# EMOTION DRIVEN

# IMAGE CAPTIONING

*Submitted in partial fulfilment of the requirements of the degree of*

*(Bachelor of Technology)*

*by*

Nabagata Saha (167140)

Yesantarao Venkata Akhila (167170)

Project Supervisor

Dr. P Radha Krishna

Professor



Department of Computer Science and Engineering,

National Institute of Technology Warangal

(2019-2020)

# APPROVAL SHEET

This Project Work entitled Emotion Driven Image Captioning by Yesantarao Venkata Akhila and Nabagata Saha is approved for the degree of Bachelor of Technology in Computer Science and Engineering at National Institute of Technology Warangal during the year 2019-2020.

**Examiners**

_____

_____

_____

**Supervisor**

_____

Dr. P Radha Krishna

Professor, CSE Department

**Chairman**

_____

Date: _____

Place: _____

# **DECLARATION**

I declare that this written submission represents my ideas in my own words and where others' ideas or words have been included, I have adequately cited and referenced the original sources. I also declare that I have adhered to all principles of academic honesty and integrity and have not misrepresented or fabricated or falsified any idea/ data/ fact/ source in my submission. I understand that any violation of the above will be cause for disciplinary action by the Institute and can also evoke penal action from the sources which have thus not been properly cited or from whom proper permission has not been taken when needed.

Name: Nabagata Saha
Roll Number: 167140
Date: 16/6/20

Name: Yesantarao Venkata Akhila
Roll Number: 167170
Date:16/6/20

# NATIONAL INSTITUTE OF TECHNOLOGY WARANGAL
## DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING

## CERTIFICATE

This is to certify that the project work entitled "EMOTION DRIVEN IMAGE CAPTIONING" is a bonafide record of work carried out by Yesantarao Venkata Akhila (167170) and Nabagata Saha (167140), submitted to the faculty of Computer Science and Engineering Department, in fulfilment of the requirements for the award of the degree of Bachelor of Technology in Computer Science and Engineering at National Institute of Technology, Warangal during the academic year 2019-2020.

Dr. P Radha Krishna,                                                    Dr. P Radha Krishna,
Project Guide,                                                    Head of the Department,
Department of CSE,                                                    Department of CSE,
NIT Warangal                                                    NIT Warangal

# **ACKNOWLEDGEMENT**

We would like to extend our gratitude to many people who helped and guided us along the course of completion of this project. First, we would like to thank our project partners for always cooperating and coordinating with us. It would not have been possible to complete this project without their deep involvement and dedication.

We are grateful to our Project Guide and Head of the Department, Dr. P Radha Krishna for guiding us through all the problems that we faced and pushing us to do our best. We would also like to thank him for providing us the opportunity to work under his guidance and being available to give us support whenever required in terms of invaluable advice and resources.

We would like to extend our thanks to the Project Evaluation Committee, for their strenuous efforts to evaluate our project.

We wish to express our regards to the Department for providing us with all the necessary resources and the support throughout the course of this work.

# ABSTRACT

Image captioning has been a challenging area, for generating captions which resemble very closely to how a human would caption an image. The state-of-the-art exists in factual captions, which is a great tool to caption images that contain inanimate objects. But captioning images with humans, using their facial expressions, remains a field that has not been tinkered into.

The current state-of-the-art architecture by Xu et al [21] called show attend and tell is an encoder-decoder network with soft attention. The visual features in an image are extracted by the encoder and next attention is applied to them so that different parts of the image are fed to the decoder at every step. This helps generate better captions since the features fed to the decoder are limited. However, the features extracted are limited to generic categories and so no features specific to human faces are learnt. Due to this reason, it is observed that the model by Xu et al [21] does not include adjectives that describe the emotion of the human being in the image.

This work proposes a novel method that realizes this task. The facial features of the human subject present in the image are extracted along with image features and fed to an image captioning model. The caption generated is more relevant and human-like. To add emotion to the generated caption, we performed facial expression recognition which is a popular computer vision problem.

The emotion is recognized by a deep learned model and image is captioned by an encoder-decoder network. A multi-level VGG network is used for FER such that it extracts facial features and Inception V3 (encoder) is used to extract the visual features. These features are fed to an attention tuned Gated Recurrent Unit (decoder), to produce the caption in a word by word manner. This helps more realistic captioning of images which in turn can be used to generate natural sounding video summaries. Moreover, customer satisfaction can be gauged from such captioned images.

This work also digs into various experiments that were conducted in the FER problem to achieve an accuracy of 73.11%. Along with this, we have added novelty to standard image captioning techniques. By including the facial features in the encoder-decoder network, facial expressions have been successfully incorporated into image captions.

# **INDEX**

# **APPENDIX**

## 1. List of Figures

## 2. List of Tables

# CHAPTER 1

# INTRODUCTION

## 1.1 Facial Expression Recognition

Facial expression is a form of conveying a human being's emotion to another human being. It is a contortion of the facial features to express emotions. They play a significant role while human beings converse with each other. It adds in a factor of empathy and helps determine what the other party truly feels in non-verbal communication. As per multiple surveys, one-third of human communication is verbal, while the rest is non-verbal. This emphasizes the importance of facial expression recognition.

The facial expression recognition (FER) problem has been an active field of interest due to two main reasons; its undeniable importance in human-computer interaction and due to the ambiguous nature of the problem. The basic expressions include Anger, Disgust, Fear, Happy, Sad, Surprise and Neutral.

The various advantages and applications of FER include:

- Recognizing a facial expression can aid a virtual assistant in understanding the mood and the needs of its user much better. This also extends to recognizing facial expression in images or live feed and utilizing this information to form more human-like sentences, in the case of image captioning.
- It will be of utmost importance while computer agents handle grievances and feedback of customers. It will be able to gauge the customer's reaction.
- Its applications cover all the activities that involve interaction between a human and a computer, it adds an element of humanity to the computer's personality which makes robots more intelligent.
- It helps understand the effectiveness of teaching methods based on the facial expressions of students. A similar logic can be applied to audience engagement in response to a video.
- Overall, the ability to recognize emotions using facial expressions can benefit a company in serving its customer base with much more intelligent and humane products or services.

These emphasize the importance of FER.

While recognizing a facial expression, there are a few expressions which cannot be uniformly classified by different persons. Recognizing expressions in natural conditions adds on to the complexity of the problem. These complications occur due to variations in head positions, illumination, and the fact that unposed expressions are subtle. This raises the issue of ambiguity in solving the problem due to which the standard accuracy obtained on a very famous dataset is limited to 67% [6], in spite of using convolutional neural networks which usually perform extremely well in image processing problems. Figure 1 shows the generic steps involved in FER using CNNs. A digital image is interpreted by a computer as a matrix of pixel values. In the lower

layers of a CNN, it learns small features like edges, but as the network goes deeper, these edges are stacked to represent abstract shapes which often materialize to form important features such as eyes, nose, mouth and so on. Their filter weights are optimized to detect such features.



*Figure 1 Basic steps in FER*

CNNs can perform much better, given there are changes made to the architecture of the network and the data-preprocessing phase. These changes will improve the dataset's versatility and the features extracted by the CNN.

## 1.2 Image Captioning

Image captioning combines concepts from two important artificial intelligence domains; computer vision and natural language processing, to describe the objects and the relationship between them in an image. It is the task of generating captions for an image while ensuring the structure of the sentence remains intact.

The applications of image captioning are wide. These include automatically captioning images uploaded on social media platforms such as Facebook and Twitter, which helps to detect posts that violate the policy of either platform. This can be used in image search, where every image in the search index is captioned and on receiving an image query, it is easier to retrieve results. Generating descriptions of videos uploaded on YouTube is another use case. By solving this problem, we can solve many more real-world problems.

To generate semantically correct captions the dataset must be large, and an encoder-decoder network must be used. The visual features are extracted using the encoder which is a CNN. An important aspect is to focus on the prominent parts of these features, understand which words in the ground truth captions describe them. Then using attention mechanism, the prominent features are focused on and fed to the decoder as shown in figure 2. Then the decoder generates the caption word by word.

*Figure 2 Basic steps in image captioning*

## 1.3 Objective and Problem Statement

The captions generated by current image captioning models are factual in nature. They do not have the emotional aspect that human beings employ while conversing with one another. No adjectives are used to describe the subjects identified. By giving due importance to facial expressions, we can generate coherent and semantically valid captions that describe the image in a better manner and sound human-like. This would improve the quality of the captions in general.

This task is very hard in nature because the image captioning model needs to first understand the image by extracting its features, identify facial expressions if any and then find words to describe these features detected. To add emotion to the captions, more attention needs to be given to the facial features and the words that denote them. This thesis focuses on solving this problem, which is to add emotions to image captions.

## 1.4 Proposed Method

The image captioning dataset is modified to accommodate emotions. This is done by adding an emotion describing the facial expression before the subject. An encoder decoder network is used to generate captions. The encoder expresses the input image as convolutional feature maps, each describing various parts of the image. Simultaneously the image is fed to a FER network that extracts the facial features. Both these features are fed to the attention-based decoder that learns

how to generate the captions, word by word. The attention mechanism focuses on different image features and facial features at every time step and the decoder uses them to generate every word.

## 1.5 Contributions of the Thesis

The main contributions in this work are:

- Various experiments are carried out on the FER problem to improve the accuracy of the expression recognition network. Multi-level architectures were used to obtain a FER network with accuracy of 73.11%.
- Modifications are made to MS COCO to accommodate emotions.
- A novel approach to add emotion to image captions is proposed. It is an end-to-end trainable neural network with various components.

## 1.6 Organization of the Thesis

The rest of the thesis is organized as follows. Chapter 2 discusses the related work in this literature, Chapter 3 describes the specifications of the datasets used. This is followed by Chapter 4 that provides details about the techniques used to implement a FER network that achieves good accuracy. Chapter 5 digs into the proposed methodology and its implementation. Results are presented and analyzed in Chapter 6 and Chapter 7 concludes the thesis and provides insights into future scope of work. The references are stated in the last section.

# CHAPTER 2

# LITERATURE REVIEW

## 2.1 Facial Expression Recognition

A large body of work has emerged on recognizing facial expressions, usually using a purportedly universal framework of a small number of basic emotions (happiness, sadness, fear, surprise, anger, and disgust) or this set including a neutral expression (Field, Woodson, Greenberg, & Cohen, 1982 [5]; Kanade, Cohn, & Tian, 2000 [23] ; Fasel & Luettin, 2003 [3]) or more fine-grained facial features such as facial action units (AUs), defined as the deformations of facial muscles (Tian, Kanade, & Cohn, 2001 [16]).

| Title | Authors | Approach |
|---|---|---|
| Recognizing AUs for Facial Expression Analysis | Tian, Kanade, & Cohn, 2001 [16] | Their system recognizes fine-grained changes in facial expressions called facial action units. Using these features, they propose a method to do FER. |
| Deep learning using linear support vector machines | Tang Y (2013) [15] | Uses CNN with linear SVM instead of softmax. Margin-based loss is optimized to maximize margins among different data points. Won FER2013 challenge. |
| Multimodal deep learning approaches for emotion recognition in video. | Kahou, Bouthillier, Lamblin, Gulcehre, Michalski, Konda, Jean, Froumenty, Dauphin, and Boulanger-Lewandowski (2016) [8] | Features extraction and transfer learning is applied to learn from diverse modalities such as video, audio and spatio-temporal information. Won EmotiW challenge. |
| Facial expression recognition using convolutional neural networks: State of the art. | Pramerdorfer and Kampel (2016) [13] | They used an ensemble of various CNNs like VGG, ResNet and Inception. |

*Table 1 A summary of related work in FER*

To find effective representations, deep learning-based methods have been recently successful in this domain. Due to their complex architectures including multiple layers, they can capture hierarchical structures from low- to high-level representations of facial expression data. Tang (2013) [15], the winner of the 2013 Facial Expression Recognition (FER) challenge (Goodfellow et al., 2013 [6]), trained a Convolutional Neural Network (CNN) with a linear support vector machine (SVM) to detect facial expressions. He replaced the softmax layer, used to generate a probability distribution across multiple classes, with a linear SVM and showed a consistent improvement compared to the previous work. Instead of cross-entropy loss, his approach optimizes a margin-based loss to maximize margins among data points belonging to diverse classes.

Convolutional Neural Networks (CNNs) are also used for feature extraction and transfer learning in this domain. Kahou, Bouthillier, Lamblin, Gulcehre, Michalski, Konda, Jean, Froumenty, Dauphin, and Boulanger-Lewandowski (2016) [8] applied a CNN model to recognize facial expressions and won the 2013 Emotion Recognition in the Wild (EmotiW) Challenge. Their approach uses a combination of deep neural networks to learn from diverse data modalities including video frames, audio data and spatio-temporal information (Kahou, Pal, Bouthillier, Froumenty, Gulcehre, Memisevic, Vincent, Courville, Bengio, & Ferrari, 2013 [9]).

