# One Stop Health

Arda Celik, Wesley Burnawan
Ontario Tech University
19 November 2024
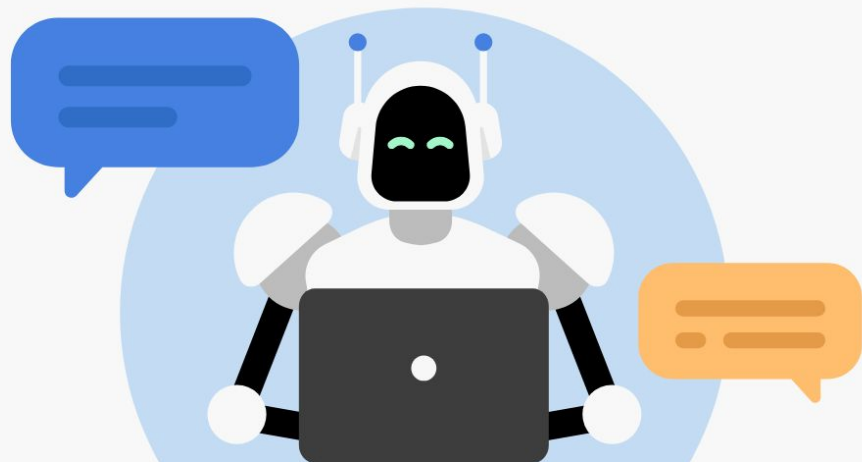
# Motivation

- 22% of Canadian adults don't have a family doctor
- 66% described Canadian healthcare system as "long wait"
- Only 26% of Canadians could get same day or next day appointment
- Google is not always reliable
- Research question: How can LLMs be optimized to provide accurate and context-aware health-related responses?

# Background

1. Chatbot with prefixed query and response
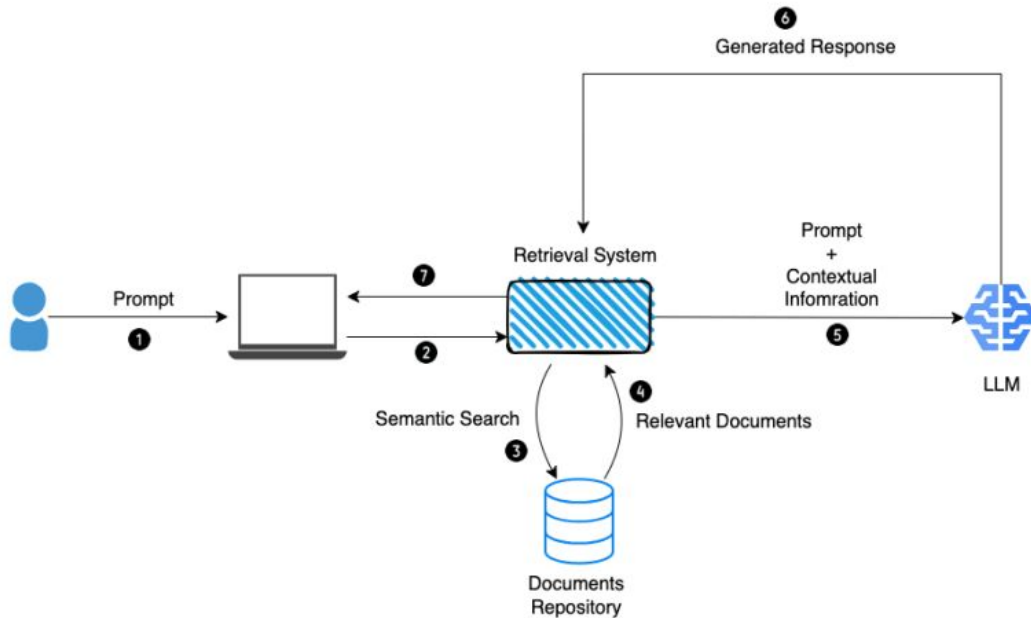2. Chatbot with NLP
3. Chatbot with LLM and sensors

