# UCF

## COLLEGE OF ENGINEERING AND COMPUTER SCIENCE

**A PROJECT PROPOSAL TO MODEL A QUESTION ANSWERING SYSTEM USING TENSORFLOW**

- TEAM 12(VHAIJAYANTHISHREE,NEETHIKA, SRIDEVI )

## 1. PROBLEM DESCRIPTION:

A Question-Answering (QA) system is a system that is designed to automatically answer questions posed by humans in a natural language. Question-Answering has received more focus in the recent times as large search engines have basically mastered general information retrieval and are starting to cover more edge cases.

Standard information retrieval statistical models like Latent Dirichlet allocation (LDA) or Latent Semantic Indexing (LSI) do not efficiently capture syntactic nuances. Deep learning models are better suited to model QA system given their ability to capture higher-order syntax and efficiency. This project focuses on improving the existing QA system based out of Deep Neural Networks, by exploring supervised learning techniques and by using more specific word vectorizations schemes.

## 2. RELATED WORK:

1) "Applying deep learning to answer selection: a study and an open task" (Feng et al., 2015).
2) "LSTM-based deep learning models for non-factoid answer selection" (Tan et al., 2016).

The above papers are recent examples that have applied deep learning to Question-Answering tasks yielding good results.

## 3. PROPOSED EXPERIMENT:

In this experiment, we intend to implement and enhance existing Question-Answering system using TensorFlow, a highly sophisticated framework in Python developed by Google for Machine Learning.

### A. DATA SET:

**bAbI:** It is a set of 20 QA tasks, each consisting of several context-question-answer triplets, prepared and released by Facebook. Each task aims to test a unique aspect of reasoning and is, therefore, geared towards testing a specific capability of QA learning models.The bAbI dataset is composed of synthetically generated stories about activity in a simulated world. (10,000 training examples). If time permits we plan to test the

existing system on MCTest dataset of Microsoft to test its ability to predict the Multiple choice questions.

## B. EXTERNAL LIBRARIES AND TOOLS:

- **NLTK** - The most famous Natural Language Processing toolkit for Python. Using this toolkit, words, lines, paragraphs, texts and many more entities can be analyzed with ease

- **TensorFlow** - A sophisticated python framework for Machine Learning which focuses predominantly on Deep Neural Network models. TensorFlow provides many APIs.

- **Stanford GloVe** - A deep learning model  for obtaining vector representations for words. This model predicts contextually similar words thus in QA system query can be expanded to word vectors of close proximity facilitating efficient query expansion.

## C. SOFTWARE DEVELOPMENT:

- **Programming Language:** Python

- **External Libraries:**  Keras, Numpy, NLTK, JSON, re, Matplotlib

## D. EXPERIMENTAL DESIGN:

a. Parse JSON and tokenize the text.

b. Create word embeddings using GloVe (linguistic or semantic similarity of the corresponding words).

c. Encode the input to fit neural network model.

d. Construct model using Keras library and train it.

e. Test the model.

**Example 1:**

**TEXT:** bill grabbed the apple there . bill got the football there . jeff journeyed to the bathroom . bill handed the apple to jeff . jeff handed the apple to bill . bill handed the apple to jeff . jeff handed the apple to bill . bill handed the apple to jeff .

**QUESTION:**  what did bill give to jeff ?

**RESPONSE:**  apple (Correct)

**EXPECTED:**  apple

**Example 2:**

**TEXT:** bill moved to the bathroom . mary went to the garden . mary picked up the apple there . bill moved to the kitchen . mary left the apple there . jeff got the football there . jeff went back to the kitchen . jeff gave the football to fred .

**QUESTION:** what did bill give to fred ?

**RESPONSE:** apple (Incorrect)

**EXPECTED:** football

## 4. DIVISION OF LABOR:

### Week 1:

**Vhaijaiyanthi:** Overview of TensorFlow - Aaron Schumacher's "Hello TensorFlow".

Installing TensorFlow, Numpy, Matplotlib, NLTK.

**Neethika:** Installing Keras from binary and understanding the language models in it.

**Divya:** Reviewing different tasks of bAbI

### Week 2:

**Divya:** Parsing Glove dataset and handling unknown words

**Neethika:** Converting the sentence to sequence and contextualizing based on known and unknown words.

**Vhaijaiyanthi:** Constructing the network by defining the hyperparameters

### Week 3:

**Neethika:** Structuring using Dynamic Memory Network

**Vhaijaiyanthi:** Working on the 4 different modules of the Dynamic Memory Network

**Divya:** Exploring supervised learning for existing unsupervised implementation by using loss function, and training the model.

### Week 4:

**Vhaijaiyanthi:** Validation - Training the network

**Neethika:** View training progress and get training speed metrics

**Divya:** Testing the model on different Task sets provided in the dataset and concluding results.

### Week 5:

This week is dedicated to working on the project report.

**5. <u>REFERENCES</u>:**

- https://codekansas.github.io/blog/2016/language.html

- https://en.wikipedia.org/wiki/Question_answering

- https://web.stanford.edu/~jurafsky/slp3/28.pdf

- https://www.oreilly.com/ideas/question-answering-with-tensorflow

- https://nlp.stanford.edu/projects/glove/