# *Knowledge Discovery and Management*

## **Final Report**

# *Semantic Search*

## *- Movie Plots*

**Prepared by**
Sai Venkatesh Gatiganti (Class ID: 08)
Chaitanya Sai Manne (Class ID: 20)
Sri Chaitanya Patluri (Class ID: 32)
Karthik Reddy Vundela (Class ID: 43)

## Introduction:

An application to search the movies based on the plot provided by the user. We have an applications based on the keywords to find details about the keyword. Now, we are developing an application where user gives a story line or plot based details and can fetch the related movie either based in the genre or topic based. It's a simple application to find the movie names based on the plot.

## Motivation:

There is an existing Semantic search engine to find the details about the movies. We can find the details by giving a key word about the movie like movie name, cast etc. We are building an application to find the details based on the plot. We can search by giving a plot of a movie and find the movies similar to that plot.

## Objective:

To design an application, Semantic Search Engine that provides search results based on the plot given by the user on movies which are obtained from DBpedia and Wikipedia.

## Expected Outcome:

To obtain search results based on context besides keywords for better accuracy. For example, if a user searches for a plot in the search engine, the movies with similar plot as entered by the user are shown as the search results.

**Project Domain:** Movies (Plot Based Semantic Search Engine)

## Datasets:
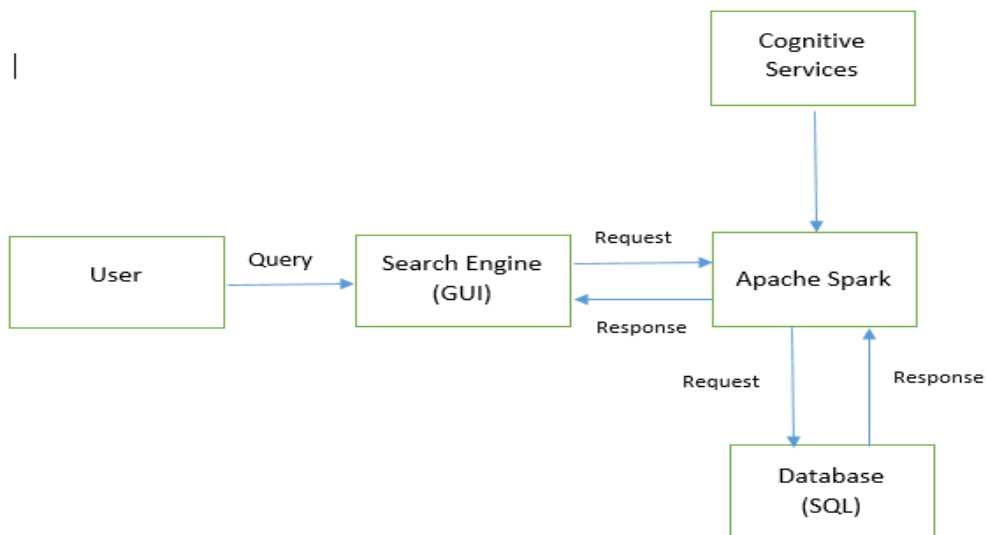
Amazon Review Data collected by Julian McAuley, UCSD -
http://jmcauley.ucsd.edu/data/amazon/links.html
***Image-based recommendations on styles and substitutes*** *J. McAuley, C. Targett, J. Shi, A. van den Hengel* SIGIR, *2015*
***Inferring networks of substitutable and complementary products*** *J. McAuley, R. Pandey, J. Leskovec, Knowledge Discovery and Data Mining, 2015*
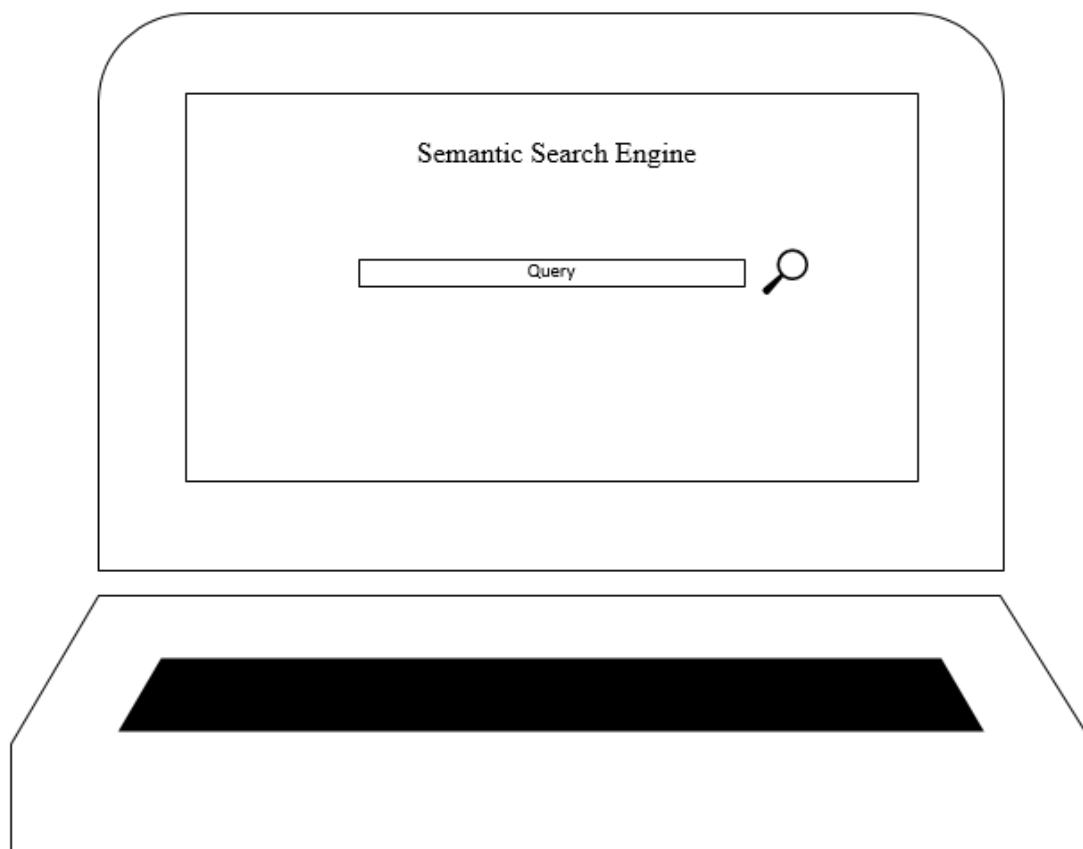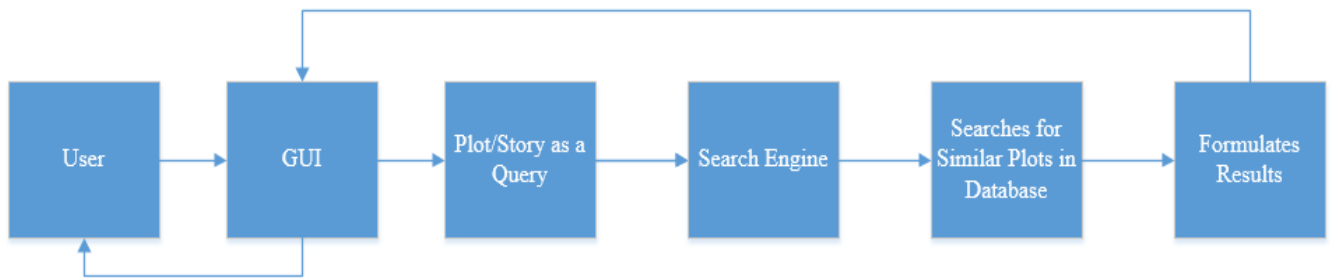
Wikipedia - https://dumps.wikimedia.org/
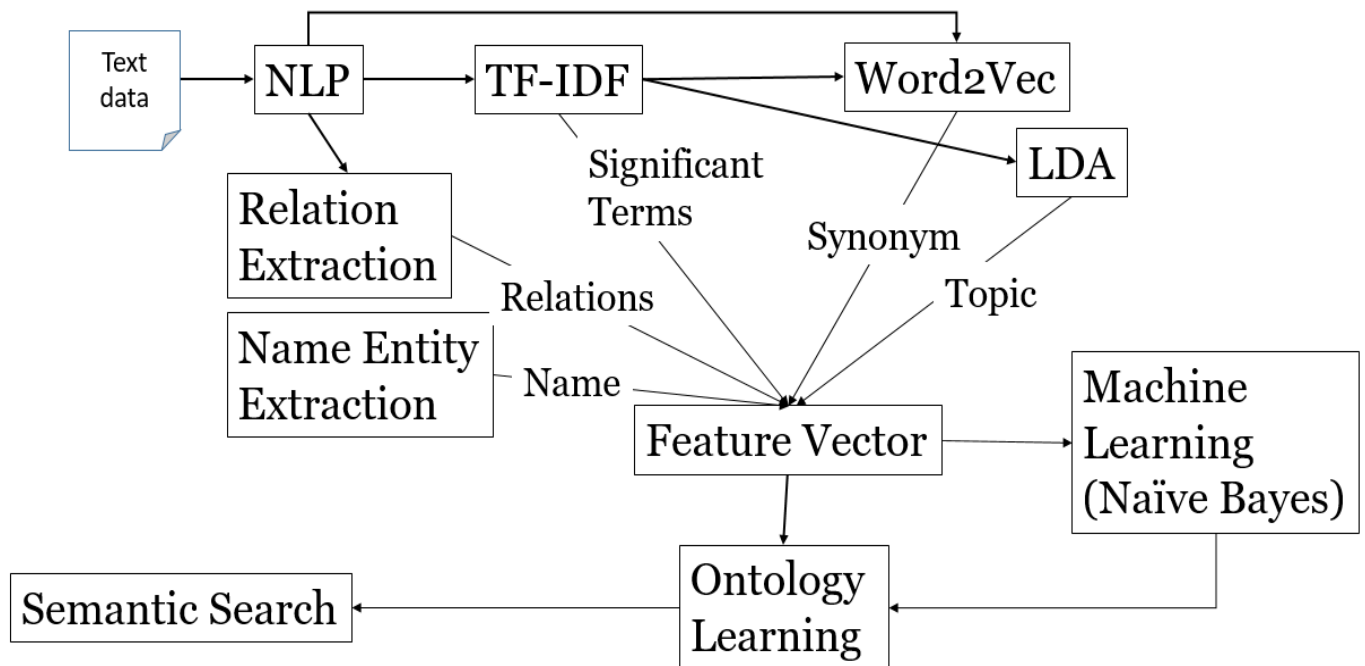
DBpedia - http://wiki.dbpedia.org/

