

# Machine Learning II

## — Visual Question Answer

### Group 7

Sipeng Wang

Jiafang Liu

Kangning Gao

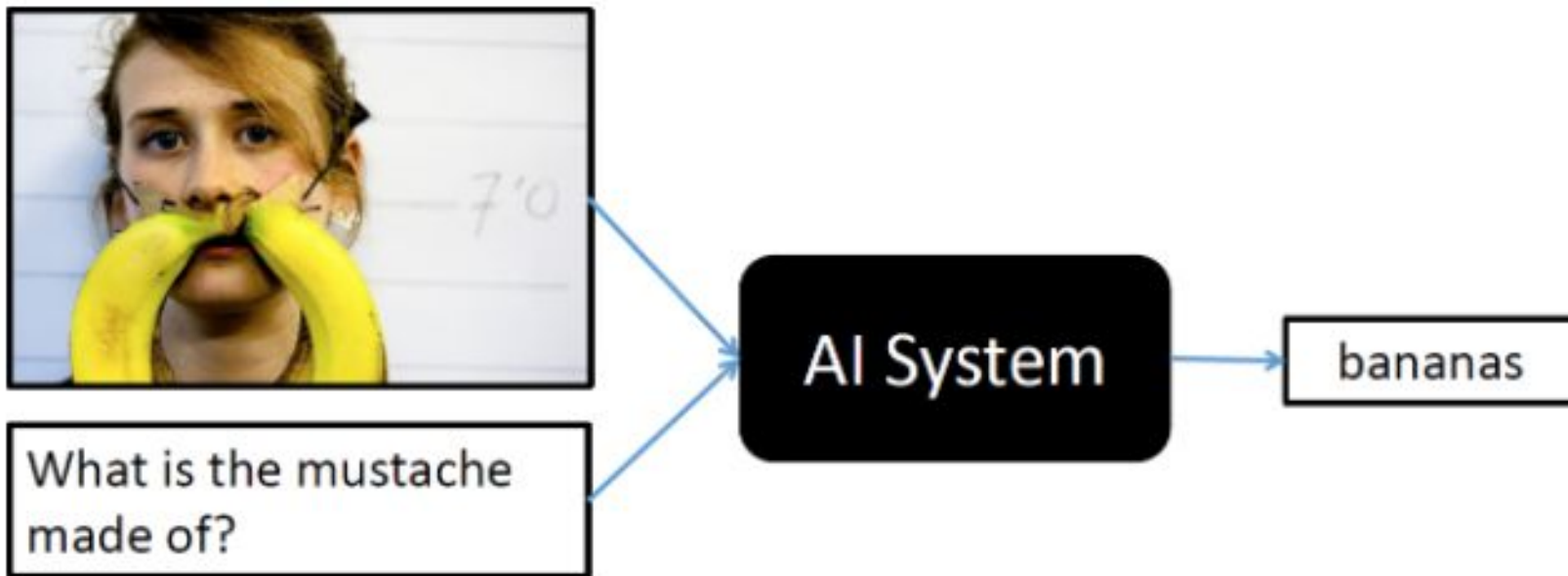
12/11/2017



# Outline

- Introduction
- Data Process
- VAQ Model
- Results
- Conclusion

# Introduction



# View Dataset - MSCOCO

## Image Set

204,721 Image quantities

Input: 4096 dimensions feature



## Question Set

1,105,904 questions

Input: 384 dimensions feature

```
[sipeng_wang@ml-class-ubuntu16:~/final/data/preprocessed]$ cat questions_val2014.txt | head -10
Where is he looking?
What are the people in the background doing?
What is he on top of?
What website copyrighted the picture?
Is this a creamy soup?
Is this rice noodle soup?
What is to the right of the soup?
What is the man doing in the street?
How many photo's can you see?
What does the truck on the left sell?
```

