

Virtual Internship - 6.0

Milestone-2

Group-4

Vaibhavi Nimba Patil

Preetham Aditya

Vardhan

Ravi Kumar Rathlavat

- **Dataset collection and Preprocessing & NER** - Vaibhavi Patil
- **RE** - Ravi Kumar Rathlavat
- **Knowledge Graph** - Vardhan
- **Semantic Search** - Preetham Aditya
- Combined All these into a single working pipeline

NER

```
1 import pandas as pd
2 import spacy
3 from tqdm import tqdm
4
5 data = pd.read_csv("KMapFinal.csv")
6 nlp = spacy.load("en_core_web_sm")
7 tqdm.pandas()
8 def extract_entities(text):
9     doc = nlp(str(text))
10    return [(ent.text, ent.label_) for ent in doc.ents]
11
12 data["entities"] = data["article"].progress_apply(extract_entities)
13
14 data.to_csv("entities_extracted.csv", index=False)
15
16 print("NER extraction complete – saved as 'entities_extracted.csv'")
```

Relation Extraction

```
1  import pandas as pd
2  import spacy
3  from tqdm import tqdm
4
5  # Load dataset
6  data = pd.read_csv("KMapFinal.csv")
7  nlp = spacy.load("en_core_web_sm")
8  tqdm.pandas()
9  # Extract relations
10 def extract_relations(text):
11     doc = nlp(str(text))
12     triples = []
13     for sent in doc.sents:
14         subject, relation, object_ = None, None, None
15         for token in sent:
16             if "subj" in token.dep_:
17                 subject = token.text
18             if token.pos_ == "VERB":
19                 relation = token.lemma_
20             if "obj" in token.dep_:
21                 object_ = token.text
22         if subject and relation and object_:
23             triples.append((subject, relation, object_))
24     return triples
25 data["relations"] = data["article"].progress_apply(lambda x: extract_relations(str(x)))
26 data.to_csv("relations_output.csv", index=False)
27 print("✅ Relation extraction complete – saved as 'relations_output.csv'")
```

output

```
rgukt-basar@ravi-kumar-rathlavath:~/MileStone-Final-Task/new$ python3 entity_triplet_extractor.py
Loaded 5 sentences from dataset.
First few rows of data:
                                text
0  Albert Einstein developed the theory of relati...
1      Marie Curie discovered radium and polonium.
2      Isaac Newton formulated the laws of motion.
3  Charles Darwin proposed the theory of evolution.
4      Galileo Galilei improved the telescope.

Entities saved to: /home/rgukt-basar/MileStone-Final-Task/new/entities_extracted.csv
Triplets saved to: /home/rgukt-basar/MileStone-Final-Task/new/triplets_extracted.csv

Sample Entities:
Sentence: Albert Einstein developed the theory of relativity.
Entity: Albert Einstein (PERSON)

Sentence: Marie Curie discovered radium and polonium.
Entity: Marie Curie (PERSON)

Sentence: Isaac Newton formulated the laws of motion.
Entity: Isaac Newton (PERSON)

Sentence: Charles Darwin proposed the theory of evolution.
Entity: Charles Darwin (PERSON)

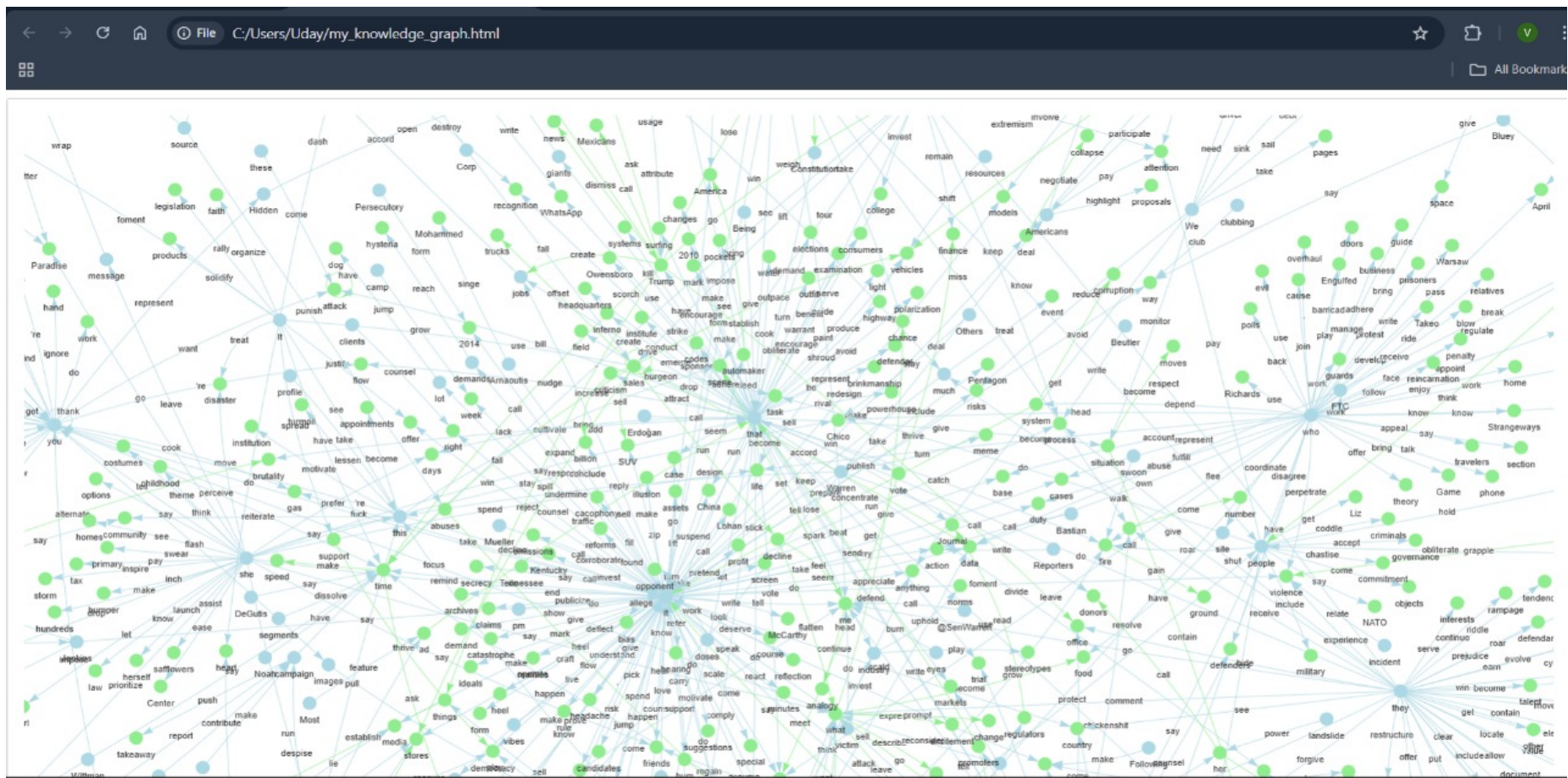
Sentence: Galileo Galilei improved the telescope.
Entity: Galileo Galilei (PRODUCT)

Sample Triplets:
Einstein develop relativity
Curie discover radium
Newton formulate motion
Darwin propose evolution
Galilei improve telescope
```

