

Question Answer System for Twitter and Wikipedia on US Election 2016

Project Report Phase 3

Submission Date: 07/25/2016

by

Ganesh Taduri (#40)

Swatvik Gunamaneni (#10)

Venkata Sasidhar kanumuri(#13)

Tejo Kumar Manchu (#19)

Motivation: There are many advantages in extracting the information from the slew of sources available on the internet. Not a decade ago, we realized that we can process this large amount of data and extract meaningful information from it. There is a growing importance on the development of information technologies that are capable of accessing and analyzing the data available via social media. Twitter is one such source of large amount of raw data which contains worldwide news, public opinion on several events which can be processed and relevant information can be extracted. On average, every second around 6000 tweets are tweeted on twitter. This data can be used to build a question answer system

Objectives:

The objective of the project is to learn machine learning algorithms, natural language processing through implementation of available services and APIs like Natural language classifier, Speech to text, speech to text etc. We plan to tap the twitter data and develop a question answering system which takes a text question as input and answers it in accordance with the data available on twitter. In addition, we wish to use Wikipedia data to answer questions on static affairs related to US election 2016

Data collection: Our data is related to US election 2016. We collected data from various sources like Wikipedia and twitter using python and their respective API. The twitter data is collected based on #election2016. The Wikipedia data is a text file but the twitter data is extracted in the form of JSON object. We extracted the tweet text from the object. The data is distributed among ten documents depending upon its content. The extracted data from Wikipedia is a meaningful one whereas most of data from twitter is inconsistent with various errors like spelling and grammatical mistakes and unrelated data. We collected data from Wikipedia pages like Donald Trump, Obama, Hillary Clinton etc.

Data collection URL: www.twitter.com & www.wikipedia.com

Input:

The text data we collected from various sources was given as input to the model.

Expected output: The expected output is a Question Answer system, where a user can ask a question in text format and the model answers the question in text format.

Sample Question1:

Who is the republican candidate for US elections 2016?

Sample Answer1:

Donald Trump

Sample Question2:

What is the public opinion on Hillary Clinton winning chances?

Sample Answer2:

She has 60% winning chance

Domain: Question answering system

Class diagram:

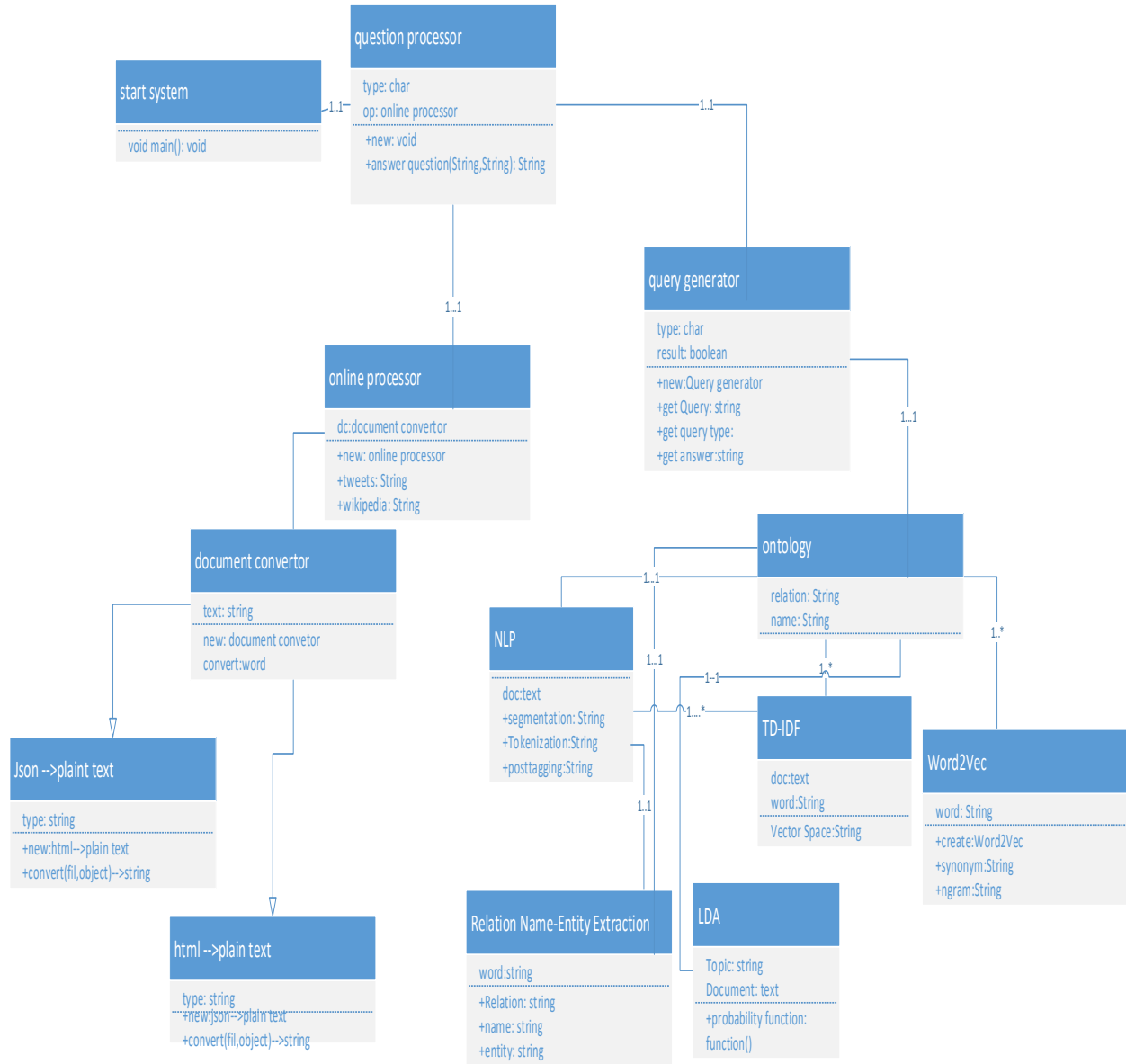
Class diagrams shows the classes in the system and their interrelationships between them

Classes involved in this system are

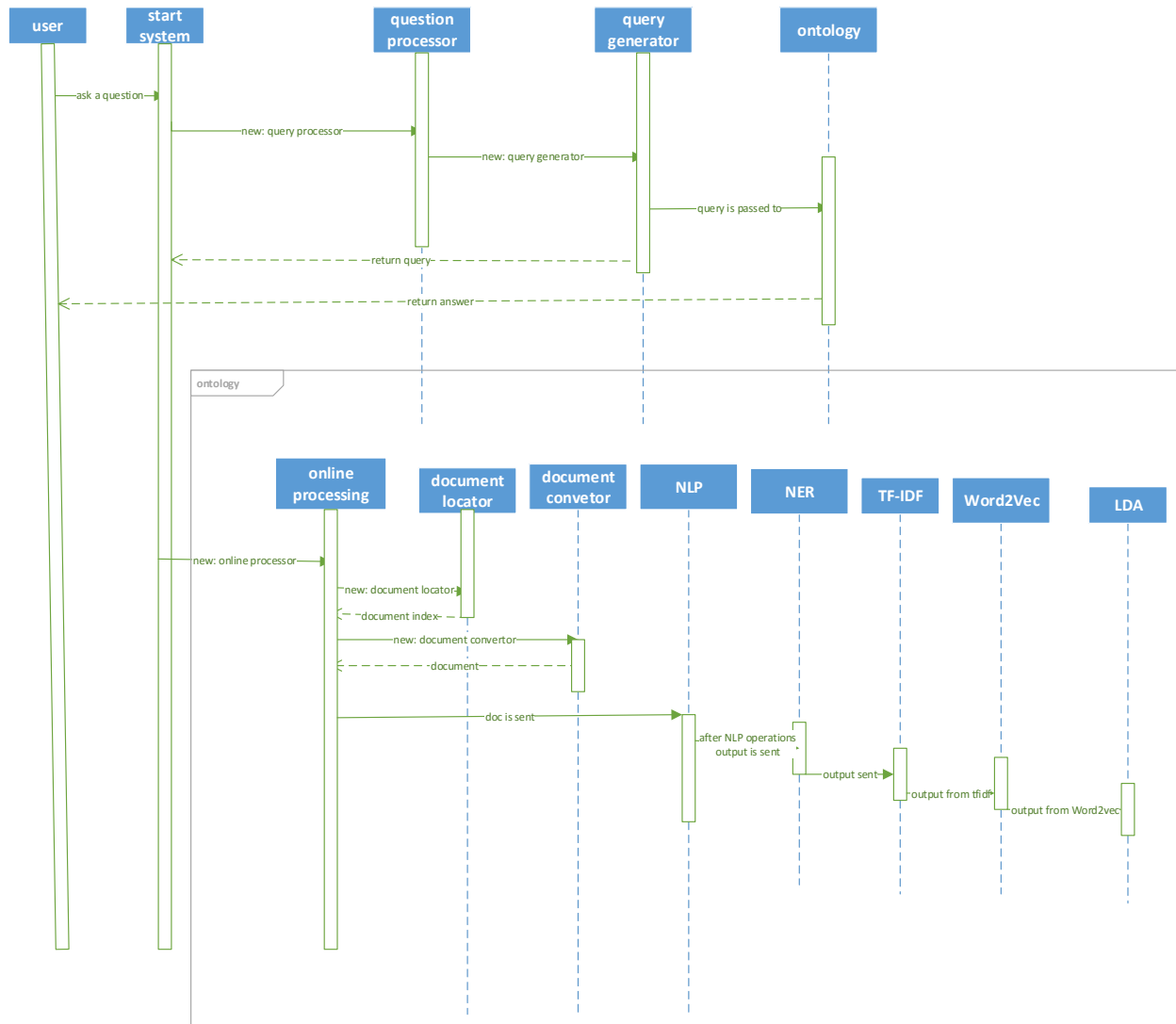
- Start system
- Online processor
- Query generator
- NLP
- Question processor
- Document convertor
- Word2Vec
- Json_plain text
- Html_plain text
- Ontology
- LDA
- Name Relation Entity Extraction

