-December, pp. 5999–6009 (Neural information processing systems foundation, 2017).

Vedantam, R., C. Lawrence Zitnick and D. Parikh, "Cider: Consensus-based image description evaluation", in "Proceedings of the IEEE conference on computer vision and pattern recognition", pp. 4566–4575 (2015).

Vinyals, O., A. Toshev, S. Bengio and D. Erhan, "Show and tell: A neural image caption generator", URL `http://arxiv.org/abs/1411.4555`, cite arxiv:1411.4555 (2014).

Wang, L., Y. Li, J. Huang and S. Lazebnik, "Learning Two-Branch Neural Networks for Image-Text Matching Tasks", Tech. rep. (2017).

Wang, X., W. Jiang and Z. Luo, "Combination of convolutional and recurrent neural network for sentiment analysis of short texts", in "Proceedings of COLING 2016, the 26th international conference on computational linguistics: Technical papers", pp. 2428–2437 (2016).

Wu, Y. and K. He, "Group normalization", in "Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)", vol. 11217 LNCS, pp. 3–19 (Springer Verlag, 2018).

Xu, K., J. L. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhutdinov, R. S. Zemel and Y. Bengio, "Show, attend and tell: Neural image caption generation with visual attention", in "32nd International Conference on Machine Learning, ICML 2015", vol. 3, pp. 2048–2057 (International Machine Learning Society (IMLS), 2015).

Zhu, Y., R. Kiros, R. Zemel, R. Salakhutdinov, R. Urtasun, A. Torralba and S. Fidler, "Aligning books and movies: Towards story-like visual explanations by watching movies and reading books", in "Proceedings of the IEEE international conference on computer vision", pp. 19–27 (2015).

APPENDIX A

ADDITIONAL SAMPLES

| Image | Generated Caption | Actual Caption |
|-------|-------------------|----------------|
|  | A man sitting at a table with a plate of food | a little girl in a blue apron and chefs hat and a girl in a red apron and hat |
|  | A zebra standing in a zoo enclosure with a zebra. | A zebra that is bending it 's neck backwards to reach it 's tail. |
|  | A black and white image of a cat sitting on a park bench. | Three dogs sitting on the levels of an empty tiered garden. |
|  | A man flying through the air while riding a skateboard | A MAN IS JUMPING ON HIS SKATE BOARD IN THE SKY |

Figure A.1: Visualizations for images to caption prediction.Qualitative samples from the results. For each example, we show the input image (left), human-annotated caption (center), model-generated caption (right). The samples are from the MS COCO captioning dataset.

| Image | Actual Caption | Generated Caption |
|---|---|---|
|  | a man wearing a traditional dress and a traditional dress | woman selling food by the railway line |
|  | the interior of the church 's largest temple complex. | our homage to the necklace and earrings won the aesthetic first prize at award. |
|  | person poses for a photo with a black suit and white shirt and white shirt. | before and after of person the mid 1800 's, carefully restored and preserved for generations to come |
|  | portrait of a boy lying on the beach with a blue sky | woman relaxing at the sea |

Figure A.2: Visualizations for images to caption prediction.Qualitative samples from the results. For each example, we show the input image (left), human-annotated caption (center), model-generated caption (right). The samples are from the Conceptual Captions dataset.