# VISUAL QEUSTION ANSWERING WITH DEEP LEARNING

TU FENGZHI          1155058610

ZHENG ZHIXUAN       1155058600

*Supervised by Prof. LYU Rung Tsong Michael*

# Table of Contents

# TABLE OF CONTENTS

# Introduction

## 1. Overview

In March 2016, a computer program named AlphaGo (AlphaGo) which is developed by Google DeepMind beat Lee Sedol in playing the board game Go and it was the first time a computer program beat human 9-dan professional in game Go because Go is one of most difficult games that have a large search space. AlphaGo uses algorithms in deep learning to "learn" from the previous game record both from top human players and computer to computer records. The success of AlphaGo makes deep learning become a popular topic in computer science and attracts more and more researcher dive into the deep learning area.

Recent years, more and more researchers in both university and industry devote them into the field of deep learning and the number of published articles about deep learning algorithms have exploded. Both the universities and companies are pushing forward to speed up the developments and research of this area since the deep learning can solve a lot of problems better than classic machine learning algorithm. Deep Neural

# INTRODUCTION

Network are widely used nowadays in our daily life such as recommendation systems, automatic translation, auto-drive car, etc.

Among all the popular deep learning problems, we are interested in the problem that joining image and text based on the classic question-answering problem. This problem combines computer vision problems such as image caption and object detection and some text process problems such as natural language processing.

The objective of our project is to use the deep learning method to solve Visual Question Answering problem. Visual Question Answering is a new dataset which include both open-end questions and multiple choice questions about images, it allows us to rise a question to an image, and the deep learning model will answer the question based on the image. Visual Question Answering can be regarded as a promotion of classic text based Question Answering problem. From the previous works, we can find that the most common method to solve this Visual Question Answering problem is to use Convolutional Neural Network to extract visual features from images

first. Then using the Long Short-Term Memory (a kind of Recurrent Neural Network) to process the embedded natural language question vectors.

## 2. Statement of Purpose

The task of our project is to come up with a new model about Visual Question Answering to enhance the accuracy of previous models since there are many limitations in previous models.
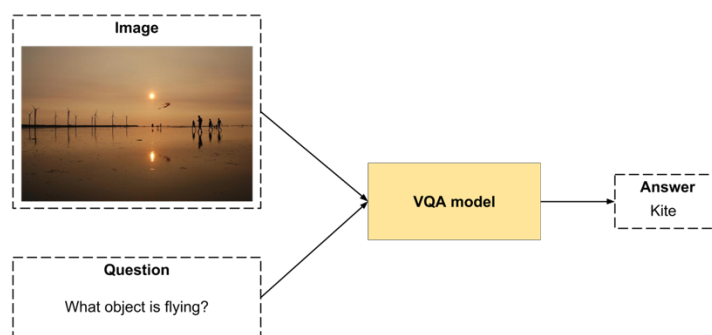


Figure 1 Example of Visual Question Answering

Based on the current work in Visual Question Answering area, the basic objectives of our project are:

- Study and explore the methods and theories used in text-based Question Answering problem.

- Study and explore the methods of extracting the features from images.

# INTRODUCTION

- Study and explore the methods of word embedding and find the best one that most suitable to our project.

- Build a model which can answer questions based on the input image.

- Come up with some models that maybe can enhance the accuracy of the Visual Question Answering and verify their feasibility.

   The advance objectives of our project:

- Try different architectures and parameters to increase model's accuracy.

- Implement a model that have a high accuracy in number-relate and logic-relate questions.

- Extend the language of the model to Chinese, since the main language we use now is English and there exists huge difference semantics between the English word and Chinese word and we should use different way to process the Chinese text.

# Study of Technologies

## 1. TensorFlow

### 1.1 Introduction to TensorFlow



Figure 2 Logo of TensorFlow (TensorFlow)

In our project, we use a famous deep learning platform TensorFlow (TensorFlow), which support various of APIs of deep learning algorithms such as CNN (Convolutional Neural Network), RNN (Recurrent Neural Network) and LSTM (Long Short-Term Memory) which are the most popular deep learning algorithm models in the fields of image processing, speech process and natural language processing. TensorFlow is an open source library for machine learning which is developed by Google brain team and it is based on DistBelief and it is named TensorFlow because of its methodology.

There are two basic but important objects in TensorFlow, one is the Tensor, another is the data flow graph. The basic principle of TensorFlow is to do operation using Tensor. Tensor is a N-dimension array and Flow means the calculation is based on the data flow graph. And the process of "TensorFlow" is the flow of tensor from a node of date flow graph to another node of graph. TensorFlow is a system that transfers the complex data structures to the artificial neural network and analysis and process the input data.

## 1.2 Tensor

Tensor is a linear form function that can be used to represent the relationship of geometric vectors, scalars and other tensors. The relationship can be inner product, outer product and linear mapping and Cartesian product.

The following figure is an example of a tensor. As we can see, the coordinate of a tensor is in a N-Dimension space and contains $n^r$ vectors.
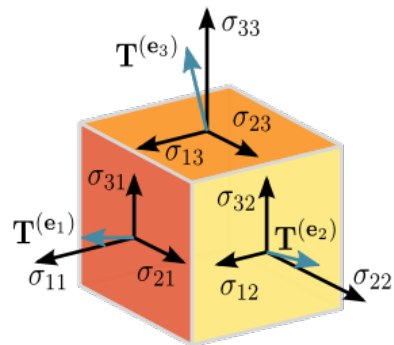
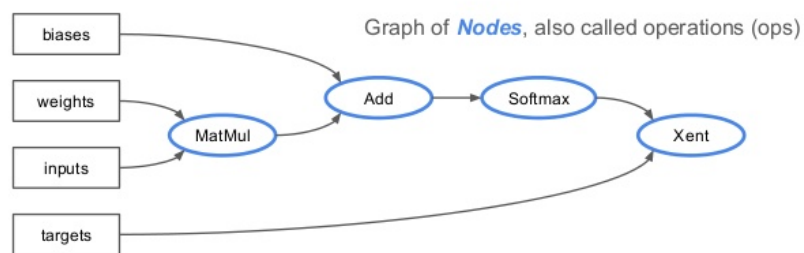**Figure 3 Tensor (htt2)**

## 1.3 Data Flow Graph



**Figure 4 Data Flow Graph (htt3)**

The data flow graph is the computation model using directed graph, nodes and edges. In this graph, nodes not only represent the mathematics

operations, but also represent the operation that read/write the global variable, endpoints for feeding data and endpoints for pushing out result. Edges represent the communication between nodes and indicates the relationship between nodes by transiting tensor, multidimensional data arrays, between nodes. When running the TensorFlow program, the nodes will be assigned to computational devices like CPU and GPU, and each node can run in parallel asynchronously and execute instructions when data of all its incoming edges is really.

### 1.4 Single-Device and Multi-Device

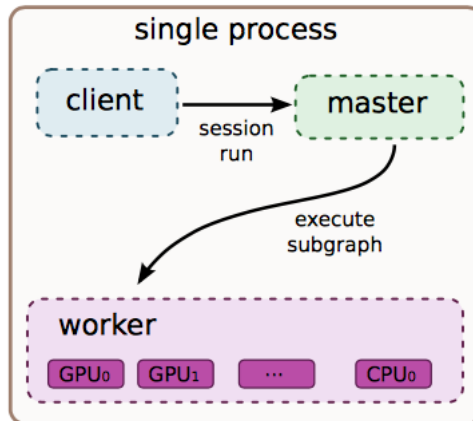TensorFlow not only supports a single machine mode but also supports a distributed system.
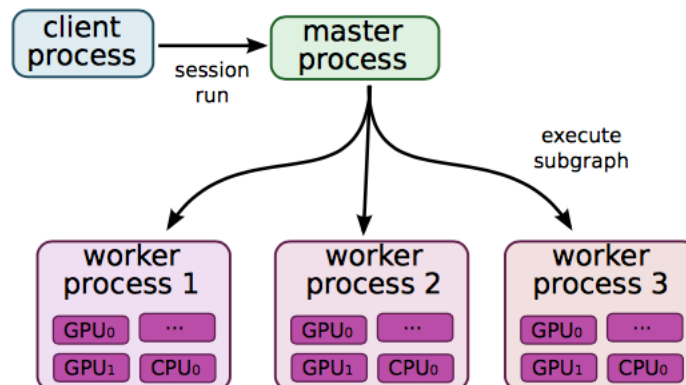
Figure 5 Single Machine (htt4)



Figure 6 Distributed System (htt4)

Thus, when the work load is too heavy for single machine, we can use a distributed system to reduce the runtime of the operations.

## 2. Natural Language Processing

### 2.1 Traditional Methods

The ability for computers to communicating with human in natural language is what human are keeping on chasing. To let computers to be able to communicate with human in natural language, computers need to be able to understand the meaning of natural language and able to express the ideas in natural language. The first one is called Natural Language Understanding and the second one is called Natural Language Generating. Therefore, Natural Language Processing basically contains two parts which are Natural Language Understanding and Natural Language Generation. Historically, scientists did more research on Natural Language Understand and less on Natural Language Generation. However, situation is changed now.

Achieving the communication between human and computers is natural language is very difficult because of:

- Difficulty to determine the boundary of words.

In speaking, there is not boundary between word and word. Therefore, computer need find out the best combination of words based on the context.

- Ambiguity of words.

One word may have different meanings. The computer need to find out the most suitable meaning based on the context.

- Ambiguity of grammar.

Sentences may have different parsing trees because of ambiguous grammar. Therefore, computer need to determine which parse tree is the best based on the context.

When neural network is not introduced in this field, conceptual dependency theory is widely used in Natural Language Processing.

Conceptual dependency theory is first introduced in 1969. In this model, two sentences with the same meaning will have the same representation even though the words in these two sentences are not the same. Four basic representational tokens are used in conceptual dependency theory model which are real word objects with some attributes, real world actions with some attributes, times, locations.

Conceptual dependency representing
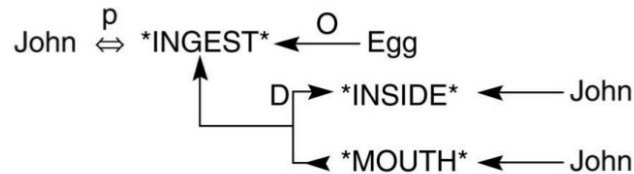"John ate the egg"(Schank and Rieger 1974).



**Figure 7 Example of Conceptual Dependency Theory**

After machine learning is introduced, A revolution occurred in this field. At the beginning, scientist uses machine learning to produce something like decision thing to represent the natural language which is like previous traditional approach. Later, scientists began to focus on develop representation algorithm, like word embedding, on statistical model.

