**Question/Answering System for Diseases**

**Team 12**

❖ Gupta, Arunit – 11
❖ Patel, Marmikkumar Navinchandra – 31
❖ McDuff, Luke Joseph – 21
❖ Ejjirothu, Manoj Prabhakar – 7

# Table of Contents

# 1. Introduction

## Question and Answering system for Diseases

The QA system is designed for users around the globe. Users can access the web page and ask their queries regarding health and diseases. The system is designed for a user to have the ability to ask a question on the webpage in a textbox and retrieve the answer for the same in an organized manner.

The system uses a corpus of unstructured information across the web and live data from twitter tweets for the knowledge base and processes the data to return results.

Input is provided by the user on the webpage and output is displayed on the same page by processing data from the knowledge base. It uses natural language processing to lemmatize the data, divide into Parts of Speech, parsing, and creating named entity relationships. The system uses machine learning principles and ontologies to process data and come up with a correct answer or closest match.

The system is useful for people to learn about a particular disease in a precise and organized format.  If we have any problem or query we usually search on the internet. We come up with multiple links where we get scattered information and need to search across links to get the answer; this system has the ability to analyze and interpret the question, and get the answer from our knowledge base after applying processing techniques. The system then returns the answer in an organized and specific format with graphical and image representation which is useful for naïve user.

# 2. Project Goals and Objectives

## 2.1 Motivation

People always have medical queries and usually they use search engines to get answers for their question, but all information is scattered on the web, we have a huge corpus of information but not in a structured manner. This system is designed for the users where they can ask their queries and get response in a precise and structured manner.

## 2.2 Objective

To design a system which can take a question from a user and give an answer as output by applying natural language and machine learning principles and display the information in a structured format.

## 2.3 Expected Outcomes

A webpage accessible by users across the world having the ability to ask a question, and get a response on the webpage in a structured format.