We have the following relations for classes

- Dependency relation shows class is dependent on or uses other class
- Generalization which shows the inheritance of classes like which class is a subset of or what is the superclass of this subclass
- Generalization is placed from document convertor to Json_plain text and html_plain text where document convertor is the superclass of both classes
- Dependency is placed between word parser (NLP, Word2vec,) and Topic discovery (LDA) as to parse the data we need text split into sentences and we get this from sentence extraction
- Other classes are linked with association relation that shows the connection of two classes



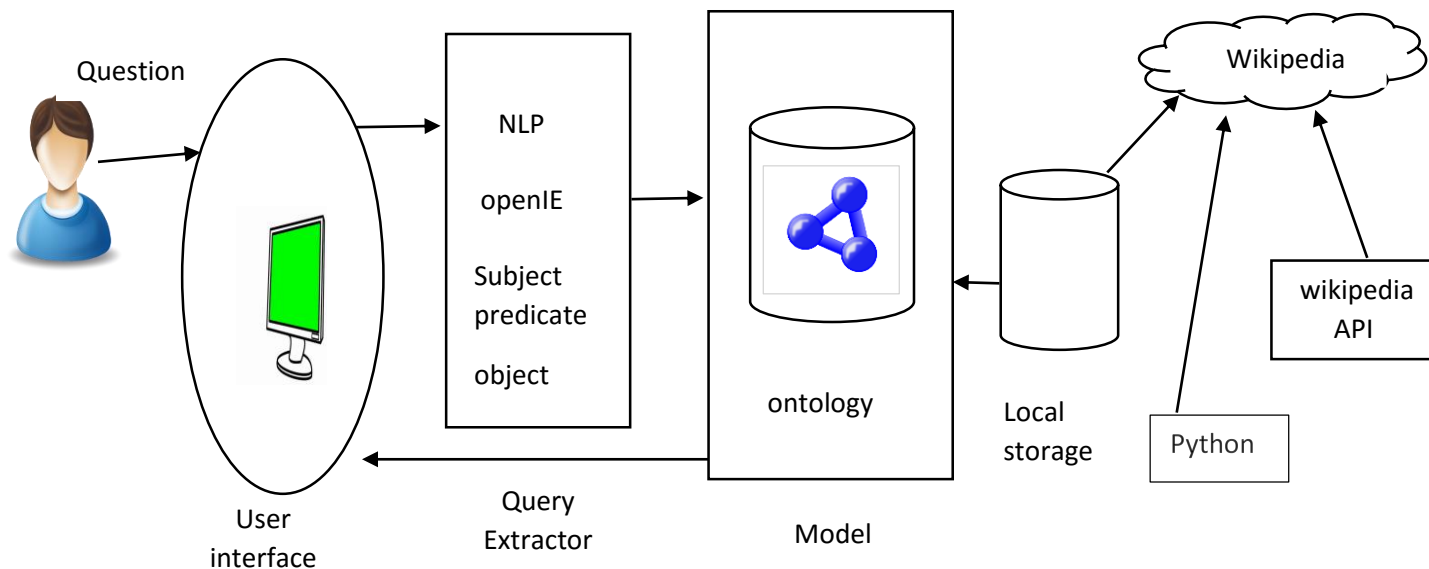
Sequence diagram: Sequence diagram gives the information about the interactions between the objects. The communication between the objects is done by passing messages from one object to another.



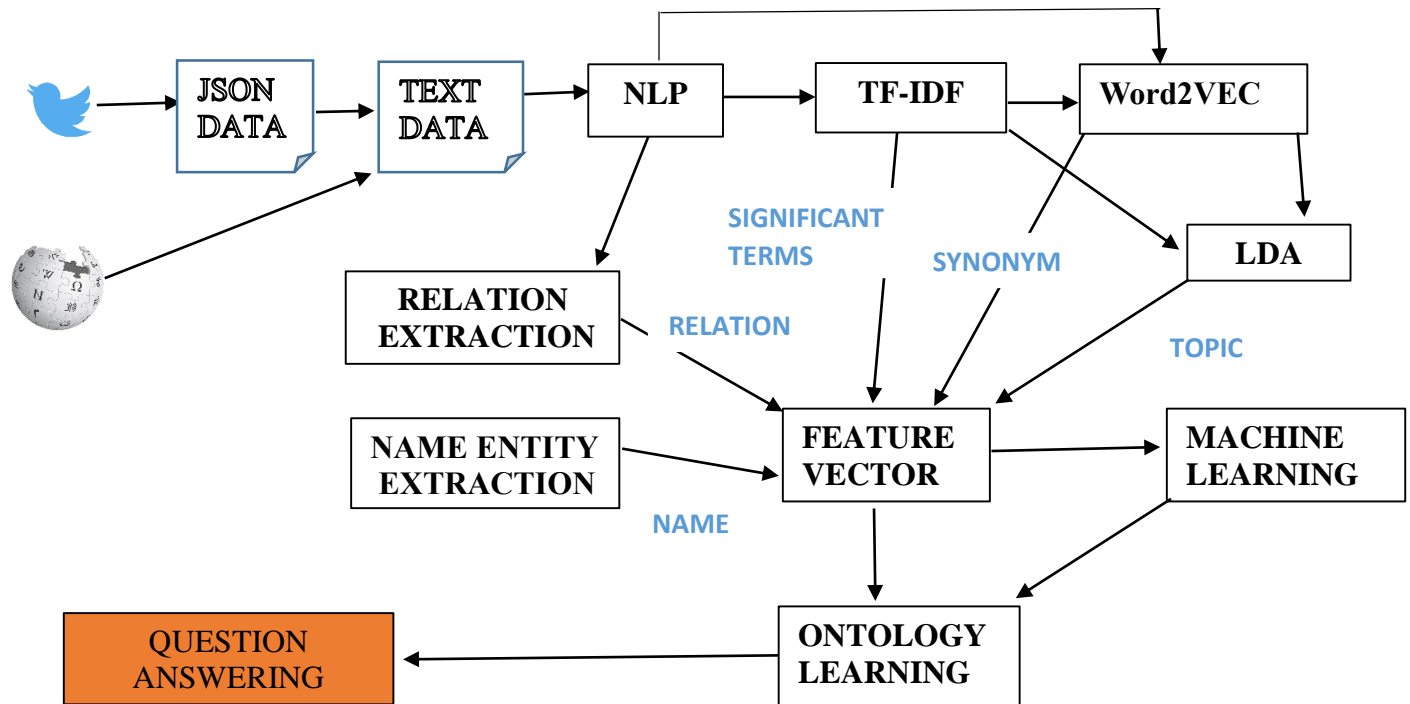
The objects involved in this model are:

