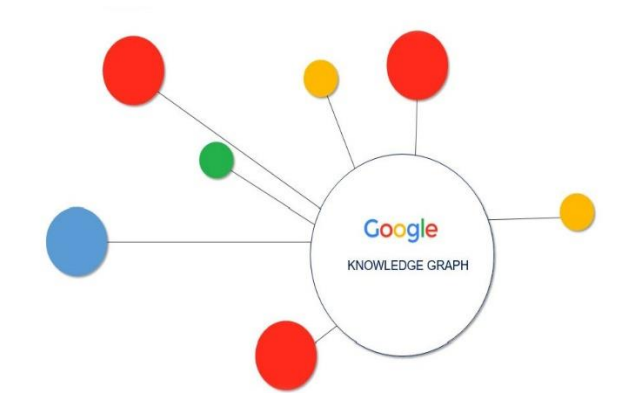# Knowledge Discovery and Management

# Project Report



# TechChamps

## Team 1

1. Jakkepalli Rama Charan Pavan (8)

2. Puthana Sujitha (24)

3. Yalamanchili Sowmya (30)

4. Nandanamudi Sreelakshmi (17)

# 1.Project Motivation, Objectives, and Significance

## 1.1 Motivation:

Data Science is a system to extract knowledge of data in various forms, either structured or unstructured from various domains, similar to Knowledge Discovery in Databases(KDD). Natural language processing is used for processing the text which is machine understandable and which will help for fast retrieval of data.

## 1.2 Specific Objectives:

1.2.1    Easy search of information from huge amount of text.

1.2.2    Helps in precise answer for customized questions.

1.2.3    Increase the knowledge management process.

## 1.3 Specific Significance:

This application helps in fetching the answer to particular questions by using NLP Process, word2vec, TF-IDF, N-gram. NLP is useful step for text processing and then we are extracting the relevant data.

# 2. Domain and Q/A application

We are taking News as our domain for our project and applying NLP operations on it and further applying question answering system for the dataset. For this question answering system, we are considering two datasets from News domain.

# 3. Specific Datasets

For our project implementation, we have considered two datasets as follows:

- ➢ WikiRef220
- ➢ BBC News-In this especially we have selected politics area and sports.

# 4. Design

## 4.1 Workflow

Step 1: Natural language processing – This process includes the identification of token, lemmatization, named entity reference(NER), co-reference resolution.

Step 2: Information Retrieval – Retrieving the information from the text. We are including the identification of the NER i.e., PERSON, LOCATION, ORGANIZATION.



Step 3: Topic Discovery – Topic discovery helps identification of the topics from the context question.

Step 4: Knowledge Graph construction – Construction of the knowledge graph from generated NER.

## 4.2 Knowledge Graph

### 4.2.1 Design workflow of knowledge Graph



### 4.2.2 Knowledge Graph for our dataset:

## 4.3 A Question-Answer Set for our Dataset.

We are designing the questions from datasets considering mainly the PERSON, LOCATION, ORGANIZATION, NUMBER entity.

1. When was Obama born?
   Born on Aug. 4, 1961.

2. Where did Obama did his schooling?
   Punahou School.

3. Who is father of Obama?
   Barack Hussein Obama.

4. Whom did Obama compete in primary race?
   Hillary Rodham Clinton.

5. What is the minimum duration for maternity leave?
   6 months.

6. What is the topic about?
   career sexism.

7. Who is the speaker?
   Ms. Hewitt.

8. What is the average pay for full-time women.
   80p

9. What is the average pay for part-time women.
   60p.

10. What is the average pay for retired women compared to men?
    Half.

# 5. Implementation

## 5.1 Output of NLP operations for our dataset

We have performed the NLP operations on the dataset which we have chosen and the result of each operation is shown in the below mentioned screenshots.

Tokenization:

Lemmatization:

POS Tagging:

NER:

Coreference Resolution:

Below is the code for all the operations of the NLP.

