# HepaBot

Group members:
Hasnain Sohail 24280029

# Scope and Deliverables

🎓 **Project Title:**
 **HepaBot – An AI-powered Voice Assistant for Liver Diagnostics and Report Generation**

📝 **Brief Description:**
 Streamlines clinical documentation by transcribing doctor-patient conversations, performing speaker and role identification, and auto-generating structured medical reports.

📊 **System Overview:**

- **Input:** Audio from real-time patient consultation or uploaded recordings

- **Processing:** Whisper for transcription → Pyannote for diarization → ClinicalBERT for role classification

- **Output:** Structured, labeled medical report (D: Doctor / P: Patient) in .txt / .csv formats

✧ **What Makes It Unique:**

- Context-aware retrieval from fragmented medical speech
- Structured doctor-patient dialogue generation
- Role classification using fine-tuned ClinicalBERT
- Granular section-level search (symptoms, diagnoses, treatments)
- Auto-generated clinical summaries from fragmented records
- AI enhanced queries

🎯 **Intended Audience:**

- Clinical QA teams

- Medical researchers

- Hospital digitization staff

# Project Progress

**📦 Deliverables & Completion:**

- ✅ Transcription (Whisper) – 100%
- 🔄 Speaker Diarization – 50%
- 🔄 PDF-to-Structured Pipeline – 70%
- 👁 Vector Search & Retrieval – 40%
- 📈 Analytics Dashboard – 70%
- 📑 Auto-Report Generation – 80%

**📌 Current Status:**

- ✅ Working: Audio pipeline, basic semantic search, role classification, report generation
- ⚠ Blockers: Custom model inference on Hugging Face, UML-based entity extraction

**🛠 Tech Stack:**

- **Frameworks:** Streamlit, LangChain, ChromaDB
- **Models:** Whisper, pyannote.audio v3.x, ClinicalBERT (fine-tuned), LLaMA-3 (via Ollama)
- **Languages & Libraries:** PyTorch, HuggingFace Transformers, Pandas, Scikit-learn
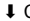
**📏 Evaluation Metrics:**

- Transcription WER (planned)
- Accuracy for role classification (TBD)
- Latency for real-time pipeline

# Live Demo Highlights

**💻 What You'll Showcase Live:**

1. Audio-to-Transcript flow using Whisper
2. Speaker Diarization results
3. Role Identification via ClinicalBERT
4. Downloadable D/P Structured Report
5. Summary Generation
6. Structured Json creation and storing in Vector DB
7. Patient Analytics
8. Semantic Search (Partially)
9. Enhanced Search capabilities with AI

**🎬 Key Outputs / UI Elements:**

- 🎙 Transcription text area
- 👥 Speaker-labeled segments
- 🧑‍⚕️ Role-tagged final dialogue (D: / P:)
- ⬇ CSV and TXT report download

**⚙ Technical Wins:**

- ClinicalBERT role detection on noisy transcripts
- Pyannote.audio v3.x migration (modern pipeline adaptation)
- Whisper performance optimization on long consultations
- Hybrid retrieval using both vector search and structured metadata
- Fine-tunning Llama2 7b on medical dataset using PEFT

# What's Next:

- Improve diarization speaker-label accuracy

- Refine prompt tuning for LLaMA-3 to enhance semantic search

- Deploy on Streamlit Cloud / Hugging Face Spaces

⚙ **Solving Blockers / Boosting Performance:**

- Enable gated model access with HuggingFace tokens

- Entity extraction via MedSpaCy or UML graph embeddings

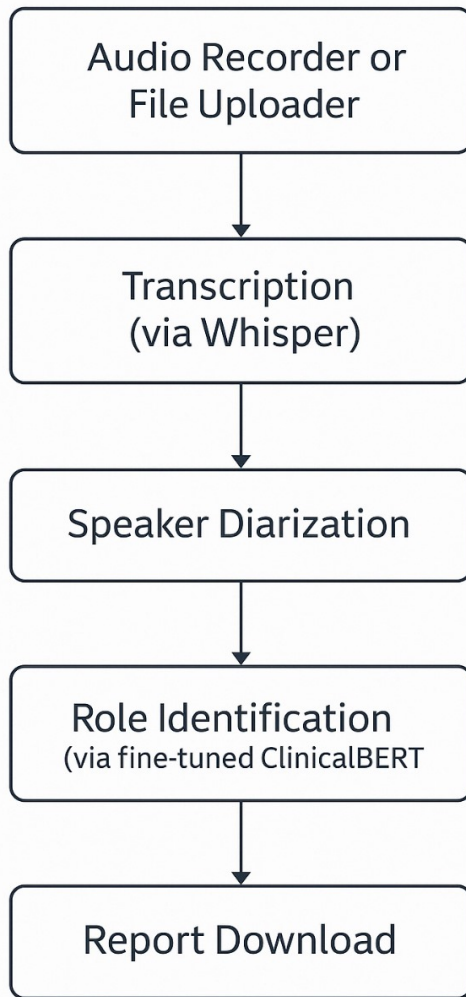- Caching for inference speed-up and latency reduction