- User: who interacts with the system by giving a question
- Start system: this takes the input from the user and interacts with our model
- Question processor: processes the question and classifies the question based on what, when, and other classifications
- Query generator: generates a machine readable query of the question from the user
- Online processing: streaming for dynamic data that comes from the twitter API
- Document locator: locates the document like this gives the index of the document in the local storage
- Document convertor: converts the html file or Json file
- Sentence extractor: gives the sentences from the paragraphs in the document
- Word parser: generates the tokens from the sentences
- Answer generator: generates the answers based on machine learning algorithms provided

System Architecture



WORKFLOW



Services used:

Natural language Classifier (NLC): this service uses the deep learning techniques and by using these techniques we make predictions about the predefined classes for sentences that we give as input for this service. This makes use of some machine learning algorithms to make predictions

Example of his service could be IBM Watson

MongoDB: Mongo DB is an open source database which has the convenience for the developers in terms of scalability and ease of use. Data is stored in the form of a Json file in mongo DB

Alchemy API: we make use of alchemy API like named entity extraction and other keyword extraction

- This twitter API forms as an interface for the developer to use data coming from twitter, here data is in the form of tweets. Data is collected in Json format.
- We will also use speech to text API to convert the speech that we give as input from the device into text format

Natural classifier demo: We plan to use Natural language classifier to determine the question type. In this section, the demo of the same classifier is shown. In this demo we train the classifier with twenty-five male names and twenty-five female names. Then we give a test name which the classifier takes as input and gives the percentage confidence of male and female classification.

Features developed for phase 1:

Cognitive services you want to use for project:

Implemented natural language classifier service from alchemy API. We are planning to use this service to classify questions into various categories. As of now we practiced implementing this service for classifying male and female names.

Screen shot of giving training data:

```
C:\>curl -i -u "f6de402c-0f6e-4c26-bbdd-b5e314767f72":"agBMTwd7DPnF" \ -F training_data=@C:/Users/tmcx4/Desktop/names.csv
//gateway.watsonplatform.net/natural-language-classifier/api/v1/classifiers"
curl: (6) Could not resolve host: \
curl: (6) Could not resolve host: \
curl: (6) Could not resolve host: \
HTTP/1.1 100 Continue
X-Note: Gateway Ack

HTTP/1.1 200 OK
X-Backside-Transport: OK OK
Connection: Keep-Alive
Transfer-Encoding: chunked
Date: Sat, 25 Jun 2016 00:06:41 GMT
Content-Type: application/json
```

Classifier response:

```
{
  "classifier_id" : "2373f5x67-nlc-7103",
  "name" : "TutorialClassifier",
  "language" : "en",
  "created" : "2016-06-25T00:06:40.763Z",
  "url" : "https://gateway.watsonplatform.net/natural-language-classifier/api/v1/classifiers/2373f5x67-nlc-7103",
  "status" : "Training",
  "status_description" : "The classifier instance is in its training phase, not yet ready to accept classify requests"
}
C:\>
```


Classifier Status:

```
C:\>
C:\>curl -u "f6de402c-0f6e-4c26-bbdd-b5e314767f72":"agBMTwd7DPnF" \ "https://gateway.watsonplatform.net/natural-lang
curl: (6) Could not resolve host: \
{
  "classifier_id" : "2373f5x67-nlc-7103",
  "name" : "TutorialClassifier",
  "language" : "en",
  "created" : "2016-06-25T00:06:40.763Z",
  "url" : "https://gateway.watsonplatform.net/natural-language-classifier/api/v1/classifiers/2373f5x67-nlc-7103",
  "status" : "Available",
  "status_description" : "The classifier instance is now available and is ready to take classifier requests."
}
C:\>
```

Classifier Output for test data: The name “Suresh” is classified as male with 85% confidence.

```
C:\>curl -G -u "f6de402c-0f6e-4c26-bbdd-b5e314767f72":"agBMTwd7DPnF" \ "https://gateway.watsonplatform.net/natural-
resh"
curl: (6) Could not resolve host: \
{
  "classifier_id" : "2373f5x67-nlc-7103",
  "url" : "https://gateway.watsonplatform.net/natural-language-classifier/api/v1/classifiers/2373f5x67-nlc-7103",
  "text" : "suresh",
  "top_class" : "male",
  "classes" : [ {
    "class_name" : "male",
    "confidence" : 0.8489180884346
  }, {
    "class_name" : "female",
    "confidence" : 0.15108191156540002
  } ]
}curl: (6) Could not resolve host: \
```

The name “Jessica” is classified as female with 95% confidence.

```
C:\>curl -G -u "f6de402c-0f6e-4c26-bbdd-b5e314767f72":"agBMTwd7DPnF" \ "https://gateway.watsonplatform.net/natural-
ssica"
curl: (6) Could not resolve host: \
{
  "classifier_id" : "2373f5x67-nlc-7103",
  "url" : "https://gateway.watsonplatform.net/natural-language-classifier/api/v1/classifiers/2373f5x67-nlc-7103",
  "text" : "jessica",
  "top_class" : "female",
  "classes" : [ {
    "class_name" : "female",
    "confidence" : 0.9571645890651067
  }, {
    "class_name" : "male",
    "confidence" : 0.042835410934893264
  } ]
}curl: (6) Could not resolve host: \
```

Word count count on the data collected:

Input:

Text data from twitter.

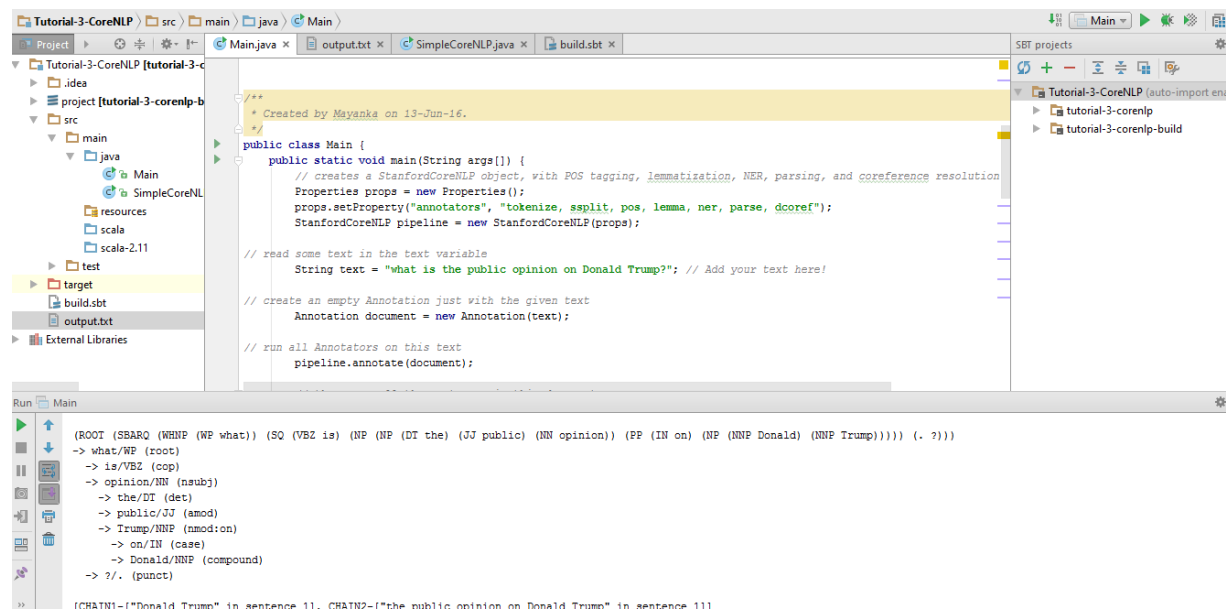
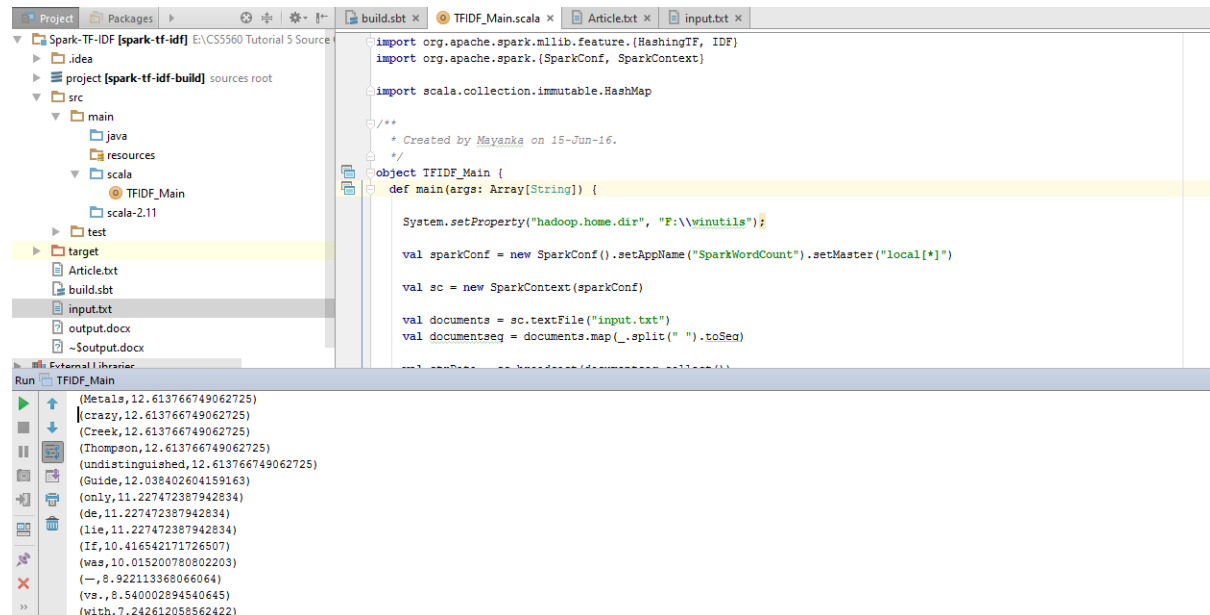
Output:

