

Infosys Springboard Virtual Internship

PROJECT: KNOWMAP CROSS DOMAIN KNOWLEDGE MAPPING

Name: Janapati Lakshmi Sai Suvarnika
Date:30-10-2025

KNOWMAP CROSS DOMAIN KNOWLEDGE MAPPING

1. Introduction

In the modern digital world, a massive amount of information is stored in unstructured formats such as text documents, research papers, and online articles. Extracting meaningful knowledge from such large content repositories is a challenging task. Knowledge graphs have emerged as a powerful solution for representing and organizing information in a structured, machine-understandable way.

This project focuses on building an **intelligent knowledge graph system** that can extract entities and their relationships from multiple data sources, store them in a graph database, and visualize the connections interactively. The system also incorporates semantic search and feedback-based refinement to continually improve accuracy and relevance. A user authentication module and dataset management UI support secure and flexible data handling. The final outcome is a fully deployable application capable of supporting real user queries for knowledge discovery.

2. Objective

- ✓ Build a secure user interface for dataset upload and selection
- ✓ Automate entity and relation extraction using NLP models (spaCy/Transformers/BERT)
- ✓ Store extracted triples (Entity → Relation → Entity) in a knowledge graph database (e.g., Neo4j)
- ✓ Provide an interactive visual explorer using NetworkX or PyVis
- ✓ Integrate semantic search using Sentence Transformers
- ✓ Enable subgraph querying and search-driven graph expansion
- ✓ Provide admin-level monitoring, data correction, and feedback mechanism
- ✓ Deploy the final system using containerized technologies (e.g., Streamlit, Hugging Face Spaces)

3. Workflow

① User Authentication & Dataset Upload

- Users create accounts and log in securely
- Upload textual data (news, research papers, articles)
- Data format validation and storage

2 NLP Processing

- Named Entity Recognition (NER) identifies important entities (people, places, concepts)
- Relation Extraction models detect relationships between entities
- Entity–Relation–Entity triples are generated

3 Knowledge Graph Construction

- Triples are structured and inserted into a graph database
- Graph schema is maintained and continuously updated

4 Visualization & Semantic Search

- Users explore the knowledge graph interactively
- Semantic search retrieves related concepts using sentence embeddings
- Subgraph queries help users focus on specific knowledge areas

5 Feedback & Admin Tools

- Admin dashboard monitors extraction and graph correctness
- Manual correction for merging or editing nodes
- User feedback loop improves relevance over time

6 Deployment

- Application containerized and deployed for real usage
- Final project documentation and presentation prepared

4.Code

```
import os
import tempfile
import shutil
from pathlib import Path
from typing import List, Tuple, Dict
from collections import deque
from datetime import datetime
import streamlit as st
import pandas as pd
import networkx as nx
from pyvis.network import Network

try:
    import spacy
except Exception:
    spacy = None

try:
    from sentence_transformers import SentenceTransformer
    import numpy as np
```

```

from sklearn.metrics.pairwise import cosine_similarity
except Exception:
    SentenceTransformer = None
    np = None
    cosine_similarity = None

st.set_page_config(page_title="Streamlit App", layout="wide")
UPLOAD_FOLDER = "uploads"
os.makedirs(UPLOAD_FOLDER, exist_ok=True)

nlp = None
if spacy:
    try:
        nlp = spacy.load("en_core_web_sm")
    except Exception:
        try:
            nlp = spacy.blank("en")
        except Exception:
            nlp = None

MODEL = None
if SentenceTransformer is not None:
    try:
        MODEL = SentenceTransformer("all-MiniLM-L6-v2")
    except Exception:
        MODEL = None

SAMPLE_SENTENCES = [
    "Albert Einstein developed the theory of relativity in 1905.",
    "Marie Curie discovered polonium and radium in the early 1900s.",
    "Isaac Newton formulated the laws of motion.",
    "Nikola Tesla worked on alternating current systems.",
    "Alexander Fleming discovered penicillin in 1928.",
    "Charles Darwin proposed the theory of evolution.",
    "Rosalind Franklin contributed to the discovery of DNA structure.",
    "Ada Lovelace wrote the first algorithm for Charles Babbage's machine.",
    "Katherine Johnson calculated trajectories for NASA spacecraft.",
    "Tim Berners-Lee invented the World Wide Web.",
    "Alan Turing broke the Enigma code during World War II.",
    "Grace Hopper developed the first compiler for a programming language.",
]

EXAMPLE_TRIPLES = [