## 2.2 Word Embedding

Word embedding is a set of techniques that maps words or phrases to vectors of real number. We will use word2vec model developed by Google

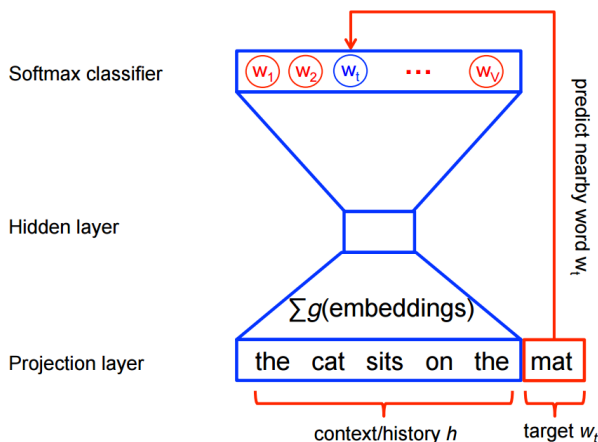Inc. (Mikolov, Sutskever and Chen) to process the natural language text of train questions and answers.



**Figure 8 Model of Word Embedding**

Word2vec [3] learning how words are used in input text by using a shallow neural network. Word2vec output a matrix and each column of this matrix represents a word in input text and provides a numerical description on how the word is used in input text. If the input train text is large enough, the vector that represent two words with similar meaning will have small vector distance. For instance, the word "woman" and word "lady" may have a small vector distance. Most application of word2vec judge the similarity of two words using cosine similarity. Empirically, Word2vec has very nice

performance when the input text is very large, consistent and without ambiguity.

To map the words or phrases to numerical vector, A technique called "skip-gram with negative sample" is used by Word2vec. This method can be roughly divided into 5 steps.

First, take a word in the input text as target and a few words that are close to this target as context.

Second, represent each word in the input text by a random numerical vector. After this step, we can get the numerical vector of the target and context.

Third, put the vector of target and context together and shorten the distance of these vectors.

Fourth, pick up the words randomly from the input text outside the context, and enlarge the distance between the vectors of these words and the vector of the target. By doing this, we can make our target be further away from the words which is rarely used in context.

As keep on applying this process on the targets and their contexts, the vectors of words which are always used together will be pulling closer and closer which the vectors of words which are rarely used together will be pushing further and further.

## 2. Image Processing

In image science, image processing represents a various of technologies that used to analysis and process the images such that the images meet the visual or mentality request or some other technical requirements. It is an application of signal processing in the domain of images. Nowadays, most of images are stored in digital format, thus image processing most refers to the digital image processing. Besides, some methods based on the optical theories still holds an importance position.

Image processing is sub-class of signal processing and has a deep connection with computer science and artificial intelligence. The traditional ways in process the one-dimension signals can be used in image processing such as quantization and noise elimination. But the images are still two-dimension signals which have their own special methods.

### 2.1 Feature Detection

When the neural network is not widely used, scientists use feature detection to gather information from images, deciding whether every single pixel belongs to an image feature. Therefore, the result of feature detection is usually a collection of sets of pixels and the sets contains isolated points, continuous lines or continuous regions.
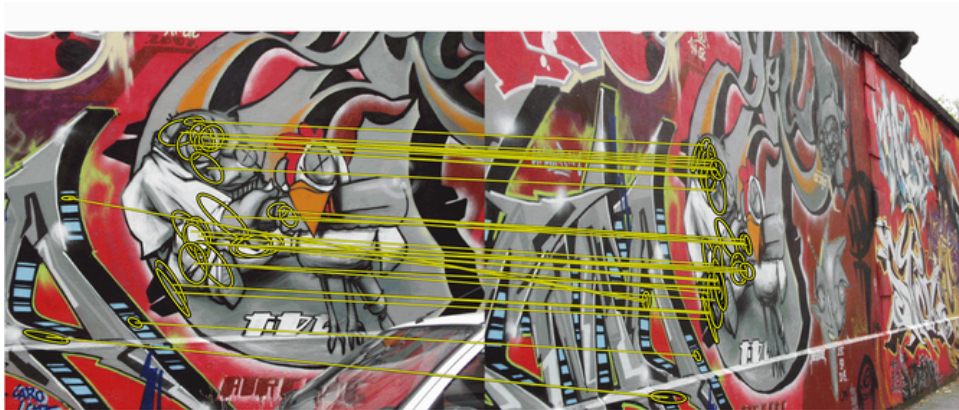


**Figure 9 Feature Detection**

Nowadays, there is not universal and accurate definition of features. Features are usually determined by the problems or the applications. Features can be the interest parts in a digital image and they are the starting point of many computational images analysis algorithm. Whether an algorithm will be successful or not depends on which features this algorithm

uses. Therefore, features must be repeatable which means the features extracted on different images in same scenario should always be the same.

Feature detection is a very fundamental algorithm in image processing, which means it is always the first processing on the images. Feature detection checks every single pixel and determine whether this pixel represents a feature. If feature detection is a part of a larger algorithm, then feature detection usually only check the feature regions of images. As the preprocessing of feature detection, the input images are usually smoothed by using Gaussian kernel.

Many different image processing algorithms use feature detection as its starting point. Therefore, many different detection algorithms are developed and the types of features, the computational complexity and repeatability are varied from each other.

### 2.1.1 Edge detection

Edge detection is a fundamental problem in computer vision and image processing and aims to find the pixel at which image brightness changes a lot in the images. The sharp change in pixels usually refers to important change

of feature. These changes include discontinuities in depth, discontinuities in surface orientation, changes in material properties and variations in scene illumination.
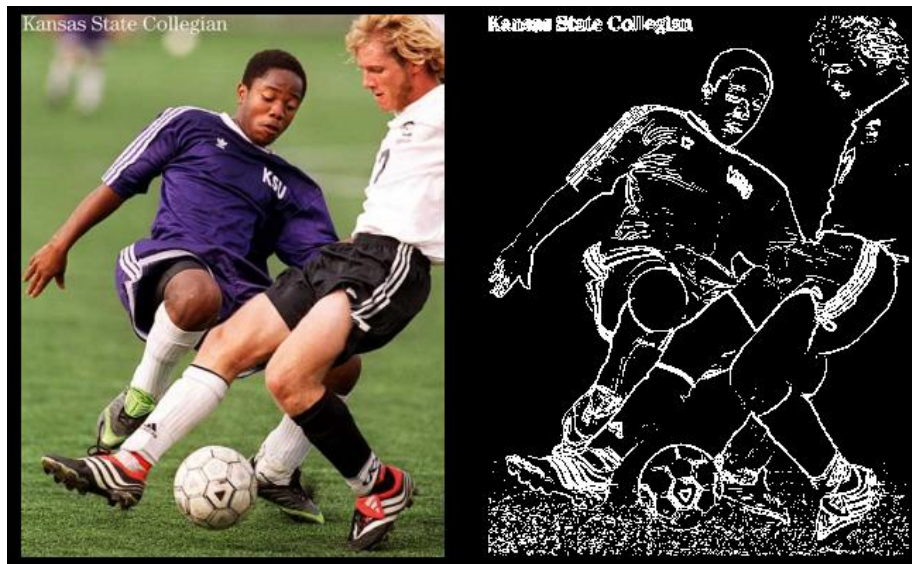
**Figure 10 Sample of Edge Detection**

After applying feature detection, the size of data is significantly reduced, some information which is considered as not important is removed and the important structural features of image are remained.

There are many different methods to perform edge detection, most of them can be categorize into two part, search-based and zero-based. The search-based method detects edges by computing the gradient of images and

find out the local directional maximum of gradient. The zero-crossing based methods search for the zero-crossing of the Laplacian of images.

## 2.2 The Shortcoming of Traditional Methods

Although some traditional method in image processing can achieve a better performance, there are still some shortcomings in these ways. As the rapidly development of the computation speed of computers, deep learning algorithms are wildly used in image processing, such as CNN (Convolutional Neural Network). The results of experiments show that deep learning can perform better than the traditional methods in image processing.



Figure 11 Enhance of Image Using Optical Theory (htt5)

In our project, we should extract the features of an image and feed the features into the neural network. According to the previous experiments,

CNN (Convolutional Neural Network) has an excellent performance in extracting the features of images. And we will introduce the Convolutional Neural Network in the next chapter.

## 2.3 Convolutional Neural Network

In machine learning, convolutional neural network, inspired by animal vision system, is a type of feed-forward neural network that is generally used for image processing because the architecture of convolutional neural network is very suitable for this kind of tasks. For example, the ImageNet (Krizhevsky, Sutskever and Hinton) model has an excellent performance on image classification.

A very important difference between standard artificial neural network and convolutional neural network is that the neurons in the layers of convolutional neural networks is arranged into three dimensions instead of two dimensions in artificial neural network. In addition, differing from the neurons in artificial neural network, the neurons in convolutional neural network only connects to a small region of neurons instead of all neurons in previous layer. (O'Shea and Nash)

### 2.3.1 Overall architecture

Convolutional neural networks are composed of three kinds of layers, which are fully-connected layers, pooling layers and convolutional layers.
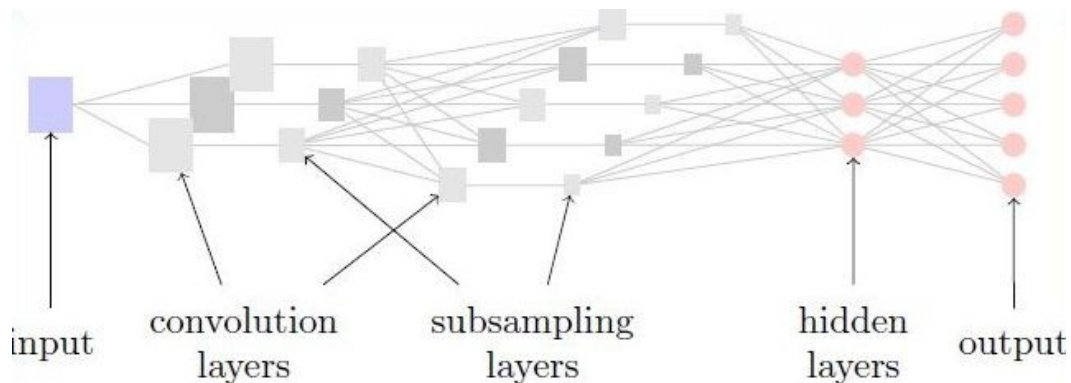


**Figure 12 A Simple Convolutional Neural Network with Five Layer**

As an example, the architecture shown above can be decomposed into several parts.

