# Assignment 4 - End-to-End-Reseach-Tool

**Dharun Karthick Ramaraj**

**Nikhil Godalla**

**Linata Deshmukh**

## Introduction

The **Assignment4-End-to-End-Research-Tool** is an advanced, AI-powered research assistant designed to automate and streamline the research workflow, from document ingestion to report generation. Built with a combination of **FastAPI** for backend operations and **Next.js** for the frontend interface, it leverages cloud services like **AWS S3**, **Snowflake**, and **Pinecone** for storage, processing, and querying large datasets. The tool incorporates **LangChain** and **CopilotKit** to manage multiple agents, each specialized in tasks such as document selection, querying, and external research, providing an intuitive, powerful platform for researchers.

### Technologies and Their Roles

- **FastAPI**: Chosen for its speed and asynchronous capabilities, FastAPI allows efficient backend processing and seamless integration with Pydantic for data validation.
- **Next.js**: Selected for its high-performance, interactive user interface, essential for real-time research interaction.
- **LangChain** and **CopilotKit**: These enable a modular, scalable multi-agent system with various agents specialized in querying, document retrieval, and interacting with language models.
- **AWS S3**: Used for cloud storage, facilitating secure, scalable access to research documents.
- **Snowflake**: Provides a robust data warehouse for storing and accessing publication metadata, enhancing retrieval speeds.
- **Pinecone**: Acts as a vector store, supporting similarity search and RAG (Retrieval-Augmented Generation) by enabling quick, relevant responses to user queries.
- **Google Docs API and Codelabs**: Supports the export of research notes to Google Docs in Codelabs format, offering a structured, shareable document style for research findings.

The main goal is to create a tool that allows users to efficiently parse, retrieve, and research document-based content, complete with export options for sharing results in professional formats (PDF, Codelabs).

# Problem Statement

The research assistant tool solves the challenge of managing large sets of documents, performing efficient queries, and producing structured, exportable outputs. Traditional research workflows are often fragmented, requiring separate tools for each step. This tool consolidates research tasks, from document selection and query processing to output generation, all within a single interface.

### Description of the Problem

The research process involves handling massive volumes of unstructured data from various sources, including academic publications, web articles, and documents. Researchers often face challenges such as:

1. **Inefficient Data Management**: Difficulty in organizing and storing research publications, metadata, and related materials for easy access.
2. **Time-Consuming Content Retrieval**: Extracting specific, relevant information from large documents is labor-intensive and prone to errors.
3. **Lack of Interactive Insights**: Limited tools to query content contextually and interactively, hindering deeper analysis.
4. **Fragmented Documentation Workflow**: Lack of seamless integration for note-taking, drafting, and exporting research in shareable formats.

### Desired Outcome

The solution aims to provide an **AI-powered research assistant** that simplifies and accelerates the research workflow by offering:

1. **Automated Data Ingestion**: A system to scrape, store, and manage research documents and metadata.
2. **Efficient Retrieval and Insights**: Tools for querying document content using RAG (Retrieval-Augmented Generation) to deliver precise and context-aware responses.
3. **Streamlined Workflow**: An intuitive interface for document selection, note-taking, and drafting, integrated with robust export options like PDF and Google Codelabs.
4. **Scalability and Reliability**: A backend architecture capable of handling large-scale data with minimal latency, ensuring consistent performance.

### Constraints and Limitations

1. **Data Source Constraints**:
    - Dependence on publicly accessible or API-supported sources like CFA publications and Arxiv.
    - Limited access to proprietary or paywalled content.
2. **System Requirements**:

- Integration with Snowflake, AWS S3, Pinecone, and Google APIs requires proper setup of credentials and configurations.
- Resource-intensive tasks such as document chunking and indexing demand optimized processing pipelines.

3. **Export Formatting**:
   - Ensuring consistent formatting for Codelabs and PDF exports may require additional validation and adjustments.
4. **Scalability Limitations**:
   - Handling increasing volumes of data and queries requires efficient database scaling and optimized API endpoints.

By addressing these challenges and working within the constraints, the project aims to provide a seamless, efficient, and user-friendly research assistant for academics, professionals, and students.

# Proof of Concept

To solve the challenges outlined, this project leverages various technologies, each serving a specific purpose. The technologies were chosen based on their efficiency, flexibility, and ability to handle large data volumes:

## Major Technologies

1. **FastAPI & Next.js**: Chosen for their speed and scalability, ensuring responsive interactions and real-time querying.
2. **Airflow**: Automates data pipeline workflows for data scraping, loading, and processing.
3. **AWS S3 & Snowflake**: S3 provides reliable cloud storage for documents, while Snowflake enables efficient retrieval of structured metadata.
4. **Pinecone**: Supports fast, context-aware document retrieval, essential for delivering relevant insights to user queries.
5. **CopilotKit**: Manages multi-agent processing for diverse tasks, such as document retrieval, RAG-based queries, and formatted exports.
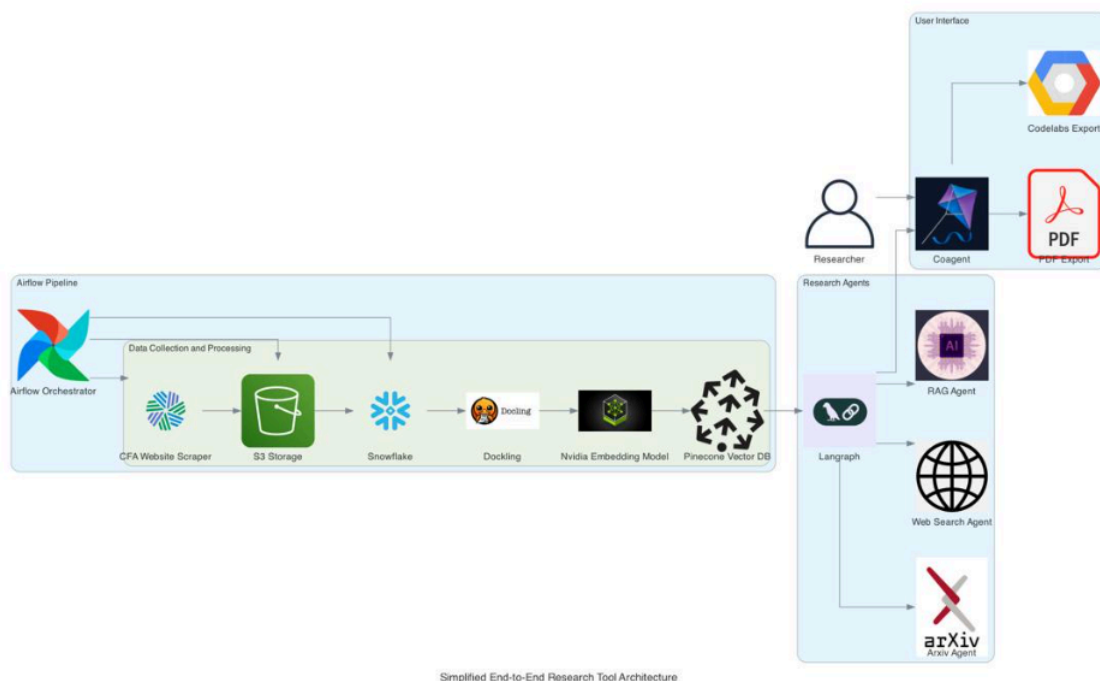
## Initial Setup and Testing

The initial setup included:

- Configuring API endpoints in FastAPI for document management.
- Setting up Airflow to automate pipeline tasks for data ingestion and loading.
- Testing Pinecone and Snowflake integration for smooth data retrieval.