```

```

        ("Barack Obama", "born_in", "Honolulu"),
        ("Barack Obama", "position_held", "44th President of the United States"),
        ("Michelle Obama", "spouse", "Barack Obama"),
        ("Honolulu", "located_in", "Hawaii"),
        ("Hawaii", "country", "USA"),
        ("Theory of Relativity", "field", "Physics"),
        ("Physics", "related_to", "Mathematics"),
        ("Mathematics", "feeds", "Computer Science"),
        ("Computer Science", "includes", "Knowledge Graphs"),
        ("Knowledge Graphs", "used_for", "Semantic Search"),
    ]
}

def sanitize_filename(name: str) -> str:
    return Path(name).name

def save_uploaded_file(file_obj, username: str) -> str:
    user_dir = Path(UPLOAD_FOLDER) / username
    user_dir.mkdir(parents=True, exist_ok=True)
    dest = user_dir / sanitize_filename(file_obj.name)
    try:
        with open(dest, "wb") as f:
            f.write(file_obj.getbuffer())
    except Exception:
        file_obj.seek(0)
        with open(dest, "wb") as f:
            shutil.copyfileobj(file_obj, f)
    return str(dest)

def extract_triples_from_text(text: str) -> List[Tuple[str, str, str]]:
    """Basic subject-verb-object extraction using spaCy where available.
    Returns list of (subj, verb, obj)."""
    triples = []
    if nlp is None:
        for sent in text.split("."):
            sent = sent.strip()
            if not sent:
                continue
            parts = sent.split()
            if len(parts) >= 3:
                triples.append((parts[0], parts[1], " ".join(parts[2:])[::60]))
    return triples

doc = nlp(text)

```

```

for sent in doc.sents:
    root = None
    for token in sent:
        if token.dep_ == "ROOT":
            root = token
            break
    if root is None:
        continue

    subj = None
    obj = None
    for ch in root.lefts:
        if ch.dep_.endswith("subj"):
            subj = ch
            break
    for ch in root.rights:
        if ch.dep_.endswith("obj") or ch.dep_ == "pobj":
            obj = ch
            break

    subj_text = subj.text if subj is not None else None
    obj_text = obj.text if obj is not None else None
    verb_text = root.lemma_ if root is not None else None

    if subj_text and verb_text and obj_text:
        triples.append((subj_text, verb_text, obj_text))
    else:
        ents = [ent.text for ent in sent.ents]
        if len(ents) >= 2:
            triples.append((ents[0], verb_text or "related_to", ents[1]))

return triples

def extract_triples_from_sentences(sentences: List[str], min_triples: int = 10) -> List[Tuple[str, str, str]]:
    all_triples = []
    for s in sentences:
        t = extract_triples_from_text(s)
        for tr in t:
            if tr not in all_triples:
                all_triples.append(tr)
    if len(all_triples) < min_triples and nlp is not None:

```

```

ents_pool = []
for s in sentences:
    doc = nlp(s)
    for ent in doc.ents:
        ents_pool.append(ent.text)
i = 0
while len(all_triples) < min_triples and i + 1 < len(ents_pool):
    a = ents_pool[i]
    b = ents_pool[i + 1]
    candidate = (a, "related_to", b)
    if candidate not in all_triples:
        all_triples.append(candidate)
    i += 1
idx = 0
while len(all_triples) < min_triples and idx < len(EXAMPLE_TRIPLES):
    if EXAMPLE_TRIPLES[idx] not in all_triples:
        all_triples.append(EXAMPLE_TRIPLES[idx])
    idx += 1

return all_triples

def embed_sentences(sentences: List[str]):
    if MODEL is None:
        return None
    return MODEL.encode(sentences, convert_to_numpy=True)

def semantic_search(query: str, sentences: List[str], embeddings, top_k=5):
    if MODEL is None or embeddings is None:
        return []
    q_emb = MODEL.encode([query], convert_to_numpy=True)
    scores = cosine_similarity(q_emb, embeddings)[0]
    idx = list(np.argsort(scores)[::-1][:top_k])
    results = [(sentences[i], float(scores[i])) for i in idx]
    return results

def build_graph_from_triples(triples: List[Tuple[str, str, str]]) -> nx.Graph:
    G = nx.Graph()
    for s, p, o in triples:
        G.add_node(s, label=s)
        G.add_node(o, label=o)
        G.add_edge(s, o, relation=p)
    return G