Word count for given input

Sample screen shots:

```
SparkWordCount - [E:\CS5560 T4-SourceCode\CS5560 Tutorial 4 Source Code\Spark WordCount] - [root] - ...\\output\\part-00000 - IntelliJ IDEA 2016.1.3
File Edit View Navigate Code Analyze Refactor Build Run Tools VCS Window Help
SparkWordCount > output > part-00000
Project Packages SparkWordCount.scala part-00000 part-00001 input PairRDDFunctions.scala build.sbt
SparkWordCount [root] E:\CS5560 T4-SourceCode\CS5560 Tutorial 4 Source Code\Spark WordCount
  .idea
  output
    _SUCCESS.crc
    _part-00000.crc
    _part-00001.crc
    _SUCCESS
    part-00000
    part-00001
  project [root-build] sources root
  src
    main
      java
      resources
      scala
        SparkWordCount
        (Uganda, 4)
        (better, 1)
        (House, 65)
        (Board, 2)
        (#auapoli2016, 1)
        (preludes, 9)
        (lifeline, 1)
        (order, 1)
        (https://t.co/4MpcMytE6U, 1)
        ('don't', 1)
        (@bernardkeane, 1)
        (Wake, 1)
        (been, 5)
        (underhanded, 1)
        (pandering, 2)
        (@S0C0R0C, 1)
        (@errollouis, 1)
        (02:41PM, 1)
        (@republican, 2)
Run SparkWordCount
16/06/24 20:03:53 INFO ContextCleaner: Cleaned accumulator 1
16/06/24 20:03:54 INFO SparkContext: Invoking stop() from shutdown hook
16/06/24 20:03:54 INFO SparkUI: Stopped Spark web UI at http://10.89.0.126:4040
16/06/24 20:03:54 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
16/06/24 20:03:54 INFO MemoryStore: MemoryStore cleared
16/06/24 20:03:54 INFO BlockManager: BlockManager stopped
16/06/24 20:03:54 INFO BlockManagerMaster: BlockManagerMaster stopped
16/06/24 20:03:54 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
16/06/24 20:03:54 INFO SparkContext: Successfully stopped SparkContext
16/06/24 20:03:54 INFO ShutdownHookManager: Shutdown hook called
16/06/24 20:03:54 INFO ShutdownHookManager: Deleting directory C:\Users\Vivek Veldandi\AppData\Local\Temp\spark-fc73ad7c-4d3c-4c8f-9f2a-c81e54c345od
16/06/24 20:03:54 INFO RemoteActorRefProvider$RemotingTerminator: Shutting down remote daemon.
Process finished with exit code 0
```

```
SparkWordCount - [E:\CS5560 T4-SourceCode\CS5560 Tutorial 4 Source Code\Spark WordCount] - [root] - ...\\output\\part-00001 - IntelliJ IDEA 2016.1.3
File Edit View Navigate Code Analyze Refactor Build Run Tools VCS Window Help
SparkWordCount > output > part-00001
Project Packages SparkWordCount.scala part-00000 part-00001 input PairRDDFunctions.scala build.sbt
SparkWordCount [root] E:\CS5560 T4-SourceCode\CS5560 Tutorial 4 Source Code\Spark WordCount
  .idea
  output
    _SUCCESS.crc
    _part-00000.crc
    _part-00001.crc
    _SUCCESS
    part-00000
    part-00001
  project [root-build] sources root
  src
    main
      java
      resources
      scala
        SparkWordCount
        (Hammer&#8230;Forget, 1)
        (Binding, 1)
        (https://t.co/WZuW0gkDJF, 2)
        (@Eamonnaro, 1)
        (fun, 1)
        (@brakeith, 1)
        (voters, 19)
        (but, 20)
        (turnout, 2)
        (comment, 7)
        ('liberal', 3)
        (me!, 1)
        (Must, 1)
        (https://t.co/dVSwzEKdbP, 1)
        (https://t.co/S8RzLpJ01, 1)
        (Election, 1)
        (raging, 2)
        ('educated', 1)
        (But, 7)
Run SparkWordCount
16/06/24 20:03:53 INFO ContextCleaner: Cleaned accumulator 1
16/06/24 20:03:54 INFO SparkContext: Invoking stop() from shutdown hook
16/06/24 20:03:54 INFO SparkUI: Stopped Spark web UI at http://10.89.0.126:4040
16/06/24 20:03:54 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
16/06/24 20:03:54 INFO MemoryStore: MemoryStore cleared
16/06/24 20:03:54 INFO BlockManager: BlockManager stopped
16/06/24 20:03:54 INFO BlockManagerMaster: BlockManagerMaster stopped
16/06/24 20:03:54 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
16/06/24 20:03:54 INFO SparkContext: Successfully stopped SparkContext
16/06/24 20:03:54 INFO ShutdownHookManager: Shutdown hook called
16/06/24 20:03:54 INFO ShutdownHookManager: Deleting directory C:\Users\Vivek Veldandi\AppData\Local\Temp\spark-fc73ad7c-4d3c-4c8f-9f2a-c81e54c345od
16/06/24 20:03:54 INFO RemoteActorRefProvider$RemotingTerminator: Shutting down remote daemon.
Process finished with exit code 0
```