Most models usually use conventional CNN architectures to report the performance on different facial expression recognition datasets including the FER-2013 dataset (Goodfellow et al., 2013 [6]). This dataset is a publicly available with a huge number of human faces collected in the wild condition. Pramerdorfer and Kampel (2016) [13] instead used an ensemble of very deep architectures of CNNs such as VGG Net, Inception and Residual Network (Res Net) by identifying the bottlenecks of the previous state-of-the-art facial expression recognition models on the FER-2013 dataset and achieving a new state-of-the-art result on the dataset. The quality of these recent models is high: it is at least as good as human performance (Goodfellow et al., 2013 [6]).

Another interesting approach to recognize facial expression is by using multi-level CNNs. The idea of applying features detected at various levels in a CNN motivates our work to make a facial expression recognition module. We use the module to extract facial features from human faces to apply in our image captioning models.

## 2.2 Image Captioning

Image captioning has received immense attention over the years and remains a problem worth solving. There are different types of models which have been implemented, such as, template-based, retrieval-based and deep-learning based.

| Title | Author | Approach |
|---|---|---|
| Every picture tells a story: Generating sentences from images | Farhadi, Hejrati, Sadeghi, Young, Rashtchian, Hockenmaier, & Forsyth, 2010 [2] | Object attributes and relations between them were detected. These are expressed as words and used to fill in the pre-defined caption template. |
| Framing image description as a ranking task: Data, models and evaluation metrics | Hodosh, Young, & Hockenmaier, 2013 [7] | They used datasets which contain images and captions. These were used to find the similar images and the most appropriate caption. |
| Unifying visual-semantic embeddings with multimodal neural language models | Kiros, Salakhutdinov, & Zemel, 2014 [11] | An encoder-decoder network was used. The visual features were extracted using the encoder and captions were generated word by word using the decoder. |
| Image Captioning with Semantic Attention | You et al [18] | A CNN is first used to detect visual features and simultaneously detect visual regions of interest. A semantic attention model is used to combine the visual features with the visual regions in an RNN. This generates the caption |
| Show, Attend and Tell: Neural Image Caption Generation with Visual Attention | Xu et al [22] | Fine-grained fragments are detected using a soft attention mechanism. Instead of feeding features corresponding to the whole image to the decoder, different regions are used. This is inserted in an encoder-decoder network. |

*Table 2 Summary of related work in Image Captioning*

Template-based ones first detect visual objects, their attributes and relations and then fill a predefined template's blank slots (Farhadi, Hejrati, Sadeghi, Young, Rashtchian, Hockenmaier, & Forsyth, 2010 [2]). Retrieval-based ones generate captions using the available captions corresponding to similar images in their corresponding datasets (Hodosh, Young, & Hockenmaier, 2013 [7]). These classical image captioning models have some limitations. For example, template-based ones cannot generate a wide variety of captions with different lengths, and retrieval-based ones are not able to generate specifically designed captions for different images. Moreover, classical models do not incorporate the detection and generation steps using an end-to-end training approach.

Because of these limitations, modern image captioning models using deep learning are currently the most popular. Modern image captioning models usually use an encoder-decoder paradigm (Kiros, Salakhutdinov, & Zemel, 2014 [11]). They apply a top-down approach where a CNN model learns the image content (encoding), followed by a LSTM generating the image caption (decoding). This follows the paradigm employed in machine translation tasks, using deep neural networks (Sutskever, Vinyals, & Le, 2014 [14]), to translate an image into a caption. This top-down mechanism directly converts the extracted visual features into image captions (Chen & Lawrence Zitnick, 2015 [1]). But by using a top-down approach, it is difficult to pay attention to fine-grained and critical elements of the image and generate better image captions. To solve this problem, a combination of a top-down approach and a bottom-up approach, inspired from the classical image captioning models, is proposed by You et al [18]. This limitation is overcome by the bottom-up approach, by generating the relevant words and phrases that are detected from images and these are combined to generate an appropriate caption.

To attend to fine-grained fragments, attention-based image captioning models have been proposed (Xu et al., 2015 [22]). These kinds of approaches usually analyze different regions of an image in different time steps of a caption generation process, in comparison to the initial encoder-decoder image captioning systems which consider only the whole image (Vinyals & Le, 2015 [17]) as an initial state for generating image captions. The spatial information of visual data can also be taken into consideration to provide the relevant words and phrases to generate the caption. The current state-of-the-art models in image captioning are attention-based systems (Xu et al [22]).

## 2.3 Summary

Various methods that have been used to solve FER and image captioning are summarized in table 1 and table 2 respectively. A clear evolution in the approach is seen as deep learning methods came into popularity.

In FER, the first method used detection of facial action units using appropriate datasets. The AUs were detailed in nature and transitions in facial expressions were also considered. As CNNs came into the picture, the detection of such AUs became inherent in deeper layers. Multiple modifications were made to standard CNNs to help improve results. Different forms of information were taken into consideration to obtain better features. Finally, ensemble methods, combining

various CNNs like VGG and ResNet were able to achieve the best accuracy. A multi-level architecture is proposed in this work.

In image captioning, the initial methods focused on filling in caption template and producing a caption based on retrieval techniques. Then we advanced to encoder-decoder networks which are deep learning based methods which generated captions with similarity to how human beings caption. Modifications were suggested to the encoder-decoder network such as combination of visual features and visual regions, and the usage of attention mechanism to generate better results.

# CHAPTER 3

# DATASETS USED

## 3.1 FER2013 Dataset

The FER2013 dataset is a large dataset that was created for the Facial Expression Recognition Challenge. It contains 35887 images for each of the seven facial expressions (Anger, Disgust, Happy, Sad, Surprise, Neutral) along with 28709 for training, 3589 for private testing and the rest for public testing. Each image is a 48x48 cropped grayscale face. One of the issues faced while using this dataset is the data unbalancing. As shown in figure 3, the number of images labelled as Disgust is relatively less as compared to the others. Figure 4 shows some examples from the training set of FER2013.



*Figure 3 Data distribution in FER2013*



*Figure 4 Examples from FER2013*

## 3.2 MS COCO Dataset

The Microsoft Common Objects in Context (MS COCO) is a large-scale dataset that is used for a variety of tasks such as object detection, segmentation, and image captioning. It consists of 3,30,000 images and each has five captions with 80 object categories. This dataset is widely used to train good quality image captioning models due to the sheer volume of data that it holds.

The 2014 dataset split is used in the study, MS COCO2014. This contains 82783 images in the training set with 413915 captions, 40504 images in the validation set with 202520 captions and 40775 images in the test set. The annotations were released for the training and validation set. The captions for each image were generated by crowdsourcing and the five ground truth captions were decided by the majority vote of experts. Figure 5 shows some of these captions, along with their images.



a guy touching a scarf around the neck of a statue

two dogs sitting in the back of a vehicle

two elephants performing with a male entertainer in between them

a man poses for a picture while sitting on a motorcycle

*Figure 5 Examples of images with captions from MS COCO*

## 3.3 Summary

The datasets used in this work have been explained in sections 3.1 and 3.2. FER2013 dataset was used to train a CNN to classify between the seven facial expressions which are Anger, Disgust, Fear, Happy, Sad, Surprise and Neutral. It contains 35887 cropped human faces in grayscale.

The MS COCO dataset was partitioned such that only images containing human faces were present. This dataset was then used to train the decoder network to generate captions. The dataset can be used for other purposes such as object detection and segmentation. It contains 330000 images each with approximately five captions.

# CHAPTER 4

# BACKGROUND AND EXPLORATION

In this chapter, we provide details about the techniques used to train a FER classifier and about the language model, which form the basis of the proposed methodology. Their implementation details are presented as well.

## 4.1 Data Preparation

Affine transformations are applied to remove any camera induced distortions. This produces an image that is invariant to which angle the picture has been captured. These include translation, rotation, isotropic shearing, and scaling.

Data augmentation is used to artificially expand the size of the training dataset. This helps the model apply what it has learnt to new images, generalize better and perform better. In this method, new training images are created from existing images by applying some transforms. Transforms include flipping, crops, shift, zooms and much more.

The various approaches to perform affine transformations and data augmentation can be classified as image manipulations. Figure 6 shows various transformations applied to the image of a young girl. These include flips, zooms, crops and rotations by specific number of degrees. Some of them are used on the training images in the FER2013 dataset to perform data augmentation. Random crop of size 44 on the original image of size 48, rotation of original image by 1 degree, zooming in by 0.1 units and horizontal flip are used. These transformations help eliminate the effect of image irregularities such as varied illumination, occlusions, and the influence of different head positions.



*Figure 6 Image manipulations such as flip, crops, zooms and rotate*

In the testing phase, a deep-learned face detection model called MTCNN (Multi-Task Cascaded Convolutional Networks) is used. MTCNN uses three neural networks to accurately detect faces and the output of one is fed to the next to eliminate false detections and to obtain face landmarks. The largest face is cropped from the given image. This is followed by converting the cropped face to grayscale and resizing it to 48x48. After the above transformations, any image can be supplied to the FER model. The FER model expects a 48×48 sized grayscale face as input.

## 4.2 Conventional Architecture

Convolutional Neural Networks (CNNs) are powerful tools in computer vision and perform satisfactorily on most image classification problems. They can differentiate between different types of objects by assigning importance to certain aspects of each object. Unlike standard tasks in image processing like edge detection where filter weights have been fixed after tremendous experimentation, CNNs learn the filters while training. Their architecture is inspired by the organization of the visual cortex.

The following section describes two CNN architectures that were used in this work. These are the base models to our proposed model.

## 4.2.1 Modified Mini Xception

The Xception architecture was released by Google and stands for extreme version of Inception. It outperformed all the previous models in the ImageNet classification problem due to the modified depthwise separable convolution. The model is represented as a stack of modules, unlike previous ones represented as a stack of convolutional layers.

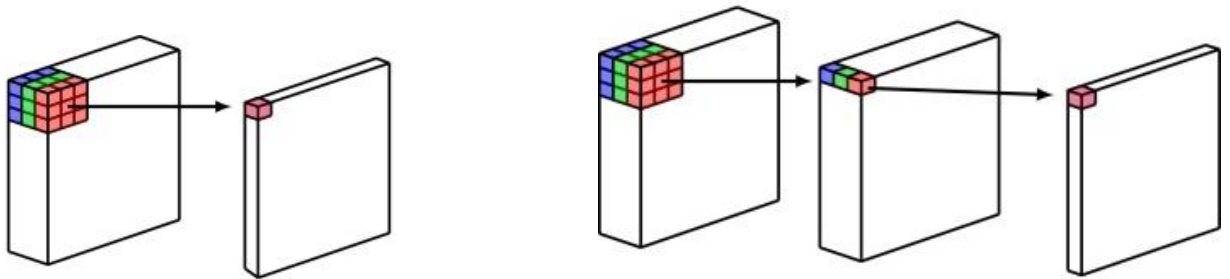

*Figure 7 Convolution and depthwise separable operations*

Figure 7 represents a simple convolution operation on an RGB image. A filter (or kernel) of dimensions n×n×3 slides along the three channels, and one cell of the output feature map is computed by multiplication and addition of pixel values with filter weights. This operation is repeated m times to obtain an output with m channels.

24

A separable convolution operation consists of a depthwise convolution followed by a pointwise convolution. The depthwise convolution is performed for each channel separately using a nxnx1 spatial filter. Then a pointwise convolution using a $1\times1\times n$ spatial filter is carried out. It can increase the number of channels to m by performing pointwise convolution using the $1\times1\times n$ filter m times. This produces the same result as a normal convolution. The advantage of using such an operation is that the number of multiplications involved is less than half. This is because in separable convolutions, the image is transformed only during the depthwise convolution. The pointwise convolution simply increases the number of channels.

This architecture uses a modified separable convolution where in, a pointwise convolution is followed by a depthwise convolution (see figure 8). The pointwise convolution is the $1\times1$ convolution to change the dimension. The depthwise convolution is channel-wise $n\times n$ spatial convolution, which makes it unnecessary to perform convolution across all channels. This makes the model light with a fewer number of connections as shown in the below diagram. Furthermore, there is no ReLU non-linearity in the middle due to the modified separable convolutions.



*Figure 8 Modified separable convolution operation*

The Xception architecture consist of 36 such convolutional layers which perform feature extraction. These modified separable convolutional layers are linearly stacked with residual connections between the pointwise convolutions and the depthwise convolution in every block as shown in figure 9.