File   Edit   View   Navigate   Code   Analyze   Refactor   Build   Run   Tools   VCS   Window   Help

Tutorial-3-CoreNLP ⟩ src ⟩ main ⟩ java ⟩ Main ⟩ Main

```java
        // a CoreLabel is a CoreMap with additional token-specific methods
        for (CoreLabel token : sentence.get(CoreAnnotations.TokensAnnotation.class)) {
            System.out.println("\n" + token);
            //    this is the text of the token
            String word = token.get(CoreAnnotations.TextAnnotation.class);
            System.out.println("Text Annotation");
            System.out.println(token + ":" + word);
            // this is the POS tag of the token
            String lemma = token.get(CoreAnnotations.LemmaAnnotation.class);
            System.out.println("Lemma Annotation");
            System.out.println(token + ":" + lemma);
            String outputPath = "output.txt";
            String x= token + ":" + lemma;
            // this is the Lemmatized tag of the token
            String pos = token.get(CoreAnnotations.PartOfSpeechAnnotation.class);
            System.out.println("POS");
            System.out.println(token + ":" + pos);
            // this is the NER label of the token
            String ne = token.get(CoreAnnotations.NamedEntityTagAnnotation.class);
            System.out.println("NER");
            System.out.println(token + ":" + ne);
            System.out.println("\n\n");
            // this is the parse tree of the current sentence
            Tree tree = sentence.get(TreeCoreAnnotations.TreeAnnotation.class);
            System.out.println(tree);
            // this is the Stanford dependency graph of the current sentence
            SemanticGraph dependencies = sentence.get(SemanticGraphCoreAnnotations.CollapsedCCProcessed
            System.out.println(dependencies.toString());
            Map<Integer, CorefChain> graph =
                    document.get(CorefCoreAnnotations.CorefChainAnnotation.class);
            System.out.println(graph.values().toString());
```

SBT projects

▼ Tutorial-3-CoreNLP (auto-import enabled
    tutorial-3-corenlp
    tutorial-3-corenlp-build

All files are up-to-date (4 minutes ago)          28:45   CRLF   UTF-8

Below are the outputs of the operations.

heat-5:heat
Lemma Annotation
heat-5:heat
POS
heat-5:NN
NER
heat-5:O

(ROOT (S (ADVP (RB Once)) (, ,) (PP (IN during) (NP (NP (DT the) (NN heat)) (PP (IN of) (NP (NP (DT the) (JJ primary) (NN race)) (PP (IN between) (NP (NP (NNP Obama)) (CC and) (NP (NNP Hillary) (NNP Rodham) (NNP Clinton)))))))) (, ,) (NP (DT a) (NN claim)) (VP (VBD came) (PP (IN from) (NP (NP (NNP Bill) (NNP Clinton)) (PP (IN that) (NP (NP (PRP he)) ('' '') (VP (VBN understood) ('' '') (NP (NNP Obama)))))))) (. .)))

-> came/VBD (root)
  -> Once/RB (advmod)
  -> ,/, (punct)
  -> heat/NN (nmod:during)
    -> during/IN (case)
    -> the/DT (det)
    -> race/NN (nmod:of)
      -> of/IN (case)
      -> the/DT (det)
      -> primary/JJ (amod)
      -> Obama/NNP (nmod:between)
        -> between/IN (case)
        -> and/CC (cc)
        -> Clinton/NNP (conj:and)
          -> Hillary/NNP (compound)
          -> Rodham/NNP (compound)
        -> Clinton/NNP (nmod:between)
  -> ,/, (punct)
  -> claim/NN (nsubj)
    -> a/DT (det)
  -> Clinton/NNP (nmod:from)
    -> from/IN (case)
    -> Bill/NNP (compound)
    -> he/PRP (nmod:that)
      -> that/IN (case)
      -> ''/'' (punct)
      -> understood/VBN (acl)
        -> ''/'' (punct)
        -> Obama/NNP (dobj)
  -> ./. (punct)

