



OncoDetectAI

Project Tracker

Project Overview

Timeline: September 13 - December 12, 2024

Team: Aditi Deshmukh, Siddharth Pawar, Pranali Chipkar

Status: ✓ Completed

Total Sprints: 7

Total Tasks: 80

⌚ Sprint 1: Foundation & Environment Setup

Duration: Sept 13 - Sept 26 (2 weeks)

Status: ✓ Completed

Goal: Establish development environment and complete initial research

Sprint 1 Tasks

Aa Task	⌚ Status	⌄ Priority	⌄ Assignee	≡ Type
<u>Research SCLC literature and identify relevant papers</u>	Done	Medium	Aditi Deshmukh	Research
<u>Review DBT documentation for data transformation</u>	Done	Medium	Pranali Chipkar	Research
<u>Review Snowflake Cortex AI documentation</u>	Done	Medium	Siddharth Pawar	Research
<u>Configure Airflow environment using Docker</u>	Done	Medium	Siddharth Pawar	Setup

Aa Task	Status	Priority	Assignee	Type
<u>Create .env configuration files and Dockerfile for Airflow</u>	Done	Low	Aditi Deshmukh	Setup

Challenges Faced:

- Docker container compatibility issues with different operating systems
- Initial Airflow configuration complexity

🚀 Sprint 2: Data Pipeline Development & Deployment

Duration: Sept 27 - Oct 10 (2 weeks)

Status:  Completed

Goal: Build and deploy multimodal data ingestion pipelines

Sprint 2 Tasks

Aa Task	Status	Priority	Assignee	Type
<u>Set up S3 buckets for PDF storage and image storage</u>	Done	Medium	Pranali Chipkar	Setup
<u>Define DDL schemas for unstructured data tables in Snowflake</u>	Done	Medium	Aditi Deshmukh	Development
<u>Develop DAG to fetch SCLC papers from PubMed and upload to S3</u>	Done	High	Siddharth Pawar	Development
<u>Develop DAG for text embedding pipeline with Snowflake Arctic</u>	Done	High	Pranali Chipkar	Development
<u>Develop DAG for image embedding pipeline with Voyage Multimodal</u>	Done	High	Aditi Deshmukh	Development
<u>Deploy Airflow DAGs to GitHub repository via VS Code</u>	Done	Medium	Siddharth Pawar	Deployment

Challenges Faced:

- Papers being fetched via Airflow took longer than expected due to large file sizes
 - XCom errors in Airflow when passing data between tasks
 - Batch processing in Airflow implementation to reduce load and processing time
 - S3 bucket permission configurations
-

Sprint 3: DBT Data Modeling & Transformation

Duration: Oct 11 - Oct 24 (2 weeks)

Status:  Completed

Goal: Build structured data transformation pipeline using DBT

Sprint 3 Tasks

Aa Task	⌚ Status	⌚ Priority	⌚ Assignee	≡ Type
<u>Create DBT Staging Models for Raw Clinical & Genomic Data</u>	Done	Medium	Pranali Chipkar	Development
<u>Implement Data Quality & Validation Tests in Staging Layer</u>	Done	Medium	Siddharth Pawar	Development
<u>Develop Intermediate Models for Biomarker & Risk Feature Engineering</u>	Done	Medium	Pranali Chipkar	Development
<u>Build Mart-Level ML Feature Tables</u>	Done	Medium	Aditi Deshmukh	Development
<u>Configure Snowflake Connections via profiles.yml</u>	Done	High	Siddharth Pawar	Setup

Challenges Faced:

- Snowflake connection authentication errors via profiles.yml
 - Data quality issues with raw clinical data containing null values
 - Complex transformations for biomarker feature engineering
-

Sprint 4: Midterm Deliverables & Subtype Classification

Duration: Oct 25 - Nov 7 (2 weeks)

Status:  Completed

Goal: Complete midterm requirements and build SCLC subtype classification feature

Sprint 4 Tasks

Aa Task	Status	Priority	Assignee	Type
<u>Midterm ppt</u>	Done	Medium	Siddharth Pawar	Documentation
<u>Midterm paper research and code</u>	Done	High	Aditi Deshmukh	Documentation
<u>Midterm demo implementation logic</u>	Done	Medium	Pranali Chipkar	Development
<u>Quick classification logic for sclc subtype</u>	Done	Medium	Aditi Deshmukh	Development
<u>Full analysis logic for sclc subtype results</u>	Done	High	Pranali Chipkar	Development
<u>Setup UI page for Subtype classification</u>	Done	Medium	Siddharth Pawar	Development