Figure 9 depicts the architecture of the modified mini Xception model. The modified mini Xception consists of two convolution layers followed by three blocks and it accepts as input a grayscale face of dimensions $48\times48\times3$. This constitutes the entry flow and is represented as block 1. In every block, there are two modified depthwise separable convolutions. Each of them is preceded by a ReLU activation function and followed by a batch normalization layer. There is a max pool layer at the end of each block. These blocks are labelled as block 2, block 3, block 4. The ReLU (Rectified Linear Unit) activation function activates only when important information

is detected, all negative values are made zero. Pooling layers help retain the important features and reduce the size of the feature maps also called downsizing.

The exit flow consists of global average pooling and a fully connected layer with 7 units to predict the probabilities for each class. The convolutional layers in the lower blocks detect low level features like edges in the image. High level features are detected as the number of blocks increase and all the lower-level features are stacked. These include eyes, nose, lips, eyebrows.

Input of 48x48x3

32x3x3 conv2d
Block 1
64x3x3 conv2d

128x3x3 sepconv2d
128x1x1 conv2d
Block 2
128x3x3 sepconv2d

256x3x3 sepconv2d
256x1x1 conv2d
Block 3
256x3x3 sepconv2d

512x3x3 sepconv2d
512x1x1 conv2d
Block 4
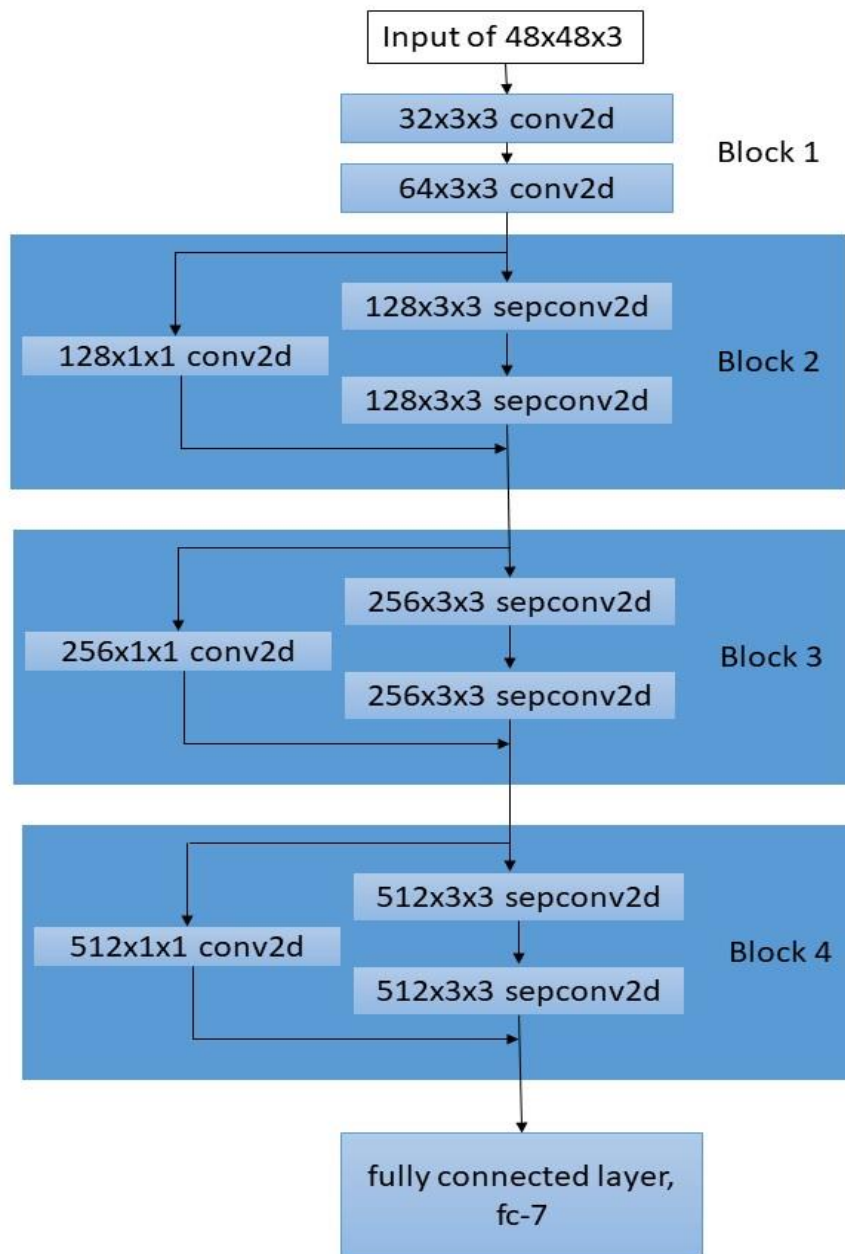512x3x3 sepconv2d

fully connected layer, fc-7

*Figure 9 Modified mini xception architecture*

## 4.2.2 VGGNet 19

The Verified Geometric Group in the Oxford University had come up with this architecture. VGGNet 19 as the name suggests, consists of 19 layers, 16 convolution layers and 3 fully connected layers. This architecture won the ILVSR 2014 competition. It was the first of its kind, a deeper network with smaller filters that adapted well to the images in the training set. It has other variations, starting from VGGNet 11, VGGNet 13 and VGGNet 16. Each had an increase in the number of layers. It was observed that adding more layers than those in VGGNet 19 resulted in some increase in error rate and thus no more layers were added.

It is a deep network and uses small filters (3×3) with a stride of 1 unlike the other CNNs designed which were shallower and had big filters (7×7). Due to stride being 1, the filters convolve at every pixel value. The usage of three 3x3 filters was equivalent to one 7×7 filter since both had the same receptive field. This resulted in fewer parameters. This led to faster convergence and reduced the problem of overfitting due to fewer number of parameters.

VGG is a standard CNN in the field of computer vision. It is often used for other computer vision tasks by the method of transfer learning where new layers are added and only their weights are learnt. Otherwise, the architecture of VGG is adopted and trained on a completely new dataset, like in our case.

This architecture can be used to solve the FER problem as well. It is shown in figure 10. It consists of 17 convolutional layers and each such layer is followed by a batch normalization layer and a ReLU activation layer. There is a max pool layer at the end of each block. This helps downsize the feature maps and makes the network more robust so that it can identify features irrespective of their location in the feature map. The entry flow consists of two convolution layers and is called block 1. Block 2, block 3, block 4 and block 5 are designed in a similar fashion with the number of channels being doubled at each. Block 2 contains three convolutional layers, while the rest contain four each. The exit flow consists of a fully connected dense layer with 7 units to predict the probabilities for each of the seven emotions. After an image of 48×48 dimension is put through the VGG network, the important features are retained while the irrelevant ones are removed via the convolutional, pooling and ReLU activation layers.

Input of 48x48x3

**Block 1**
64x3x3 conv
64x3x3 conv

**Block 2**
128x3x3 conv
128x3x3 conv
128x3x3 conv

**Block 3**
256x3x3 conv
256x3x3 conv
256x3x3 conv
256x3x3 conv

**Block 4**
256x3x3 conv
256x3x3 conv
256x3x3 conv
256x3x3 conv

**Block 5**
512x3x3 conv
512x3x3 conv
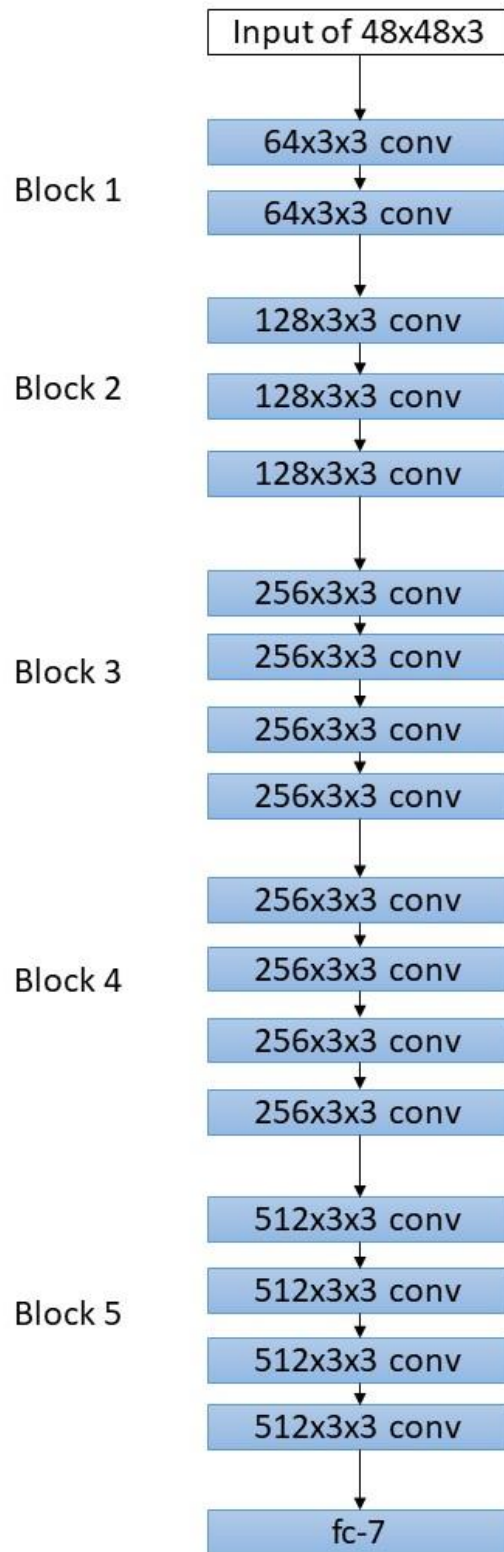512x3x3 conv
512x3x3 conv

fc-7

*Figure 10 VGGNet 19 architecture*

## 4.3 Multi-Level Architecture

An interesting approach to enhance the features extracted from an image in convolutional neural networks is to use a multi-level CNN. As an image passes through various layers in the CNN, the features extracted from the previous modules are further streamlined and some information is dropped out at each module due to the max pool and activation layers. Thus, with higher number of convolutional layers, better features are extracted. However, the features extracted after each module represent essential information about the image. They can be classified as low-level, medium-level and high-level features. Low-level features include edges, background whereas the high level-features would include lips, eyes, facial muscles in our case. In a multi-level CNN all these feature maps are considered while finally being used to learn the expression corresponding to a face.

Gradient Class Activation Mappings (grad-cam) can be used to visualize the regions of interest to a network with attention levels. It helps understand at which part of the image the CNN places its attention in every layer or block of the network, so that CNNs are not hard to interpret anymore. It uses the gradient information flowing into a convolution layer during backpropagation to understand which neurons have made which decision.

### 4.3.1 Multi-Level Mini Xception



*Figure 11 Gradcam visualization for mini xception*

The grayscale image of a man smiling was fed to the modified mini xception model. The grad-cam representations from block 2 and block 3 were obtained and are shown in figure 11. In the heatmap, red represents area with highest attention, followed by yellow and then blue. The grad-cam representation of the features identified in block 2 is quite interesting. It identifies the arc line of the smile as a medium attention level feature. A spot of red is visible at the chin, at the end of the arc line on the left side of the face. However, it finds some important information at the nose as well. The grad-cam representation of the features identified in block 3, explain the previous features better. The right side of the face containing the smile line is considered important and the

teeth are given adequate attention. Thus, it shows us that the network focuses on important features of the face that relate to which expression is being portrayed.