- Same as traditional artificial neural network, the input layer of Convolutional Neural Network holds the input data.

- Every neuron in the convolutional layers generate its own output by calculating the scalar sum of the products of weights and inputs from the neuron which it connected.

- The pooling layers apply down-sampling to reduce the dimension of input from previous layer. Therefore, the number of neurons, such as, neurons in preceding convolutional layer, and the number of parameters, such as, weights hold by neurons in preceding convolution layer, will be reduced. Therefore, the total model complexity is reduced due to pooling layers.

- The full-layers have the same functionality as in artificial neural network which is calculating the class scores from the activations, which will be used for classification.

According to this small sample, convolutional neural networks can transform the input using several technologies, like convolution, down-sampling layer by layer to produce class score for classification or other purposes.

### 2.3.2 Convolution Layer

Convolution layer is the most important layer of convolutional neural network and the parameter in this kind of layer is the learnable kernels.

These kernels are commonly very small in terms of dimension comparing with input. In convolutional layer, when the input is ready, the input convolves with the small kernel to create a 2D activation map.
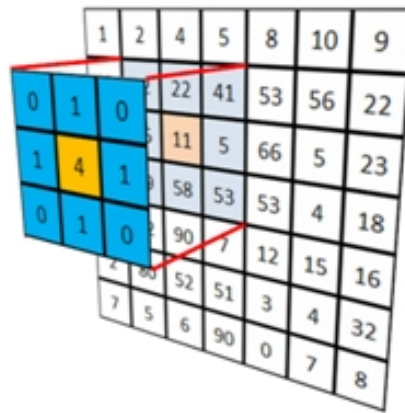


Figure 13 Gaussian Kernel

As we scan the input, the convolution result will be calculated using the kernels for each portion of input. Kernels will learn whether a specific feature has been sawn at a certain position of the input from the network and these processes are usually called activations. Each kernel has its own activation map based on the input, and these activation maps will be stack together to form a 3D-array as output of convolutional layers.
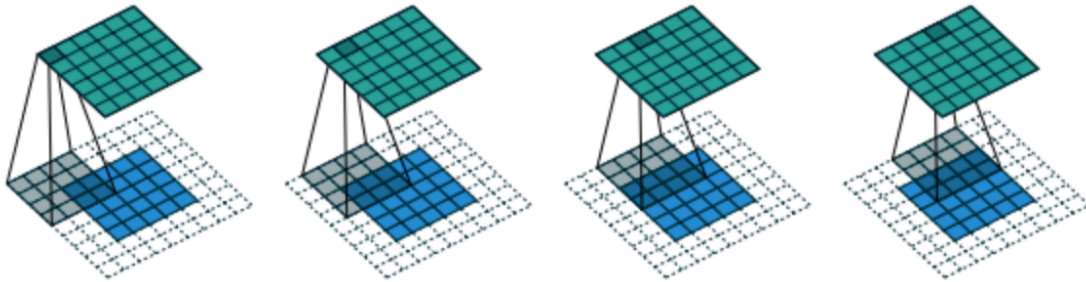
**Figure 14 Convolution Process of An Image**

As we all known, artificial neural networks which are been trained for image processing are always too large due to fully-connected manner of standard artificial neural network, therefore this kind of artificial neural networks cannot be efficient enough for training. To solve this problem, neurons in convolutional layers only connect to a small region of neurons in previous layer. The size of this region is called receptive field size.

For instance, if the input image is RBG image with size 100 * 100 * 3 and receptive field size is set to 5 * 5, each neuron in convolution layers should only maintain 75 weights comparing with 30000 weights for neurons in artificial neural networks.

Convolution layers also can optimize their output to reduce the model complexity significantly. There are three basic hyper-parameters for optimization which are the depth, the strike, and the size of zero-padding.

Depth is the third dimension of the output and it controls the number of neurons connect to the same region of input. Reducing the depth can significantly reduce the neurons in convolutional layers.

Stride controls the height and weight of output. If the stride is set to 1, the receptive fields will be heavily overlapped and the output will be very large. So increasing the stride can reduce the overlapping of receptive fields and the size of output.

Zero-padding is the process that padding on the border of the input. Therefore, the number of zeros that zero-padding applies can further control the size of output.

It is important to understand these three hyper-parameters and we can calculate the size of output based on the formula $(N - F + 2 * Z) / (S + 1)$ where N is the size of input, F is receptive field size, Z is the size of zeros padding on the border and S is the stride. If the result is not an integer, the

stride is set incorrectly because the neurons cannot fit in this configuration on the input.

### 2.3.3 Pooling layers

Pool layers is used to reduce the size of data and therefore reduce the parameters (weights) of each neuron and model complexity. The pool layers take activation maps as input and reduce their size by using Max function. Since the pool layers use Max function, this kind of pool layers are usually referred to max-pooling layers. For instance, if a kernel with size 2 * 2 is applied with stride 2 to the input, the height and width will all be reduced to the half of their counterpart of input. Therefore, the size of output is reduced to only 25% of the size of input with the depth unchanged.
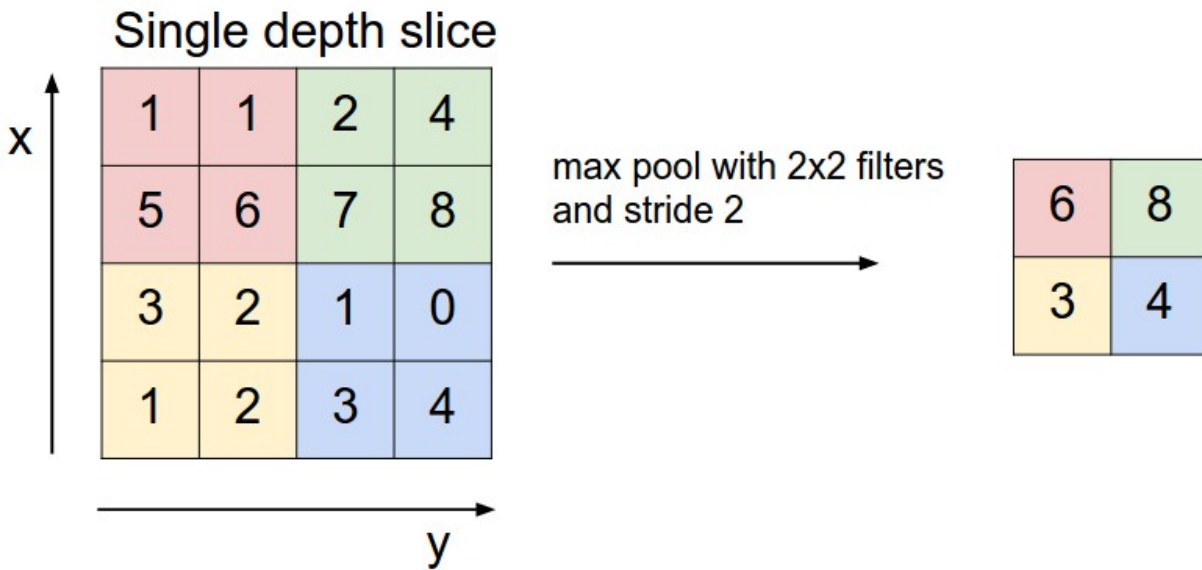
**Figure 15 Pooling Layer**

As the pooling layers remove information from input, there are only two ways to implement max-pooling. The first one is already described in example above. The kernel is set to 2 * 2 and stride is set to 2. Another one is called overlapping-pooling, where the kernel is set to 3 * 3 and stride is set to 2. In general, the size of kernel cannot exceed 3, because it will remove too much information.

### 2.3.4 Fully-connected Layers

each neuron in fully-connected layer connects to all the neurons in previous layer and preceding layer. The functionality of fully-connected

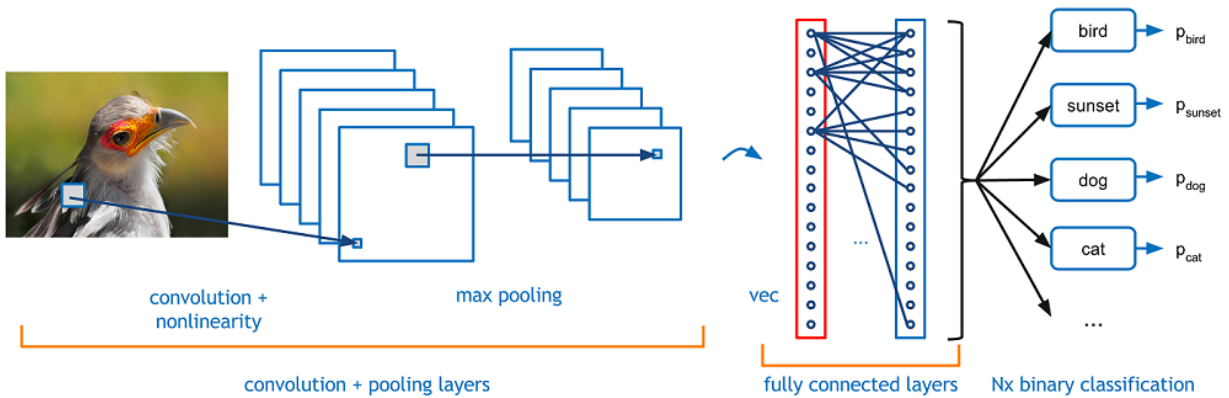layers in convolutional neural networks is the same as theirs in artificial neural network.



**Figure 16 Fully Connected Layers In CNN**

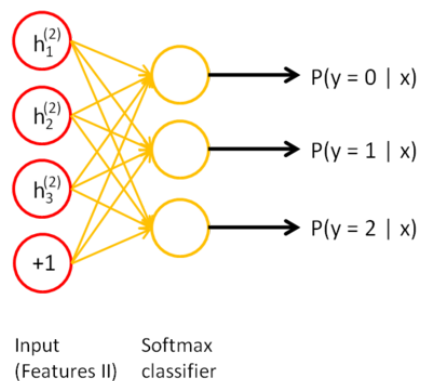## *2.3.5 Softmax Layer (Output Layer)*



**Figure 17 Softmax Layers**

Softmax function is an extension of logistic function. The logistic regression is to solve the binary problem which are labeled {0, 1} .