🌀 **Hosting Plan:**

- Target: Streamlit Cloud for lightweight deployment

- Backup: Local Dockerized version for hospitals without public cloud access

☐ **Dependencies / Support Needed:**

- HuggingFace token access to pyannote/speaker-diarization

- Guidance on advanced speaker-role assignment logic (multi-speaker cases)

# SCREENSHOTS

```json
{
    {
        "source_file": "history_physical_30.pdf",
        "raw_text": "Based on the conversation data, I extracted the following information:\n\n**Patient Demographics**\n\n* Gender: Male\n* Age: 13\n* MRN: 738879\n* Dia
        "metadata": {
            "gender": "Male",
            "age": "13",
            "mrn": "738879",
            "diagnosis": "Budd-Chiari syndrome"
        }
    },
    {
        "source_file": "history_physical_24.pdf",
        "raw_text": "Based on the conversation data, I extracted the following information:\n\n**Patient Demographics**\n\n* Gender: Male\n* Age: 17\n* MRN: 776512\n* Dia
        "metadata": {
            "gender": "Male",
            "age": "17",
            "mrn": "776512",
            "diagnosis": "Budd-Chiari syndrome"
        }
    },
    {
        "source_file": "history_physical_18.pdf",
        "raw_text": "Based on the conversation, I extracted the following information:\n\nPatient Demographics:\n- Gender: Male\n- Age: 13\n- MRN: 858063\n- Diagnosis: He
        "metadata": {
            "gender": "Male",
            "age": "13",
            "mrn": "858063",
            "diagnosis": "Hepatic fibrosis"
```

Terminal

```
streamlit run app.py

  You can now view your Streamlit app in your browser.

  Local URL: http://localhost:8501
  Network URL: http://192.168.50.129:8501

  For better performance, install the Watchdog module:

  $ xcode-select --install
  $ pip install watchdog


/Users/aks/Documents/ollama-models-llm/ollama-fundamentals/ollama-fundamentals/after-mid/app.py:45: LangChainDeprecationWarning: The class `Chroma` was deprecated in LangChain 0.2.9 and will be removed in
1.0. An updated version of the class exists in the :class:`~langchain-chroma package and should be used instead. To use it run `pip install -U :class:`~langchain-chroma` and import as `from
:class:`~langchain_chroma import Chroma``.
  vectordb = Chroma(
```

# 🏥 Medical Records Semantic Search

## Semantic Search

Search across all patient records using natural language queries

Enter your search query:

E.g., Patients with kidney stones and HCV

☑ Enhance query with AI

Search

Advanced Filters

---

### Navigation

Semantic Search

Patient Browser

Analytics

Generate Report

This application provides semantic search capabilities for medical records. It uses vector embeddings to find relevant patient records based on your query.

Deploy

## Navigation

**Semantic Search**

**Patient Browser**

**Analytics**

**Generate Report**

This application provides semantic search capabilities for medical records. It uses vector embeddings to find relevant patient records based on your query.

management of of bleeding and monitoring of liver function due to hepatic fibrosis. Further follow-up appointments are scheduled to assess treatment response and address any emerging symptoms.

## Search Results (10 found)

Patient: 753394 - Budd-Chiari syndrome (Diagnostic Conclusions)

Patient: 708087 - Hepatic fibrosis (Summary Narrative)

Patient: 708087 - Hepatic fibrosis (Diagnostic Conclusions)

Patient: 858063 - Hepatic fibrosis (Summary Narrative)

Patient: 708087 - Hepatic fibrosis (FULL_TEXT)

Patient: 708087 - Hepatic fibrosis (Summary Narrative)

Patient: 753394 - Budd-Chiari syndrome (FULL_TEXT)

Patient: 708087 - Hepatic fibrosis (FULL_TEXT)

Patient: 858063 - Hepatic fibrosis (Negative Findings)

Patient: 753394 - Budd-Chiari syndrome (Diagnostic Conclusions)

localhost

Aizelsheikh/llama2-finetuned - Hugging Face | slides report - Google Docs | (2) WhatsApp | Medical Records Semantic Search

RUNNING...  Stop  Deploy

# 🏥 Medical Records Semantic Search

## Navigation

Semantic Search

Patient Browser

Analytics

Generate Report

This application provides semantic search capabilities for medical records. It uses vector embeddings to find relevant patient records based on your query.