Challenges Faced:

- Subtype classification logic required multiple iterations
- Code refactoring needed for demo-ready prototype

Sprint 5: Streamlit Application - User Management & Risk Prediction

Duration: Nov 8 - Nov 21 (2 weeks)

Status:  Completed

Goal: Build user authentication, risk prediction features, and Maps integration

Sprint 5 Tasks

Aa Task	Status	⌚ Priority	👤 Assignee	≡ Type
<u>Create user sign up and login page</u>	Done	Medium	Pranali Chipkar	Development
<u>Store user authentication details in SF (User management schema)</u>	Done	High	Pranali Chipkar	Development
<u>Setup feedback page</u>	Done	Medium	Siddharth Pawar	Development
<u>Store feedback data in a separate table</u>	Done	High	Siddharth Pawar	Development
<u>Create profile page for user including main buttons</u>	Done	Medium	Aditi Deshmukh	Development
<u>Create profile page for user including the product analytics button</u>	Done	Medium	Aditi Deshmukh	Development
<u>Quick prediction logic for risk</u>	Done	Medium	Pranali Chipkar	Development
<u>Full AI analysis logic for risk prediction</u>	Done	High	Pranali Chipkar	Development
<u>Integrating in-house embeddings to enhance AI analysis for risk prediction</u>	Done	High	Aditi Deshmukh	Development
<u>Setup UI for risk prediction page</u>	Done	Medium	Siddharth Pawar	Development
<u>Setup logic for risk categories</u>	Done	High	Aditi Deshmukh	Development
<u>Map agent for high risk patients</u>	Done	High	Siddharth Pawar	Development
<u>Setup network integrations for google maps api</u>	Done	High	Pranali Chipkar	Setup

Aa Task	Status	Priority	Assignee	Type
<u>Maps logic based on location</u>	Done	High	Siddharth Pawar	Development

Challenges Faced:

- Network integrations for Google Maps API setup and configuration
- Integrating and fine-tuning embedding-based context for Full AI Analysis in risk prediction
- Streamlit session state management for user authentication
- UI responsiveness issues when loading large datasets

Sprint 6: Research Assistant - Multi-Agent RAG System

Duration: Nov 22 - Dec 5 (2 weeks)

Status:  Completed

Goal: Build intelligent research assistant with multi-agent architecture

Sprint 6 Tasks

Aa Task	Status	Priority	Assignee	Type
<u>Work on RAG agent logic for research assistant</u>	Done	High	Aditi Deshmukh	Development
<u>Work on Arxiv agent logic for research assistant</u>	Done	High	Siddharth Pawar	Development
<u>Store serp api key into config table</u>	Done	High	Pranali Chipkar	Setup
<u>Work on websearch agent logic for research agent</u>	Done	High	Aditi Deshmukh	Development
<u>Work on logic for saving research notes</u>	Done	High	Siddharth Pawar	Development

Aa Task	Status	Priority	Assignee	Type
<u>Add translation, summarization features to the chat</u>	Done	Medium	Pranali Chipkar	Development
<u>Decide on threshold to call the 3 agents as per response confidence</u>	Done	High	Aditi Deshmukh	Development
<u>Create networks for arxiv and websearch access integrations</u>	Done	High	Siddharth Pawar	Setup
<u>UI Setup for research assistant page</u>	Done	High	Pranali Chipkar	Development
<u>Making UI uniform and clean across all pages</u>	Done	Medium	Aditi Deshmukh	Development

Challenges Faced:

- Network integrations for SERP API access and authentication
- Arxiv library access setup and integration
- Agent orchestration and determining optimal confidence thresholds
- Balancing response time with agent accuracy



Sprint 7: Product Analytics & Final Documentation

Duration: Dec 6 - Dec 12 (1 week)

Status: Completed

Goal: Implement Kafka-based analytics pipeline and complete all documentation

Sprint 7 Tasks

Aa Task	Status	Priority	Assignee	Type
<u>Setup Confluent Cloud</u>	Done	Medium	Aditi Deshmukh	Setup
<u>Setup topic and event structure</u>	Done	High	Aditi Deshmukh	Development