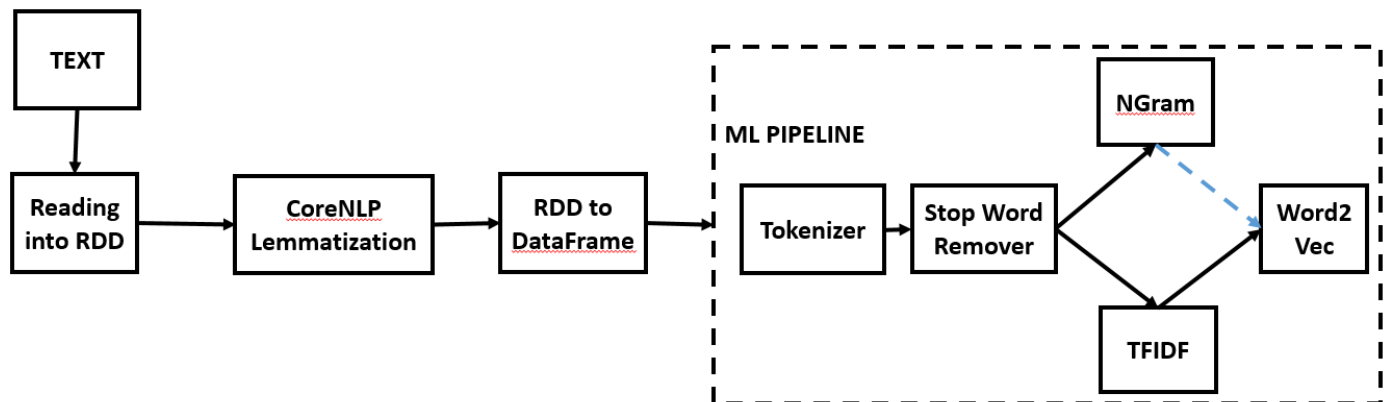


Features developed for phase 2:

NGram and Word2Vec Exercise on the data:

For the data we collected from Wikipedia and twitter we collected NGram and Word2Vec. The NGram data frame we took was 5 (Greater the number of data frame better is the sense of the output). We also collected top TF-IDF words and found synonyms for them using Word2Vec.

workflow



Input:

Text data collected from Wikipedia and twitter.

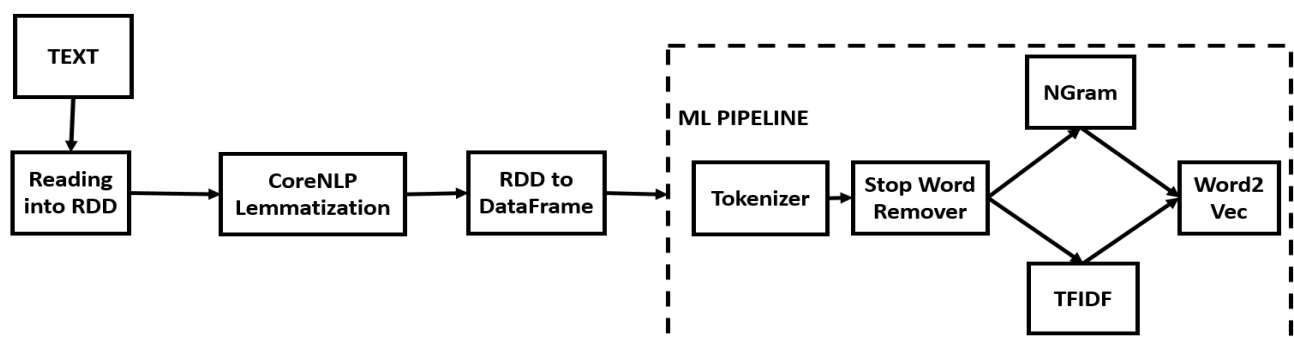
Output:

N Grams with DataFrame = 5 and synonyms for top TF-IDF words.

Name Entity Extraction and Relation Extraction on our Data:

Next we extracted structured relation triples Using OpenIE. Then by using ConceptNet5 we tried to retrieve some concepts and relationships for a few of the wiki pages (The pages where we collected our raw data).

Workflow:



The Sample output we got using OpenIE with coreNLP on Twitted data is

[(This,is,brilliant,1.0)]

[(they,call,election,1.0)]

Usage of WordNet on the data:

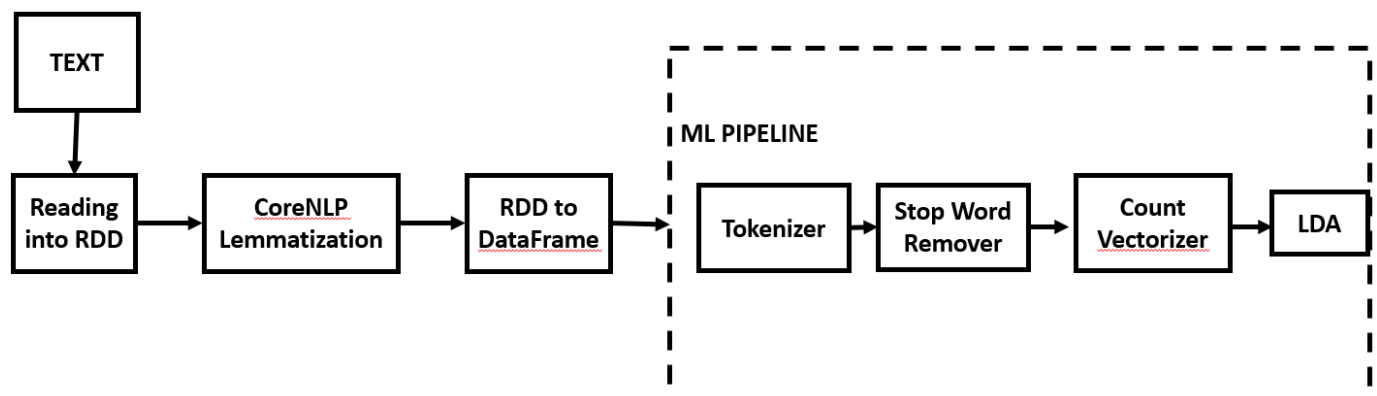
Then with the help of WordNet we tried to extract some Synonyms for some of the sample words in our data. By comparing these Synonyms with the synonyms extracted from Word2Vec we identified that synonyms from WordNet are more accurate and related than that are form Word2Vec.

Topic Discovery on the data:

We collected around 9 documents of data from Wikipedia and one document of twitter data and we classified all the words in all the documents into 4 topics using Latent Dirichlet allocation (LDA) method.

Later, we also used K- means algorithm to classify the words into groups again.

Workflow for LDA:



Input:

All the documents of data from Wikipedia and twitter.

Output:

Output has group of 4 topics each topic has words and its respective probability value for belonging to that group.

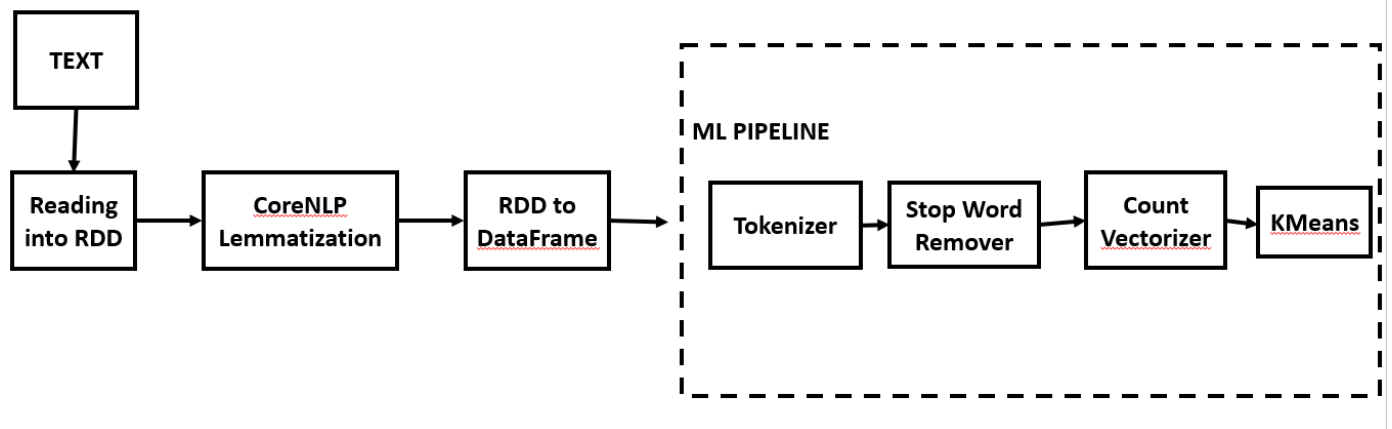
Summary of LDA execution:

```
Corpus summary:
  Training set size: 2473 documents
  Vocabulary size: 8006 terms
  Training set size: 45460 tokens
  Preprocessing time: 27.113088039 sec
```

```
Finished training LDA model. Summary:
  Training time: 29.763379458 sec
  Training data average log likelihood: -178.32006826792988
```

K – Means Clustering:

Now we applied K -means Clustering method to group the data into clusters using K – Means algorithm.



