

# Table of Contents

<b>Table of Contents</b>	<b>1</b>
<b>Project Overview</b>	<b>1</b>
<b>Milestone 2</b>	<b>2</b>
Dataset	2
Virtual Environment Setup Process	2
Screenshots of the running instances and logs:	2
Containerized Components	5
1. Container for data scraping/fetching:	6
Instructions for Running:	6
2. Container for data embedding, storing, and other preprocessing:	6
Instructions for Running:	7
3. Container for user prompt processing through a RAG model:	7
Instructions for Running:	8
Orchestration Using Docker-Compose	8
Versioned Data Strategy (Planned, not yet done)	9
Utilizing LLM (Large Language Model)	9
Mock-up of the Application	11

# Project Overview

**Title:** Global AI Colab For Good

## **Members:**

Rama Edlabadkar (ramasandeepedlabadkar@g.harvard.edu),  
Shanzeh Batool (shanzehbatool@g.harvard.edu),  
Labdh Gandhi (labdhigandhi@g.harvard.edu),  
Hinal Jajal (hjajal@g.harvard.edu)

**Project:** The scope of this project is to build a global platform that links AI research groups with organizations aiming to solve social issues using AI. The platform will have a search interface for organizations to look for AI research papers relevant to their social cause. A dashboard will provide a curated list of relevant research to the user prompt, the research groups, and how the research work relates to the user's problem prompt. The platform will be designed to support a growing number of research groups and global organizations. We process a large corpus of AI

research papers & social issue descriptions and train LLMs for information retrieval and matching between research and real-world problems.

## Milestone 2

### Dataset

We utilized textual data fetched from social impact-related papers obtained using the ArXiv API, then embedded it using the Hugging Face library and stored it in the vector database. Currently, we fetch and store the manuscripts (both raw and embedded) for around 30 papers.

### Virtual Environment Setup Process

We use Docker containers to separate tasks like data scraping, preprocessing, embedding generation, and user prompt processing via the RAG model. This setup allows for easy management and scalability of different services. Each container runs specific processes, which can be updated, modified, or scaled independently based on requirements. We have set up the virtual environment using Docker to support containerized components. This will ensure that all project elements, from data scraping to large language model (LLM) processing, remain isolated and easy to integrate.

### Screenshots of the running instances and logs

#### **Container for data scraping/fetching (retrieve\_papers)**

```

Locking [dev-packages] dependencies...
Updated Pipfile.lock (073193008382e006cc494452dfcf769de6e5f7849cd8b28bcbf0e1b27a69309) !
Locking [build-dependencies] dependencies...
Updated Pipfile.lock (073193008382e006cc494452dfcf769de6e5f7849cd8b28bcbf0e1b27a69309) !
[+] Building 30s (11/12) FINISHED
   =+ [internal] load build definition from Dockerfile
   => [internal] transfering dockerfile: 878B
   => [internal] load metadata for docker.io/library/python:3.12-slim-bookworm
   => [auth] library/python:pull token for registry-1.docker.io
   => [internal] load dockerignore
   => [internal] load context
   => [internal] CACHED [1/6] FROM docker.io/library/python:3.12-slim-bookworm@sha256:032c52613401895aa3d418a4c563d2d05f993bc3ecc065c8f4e2280978acd249
   => [internal] load build context
   => transferring context: 36.38kB
   => [2/6] RUN set -ex; apt-get update && apt-get upgrade -y && apt-get install -y --no-install-recommends build-essential git ffmpeg && pip install --no-cache-dir --upgrade pip &&
   docker:desktop-linux
   => [3/6] WORKDIR /app
   => [internal] ADD Pipfile Pipfile.lock /app/
   => [5/6] COPY . /app
   => [6/6] ADD /app
   => exporting to image
   => exporting layers
   => writing image sha256:bca678aa590daf9e96596b9ec65a4fb9eacaeedc1ead2b553105e7655a66aa1
   => naming to docker.io/library/retrieve_papers
View build details: docker-desktop://dashboard/build/desktop_linux/desktop-linux/riyekxoiuhxyujchcc6ne4np

1 warning found (use docker --debug to expand):
- SecretsUsedInArgOrEnv: Do not use ARG or ENV instructions for sensitive data (ENV "GOOGLE_APPLICATION_CREDENTIALS") (line 8)

What's next:
  View a summary of image vulnerabilities and recommendations > docker scan quickview
Lambdigandhi@dhcp-10-250-43-94: ~retrieve_papers % docker run --rm -ti -v $(pwd):/app embed_papers
/bin/bash: ./docker-entrypoint.sh: No such file or directory
Lambdigandhi@dhcp-10-250-43-94: ~retrieve_papers % docker run --rm -ti -v $(pwd):/app retrieve_papers
Launching subshell in virtual environment...
root@419b27dca38:/app# source /root/.local/share/virtualenvs/app-4PlAip00/bin/activate
(app) root@419b27dca38:/app# python3 -m embed_papers.py
File manuscript_texts_to_retrieve.txt uploaded to paper-rec-bucket.
Processing http://arxiv.org/src/astro-ph/9704140v2...
Downloaded .tar file to downloads/astro-ph_9704140v2.tar
/app/retrieve_papers.py:77: DeprecationWarning: Python 3.14 will, by default, filter extracted tar archives and reject files or modify their metadata. Use the filter argument to control this behavior.
tar.extractall(path=extract_path)
Extracted contents to downloads/astro-ph_9704140v2/extracted
Found .tex files in downloads/astro-ph_9704140v2/extracted/G347_5_8_S_accept_1.tex
Saved .tex file as manuscript_texts_to_retrieve/astro-ph_9704140v2.txt
File manuscript_texts_to_retrieve/astro-ph_9704140v2.txt uploaded to paper-rec-bucket.
Processing http://arxiv.org/src/cs/9903013v1...
Downloaded .tar file to downloads/cs_9903013v1.cs_9903013v1.tar
Extracted contents to downloads/cs_9903013v1/extracted
Found .tex file: downloads/cs_9903013v1/extracted/text.tex
Saved .tex file: downloads/cs_9903013v1/extracted/text.tex
File manuscript_texts_to_retrieve/cs_9903013v1.txt uploaded to paper-rec-bucket.
Processing http://arxiv.org/src/quant-ph/0003132v1...
Downloaded .tar file to downloads/quant-ph_0003132v1/quant-ph_0003132v1.tar
Extracted contents to downloads/quant-ph_0003132v1/extracted