[CHAIN1-["Obama" in sentence 1, "Himself" in sentence 1, "His" in sentence 1, "Barry Obama" in sentence 1, "he" in sentence 2, "he" in sentence 2, "Barack Obama" in sentence 6, "Obama" in sentence 8, "Obama" in sentence 10, "young Obama" in sentence 13, "his" in sentence 13, "He" in sentence 14, "himself" in sentence 14, "he" in sentence 15, "his" in sentence 15, "Barack Obama 's" in sentence 16, "Barack Hussein Obama" in sentence 25, "Obama 's" in sentence 31, "his" in sentence 32, "him" in sentence 32, "Obama" in sentence 34, "Obama" in sentence 35, "he '' understood '' Obama" in sentence 35, "Obama" in sentence 35, "he" in sentence 36], CHAIN3=["10th" in sentence 1], CHAIN4=["S. Beretania" in sentence 1, "Beretania" in sentence 28], CHAIN5=["Waikiki Beach" in sentence 1], CHAIN7=["It" in sentence 1, "That" in sentence 1], CHAIN8=["Hawaii That Made His Rise Possible On weekday mornings as a teenager , Barry Obama left his grandparents ' apartment on the 10th floor of the 12-story high-rise at 1617 S. Beretania , a mile and a half above Waikiki Beach , and walked up Punahou Street in the shadows of capacious banyan trees and date palms" in sentence 1, "Punahou" in sentence 28], CHAIN10=["His Rise Possible On weekday mornings as a teenager , Barry Obama left his grandparents ' apartment on the 10th floor of the 12-story high-rise at 1617 S. Beretania , a mile and a half above Waikiki Beach , and walked up Punahou Street in the shadows of capacious banyan trees and date palms" in sentence 1, "his rise" in sentence 6, "it" in sentence 6, "a geographical truth" in sentence 6, "his rise" in sentence 15], CHAIN12=["weekday mornings as a teenager" in sentence 1], CHAIN13=["a teenager" in sentence 1, "his" in sentence 1], CHAIN14=["his grandparents ' apartment" in sentence 1], CHAIN15=["his grandparents '" in sentence 1], CHAIN17=["the 10th floor of the 12-story high-rise at 1617 S. Beretania , a mile and a half" in sentence 1], CHAIN18=["the 12-story high-rise at 1617 S. Beretania , a mile and a half" in sentence 1, "the high-rise" in sentence 28], CHAIN19=["1617 S. Beretania , a mile and a half" in sentence 1], CHAIN20=["1617 S. Beretania" in sentence 1], CHAIN21=["a mile" in sentence 1], CHAIN22=["a half" in sentence 1], CHAIN23=["the shadows of capacious banyan trees and date palms" in sentence 1], CHAIN24=["the shadows of capacious banyan trees and date palms" in sentence 1], CHAIN25=["Punahou Street in the shadows of capacious banyan trees and date palms" in sentence 1],

---

30], CHAIN334=["a human being" in sentence 30], CHAIN335=["a politician , without his mother 's sensibility , naive or adventurous or both" in sentence 30, "They" in sentence 32, "they" in sentence 33, "they" in sentence 33], CHAIN336=["a politician , without his mother 's sensibility" in sentence 30], CHAIN339=["his mother 's sensibility" in sentence 30], CHAIN342=["both" in sentence 30], CHAIN346=["the same roof for only perhaps 12 years" in sentence 32], CHAIN347=["only perhaps 12 years" in sentence 32], CHAIN348=["his adolescence" in sentence 32], CHAIN350=["her lessons and judgments" in sentence 32], CHAIN353=["some sense" in sentence 33], CHAIN356=["just 18 years" in sentence 33], CHAIN357=["each following a singular path toward maturity" in sentence 33, "him" in sentence 34], CHAIN358=["a singular path toward maturity" in sentence 33], CHAIN359=["Bill Clinton" in sentence 34], CHAIN361=["many presidential aspirants before him , and perhaps most like Bill Clinton" in sentence 34], CHAIN362=["many presidential aspirants before him" in sentence 34], CHAIN364=["most like Bill Clinton" in sentence 34], CHAIN365=["strong women , the male figures either weak or absent" in sentence 34], CHAIN366=["strong women" in sentence 34], CHAIN367=["the male figures either weak or absent" in sentence 34], CHAIN368=["the male figures" in sentence 34], CHAIN370=["Hillary Rodham Clinton" in sentence 35], CHAIN372=["the heat of the primary race between Obama and Hillary Rodham Clinton" in sentence 35], CHAIN373=["the primary race between Obama and Hillary Rodham Clinton" in sentence 35], CHAIN374=["Obama and Hillary Rodham Clinton" in sentence 35], CHAIN375=["a claim" in sentence 35], CHAIN376=["Bill Clinton that he '' understood '' Obama" in sentence 35], CHAIN378=["different as their backgrounds and families" in sentence 36, "it" in sentence 36, "this strong-female , weak-male similarity" in sentence 36, "this strong-female" in sentence 36], CHAIN379=["their backgrounds and families" in sentence 36]]