Input:

1. All the 9 documents containing Wikipedia data
2. Document containing tweets from the twitter

Output:

The K-means clustering has divided all the words in all the documents into 4 groups. The output has 4 vectors with their group words and their respective cluster IDs.

Sample Corpus Summary from K- means:

```
Corpus summary:
  Training set size: 9 documents
  Vocabulary size: 6889 terms
  Preprocessing time: 41.883770357 sec
```

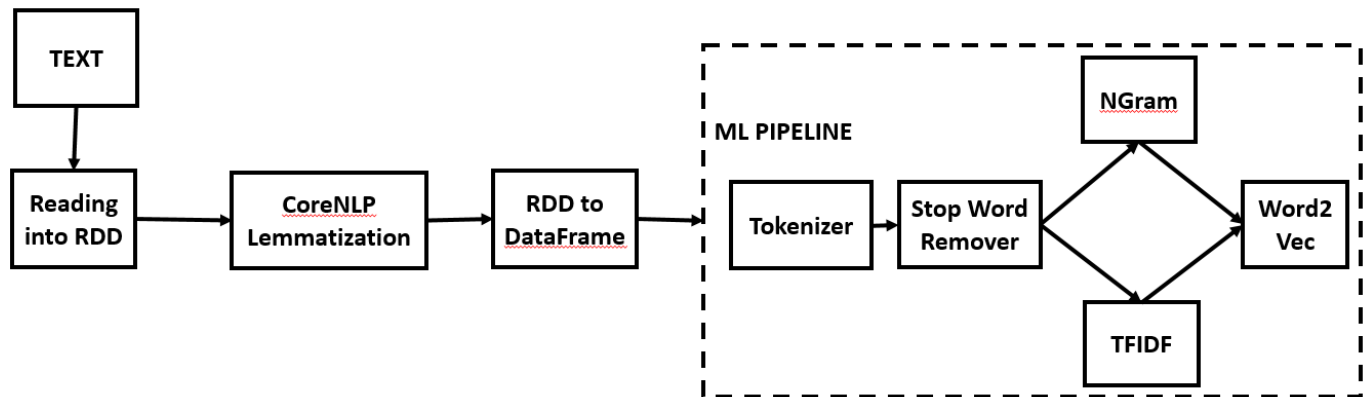
```
Finished training KMeans model. Summary:
  Training time: 2.212250886 sec
```

Feature Vector Generation: We generated two feature vectors for machine learning one with TFIDF and other with TFIDF + Word2vec.

Input: set of documents containing the text data.

Output: Feature vectors

Work flow:

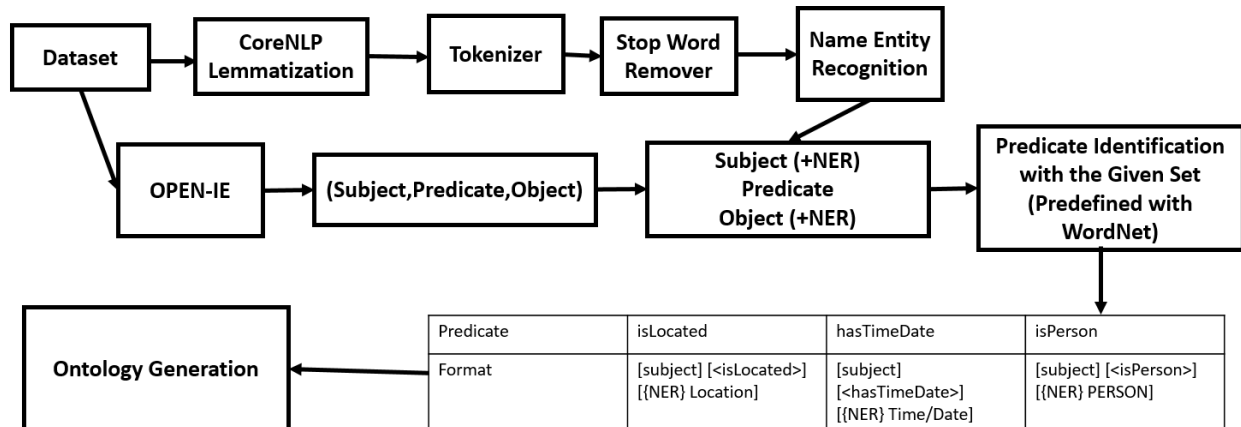


Features developed for phase 3:

We created ontology using our data both statically and dynamically

Dynamically developed Ontology: To create dynamic ontology we pruned our data to maximum extent to make the ontology meaningful. We developed two ontologies dynamically using twitter and Wikipedia data. Of the two, ontology based on Wikipedia data is more meaningful. Even though the ontology is meaningful it is not useful enough to extract sensible data.

The workflow used to create dynamic ontology is:



Screenshots Visualization using VOWL:

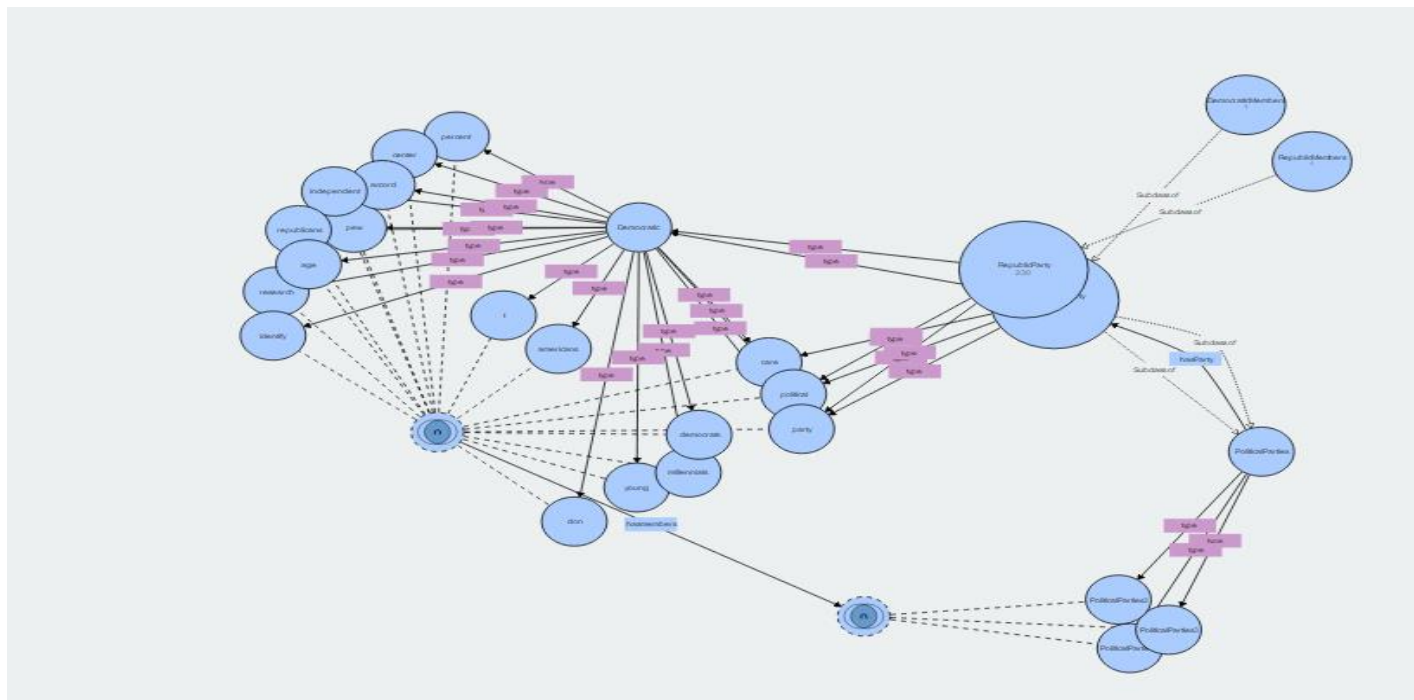


Fig: *Screenshot of Dynamic Ontology*

Note: The ontology created using the whole data could not be loaded into VOWL tool as the file was very huge. So, only part of data is used to develop and visualize the ontology

Statically developed ontology: we used protégé tool to develop static ontology.

