# Visual Question Answering: Experiments with deep learning and text features

Vien Nguyen

Data Science Retreat

*trucvien.nguyen@yahoo.com*

July 16, 2016

# Overview

- Introduction to the task Visual Question Answering (*VQA*)
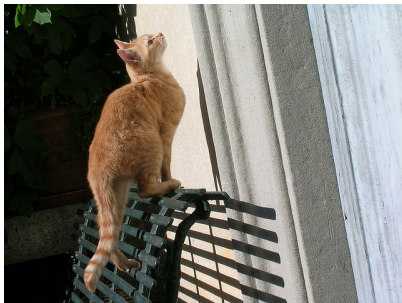- Method and Results
- Conclusion

## Introduction

- Given an image and a question related to this image, the system will automatically learn to generate an answer for this question.
- Use the image to generate visual features with Convolution Neural Network.
- Use the text of the question to generate "bag-of-words" features.
- Use machine learning to learn the answer.

# Examples (1)

Given an image and a question related to this image, the system will automatically learn to generate an answer for this question.
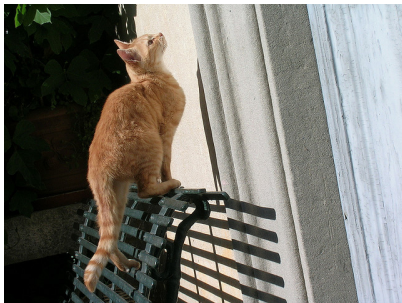
**Image**



**Question**

What animal is in this picture?

# Examples (1)

**Image**



**Question**

What animal is in this picture?

**Answer**

cat

**Image**



**Question**

How many chairs are in this shot?

# Examples (2)

**Image**



**Question**

How many chairs are in this shot?

**Answer**

3

**Image**



**Question**

What season is it?

**Image**



**Question**
What season is it?

**Answer**
summer

## Method

We use a machine learning approach.

- Extract features from the question
- Extract features from the image
- Combine features
- Learn the model

## Dataset

**VQA Visual Question Answering**: http://visualqa.org/

|            | Images | Questions | Answers   |
|------------|--------|-----------|-----------|
| Training   | 82,783 | 248,349   | 2,483,490 |
| Validation | 40,504 | 121,512   | 1,215,120 |

Table: Data statistics

# Pre-processing (1)

- **Part-of-speech tagging**:
  Example:

| What | animal | is | in | this | picture | ? |
|------|--------|-----|-----|------|---------|---|
| WDT | NN | VBZ | IN | DT | NN | . |

Table: POS-tagging example

**Description of pos-tag labels**

WDT: Wh-determiner

NN: Noun, singular or mass

VBZ: Verb, 3rd person singular present

IN: Preposition or subordinating conjunction

DT: Determiner

# Pre-processing (2)

- **Question**: Extract n-gram features of words ($n \leq 2$). Example: What animal is in this picture?
  **unigram ($n = 1$)**: What, animal, is, in, this, picture, ?
  **bigrams ($n = 2$)**: What animal, animal is, is in, in this, ...
- **Question**: Extract n-gram features of pos-tags ($n \leq 2$). Example: WDT NN VBZ IN DT NN .
  **unigram**: WDT, NN, VBZ, IN, DT, NN, .
  **bigrams**: WDT NN, NN VBZ, VBZ IN, IN DT, DT NN, NN .
- **Image**: Resize to 64*64.

# Experiments

- Random Forest Classifier on n-gram features.
- Deep Learning on image features.
- Apply convolutional neural networks (CNN) layer on images. Combine features from images and text. Apply another layer on the combined features.

**Random Forest**

- Based on Decision Trees.
- From multiple trees, select the one with highest frequency.

**Deep Learning**

- A neural network with many layers ($\geq 3$), including CNN layer(s).
- Convolutional neural networks: not use pre-defined funtion like a normal neural network, but instead learn a function from the data.

# Technologies

- **Amazon Web Service**: Ubuntu Server 14.04 LTS (HVM), 64-bit, GPU g2.8xlarge.
- 25 minutes for pre-processing image features, 10 minutes for pre-processing text features.
- 10 hours training on 50% images of the training set.
- 20 minutes training on 70% questions of the training set.
- **Softwares**: Python, NLTK, Theano, scikit-learn.

# Results (1)

| | |
|---|---|
| Image features | 28.38 (50% data) |
| Text features | 42.46 (70% data) |

Table: Results

# Results (2)

| Per Question Type Accuracy | |
| --- | --- |
| are these | 67.56 |
| is there a | 85.01 |
| how many | 33.12 |
| what animal is | 17.17 |
| what | 17.19 |
| does the | 75.50 |
| could | 87.36 |

Table: Results Per Question Type

# Results (3)

| Per Answer Type Accuracy | |
|---|---|
| number | 27.81 |
| other | 21.66 |
| yes/no | 75.16 |

Table: Results Per Answer Type

# Examples (1)

**Image**



**Question**
What animal is in this picture?

**Question type**
what animal is

**Answer type**
other

**Answer**
cat

# Examples (2)

**Image**



**Question**

How many chairs are in this shot?

**Question type**

how many

**Answer type**

number

**Answer**

3

**Image**



**Question**
What season is it?

**Question type**
what

**Answer type**
other

**Answer**
summer

**Image**



**Question**

Does the weather appear rainy?

**Question type**

does the

**Answer type**

yes/no

**Answer**

yes

**Image**

**Question**

Could the items in this picture be used for sewing?

**Question type**

could

**Answer type**

yes/no

**Answer**

yes

# Results per question type
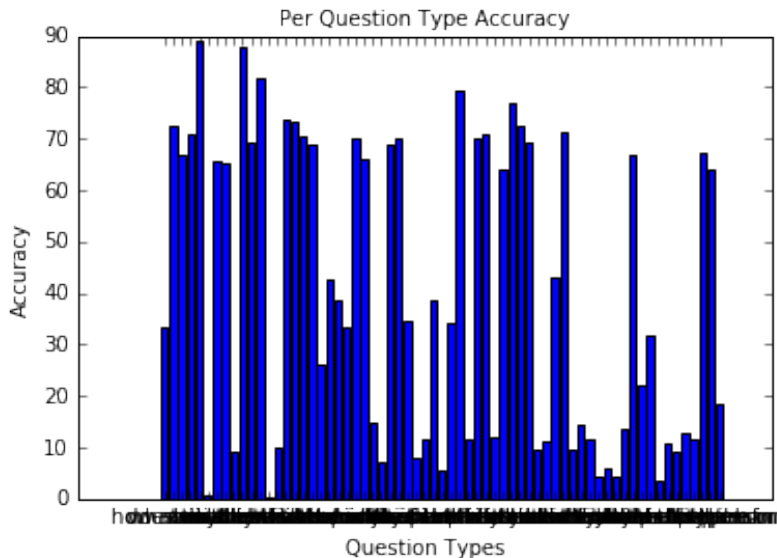


Per Question Type Accuracy

# Conclusion

- Visual Question Answering is a new research direction.
- Results depend on both visual and textual features.
- Requiring techniques in computer vision, language, integrating vision + language.

The deep learning framework here is based on the tutorial of Michael Nielsen: http://neuralnetworksanddeeplearning.com

# Thank you!