# Visual Question Answering

#5 Jada, Santosh Yadav| #11 Mehta, Avni | #14 Muppala, Raji | #19 Patel, Harshil Lavjibhai

## Abstract

In the world of computer vision, there are many projects on visual question answering, but we put forth this novel methodology to get answers from images and also tried to incorporate a visual question answering system with the help of bounding box concept. we are taking MS COCO data which got 80,000 images with different domain. We could use this system to help visually impaired people to identify the objects in a given image. This is an application that enables the user to either select an image from the available images or upload an image. The user can ask a question about the selected image. The main purpose of this application is to provide a natural language answer in real-time.

## 1. Introduction

Image processing and computer vision has become popular over the last few years and research is growing exponentially. In this project we are using MS COCO data set and limiting our VQA model to interior (bathroom, bedroom, kitchen, living room). We already built models using Machine Learning and Deep learning algorithms like CNN. We are using bounding box to identify the object accurately. This is done when we use Mask R-CNN with python and tensorflow. We are creating a box by picking a smallest box that captures all the pixels of the mask as the bounding box.

## 2. Related Work
- **Open Source Projects**

There have been several papers that are working on tasks of image tagging, image captioning and text-based Q&A. A team from MIT has worked on 'Simple Baseline for Visual Question Answering'.
Show and Tell: Lessons learned from the 2015 MS-COCO Image Captioning Challenge.

- **Literature Reviews**

We found the following scholarly papers related to Visual Question Answering topic:
[01] Learning to Answer Questions from Image Using Convolutional Neural Network
[02] An attention based convolutional neural network for visual question answering
[03] Exploring Models and Data for Image Question Answering
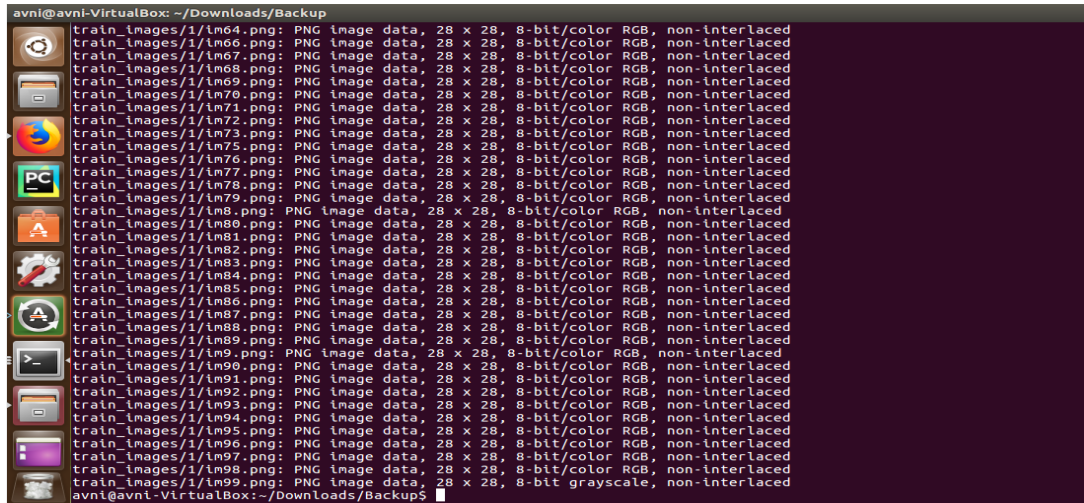
## 3. Approach
- Segregate Interiors (bathroom, bedroom, kitchen, living) from MSCOCO data.
- Build Bayes Machine learning models Decision Tree, Random Forest.
- Compare the results. Run test data through Clarify, build deep learning algorithms - SoftMax, CNN.
- Compare the accuracy between all the models. Extract VQA related to interiors and build VQA.

- Implement bounding box by using mask R-CNN with python and tensorflow.