## Anticipated Challenges and Solutions

1. **Multi-Agent Coordination**: Managing consistent state across multiple agents. *Solution*: Use CopilotKit for efficient state management.
2. **Handling Large Documents**: Processing large PDFs and enabling efficient retrieval. *Solution*: Use Docling to chunk documents for faster search and retrieval.
3. **Codelabs Export Formatting**: Ensuring research drafts export correctly to Codelabs format. *Solution*: GPT-4 structures and formats drafts, and Google Docs API uploads them for easy sharing.

# Architecture Diagram



Simplified End-to-End Research Tool Architecture

**1. Frontend (UI)**

- **Technology**: Built with **Next.js** and styled using **TailwindCSS**.
- **Role**: Provides an intuitive interface for researchers to interact with documents, ask questions, manage resources, and generate reports.
- **Key Features**:
    - **Query Interface**: Allows users to input research questions or specific queries.
    - **Document Selection**: Enables users to add or remove resources for analysis.
    - **Research Drafting and Export Options**: Provides options to export research drafts as PDFs or Codelabs, integrating with backend APIs for seamless document generation.
- **Interaction**:

- ○ Communicates with the backend for querying, retrieving resources, and exporting results.
- ○ Triggers document processing actions based on user input.

## 2. Backend (API Layer)

- **Technology**: **FastAPI** for creating RESTful API endpoints.
- **Role**: Serves as the central logic layer, connecting the frontend to the data pipelines, database, and AI services.
- **Key Features**:
  - ○ **Document Management**: Handles routes for querying documents, generating summaries, and managing document metadata.
  - ○ **Export Functionality**: Generates structured outputs for PDF and Codelabs formats.
  - ○ **Processing Orchestration**: Interacts with the pipeline orchestrator (Airflow) for automated data collection and processing tasks.
- **Code Reference**:
  - ○ document_selection.py: Manages document retrieval and metadata handling from Snowflake.
  - ○ rag_router.py: Manages RAG (Retrieval-Augmented Generation) queries, leveraging embeddings in Pinecone.
  - ○ export_router.py: Facilitates export requests for generating PDFs or Codelabs.
- **Interaction**:
  - ○ Communicates with Airflow for data ingestion and processing.
  - ○ Retrieves structured metadata from Snowflake and document embeddings from Pinecone.
  - ○ Connects with AI services for summarization and query processing.

## 3. Pipeline Orchestration (Airflow)

- **Technology**: **Apache Airflow**.
- **Role**: Automates and schedules data ingestion and processing tasks.
- **Key Pipelines**:
  - ○ **Scraping Pipeline**: Scrapes publications from the CFA website and uploads metadata to **AWS S3**.
  - ○ **Processing Pipeline**: Processes PDF content using **Docling**, splitting it into chunks suitable for embedding.
  - ○ **Snowflake Load Pipeline**: Loads structured data into **Snowflake** for querying.
- **Interaction**:
  - ○ Runs periodically or on-demand to keep document data up-to-date in S3 and Snowflake.
  - ○ Triggers processing tasks that store document chunks in Snowflake and embeddings in Pinecone.

## 4. Database (Snowflake)

- **Technology**: **Snowflake** Cloud Data Platform.
- **Role**: Stores metadata and document-related information, enabling fast and efficient querying.
- **Interaction**:
    - Metadata is loaded by Airflow and queried by the backend.
    - Backend fetches document metadata based on user queries.

## 5. Cloud Storage (AWS S3)

- **Technology**: **Amazon S3**.
- **Role**: Stores raw PDF files and processed data, including publication metadata.
- **Interaction**:
    - Airflow pipelines upload documents and metadata to S3.
    - Backend generates presigned URLs for secure access to stored files.

## 6. Vector Database (Pinecone)

- **Technology**: **Pinecone Vector Database**.
- **Role**: Stores embeddings of document chunks, enabling fast and scalable similarity-based retrieval during RAG queries.
- **Interaction**:
    - Document embeddings are generated by NVIDIA AI endpoints and indexed in Pinecone.
    - Backend retrieves relevant chunks based on similarity matching for user queries.

## 7. AI Services

- **Technologies**: **OpenAI** (for GPT-4), **NVIDIA AI endpoints**, **LangChain**.
- **Role**: Powers intelligent querying, summarization, and context-aware responses.
- **Interaction**:
    - Backend uses OpenAI for generating responses to user questions.
    - NVIDIA models create embeddings for document chunks, which are stored in Pinecone for RAG.
    - LangChain facilitates multi-modal RAG querying across text and document content.

## 8. Export Integration

- **Technologies**: **Google Drive and Docs APIs** for Codelabs export, **PDFKit** for PDF generation.
- **Role**: Provides export functionalities to save research drafts in user-friendly formats.
- **Interaction**:
    - Backend routes process export requests, converting research notes into structured Google Docs (Codelabs format) or downloadable PDFs.

### Data Flow

#### Data Ingestion

1. Airflow scrapes CFA publications, storing metadata and PDFs in S3.
2. Metadata is loaded into Snowflake for structured querying.

#### Processing and Indexing

1. PDFs are processed using Docling, with chunks embedded using NVIDIA AI endpoints.
2. Document embeddings are indexed in Pinecone for efficient retrieval.

#### Frontend Interaction

1. Users query documents via the Next.js frontend.
2. Backend retrieves data from Snowflake or Pinecone, processes it with AI models, and generates responses.
3. Users can export results to Google Docs (Codelabs) or PDFs.

#### Backend Communication

- API endpoints facilitate interaction between the frontend, pipelines, and cloud services, ensuring seamless integration across the system.

## Component Interaction Summary

- **Frontend ↔ Backend**: API endpoints handle document selection, querying, and export functionalities.
- **Backend ↔ Pipelines**: Backend triggers or retrieves results from Airflow pipelines for data ingestion and processing.
- **Backend ↔ Databases**: Fetches metadata from Snowflake and embeddings from Pinecone.
- **Backend ↔ AI Services**: Utilizes AI models for summarization, querying, and insight generation.
- **Backend ↔ Export Tools**: Generates outputs for PDFs and Google Docs.

# Walkthrough of the Application

## Application Overview

The **Research Helper** application provides a streamlined interface for conducting research and managing resources effectively. The main page is divided into three functional sections:

1. **Chat Interface**
   - On the right side of the page, there is a conversational chat interface where users can interact with the AI assistant. Users can:
   - Ask questions related to specific research publications. The AI will first search within the provided publications to answer the question. If the relevant information isn't found in the existing resources, the AI will broaden its search to the web.
   - View publications or resources available for research. Each publication is displayed with a title, ID, and a Download Link for easy access.
2. **Resources Management**:
   - In this area, users can add, edit, or delete resources relevant to their research. The resources that are retrieved in response to queries are automatically populated here, allowing users to manage and curate information conveniently.
3. **Research Draft**:
   - This is where users can draft and compile their research findings. Users can manually add content or import information directly from the responses generated by the AI.
   - Once the draft is complete, users have the option to **Export as PDF** or **Export as Codelabs**, making it easy to share, save, or continue working on the document in other formats.