```

```

def graph_to_pyvis(_G: nx.Graph, highlight_nodes=None, size_map=None) -> str:
    net = Network(height="650px", width="100%", bgcolor="#ffffff", font_color="black",
notebook=False)
    net.barnes_hut()
    for node in _G.nodes():
        title = node
        size = 15
        if size_map and node in size_map:
            size = size_map[node]
        color = "#7FB3FF"
        if highlight_nodes and node in highlight_nodes:
            color = "#FF6B6B"
            size = max(size, 25)
        net.add_node(node, label=title, title=title, color=color, size=size)
    for u, v, data in _G.edges(data=True):
        net.add_edge(u, v, title=data.get("relation", ""), label=None)
    net.set_options("""
var options = {
    "nodes": {"borderWidth": 2, "shape": "dot", "font": {"size": 14}},
    "physics": {"stabilization": {"enabled": true, "iterations": 200}}
}
""")
    tmp = tempfile.NamedTemporaryFile(delete=False, suffix=".html")
    net.write_html(tmp.name)
    return tmp.name

```

st.title("👁️ know map cross domian knowledge mapping using AI ")

```

if "logged_in" not in st.session_state:
    st.session_state["logged_in"] = None

```

```

tab1, tab2, tab3, tab4 = st.tabs(["User Workspace", "Extraction (Milestone-2)", "Visualization & Search (Milestone-3)", "Admin"])

```

with tab1:

```

st.header("🔗 Login & Dataset Management")
col1, col2 = st.columns(2)
with col1:
    username = st.text_input("Username", key="ui_username")
    password = st.text_input("Password", type="password", key="ui_password")
    if st.button("Login"):
        # lightweight local auth: accept any non-empty username for demo
        if username:

```

```

        st.session_state["logged_in"] = username
        st.success(f"Logged in as {username}")
    else:
        st.error("Enter a username to login")
    if st.button("Logout"):
        st.session_state.pop("logged_in", None)
        st.success("Logged out")
with col2:
    st.markdown("**Quick status**")
    st.write(f"Logged in as: {st.session_state.get('logged_in')}")

st.markdown("---")
st.subheader("📁 Upload Dataset (optional)")
up_col1, up_col2 = st.columns([2, 1])
with up_col1:
    uploaded = st.file_uploader("Upload CSV (for dataset validation only)", type=["csv"],
key="uploader_csv")
    with up_col2:
        up_user = st.text_input("Your username (same as login)",
value=st.session_state.get("logged_in") or "", key="up_user")
        if st.button("Upload File"):
            if not up_user:
                st.error("Provide username to save file")
            elif not uploaded:
                st.error("Attach a CSV file before uploading")
            else:
                path = save_uploaded_file(uploaded, up_user)
                st.success(f"File saved: {path}")

st.subheader("☑ Validate / 📺 Preview (optional)")
if st.button("Run Validation"):
    p = None
    if "last_uploaded_path" in st.session_state:
        p = st.session_state["last_uploaded_path"]
    if uploaded is not None:
        try:
            df = pd.read_csv(uploaded)
            if "grade" in df.columns:
                invalid = df[(df["grade"] < 0) | (df["grade"] > 100)]
                if not invalid.empty:
                    st.error("Validation error: grade out of range")
                else:
                    st.success("Dataset passed validation")
        except:
            st.error("Error reading the uploaded CSV file")

```

```

else:
    st.info("No 'grade' column detected — basic checks passed")
except Exception as e:
    st.error(f"Read error: {e}")
else:
    st.info("No CSV uploaded")

```

with tab2:

```

st.header("⌚ Milestone-2 — Extract Triples (min 10)")
st.markdown("This section uses a built-in set of sample sentences and extracts at least 10 triples from them.")

```

```

st.subheader("Sample Sentences (editable)")
sentences_text = st.text_area("Sentences (one per line)",
value="\n".join(SAMPLE_SENTENCES), height=220)

```

```

if st.button("Run Extraction"):
    user_sentences = [s.strip() for s in sentences_text.splitlines() if s.strip()]
    triples = extract_triples_from_sentences(user_sentences, min_triples=10)
    st.session_state["extracted_triples"] = triples
    st.success(f"Extracted {len(triples)} triples")
    st.write(triples)

```

```

if "extracted_triples" in st.session_state:
    st.info(f"Currently stored triples: {len(st.session_state['extracted_triples'])}")

```

with tab3:

```

st.header("🕸️ Milestone-3 — Visualization & Semantic Search (no CSV upload)")
st.markdown("Milestone-3 now uses only built-in/example triples derived from Milestone-2 extraction. CSV upload option has been removed from this tab.")