of-6
Text Annotation
of-6:of
Lemma Annotation
of-6:of
POS
of-6:IN
NER
of-6:O

(ROOT (S (ADVP (RB Once)) (, ,) (PP (IN during) (NP (NP (DT the) (NN heat)) (PP (IN of) (NP (NP (DT the) (JJ primary) (NN race)) (PP (IN between) (NP (NP (NNP Obama)) (CC and) (NP (NNP Hillary) (NNP Rodham) (NNP Clinton)))))))) (, ,) (NP (DT a) (NN claim)) (VP (VBD came) (PP (IN from) (NP (NP (NNP Bill) (NNP Clinton)) (PP (IN that) (NP (NP (PRP he)) ('' '') (VP (VBN understood) ('' '') (NP (NNP Obama)))))))) (. .)))

-> came/VBD (root)
  -> Once/RB (advmod)
  -> ,/, (punct)
  -> heat/NN (nmod:during)
    -> during/IN (case)
    -> the/DT (det)
    -> race/NN (nmod:of)
      -> of/IN (case)
      -> the/DT (det)
      -> primary/JJ (amod)
      -> Obama/NNP (nmod:between)
        -> between/IN (case)
        -> and/CC (cc)
        -> Clinton/NNP (conj:and)
          -> Hillary/NNP (compound)
          -> Rodham/NNP (compound)
        -> Clinton/NNP (nmod:between)
  -> ,/, (punct)
  -> claim/NN (nsubj)
    -> a/DT (det)
  -> Clinton/NNP (nmod:from)
    -> from/IN (case)
    -> Bill/NNP (compound)
    -> he/PRP (nmod:that)
      -> that/IN (case)
      -> ''/'' (punct)
      -> understood/VBN (acl)
        -> ''/'' (punct)
        -> Obama/NNP (dobj)

```
Chicago" in sentence 18, "Chicago" in sentence 19, "it" in sentence 19, "Chicago" in sentence 21], CHAIN158-["Chicago , the antipode of remote Honolulu , deep in the fold of the mainland" :
, CHAIN185-["Hawaii and Chicago" in sentence 16, "the two main threads weaving through the cloth of Barack Obama 's life" in sentence 16], CHAIN187-["the cloth of Barack Obama 's life" in :
ner conflicts" in sentence 19], CHAIN212-["the subtle , coolly ambitious persona" in sentence 19], CHAIN214-["the presidential election" in sentence 19], CHAIN214-["first" in sentence 20],
sentence 24, "It" in sentence 25, "It" in sentence 26, "It" in sentence 27, "that" in sentence 28], CHAIN232-["community work" in sentence 24], CHAIN233-["lives of public service" in senter
HAIN255-["his grandparents , Madelyn and Stan Dunham , Toot and Gramps , the white couple with whom he lived for most of his teenage years , she practical and determined , he impulsive , he
ve , hokey , well-intentioned and , by his grandson 's account , burdened with the desperate lost hopes of a Willy Loman-style salesman" in sentence 26], CHAIN261-["the white couple with wh
HAIN284-["the West Coast to Hawaii" in sentence 27], CHAIN286-["52" in sentence 28], CHAIN291-["their daughter , who followed the Pacific farther to Indonesia" in sentence 28], CHAIN292-["
in sentence 29], CHAIN310-["1995" in sentence 29], CHAIN313-["his debut" in sentence 29], CHAIN315-["the national stage" in sentence 29], CHAIN316-["a book about himself that searched for 1
a politician , without his mother 's sensibility , naive or adventurous or both" in sentence 30, "They" in sentence 32, "they" in sentence 33, "they" in sentence 33], CHAIN336-["a politicia
re him , and perhaps most like Bill Clinton" in sentence 34], CHAIN362-["many presidential aspirants before him" in sentence 34], CHAIN364-["most like Bill Clinton" in sentence 34], CHAIN36
ale" in sentence 36], CHAIN379-["their backgrounds and families" in sentence 36]]


heat-5
Text Annotation
heat-5:heat
Lemma Annotation
heat-5:heat
POS
heat-5:NN
NER
heat-5:O


(ROOT (S (ADVP (RB Once)) (, ,) (PP (IN during) (NP (NP (DT the) (NN heat)) (PP (IN of) (NP (NP (DT the) (JJ primary) (NN race)) (PP (IN between) (NP (NP (NNP Obama)) (CC and) (NP (NNP Hil
-> came/VBD (root)
  -> Once/RB (advmod)
  -> ,/, (punct)
  -> heat/NN (nmod:during)
    -> during/IN (case)
    -> the/DT (det)
    -> race/NN (nmod:of)
      -> of/IN (case)
      -> the/DT (det)
      -> primary/JJ (amod)
      -> Obama/NNP (nmod:between)
        -> between/IN (case)
        -> and/CC (cc)
        -> Clinton/NNP (conj:and)
          -> Hillary/NNP (compound)
          -> Rodham/NNP (compound)
        -> Clinton/NNP (nmod:between)
  -> ,/, (punct)
```