This conversational approach makes it easy to ask follow-up questions, request summaries, and gather insights.

## Image1: Asking a Question About Publications

In this image, the user initiates a query by typing in the right-side chat interface, asking for available publications. The AI assistant responds by listing relevant publications, complete with titles, IDs, and download links. This setup allows users to quickly access and download documents for their research. The left side of the screen remains empty initially, with placeholders indicating where research question input, resources, and draft notes will be displayed as the session progresses.
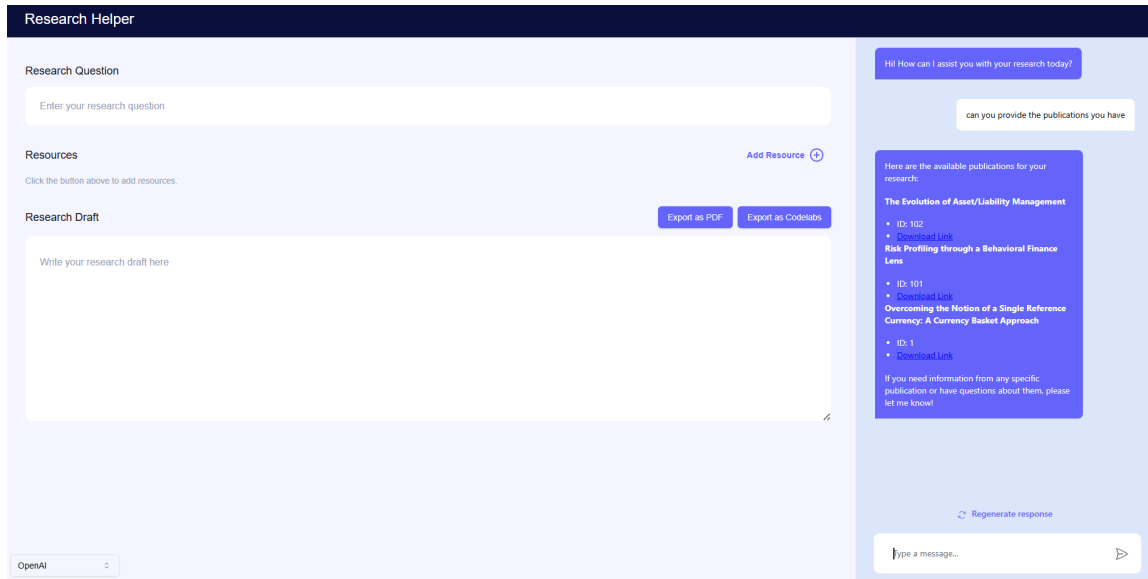
## Image 2: Adding a Research Question

In this image, the user inputs a specific research question in the **Research Question** field on the left side. This step allows the AI assistant to better understand the user's focus, guiding the interaction with relevant information and resources. The question also serves as a reference point for future interactions, helping structure the research process and responses from the assistant.
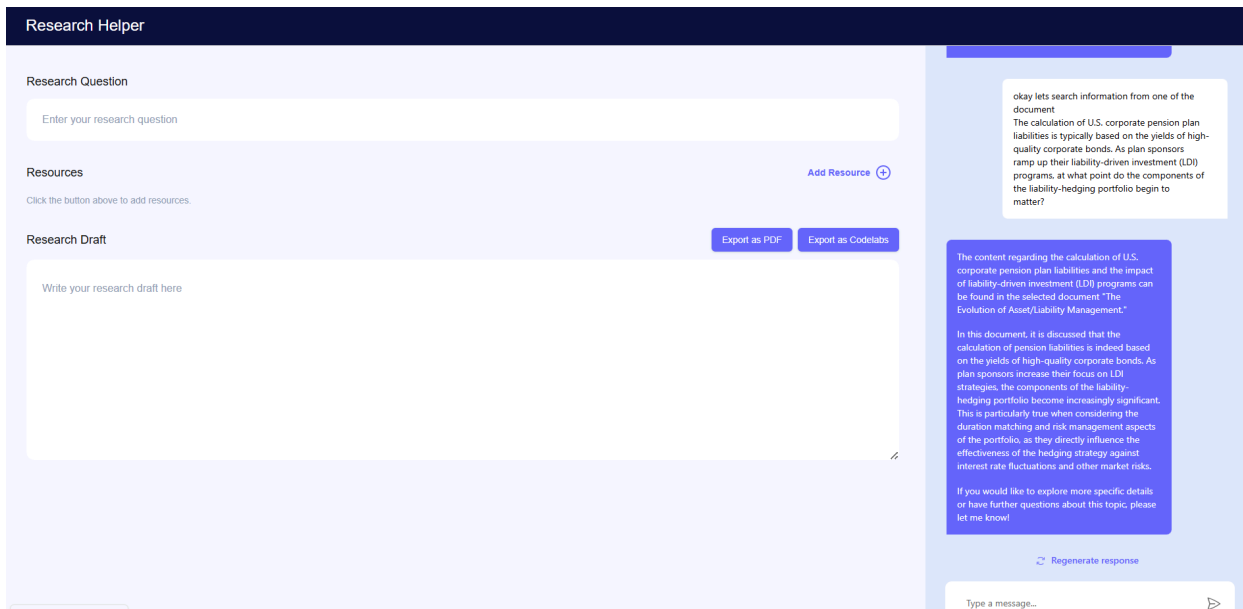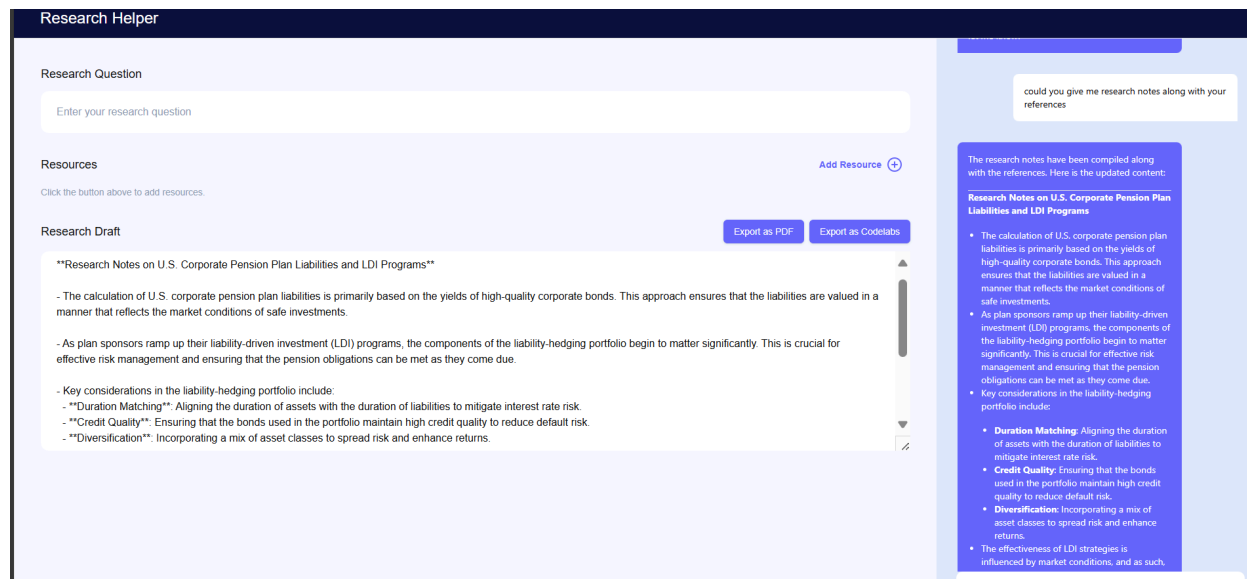