The block 2 features can be considered as medium-level features and the block 3 features can be considered as high-level features. By concatenating both these features, attention is given to the left side of the face as well, strengthening the model's prediction that the man is happy. Thus, the features obtained after block 2 and block 3 can help the model choose better features to base its prediction on.
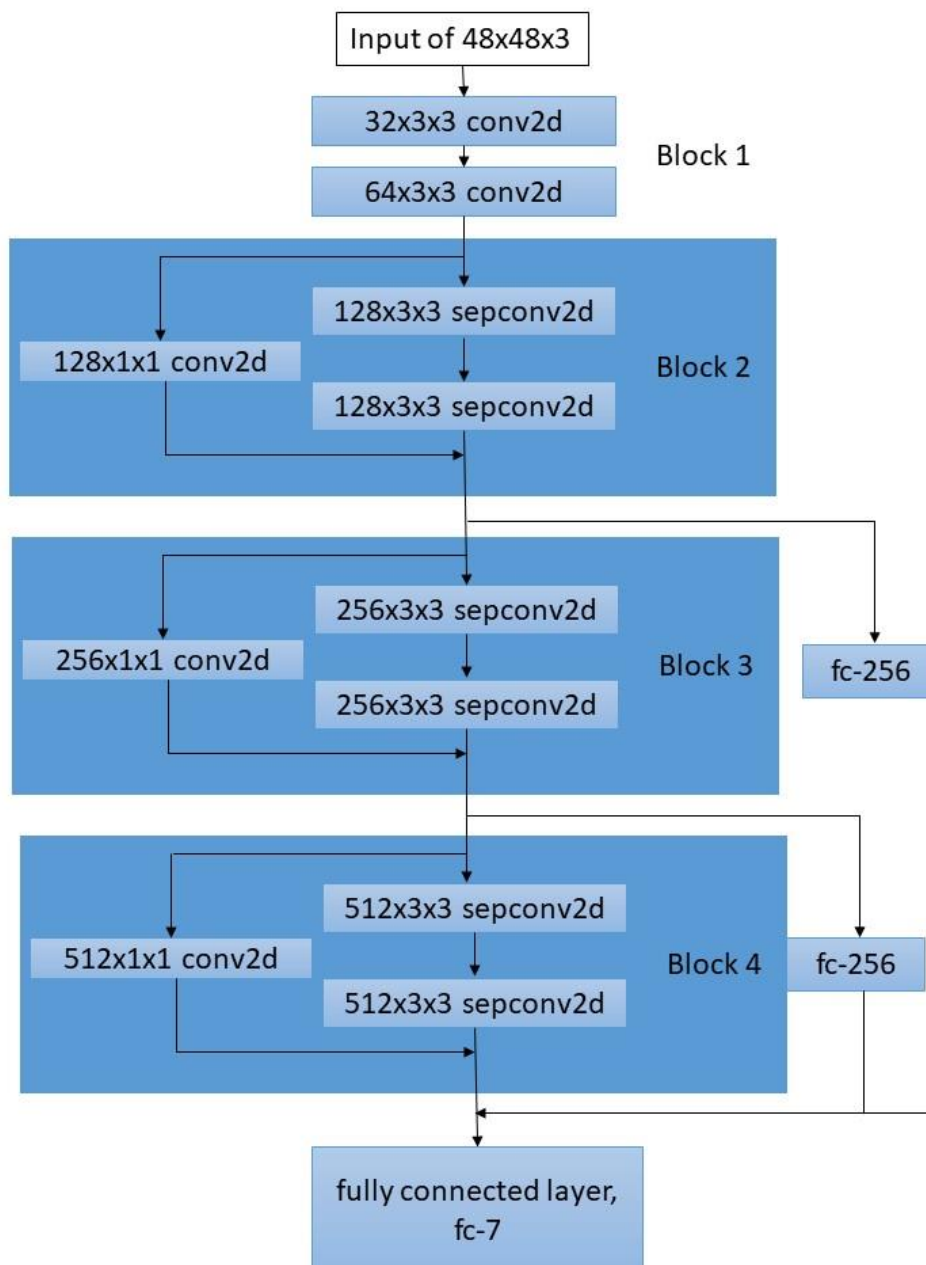


*Figure 12 Multi level mini xception architecture (MLMX)*

30

Figure 12 shows how the Mini Xception model is modified to produce a multi-level network. The feature maps from each of the blocks except the first block are put through a fully connected layer of 256 units. This is done to reserve significant filters and remove the irrelevant ones. The features identified after block 1 are not considered since those are low-level features which may further confuse the model. All these feature maps are concatenated and then put through a fully connected layer with 7 units.

## 4.3.2 Multi-Level VGG



*Figure 13 Gradcam visualization for VGGNet 19*

The image of the woman with a smile was given as input to the VGGNet 19 model. Figure 13 shows the input image, the grad-cam generated from block 1 and the grad-cam generated from block 3. The grad-cam from block 1 shows that the model was focusing on irrelevant features like the background and the hair, which are the low-level features. This justifies not using features from block 1 in the multi-level architecture. The grad-cam from block 3 shows that the model was focusing on the facial muscles near the smile. This reinstates the fact that by concatenating the medium-level and high-level features, the model can classify the image more accurately.

The grad-cam representations represent all the feature maps produced because of a block of operations in the network. If the number of output channels in a convolution layer is x, then the number of feature maps produced is also x. Every feature map may not contain relevant information, and so in our architecture, the feature maps are put through a fully connected layer with 256 units. This ensures that only 256 feature maps containing relevant information are retained. So, the feature maps obtained after every block are more informative and relevant than before.

Figure 14 shows the multi-level architecture of VGG19. Features from every block except the first block are taken into consideration. This is because the feature maps after the first block usually focus on areas in the background which are not useful in facial expression recognition as illustrated in the grad-cam representations. Like the multi-level tiny xception model, the feature maps from every other block are put through a fully connected layer with 256 units to reserve the important

filters as explained above. Then these feature maps are concatenated before being fed into another fully connected layer with 7 units to predict the probability for each of the seven emotions.
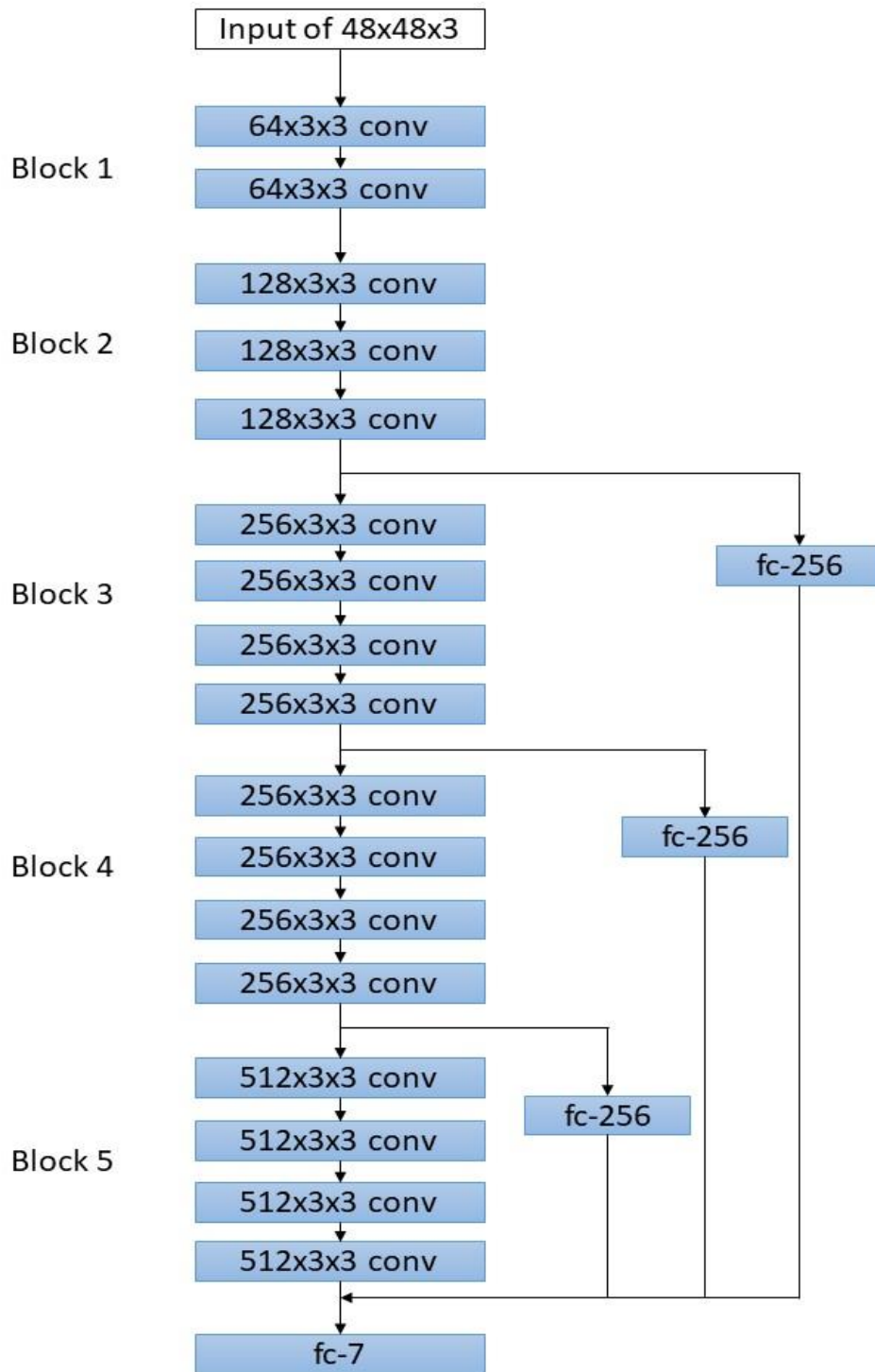


*Figure 14 Multi level VGGNet 19 architecture (MLVGG)*

## 4.4 Implementation Details

All the models were implemented in PyTorch. The mini xception model used batch size of 128 and 250 epochs to obtain the best performance. The VGG19 model used batch size of 96 and 100 epochs to do the same.

To train the multi-level tiny xception model, the FER dataset was first preprocessed as described in chapter 3. It used a batch size of 96 and 250 epochs. We found better results using the conventional model weights to initialize the multi-level xception model and then to further train the model, which is called pre-training. A similar approach was adopted for the VGG19 model as well. The model without transfer learning used batch size of 128 and 100 epochs, while the one with transfer learning used the same batch size and 150 epochs.

The learning rate was initialized to 0.001 and after reaching a threshold epoch, the rate would start to decrease by a factor of 0.9. This was done to ensure the model does not skip the convergence due to a high learning rate. The hyperparameters for each model were finalized after a lot of experimentation. For the base models, hyperparameter tuning methods like grid search with cross validation were used as well.

The network is expected to classify the image into one of seven expressions and thus categorical cross entropy metric is used to measure the loss. The logarithm of the predicted softmax value for every class is multiplied with the ground truth label (0/1) and all of them are added to obtain the total loss value for an image. Stochastic Gradient Descent (SGD) algorithm is used to help the model reach convergence iteratively. For every image in the training set, the categorical cross entropy loss is calculated and using stochastic gradient descent backpropagation is done. All the filter weights are updated after a few samples so that the loss is minimized.

The multi-level CNN models were trained in Keras as well for compatibility with the Image Captioning code.


## 4.5 Language Model

The main idea of machine learning in NLP is to train a model to understand a text corpus well enough so that it can automatically perform tasks such as text classification or text generation. In order to do this, a machine learning model must be able to quantify what counts as a "good" or "bad" sequence of words. Like every machine learning task, the task of quantifying text quality is centered around calculating probabilities. In this case, to find the words describing a person, we would need to find the relevant keywords. We can perform this task using a language model.

A language model can tell us the likelihood of each word in a given sentence or text passage based on the words that came before it. We can then determine how likely a sentence or text passage is by aggregating its individual word probabilities.

### 4.5.1 Language model tasks

Language models are useful for both text classification and generation. In text classification, we can use the language model's probability calculations to separate texts into different categories. For example, if we trained a language model on spam email subject titles, the model would likely give the subject "CLICK HERE FOR MONEY" a relatively high probability of being spam. In our case, we classified the person entities from the captions of the images which was later fed to the Decoder network.

In text generation, a language model completes a sentence by generating text based on the incomplete input sentence. This is the idea behind the autocomplete feature when texting on a phone or typing in a search engine. The model will give suggestions to complete the sentence based on the words it predicts with the highest probabilities.

### 4.5.2 Word probabilities

The purpose of a language model is to assign probabilities to words in sequences of text. The probability for each word is conditioned on the words that appear before it in the sequence. Though it may not seem like it at first glance, the task of calculating a word probability based on the previous words in a sequence is essentially multi-class classification.

Since each word in a sequence must come from the text corpus' vocabulary, we consider each vocabulary word as a separate class. We then use the previous sequence words as input into our language model to calculate probabilities for each vocabulary class.

## he went to the ____

| Vocabulary | Probability |
| --- | --- |
| he | 0.02 |
| went | 0.02 |
| to | 0.01 |
| the | 0.0 |
| coffee | 0.15 |
| store | 0.75 |
| wrist | 0.05 |

### 4.5.3 Inputs and Targets

The main idea behind training a language model is to predict each word in a sequence based on the words that come before it. However, rather than making training pairs where the target is a single word, we instead make training pairs consisting of equal length input and target sequences.

Original Sequence: ["she", "bought", "a", "book", "from", "me"]

Input Sequence: ["she", "bought", "a", "book", "from"]

Target Sequence: ["bought", "a", "book", "from", "me"]

The language model attempts to correctly predict each word in the target sequence based on its corresponding prefix in the input sequence. In the example above, the input prefix-target word pairs are:

- (["she"], "bought")
- (["she", "bought"], "a")
- (["she", "bought", "a"], "book")
- (["she", "bought", "a", "book"], "from")
- (["she", "bought", "a", "book", "from"], "me")

For each pair, we would ideally want the language model to calculate a very high probability for the target word based on the input prefix. This is the same goal as training a multi-class classification model.