## 5.2 Question Answering for our dataset

After performing the NLP operations, we have taken the post processed dataset and we have separately stored the result of NER output like from the NER result we have the all person related entities to one file and similarly we have done for every group and based on that we have generated answers for the questions we choose. The below screenshots will depict the same.



```java
public class NewNLP {

    public static void main(String args[]) throws IOException {

        Properties props = new Properties();
        props.setProperty("annotators", "tokenize, ssplit, pos, lemma, ner, parse, dcoref");
        StanfordCoreNLP pipeline = new StanfordCoreNLP(props);
        Set personSet = new HashSet();
        Set locSet = new HashSet();
        Set orgSet = new HashSet();

        String text = readFile("input.txt");


        Annotation document = new Annotation(text);

        pipeline.annotate(document);

        List<CoreMap> sentencess = document.get(CoreAnnotations.SentencesAnnotation.class);

        for (CoreMap sentence : sentencess) {
            for (CoreLabel token : sentence.get(CoreAnnotations.TokensAnnotation.class)) {
                String nameAndEntity = token.get(CoreAnnotations.NamedEntityTagAnnotation.class);

                if (nameAndEntity.equals("PERSON")) {
                    personSet.add(token);
                }
                if (nameAndEntity.equals("LOCATION")) {
                    locSet.add(token);
                }
                if (nameAndEntity.equals("ORGANIZATION")) {
                    orgSet.add(token);
                }
```
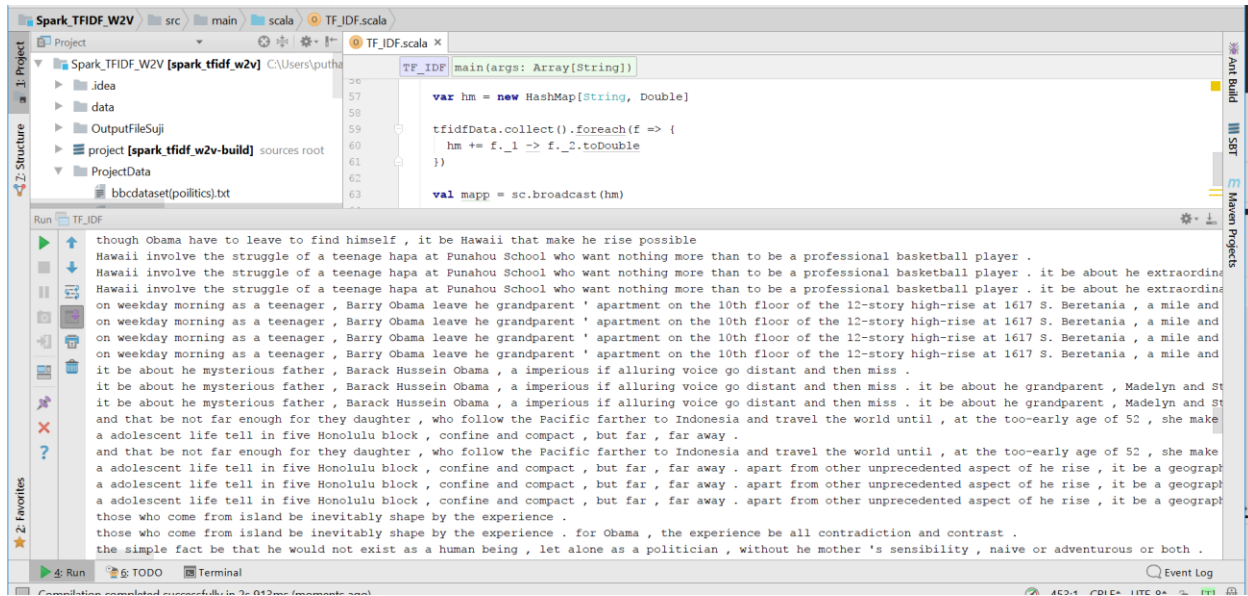
```java
System.out.println("------Named Entity Recognition----");
System.out.println("\n\n");
System.out.println("People");
System.out.print(personSet);
System.out.println("\n\n");
System.out.println("Location");
System.out.print(locSet);
System.out.println("\n\n");
System.out.println("Organization");
System.out.print(orgSet);
System.out.println("\n\n");
System.out.println("Enter question based on above data");