```

```

Saved .tex file as manuscript_texts_to_retrieve/cond-mat_0503607v2.txt
File manuscript_texts_to_retrieve/cond-mat_0503607v2.txt uploaded to paper-rec-bucket.
Processing http://arxiv.org/src/cs/0502056v2...
Downloaded .tar file to downloads/cs_0502056v2/cs_0502056v2.tar
Extracted contents to downloads/cs_0502056v2/extracted
Found .tex file: downloads/cs_0502056v2/extracted/jpm-cauthor.tex
Saved .tex file as manuscript_texts_to_retrieve/cs_0502056v2.tex
File manuscript_texts_to_retrieve/cs_0502056v2.txt uploaded to paper-rec-bucket.
Processing http://arxiv.org/src/physics/0508029v1...
Downloaded .tar file to downloads/physics_0508029v1/physics_0508029v1.tar
Extracted contents to downloads/physics_0508029v1/extracted
Found .tex file: downloads/physics_0508029v1/extracted/contrarians_v6.tex
Saved .tex file: downloads/physics_0508029v1/extracted/physics_0508029v1.txt
File manuscript_texts_to_retrieve/physics_0508029v1.txt uploaded to paper-rec-bucket.
Processing http://arxiv.org/src/cs/0409019v2...
Failed to download the file from http://arxiv.org/src/cs/0409019v2. HTTP Status Code: 403
Skipping http://arxiv.org/src/cs/0409019v2 due to download error.
Skipping upload for http://arxiv.org/src/cs/0409019v2 due to processing errors.
Processing http://arxiv.org/src/cs/0511011v1...
Downloaded .tar file to downloads/cs_0511011v1/cs_0511011v1.tar
Extracted contents to downloads/cs_0511011v1/extracted
Found .tex file: downloads/cs_0511011v1/extracted/link8impact.tex
Saved .tex file as manuscript_texts_to_retrieve/cs_0511011v1.txt
File manuscript_texts_to_retrieve/cs_0511011v1/cs_0601039v1.tar
Processing http://arxiv.org/src/cs/0601039v1...
Downloaded .tar file to downloads/cs_0601039v1/cs_0601039v1.tar
Extracted contents to downloads/cs_0601039v1/extracted
Found .tex file: downloads/cs_0601039v1/extracted/q-bio_0601039v1.tex
Saved .tex file as manuscript_texts_to_retrieve/cs_0601039v1.txt
File manuscript_texts_to_retrieve/cs_0601039v1/cs_0601032v1.tar
Processing http://arxiv.org/src/q-bio/0601032v1...
Saved .tex file as manuscript_texts_to_retrieve/cs_0601032v1.txt
File manuscript_texts_to_retrieve/cs_0601032v1.txt uploaded to paper-rec-bucket.
Processing http://arxiv.org/src/q-bio/0601032v1...
Downloaded .tar file to downloads/q-bio_0601032v1/q-bio_0601032v1.tar
Extracted contents to downloads/q-bio_0601032v1/extracted
Found .tex file: downloads/q-bio_0601032v1/extracted/ising_jn.tex
Saved .tex file: downloads/q-bio_0601032v1/extracted/physics_0603248v1.txt
File manuscript_texts_to_retrieve/q-bio_0601032v1.txt
File manuscript_texts_to_retrieve/q-bio_0601032v1.txt uploaded to paper-rec-bucket.
Processing http://arxiv.org/src/physics/0603248v1...
Failed to download the file from http://arxiv.org/src/cs/0802018v2. HTTP Status Code: 403
Skipping http://arxiv.org/src/cs/0802018v2 due to download error.
Skipping upload for http://arxiv.org/src/cs/0802018v2 due to processing errors.
Processing http://arxiv.org/src/physics/0605129v1...
Downloaded .tar file to downloads/physics_0605129v1/physics_0605129v1.tar
downloads/physics_0605129v1/physics_0605129v1.tar is not a valid tar file.
Skipping http://arxiv.org/src/physics/0605129v1 due to extraction error.
Skipping upload for http://arxiv.org/src/physics/0605129v1 due to processing errors.
Processing http://arxiv.org/src/cs/0306128v3...
Downloaded .tar file to downloads/cs_0306128v3/cs_0306128v3.tar
downloads/cs_0306128v3/cs_0306128v3.tar is not a valid tar file.
Skipping http://arxiv.org/src/cs/0306128v3 due to extraction error.
Skipping upload for http://arxiv.org/src/cs/0306128v3 due to processing errors.
Processing http://arxiv.org/src/physics/0608296v1...
Downloaded .tar file to downloads/physics_0608296v1/physics_0608296v1.tar
Extracted contents to downloads/physics_0608296v1/extracted
Found .tex file: downloads/physics_0608296v1/extracted/L5.tex
Saved .tex file as manuscript_texts_to_retrieve/physics_0608296v1.txt
File manuscript_texts_to_retrieve/physics_0608296v1.txt uploaded to paper-rec-bucket.
(=) root@419b27dca38:/app#
```