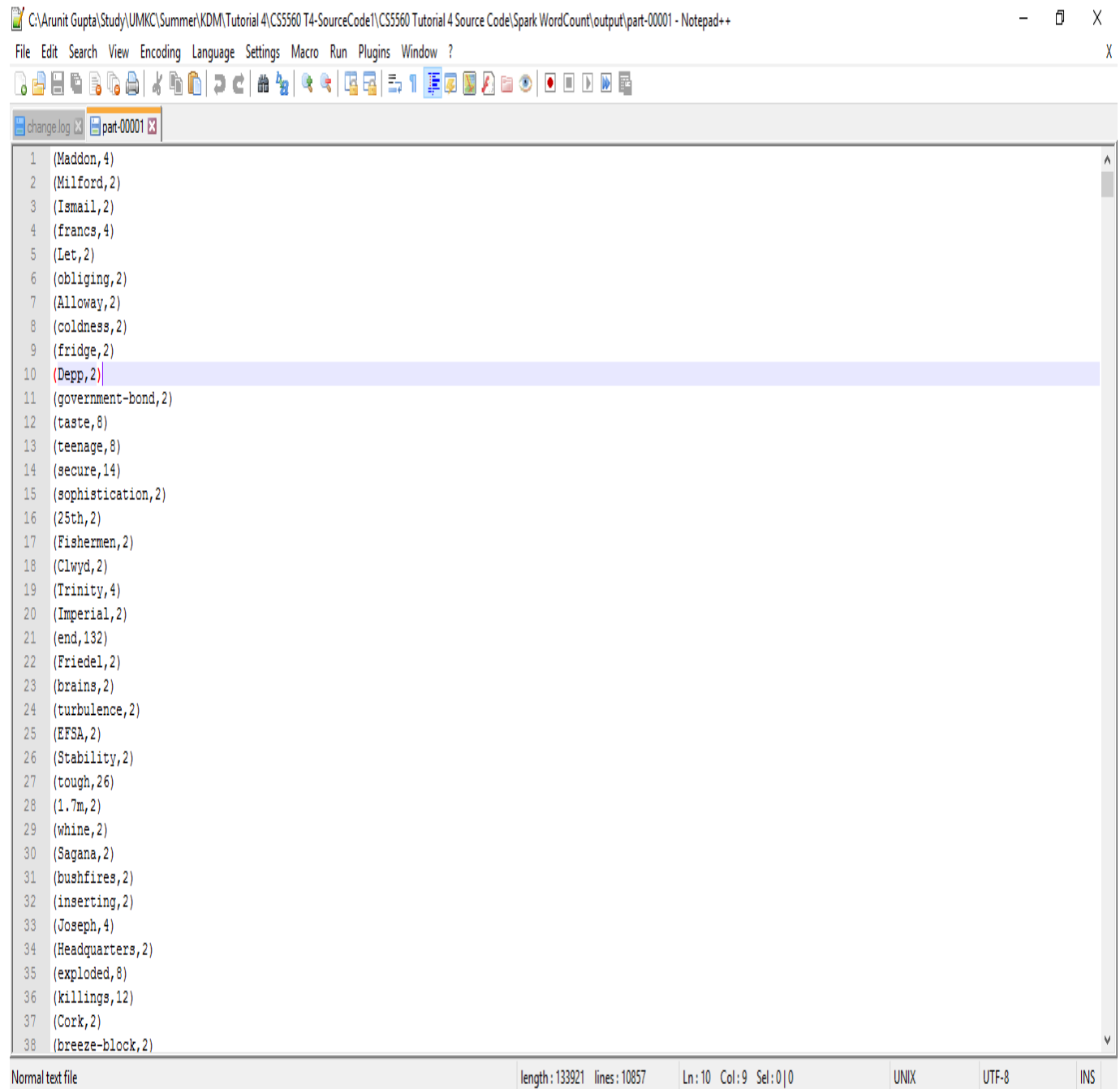# 3. Project Domain & Datasets

## 3.1 Project Domain

Healthcare - Consisting of Diseases, health issues, drug information, trending viruses, allergies etc.

## 3.2 Dataset Links

❖ https://dev.twitter.com/overview/api  -  Collected 200K tweets on health and Diseases.

❖ http://dumps.wikimedia.org/enwiki/latest/enwiki-latest-pages-articles.xml.bz2

❖ http://www.statmt.org/lm-benchmark/1-billion-word-language-modeling-benchmark-r13output.tar.gz

❖ http://ebiquity.umbc.edu/redirect/to/resource/id/351/UMBC-webbase-corpus

❖ http://blog.wolframalpha.com/category/health-med/

❖ http://archive.ics.uci.edu/ml/datasets/Hepatitis

# 4. Features Implementation

## 4.1 Word Count

## 4.2 NLP Processing

### Lemmatization

## POS

## 4.3 Information Retrieval/Extraction Technologies

# 5. Implementation Specification

## 5.1 Software Architecture Diagram

## 5.2 Class Diagram

| Data |
| --- |
| |
| DataLoader() |

| lemmatiziation |
| --- |
| definition: String[] |
| |

| Documment Procsessing |
| --- |
| |
| DocumentExtraction()<br>DocumentConvertion() |

| POS tagging |
| --- |
| name : String[] |
| |

| Twitter Data |
| --- |
| |
| StreamTwitterData() |

| DomainClassification |
| --- |
| |
| FilterWords() |

| FeatureVector Encoder |
| --- |
| |
| WordValueEncoder() |

| Home |
| --- |
| |
| InputDomain()<br>InputQuestions() |

| Machine Processing |
| --- |
| |
| TrainingData()<br>AlgorithmLearning()<br>Optimization()<br>FinalModelCalssiication() |