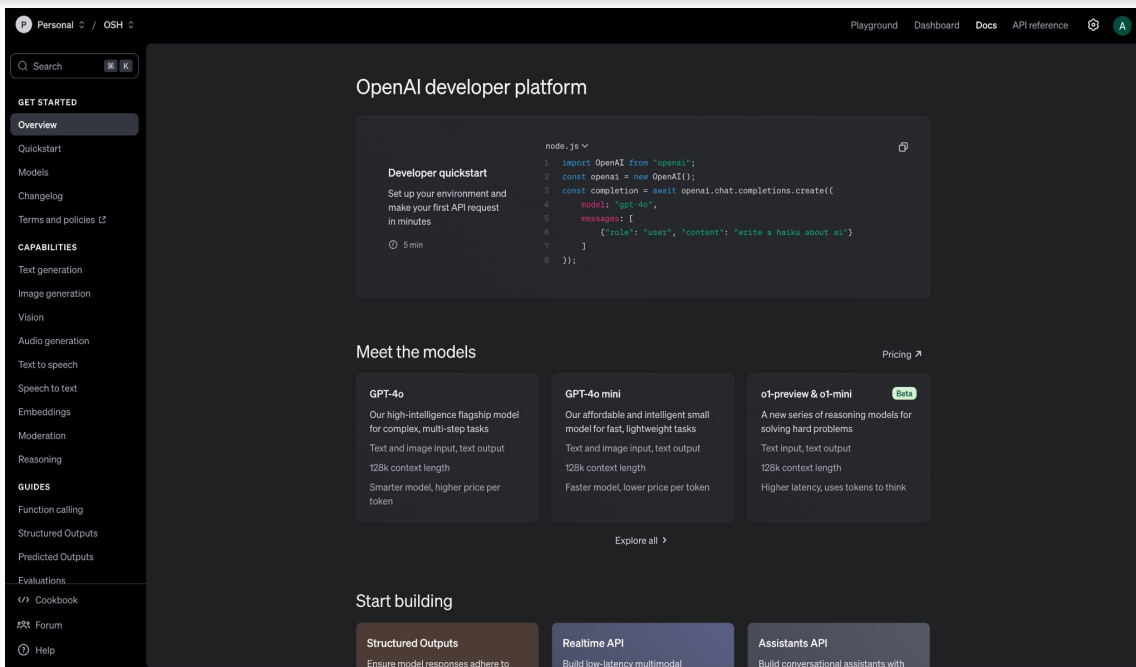Our version: RAG Chatbot

# Design and Implementation

App components:

- UI
- API
- LLM
- MongoDB
- Pinecone (vector database)
- CI/CD
- Google Cloud

# LLM Integration ⊛ OpenAI

- OpenAI's API platform
  - Models: GPT-4o, GPT-4o-mini, text-embedding-3-small
  - Capable models
  - Easy integration