## Container for data embedding, storing, and other preprocessing (embed\_papers)

```

labdhigandhi@dhcp-10-250-43-94 ~% docker build --no-cache -t embed_papers .
+ Building 79.0s (11/11) FINISHED
--> [internal] load build definition from Dockerfile
--> [internal] transferring dockerfile: 1.60kB
--> [internal] load metadata for docker.io/library/python:3.12-slim-bookworm
--> [internal] load .dockerignore
--> [internal] transferring context: 28
--> CACHED [1/6] FROM docker.io/library/python:3.12-slim-bookworm@sha256:032c52613401895aa3d418a4c563d2d05f993bc3ecc065c8f4e2280978acd249
--> [internal] load build context
--> [internal] transferring context: 262.35kB
--> [5/6] RUN    for i in $(seq 1 8); do mkdir -p "/usr/share/man/man${i}"; done && apt-get update && apt-get upgrade -y && apt-get install -y --no-install-recommends build-essential
--> [6/6] WORKDIR /app
--> [4/6] ADD --chown=app:app Pipfile Pipfile.lock /app
--> [5/6] RUN pipenv sync
--> [6/6] ADD --chown=app:app . /app
--> exporting layers
--> writing image sha256:31c7b663d84dafe524ab22a76a52a5626a0fe5deb3e0aa725cad5ea3abb73da
--> naming to docker.io/library/embed_papers

View build details: docker-desktop://dashboard/build/desktop-linux/desktop-linux/vxzhpl5shpfm0scf96pxpwr91

What's next:
  View a summary of image vulnerabilities and recommendations + docker scout quickview
labdhigandhi@dhcp-10-250-43-94 ~% docker run --rm -ti -v "$PWD":/app embed_papers
Container is running!!!

Launching subshell in virtual environment...
app@9bb0012bb05a:/app$ source /home/app/.local/share/virtualenvs/app-4PlAip0Q/bin/activate
(app) app@9bb0012bb05a:/app$ python embed_papers.py
Error: failed to open file PR_SVE_GET_VL failed
modules.json: 100%
config_sentence_transformers.json: 100%
README.md: 100%
sentence_bert_config.json: 100%
config.json: 100%
model_sentence_transformers: 100%
tokenizer_config.json: 100%
vocab.txt: 100%
tokenizer.json: 100%
special_tokens_map.json: 100%
l_Pooling/config.json: 100%
Downloaded manuscript_texts_to_retrieve/astro-ph_9704140v2.txt from paper-rec-bucket to /tmp/astro-ph_9704140v2.txt
Deleted manuscript_texts_to_retrieve/astro-ph_9704140v2.txt
Created a chunk of size 1088, which is longer than the specified 1000
Created a chunk of size 1341, which is longer than the specified 1000
Created a chunk of size 1480, which is longer than the specified 1000
Created a chunk of size 2282, which is longer than the specified 1000
Created a chunk of size 1438, which is longer than the specified 1000
Created a chunk of size 2470, which is longer than the specified 1000
Created a chunk of size 1080, which is longer than the specified 1000
Created a chunk of size 1054, which is longer than the specified 1000
Moved manuscript_texts_to_retrieve/astro-ph_9704140v2.txt to manuscript_texts_done/astro-ph_9704140v2.txt in the cloud
Deleted manuscript_texts_to_retrieve/astro-ph_9704140v2.txt from Google Cloud Storage
Deleted local temp file /tmp/astro-ph_9704140v2.txt
Downloaded manuscript_texts_to_retrieve/cond-mat_00804026v1.txt from paper-rec-bucket to /tmp/cond-mat_00804026v1.txt
Moved manuscript_texts_to_retrieve/cond-mat_00804026v1.txt to manuscript_texts_done/cond-mat_00804026v1.txt in the cloud
Deleted manuscript_texts_to_retrieve/cond-mat_00804026v1.txt from Google Cloud Storage
Deleted local temp file /tmp/cond-mat_00804026v1.txt
Downloaded manuscript_texts_to_retrieve/cond-mat_0111070v1.txt from paper-rec-bucket to /tmp/cond-mat_0111070v1.txt
Created a chunk of size 1082, which is longer than the specified 1000
Created a chunk of size 1528, which is longer than the specified 1000
Created a chunk of size 1166, which is longer than the specified 1000
Created a chunk of size 1198, which is longer than the specified 1000
[green checkmark] Scripts have been embedded and uploaded to the cloud.

[green checkmark] app@9bb0012bb05a:/app$ █
```