Knowledge Graph

```
1  import pandas as pd
2  import networkx as nx
3  from pyvis.network import Network
4  import webbrowser
5  # Load relation data
6  data = pd.read_csv("relations_output.csv").head(50)
7  # Extract triples
8  triples_list = []
9  for cell in data.iloc[:, -1]:
10     try:
11         tuples = eval(cell)
12         triples_list.extend(tuples)
13     except:
14         pass
15  print(f"✅ Extracted {len(triples_list)} triples from relations_output.csv")
16  # Build directed graph
17  G = nx.DiGraph()
18  for subj, rel, obj in triples_list:
19     G.add_node(subj, color="lightblue")
20     G.add_node(obj, color="lightgreen")
21     G.add_edge(subj, obj, label=rel)
22  print(f"✅ Graph built with {len(G.nodes())} nodes and {len(G.edges())} edges.\n")
23  # Visualize
24  net = Network(height="750px", width="100%", directed=True, notebook=False)
25  net.from_nx(G)
26  net.toggle_physics(True)
27  net.write_html("my_knowledge_graph.html")
28
29  print("🎉 Interactive Knowledge Graph saved as 'my_knowledge_graph.html'")
30  webbrowser.open(["my_knowledge_graph.html"])
```


output



Semantic Search

```
1  import pandas as pd
2  from sentence_transformers import SentenceTransformer, util
3
4  # Load dataset
5  data = pd.read_csv("KMapFinal.csv")
6  # Load model
7  model = SentenceTransformer('all-MiniLM-L6-v2')
8
9  # Encode sentences
10 sentences = data["article"].dropna().tolist()
11 embeddings = model.encode(sentences, convert_to_tensor=True)
12
13 # Query input
14 query = input("\n🗨 Enter a search query (e.g., 'AI technology'): ")
15 query_embedding = model.encode(query, convert_to_tensor=True)
16
17 # Compute similarity
18 cosine_scores = util.cos_sim(query_embedding, embeddings)
19 top_results = cosine_scores[0].argsort(descending=True)[:5]
20 |
21 # Display results
22 print("\n🔍 Top 5 Semantic Matches:")
23 for idx in top_results:
24     print(f"➡ {sentences[idx]} | Score: {float(cosine_scores[0][idx]):.3f}")
25
26 print("\n✅ Semantic search completed successfully!")
```