Aa Task	Status	Priority	Assignee	Type
<u>Snowflake sink connector to Confluent Cloud</u>	Done	High	Siddharth Pawar	Development
<u>Store api keys, bootstrap url and rest api url into config table</u>	Done	High	Pranali Chipkar	Setup
<u>Trigger Kafka events on user_authentication success and failures</u>	Done	High	Aditi Deshmukh	Development
<u>Trigger Kafka events on main buttons clicks</u>	Done	High	Siddharth Pawar	Development
<u>Create flattened view of the event data received from Kafka</u>	Done	High	Pranali Chipkar	Development
<u>Setup network access to Confluent Cloud</u>	Done	High	Aditi Deshmukh	Setup
<u>Create UI for analytics query chatbot</u>	Done	Medium	Siddharth Pawar	Development
<u>Work on logic for fetching text to sql function</u>	Done	High	Pranali Chipkar	Development
<u>Integrate LLM to choose relevant tables for querying data</u>	Done	High	Aditi Deshmukh	Development
<u>Including AI insights and choice of visualization on analytics page</u>	Done	High	Siddharth Pawar	Development
<u>Deploy streamlit code to github</u>	Done	High	Pranali Chipkar	Deployment
<u>Upload snowflake ddls to github</u>	Done	High	Aditi Deshmukh	Deployment
<u>Documentation: poster</u>	Done	High	Siddharth Pawar	Documentation
<u>Documentation: github readme</u>	Done	High	Pranali Chipkar	Documentation
<u>Documentation: ppt file</u>	Done	Low	Aditi Deshmukh	Documentation

Aa Task	Status	Priority	Assignee	Type
<u>Documentation: architecture diagram</u>	Done	High	Siddharth Pawar	Documentation
<u>Updating Notion for the work logs</u>	Done	Medium	Aditi Deshmukh	Documentation
<u>Documentation: YouTube video</u>	Done	Medium	Aditi Deshmukh	Documentation

Challenges Faced:

- Connection of Confluent Cloud to Snowflake sink connector
- Decisiveness in event structure for product analytics
- Network access configuration between services
- Text-to-SQL function accuracy improvements
- Time management for comprehensive documentation

Overall Project Summary

Aa Sprint	Duration	# Tasks	Status
<u>Sprint 1: Foundation & Environment Setup</u>	Sept 13 - Sept 26	5	Done
<u>Sprint 2: Data Pipeline Development & Deployment</u>	Sept 27 - Oct 10	6	Done
<u>Sprint 3: DBT Data Modeling & Transformation</u>	Oct 11 - Oct 24	5	Done
<u>Sprint 4: Midterm & Subtype Classification</u>	Oct 25 - Nov 7	6	Done
<u>Sprint 5: User Management & Risk Prediction</u>	Nov 8 - Nov 21	14	Done
<u>Sprint 6: Research Assistant Multi-Agent System</u>	Nov 22 - Dec 5	10	Done
<u>Sprint 7: Product Analytics & Documentation</u>	Dec 6 - Dec 12	20	Done

Sprint Overview

Task Distribution by Member

- **Aditi Deshmukh:** 26 tasks
- **Siddharth Pawar:** 28 tasks
- **Pranali Chipkar:** 26 tasks

Task Distribution by Type

- **Development:** 54 tasks
- **Setup:** 12 tasks
- **Documentation:** 9 tasks
- **Research:** 4 tasks
- **Deployment:** 1 task

Completion Metrics

- **Total Tasks:** 80
 - **Completed:** 80
 - **Completion Rate:** 100%
 - **Project Duration:** 13 weeks (Sept 13 - Dec 12)
-

Key Achievements

- Built multimodal RAG system with text and image embeddings
 - Implemented 7 specialized AI agents
 - Developed comprehensive risk prediction and subtype classification
 - Created modern BI analytics pipeline with Kafka
 - Delivered production-ready Streamlit application
 - Completed all documentation and demo materials
-

Technical Stack

Data Engineering: Airflow, DBT, Snowflake, AWS S3

ML/AI: Snowflake Cortex, Claude Code

Application: Streamlit, Python

Analytics: Kafka, Confluent Cloud

DevOps: Docker, GitHub, VS Code

APIs: PubMed, ArXiv, SerpAPI, Google Maps API