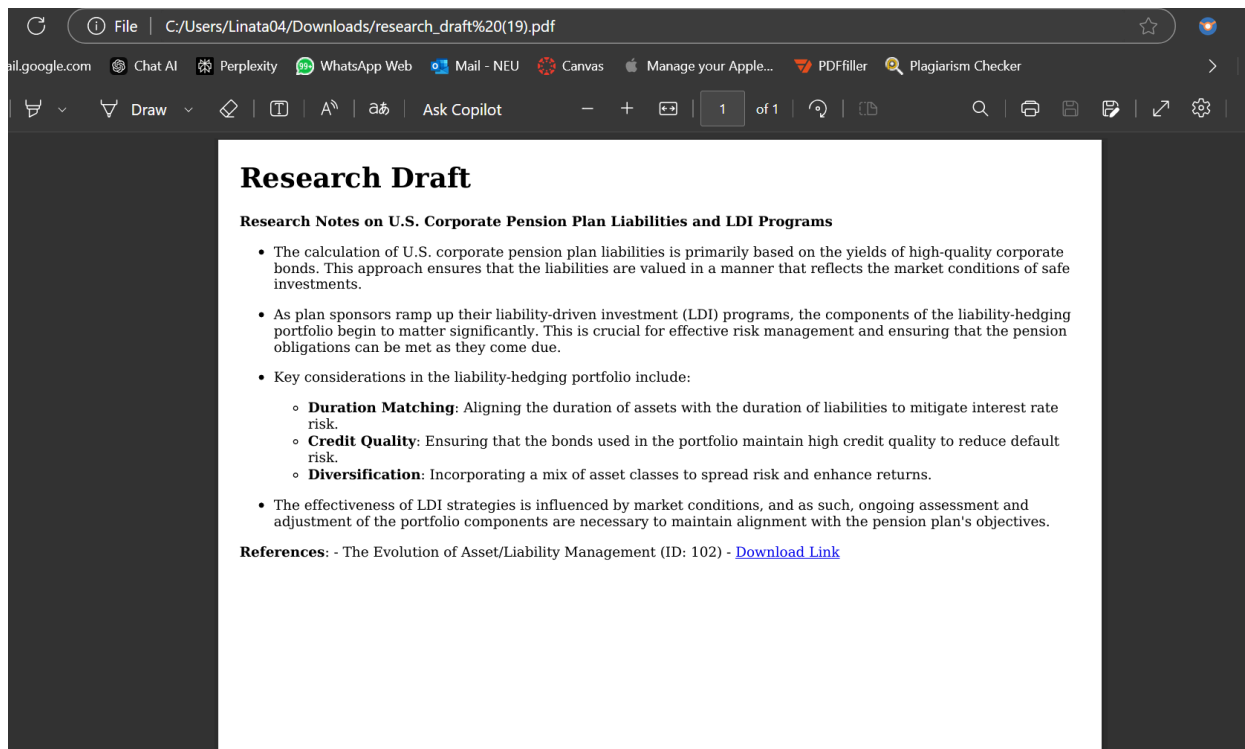## Image 3: Adding Research Notes to the Draft

In this image, the user has asked the AI assistant on the right to generate research notes along with relevant references. The assistant responds with a well-organized summary, listing key insights and recommendations related to U.S. corporate pension plan liabilities and liability-driven investment (LDI) programs. This content is then transferred to the **Research Draft** section on the left, where it is displayed in a structured format.

The **Research Draft** area now includes detailed notes with bolded headings and bullet points for easy reading. These notes provide a concise overview of the requested topic, complete with sections on **Duration Matching**, **Credit Quality**, and **Diversification**, which are crucial considerations in liability-hedging portfolios.
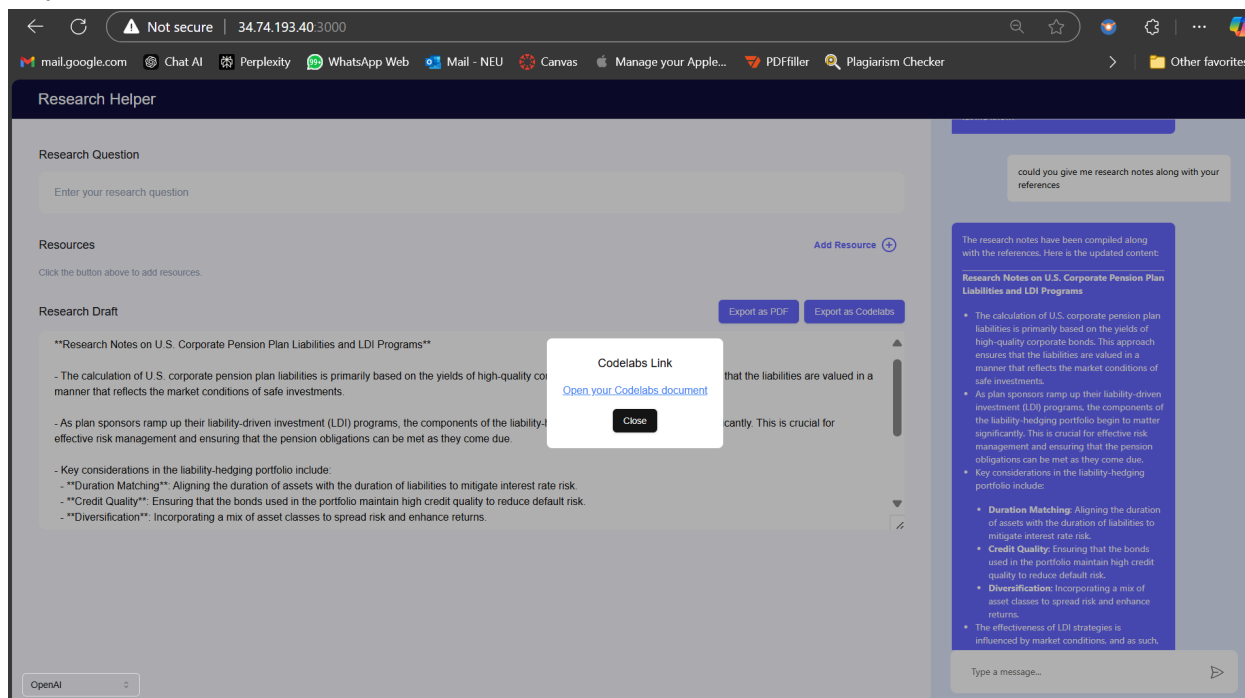
At this stage, the user has the option to review the notes, make edits if needed, and proceed to export the document either as a PDF or in Codelabs format using the buttons provided. This seamless integration of research notes into the draft allows users to compile information quickly and efficiently within the app.