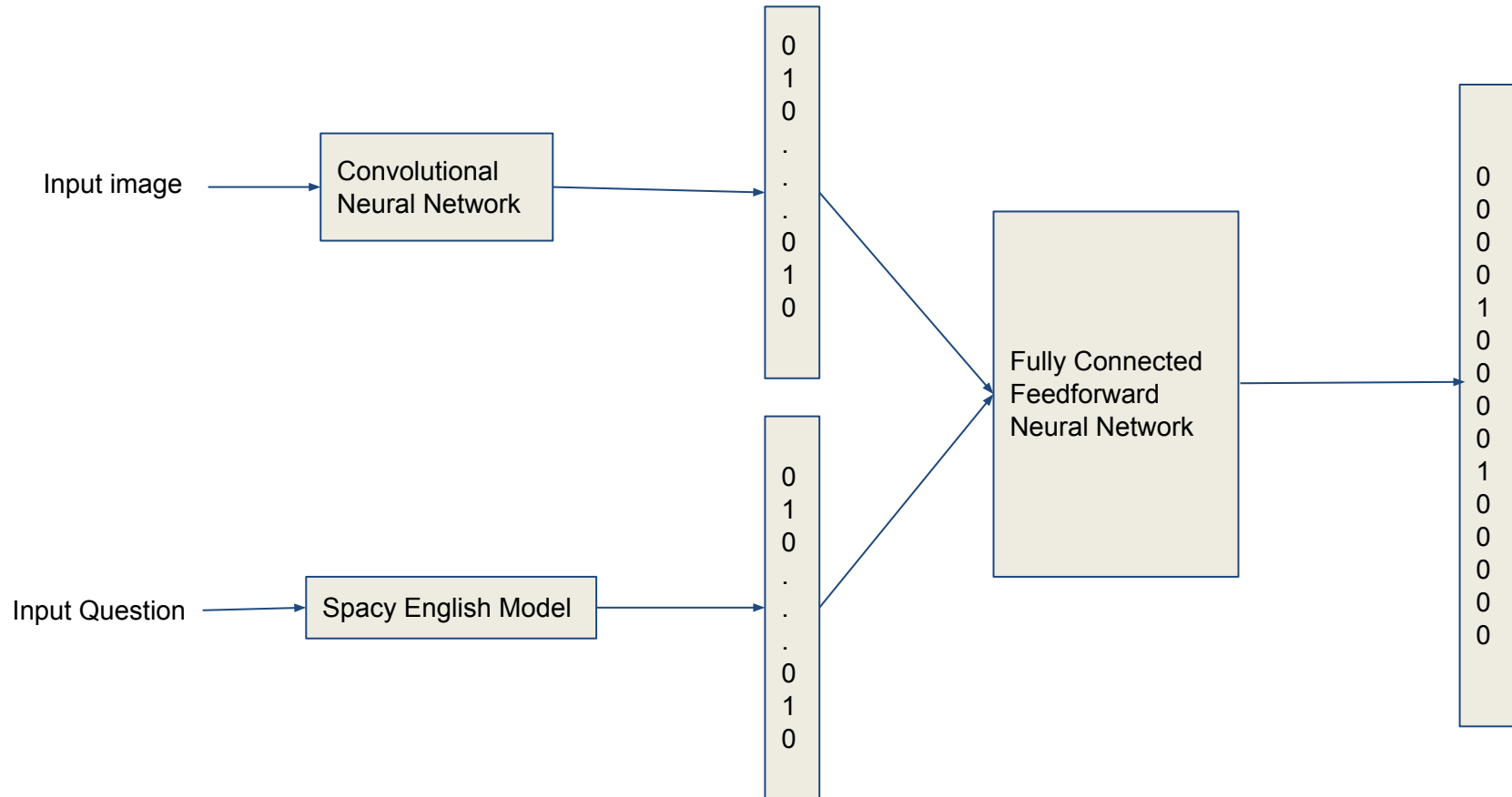
## Answer Set

Top 1000 answers from

11,059,040 ground truth answer

```
[sipeng_wang@ml-class-ubuntu16:~/final/data/preprocessed]$ cat answers_val2014_model.txt | head -10
down
watching
picnic table
foodiebakercom
no
yes
chopsticks
walking
1
ice cream
```