  

```

Created a chunk of size 1785, which is longer than the specified 1000
Created a chunk of size 1080, which is longer than the specified 1000
Created a chunk of size 1049, which is longer than the specified 1000
Created a chunk of size 2956, which is longer than the specified 1000
Created a chunk of size 2476, which is longer than the specified 1000
Created a chunk of size 1235, which is longer than the specified 1000
Created a chunk of size 3222, which is longer than the specified 1000
Created a chunk of size 1179, which is longer than the specified 1000
Created a chunk of size 1080, which is longer than the specified 1000
Created a chunk of size 1087, which is longer than the specified 1000
Created a chunk of size 1912, which is longer than the specified 1000
Created a chunk of size 1134, which is longer than the specified 1000
Moved manuscript_texts_to_retrieve/physics_0502047v1.txt to manuscript_texts_done/physics_0502047v1.txt in the cloud
Deleted manuscript_texts_to_retrieve/physics_0502047v1.txt from Google Cloud Storage
Deleted local temp file /tmp/physics_0502047v1.txt
Downloaded manuscript_texts_to_retrieve/physics_05080029v1.txt from paper-rec-bucket to /tmp/physics_05080029v1.txt
Created a chunk of size 1194, which is longer than the specified 1000
Created a chunk of size 1026, which is longer than the specified 1000
Created a chunk of size 1163, which is longer than the specified 1000
Moved manuscript_texts_to_retrieve/physics_05080029v1.txt to manuscript_texts_done/physics_05080029v1.txt in the cloud
Deleted manuscript_texts_to_retrieve/physics_05080029v1.txt from Google Cloud Storage
Deleted local temp file /tmp/physics_05080029v1.txt
Downloaded manuscript_texts_to_retrieve/physics_0603248v1.txt from paper-rec-bucket to /tmp/physics_0603248v1.txt
Created a chunk of size 1339, which is longer than the specified 1000
Created a chunk of size 3539, which is longer than the specified 1000
Moved manuscript_texts_to_retrieve/physics_0603248v1.txt to manuscript_texts_done/physics_0603248v1.txt in the cloud
Deleted manuscript_texts_to_retrieve/physics_0603248v1.txt from Google Cloud Storage
Deleted local temp file /tmp/physics_0603248v1.txt
Downloaded manuscript_texts_to_retrieve/physics_0608296v1.txt from paper-rec-bucket to /tmp/physics_0608296v1.txt
Created a chunk of size 1195, which is longer than the specified 1000
Created a chunk of size 1026, which is longer than the specified 1000
Created a chunk of size 1211, which is longer than the specified 1000
Created a chunk of size 1098, which is longer than the specified 1000
Created a chunk of size 1558, which is longer than the specified 1000
Moved manuscript_texts_to_retrieve/physics_0608296v1.txt to manuscript_texts_done/physics_0608296v1.txt in the cloud
Deleted manuscript_texts_to_retrieve/physics_0608296v1.txt from Google Cloud Storage
Deleted local temp file /tmp/physics_0608296v1.txt
Downloaded manuscript_texts_to_retrieve/q-bio_0601032v1.txt from paper-rec-bucket to /tmp/q-bio_0601032v1.txt
Created a chunk of size 4180, which is longer than the specified 1000
Created a chunk of size 1443, which is longer than the specified 1000
Created a chunk of size 1092, which is longer than the specified 1000
Created a chunk of size 1200, which is longer than the specified 1000
Created a chunk of size 1492, which is longer than the specified 1000
Created a chunk of size 1727, which is longer than the specified 1000
Created a chunk of size 1954, which is longer than the specified 1000
Created a chunk of size 2349, which is longer than the specified 1000
Created a chunk of size 1103, which is longer than the specified 1000
Created a chunk of size 1333, which is longer than the specified 1000
Created a chunk of size 1215, which is longer than the specified 1000
Created a chunk of size 2407, which is longer than the specified 1000
Moved manuscript_texts_to_retrieve/q-bio_0601032v1.txt to manuscript_texts_done/q-bio_0601032v1.txt in the cloud
Deleted manuscript_texts_to_retrieve/q-bio_0601032v1.txt from Google Cloud Storage
Deleted local temp file /tmp/q-bio_0601032v1.txt
Uploaded paper_vector_db/36fbff9c5-4be9-40fc-a652-a9fa85baa2c1/data_level0.bin to paper-rec-bucket
Uploaded paper_vector_db/36fbff9c5-4be9-40fc-a652-a9fa85baa2c1/length.bin to paper-rec-bucket
Uploaded paper_vector_db/36fbff9c5-4be9-40fc-a652-a9fa85baa2c1/link_lists.bin to paper-rec-bucket
[green checkmark] Scripts have been embedded and uploaded to the cloud.

[green checkmark] app@9bb0012bb05a:/app$ █
```