The hypothesis function of logistic regression is:

$$h_\theta(x) = \frac{1}{1 + \exp(-\theta^T x)}$$

The cost function of logistic regression is:

$$J(\theta) = -\frac{1}{m} \left[ \sum_{i=1}^{m} y^{(i)} \log h_\theta(x^{(i)}) + \left(1 - y^{(i)} \log\left(1 - h_\theta(x^{(i)})\right)\right) \right]$$

The train process is to minimize the cost function $J(\theta)$. After we obtain

the optimal parameter, we can use the logistic function to predict the result

of input data. But the logistic function can only be used in binary problem,

when the labels of data is larger than two, we need to use k logistic classifiers

and it is fussy. In this case, we can use the softmax function to solve the

problem.

The hypothesis function of softmax function is:

$$h_\theta(x^{(i)}) = \begin{bmatrix} p(y^{(i)} = 1 | x^{(i)}; \theta) \\ p(y^{(i)} = 2 | x^{(i)}; \theta) \\ \vdots \\ p(y^{(i)} = k | x^{(i)}; \theta) \end{bmatrix} = \frac{1}{\sum_{j=1}^{k} e^{\theta_j^T x^{(i)}}} \begin{bmatrix} \theta_1^T x^{(i)} \\ \theta_2^T x^{(i)} \\ \vdots \\ \theta_k^T x^{(i)} \end{bmatrix}$$

The cost function of logistic regression is:

$$J(\theta) = -\frac{1}{m}\left[\sum_{i=1}^{m}\sum_{j=1}^{k}1\{y^{(i)}=j\}log\frac{\theta_j^T x^{(i)}}{\sum_{l=1}^{k}\theta_l^T x^{(i)}}\right]$$

The train process of the softmax function is like to the logistic function, and finally we can use the function to predict the result:

$$P(y=j|x) = \frac{e^{x^T w_j}}{\sum_{k=1}^{K}e^{x^T w_k}}$$

# 3. Long Short-Term Memory (Recurrent Neural Network)

## 3.1 Recurrent Neural Network

When human beings read articles, they comprehend the whole article by understanding each word based on the meaning of previous words since a single word doesn't have any meaning.

In traditional Neural Network, there have no connection between the nodes in the same layer. That means traditional Neural Network cannot preserve the relation between the original input data and this is a major defect of traditional full connected neural network.
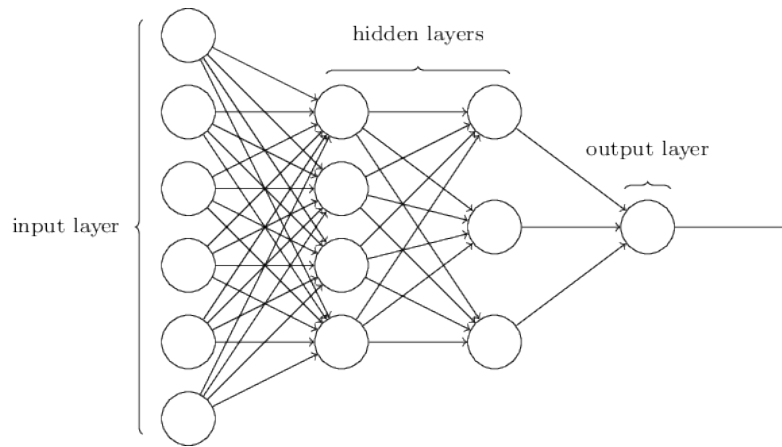
**Figure 18 Structure of Traditional Neural Network**

Recurrent Neural Network is designed to preserve the relationship between the node and its previous node. There are connections between nodes in the same layer of Recurrent Neural Network and this property makes Recurrent Neural Network becomes applicable to solve the problems like handwriting and natural language processing.
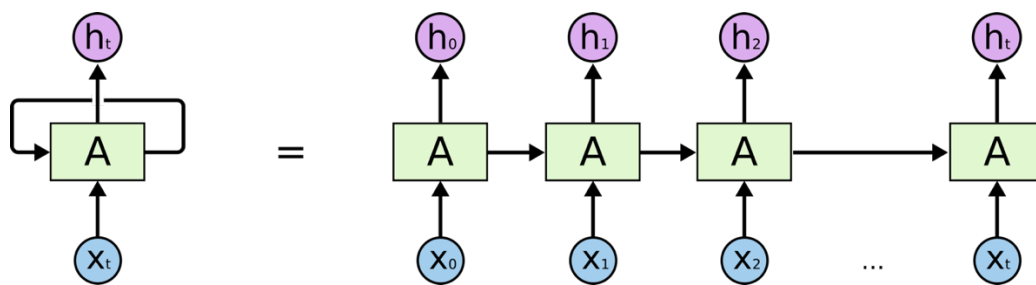


**Figure 19 A Single Layer of Recurrent Neural Network (Understanding LSTM Networks)**

In the above figure, A is a single layer in a Recurrent Neural Network while it has input parameter X and output value h and the value of information can pass from the previous nodes to next nodes in the same layer, that is from Xt to Xt-1.

Recurrent Neural Network are widely used in recent years since it has been proved to have an excellent performance in solving a various of problems like speech recognition, language translation and so on.

### 3.2 Long Short-Term Memory

Although Recurrent Neural Network have a good performance in some area, but it still has some shortcomings. Sometimes, we only need to loop back to just a few steps, for example, when we want to predict the next word the user may input, we may only look up the previous several words that has been inputted. In the above case, Recurrent Neural Network can have an excellent performance. But if we want to solve the Question Answering problem, looking up only several previous words doesn't work anymore. We need the content from further back. In practice, Recurrent Neural Network doesn't seem to solve this "Long-Term dependencies" problem well.

Long Short-Term Memory is explicitly designed to solve the "Long-Term dependencies" problem (Understanding LSTM Networks) and it is kind of Recurrent Neural Network which proposed in 1997.
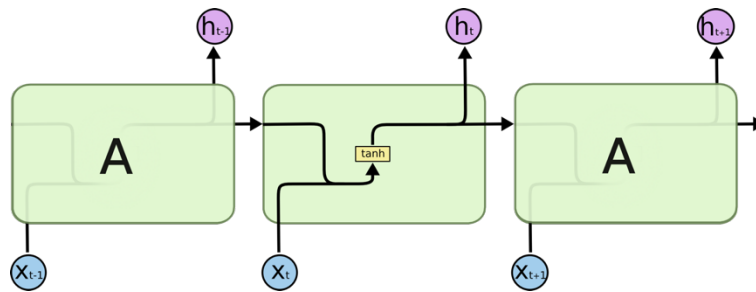


**Figure 20 The Structure of a Single Node of Standard Recurrent Neural Network (Understanding LSTM Networks)**
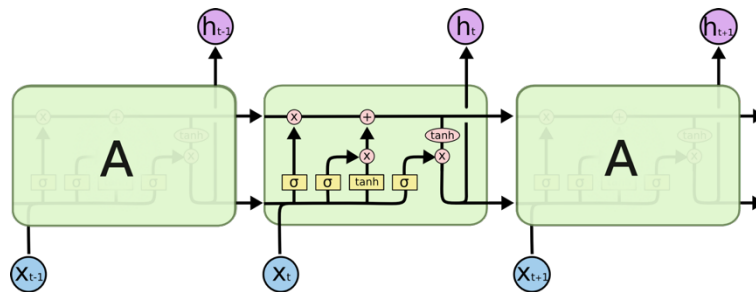


**Figure 21 The Structure of a Single Layer of Long Short-Term Memory (Understanding LSTM Networks)**
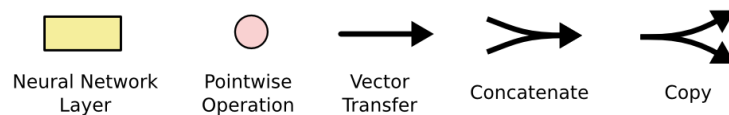


**Figure 22 The Meaning of Symbols (Understanding LSTM Networks)**

Long Short-Term Memory have the same chain structure like Recurrent Neural Network, but they have a different repeat module in a single node. In General, a standard Recurrent Neural network only has single neural network layer while Long Short-Term Memory has four.

# Methodology

In general, an image can be described as a row * column * 3-dimension matrix which in a high dimension and text can be represented as a vector which in a low dimension. What we need to is to reduce the dimension of images or increase the dimension of text so that these two items can be combined and feed to neural network.

In the first, we consult to the structure of Question-Answering model based on text, then extent the text-base to image-base structure.

## 1. Text-Based QA Model

### 1.1 Word Embedding

First, we perform the word embedding to convert word strings into vectors. The method to perform the word embedding is call Skip-gram (Mikolov, Sutskever and Chen). First, we will traverse all the words in the training set and get a unique word dictionary vector. Second, we can replace words in all the sentence in the train set with the index of the word in the dictionary vector. After converting sentence in the train set into vectors, we

can use Long Short-Term Memory (Understanding LSTM Networks) (LSTM) to capture more information about the sentence.

## 1.2 End-to-End Memory Network (Sukhbaatar, Szlam and Weston)

After converting the words to vectors, end-to-end memory network is used to solve question answering problem, this network is a kind of recurrent neural network where the recurrence reads from a large data before it produces the result.
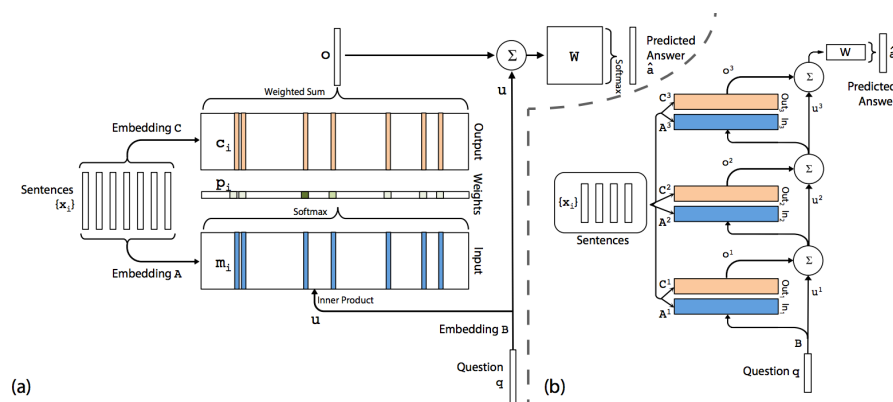


**Figure 23 (a) A Single Layer Version (b) A Three Layer Version (Sukhbaatar, Szlam and Weston)**

End-To-End Memory Network can divide into the following parts:

- Input representation:

In this stage, the model maintains several embedding matrices with the same dimension. We use a certain embedding matrix to convert the input vectors to memory vectors and convert the query to an internal state. After that, we compute the match, which is a possibility, of the memory vectors and the internal state by using SoftMax function on the dot product of these vectors.