```

```

triples = st.session_state.get("extracted_triples") or EXAMPLE_TRIPLES
st.write(f"Using {len(triples)} triples for graph & search")

```

```

G = build_graph_from_triples(triples)
c1, c2, c3 = st.columns(3)
c1.metric("Nodes", G.number_of_nodes())
c2.metric("Edges", G.number_of_edges())
c3.metric("Avg Degree", round(sum(dict(G.degree()).values()) / max(1, G.number_of_nodes()), 2))

st.markdown("---")

```

```

sentences = [f"{s} {p} {o}" for s, p, o in triples]
embeddings = None
if MODEL is not None:
    embeddings = embed_sentences(sentences)
else:
    st.warning("SentenceTransformer not available — semantic search disabled")

query = st.text_input("Enter a natural language query to search triples:")
top_k = st.slider("Top K", 1, 10, 5)
depth = st.slider("Subgraph depth for visualization", 0, 3, 1)

highlight_nodes = []
size_map = {}

vis_path = None
if query and MODEL is not None and embeddings is not None:
    results = semantic_search(query, sentences, embeddings, top_k=top_k)
    st.subheader("Search Results")
    if results:
        df = pd.DataFrame(results, columns=["Triple Text", "Score"]) if isinstance(results[0], tuple) else pd.DataFrame(results)
        st.table(df)
        highlight_trip_nodes = [r[0].split()[0] for r in results] # rough pick of subjects
        highlight_nodes = highlight_trip_nodes
        size_map = {r[0].split()[0]: 30 - i * 4 for i, r in enumerate(results)}
        # auto subgraph around first match
        seed_label = results[0][0].split()[0]
        # generate subgraph
        def generate_subgraph(_G: nx.Graph, seeds: List[str], depth: int = 1, max_nodes: int = 200) -> nx.Graph:
            visited = set()
            q = deque([(s, 0) for s in seeds if s in _G])
            edges = []
            while q:
                node, d = q.popleft()
                if d >= depth:
                    continue
                for nbr in _G.neighbors(node):
                    edges.append((node, nbr))
                    if nbr not in visited:
                        visited.add(nbr)
                        q.append((nbr, d + 1))
            if len(visited) >= max_nodes:
                break
            return nx.Graph(edges)
        G = generate_subgraph(nx.Graph(), highlight_nodes, depth=depth)
        st.write(G)
    else:
        st.write("No results found for the given query.")



```

```

        break
    sub = nx.Graph()
    for u, v in edges:
        sub.add_node(u, **_G.nodes[u])
        sub.add_node(v, **_G.nodes[v])
        sub.add_edge(u, v, **_G.edges[u, v])
    return sub
subG = generate_subgraph(G, [seed_label], depth)
vis_path = graph_to_pyvis(subG, highlight_nodes, size_map)
else:
    st.info("No semantic results (or embeddings missing). Showing full graph.")
    vis_path = graph_to_pyvis(G)
else:
    vis_path = graph_to_pyvis(G)

try:
    with open(vis_path, "r", encoding="utf-8") as f:
        html = f.read()
    st.components.v1.html(html, height=700, scrolling=True)
except Exception as e:
    st.error(f"Unable to render graph: {e}")
finally:
    try:
        os.unlink(vis_path)
    except Exception:
        pass

```

with tab4:

```

st.header("⚡ Admin Dashboard")
stats = {
    "Total Entities": G.number_of_nodes(),
    "Total Relations": G.number_of_edges(),
    "Data Sources": 1,
    "Extraction Accuracy": "N/A"
}
col1, col2, col3, col4 = st.columns(4)
for col, (k, v) in zip([col1, col2, col3, col4], stats.items()):
    col.metric(label=k, value=v)

st.markdown("---")
st.subheader("Manual Graph Control")
node = st.text_input("Enter node name to edit (exact)", key="admin_node")
new_name = st.text_input("Enter new name", key="admin_new_name")

```

```

if st.button("Update Node"):
    if node in G:
        nx.relabel_nodes(G, {node: new_name}, copy=False)
        st.success(f"Node '{node}' renamed to '{new_name}'")
    else:
        st.error("Node not found")

    st.write("---")
    merge1 = st.text_input("Merge Node A", key="merge_a")
    merge2 = st.text_input("Merge Node B", key="merge_b")
    if st.button("Merge Nodes"):
        if merge1 in G and merge2 in G:
            for nbr in list(G.neighbors(merge2)):
                if not G.has_edge(merge1, nbr) and nbr != merge1:
                    G.add_edge(merge1, nbr, **G.edges[merge2, nbr])
            G.remove_node(merge2)
            st.success(f"Merged '{merge2}' into '{merge1}' ")
        else:
            st.error("One or both nodes not found")

    st.markdown("---")
    st.subheader("Feedback")
    rating = st.slider("Graph Relevance", 1, 5, 3, key="fb_rating")
    comment = st.text_input("Comment", key="fb_comment")
    if st.button("Submit Feedback"):
        st.success("Feedback submitted ")

```

5. Code Explanation

This code creates a Streamlit web application for:

- Logging in users
- Uploading datasets
- Extracting knowledge (triples) from sentences
- Building and visualizing a Knowledge Graph
- Allowing semantic search on extracted knowledge
- Admin tools for editing graph nodes

6.Output Screenshots:

🔑 Login & Dataset Management

Username
suvarnika

Password
....