After exporting to PDF

⊔ ⌄ | ▽ Draw ⌄ | ◇ | ⊤ | Aᴺ | аӡ | Ask Copilot — + ⊡ | 1 | of 1 | ↻ | ⟋ | 🔍 🖨 💾 📝 | ↗ ⚙ |

# Research Draft

**Research Notes on U.S. Corporate Pension Plan Liabilities and LDI Programs**

- The calculation of U.S. corporate pension plan liabilities is primarily based on the yields of high-quality corporate bonds. This approach ensures that the liabilities are valued in a manner that reflects the market conditions of safe investments.
- As plan sponsors ramp up their liability-driven investment (LDI) programs, the components of the liability-hedging portfolio begin to matter significantly. This is crucial for effective risk management and ensuring that the pension obligations can be met as they come due.
- Key considerations in the liability-hedging portfolio include:
  - **Duration Matching**: Aligning the duration of assets with the duration of liabilities to mitigate interest rate risk.
  - **Credit Quality**: Ensuring that the bonds used in the portfolio maintain high credit quality to reduce default risk.
  - **Diversification**: Incorporating a mix of asset classes to spread risk and enhance returns.
- The effectiveness of LDI strategies is influenced by market conditions, and as such, ongoing assessment and adjustment of the portfolio components are necessary to maintain alignment with the pension plan's objectives.

**References**: - The Evolution of Asset/Liability Management (ID: 102) - [Download Link](Download Link)

After clicking the **Export as Codelabs** button, a popup appears in the center of the screen, providing a link to view the exported document in Codelabs format. This popup contains a clickable link labeled "Open your Codelabs document," which redirects the user to the Codelabs platform, where the document is now formatted and accessible for further editing or sharing. The "Close" button on the popup allows the user to return to the main application interface once they have noted or accessed the Codelabs link.

**Research Helper**

**Research Question**

Enter your research question

could you give me research notes along with your references

**Resources**                                    Add Resource ⊕

Click the button above to add resources.

The research notes have been compiled along with the references. Here is the updated content:

**Research Draft**              Export as PDF | Export as Codelabs

**Research Notes on U.S. Corporate Pension Plan Liabilities and LDI Programs**

\*\*Research Notes on U.S. Corporate Pension Plan Liabilities and LDI Programs\*\*

- The calculation of U.S. corporate pension plan

- The calculation of U.S. corporate pension plan liabilities is primarily based on the yields of high-quality co[...] that the liabilities are valued in a manner that reflects the market conditions of safe investments.

**Codelabs Link**

[Open your Codelabs document](Open your Codelabs document)

Close

- As plan sponsors ramp up their liability-driven investment (LDI) programs, the components of the liability-[...] cantly. This is crucial for effective risk management and ensuring that the pension obligations can be met as they come due.

- Key considerations in the liability-hedging portfolio include:
  - \*\*Duration Matching\*\*: Aligning the duration of assets with the duration of liabilities to mitigate interest rate risk.
  - \*\*Credit Quality\*\*: Ensuring that the bonds used in the portfolio maintain high credit quality to reduce default risk.
  - \*\*Diversification\*\*: Incorporating a mix of asset classes to spread risk and enhance returns.

liabilities is primarily based on the yields of high-quality corporate bonds. This approach ensures that the liabilities are valued in a manner that reflects the market conditions of safe investments.
- As plan sponsors ramp up their liability-driven investment (LDI) programs, the components of the liability-hedging portfolio begin to matter significantly. This is crucial for effective risk management and ensuring that the pension obligations can be met as they come due.
- Key considerations in the liability-hedging portfolio include:
  - **Duration Matching**: Aligning the duration of assets with the duration of liabilities to mitigate interest rate risk.
  - **Credit Quality**: Ensuring that the bonds used in the portfolio maintain high credit quality to reduce default risk.
  - **Diversification**: Incorporating a mix of asset classes to spread risk and enhance returns.
- The effectiveness of LDI strategies is influenced by market conditions, and as such,
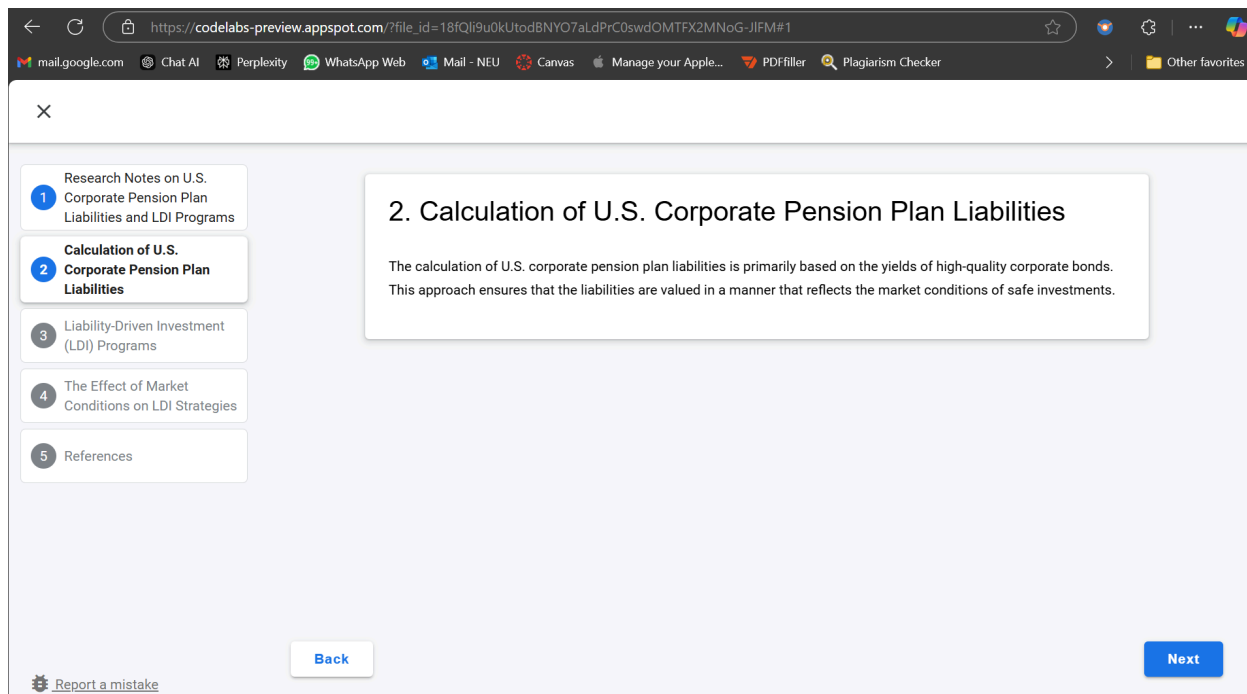
OpenAI ⌄

Type a message... ➤

## Image 4: Asking a Question Outside of Existing Publications

In this image, the user has asked a question that goes beyond the information available in the existing publications. As a result, the system initiates an external search to gather relevant resources. The "Resources" section now displays various external articles related to the query, each summarized and with clickable links for deeper exploration.