### 4.5.4 Maximum length

When creating the input-target sequences for training a language model, a big factor to consider is the length of the sequences. Depending on the text corpus, we can choose to limit our training pair sequences to a fixed maximum length. Using a fixed max sequence length can increase the training speed and also help the language model avoid over-fitting uncommon text dependencies which are sometimes found in long, run-on sentences.

Original Sequence: ["they", "went", "to", "the", "big", "red", "truck"]

Input Sequence: ["they", "went", "to"]

Target Sequence: ["went", "to", "the"]

### 4.5 Summary

This section described various CNN architectures like Mini Xception and VGGNet 19 that were used in this work for facial expression recognition. Multi-level modifications were made to both architectures to help extract better features for classification and the implementation details were mentioned. The working of the language model was also described.

# CHAPTER 5

# PROPOSED METHODOLOGY

The following chapter describes the end-to-end method that is proposed to train the decoder network as shown in figure 15. This in turn add emotions to generated captions.
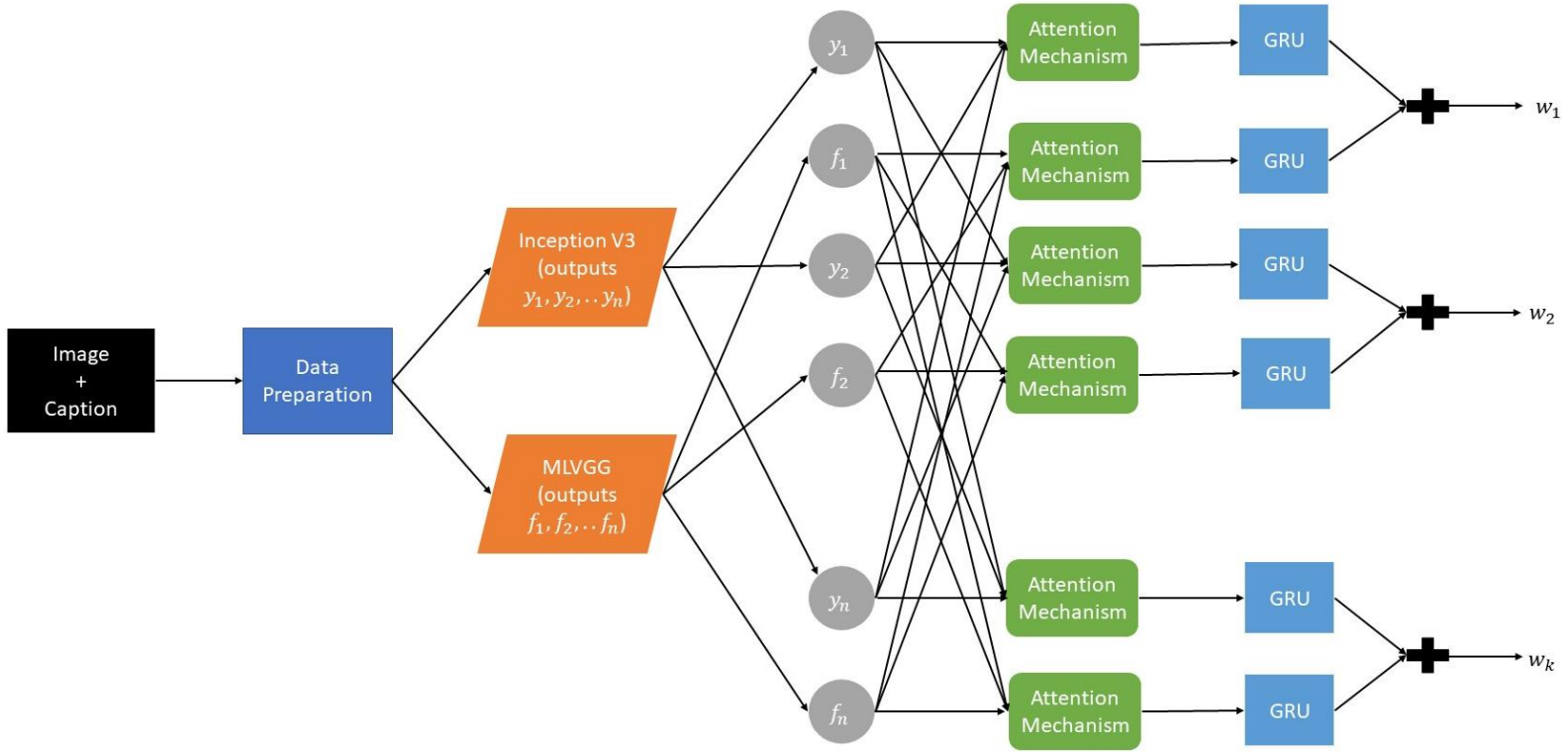


*Figure 15 End-to-end pipeline for facial expression driven image captioning*

The data in the MS COCO dataset is firstly preprocessed. Then the image is fed to the inception network and MLVGG network simultaneously, to extract image features $(y_1, y_2, .., y_n)$ and facial features $(f_1, f_2, .., f_n)$ respectively. These are fed to an attention model followed by the GRU one after the other. The outputs generated are combined to predict the word in every time step.

## 5.1 Dataset Preparation

Since facial expression recognition, which is used to add emotion to the captions, is applicable to human faces, we first apply a face detection algorithm on the training set of MS COCO and obtain 26,236 images with 112462 captions. The MTCNN network is used to do so which uses three neural networks to obtain the face landmarks. All the faces present in the image are fed to the FER network, which predicts the expressions and gives the most occurring expression as output.

The captions in the training set are processed so that they contain an adjective which describes the facial expression before the subject. A natural language processing language model is trained to identify the entities in a sentence. The ground truths of the dataset are modified to accommodate emotion.

The language model is a named entity recognition (NER) model which parses through bodies of text and identifies specified categories. These categories include countries, companies, organizations, persons, animals, locations and so on. This involves the implementation of two tasks, tokenizing the text and detecting names, and classifying these names into pre-defined categories.

For our purpose, the person entity is required to be extracted and a customized version of Spacy's NER model is used. Spacy is an open source library for NLP tasks published under an MIT license. Its second version contains neural models to solve problems like NER, text segmentation and parsing, using neural networks with residual connections to give much better latency. The NER model is customized by adding a small amount of training data. The data contains the sentence and the starting and ending position of the entity to be detected and the name of the entity. The model is trained on this and other specified categories in the language model are disabled.

This trained language model is used to find subjects, which is the person entity, in every caption in the training set of MS COCO. Since the expression is already detected, a word describing the expression is added before the subject. These words include happy, sad, scared, surprised, angry, disgusted, and neutral. In this way, a new dataset for the purpose of training a more human-like image captioning model was created. For example:

- `<start> a couple of happy people are sitting on a bench <end>`
- `<start> a sad woman in a shirt and tie sitting down <end>`
- `<start> a neutral man is standing next to a small clock <end>`
- `<start> an angry baseball player dropping his bat and starting to run <end>`
- `<start> some surprised players in action on the baseball field <end>`
- `<start> a disgusted woman is showing a pepperoni deep pizza <end>`
- `<start> a scared man riding a skateboard up the edge of the skate board park <end>`


For the task of image captioning, further data processing is required. The captions in the training set are fed to a tokenizer where every unique word is extracted and arranged in descending order of their frequency. These words form the vocabulary of the network and two lists are created. The word to index maps each word to the corresponding index and index to word does vice versa. Then every word in the captions is substituted by its corresponding index. After this all captions are padded to have the same length.

Steps involved in preprocessing the captions, given the corresponding image:

- The image is fed to MTCNN (face detector) and the largest face is cropped and resized to $48 \times 48 \times 3$
- The facial expression corresponding to the cropped face is recognized by feeding it to MLVGG.
- The subject is identified in the caption and the adjective describing the facial expression is added before it.
- The caption is tokenized, and each word is replaced by its index in the word to index list.
- The caption is padded to increase its length, so that all captions have the same length.

## 5.2 Encoder Network

An encoder is used to capture the information in the input and represent it in a certain format. In the end-to-end method, the encoder has the same functionality. It is used to extract features from an image. These can include persons, animals, objects, and locations detected.

For this purpose, we use Inception V3 to obtain the image features. This is the refined version of Inception V1, otherwise called GoogleNet. It was the first runner-up in the ILSVRC 2015 (Imagenet Large Scale Visualization Competition) and is conceived by Google. Inception V3 is a convolutional neural network and is specifically designed to solve the problem of large-scale object detection. It is trained on the Imagenet dataset that consists of around 15 million high resolution images which contain labelled objects in around 22,000 categories.

Since Inception is trained on such a large dataset and is used to classify images into one of the 22,000 categories, it must be able to differentiate between all these categories well. The ability to distinguish rises from the model's ability to identify good differentiating features. Thus, Inception V3 is capable to extract features containing lots of information from an image and these coincide with our requirements of an encoder.

The Xception and VGG architectures were described in the previous chapter and their modified versions were used to train on a different dataset for a different task. However, in this case the pre-trained weights of Inception V3 are used. Models such as ResNet, VGG, Xception, Inception have become universal and so, many deep learning platforms provide their architecture and weights in their packages.
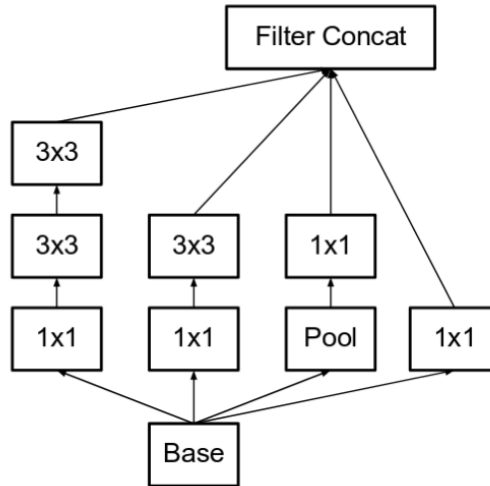
*Figure 16 Inception Module A*

Inception V3 is a 42-layered deep convolutional neural network which achieved computational efficiency, lesser number of parameters and a low error rate. It uses asymmetric convolutions to further reduce the number of parameters. The first type of module called Inception module A as shown in figure 16, uses VGG's technique of replacing a 5×5 filter by two 3×3 ones.



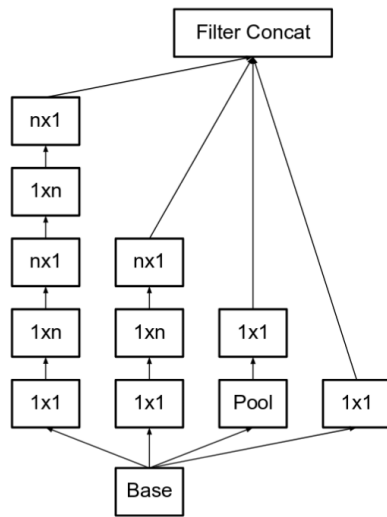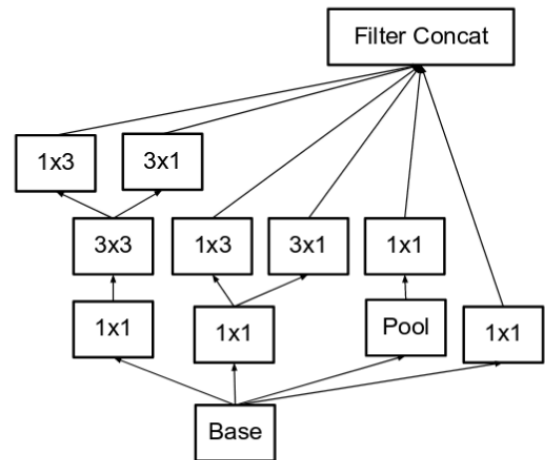*Figure 17 Inception Module B*



*Figure 18 Inception Module C*

In Inception module B, a 3×3 filter is replaced by a 1×3 filter and a 3×1 filter. This reduces the number of parameters from 9 to 6, that is by 33%, as shown in figure 17. Inception module C (figure 18) is used for high dimensional representations. As the number of parameters are reduced the model can go deeper without the fear of overfitting.
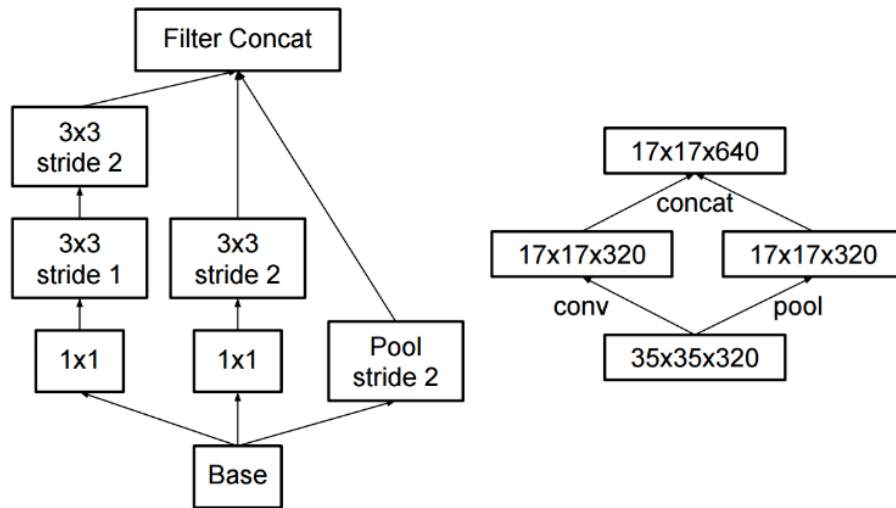
*Figure 19 Grid Size Reduction*

Another important feature is the grid size reduction. The downsizing of feature maps is usually done by max pooling in VGG and AlexNet architectures. If a max pool layer is added before the convolutional layer, it becomes too greedy and some important information might get lost. On the other hand, if the max pool layer is added after the convolutional layer, it becomes computationally expensive. To overcome this issue, Inception V3 implements grid size reduction as shown in figure 19. Here the feature maps are fed parallelly to a convolution layer with stride 2 and a max pool layer. Then both outputs are concatenated, this solves the previous two shortcomings of the position of the max pool layer.