In order to perform the SoftMax and CNN on our project dataset, we had to perform the following pre-processing steps:

1. Resize all the training and test images to 28x28 pixels.



2. Change the image color from RGB to black and white.
3. Convert the resized images to MNIST format

| | | | |
|---|---|---|---|
| t10k-images-idx3-ubyte | 3/13/2018 1:28 PM | WinRAR archive | 23 KB |
| t10k-labels-idx1-ubyte | 3/13/2018 1:28 PM | WinRAR archive | 1 KB |
| train-images-idx3-ubyte | 3/13/2018 1:28 PM | WinRAR archive | 225 KB |
| train-labels-idx1-ubyte | 3/13/2018 1:28 PM | WinRAR archive | 1 KB |

- **Analytical Tools**

  We have used the following tools for this increment:
  1. Shallow learning – Decision Tree, Naïve Bayes, Random Forest Model
  2. Clarify API, Spark API
  3. Deep Learning – SoftMax
  4. Deep Learning – CNN + Mask R-CNN

- **Analytical Task**

  1. Build Decision Tree, Naïve Bayes, Random Forest models. Out of all Random Forest got good accuracy.
  2. Object detection using Random Forest.
  3. Found accuracy and object detection using clarify.
  4. Convert data set to ubyte.
  5. Found accuracy using SoftMax and CNN.
  6. Implementation of bounding box using Mask R-CNN.


- **Expected Inputs/Outputs**

  Input: Images belonging to one of the four Interiors (bathroom, bedroom, kitchen, living).For each category we had 50 images for training and 10 images for testing.

  Expected Output:

  Find the accuracy using Random Forest, Clarify, SoftMax and CNN compare the results.
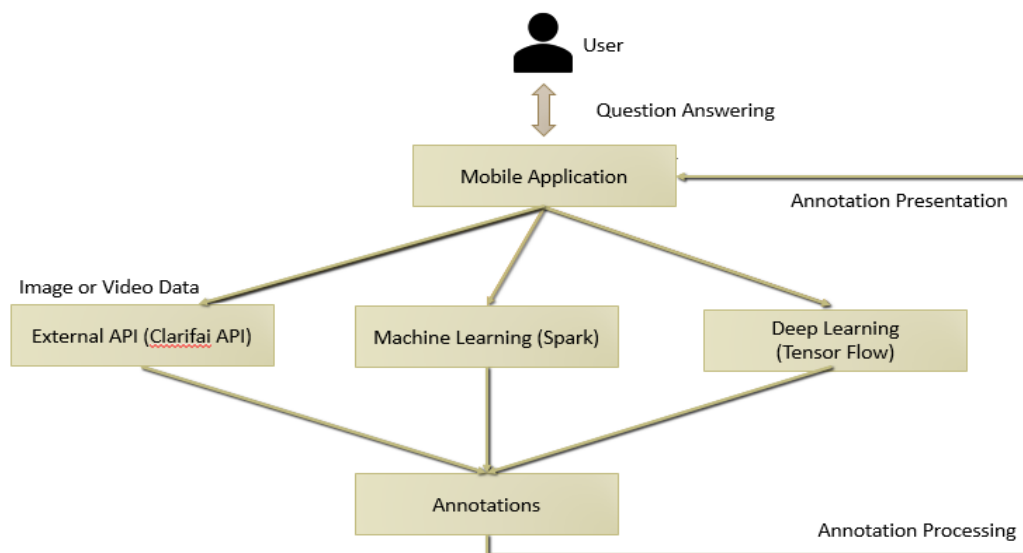
  Find out different objects using bounding box

  Answer the question asked by the user.

## 4. Application Specification & Implementation

### System Specification

- **Software Architecture**



The above figure represents the software architecture of this project. There will be a user interface in form of a web application where in the user can select image and ask the natural language question.

We have built following models Machine learning – Decision Tree, Naïve Baiyes, Random Forest Model, Deep Learning – SoftMax, CNN and also Mask R-CNN.