## Semantic Search

Search across all patient records using natural language queries

Enter your search query:

abdominal pain

☑ Enhance query with AI

Search

Advanced Filters

Enhanced query: Abdominal pain (dyspepsia) with nausea and vomiting, accompanied by fever and chills, in a patient with a history of gastrointestinal issues (e.g., IBS, Crohn's disease). Include symptoms such as weight loss, diarrhea, or constipation. Also, consider related terms like abdominal tenderness, guarding, or rebound tenderness. Use MeSH terms: 741.5 (abdominal pain), 382.4 (gastrointestinal disorders), and 111.8 (inflammatory bowel disease).

⟳ Generating clinical summary...

Chroma` was deprecated in ed instead. To use it run

Python 3.10 (ollama-models-llm)

| | | | | | |
|---|---|---|---|---|---|
| 8 | Male | 55 | 798650 | Chronic viral hepatitis C | history_physical_8.pdf |
| 9 | Male | 56 | 796637 | Chronic viral hepatitis C | history_physical_9.pdf |

## Patient Detail View

Select Patient MRN

708087 ⌄

| MRN | Age | Gender |
|---|---|---|
| **708087** | **5** | **Male** |

## Diagnosis: Hepatic fibrosis

Full Clinical Summary

Patient Demographics:
- Gender: Male
- Age: 5
- MRN: 708087
- Diagnosis: Hepatic fibrosis

Clinical Summary:

Active Symptoms:
- Abdominal pain (duration: 2 days): described as sharp and intermittent
- Fatigue (duration: 3 days): reported by the patient's mother

Negative Findings:
- No fever or chills
- No recent travel history
- No new symptoms in the past week

Diagnostic Conclusions:

# 🏥 Medical Records Semantic Search

## Patient Record Browser

| | gender | age | mrn | diagnosis | source_file |
|---|---|---|---|---|---|
| 0 | Male | 13 | 738879 | Budd-Chiari syndrome | history_physical_30.pdf |
| 1 | Male | 17 | 776512 | Budd-Chiari syndrome | history_physical_24.pdf |
| 2 | Male | 13 | 858063 | Hepatic fibrosis | history_physical_18.pdf |
| 3 | Female | 12 | 661446 | Hepatic fibrosis | history_physical_19.pdf |
| 4 | Female | 12 | 691761 | Budd-Chiari syndrome | history_physical_25.pdf |
| 5 | Female | 5 | 709356 | Budd-Chiari syndrome | history_physical_31.pdf |
| 6 | Male | 7 | 753394 | Budd-Chiari syndrome | history_physical_27.pdf |
| 7 | Male | 12 | 698512 | Budd-Chiari syndrome | history_physical_33.pdf |
| 8 | Male | 55 | 798650 | Chronic viral hepatitis C | history_physical_8.pdf |
| 9 | Male | 56 | 796637 | Chronic viral hepatitis C | history_physical_9.pdf |

## Patient Detail View

Select Patient MRN

602927

MRN
## 602927

Age
## 51

Gender
## Female

## Diagnosis: Hepatic fibrosis

Deploy

# 🏥 Medical Records Semantic Search

## Medical Records Analytics

### Patient Demographics



### Diagnosis Distribution



### Age Distribution

# Model Card for Model ID

## Model Details

### Model Description

- **Developed by:** [Rabia Aslam]
- **Funded by [optional]:** [More Information Needed]
- **Shared by [optional]:** [More Information Needed]
- **Model type:** [More Information Needed]
- **Language(s) (NLP):** [More Information Needed]
- **License:** [More Information Needed]
- **Finetuned from model [optional]:** [More Information Needed]

## Model Sources [optional]

- **Repository:** [More Information Needed]
- **Paper [optional]:** [More Information Needed]
- **Demo [optional]:** [More Information Needed]

Uses