On the right side, the assistant has provided a comprehensive response based on the aggregated information from the retrieved articles. This functionality demonstrates the tool's capability to extend beyond internal resources and gather real-time information, enhancing the research process with broader insights and context. This external search feature helps users get a comprehensive overview when in-house publications don't have the required information.

## Image 5: Adding Information to Research Draft

In this screenshot, the user has asked the assistant to add the retrieved information on photosynthesis to the research draft. The assistant has successfully appended the details to the "Research Draft" section, where it provides a structured summary of how plants produce oxygen through photosynthesis, broken down into key points.

This screenshot displays the exported PDF of the research draft. The draft also provides additional resources as clickable links, allowing the reader to explore more in-depth information on related topics. This exported document serves as a formalized output, ideal for presentations, reports, or further research reference, demonstrating the tool's capability to produce publication-ready content. In addition to PDF export, this tool allows you to generate a structured Codelabs document from the research draft.



## Application Workflow (Data Engineering Work + Code Explanation)

**1. Full Data Flow: From Frontend Input to Backend Processing and Final Output**

**Frontend Input**:

- **User Actions**: The user interacts with the UI by entering a research question, adding resources, or asking specific questions related to preloaded publications.
- **State Management**: The application's state (managed by **CopilotKit** and **React hooks**) ensures that each interaction is processed and maintained per user session. This includes handling the current research question, selected resources, and logs of the ongoing processes.

**Data Flow from Frontend to Backend**:

- **API Requests**: When a user asks a question or requests a resource addition, the frontend sends an API request to the backend using the configured endpoints.

- **Backend Processing**: The backend processes the request based on the specified endpoint:
  - **Document Retrieval**: If the question is related to a publication, it checks the indexed documents first.
  - **Web Search**: If there's no match in the publications, it initiates an agent search to retrieve relevant resources online, using APIs integrated into the backend.
- **Response**: The backend returns either a document-based answer or web-based search results. This response is rendered in the chat window for the user to view, with options to add specific findings to the research notes.

**Final Output**:

- **Research Draft Generation**: Users can compile all relevant responses and notes into the research draft section. Once finalized, they can export this draft as either a **PDF** or **Codelabs** document for easy reference or sharing.

```tsx
Main.tsx    ✕

ui > src > app > Main.tsx > Main
1    import { ResearchCanvas } from "@/components/ResearchCanvas";
2    import { useModelSelectorContext } from "@/lib/model-selector-provider";
3    import { AgentState } from "@/lib/types";
4    import { useCoAgent } from "@copilotkit/react-core";
5    import { CopilotChat } from "@copilotkit/react-ui";
6
7    export default function Main() {
8      const { model, agent } = useModelSelectorContext();
9      const { state, setState } = useCoAgent<AgentState>({
10       name: agent,
11       initialState: {
12         model,
13         research_question: "",
14         resources: [],
15         report: "",
16         logs: [],
17       },
18     });
19
20     return (
21       <>
22         <h1 className="flex h-[60px] bg-[#0E103D] text-white items-center px-10 text-2xl font-med:
23           Research Helper
24         </h1>
25
26         <div
27           className="flex flex-1 border"
28           style={{ height: "calc(100vh - 60px)" }}
29         >
30           <div className="flex-1 overflow-hidden">
31             <ResearchCanvas />
32           </div>
33           <div
```

The Main component integrates ResearchCanvas for managing research workflows and CopilotChat for AI-powered assistance. It initializes the app state using useCoAgent, managing

fields like research_question and logs. The layout features a resource management area and a chat assistant, enabling seamless interaction with dynamic state updates and AI support.

```python
export_router.py 8 ×
backend > research_canvas > export_router.py > ...
42
43    @router.post("/export/pdf")
44    async def export_pdf(request: DraftRequest):
45        """Convert research draft to PDF and return it as a downloadable file."""
46        try:
47            # Convert draft Markdown to HTML
48            markdown_content = f"# Research Draft\n\n{request.draft}"
49            html_content = markdown2.markdown(markdown_content)
50
51            # Define file paths
52            html_file_path = "/tmp/research_draft.html"
53            pdf_file_path = "/tmp/research_draft.pdf"
54
55            logger.info("Writing HTML content for PDF generation.")
56            # Save HTML content as a file for PDF generation
57            with open(html_file_path, "w") as file:
58                file.write(html_content)
59
60            logger.info("Attempting to generate PDF from HTML content.")
61            # Generate PDF from the HTML content
62            pdfkit.from_file(html_file_path, pdf_file_path)
63            return FileResponse(pdf_file_path, media_type="application/pdf", filename="research_dra
64        except Exception as e:
65            logger.error(f"PDF generation failed: {e}")
66            raise HTTPException(status_code=500, detail=f"PDF generation failed: {e}")
67
68    @router.post("/export/codelabs")
69    async def export_codelabs(request: DraftRequest):
70        """Format research draft as Codelabs using GPT-4 and upload it to Google Docs."""
71        prompt = (
72            "You are a Codelabs formatting assistant. Format the given document in Codelabs style w
73            "- Use `#` only once, for the main title of the document.\n"
74            "- Use `##` for each major section that should appear in the Table of Contents.\n"
75            "- Use `###` for any finer details or sub-sections within each maior section. where appi
```

The export_router.py script provides FastAPI endpoints to export research drafts as PDFs or Codelabs documents. The /export/pdf endpoint converts Markdown drafts to PDFs using markdown2 and pdfkit, while /export/codelabs uses GPT-4 to format drafts in Codelabs style and uploads them to Google Docs via the Drive API. It returns links to the Google Doc and its Codelabs preview. The script integrates OpenAI and Google APIs for seamless formatting, document creation, and sharing, with robust logging and error handling throughout.

**2. Pipeline Orchestration with Airflow DAGs (Data Extraction, Transformation, and Loading)**

The data engineering pipeline is orchestrated using **Apache Airflow**, which manages the **Extraction**, **Transformation**, and **Loading (ETL)** process across multiple DAGs. Here's a breakdown of each key pipeline:

- **DAG 1**: scrape_cfa_publications_dag.py

- ○ **Purpose**: Scrapes publications from the CFA website, extracts key metadata (title, summary, publication date, etc.), and uploads this data to Amazon S3.
  - ○ **Key Task Scheduling**:
    - ■ scrape_task: A task that scrapes data from the CFA site.
    - ■ upload_task: Uploads extracted metadata and publications to S3 for storage.
  - ○ **Data Flow**: The scraped data is stored as a CSV in S3, which the other pipelines can access.
- ● **DAG 2**: snowflake_setup_dag.py
  - ○ **Purpose**: Sets up a **Snowflake** database, creating necessary tables and schemas for storing publication data.
  - ○ **Key Task Scheduling**:
    - ■ create_table_task: Creates tables for storing publication metadata.
    - ■ config_task: Configures the database and sets permissions.
  - ○ **Data Flow**: This DAG prepares the Snowflake environment for storing and querying publication data.
- ● **DAG 3**: snowflake_load_dag.py
  - ○ **Purpose**: Loads publication metadata from the S3 CSV into the Snowflake database, allowing efficient querying.
  - ○ **Key Task Scheduling**:
    - ■ download_from_s3_task: Downloads data from S3.
    - ■ load_to_snowflake_task: Loads data into Snowflake.
  - ○ **Data Flow**: Transforms data from the S3 bucket and imports it into Snowflake.
- ● **DAG 4**: pdf_processing_pipeline_dag.py
  - ○ **Purpose**: Processes PDFs by converting them to Markdown, chunking the content, and indexing in Pinecone for quick retrieval.
  - ○ **Key Task Scheduling**:
    - ■ pdf_to_md_task: Converts PDFs to Markdown.
    - ■ chunk_and_index_task: Splits content into chunks and indexes in Pinecone.
  - ○ **Data Flow**: Processed document data is stored in Pinecone, enabling quick retrieval based on document context.

```python
26    def init_driver():
35        driver = webdriver.Chrome(service=service, options=chrome_options)
36        return driver
37
38    # Function to download files and upload them to S3 using an in-memory buffer
39    def download_and_upload_file(url, s3_dir, s3_bucket_name, aws_region, s3):
40        if url:
41            file_name = os.path.basename(url.split('?')[0])
42            response = requests.get(url, stream=True)
43            if response.status_code == 200:
44                file_buffer = BytesIO(response.content)
45                s3_key = f"{s3_dir}/{file_name}"
46                s3.upload_fileobj(file_buffer, s3_bucket_name, s3_key)
47                s3_url = f"https://{s3_bucket_name}.s3.{aws_region}.amazonaws.com/{s3_key}"
48                return s3_url
49        return None
50
51    # Function to handle PDF extraction
52    def extract_pdf_link(pdf_soup):
53        primary_link = pdf_soup.find('a', class_='content-asset--primary', href=True)
54        if primary_link and '.pdf' in primary_link['href']:
55            return primary_link['href'] if primary_link['href'].startswith('http') else f"https://r
56
57        secondary_pdf_tag = pdf_soup.find('a', class_='items__item', href=True)
58        if secondary_pdf_tag and '.pdf' in secondary_pdf_tag['href']:
59            return secondary_pdf_tag['href'] if secondary_pdf_tag['href'].startswith('http') else f
60
61        return None
62
63    # Function to scrape publications using Selenium and save data as a pandas DataFrame
64    def scrape_publications_with_selenium(**kwargs):
65        print("Starting publication scraping process...")
66        s3 = boto3.client('s3', aws_access_key_id=aws_access_key, aws_secret_access_key=aws_secret_
67        driver = init driver()
```

Ln 2, Col 10 (399 selected)    Spaces: 4    UTF-8    LF    {} Python    3.12.3 64-bit

The scrape_cfa_publications_dag.py defines an Airflow DAG to scrape publication data from the CFA Institute website using Selenium and BeautifulSoup, process it, and upload it to AWS S3. It extracts metadata, cover images, and PDFs, saves the data as a CSV using Pandas, and uploads it to S3. The workflow involves two tasks: scraping publications and saving/uploading the CSV, executed sequentially with data passed via XCom.

```python
33    )
34    pc = PineconeClient(api_key=os.getenv('PINECONE_API_KEY'))
35
36    default_args = {
37        'owner': 'airflow',
38        'depends_on_past': False,
39        'retries': 1,
40        'retry_delay': timedelta(minutes=5),
41    }
42
43    def fetch_pdf_data_from_snowflake(**kwargs):
44        conn = snowflake.connector.connect(
45            user=os.getenv("SNOWFLAKE_USER"),
46            password=os.getenv("SNOWFLAKE_PASSWORD"),
47            account=os.getenv("SNOWFLAKE_ACCOUNT"),
48            warehouse=os.getenv("SNOWFLAKE_WAREHOUSE", "WH_PUBLICATIONS_ETL"),
49            database=os.getenv("SNOWFLAKE_DATABASE", "RESEARCH_PUBLICATIONS"),
50            schema=os.getenv("SNOWFLAKE_SCHEMA", "RESEARCH_PUBLICATIONS"),
51            role=os.getenv("SNOWFLAKE_ROLE")
52        )
53        cursor = conn.cursor()
54
55        table_name = os.getenv("SNOWFLAKE_TABLE", "PUBLICATION_LIST")
56
57        query = f"SELECT id, title, pdf_link FROM {table_name};"
58        cursor.execute(query)
59        records = cursor.fetchall()
60
61        cursor.close()
62        conn.close()
63        # Push records to XCom for use in downstream tasks
64        kwargs['ti'].xcom_push(key='pdf_data', value=records)
65
```

The pdf_processing_pipeline_dag.py defines an Airflow DAG to process, chunk, and index PDFs for efficient retrieval. It fetches metadata from Snowflake, downloads PDFs from S3, converts them to Markdown, chunks the content using Docling, and stores it back in S3. The chunks are embedded with NVIDIA's embeddings API and indexed in Pinecone for searchability. Tasks are executed sequentially: fetching data, processing PDFs, and indexing chunks, integrating Snowflake, AWS S3, NVIDIA, and Pinecone for an automated document processing workflow.

### 3. Backend Architecture and API Layer

The backend is built with **FastAPI**, providing an API layer that integrates with AWS S3, Snowflake, Pinecone, and NVIDIA's RAG services. Here's an overview of its main components:

- **API Endpoints**:
  - **Document Retrieval** (document_selection.py): Fetches publications stored in Snowflake and provides pre-signed URLs for S3-hosted PDFs.
  - **Q/A and RAG-Based Responses** (rag_router.py): Handles queries and searches for document content using NVIDIA's RAG model integrated with Pinecone.
  - **Arxiv and Web Search** (arxiv_search.py, search.py): Fetches academic articles and broader web content to expand research resources.
  - **Export Functions** (export_router.py): Exports research drafts to either PDF or Codelabs format.
- **Database Interactions**:
  - **Snowflake**: Stores publication metadata, and FastAPI endpoints query it to retrieve document information or save processed data.
  - **AWS S3**: Stores raw data and processed PDFs; the backend fetches and provides links to this data as required.
  - **Pinecone**: Hosts the indexed content from PDFs, enabling RAG model queries to retrieve relevant document content.

```python
27    async def arxiv_search_node(state: AgentState, config: RunnableConfig):
39                    "message": f"Searching Arxiv for papers related to '{query}'",
40                    "done": False
41                })
42
43        await copilotkit_emit_state(config, state)
44
45        arxiv_results = []
46
47        for i, query in enumerate(queries):
48            try:
49                print(f"Executing Arxiv search with query: {query}")
50                search = arxiv.Search(
51                    query=query,
52                    max_results=5,
53                    sort_by=arxiv.SortCriterion.Relevance
54                )
55                results = [
56                    {
57                        "title": result.title,
58                        "url": result.pdf_url,
59                        "description": result.summary
60                    } for result in search.results()
61                ]
62                arxiv_results.extend(results)
63                state["logs"][i]["done"] = True
64                await copilotkit_emit_state(config, state)
65            except Exception as e:
66                print(f"An unexpected error occurred during Arxiv search for query '{query}': {e}")
67                state["logs"][i]["message"] = f"Unexpected error: {str(e)}"
68                state["logs"][i]["done"] = True
69                await copilotkit_emit_state(config, state)
70
```
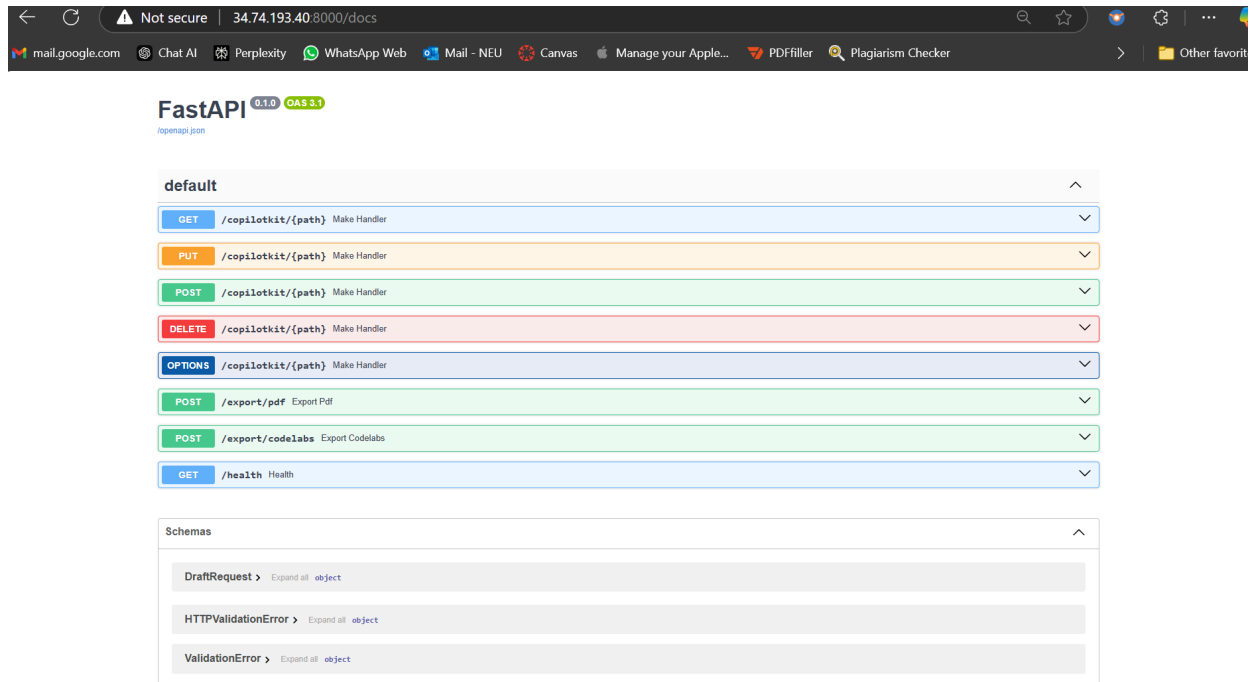
The arxiv_search.py script facilitates targeted research by integrating Arxiv searches into a research assistant. It processes user queries, retrieves up to 5 relevant papers per query using the arxiv library, and filters results through the ExtractArxivResources tool. It updates logs, resources, and messages in the application state with copilotkit_emit_state while handling errors robustly, enabling efficient Arxiv-based research workflows.

```python
16          aws_secret_access_key=os.getenv('AWS_SECRET_ACCESS_KEY'),
17          region_name=os.getenv("AWS_REGION"),
18          config=Config(signature_version='s3v4')
19      )
20
21  def extract_bucket_and_key(s3_url):
22      """Extract bucket name and key from an S3 URL."""
23      match = re.match(r'https://(.+?)\.s3\..*?\.amazonaws\.com/(.+)', s3_url)
24      if match:
25          bucket = match.group(1)
26          key = match.group(2)
27          return bucket, key
28      else:
29          raise ValueError("Invalid S3 URL format")
30
31  async def document_selection_agent(state: AgentState, config):
32      """
33      Retrieves documents from Snowflake and generates presigned URLs for them.
34      """
35
36      try:
37          # Connect to Snowflake
38          conn = snowflake.connector.connect(
39              user=os.getenv("SNOWFLAKE_USER"),
40              password=os.getenv("SNOWFLAKE_PASSWORD"),
41              account=os.getenv("SNOWFLAKE_ACCOUNT"),
42              warehouse=os.getenv("SNOWFLAKE_WAREHOUSE", "WH_PUBLICATIONS_ETL"),
43              database=os.getenv("SNOWFLAKE_DATABASE", "RESEARCH_PUBLICATIONS"),
44              schema=os.getenv("SNOWFLAKE_SCHEMA", "RESEARCH_PUBLICATIONS"),
45              role=os.getenv("SNOWFLAKE_ROLE")
46          )
```

The document_selection_agent retrieves document metadata from Snowflake, generates secure presigned URLs for S3 PDFs, and updates the application state with the document list. It uses regex to extract S3 bucket and key details, ensures URLs are valid for 1 hour, and handles errors gracefully, enabling secure and efficient document access.

```python
 9      # Initialize Pinecone and NVIDIA Embeddings clients
10      pc = PineconeClient(api_key=os.getenv('PINECONE_API_KEY'))
11      embedding_client = NVIDIAEmbeddings(
12          model="nvidia/nv-embedqa-e5-v5",
13          api_key=os.getenv("NVIDIA_API_KEY"),
14          truncate="END"
15      )
16
17      def rag_agent(state: AgentState):
18          """
19          Queries a specific Pinecone index for a document based on publication_id.
20          Updates the state with the RAG query results.
21          """
22          # Retrieve query and publication_id from state
23          query = state.get("query")
24          publication_id = state.get("publication_id")
25
26          if not query or not publication_id:
27              state["rag_query_result"] = {"error": "Missing query or publication ID."}
28              return state
29
30          index_name = f"pdf-index-{publication_id}"
31
32          # Check if the index exists
33          if index_name not in pc.list_indexes().names():
34              state["rag_query_result"] = {"error": f"Index {index_name} does not exist in Pinecone."}
35              return state
36
37          # Initialize Pinecone index
38          pinecone_index = pc.Index(index_name)
39
40          # Initialize Pinecone vector store with the embedding client and Pinecone index
41          vector_store = PineconeVectorStore(index=pinecone_index, embedding=embedding_client)
42
```

The rag.py script implements a Retrieval-Augmented Generation (RAG) agent to query Pinecone indexes for document chunks based on user queries. It retrieves the query and publication_id from the state, generates embeddings using NVIDIA's embedding client, and performs a similarity search in the relevant Pinecone index. Top matches are processed and added to the state, with error handling for missing inputs or non-existent indexes.

# References

**Next.js Documentation:** https://nextjs.org/docs

**FastAPI Documentation:** https://fastapi.tiangolo.com

**Apache Airflow Documentation:** https://airflow.apache.org/docs

**Snowflake Documentation:** https://docs.snowflake.com

**AWS S3 - Getting Started:** https://aws.amazon.com/s3/getting-started

**Pinecone Vector Database Documentation:** https://docs.pinecone.io

**OpenAI API Documentation:** https://platform.openai.com/docs

**LangChain Documentation:** https://langchain.readthedocs.io

**Google Drive API Documentation:** https://developers.google.com/drive

**PDFKit Documentation:** https://pdfkit.org