## Container for user prompt processing through a RAG model (perform\_rag)

```

labdhigandhi@dhcp-10-250-43-94 perform_rag % pipenv lock --clear
Creating a virtualenv for this project
Pipfile.lock is being loaded from /Users/labdhigandhi/.local/share/virtualenvs/perform_rag/lib/pipfile
Using /Library/Frameworks/Python.framework/Versions/3.12/bin/python3.12.2 to create virtualenv...
Creating virtual environment...created virtual environment CPython3.12.2.final.0-156ms
creator CPython3Posix(dest=/Users/labdhigandhi/.local/share/virtualenvs/perform_rag-ariaCtcl, clear=False, no_vcs_ignore=False, global=False)
seeded FromAppData(download=False, pip=bundle, via=copy, app_data_dir=/Users/labdhigandhi/Library/Application Support/virtualenv)
added seed packages: pip==24.2
activators BashActivator,CShellActivator,FishActivator,NushellActivator,PowerShellActivator,PythonActivator

✓ Successfully created virtual environment!
Virtualenv location: /Users/labdhigandhi/.local/share/virtualenvs/perform_rag-ariaCtcl
Locking [packages] dependencies...
Building requirements...
Resolving dependencies...
└ Success!
Locking [dev-packages] dependencies...
Updated Pipfile.lock (d374ef7193d3ca2c51df33cb832b40073fd80579db67fa4b762877506be07e)!

labdhigandhi@dhcp-10-250-43-94 perform_rag % docker build --no-cache -t perform_rag .
[+] Building 82.1s (12/12) FINISHED
   => [internal] load build definition from Dockerfile
   => transferring dockerfile: 88B
   => [internal] load metadata for docker.io/library/python:3.12-slim-bookworm
   => [auth] library:authenticating token for registry-1.docker.io
   => [internal] load dockerignore
   =>> transferring context: 2B
   => CACHED [1/6] FROM docker.io/library/python:3.12-slim-bookworm@sha256:032c52613401895aa3d418a4c563d2d085f993bc3ecc065c8f4e2280978acd249
   => [internal] load .envfile context
   =>> transferring context: 2B+48KB
   => [1/6] RUN apt-get update && apt-get upgrade -y && apt-get install -y --no-install-recommends build-essential git ffmpeg && pip install --no-cache-dir --upgrade pip &&
      32.4s
   => [2/6] RUN apt-get update && apt-get upgrade -y && apt-get install -y --no-install-recommends build-essential git ffmpeg && pip install --no-cache-dir --upgrade pip &&
      0.0s
   => [3/6] RUN apt-get update && apt-get upgrade -y && apt-get install -y --no-install-recommends build-essential git ffmpeg && pip install --no-cache-dir --upgrade pip &&
      0.0s
   => [4/6] ADD Pipfile Pipfile.lock /app/
   => [5/6] RUN pipenv sync
   => [6/6] ADD . /app
   => exporting to image
   => exporting layers
   =>> writing image sha256:cd4a92d0d83c33bab0cf26f474f7614973bcd02f3e8c136316048ee81b521
   =>> naming to docker.io/library/perform_rag

View build details: docker-desktop://dashboard/build/desktop-linux/desktop-linux/@kuuijsoh538139q3zajh7hou

1 warning found (use docker --debug to expand):
- SecretsUsedInArgOrEnv: Do not use ARG or ENV instructions for sensitive data (ENV "GOOGLE_APPLICATION_CREDENTIALS") (line 8)

What's next:
View a summary of image vulnerabilities and recommendations → docker scout quickview
labdhigandhi@dhcp-10-250-43-94 perform_rag % docker run --rm -ti -v "$PWD":/app perform_rag
Launching subshell in virtual environment...

