# ChemNovus HR Assistant chatbot

May / 2025

Presented by

Eugene Hao

# Content

# What is built

- Gradio web interface for HR policy Q&A

- Powered by OpenAI GPT-4o with RAG architecture

- Ingest DOCX, PDF, PPTX, XLSX files as the ground truth

- Region-aware responses (Global/US/EU)

- Conversational history support

# Tech stack

**1** ⊙

## Gradio App

open-source Python package to build web app for ML model, API

**2** ⊙

## LLM

OpenAI GPT-4o

**3** ⊙

## Embedding

text-embedding-ada-002

**4** ⊙

## Search

Cosine similarity for KNN

# RAG workflow

**1**

**2**

**3**

**4**

## Ingestion

- Process documents by region (Global / US / EU)

- Chunk text (2000 chars), stored in memory

## Embedding

- Generate ADA-002 embeddings

- Store the embedding in memory

## Query

- Embed user question

- Find top-20 similar chunks

## Generation

- Inject chunks + history into GPT-4o prompt

- Return formatted response

# Prompt Engineering

**GPT-4o**

```
"role": "system",
"content":
"You are an HR assistant for ChemNovus
Incorporated."
● Answer questions using ONLY the
  provided documents.
● Pay special attention to
  region-specific information (Global,
  US, Europe).
● When information differs between
  regions, present it in a clear
  comparison format.
● For global policies, indicate they
  apply worldwide.
● Format responses with clear section
  headers and bullet points.
```
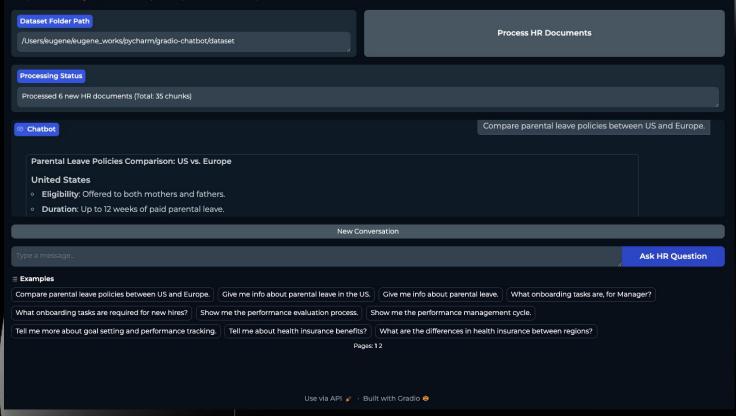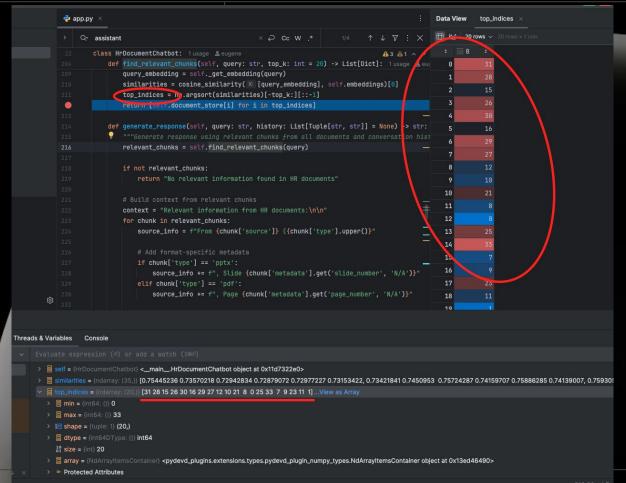
# Demo

📁 **ChemNovus Incorporated HR Assistant**

Ask questions about global, US, and European HR policies, benefits, and procedures

**Dataset Folder Path**

/Users/eugene/eugene_works/pycharm/gradio-chatbot/dataset

**Process HR Documents**

**Processing Status**

Processed 6 new HR documents (Total: 35 chunks)

💬 **Chatbot**

Compare parental leave policies between US and Europe.

**Parental Leave Policies Comparison: US vs. Europe**

**United States**
- **Eligibility**: Offered to both mothers and fathers.
- **Duration**: Up to 12 weeks of paid parental leave.

**New Conversation**

Type a message...                                                              **Ask HR Question**

☰ **Examples**

Compare parental leave policies between US and Europe.    Give me info about parental leave in the US.    Give me info about parental leave.    What onboarding tasks are, for Manager?

What onboarding tasks are required for new hires?    Show me the performance evaluation process.    Show me the performance management cycle.

Tell me more about goal setting and performance tracking.    Tell me about health insurance benefits?    What are the differences in health insurance between regions?

Pages: 1 2

Use via API 🚀   ·   Built with Gradio 🟠

● Company Name                                    ■ Quarter Month                                                        Year

# Demo

# Productionize Roadmap

## Infra

- Dockerize app + Kubernetes deployment
- Persistent vector database (Pinecone/Weaviate)
- Redis for conversation history

## Performance

- Cache common queries
- Async document processing
- Rate limiting
- PCA for embedding (shorter vector)

## Enhancements

- PDF text extraction improvement (OCR for scans)
- User history and recommendation
- Feedback mechanism (thumbs up/down)
- Doc gen: work verification
- ElasticSearch to replace in-mem storage
  - disk
  - ANN (faster than sklearn linear knn scan)
  - filter by region
- CI / CD:
  - Regression test
  - versioned release
  - Deploy (k8s rollout)

## Monitoring

- Log LLM token usage/costs
- Track unanswered questions
- Alert for document changes
- 4xx / 5xx error
- CPU / Mem usage

# Key Challenges

## Chunk size

- chunk size 2000 chars works ok
- ada-002 token limit: 8191 tokens
  - avg 4 chars / token
  - 2000 / 4 = 500 token per chunk

## Top k

- top_k = 20
- when used top_k = 3, missed some answers for question *"What onboarding tasks are required for new hires?"*

## Temperature

- temperature=0.3
- Lower values (e.g., 0.3) make the output more focused and deterministic.
- Higher values (e.g., 1.0 or above) make the output more random and creative.

## Region

- should categorize the region for the ingested unstructured documents

# Business impact

- Huge reduction in HR policy human communication

- Instant access to latest documents

- Quick response then email HR (24 x 7)

Q & A