# VQA Model

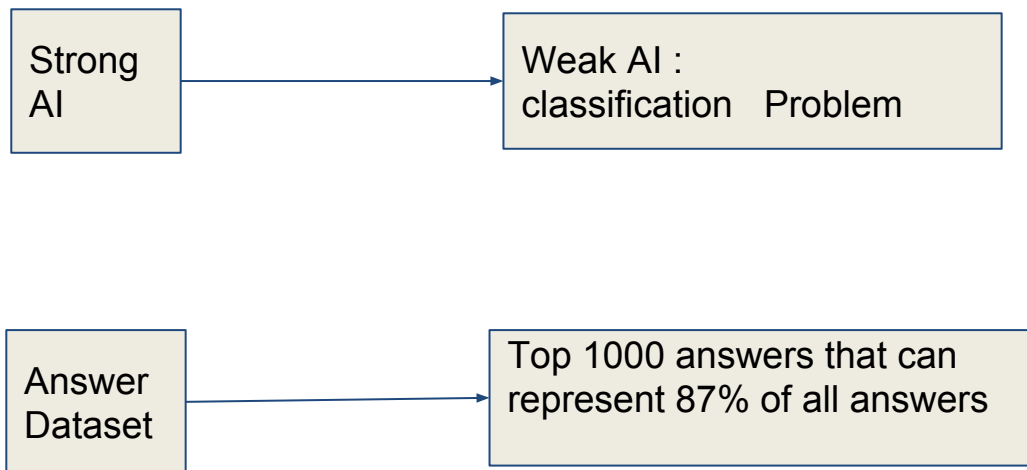


# Data Pre-process

- Extract data
- Generate Answer
- Categorize answer data

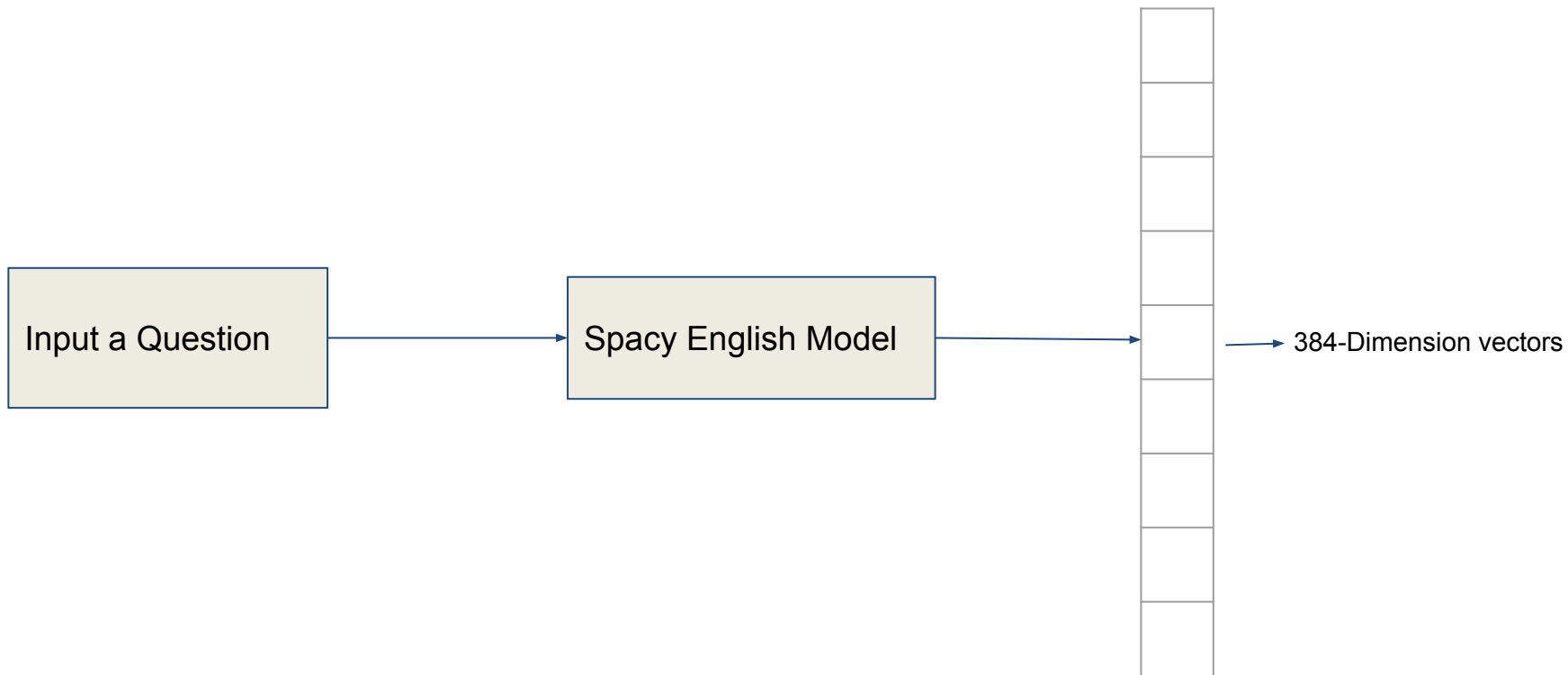
# VQA Model

- Generate Answers



# VQA Model

- Flatten Questions





# VQA Model

- Flatten Pictures

Input picture



CNN



4096-dimension Vectors

# VQA Model

## VGG 16

```
def VGG_16(weights_path=None):
    model = Sequential()
    model.add(ZeroPadding2D((1,1),input_shape=(3,224,224)))
    model.add(Convolution2D(64, 3, 3, activation='relu'))
    model.add(ZeroPadding2D((1,1)))
    model.add(Convolution2D(64, 3, 3, activation='relu'))
    model.add(MaxPooling2D((2,2), strides=(2,2)))

    model.add(ZeroPadding2D((1,1)))
    model.add(Convolution2D(128, 3, 3, activation='relu'))