- Output representation:

In this stage, we use the output vectors, which are calculated for each input vector by using another embedding matrix, to compute to the response vector of the memory. The response vector is simply the sum of the products of each output vector and its corresponding possibility which is computed in last stage.

- Find out the prediction:

In the single layer case, the prediction of the network is to sum all response vector and the internal state and then pass this sum through the weight matrix. Finally, we can use Softmax function to get the result.

## 2. From Text to Image

### 2.1 Extract Image Features.

In our application, we used the VGG model (K. Simonyan) to extract the feature from training set.

VGG model contains a stack of convolutional layers with kernel size of 3 * 3. The convolutional stride is 1 and the zero-padding size is also 1. For pooling layers, VGG model contains five max-pooling layers with 2 * 2 kernel and stride equal to 2. Following the convolutional layers, there are three Fully-connected layers. The first and second fully-connected layer contains 4096 neurons each while the third layer has only 1000 neurons because this layer is used to perform ILSVRC classification. Finally, we add a soft-max layer as final layer.

Instead of training our own VGG model, we use the VGG provided Keras (VGG16 model for Keras). After deploy the model in our application, we extract the feature from the images in MSCOCO dataset. In our program for extracting feature, we set parameter batch size to 10, which means the program will process 10 images in one turn. In each turn, we extract feature

of these ten images, which are vectors, using the VGG model. After the features are produced by the VGG model, we store the features into a file as the input of LSTM.

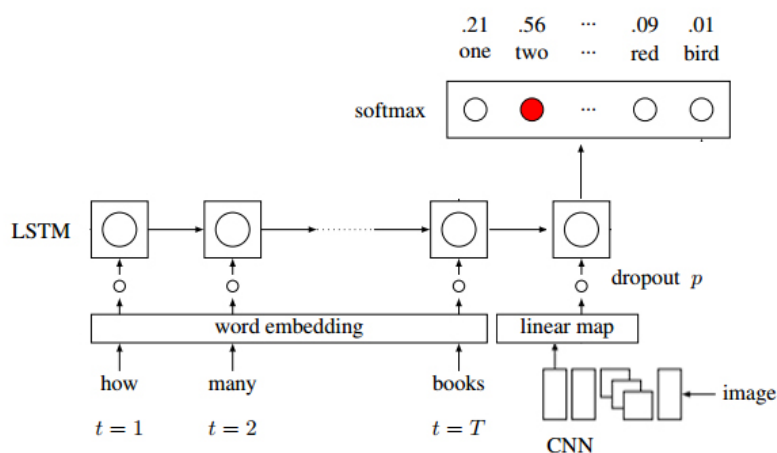## 2.2 Feed the Image Features to LSTM



Figure 24 Whole Model of Visual Question Answering

After using Convolutional Neural Network to extract the features from the image train set, we feed these features with the text vectors to the Long Short-Term Memory to train the model.

# 3. Hypothetical Model

To improve the accuracy of number-relative question, we decide to use some ways to process the image such that we can count all objects in an image and feed the result into the neural network as image features.

## 3.1 Object Counting Model

To improve the performance of our model for visual question answering, we build a sub-model for objection detection which may have benefit to our model.

This object detection is referred to the paper (Girshick), which describe a kind of neural network called Fast Region-based Convolutional Network. Comparing with other neural network, this neural network has faster training speed and more accurate detection.
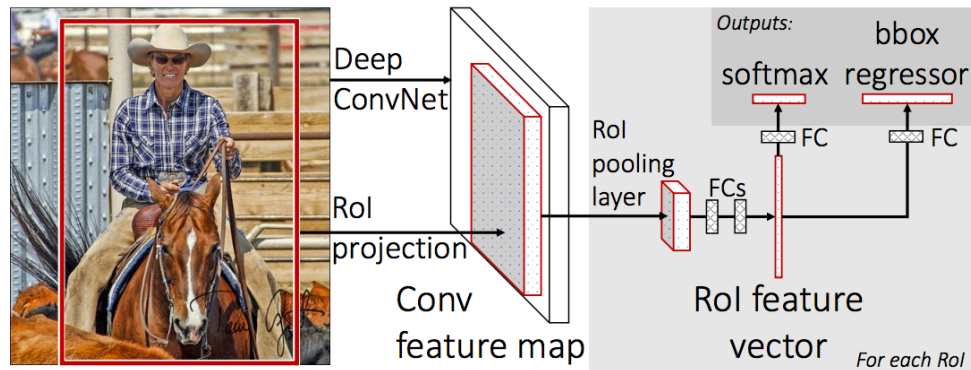
**Figure 25 Fast R-CNN Architecture (DIY Deep Learning for Vision: a Hands-On Tutorial with Caffe)**

Fast R-CNN (Girshick) uses the whole image and its object proposals as input. First, Fast R-CNN passes the input through several convolutional layers and max-pooling layer to generate the convolutional feature map of this image. After that, the convolutional feature map is passed to a region of interest (RoI) layer to generate a fixed-length feature vector for each proposal. Finally, these feature vectors are passed through two separate fully-connected layer. For the first fully-connect layer which perform softmax function, this layer produces a numerical number for each proposal which indicate the possibility that this object is in this image. For the second fully-connected layer, this layer output four numbers for each proposal which indicate the location of the object in this image.

In our model, as the accuracy of answering question like predicting number of objects is very low, I think Fast R-CNN is very suitable to our model. As mentioned earlies, the output of Fast R-CNN can be considered as a set of regions of interest (rectangles) indicating objects and the corresponding label for each object in these regions of interest. Therefore, by reviewing the output of Fast R-CNN, we can find out the numbers of different objects in input image.

For example, the input image is shown below.

In this image, there are a female and a horse. After the processing of Fast R-CNN, we can toughly consider the output as image below.
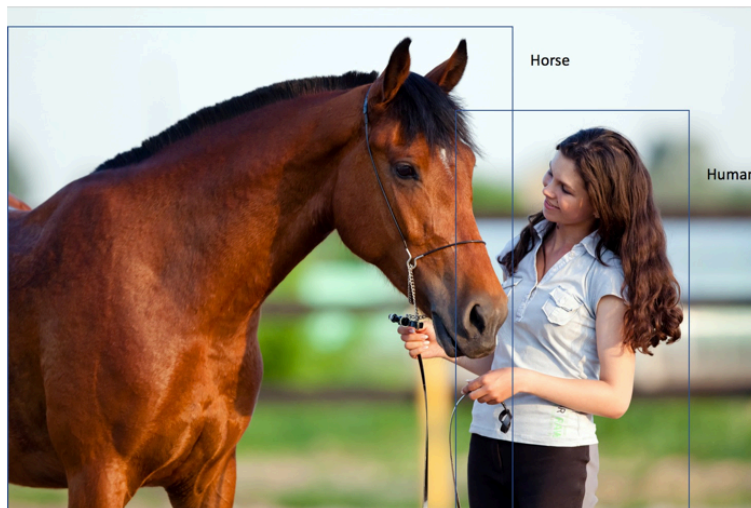
**Figure 27 Output of Fast R-CNN**

In the output of Fast R-CNN, there are two regions of interest (rectangles) and a region of interest contains a female while the other region of interest contains a horse.

Therefore, by pushing these kinds of information to the LSTM for question answering. It is possible that the model can answer the question about counting, like "How many people are there in this image" or "Is there any horse in this image", more accurate. Therefore, this model can help us improve the accuracy of our visual question answering model.
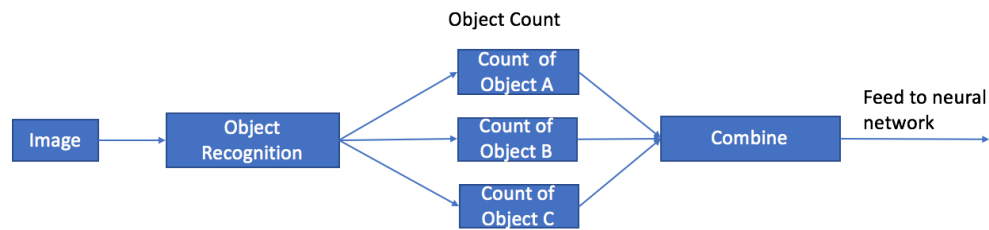
**Figure 28 Count Model**

In this model, we want to use object recognition to count each object in an image and together with the image features to enhance the accuracy of model.

## 3.2 Image Caption Model

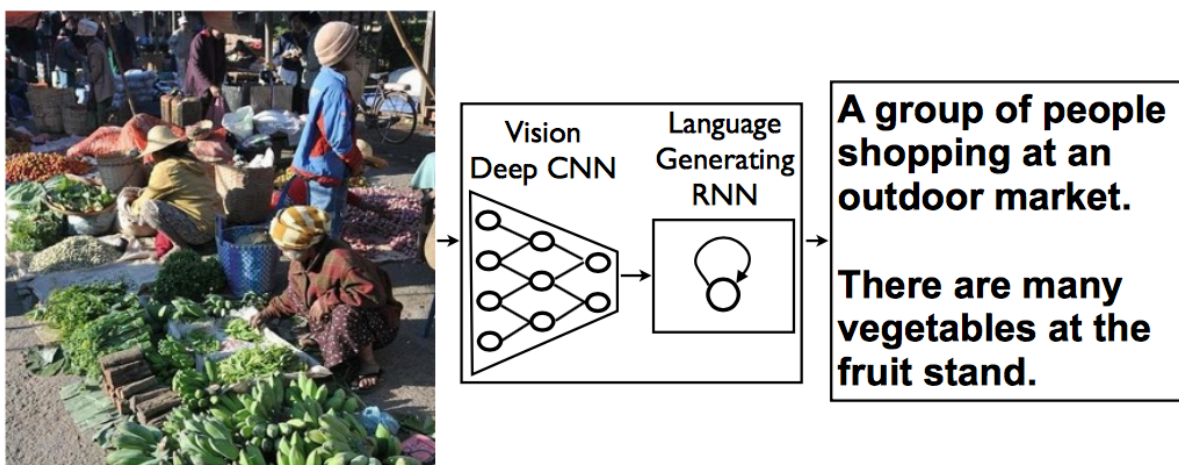### 3.2.1 NIC Model (Show and Tell: A Neural Image Caption Generator)



**Figure 29 Sample Output**

The score of this model is as following:

| Metric | BLEU-4 | MeTEOR | cider |
|--------|--------|--------|-------|
| NIC | 27.7 | 23.7 | 85.5 |
| Random | 4.6 | 9.0 | 5.1 |
| Neighbor | 9.9 | 15.7 | 36.5 |
| Human | 21.7 | 25.2 | 85.4 |

## 3.2.2 Recurrent Visual Representation (Mind's Eye: A Recurrent Visual Representation for Image Caption Generation)
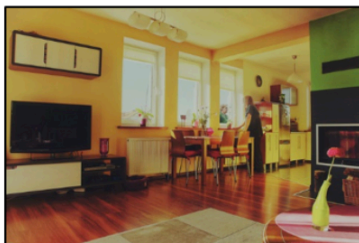


A table topped with plates of food and bowls of food.
This table is filled with a variety of different dishes.