*Figure 20 Inception V3 Architecture*

Figure 20 depicts the architecture of Inception V3. An input image of size 229×229×3 is fed to the network and is first put through five convolution layers. Then, the first three blocks are module A type, followed by a grid size reduction. Next four blocks are of module B type again followed by grid size reduction. The last two blocks are of module C type to handle high dimensionality feature

40

maps. After all these operations we obtain feature maps of size 8×8×2048. Then a fully connected layer is applied, followed by softmax to predict the probabilities of the Imagenet classes.

For using Inception V3 in the encoder network, features are extracted from the second last layer, which is before the fully connected layer. The feature maps are of size 8×8×2048. The encoder network contains another fully connected layer whose output number of channels is the same as the embedding size used in the decoder. In this case we use 256. Therefore, the dimensions of the encoded features are x×64×256, where x is the batch size. Each feature map contains features corresponding to an object identified in the image.

## 5.3 FER Network

The best model from the previous chapter is used for extracting facial features. The multi-level VGG model is used and features are extracted after the concatenation layer before feeding to the fully connected layer. The features are rich in their content since the features extracted from the medium level blocks are concatenated. The dimensions of the facial features extracted are x×9×256, where x is the batch size. There is no requirement to modify the dimensionality of these features since they coincide with the embedding dimension of the decoder. This happens since the features extracted from the medium level blocks are fed into a fully connected layer with 256 units to retain the important filters, as shown in figure 14.

## 5.4 Attention Mechanism

The attention mechanism is used to ensure the decoder understands which part of the feature map it should focus on while generating a word. In the absence of this mechanism, the decoder is always fed the same feature maps and the hidden vector from the previous iteration, which correspond to the image as a whole and the previous state of the decoder. This leads confusion within the decoder, while decoding the features. Attention model helps decide which part of the encoded vector corresponds to each word in the output while training and how to use relevant information to select the output word while testing. These two processes are called align and translate.
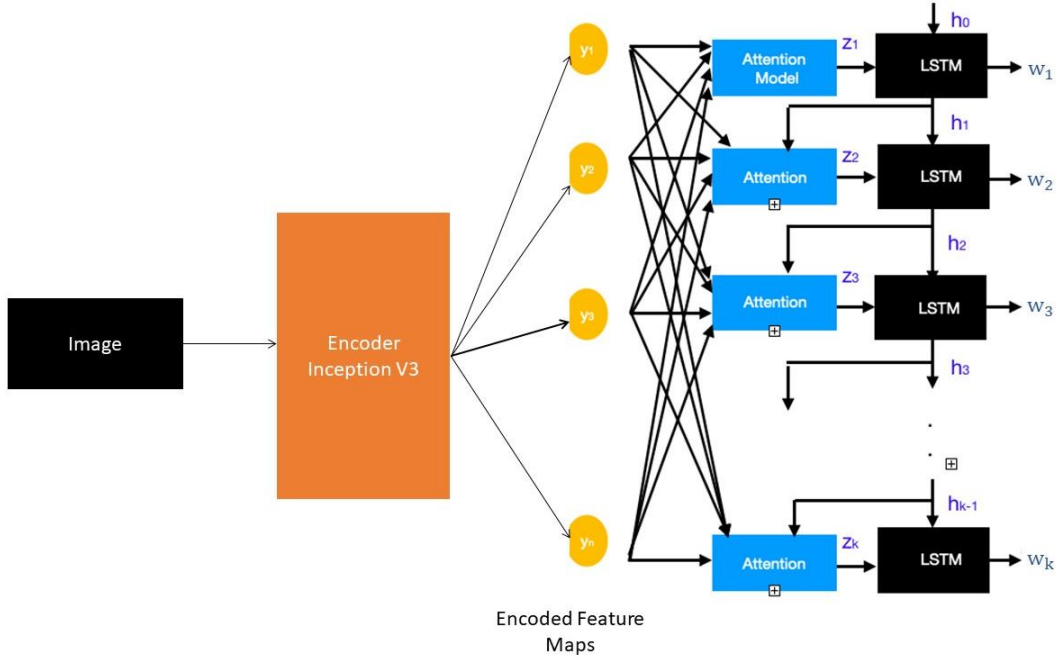
*Figure 21 Attention Mechanism*

As shown in figure 21, the attention model takes the encoded feature maps from the Inception V3 $(y_1, y_2, ..., y_n)$ where the number of feature maps, $n$, corresponds to the dimension of the embedding which is 256 in this case. The feature maps are put through a CNN denoted as attention model in figure 21, to focus on the relevant parts of the encoded feature map and produce a context vector for each output time step (every word generated). In an encoder-decoder network without attention a constant fixed length context vector was used to generate every word in the caption. However, due to the attention model the context vector varies in every time step to reflect relevant information corresponding to each word and the previous hidden state of the decoder. The decoder uses the attention model's output called the context vector $z_t$, the hidden state $h_t$ and the previous word $w_{(t-1)}$ to predict a word $w_t$ in the time step $t$.

A mechanism $\phi$ is used to compute the context vector $z_i$ from the encoded vectors $(y_1, y_2, ...., y_n)$ which correspond to features extracted from different locations of the image. For every location $i$ it computes an attention weight $\alpha_i$, which denotes the importance of that location. The attention weight $\alpha_i$ corresponding to each vector $y_i$ is calculated using an attention mechanism denoted by $f_{(att)}$. This is done by using the hidden vector from the previous state $h_{(t-1)}$.

$$e_{(ti)} = f_{(att)}(y_i, h_{(t-1)})$$

$$\alpha_{(ti)} = \frac{exp(e_{(ti)})}{\sum_{k=1}^{n} exp(e_{(tk)})}$$

After computing the attention weights, the context vector is computed as

$$z_t = \phi(y_i, \alpha_i)$$

Given the set of encoded vectors and their attention weights, $\phi$ produces a single vector $z_t$. There are two ways to compute the attention weights. The following sub sections describe them.

## 5.4.1 Stochastic Hard Attention

Let $s_{(t,i)}$ denote the location that the model chooses to look at while generating the word $w_t$ in the time step $t$. The variable $s_{(t,i)}$ is a one-hot variable which is set to 1 if the model considers $y_i$ to extract the image features. Hard attention uses a multinouilli distribution parameterized by $\alpha_i$ to sample whether location $y_i$ is selected which is reflected in the value of $s_{(t,i)}$. Here, $\alpha_i$ functions at the sample rate.

$$p\big(s_{(t,i)} = 1 \big| s_{(j<t)}, y\big) = \alpha_{(t,i)}$$

$$z_t = \sum s_{(t,i)} * y_i$$

The Monte Carlo method is used in backpropagation to ensure the gradient descent occurs correctly. It performs end-to-end episodes to calculate an average for all sampling results. A Monte Carlo based sampling approximation of the gradient with respect to the model parameters is used. The accuracy depends on number of samplings performed and the quality of these samplings.

Using this method, every location in the encoded vectors $y_i$ is allotted a weight $s_{(t,i)}$ in the time step $t$ where the value of $s_{(t,i)}$ is sampled from a multinouilli distribution parameterized by $\alpha_i$.

## 5.4.2 Deterministic Soft Attention

This attention model was introduced by Bahdanau et al. wherein the attention weights $\alpha_i$ corresponding to each location $y_i$ are computed in a deterministic method unlike the sampling method adopted in the previous section. This is done by computing a soft attention weighted vector,

$$\phi(\{y_i\}, \{\alpha_i\}) = \sum_{i=1}^{n} \alpha_i * y_i$$

Since the whole model is differentiable and smooth under deterministic attention, end-to-end training can proceed as normal by using standard backpropagation. As soft attention is being computed the variable $s_{(t,i)}$ is not one-hot anymore. Instead it represents the attention score corresponding to a location $y_i$ in the time step $t$.

$$s_{(t,i)} = tanh\left(W_c h_{(t-1)} + W_y y_i\right)$$

Then $s_{(t,i)}$ is put through softmax normalization to compute the attention weight $\alpha_i$.

$$\alpha_i = softmax\left(s_{(t,1)}, s_{(t,2)}, \ldots, s_{(t,n)}\right)$$

Due to the softmax normalization, the values of $\alpha_i$ add up to a value of one. In this way the function $\phi$ is calculated and the context vector $z_t$ is produced for every time step $t$.



*Figure 22 Caption generation using hard and soft attention*

Soft attention is found to perform better than hard attention as per various experiments conducted. The result seems obvious since hard attention relies on sampling to decide whether a location should contribute to the context vector, unlike the deterministic method of soft attention.

Figure 22 further illustrates the benefit of using the soft attention mechanism over the hard attention mechanism. It is seen that soft attention (the bottom row) is able to focus on the appropriate location in the image and the attention is dispersed in nature due to the deterministic aspect. Whereas in hard attention (the top row), a circular shape of same size is seen due to the stochastic nature of the algorithm. In the proposed methodology, soft attention mechanism is used.

Sections 5.4.1 and 5.4.2 described how attention weights corresponding to features obtained from the encoder are calculated. To incorporate the facial expression features while generating the words, the pipeline is modified to accommodate these features. First, the image features are fed to the attention model to obtain the attention weights and the context vector. Then the facial features are also fed to the attention model. A separate context vector is created that emphasizes which locations of the face contribute to a word. Ideally, the context vector corresponding to the face must contain information only when a word describing the face is found in the ground truth tokenized captions.
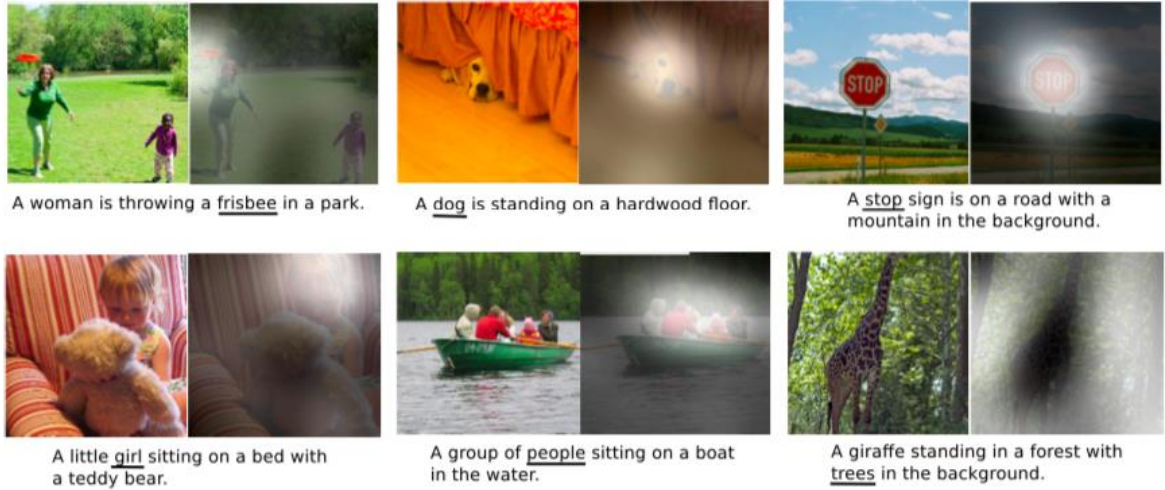


*Figure 23 Impact of attention on generated captions*

Figure 23 clearly shows how adding the attention model is useful in generating better captions. Objects such as the frisbee, dog, stop sign etc are being clearly detected. This happens since at a time step, the context vector contains features describing each of these objects fed to the decoder, which in turn generates the most relevant word.