| Output |
| --- |
| |
| DisplayText() |

| Output |
| --- |
| |
| DisplayGraph() |

| Output |
| --- |
| |
| DispalyImage() |

## 5.3 Sequence Diagram

## 5.4 Workflow



## 5.5 Existing services/features used

AWS, Alchemy API

# 6. Project Management

## 6.1 Contribution of each member

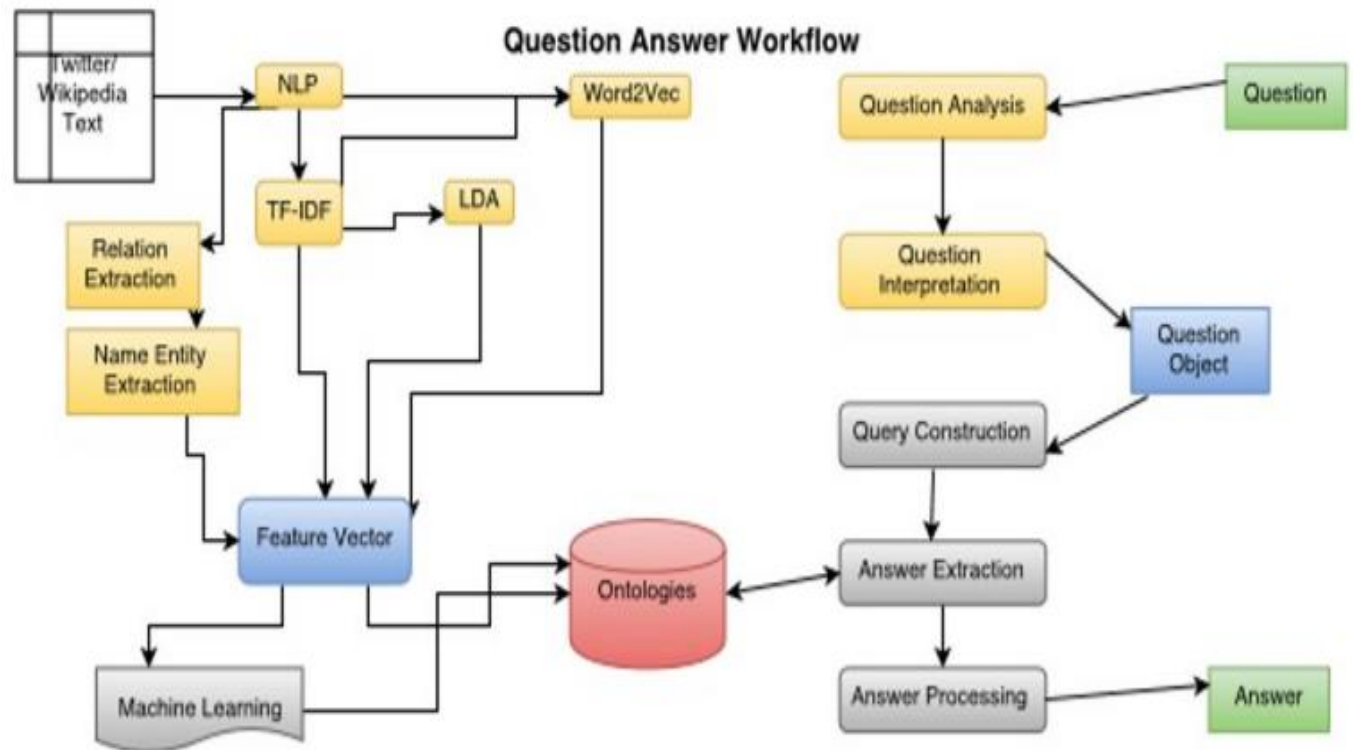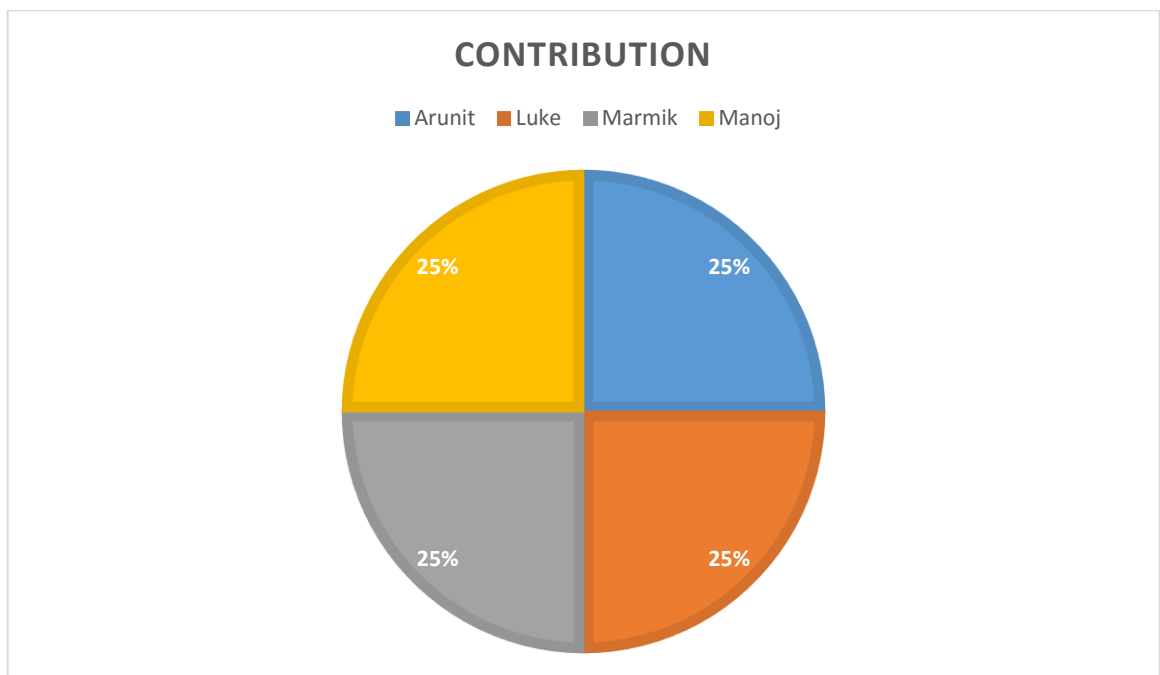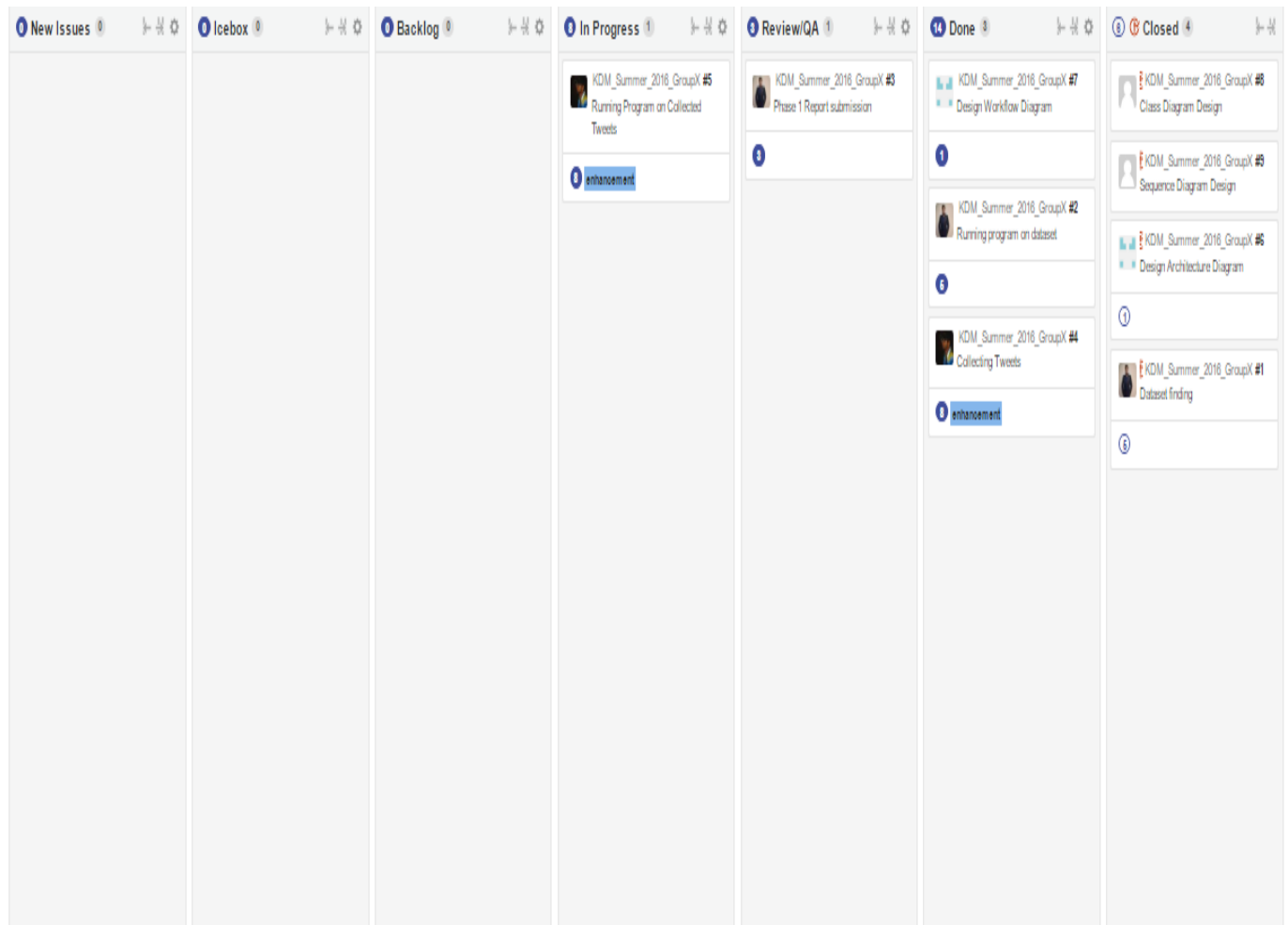- ❖ **Gupta, Arunit** – Dataset Findings, Running dataset on Program, Report compilation
- ❖ **Patel, Marmikkumar Navinchandra** – collection of tweets, running program on tweets collected