```

  

```

/app$ root@ad71ca2cfa0f:/app# ./app.py
Downloaded paper_vector_db/36fb9c5-4be9-4fc-a652-a9fa85baa2c1/data_level0.bin to paper_vector_db/36fb9c5-4be9-4fc-a652-a9fa85baa2c1/header.bin
Downloaded paper_vector_db/36fb9c5-4be9-4fc-a652-a9fa85baa2c1/link.length.bin to paper_vector_db/36fb9c5-4be9-4fc-a652-a9fa85baa2c1/length.bin
Downloaded paper_vector_db/36fb9c5-4be9-4fc-a652-a9fa85baa2c1/link.lists.bin to paper_vector_db/36fb9c5-4be9-4fc-a652-a9fa85baa2c1/link_lists.bin
Downloaded paper_vector_db/chroma.sqlite3 to paper_vector_db/chroma.sqlite3
Error in cpufifo: prctl(PR_SVE_GET_VL) failed
modules.json: 100%
config_transformers.json: 100%
Readme: 100%
preference_hart_config.json: 100%
config.json: 100%
model_safetensors: 100%
tokenizer_config.json: 100%
vocab.txt: 100%
tokenizer.json: 100%
special_token_map.json: 100%
pretrained/config.json: 100%
While the provided papers focus on social network analysis and complex systems, they offer a tangential relevance to "AI for social impact" by highlighting the complexity of social interactions and the importance of understanding social dynamics for designing effective AI applications. Here's how:
* **Modeling Social Phenomena:** The papers demonstrate the use of models from theoretical physics to study complex social phenomena like opinion formation, disease spreading, and population dynamics. This approach emphasizes the importance of understanding the underlying mechanisms of social interaction, which is crucial for developing AI systems that can effectively interact with and influence human behavior.
* **The Importance of Social Networks:** The papers stress the role of social networks in shaping individual and collective behavior. They illustrate the impact of key actors and their influence on the spread of information and opinions within a network. This understanding is crucial for designing AI systems that can navigate and leverage social networks for social impact.
* **Multi-Contextual Interactions:** One paper emphasizes the importance of understanding multi-contextual interactions in social networks. This perspective suggests that AI systems designed for social impact need to consider the diverse contexts in which people interact and the different roles they play within those contexts.
* **The Realist Approach:** Another paper discusses the importance of considering the subjective experiences and perceptions of individuals within a social network. This highlights the need for AI systems that can account for the complexities of human perception and interpretation, particularly when aiming for social impact.

Overall, while the papers don't explicitly address "AI for social impact," they provide valuable insights into the complex dynamics of social interactions and the need to understand these dynamics for developing effective AI applications. The papers suggest that AI systems designed for social impact should incorporate models of social networks, account for multi-contextual interactions, and consider the subjective experiences of individuals within those networks.
While the provided papers focus on social network analysis and complex systems, they offer a tangential relevance to "AI for social impact" by highlighting the complexity of social interactions and the importance of understanding social dynamics for designing effective AI applications. Here's how:
* **Modeling Social Phenomena:** The papers demonstrate the use of models from theoretical physics to study complex social phenomena like opinion formation, disease spreading, and population dynamics. This approach emphasizes the importance of understanding the underlying mechanisms of social interaction, which is crucial for developing AI systems that can effectively interact with and influence human behavior.
* **The Importance of Social Networks:** The papers stress the role of social networks in shaping individual and collective behavior. They illustrate the impact of key actors and their influence on the spread of information and opinions within a network. This understanding is crucial for designing AI systems that can navigate and leverage social networks for social impact.
* **Multi-Contextual Interactions:** One paper emphasizes the importance of understanding multi-contextual interactions in social networks. This perspective suggests that AI systems designed for social impact need to consider the diverse contexts in which people interact and the different roles they play within those contexts.
* **The Realist Approach:** Another paper discusses the importance of considering the subjective experiences and perceptions of individuals within a social network. This highlights the need for AI systems that can account for the complexities of human perception and interpretation, particularly when aiming for social impact.

Overall, while the papers don't explicitly address "AI for social impact," they provide valuable insights into the complex dynamics of social interactions and the need to understand these dynamics for developing effective AI applications. The papers suggest that AI systems designed for social impact should incorporate models of social networks, account for multi-contextual interactions, and consider the subjective experiences of individuals within those networks.

```

## Containerized Components

### Purpose of different containers used (for data scraping, data preprocessing, RAG model) and instructions for running them:

This milestone covers the development and setup of a virtual environment using containers for different stages of the project, including data scraping, data preprocessing, and running the Retrieval-Augmented Generation (RAG) model. Each container is designed to handle specific tasks efficiently, ensuring modularity, scalability, and ease of deployment. Broadly, containers (1) and (2) handle the creation and updating of the vector database of research papers, and container (3) handles the RAG part given a user query.

## 1. Container for data scraping/fetching:

### **Purpose:**

- This container handles getting academic research papers from ArXiv. It retrieves AI-related research papers based on predefined categories or search queries provided.
- Currently, we only fetch research papers from ArXiv with the query “social impact” for the first baseline iteration.
- [Next steps] In the future, we will consider focusing on papers from specific conferences and tracks and limiting the domain to computer science to better filter the papers.
- [Next steps] We will scrape social issue descriptions from non-profit organizations and research databases.

### **Key Components:**

- Libraries/Frameworks: Requests for API-based data retrieval (ArXiv API). In the future, we might use Python's BeautifulSoup, Selenium, or Scrapy for scraping research group descriptions or non-profit use cases.
- Tasks performed:
  1. Query ArXiv API for papers on “social impact” and fetch metadata for the top 30 results.
  2. For each of the 30 results returned by the API, perform string manipulation to get the link of the .tex file for the corresponding folder of the manuscript and download it. Untar the folder, obtain the .tex file and convert it to a .txt file.
  3. Save all manuscript .txt files to the Google Cloud bucket within manuscript\_texts\_to\_retrieve.