A man that is jumping in the air while riding a skateboard.
A man on a skateboard is performing a trick at the park.

A brown and white dog sitting on top of a street .
A picture of a dog laying on the ground.

A large living room filled with furniture and a flat screen tv.
A woman stands in the dining area at the table.

A group of motorcycles parked on the side of a road.
A motorcycle parked in a parking space next to another motorcycle.

A close up of a sink in a bathroom.
A faucet running next to a dinosuar holding a toothbrush.

**Figure 30 Sample Result**

The overall score of this model:

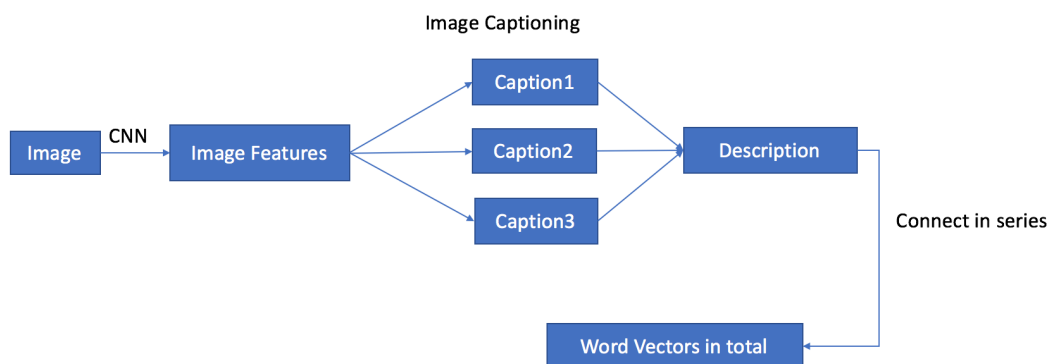| | Flickr 8K | | | Flickr 30K | | | MS COCO Val | | | MS COCO Test | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | PPL | BLEU | METEOR | PPL | BLEU | METEOR | PPL | BLEU | METEOR | BLEU | METEOR | CIDEr |
| RNN | 17.5 | 4.5 | 10.3 | 23.0 | 6.3 | 10.7 | 16.9 | 4.7 | 9.8 | - | - | - |
| RNN+IF | 16.5 | 11.9 | 16.2 | 20.8 | 11.3 | 14.3 | 13.3 | 16.3 | 17.7 | - | - | - |
| RNN+IF+FT | 16.0 | 12.0 | 16.3 | 20.5 | 11.6 | 14.6 | 12.9 | 17.0 | 18.0 | - | - | - |
| RNN+VGG | 15.2 | 12.4 | 16.7 | 20.0 | 11.9 | 15.0 | 12.6 | 18.4 | 19.3 | 18.0 | 19.1 | 51.5 |
| Our Approach | 16.1 | 12.2 | 16.6 | 20.0 | 11.3 | 14.6 | 12.6 | 16.3 | 17.8 | - | - | - |
| Our Approach+FT | 15.8 | 12.4 | 16.7 | 19.5 | 11.6 | 14.7 | 12.0 | 16.8 | 18.1 | 16.5 | 18.0 | 44.8 |
| Our Approach+VGG | 15.1 | 13.1 | 16.9 | 19.1 | 12.0 | 15.2 | 11.6 | 18.8 | 19.6 | 18.4 | 19.5 | 53.1 |
| Human | - | 20.6 | 25.5 | - | 18.9 | 22.9 | - | 19.2 | 24.1 | 21.7 | 25.2 | 85.4 |

### 3.2.3 Overall Structure



**Figure 31 Caption Model**

In our first model, we decide to using several different ways to do generate image captions and combine them into a sentence which represent the image. And in this case, we convert an image into sentence so that the Visual Question Answering problem is converted to a Text-Based Question Answering problem.

But this model still has many defects. First, current models in image caption field cannot describe the content in image well and too many information will lose in this process.

## 4. Dataset

Since we need a large amount of data to train a model, and the data set should contain image, question and answer pairs.

The training and testing dataset we used in our project is provided by the organizers of Visual Question Answering challenge which is sponsored by MS COCO website. This dataset contains three subsets which are train set, validation set and test set. Inside the dataset, there are two types of questions and first one is Open-End questions and another one is Multiple-Choice questions. The detail is as following (16ht):

Real Images:

- 82783 MS COCO training images
- 40504 MS COCO validation images
- 81434 MS COCO testing images

- 248349 questions for training

- 121,512 questions for validation

- 244,302 questions for testing (3 per image)

- 2483490 answers for training

- 1,215,120 answers for validation (10 per question)

Abstract Images:

- 20000 training images

- 10,000 validation images

- 20,000 MS COCO testing images

- 60000 questions for training

- 30000 questions for validation

- 60000 questions for testing (3 per image)

- 600000 answers for training

- 300000 answers for validation (10 per question)

# Implementation Detail

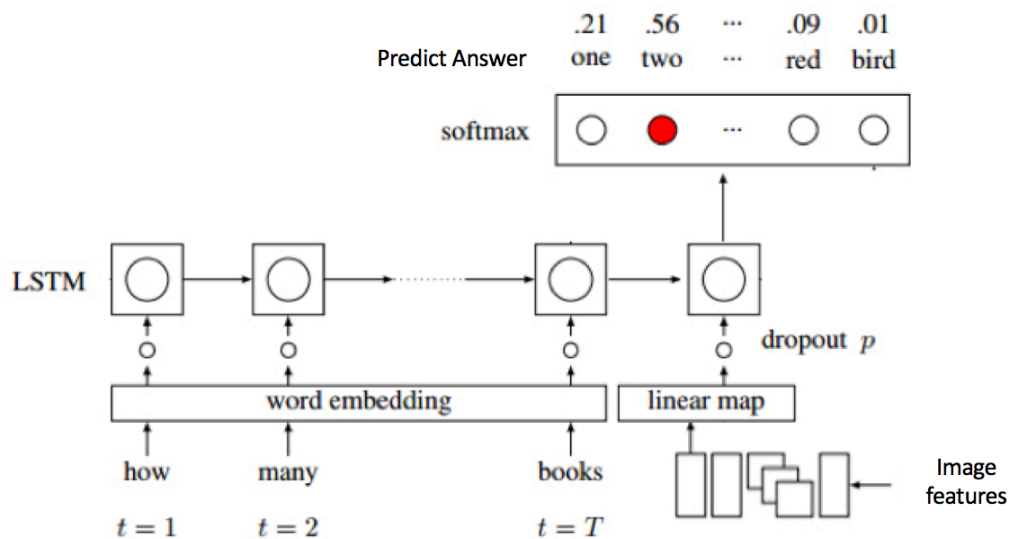The overall architecture of our model is like:



Figure 32 Overall Architecture

In the first, we used a preprocess program to extract the train images'
features and using HDF5 binary data format to store it. Then we used the
word embedding to convert the words to vectors and combining the question
vectors and the image features, and last we feed the combined data to the
LSTM (Long Short-Term Memory) to train the model.

# 1. Data Loader

Since the training images and the questions are separated in different files and in the format of JSON.

The image file name is like COCO_train2014_[image id].jpg. For example, there is an image whose name is COCO_train2014_000000348957.jpg, then the image is 348957.

In the train question file, the format is as following and it contains an image id and a question id.

{

"question": "How many windows can you see?",

"image_id": 434410,

"question_id": 4344102

}

The format of answer is as following, and it contains a question id and several choices while each choice has its own answer id.

{

"answer": "kettles",

"answer_confidence": "yes",

"answer_id": 6

}

Above all, the image id is not and question id may not continuously. In this case, we have a function to load the image, question and answer and binding them together.
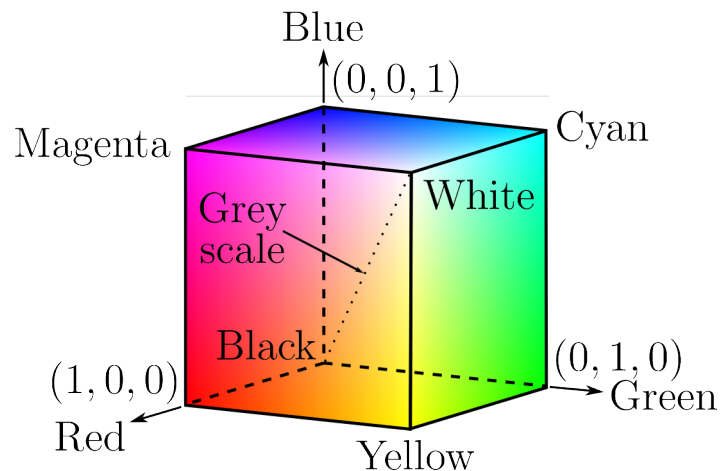
## 2. Extract Image Features



Figure 33 RGB Channels

After we can load the train data, we need to using Convolutional Neural Network to extract the feature from the images. The VGG-16 model accepts R, G, B channels that between 0 and 1. Then for each image we split it to R channel, G channel and B channel respectively. And then resize it to 224 * 224 dimension.

For example, in the following image, we first separate it into its three color channel, and then resize its dimension to (224, 224).



**Figure 34 Train Image**

**Figure 35 R Channel**



**Figure 36 G Channel**



**Figure 37 B Channel**

```
def load_image_array(image_file):
    img = misc.imread(image_file)
    # GRAYSCALE
    if len(img.shape) == 2:
        img_new = np.ndarray( (img.shape[0], img.shape[1],
3), dtype = 'float32')
        img_new[:,:,0] = img
        img_new[:,:,1] = img
        img_new[:,:,2] = img
        img = img_new

    img_resized = misc.imresize(img, (224, 224))
    return (img_resized/255.0).astype('float32')
```

After preprocessing the train images files, we feed them into the VGG-16 Model to extract their features.