    model.add(ZeroPadding2D((1,1)))
    model.add(Convolution2D(128, 3, 3, activation='relu'))
    model.add(MaxPooling2D((2,2), strides=(2,2)))

    model.add(ZeroPadding2D((1,1)))
    model.add(Convolution2D(256, 3, 3, activation='relu'))
    model.add(ZeroPadding2D((1,1)))
    model.add(Convolution2D(256, 3, 3, activation='relu'))
    model.add(ZeroPadding2D((1,1)))
    model.add(Convolution2D(256, 3, 3, activation='relu'))
    model.add(MaxPooling2D((2,2), strides=(2,2)))

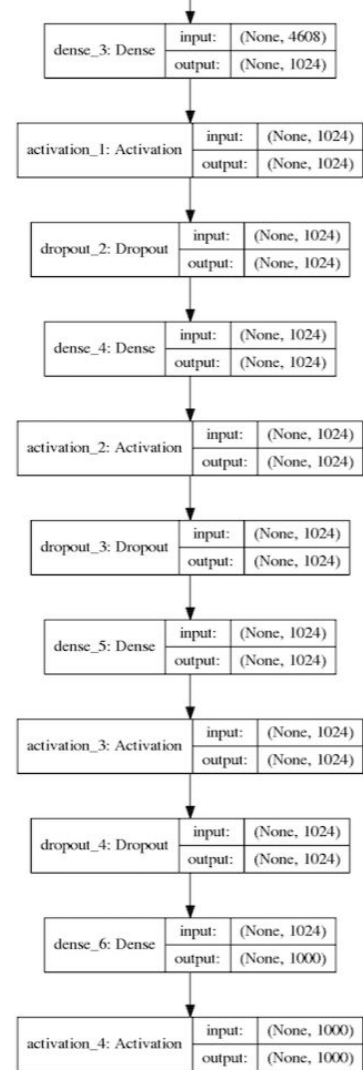
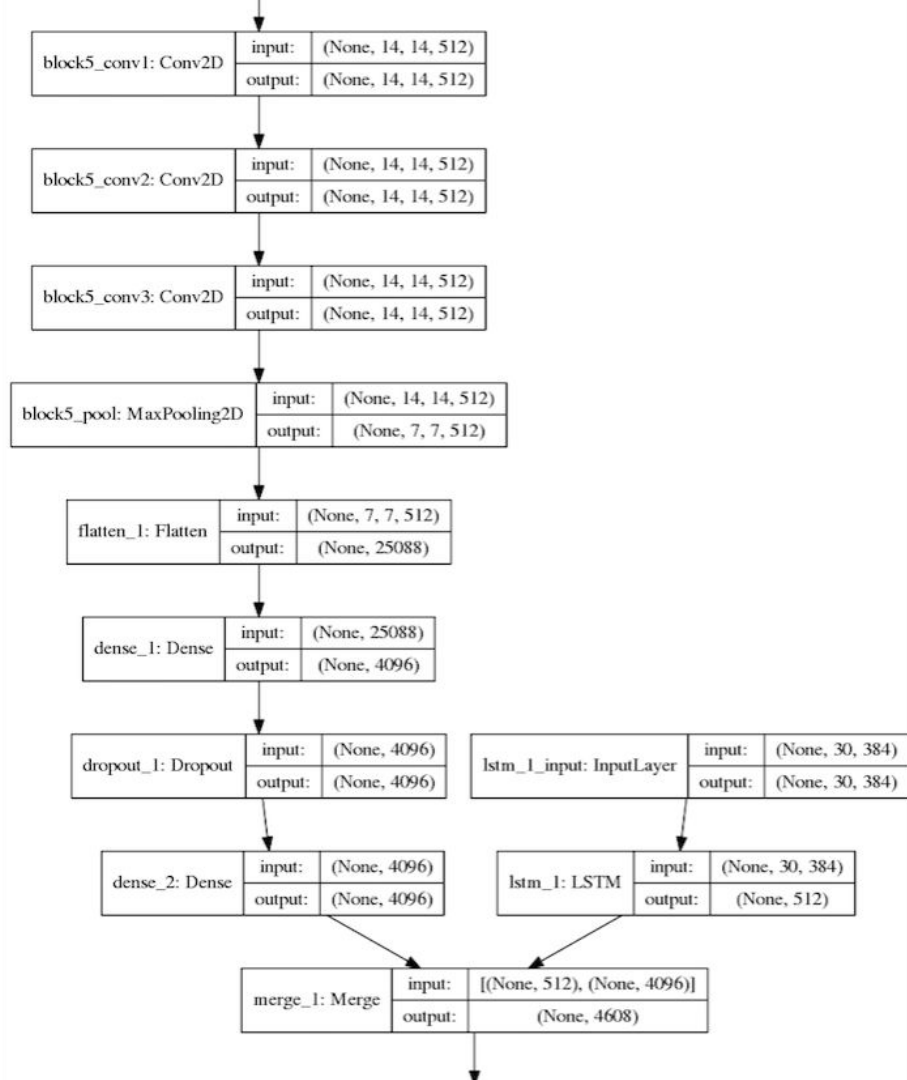
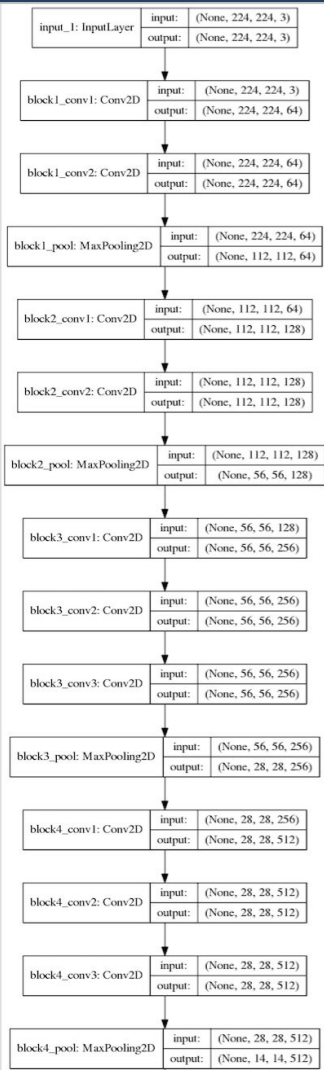
    model.add(ZeroPadding2D((1,1)))
    model.add(Convolution2D(512, 3, 3, activation='relu'))
    model.add(ZeroPadding2D((1,1)))
    model.add(Convolution2D(512, 3, 3, activation='relu'))
    model.add(ZeroPadding2D((1,1)))
    model.add(Convolution2D(512, 3, 3, activation='relu'))
    model.add(MaxPooling2D((2,2), strides=(2,2)))

    model.add(ZeroPadding2D((1,1)))
    model.add(Convolution2D(512, 3, 3, activation='relu'))
    model.add(ZeroPadding2D((1,1)))
    model.add(Convolution2D(512, 3, 3, activation='relu'))
    model.add(ZeroPadding2D((1,1)))
    model.add(Convolution2D(512, 3, 3, activation='relu'))
    model.add(MaxPooling2D((2,2), strides=(2,2)))

    model.add(Flatten())
    model.add(Dense(4096, activation='relu'))
    model.add(Dropout(0.5))
    model.add(Dense(4096, activation='relu'))
    model.add(Dropout(0.5))
    model.add(Dense(1000, activation='softmax'))

    if weights_path:
        model.load_weights(weights_path)

    return model
```