## System Architecture:



## Wireframe of Search Engine:

**Workflow Diagram:**

`



**General Work-Flow:**

## UML Class Diagram:



## Working:

## Front-End of Search Engine:

## Code Sample:



## Input Data (Categories):

## Results:



## Movie names:

## Evaluation:

```
Accuracy: 0.31686810299527063
Weighted Precision: 0.19272083267613863
Weighted Recall: 0.3168681029952707
Weighted FMeasure: 0.22625788172691708
Precision for class index 36: 0.0
```

## Created Ontology:

## Zen-hub Project Management:

## Issues Board:



## Milestones:

**Issues & Burndown:**



**Contributions:**

Sai Venakatesh Gatiganti – Naïve Bayes Classifier, Latent Dirichlet Allocation (LDA)

Karthik Reddy Vundela – SparQL & SWRL

Chaitanya Sai Manne – Front-End, NLP, TF-IDF & ppt.

Sri Chaitanya Patluri – Report, Video & Feature vector.
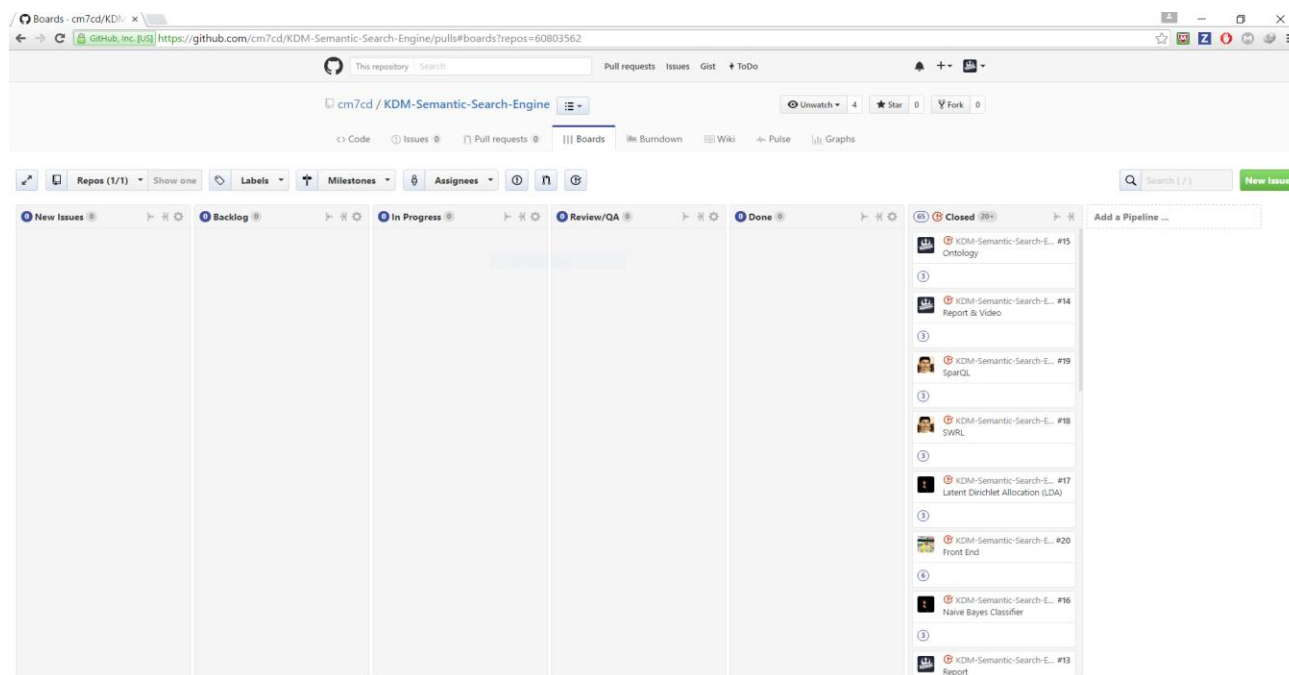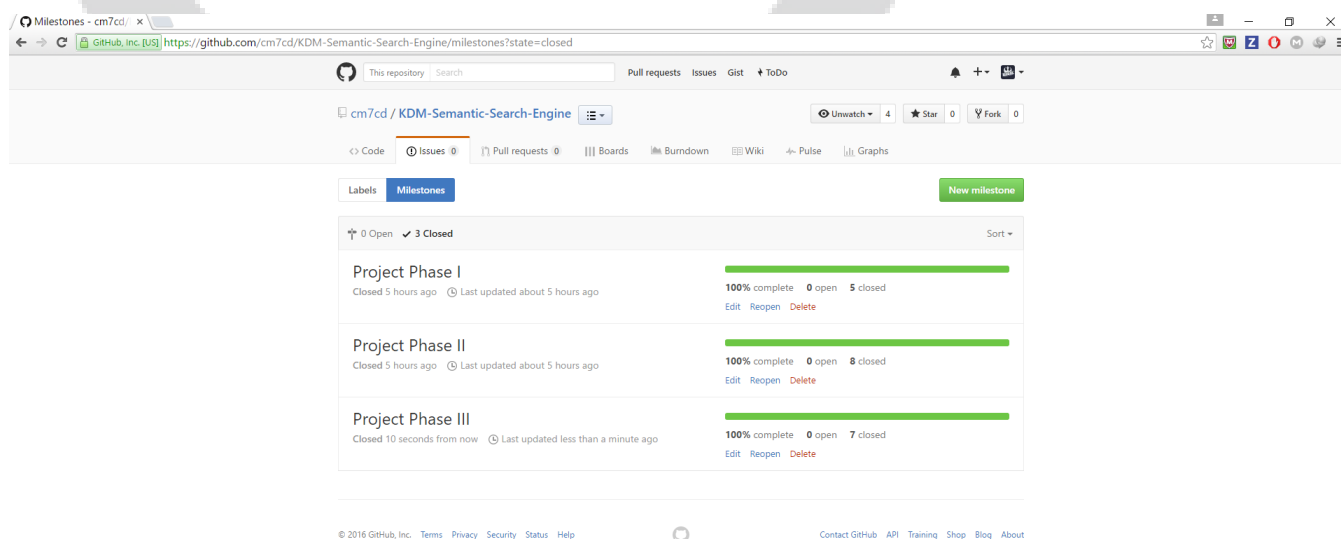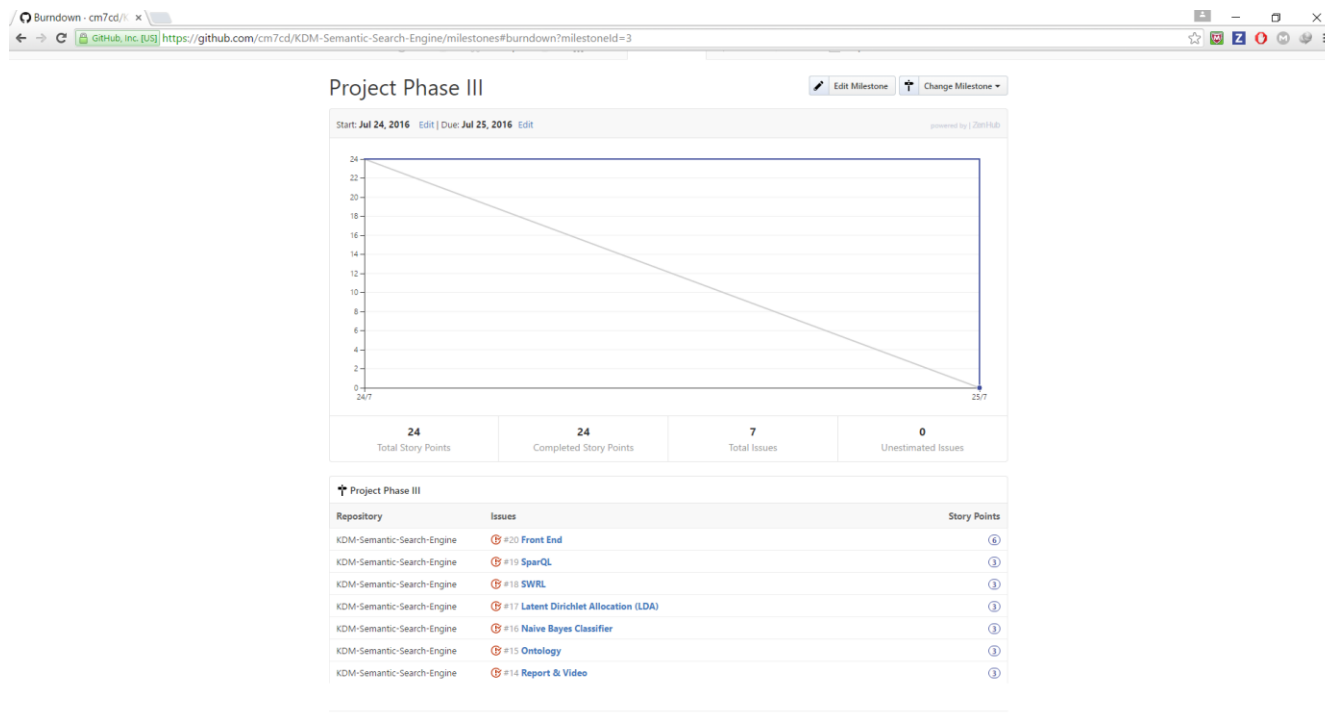
**Future Work:**

Our Future Work for the project includes including various other domains like books and research papers for the search engine and adding a dynamic web-crawler which crawls the web for related plots and displaying the results dynamically in real time.

**YouTube URL:** https://www.youtube.com/watch?v=H3vsaDTTA0M

**GitHub URL:** https://github.com/cm7cd/KDM-Semantic-Search-Engine/

**Presentation URL:** https://github.com/cm7cd/KDM-Semantic-Search-Engine/blob/master/Documentation/SemanticSearch.pptx

**Bibliography:**

- http://jmcauley.ucsd.edu/data/amazon/links.html
- http://nlp.stanford.edu/nlp/
- https://en.wikipedia.org/
- http://wiki.dbpedia.org/