- ❖ **McDuff, Luke Joseph** – Architecture Diagram design and Workflow design
- ❖ **Ejjirothu, Manoj Prabhakar** – Class Diagram and Sequence Diagram

**CONTRIBUTION**

■ Arunit  ■ Luke  ■ Marmik  ■ Manoj

25%  25%  25%  25%

## 6.2 Zen hub and GitHub URL/Statistics

GitHub link: https://github.com/ljm7b2/KDM_Summer_2016_GroupX

## Zen hub

## Increment 1



### 6.3 Concerns/Issues

Relating all the technologies and merging them to get final output

### 6.4 Future Work

Implementing the datasets and live data on machine learning and ontologies programs, implementing voice to text search

# 7. References

- https://github.com/DiseaseOntology/HumanDiseaseOntology
- In Class  Tutorials
- https://dev.twitter.com/overview/api
- http://dumps.wikimedia.org/enwiki/latest/enwiki-latest-pages-articles.xml.bz2
- http://www.statmt.org/lm-benchmark/1-billion-word-language-modeling-benchmark-r13output.tar.gz
- http://ebiquity.umbc.edu/redirect/to/resource/id/351/UMBC-webbase-corpus
- http://blog.wolframalpha.com/category/health-med/
- http://archive.ics.uci.edu/ml/datasets/Hepatitis