## 5.5 Decoder Network

The context vector $z_t$, along with the previous hidden state $h_{(t-1)}$ and the previous word generated $w_{(t-1)}$ is fed to the decoder network to generate the word $w_t$ in the time step $t$. Similarly, each word of the caption is generated word by word and as explained in the previous section on attention, the context vector is dynamic in nature.

Recurrent neural networks (RNNs) are used to decode the context vector. They have the property to remember the information learnt from previous inputs and use this information to generate the next output. The context extracted from prior inputs and outputs is represented by a hidden vector. Thus, the same input can generate a different output based on the previous hidden state.
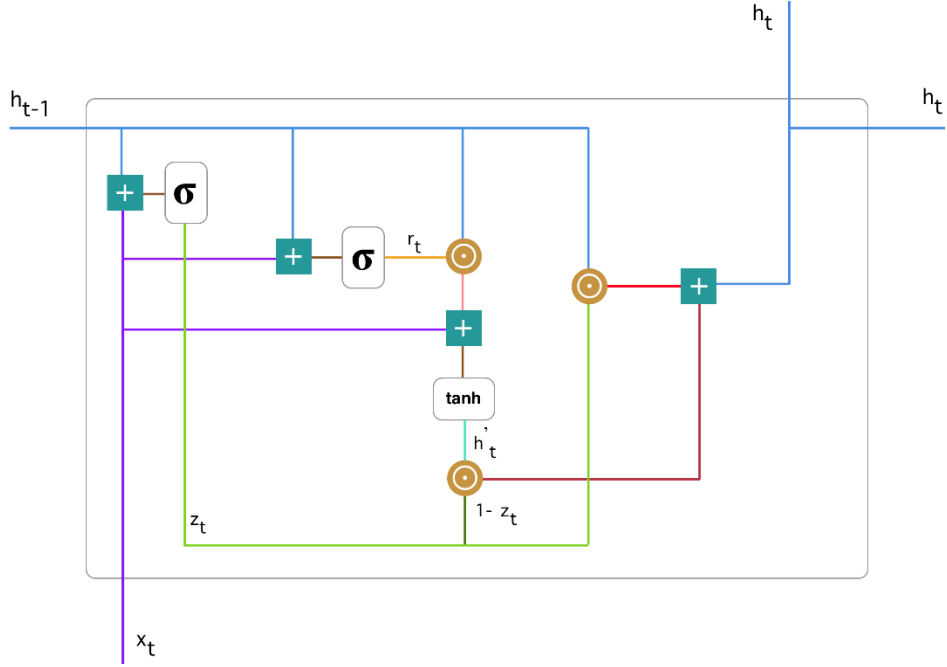
*Figure 24 Working of GRU*

A gated recurrent unit (GRU) as shown in figure 24, is an improved version of a standard RNN, which solves the issue of the vanishing gradient problem. It uses the reset gate and the update gate to do so. The update gates help the network decide the information collected from the previous time steps should be passed on to the output. The reset gate is used to decide how much of the past information should be forgotten.

The following equations represent the mathematical working of the GRU. The update gate $z_t$ in the time step $t$ is calculated using the following formula,

$$z_t = \sigma\left(W^z x_t + U^z h_{(t-1)}\right)$$

$x_t$ represents the input being fed into the network and $h_{(t-1)}$ denotes the previous hidden state. $W^z$ and $U^z$ denote the weights corresponding to the input vector and the hidden state in the update gate which are learnt during the training phase. The sigmoid function produces a value between 0 and 1. This helps the network retain information from the previous hidden state and the current input. It can decide to retain all the information from the past and this combats the vanishing gradient problem where an RNN tends to forget information not from the recent past.

The reset gate $r_t$ in the time step $t$ is calculated as,

$$r_t = \sigma\left(W^r x_t + U^r h_{(t-1)}\right)$$

with $W^r$ and $U^r$ representing the weights corresponding to the input vector and the previous hidden state in the reset gate. The sigmoid function produces the same effect as before. However, this gate is used to determine how much of the past information is to be forgotten. The difference in functionality between the update and reset gate is achieved by their respective weights, in spite of having the same formula.

46

The current memory content is denoted by $h'_t$

$$h'_t = tanh(Wx_t + r_t \odot Uh_{(t-1)})$$

An element-wise product operation is done between the reset gate and the previous hidden state. This determines what part of the hidden state is retained and what is forgotten while computing the current memory content. As usual, $W$ and $U$ denote the weights. The non-linear activation function, $tanh$ is used to compute $h'_t$.

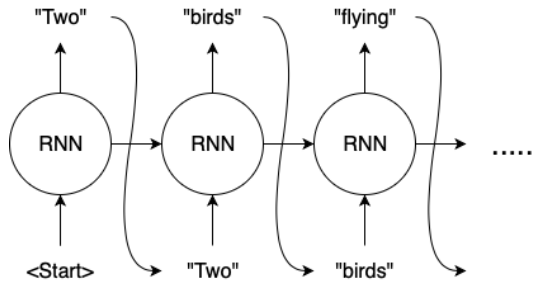The current hidden state $h_t$ is computed as follows,

$$h_t = z_t \odot h_{(t-1)} + (1 - z_t) \odot h'_t$$

It is the sum of the element-wise product of the update gate and the previous hidden state, and the element-wise product of the additive inverse of the update gate and the current memory content. Since $h_t$ denotes the current hidden state and it must express some amount of information from the previous hidden state and some amount of information from the current memory cell. $z_t$ denotes how much of the previous hidden state should be retained and so is multiplied with the same. $(1 - z_t)$ will denote how much of the current memory cell should be retained. And so, the current hidden state is calculated for the time step $t$. These equations describe how a gated recurrent unit can generate the next output based on the input and the previous information that it has collected.
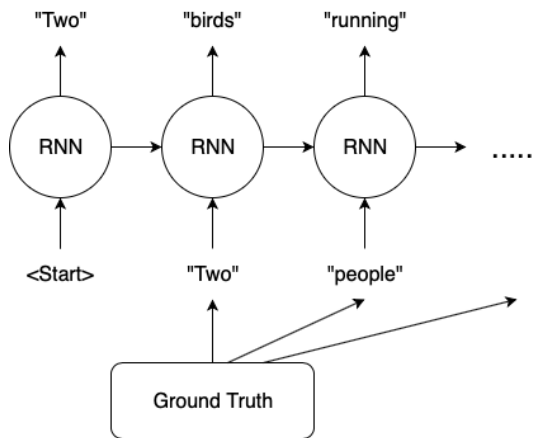
Teacher forcing is a technique used during the training phase of RNNs to increase the computational speed and so that the model reaches convergence quickly. In this technique, the target word is passed as the next input to the decoder.

Without teacher forcing, the RNN learns using what it predicted in the previous time step. In figure 25, the RNN predicts the word "birds" in the second time step while the ground truth is "people". In this case, the model learns weights corresponding to the wrong prediction. With teacher forcing, the ground truth caption is fed to the RNN at every time step. This helps the decoder reach convergence sooner because it is trained on the correct input word.

The decoder input is initialized to the start token which is all zeroes and the hidden vector is also initialized similarly. The features are obtained from the encoder network and these three inputs are fed to the decoder network. The attention model is used to obtain the context vector and the attention weights corresponding to the features. The context vector contains the important features located in the input image and the decoder input contains the index of the previous word in the ground truth, which enforces teacher forcing. Both these vectors are concatenated to produce the input to be fed to the GRU. The GRU returns the output and the hidden state. The output vector is put through two fully connected layers to obtain a vector with the same number of rows as the number of words in the vocabulary. Each entry in the vector denotes the probability with which the corresponding word can be generated by the decoder network.

*Figure 25 Example of training an RNN with and without Teacher Forcing*

The above process is repeated for every word in the ground truth caption, the loss is calculated, and the weights are updated using backpropagation. In this way the weights in the attention model and the GRU are learnt.
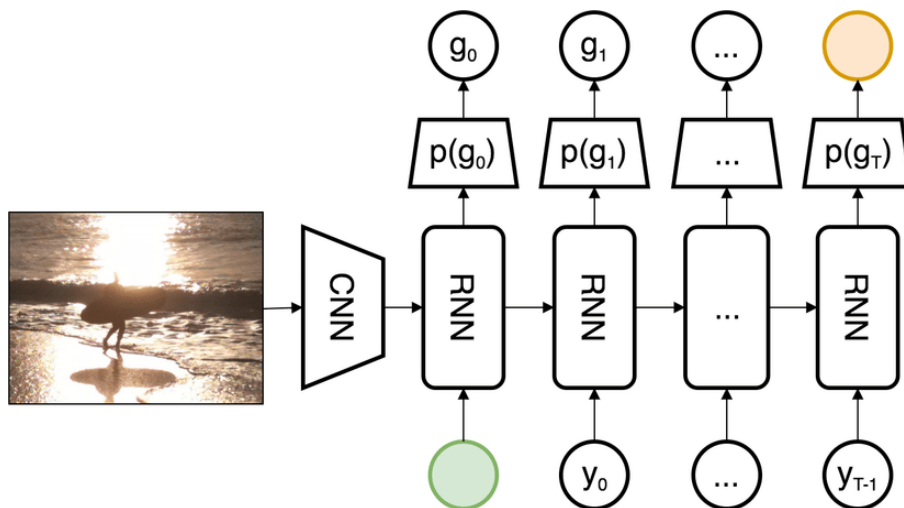


*Figure 26 Working of decoder network*

For every image in the training set, it is put through the encoder network and the FER network to obtain image features and facial features, respectively. For every word in the ground truth caption, both the features are fed into the decoder one after the other. The decoder returns the hidden state, the predictions, and the attention weights for each set of features. To add emotion to the captions, the predictions generated are added giving the facial features and image features equal importance. The loss is calculated, the decoder input is initialized accordingly, and backpropagation occurs. This process is repeated for every image in the training set and the training weights are learnt. In the caption generation mode, the caption is predicted one word at a time as shown in figure 26. The decoder network returns the probabilities corresponding to each word and the word with the highest probability is the output.

## 5.6 Implementation Details

The entire pipeline was implemented in TensorFlow. The pre-trained model of Inception V3 was used in the encoder network. The best performing model from the experiments on FER2013 was used in the FER network.

To train the decoder network, sparse categorical cross entropy was used to calculate loss. This loss is the best suited in this scenario due to the presence of multiple labels which makes it hard to represent them as one hot encoded vector. The Adam optimizer, which is a variant of stochastic gradient descent is used. It does not require much memory unlike the gradient descent which requires an entire epoch to take place before updating the weights of the model. Furthermore, Adam is well suited to handle training of models with a large number of parameters.

The embedding dimension, the number of words in the vocabulary, the batch size and the number of epochs are a few of the hyperparameters. In this case we use 256, 1000, 64 and 15, respectively. It was found that the model began to overfit beyond the fifteenth epoch.

In the training phase, every image and caption is processed. Then the image is resized and fed to Inception V3. Next, the same image is fed to MTCNN and the largest face is detected. The face is cropped out and resized, which is then fed to MLVGG. Once the image features and facial features are obtained, they are fed to the Bahdanau attention network one by one and then to the GRU. Corresponding to each image, two predictions are obtained from the GRU, these are added giving equal weightage to both. Then the loss is calculated and backpropagated.

## 5.7 Summary

This section described the end-to-end pipeline of the proposed methodology in a sequential manner. Every block present in figure 12, such as data preparation, Inception V3, MLVGG, attention mechanism and GRU, is discussed at length and architecture decisions are explained. The implementation is explained as well.