Scanner sc = new Scanner(System.in);
String ques = sc.nextLine();
if(ques.equalsIgnoreCase("Web Blog community is called into support which Bloggers")){
    System.out.println("Iranian bloggers");
}
if(ques.equalsIgnoreCase("What is Obama inspiration")){
    System.out.println("The first Obama-inspired trend is Total Absorption");
}

}
public static String readFile(String file) throws IOException {
    BufferedReader bReader = new BufferedReader(new FileReader(file));
    try {
        StringBuilder stB = new StringBuilder();
        String line = bReader.readLine();

        while (line != null) {
            stB.append(line);
            stB.append("\n");
            line = bReader.readLine();
        }
```

The below screenshot shows the output of our dataset which was categorized into different NER tags. When we compose a question, the corresponding answer will be taken based on the list of the NER entities.

```java
String nameAndEntity = token.get(CoreAnnotations.NamedEntityTagAnnotation.class);

if (nameAndEntity.equals("PERSON")) {
    personSet.add(token);
}
if (nameAndEntity.equals("LOCATION")) {
    locSet.add(token);
}
if (nameAndEntity.equals("ORGANIZATION")) {
```

```
People
[Friedel-15, Stapleton-32, Van-3, Wilson-32, Andrews-4, Sayer-109, Neill-78, Chris-8, Modric-122, Mark-31, Lennon-115, Raymond-1, Dicks-24, Behrami-11, Nsereko-93, Dixon-32, Jona

Location
[UK-26, South-18, Hampstead-106, Urmston-111, France-96, Petersburg-25, Northern-10, Ohio-36, Wales-10, UK-46, UK-29, Farnworth-56, Glasgow-58, UK-77, Hillingdon-145, Quito-31, S

Organization
[Crystal-11, Southampton-35, UNITED-8, Leeds-117, St-71, Middlesbrough-109, Chelsea-9, Coventry-100, Fulham-3, Liverpool-13, Leamington-70, McDonald-43, Wigan-64, Barnsley-75, Ha

Enter question based on above data
What are the three imports of West Africa
Emmanuel Eboue, Johan Djourou and Alex Song

Process finished with exit code 0
```

Compilation completed successfully in 2s 462ms (48 minutes ago)

## 5.3 TF-IDF for our dataset

Generating the term frequency for the words in the dataset.