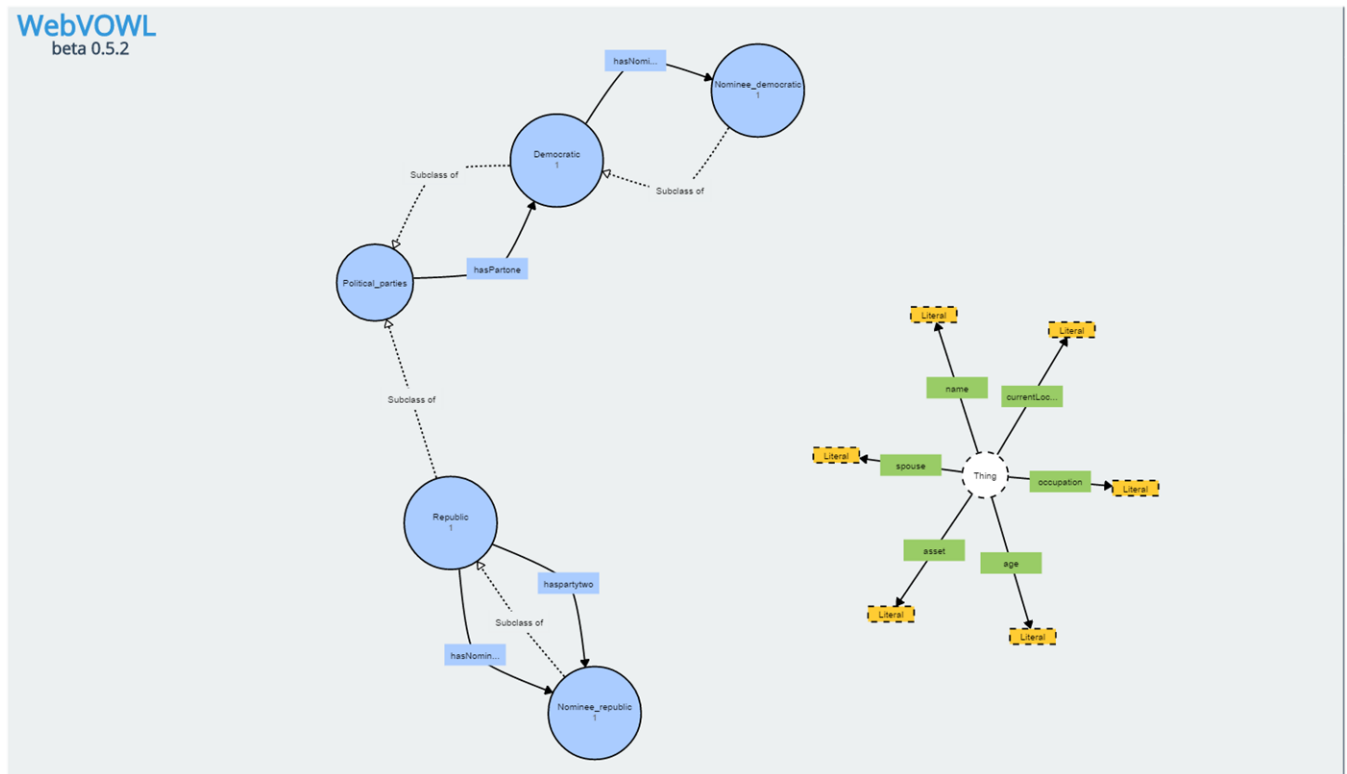


Fig: Screenshot of static Ontology using VOWL



Fig: Screenshot Data Property Hierarchy



Fig: Screenshot Object Property Hierarchy



Fig: Screenshot class Hierarchy

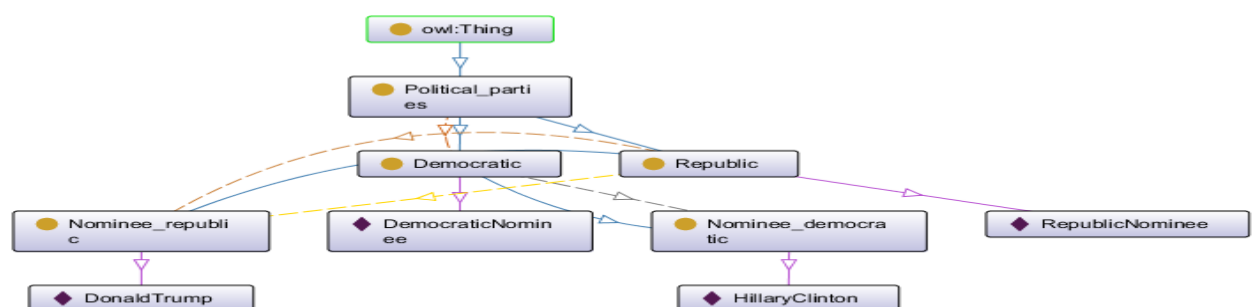


Fig: Screenshot of static Ontograph using protege

Question Answering System

Workflow of the system is as shown below:

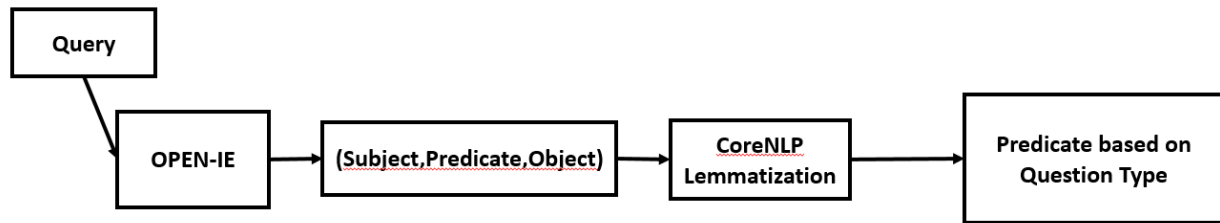


Fig: Workflow of Question answering system

From the developed ontology, we extracted data by asking some questions such as:

who is RepublicNominee?

Who is DemocraticNominee?

Where is DonaldTrump? Etc.

These questions are converted to SPARQL queries as per the workflow diagram. Sample converted queries and their respective answers

Question:Who is RepublicNominee?

Answer: Donald Trump

Converted Query(Screenshot)

```
1 PREFIX elec: <http://www.kdm.com/OWL/elections2016#>
2 SELECT ?name
3 WHERE {
4   elec:RepublicNominee elec:name ?name
5 }
6
```

Answer Screenshot:

QUERY RESULTS	
	<div>Table</div> <div>Raw Response</div> <div></div>
Showing 1 to 1 of 1 entries	
name	
1	"Donald Trump"
Showing 1 to 1 of 1 entries	

Question: Who is DemocraticNominee?

Answer: Hillary Clinton

Converted Query(Screenshot)

1	PREFIX elec: <http://www.kdm.com/OWL/elections2016#>
2	SELECT ?name
3	WHERE {
4	elec:DemocraticNominee elec:name ?name
5	}
6	

Answer Screenshot

QUERY RESULTS	
	<div>Table</div> <div>Raw Response</div> <div></div>
Showing 1 to 1 of 1 entries	
name	
1	"HillaryClinton"
Showing 1 to 1 of 1 entries	

Question: Where is DonaldTrump?