# CHAPTER 6

# RESULTS AND ANALYSIS

## 6.1 Performance on the FER Problem

Table 3 shows the accuracy obtained by all the described models on the test set of FER2013.

| Model | Architecture | Accuracy (in %) |
|---|---|---|
| Mini Xception | 12 layers | 68.15 |
| MLMX | multi-level mini xception | 69.38 |
| VGG | 18 layers | 70.30 |
| MLVGG 1 | multi-level VGG | 71.27 |
| MLVGG 2 | pre-training on multi-level VGG | 73.11 |

*Table 3 Performance of models*

The conventional mini xception model with an accuracy of 68.153% performs worse than the conventional VGG 19 model. By introducing multi-level architecture, that is using mid-level features in both mini-xception and VGG models, we notice a one-point increase in performance. Since VGG performs much better than mini xception, pre-training was implemented in this case only. A clear improvement to 73.11% is seen in the case of MLVGG 2. It is inferred that VGG architecture is more well suited to perform facial expression recognition than the mini-xception architecture. This is because of an increased number of layers as compared to mini xception, which helps the model extract better features. The best performance is achieved by using mid-level features via multi-level CNNs and pre-training the model. Table 4 shows some of the state-of-the-art architectures and their performance. It is seen that our model's performance is on par with the hybrid CNN-SIFT aggregator. It consistently outperforms all the previous models by a significant factor. Therefore, MLVGG 2 is used in the image captioning pipeline to help improve the results.

| Method | Accuracy (in %) |
|---|---|
| Inception | 66.40 |
| FER2013 Winner | 71.20 |
| Multi-scales CNNs | 71.80 |
| MNL | 72.08 |
| Hierarchical committee of CNNs | 72.72 |
| Multi-scale CNN | 72.82 |
| **MLVGG 2** | **73.11** |
| Hybrid CNN-SIFT aggregator | 73.40 |

*Table 4 Performance of state-of-the-art architectures on FER2013*

## 6.2 Extracting Facial Expressions



*Figure 27 Predictions by MLVGG*

Figure 27 shows some predictions by the multi-level VGG model, in each category of expression. The predictions are generated on a few images from the FER2013 test set. The model has successfully predicted the correct facial expression. The distinguishing facial features for every expression can be recognized. For example, happy includes a smile, with or without teeth. Anger would contain scrunched-up forehead. Disgust entails a scrunched-up face. Sad involves a frown and small eyes and so on. However, in case of some expressions, when the human being is not loud about his or her expression, a factor of subtlety is introduced. This makes distinguishing between sad and angry, and fear and surprise, a bit confusing for the model.



*Figure 28 Confusion in some predictions*

Consider the images in figure 28, the model recognizes an angry boy as sad and a scared boy as surprised. In the first example, the expression is very subtle, no frown is clearly visible. In the second example, the boy is scared. But the mouth open in a 'o' shape is regarded as a distinguishing feature for both these expressions.

The above reasons cause confusion in solving the FER problem and may be a plausible cause for accuracy not exceeding 73.11%. However, the model did perform satisfactorily well, and we have been able to increase the accuracy from the standard CNN configuration by 2.81%.


## 6.3 Generating Captions

Table 5 shows the BLEU-1,2,3,4/METEOR metrics of the Image Captioning model compared to other methods.

| Dataset | Model | BLEU | | | | METEOR |
| | | BLEU-1 | BLEU-2 | BLEU-3 | BLEU-4 | |
|---|---|---|---|---|---|---|
| Flickr8k | Google NIC(Vinyals et al., 2014)[†Σ] | 63 | 41 | 27 | — | — |
| | Log Bilinear (Kiros et al., 2014a)[○] | 65.6 | 42.4 | 27.7 | 17.7 | 17.31 |
| | Soft-Attention | **67** | 44.8 | 29.9 | 19.5 | 18.93 |
| | Hard-Attention | **67** | **45.7** | **31.4** | **21.3** | **20.30** |
| Flickr30k | Google NIC[†○Σ] | 66.3 | 42.3 | 27.7 | 18.3 | — |
| | Log Bilinear | 60.0 | 38 | 25.4 | 17.1 | 16.88 |
| | Soft-Attention | 66.7 | 43.4 | 28.8 | 19.1 | **18.49** |
| | Hard-Attention | **66.9** | **43.9** | **29.6** | **19.9** | 18.46 |
| COCO | CMU/MS Research (Chen & Zitnick, 2014)[a] | — | — | — | — | 20.41 |
| | MS Research (Fang et al., 2014)[†a] | — | — | — | — | 20.71 |
| | BRNN (Karpathy & Li, 2014)[○] | 64.2 | 45.1 | 30.4 | 20.3 | — |
| | Google NIC[†○Σ] | 66.6 | 46.1 | 32.9 | 24.6 | — |
| | Log Bilinear[○] | 70.8 | 48.9 | 34.4 | 24.3 | 20.03 |
| | Soft-Attention | 70.7 | 49.2 | 34.4 | 24.3 | **23.90** |
| | Hard-Attention | **71.8** | **50.4** | **35.7** | **25.0** | 23.04 |

*Table 5 Performance of the state-of-the-art architectures of Image Captioning models*

Figure 29 shows some of the captions generated by our model. It can be seen that the structure of the sentence is altered to generate an adjective describing the facial expression of the subject. These adjectives help us understand that our objective of adding emotion to the caption has been fulfilled. In three of the examples, it is seen that the emotion generated by the decoder network is different from the ground truth caption. However, in this case the predicted caption seems to capture the essence of the image better than the ground truth. In the first example, the boy's expression is more neutral than happy. It can be inferred that the decoder is able to extract the meaning of the facial features appropriately. Moreover, as explained in the previous section, some of the expressions are easily confused and so they seem to explain the image well.

Real Caption: a neutral person is riding on a
wake board while the waves roll in.
Prediction Caption: young happy boy riding a
surfboard down a surf



Real Caption: a sad snow skier is sitting on
the white snow
Prediction Caption: a neutral man sitting in
the snow with trees in the snow



Real Caption: a happy chef in a kitchen
preparing a meal
Prediction Caption: a neutral man cooking some
food in the kitchen



Real Caption: a happy person standing on skis
in the snow
Prediction Caption: a happy man standing on a
snowy slope in the snow

*Figure 29 Examples of predicted captions*

```
Real Caption: two angry people are riding an
elephant at a wildlife <unk>
Prediction Caption: a happy person riding on
the back of an elephant
```

*Figure 30 An odd case*

Figure 30 presents an interesting scenario. Due to an error in the FER model prediction the ground truth caption contains an incorrect emotion adjective. Inspite of this mistake, the decoder network is able to overcome this and predict a suitable emotion ajective. Furthermore the ground truth caption contains an unknown token. In this case the predicted caption is of better quality than the ground truth caption.

On the other hand, there are some lapses in the predictions made by the image captioning model. Figure 31 shows these mistakes. The first image clearly shows a woman but the model confuses the subject to be a man. In the second image, the predicted caption is not completely incorrect but

it focuses on the clothing of the boy rather than the action. Moreover, happy does not describe the facial expression of the boy accurately.



```
Real Caption: a happy woman sitting on a couch
using a laptop computer
Prediction Caption: a happy man sitting in a
room using a laptop in a chair
```

```
Real Caption: a neutral person catching a
soccer ball in front of a row of parked cars
Prediction Caption: a happy man wearing green
and black soccer uniform
```

*Figure 31 Incorrect predictions*

Overall, the image captioning model is able to change the format of the sentence to include an emotion adjective. But it also ensures that a face must be present in the input image to do so. There is some confusion regarding certain objects while decoding and a few ground truth captions have been incorrectly modified. Despite these setbacks, the model performs satisfactorily, and it is able to learn something novel.

## 6.4 Summary

This section compares the accuracy of our FER model (73.11%) and Image Captioning model to the current state-of-the-art architectures. Then examples are given for the FER model and the Image Captioning model. The results generated by both models are discussed thoroughly. Confusions and mistakes in predictions are also addressed and analyzed.

# CHAPTER 7

# CONCLUSION AND FUTURE WORK

The multi-level architecture approach along with improvements in data preprocessing and pre-training the model has helped us achieve good accuracy in comparison to other state-of-the-art methods. The accuracy can be further improved by using an ensemble of FER models. The task of predicting facial expressions for a single frame was addressed in this report. We can extend support to videos by integrating per-frame results using graphical models.

It can also be concluded that the proposed method to incorporate facial expressions while generating image captions has succeeded in doing so. As observed in the previous chapter, the semantic structure of the caption has been altered to accommodate emotion. The end-to-end neural network system presented works well for images having one human face at least. It is clear from the experiments conducted that by increasing the size of the training dataset, the performance will also improve.

To understand the effectiveness of the caption generated in a statistical format, a new similarity measuring technique must be used on the ground truth caption and the predicted caption since even though the predicted caption does not have many word in common with the real caption, it is still apt to describe the image. Further support can be extended to detecting multiple faces and incorporating their expressions into the caption generated. These improvements will help us achieve more human-like captions and this can be extended to video summarization as well.

# REFERENCES

1. Chen, X., & Lawrence Zitnick, C. (2015). Mind's eye: A recurrent visual representation for image caption generation. In CVPR, pp. 2422–2431. IEEE.

2. Farhadi, A., Hejrati, M., Sadeghi, M. A., Young, P., Rashtchian, C., Hockenmaier, J., & Forsyth, D. (2010). Every picture tells a story: Generating sentences from images. In ECCV, pp. 15–29. Springer.

3. Fasel, B., & Luettin, J. (2003). Automatic facial expression analysis: a survey. Pattern recognition, 36(1), 259–275.

4. Ferguson, Max & ak, Ronay & Lee, Yung-Tsun & Law, Kincho. (2017). Automatic localization of casting defects with convolutional neural networks.

5. Field, T. M., Woodson, R., Greenberg, R., & Cohen, D. (1982). Discrimination and imitation of facial expression by neonates. Science, 218(4568), 179–181.

6. Goodfellow, I. J., Erhan, D., Carrier, P. L., Courville, A., Mirza, M., Hamner, B., Cukierski, W., Tang, Y., Thaler, D., & Lee, D.-H. (2013). Challenges in representation learning: A report on three machine learning contests. In ICONIP, pp. 117–124. Springer

7. Hodosh, M., Young, P., & Hockenmaier, J. (2013). Framing image description as a ranking task: Data, models and evaluation metrics. Journal of Artificial Intelligence Research, 47, 853–899.

8. Kahou, S. E., Bouthillier, X., Lamblin, P., Gulcehre, C., Michalski, V., Konda, K., Jean, S., Froumenty, P., Dauphin, Y., & Boulanger-Lewandowski, N. (2016). Emonets: Multimodal deep learning approaches for emotion recognition in video. Journal on Multimodal User Interfaces, 10(2), 99–111.

9. Kahou, S. E., Pal, C., Bouthillier, X., Froumenty, P., Gulcehre, C., Memisevic, R., Vincent, P., Courville, A., Bengio, Y., & Ferrari, R. C. (2013). Combining modality specific deep neural networks for emotion recognition in video. In ICMI, pp. 543–550. ACM.

10. King, D.E.: Dlib-ml: A machine learning toolkit. Journal of Machine Learning Research 10(Jul), 1755–1758 (2009)

11. Kiros, R., Salakhutdinov, R., & Zemel, R. S. (2014). Unifying visual-semantic embeddings with multimodal neural language models. arXiv preprint arXiv:1411.2539.

12. Lin, C.Y.: Rouge: A package for automatic evaluation of summaries. Text Summarization Branches Out (2004)

13. Pramerdorfer, C., & Kampel, M. (2016). Facial expression recognition using convolutional neural networks: State of the art. arXiv preprint arXiv:1612.02903.

14. Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. In NIPS, pp. 3104–3112.

15. Tang, Y.: Deep learning using linear support vector machines. arXiv preprint arXiv:1306.0239 (2013)

16. Tian, Y.-I., Kanade, T., & Cohn, J. F. (2001). Recognizing action units for facial expression analysis. IEEE Transactions on pattern analysis and machine intelligence, 23(2), 97–115.

17. Vinyals, O., Toshev, A., Bengio, S., Erhan, D.: Show and tell: A neural image caption generator. In: CVPR. pp. 3156–3164. IEEE (2015)

18. You, Q., Jin, H., Wang, Z., Fang, C., & Luo, J. (2016a). Image captioning with semantic attention. In CVPR, pp. 4651–4659. IEEE.

19. Young, P., Lai, A., Hodosh, M., Hockenmaier, J.: From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. Transactions of the Association for Computational Linguistics 2, 67–78 (2014)

20. Nguyen, Hai-Duong & Yeom, Soonja & Lee, Guee-Sang & Yang, Hyung-Jeong & Na, In & Kim, S.H.. (2018). Facial Emotion Recognition Using an Ensemble of Multi-Level Convolutional Neural Networks. International Journal of Pattern Recognition and Artificial Intelligence.

21. Octavio Arriaga and Matias Valdenegro-Toro and Paul Plöger: Real-time Convolutional Neural Networks for Emotion and Gender Classification. (2017)

22. Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Richard Zemel, Yoshua Bengio: Show, Attend and Tell: Neural Image Caption Generation with Visual Attention (2016) arXiv preprint arXiv:1502.03044

23. Kanade, T., Cohn, J.F. and Tian, Y. (2000) Comprehensive Database for Facial Expression Analysis. Proceedings of 4th IEEE International Conference on Automatic Face and Gesture Recognition, Washington DC, 28-30 March 2000, 46-53.