## 1.1 VGG-16 Model

The preprocess program use the VGG-16 model as the basic model, the structure is like following.

```
layer {
  name: "conv1_1"
  type: "Convolution"
  bottom: "data"
  top: "conv1_1a"
  convolution_param {
    num_output: 64
    pad: 1
    kernel_size: 3
```

```
  }
}
layer {
  name: "relu1_1"
  type: "ReLU"
  bottom: "conv1_1a"
  top: "conv1_1"
}
.
.
.
layer {
  name: "drop7"
  type: "Dropout"
  bottom: "fc7"
  top: "fc7"
  dropout_param {
    dropout_ratio: 0.5
  }
}
layer {
  name: "fc8"
  type: "InnerProduct"
  bottom: "fc7"
  top: "fc8"
  inner_product_param {
    num_output: 1000
  }
}
layer {
  name: "prob"
  type: "Softmax"
  bottom: "fc8"
  top: "prob"
}
```

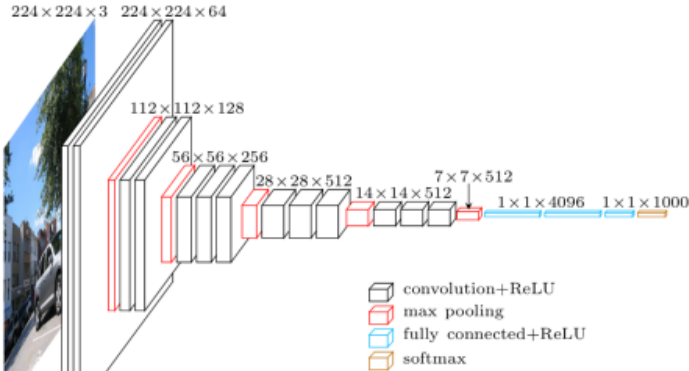The rough structure can describe as following:

Figure 38 The Structure of VGG-16 Model (htt6)

## 1.2 Features Format

The feature of an image is a vector whose shape if 4096 * 1.



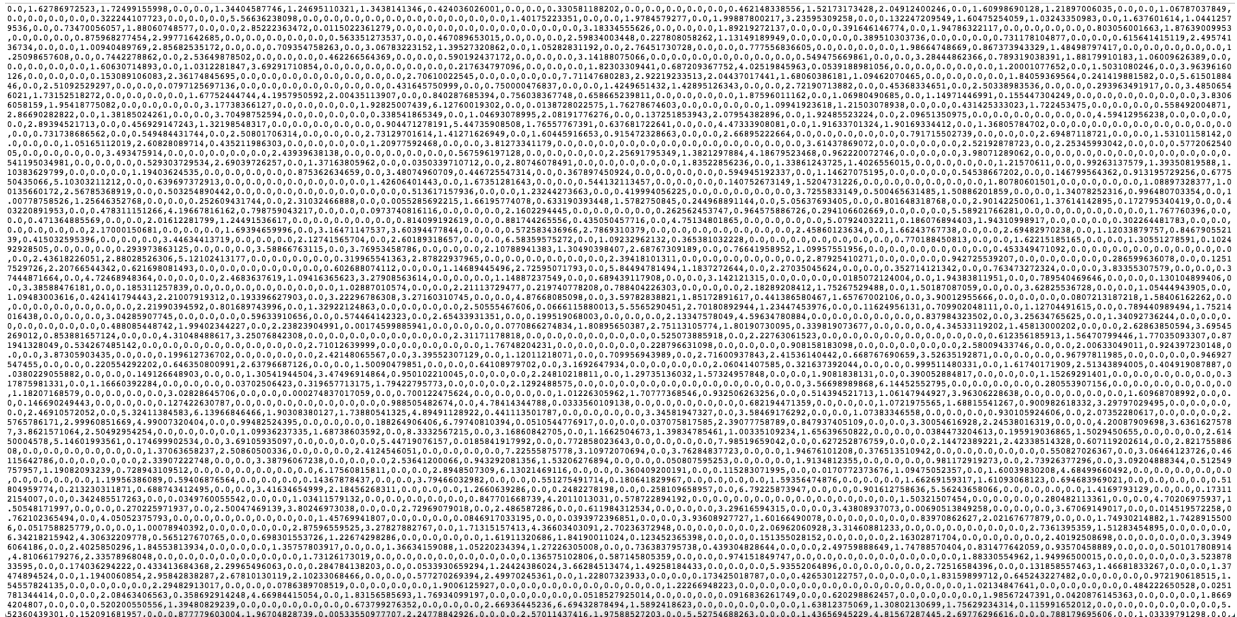Figure 39 Feature of An Image

Page 57

For each train images, we have a vector whose shape is 4096 * 1, after we finishing extracting the image features, we can obtain a matrix whose shaper is number of images * 4096, and this large matrix will store in HDF5 binary data format file. Since HDF5 has a smaller size when we store the data.

## 3. Training

We build up a LSTM Model. This LSTM model contains 2 hidden layers. and each hidden layer has 512 neurons. Therefore, for each hidden layer, it will output a 512-length vector.

For Input layer, it maintains three parameters which are Wimg, Wemb, bimg. Wimg is a matrix [4096, 512] and all elements in this matrix are random values in range [-1, 1] initially. Wemb is a matrix [q_vob_len, 512] where q_vob_len is the number of words in all questions and all the elements in this matrix are 0 initially. bimg is a 512-length vector which represents the bias. The usages of three parameters will be described below. When the data batch feed in, since the dimension of image features and words is not the same, we need to regularize these dimensions. For words, we do word

embedding to generate a 512-length vector to represents words using Wemb. For image, we multiply the Wimg and 4096-length image feature vector and add the bimg to generate another 512-length vector. Therefore, the dimensions of image feature vector and word vector are regularized.

## 4. Predicting

We first load the pre-trained QA model and VGG model from local. Then we use VGG to produce the image feature vector of input image. After that, we use pre-known knowledge to embedding the words in question and reduce the dimension of image feature vector. We concatenate these vectors and feed resultant matrix into pre-trained QA and generate the word vector representation of answer. Finally, we look up this vector in our pre-known knowledge and find out the answer.

# Results and Analysis

## 1. Training Process

The train log can help us to improve our models, and the following figures are showing the Training Accuracy and Losses during each epoch.
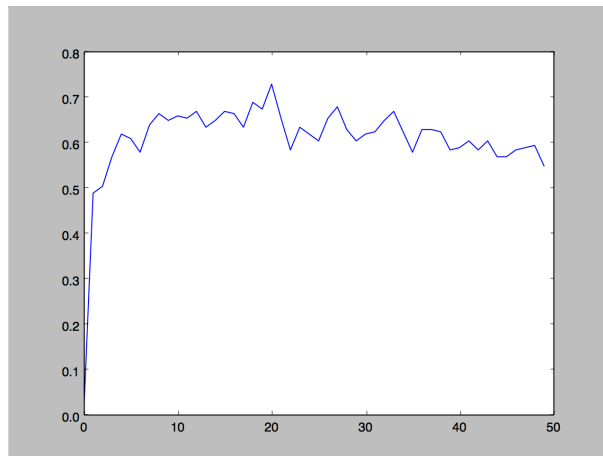


**Figure 40 Training Accuracy**

As we can from the above figure, the training accuracy go up quickly in the beginning of train process, then in maintain an accuracy between 50% to 60% and trend to converge in the middle of training process. And after finishing the whole training, the final accuracy is about 55%.
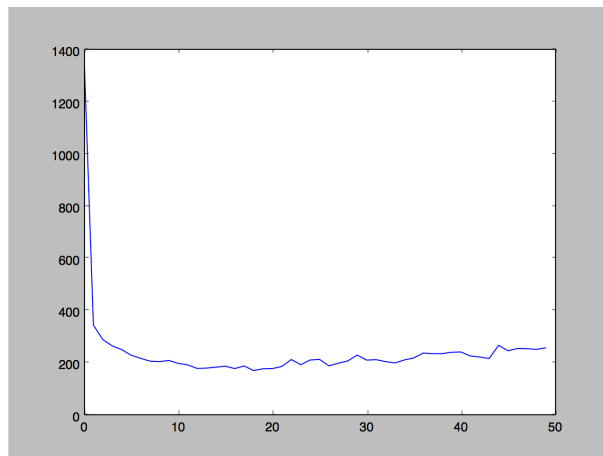
**Figure 41 Training Losses**

From the above figure, we can know that the train process was about to converge quickly. The training accuracy and the training losses didn't change a lot after epoch 10. And the total training loss is about 300 eventually.

## 2. Evaluation

### 2.1 Evaluation Metric

We use the evaluation code which offer by the Visual Question Answering challenge organizers. In this evaluation metric, the accuracy of machine will be averaged over all 10 choose 9 sets of human annotators to enhance the robustness of accuracy.

The accuracy formula is as follow:

$$Acc(ans) = \min\{\frac{\#\ humans\ that\ said\ ans}{3}, 1\}$$

## 2.2 Evaluation Dataset

The dataset we used for evaluation is mentioned in Chapter

Methodology.
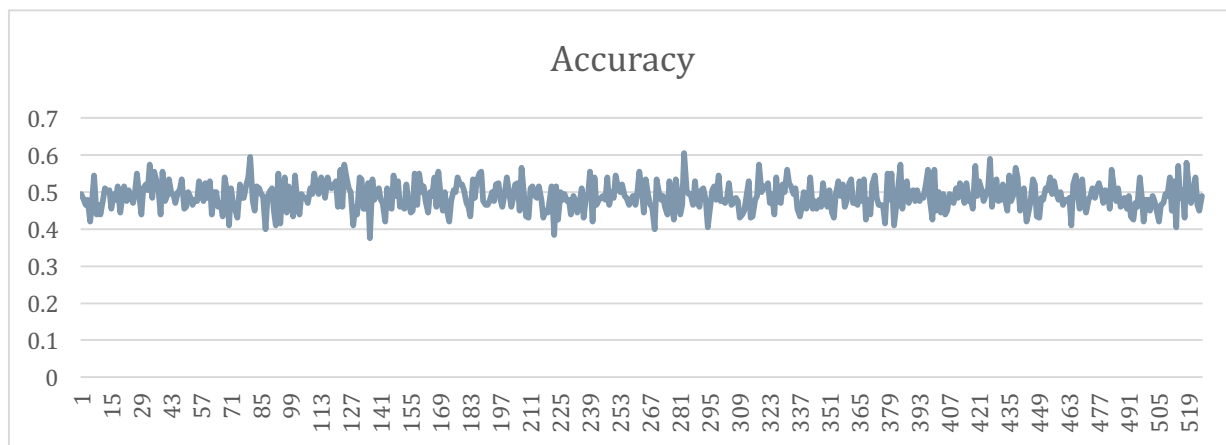
# 3. Accuracy

## 3.1 Overall Test Result



**Figure 42 Overall Accuracy of Current Model**

## 3.2 Detail Accuracy

**Table 1 Accuracy of Our Model**

| Model | Test Set | | | |
|---|---|---|---|---|
| | Yes / No | Number | Other | Overall |

| | | | | |
|---|---|---|---|---|
| Baseline All Yes | 70.97 | 0.36 | 1.21 | 29.88 |
| Baseline Prior Per Question Type | 71.40 | 34.90 | 8.96 | 37.47 |
| Baseline Nearest Neighbor | 71.89 | 24.23 | 22.10 | 42.85 |
| UC Berkeley & Sony | 83.62 | 39.47 | 58.00 | 49.12 |
| Our Model | 74.62 | 31.76 | 31.32 | 49.12 |

# 4. Samples

Question: What animal is this?