Login

Logged in as suvarnika

Quick status
Logged in as: suvarnika

Sample Sentences (editable)

Sentences (one per line)

Alexander Fleming discovered penicillin in 1928.
Charles Darwin proposed the theory of evolution.
Rosalind Franklin contributed to the discovery of DNA structure.
Ada Lovelace wrote the first algorithm for Charles Babbage's machine.
Katherine Johnson calculated trajectories for NASA spacecraft.
Tim Berners-Lee invented the World Wide Web.
Alan Turing broke the Enigma code during World War II.
Grace Hopper developed the first compiler for a programming language.

Run Extraction

Extracted 10 triples

📁 Upload Dataset (optional)

Upload CSV (for dataset validation only)



Drag and drop file here

Limit 200MB per file • CSV

Browse files



menu.csv 29.3KB

X

```

▼ 0 : [
  0 : "Einstein"
  1 : "develop"
  2 : "theory"
]
▼ 1 : [
  0 : "Curie"
  1 : "discover"
  2 : "polonium"
]
▼ 2 : [
  0 : "Newton"
  1 : "formulate"
  2 : "laws"
]
▼ 3 : [
  0 : "Fleming"
  1 : "discover"
  2 : "penicillin"
]
▼ 4 : [ 🌐
  0 : "Darwin"
  1 : "propose"
  2 : "theory"
]

```

🕸 Milestone-3 – Visualization & Semantic Search (no CSV upload)

Milestone-3 now uses only built-in/example triples derived from Milestone-2 extraction. CSV upload option has been removed from this tab.

Using 10 triples for graph & search

Nodes	Edges	Avg Degree
19	10	1.05

Enter a natural language query to search triples:

theory

Top K

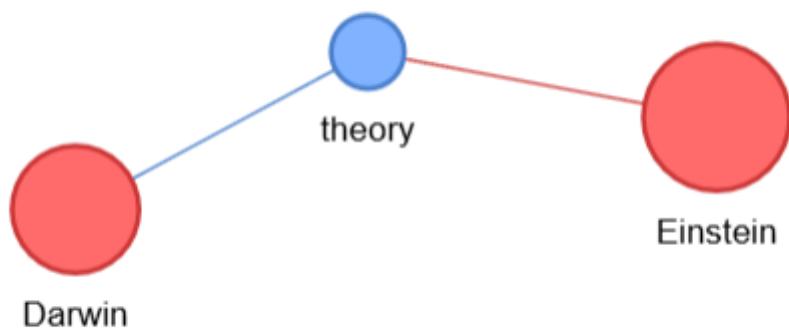
5

Subgraph depth for visualization

3

Search Results

	Triple Text	Score
0	Einstein develop theory	0.5319
1	Darwin propose theory	0.5266
2	Newton formulate laws	0.2308
3	Lee invent Web	0.2180
4	Curie discover polonium	0.1547



Manual Graph Control

Enter node name to edit (exact)

Einstein

Enter new name

theory

[Update Node](#)

Node 'Einstein' renamed to 'theory'

Merge Node A

Albert

Merge Node B

Einstein

Merge Nodes

Feedback

Graph Relevance



Comment

theory

Submit Feedback

Feedback submitted

🛠 Admin Dashboard

Total Entities

19

Total Relations

10

Data Sources

1

Extraction Accuracy

N/A

7. Conclusion:

The Known Map Cross-Domain Mapping System successfully transforms unstructured text into structured, meaningful knowledge. By extracting Subject-Relation-Object triples and converting them into a visual knowledge graph, the system provides a more intuitive way to understand relationships between concepts. The integration of semantic search enables users to retrieve information based on meaning rather than exact keywords, improving the accuracy and usability of knowledge retrieval.

The system also includes administrative tools to refine and maintain the correctness of mapped knowledge, making it scalable for various academic and industrial applications. Overall, this project demonstrates an effective method for cross-domain knowledge representation, intelligent information retrieval, and interactive visualization, contributing towards smarter knowledge management solutions.