# Neuraltalk2



a man is playing tennis on a tennis court



a train is traveling down the tracks at a train station

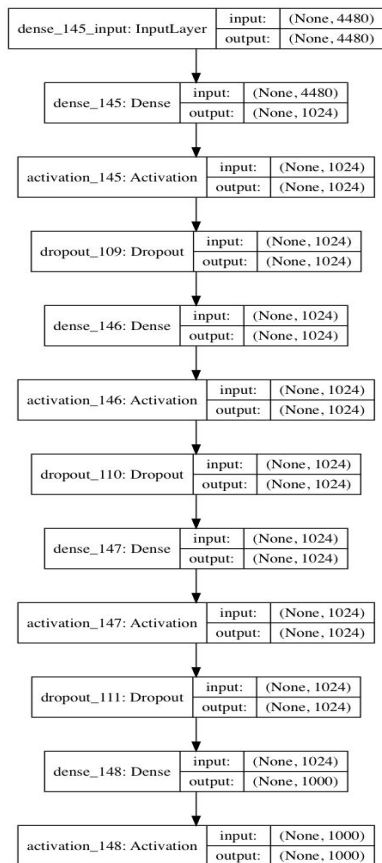


a cake with a slice cut out of it



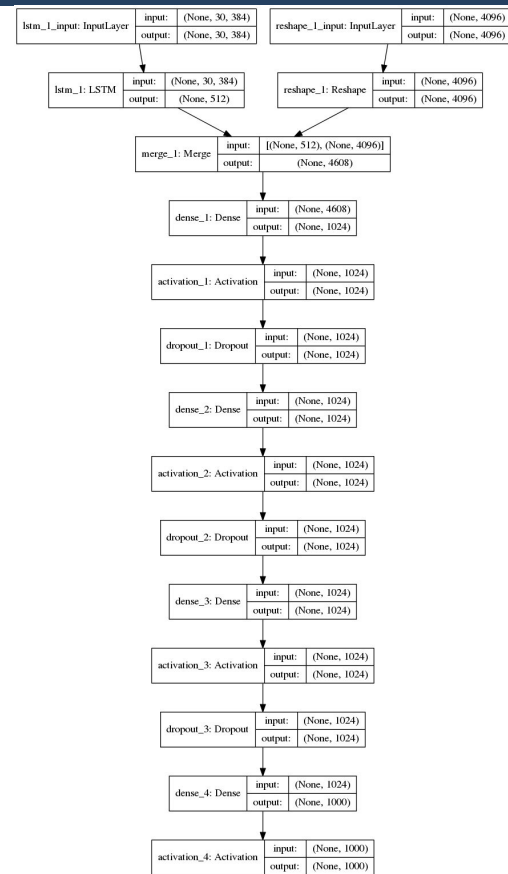
a bench sitting on a patch of grass next to a sidewalk

# Model



- Multi-layer-Perceptron(MLP)

- Long-Short-Term-Memory(LSTM)



# Training

## MLP

Epochs number	1	5	10	...	20
Time(min)	36.5	182.5	420	...	1680
Loss(lr = 0.01)	5.25	3.89	3.08	...	3.07

## LSTM

Epochs Number	1	5	10	...	20
Time(min)	63	315	630	...	2520
Loss(lr = 0.01)	3.52	3.01	2.90	...	3.02



# Accuracy

Our MLP:24.0%

Our LSTM:36.5%

2017 challenge 1st: 69%

2017 challenge 27th: 37.3%



what color is that train?

Loaded

[/Users/zack\\_wang/a](#)

UserWarning)

42.28 % blue

20.94 % red

07.67 % orange

07.64 % silver

07.37 % yellow

# Problems & Future work

- The VGG Model training implement was limited by the computational power.
- The loss of both MLP/LSTM model went smooth at around 3 after 10 epochs.