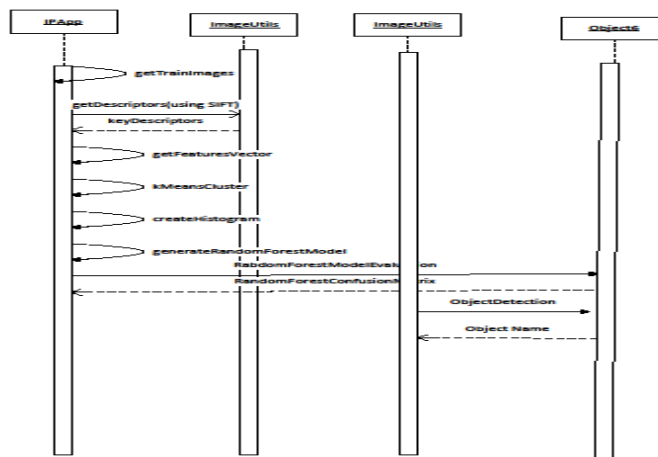- **Features, Workflow, Technologies**

1. Activity Diagram
   The activity diagram tells the process of how searching the dataset, preprocessing stage, refining the dataset, training the model and then building the model, testing the model and finally tuning the parameters to get the accuracy.
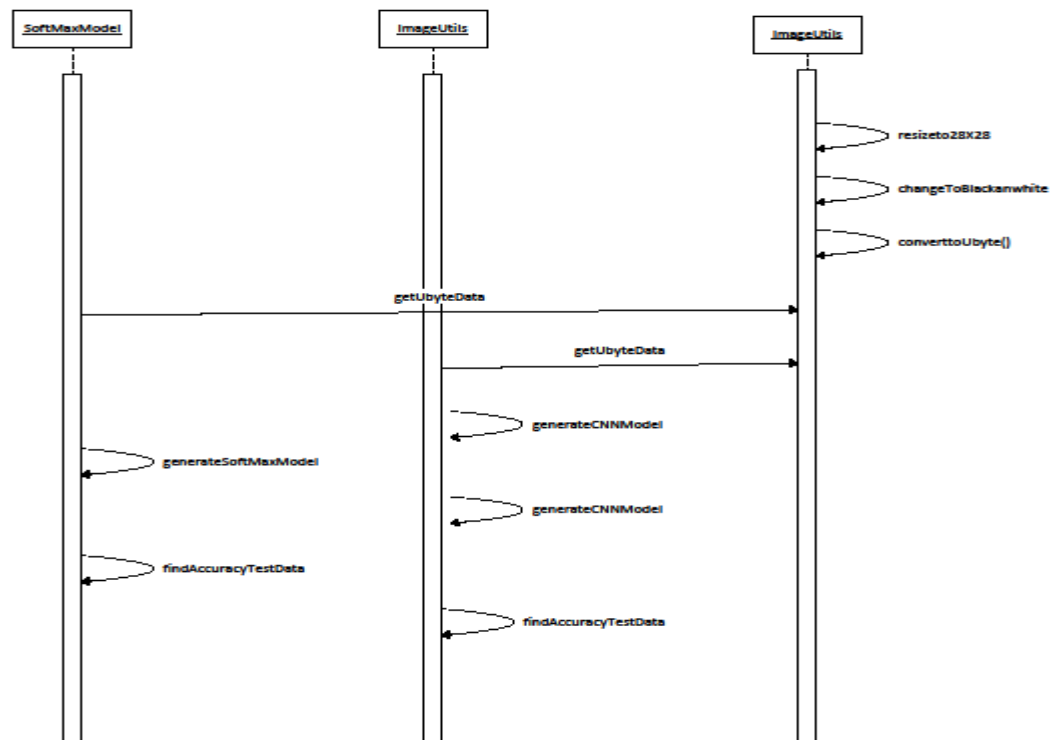
| Search Dataset | → | Preprocess Data | → | Refined Datset | → | Train and Build model | → | Test the Model | → | Get the accuracy |
|---|---|---|---|---|---|---|---|---|---|---|

2. Sequence Diagram:

   This is the sequence diagram when we use Spark :



   This is the sequence diagram when we use deep learning:

- **Feature Specification**

1. In this project increment, we have mainly focused on the Deep learning models (SoftMax, CNN) for our for interiors (bedroom, bathroom, kitchen, living) data.
2. Converted data to ubyte.
3. Compared accuracy with Random Forest Model, Clarify, SoftMax, CNN with different number of epoch and CNN AdamOptimizer, AdagradOptimizer
4. Usage of bounding box using Mask R-CNN. It is done by using python and tensor flow. This model generates bounding boxes for objects in the image.
5. The weight histograms generated would be helpful in the identification. We started of with anchor sorting and filtering, then refinement is applied on the bounding box. Then the images are scaled.
6. It also shows the percentage of the prediction.