Top 5 Answers: horse, cow, elephant, dog, zebra

Question: What color is this animal?

Top 5 Answers: brown, white and brown, white, black and white, black

Question: What are they doing now?

Top 5 Answers: walking, talking, eating, standing, taking picture



**Figure 44 Sample 2**

Question: What is this that in the front?

Top 5 Answers: train, building, truck, bikes, teddy bears

Question: How about the weather now?

Top 5 Answers: cloudy, cold, rainy, wet, evening

Question: How many cars are in there?

Top 5 Answers: 2, 3, 1, 5, 4



Figure 45 Sample 3

Question: What is this man doing?

Top 5 Answer: skateboarding, jumping, snowboarding, sitting, taking picture

Question: How about the weather?

Top 5 Answer: rainy, sunny, cloudy, clear, overcast

Question: What is the gender of the person?

Top 5 Answer: male, female, child, boy, women



**Figure 46 Sample 4**

Question: what are sitting on the counter in different stages of cutting with a knife?

Top 5 Answers: napkin, chocolate, glass, bread, fries

Question: What is the color of this vegetables?

Top 5 Answer: orange, yellow, green, blue, brown

Question: What is the name of this vegetable?

Top 5 Answer: <span style="color:red">carrot</span>, parsley, apple, bananas, carrots



**Figure 47 Sample 5**

Question: What is the color of this animal?

Top 5 Answer: <span style="color:red">black and white</span>, white, brown, black, gray

Question: What is this animal?

Top 5 Answer: <span style="color:red">zebra</span>, giraffe, horse, cow, zebras

Question: How many animals are there?

Top 5 Answer: 2, 3, <span style="color:red">4</span>, 1, 5

**Figure 48 Sample 6**

Question: How many animals are there?

Top 5 Answer: 2, 1, 3, 4, 5

Question: What is the name of this animal?

Top 5 Answer: bird, bear, elephant, cat, dog

Question: what are sitting down on the ground?

Top 5 Answer: teddy bear, bear, snow, stick, bird

**Figure 49 Sample 7**

Question: what is playing with the large UNK of ice?

Top 5 Answer: surfer, surfing, surfboard, water, boat

Question: What is the color of this animal?

Top 5 Answer: white, brown, black, white and brown, gray

Question: How many animals are there in the image?

Top 5 Answer: 2, 1, 4, 3, 6

Figure 50 Sample 8

Question: How many animals are there in the image?

Top 5 Answer: 2, 3, 1, 4, 5

Question: What is the name of object that in the middle?

Top 5 Answer: remote, laptop, toy, cat, books

Question: What is the name of object that in the left?

Top 5 Answer: laptop, keyboard, remote, phone, bowl

**Figure 51 Sample 9**

Question: what are flying through the sky?

Top 5 Answer: kites, plane, kite, clouds, airplane

Question: What is the color of background?

Top 5 Answer: blue, red, green, orange, yellow

Question: How many objects in the sky?

Top 5 Answer: 13, 10, 4, 5, 1

**Figure 52 Sample 10**

Question: what is the black dog holding?

Top 5 Answer: frisbee, kite, bat, carrot, dog

Question: What is the color of this dog?

Top 5 Answer: black, brown, black and white, white and brown, white

Question: what is the color of the object that the black dog holding?

Top 5 Answer: yellow, blue, white, pink, red

**Figure 53 Sample 11**

Question: What is this woman doing?

Top 5 Answer: cooking, eating, taking picture, smiling cutting

Question: What is this woman holding?

Top 5 Answer: banana, fork, plate, knife, pizza

Question: what color is this woman's eyes?

Top 5 Answer: blue, brown, black, green, red

## 5. Analysis of Result

The model we used now get an accuracy of 49.12% over the test set and we present the accuracy of baseline model and the model from UC Berkeley & Sony. The accuracy of model from UC Berkeley & Sony represent the highest accuracy of current Visual Question Model, we can see that the Visual Question Answering is still a hard problem over the world since the accuracy of number-relative question is much low than the accuracy of Yes / No questions.

In our opinion, there are several reasons that the accuracy is not so high. First, the Convolutional Neural Network we use is mainly used to classify the objects, which means that the image features that we extracted contained the class of each object in the image but didn't contain the count of each objects. In this case, the model is to use "common sense" to answer the number-relative questions. Thus, the model has much high accuracy on Yes / No question than the number-relative question.

In some extent, the current model has a high performance in recognizing the type of questions. If the question is asking the number of an

object, the top 5 answers will be a number, and if the question is asking about the weather, the top answers will also show some kinds of weathers. In this case, we can say that our current model has high accuracy of answering the question-type but may have low accuracy in answering the question.

We will focus on all the above bottlenecks in the next term and to enhance the whole accuracy of the model.

# Conclusion

## 1. Term Review

When we start our final year project, we know little about neural networks or even about machine learning. Therefore, we hope that we can learn something about machine learning, and then have better understanding on machine learning.

After several weeks of studying and researching, looping back to our works, we think we achieved what we want. Although we lose the direction and consider whether we needed to change the topic of the final year project, we did not give up and kept on doing research on this topic.

At the beginning, we knew nothing about Q&A model or even neural network. We set up a small millstone that we implemented a text-based Q&A first. To achieve this millstone, we read some paper about LSTM, RNN. In this step, the support from Professor Michael R. Lyu and PhDs helped us a lot.

Starting from the text-based Q&A model, we started to find a method can merge the information in image and information in text and use this information to predict the answer. Since the in the text-based Q&A model, the

words are represented by vectors, it is obvious that we also can convert images into vectors. There are several ways for converting images into vectors, such as, using CNN to extract feature vectors from images, using R-CNN for find out attention and embedding the captions, which are produced by caption generator, of images to vectors. At the end, we chose the CNN to extract the feature.

Finally, based on these knowledge and technologies, we are able to implement the first version of Visual Question Answering.

## 2. Shortcoming and further work

### 2.1 The answer is not accuracy when answering number-related question.

This is the common problem for every Visual Question Answering model using CNN because it is difficult for CNN to find out the number of objects in an image. Therefore, we will deploy the R-CNN into our model to improve the accuracy of answering number-related question. The architecture of R-CNN and the principal of how R-CNN words have already been mentioned above.

**2.2 The model cannot fully comprehend some questions when the questions contain logic judgement.**

This is the common problem for every Visual Question Answering model because there is no question about logic judgement in training dataset and then the trained model cannot recognize this kind of problems. In term of this shortcoming, we did not find out an efficient solution and we will keep on doing research on this issue.

**2.3 The current model can only support English question.**

Since English is the most popular language in the world, most of the research is on English. However, as the Chinese is becoming more and more important in the world, we think it is necessary to develop a Visual Question Answering model which is consistent of Chinese. The difficulty on implement a model consistent of Chinese is that we need to embed Chinese characters to vectors and we need Dataset of Chinese Visual Question Answering. After we did some research on this topic, we found that, in term of dataset, Baidu provides a dataset for Chinese Visual Question Answering and, in term of

# CONCLUSION

Chinese character embedding, Word2vec is sufficient for embedding Chinese

which is the same as embedding Latin characters.

# Acknowledgement

We would like to express our gratitude to our supervisor Prof. Michael Lyu for providing a lot suggestions and opinions. Secondly, we would also like to thank Mr. Edward Yau for his technical support. Last but not least, we would like to express our special thanks to Li Jian, Su Yuxin, Zeng Jichuan who gave us a lot of ideas to do this wonderful project. Without their help, we would encounter more difficulties in this project.

# Reference

27 8 2015. 29 11 2016. <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>.

n.d. 29 11 2016. <http://www.visualqa.org>.

n.d. <https://gist.github.com/baraldilorenzo/07d7802847aaad0a35d3>.

n.d. <http://tutorial.caffe.berkeleyvision.org>.

n.d. <https://upload.wikimedia.org/wikipedia/commons/thumb/4/45/Components_stress_tensor.svg/300px-Components_stress_tensor.svg.png>.

n.d. <http://image.slidesharecdn.com/april131700googleramanathan-160425195435/95/machine-intelligence-at-google-scale-tensorflow-13-638.jpg?cb=1461614089>.

n.d. <http://download.tensorflow.org/paper/whitepaper2015.pdf>.

n.d. <https://zh.wikipedia.org/wiki/图像处理>.

n.d. <https://www.cs.toronto.edu/~frossard/post/vgg16/vgg16.png>.

*AlphaGo*. n.d. Wikipedia. 29 11 2016.

   <https://en.wikipedia.org/wiki/AlphaGo>.

Girshick, Ross. "Fast R-CNN." (n.d.).

*GitHub*. n.d. <https://github.com/paarthneekhara/neural-vqa-tensorflow>.

K. Simonyan, A. Zisserman. "Very Deep Convolutional Networks for Large-
   Scale Image Recognition." (n.d.).

Krizhevsky, Alex, Ilya Sutskever and Geoffrey E Hinton. "ImageNet
   Classification with Deep Convolutional Neural Networks." (n.d.).

Mikolov, Tomas, et al. "Distributed Representations of Words and Phrases
   and their Compositionality." (n.d.).

"Mind's Eye: A Recurrent Visual Representation for Image Caption
   Generation." (n.d.).

O'Shea, Keiron Teilo and Ryan Nash. "An Introduction to Convolutional Neural
   Networks." (n.d.).

"Show and Tell: A Neural Image Caption Generator." (n.d.).

Sukhbaatar, Sainbayar, et al. "End-To-End Memory Networks." (n.d.).

*TensorFlow*. n.d. 29 11 2016. <https://www.tensorflow.org>.