# Visual Question Answering
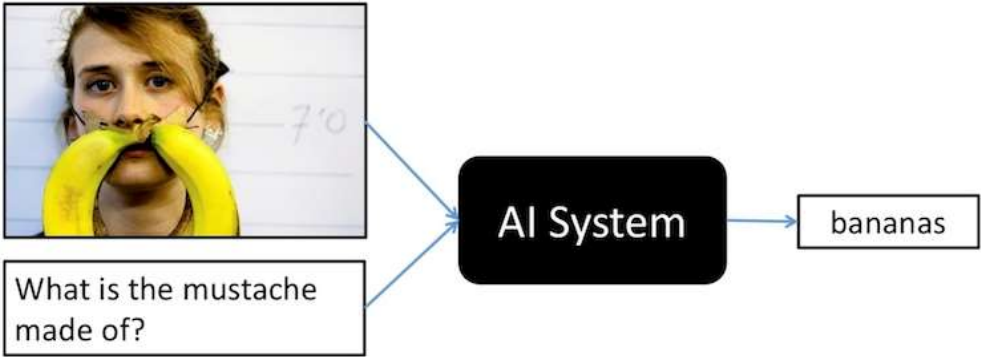
Qilong Wang, Yi Yang

Team#A2_16

## Purpose

This project aims to use artificial intelligence to help blind and other visually impaired groups deal with some of the problems they encounter in life.

The VQA task is to provide an accurate natural language answer to the given question and picture.



## Result

The model performs very well on the VQA v2.0 validation set. The result is shown below.



| Question type | Accuracy |
|---|---|
| Other | 73.90 |
| Yes/no | 95.69 |
| Count | 67.26 |
| Overall | 81.22 |

## Demo

We also make a simple UI for the project which can output an answer according to the question and picture given by the user.



How many trains are in the picture
Pred: 2, Ans: 2

How many animals are there
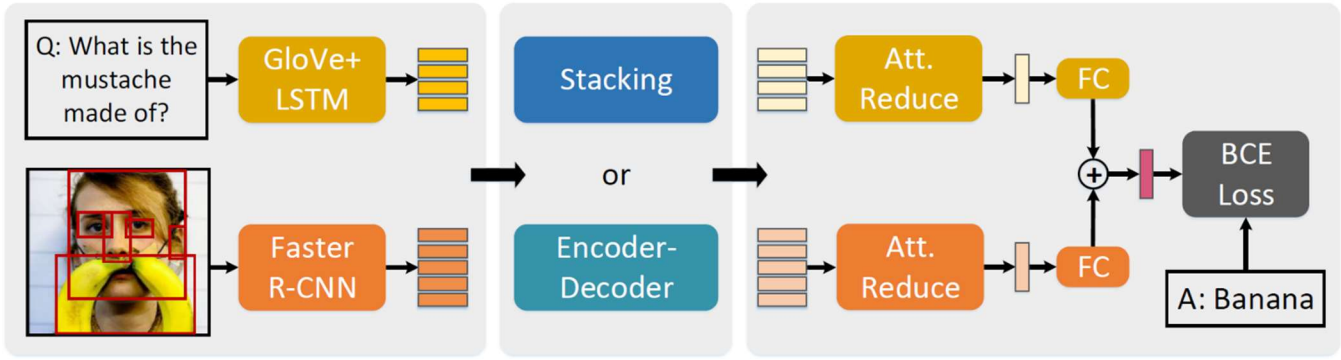Pred: 8, Ans: 8

## Model

The attention mechanism has shown very good results so far, and collaborative attention has also performed well in image and text representations.

BAN and DCN can solve the problem that co-attention cannot fully interact between multiple modalities.

Two general attention units: a self-attention (SA) unit for intramodal interaction and a guided-attention (GA) unit for modality interaction between states.

Connecting multiple module layers in series to form a MCAN network (Modular Co-Attention Network).



## Dataset

a new dataset containing open-ended questions about images. These questions require an understanding of vision, language and commonsense knowledge to answer.

## Future Work

Speech-to-text and text-to speech function: It would make the blind people use it conveniently.

Advice system: Give user some demands of picture so that the model can recognize pictures accurately.