- **Object Detection**

  Spark API Correct object Detection as Bedroom
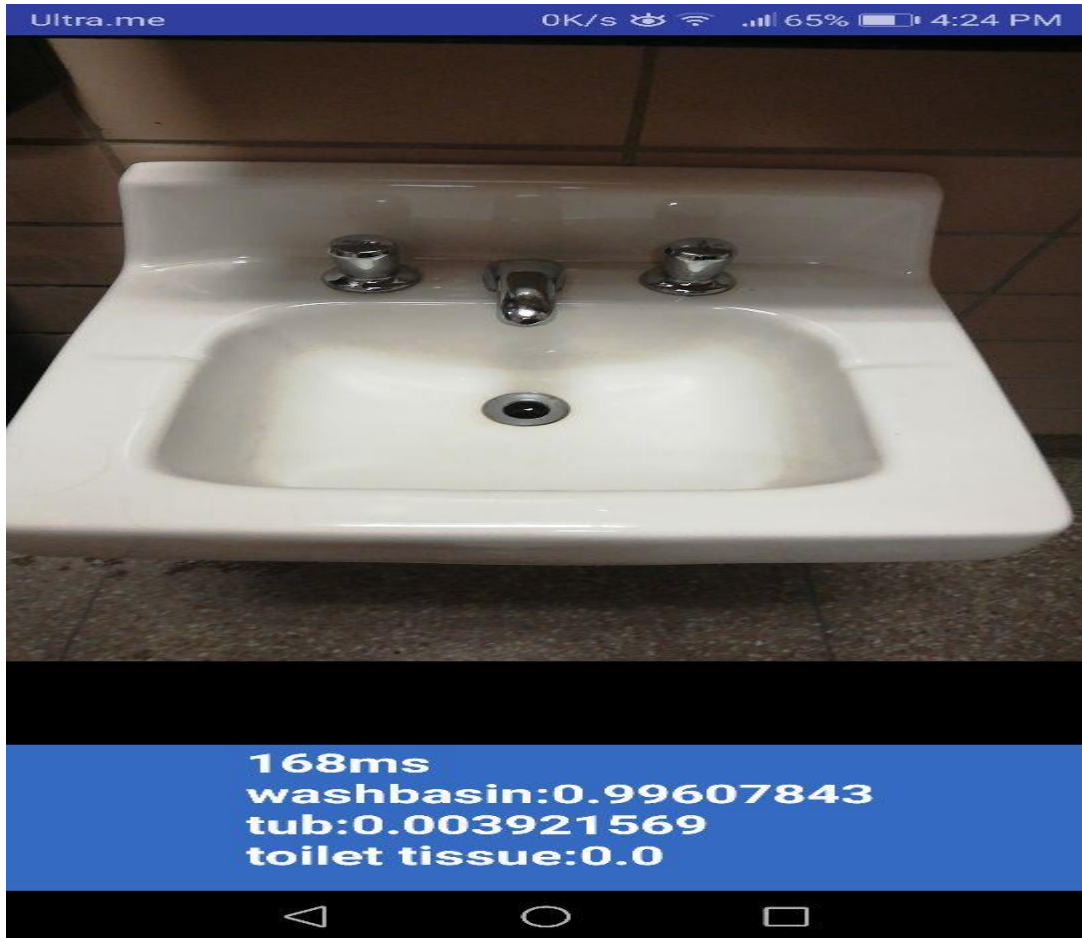


- **Identification using Clarifai:**

- **Bounding box:**

- **Tensor Flow Lite:**
  We have implemented the object prediction on tensor flow lite for our category (interiors)



- **Web application**

  1. The home page of the website looks this where any user can navigate and see different tabs associated. The Main Page talks about the project and the other tabs are self-explanatory.
  2. The image classification tab talks about the different results obtained.
  3. The image prediction tab shows the predicted answer for the pic which is chosen.
  4. In image prediction we have three methods to do it namely Spark API, Clarifai API and Tensor Flow.
  5. Integration of the website using the python flask library.