Generated Output.

## Generation N-gram for the dataset



## Generating the TF_IDF for N-gram output.

# 6. Project Management

## Programming Language Used:

We have collaborated various languages in the development of the project and in building the application. Some of them are,

- Java

- Scala

- Spark
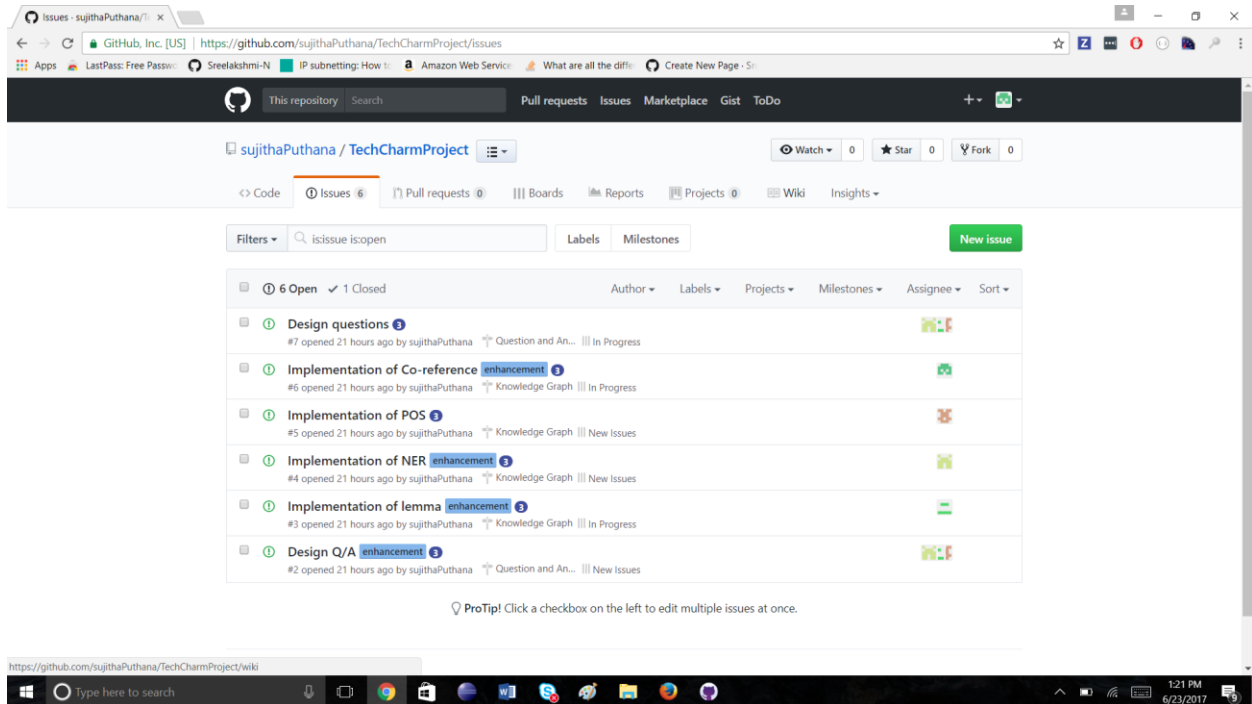
## IDE Used:

- IntelliJ

- PyCharm

## 6.1 Contributors

- Jakkepalli Rama Charan Pavan - **25%**

- Puthana Sujitha - **25%**

- Yalamanchili Sowmya **- 25%**

- Nandanamudi Sreelakshmi **- 25%**

Below is the bar graph that represents contribution of each person in the team towards project,

## 6.2 Zen-Hub Screenshots

For the first increment, we had issues regarding the working of the questions and answers section and generating the NLP output for the dataset we have chosen as the size of the dataset is larger.