### **Instructions for Running:**

- `docker build -t retrieve_papers .`
- `docker run --rm -ti -v "$(pwd)":/app retrieve_papers`

## 2. Container for data embedding, storing, and other preprocessing:

### **Purpose:**

- This container is responsible for processing the collected data and converting text into embeddings that can be used for matching research papers to social problems.

- It also handles tasks like text cleaning, tokenization, and storing embeddings in a vector database.

### **Key Components:**

- Libraries/Frameworks: all-MiniLM-L6-v2 embedding from sentence-transformers HuggingFace library, ChromaDB for vector storage, pandas, and NumPy for data manipulation.
- Tasks performed:
  1. Load the manuscripts (which are in .txt format) from the `manuscript_texts_to_retrieve` directory.
  2. Perform chunking using the character text splitting strategy and 1000 characters.
  3. Embed each chunk using the sentence transformer model `all-MiniLM-L6-v2` from the HuggingFace library. We opt for this model because it is lightweight and still ranks relatively high on well-known benchmarks.
  4. Store the embeddings in a ChromaDB vector database for later retrieval. We initiate the database if it doesn't exist and, if it exists, the code adds to the existing database. We chose to upload the vector database to the cloud for scalability as we will continue to add more papers. Apart from scalability, adding to the cloud also helps avoid redundancy (as we do not have to repeat the embedding process for papers that were processed on another team member's local computer) and makes our data storage future-proof (since GCS automates backups).
  5. For the manuscripts that have been chunked and embedded, we move the manuscripts to `manuscript_texts_done` and delete them from `manuscript_texts_to_retrieve`. This makes sure that when we run the `embed_papers` container again for a new set of retrieved papers, the embedding process is not repeated for the older papers.
  6. Currently, the only metadata we save along with the chunk embedding to the vector database is the paper source. [Next steps] In the future, we will also work on attaching other metadata to the embedding such as the title of the paper, the abstract, the publish date, and authors.

### **Instructions for Running:**

- `docker build -t embed_papers .`
- `docker run --rm -ti -v "$(pwd)":/app embed_papers`

### **3. Container for user prompt processing through a RAG model:**

#### **Purpose:**

- This container manages the retrieval of relevant research papers and generates responses for user queries using a Retrieval-Augmented Generation (RAG) model.
- It integrates the stored embeddings and research papers, processes user prompts, retrieves relevant information, and generates human-readable output for users.

### **Key Components:**

- Libraries/Frameworks: LangChain for RAG models, HuggingFace for embedding search, VertexAI for generating a response using the Gemini-1.5-Flash API, and (in the future) Flask/FastAPI for serving user prompts.
- Tasks performed:
  1. Load the existing ChromaDB from GCP onto our local instance. We believe that this is not necessarily the most efficient way to interact with a database on the cloud. However, ChromaDB currently has no support for deploying to GCP [[source](#)]. Another alternative we are considering is the [VertexAI built-in vector search](#); however, we are evaluating whether this is the best use of our GCP credits.
  2. Accept user query, embed the query, and perform vector search in the vector database with the research papers created via container (2).
  3. Feed the retrieved research paper chunks with the user query and generate a response for the user explaining how the papers relate to their query.
  4. [Next steps] Provide a dashboard for displaying relevant research to the user prompt, the research groups, and how the research work relates to the user's problem prompt. Currently, we do not have a front-end dashboard but we plan to do this in the coming milestones.

### **Instructions for Running:**

- `docker build -t perform_rag .`
- `docker run --rm -ti -v "$(pwd)":/app perform_rag`

## **Orchestration Using Docker-Compose**

Since these containers need to work sequentially, we created a `docker-compose.sh`, file to orchestrate them. The shell script is included to automate the running of all containers sequentially.

## Versioned Data Strategy (Planned, not yet done)

To ensure reproducibility and portability of data pipelines, we will implement a data versioning strategy using DVC (Data Version Control). This will allow us to track the data from different versions of research papers and social problems, ensuring consistency when embedding and storing data and fine-tuning models.

- **Strategy:** We will choose DVC for its seamless integration with Git, enabling tracking of data files alongside model code. This also provides an efficient way to manage large datasets. The versioned data pipeline will be containerized, ensuring easy integration into any environment.
- **Version Control History:** We will track dataset versions, commits, and logs to ensure the exact same datasets can be reproduced at any point.

## Utilizing LLM (Large Language Model)

We utilize LLMs in our RAG model to match AI research with societal issues. To do so, we implemented a workflow that chunked the research papers and stored the vectors in a vector database. When a user enters a prompt, we retrieve relevant research papers from the database and generate a natural language response for the user query using an LLM.

### 1. RAG Pipeline:

- **Retrieval:** Given a user query, we search our vector database of research papers and retrieve the top 5 chunks. Note that at this stage we only retrieve the relevant parts of the papers (and not the whole papers). Also, recall that we are using ChromaDB to store the chunked research papers.
- **LLM Response Generation:** The user query and retrieved relevant chunks are fed into an LLM, which is prompted to explain how the retrieved papers are related to the user's query/use case.