# Home



# Image Classification

Visual Question Answering



Upload the image:
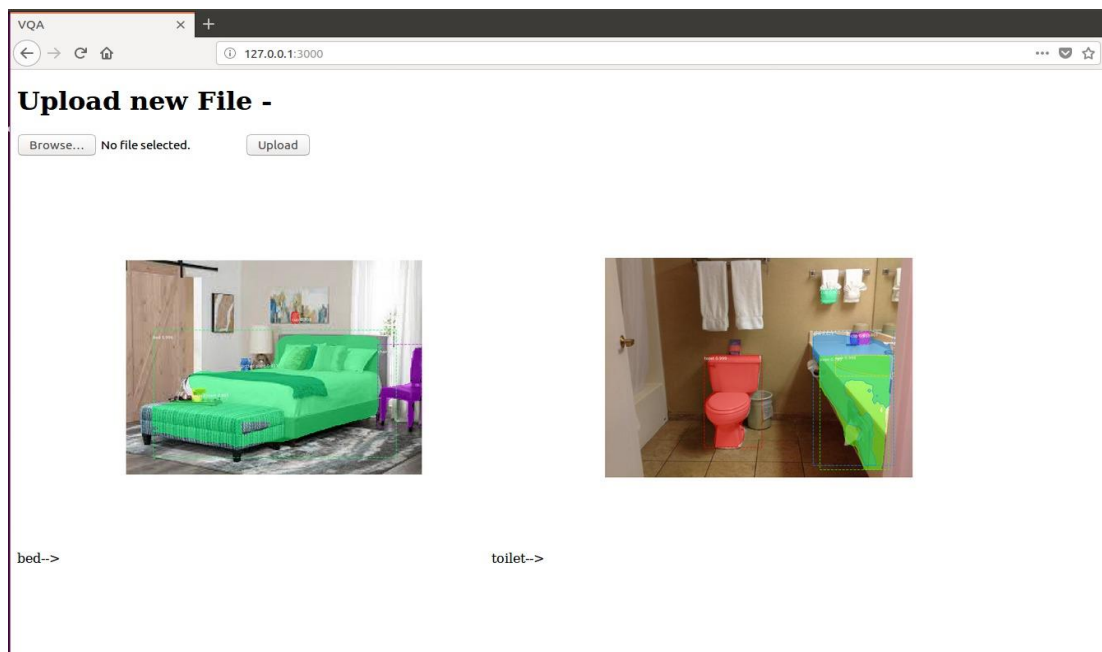
## 5.  Results:

### Data Source:
We have downloaded the dataset from http://mscoco.org/dataset/#download.
COCO is a large-scale object detection and segmentation dataset.
The following table describes the data:

|  | Training Set | Validation Set | Testing Set |
|---|---|---|---|
| **Images** | 82,783 | 40,504 | 81,434 |
| **Questions** | 2443,757 | 214,354 | 447,793 |
| **Answers** | 4,437,570 | 2,143,540 | - |

Segregate Interiors (bathroom, bedroom, kitchen, living room) from MSCOCO data, 50 training images of each category and 5 test images.

## 6. Conclusion:

We created a system which gives yes/no answer for the question asked based on the image. This system takes a picture as input form the user through the website and processes it. When the user asks a question, then the system answers it in yes/no. Implementation of more complex answers and questions would take GPUs for processing. The outputs of different models have been compared and presented during this project. This project could work as a base for more complex projects.

## 7. Appendix:

1. Introduction
2. Related Work
3. Approach
4. Application specification & Implementation
5. Results
6. Conclusion
7. Appendix
8. References
9. Project Management

## 8. References:

1. https://github.com/matterport/Mask_RCNN

2.  http://flask.pocoo.org/

3. http://www.visualqa.org/

4. https://cs.stanford.edu/people/karpathy/rcnn/

## 9. Project Management:

1. **Timelines:**
   Final Project Video/PPT: 4/30/2018
   Project Demo: 5/3/2018
   Final Project Package: 5/7/2018

2. **Team Members:**
   Santosh Yadav Jada
   Harshil Lavjibhai Patel
   Trinadha Rajeswari Muppala
   Avni Mehta