## Project Timeline, Members, and Task Responsibility

The issues that are registered and current one's which we are working are updated and can be viewed in GitHub repository. The below screenshot will show you the issues and their respective categorization's i.e. New issues, Icebox, Backlog, In Progress.
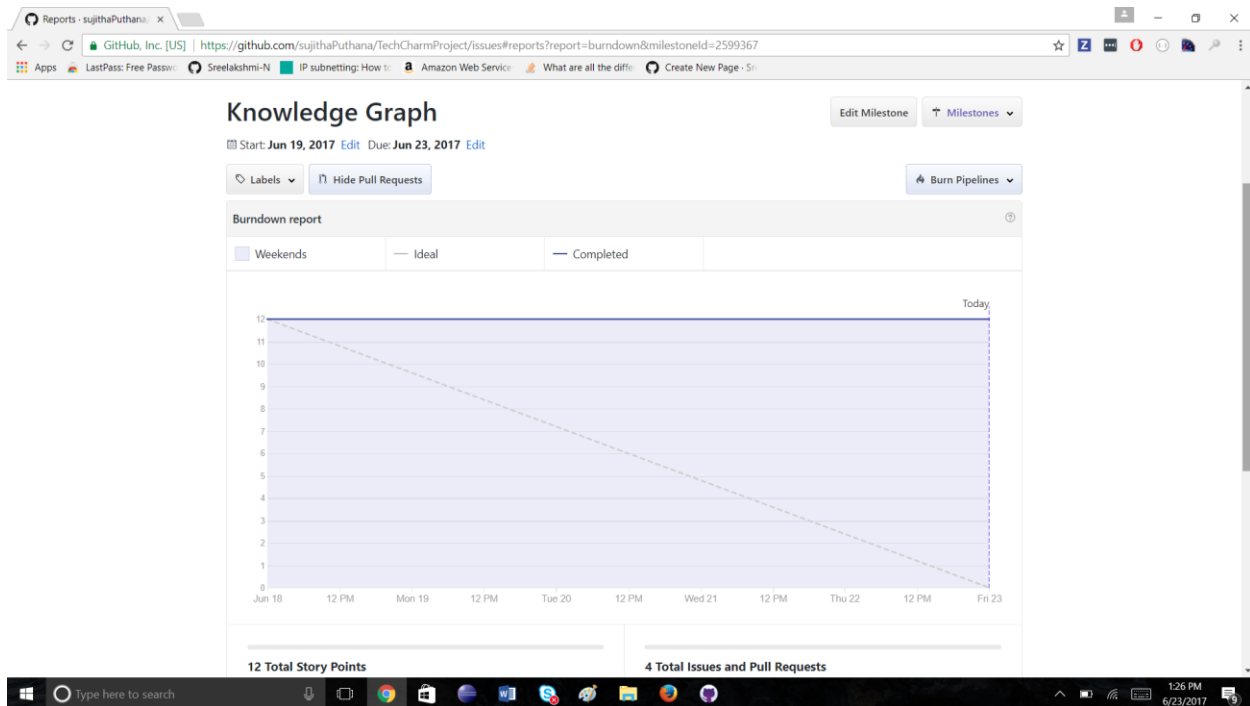
## Burn-Down Chart:

Burn-Down chart is created for the above issues via Milestones in GitHub. Below is the screenshot for more information,

# GitHub Wiki Page

The GitHub wiki page URL for the screenshots and the process flow is updated in the following link

- • https://github.com/sujithaPuthana/TechcharmProject

## 6.3 Work Completed

The completed tasks in this increment are,

- • Performed the NLP operations on the dataset.

- • Designing the question and answers for the NLP output.

- • TF-IDE and N-gram analysis

## 6.3 Future Work

❖ We need to implement the question and answer approach using the TF-IDF integrated with the NLP operations.

❖ Need to integrate the TF-IDF approach with the N-GRAM and Word2Vec for better performance.

## Bibliography

1. https://blog.algorithmia.com/introduction-natural-language-processing-nlp/

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

3. https://nlp.stanford.edu/