### 2. Fine-tuning:

We outline our plan for fine-tuning. We base our plan based on the existing efforts that address our use case, and how we can add unique value. We notice that there is existing research towards information retrieval from scientific literature (for e.g., embedding models fine-tuned on scientific abstracts). However, work towards matching these papers to relevant non-profit use cases is limited. Thus, we will fine-tune a reranker model that assigns relevance scores to papers based on their helpfulness in addressing the user-query. The steps for fine-tuning this model are below.

- Annotate a set of 500 examples of (query, paper, relevance) triplets. We will likely perform this annotation by first generating examples of user non-profit queries, then finding relevant papers and ranking them high, and finding “irrelevant” papers and ranking them low. Thus, we also make sure to have hard negatives. Because of this style of annotation, there will likely be multiple samples per query.
- Some candidate base models we consider for our reranker are bge-reranker-base and CohereRerank.

- Split the annotated dataset into 70-15-15 train, validation, and test sets. Use validation set for hyperparameter tuning and model selection, and test set for evaluation.
- Depending on how time-consuming data annotation is and how well our initial fine-tuning performs, we will consider generating synthetic data for this task based on our initial seed of annotated examples. Since the base reranker models are already highly performant, we believe that having only a few hundred training samples could be sufficient for teaching the model our task.

### 3. Experiment details/Logs:

Google Cloud Logs Explorer

VM Instance All log names All severities Correlate by

Log fields

SEARCH fields and values

RESOURCE TYPE: VM Instance

SEVERITY: Notice (37), Info (23), Error (6), Default (4)

LOG NAME: GCEGuestAgent (27), clouaudit.googleapis.com/activity (23), compute.googleapis.com/shielded\_vm... (8), diagnostic-log (4), OSConfigAgent (2), PROJECT ID: ai-research-for-good (64), INSTANCE ID: instance-20241017-173059 (36), instance-20241018-220946 (28), ZONE

Timeline: Oct 18, 5:47 PM - Oct 18, 6:48 PM

64 results

Showing logs for last 1 hour from 10/18/24, 7:54PM to 10/18/24, 8:54PM

Severity	Time	Summary
Info	Oct 18, 18:14:32.454	compute.googleapis.com [beta.compute.instances.insert] ... OSConfig Agent (version=28240926.03.g1) started.
Info	Oct 18, 18:14:33.835	compute.googleapis.com [instance-28240918-2] ... Starting the scheduler to run jobs
Info	Oct 18, 18:14:33.835	compute.googleapis.com [instance-28240918-2] ... ("logging.googleapis.com/diagnostic":{...})
Info	Oct 18, 18:14:33.837	compute.googleapis.com [instance-28240918-2] ... Scheduling job: telemetryJobID
Info	Oct 18, 18:14:33.837	compute.googleapis.com [instance-28240918-2] ... Scheduling job "telemetryJobID" to run at 24.000000 hr interval
Info	Oct 18, 18:14:33.838	compute.googleapis.com [instance-28240918-2] ... Successfully scheduled job telemetryJobID
Info	Oct 18, 18:14:33.838	compute.googleapis.com [instance-28240918-2] ... Invoking job "telemetryJobID"
Info	Oct 18, 18:14:33.840	compute.googleapis.com [instance-28240918-2] ... Scheduler added: [now 2824-10-18 22:14:33.04002644 +0000 UTC entry 1 next 2824-10-19 22:14:33 +0000 UTC]
Info	Oct 18, 18:14:39.458	compute.googleapis.com [v1.compute.instances.attachDisk] ... [all->instances/instance-28240917-173959 shanzhenbatool@harvard.edu]
Info	Oct 18, 18:14:43.732	compute.googleapis.com [v1.compute.instances.attachDisk] ... [all->instances/instance-28240917-173959 shanzhenbatool@harvard.edu]
Info	Oct 18, 18:16:49.557	compute.googleapis.com [v1.compute.instances.setMetadata] ... [all->instances/instance-28240918-220946 ramasandeepelabkar@ig.h...
Info	Oct 18, 18:16:51.285	compute.googleapis.com [instance-28240918-2] ... Created google subresource file

Google Cloud Cloud Storage Bucket details

Buckets paper-rec-bucket

Location: us (multiple regions in United States) | Storage class: Standard | Public access: Not public | Protection: Soft Delete

OBJECTS CONFIGURATION PERMISSIONS PROTECTION LIFECYCLE OBSERVABILITY INVENTORY REPORTS OPERATIONS

Folder browser

Name	Type	Created	Storage class	Last modified
arxiv_social_impact_papers.txt	text/plain	Oct 19, 2024, 12:19:12 PM	Standard	Oct 19, 2024, 12:19:12 PM
manuscript_texts_done/	Folder	—	—	—
paper_vector_db/	Folder	—	—	—

## Mock-up of the Application

We created a wireframe of the final platform, integrating the search interface, dashboard, and backend components for matching research papers to social issues.

- **Prototype:** The user interface will allow organizations to input prompts and receive relevant research papers / suggestions. The dashboard shows research groups working on similar problems and provides detailed insights.
- **Wireframe:** The wireframe below illustrates how the UI interacts with back-end components, displaying user prompts, search results, and research matches.