Answer: New York

Converted Query(Screenshot)

```
1 PREFIX elec: <http://www.kdm.com/OWL/elections2016#>
2 SELECT ?city
3 WHERE {
4   elec:DonaldTrump elec:currentLocation ?city
5 }
```

Answer Screenshot

QUERY RESULTS

 **Table**  

Showing 1 to 1 of 1 entries

	city
1	"NewYork"

Sample output screenshots:

```
16/07/25 22:38:13 INFO BlockManagerMaster: Registered BlockManager
Please enter your question
WR is RepublicanNominee
Reading POS tagger model from edu/stanford/nlp/models/pos-tagger/english-left3words/english-left3words-distsim.tagger ... done [1.8 sec].
Loading depparse model file: edu/stanford/nlp/models/parser/nndep/english_UD.gz ...
PreComputed 100000, Elapsed Time: 4.293 (s)
Initializing dependency parser done [8.6 sec].
Loading clause searcher from edu/stanford/nlp/models/naturalli/clauseSearcherModel.ser.gz...done [0.143 seconds]
Donald Trump
Process finished with exit code 0
```

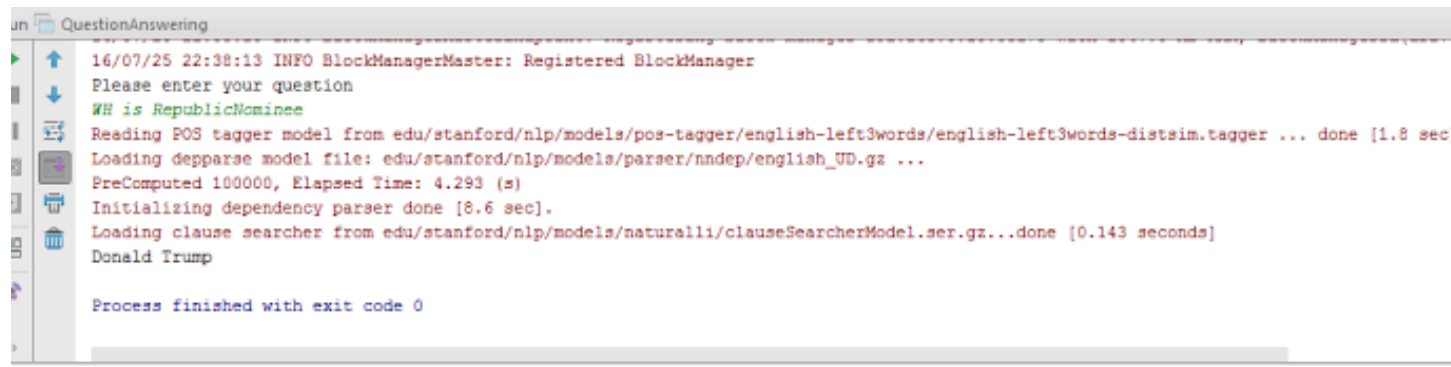
```
Tutorial 16 [naga] E:\CS5560-T16-SourceCode (2)\Source
  .idea
  data
    Election216.owl
    owl.owl
    sentenceSample
    uselection1.owl
  project [naga-build] sources root
  src
    main
    test
  target
  build.sbt
  nbst

QuestionAnswering main(args: Array[String])
7
8  */
9  object QuestionAnswering {
10    def main(args: Array[String]) {
11      val sparkConf = new SparkConf().setAppName("QuestionAnswering").setMaster("local[*]")
12
13      val sc = new SparkContext(sparkConf)
14
15      val spark = SparkSession
16        .builder
17        .appName("QuestionAnswering")
18        .master("local[*]")
19        .getOrCreate()
20      // Turn off Info Logger for Console
21
16/07/25 22:27:22 INFO SparkUI: Bound SparkUI to 0.0.0.0, and started at http://192.168.0.13:4040
16/07/25 22:27:22 INFO Executor: Starting executor ID driver on host localhost
16/07/25 22:27:22 INFO Utils: Successfully started service 'org.apache.spark.network.netty.NettyBlockTransferService' on port 58051.
16/07/25 22:27:22 INFO NettyBlockTransferService: Server created on 192.168.0.13:58051
16/07/25 22:27:22 INFO BlockManagerMaster: Trying to register BlockManager
16/07/25 22:27:22 INFO BlockManagerMasterEndpoint: Registering block manager 192.168.0.13:58051 with 1047.0 MB RAM, BlockManagerId(driver, 192.168.0.13, 58051)
16/07/25 22:27:22 INFO BlockManagerMaster: Registered BlockManager
Please enter your question
WR is DonaldTrump?
Reading POS tagger model from edu/stanford/nlp/models/pos-tagger/english-left3words/english-left3words-distsim.tagger ... done [1.8 sec].
Loading depparse model file: edu/stanford/nlp/models/parser/nndep/english_UD.gz ...
PreComputed 100000, Elapsed Time: 3.168 (s)
Initializing dependency parser done [7.5 sec].
Loading clause searcher from edu/stanford/nlp/models/naturalli/clauseSearcherModel.ser.gz...done [0.128 seconds]
NewYork
Process finished with exit code 0
```

Fig: Output for question “Where is DonaldTrump”

```
QuestionAnswering
16/07/25 22:33:42 INFO SparkUI: Bound SparkUI to 0.0.0.0, and started at http://192.168.0.13:4040
16/07/25 22:33:42 INFO Executor: Starting executor ID driver on host localhost
16/07/25 22:33:42 INFO Utils: Successfully started service 'org.apache.spark.network.netty.NettyBlockTransferService' on port 58130.
16/07/25 22:33:42 INFO NettyBlockTransferService: Server created on 192.168.0.13:58130
16/07/25 22:33:42 INFO BlockManagerMaster: Trying to register BlockManager
16/07/25 22:33:42 INFO BlockManagerMasterEndpoint: Registering block manager 192.168.0.13:58130 with 1047.0 MB RAM, BlockManagerId(driver, 192.168.0.13, 58130)
16/07/25 22:33:42 INFO BlockManagerMaster: Registered BlockManager
Please enter your question
WR is DemocraticNominee
Reading POS tagger model from edu/stanford/nlp/models/pos-tagger/english-left3words/english-left3words-distsim.tagger ... done [1.8 sec].
Loading depparse model file: edu/stanford/nlp/models/parser/nndep/english_UD.gz ...
PreComputed 100000, Elapsed Time: 4.337 (s)
Initializing dependency parser done [7.7 sec].
Loading clause searcher from edu/stanford/nlp/models/naturalli/clauseSearcherModel.ser.gz...done [0.147 seconds]
HillaryClinton
Process finished with exit code 0
```

Fig: Output for question “Who is DemocraticNominee”



```
un QuestionAnswering
16/07/25 22:38:13 INFO BlockManagerMaster: Registered BlockManager
Please enter your question
Who is RepublicNominee
Reading POS tagger model from edu/stanford/nlp/models/pos-tagger/english-left3words/english-left3words-distsim.tagger ... done [1.8 sec]
Loading depparse model file: edu/stanford/nlp/models/parser/nndep/english_UD.gz ...
PreComputed 100000, Elapsed Time: 4.293 (s)
Initializing dependency parser done [8.6 sec].
Loading clause searcher from edu/stanford/nlp/models/naturali11/c1auseSearcherModel.ser.gz...done [0.143 seconds]
Donald Trump

Process finished with exit code 0
```

Fig: Output for question “Who is RepublicNominee”

Contribution of Team members:

Ganesh Taduri (Team Leader) & Tejo Kumar: Handled the data collection part. The data is collected using Twitter API, Tweepy from Twitter and using Wikipedia API from Wikipedia. The twitter data is in Json format. We then extracted the tweet text using Json parser. We tried using the natural language classifier for the sample data. Dynamic Ontology development.

Swatwik & Sasidhar: Conceived the project idea, designed the project architecture and necessary UML diagrams. Handled the documentation part. Implemented the tutorial concepts on the collected data. Static ontology development using Protégé tool

All the teammates are involved in developing features for this phase project

Project GitHub Link: <https://github.com/ganeshtaduri/KDM>

Concerns/Issues:

- Although we started with the idea of collecting data from Twitter and LinkedIn, we had major glitches while collecting data. We had trouble in getting authorization keys.
- The execution of a few programs from tutorials is constrained by system memory and would occasionally throw errors like ‘Job aborted due to stage failure’.

Future Work:

- Data extraction from Wikipedia using DBpedia framework helps in extracting the data from info boxes
- Online twitter streaming for real time data can be incorporated
- The ontology size can be scalable
- Developing GUI for the application
- Implementing machine learning algorithms for question classification and query extraction

References:

[Class Lectures and Tutorial demos](#)

<https://docs.google.com/spreadsheets/d/1TjrVH8BxRmmvMlm3ig3X4DnwrIF3QR7ZpyUzjPbGCQ0/edit#gid=742202525>

<http://natural-language-classifier-demo.mybluemix.net/>

<http://adilmoujahid.com/posts/2014/07/twitter-analytics/>

www.developers.twitter.com

www.stackoverflow.com

www.wikipedia.com