# LLM Alternatives

**BioMistral: A Collection of Open-Source Pretrained Large Language Models for Medical Domains**

### BioGPT

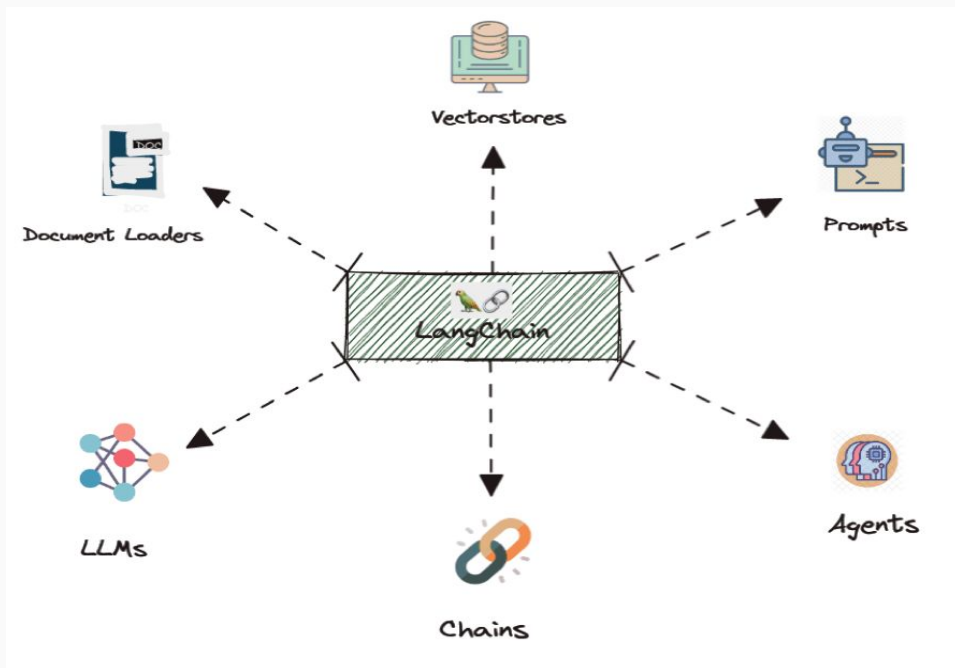Pre-trained language models have attracted increasing attention in the biomedical domain, inspired by their great success in the general natural language domain. Among the two main branches of pre-trained language models in the general language domain, i.e. BERT (and its variants) and GPT (and its variants), the first one has been extensively studied in the biomedical domain, such as BioBERT and PubMedBERT. While they have achieved great success on a variety of discriminative downstream biomedical tasks, the lack of generation ability constrains their application scope. In this paper, we propose BioGPT, a domain-specific generative Transformer language model pre-trained on large-scale biomedical literature. We evaluate BioGPT on six biomedical natural language processing tasks and demonstrate that our model outperforms previous models on most tasks. Especially, we get 44.98%, 38.42% and 40.76% F1 score on BC5CDR, KD-DTI and DDI end-to-end relation extraction tasks, respectively, and 78.2% accuracy on PubMedQA, creating a new record. Our case study on text generation further demonstrates the advantage of BioGPT on biomedical literature to generate fluent descriptions for biomedical terms.

You can use this model directly with a pipeline for text generation. Since the generation relies on some randomness, we set a seed for reproducibility:
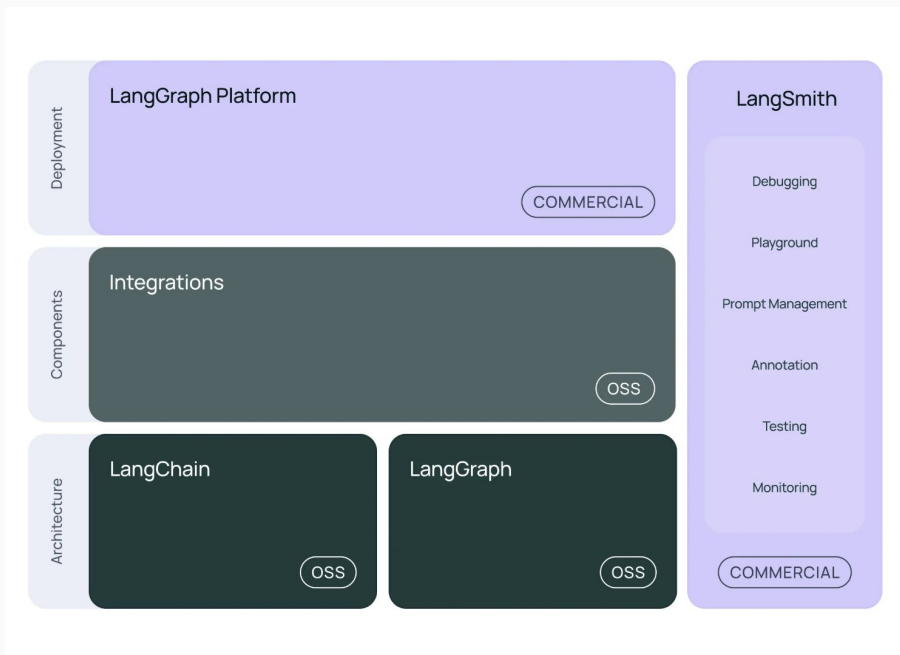
# Langchain 🦜🔗

- Langchain (JS)
  - Allows easy integration with LLMs and other tools
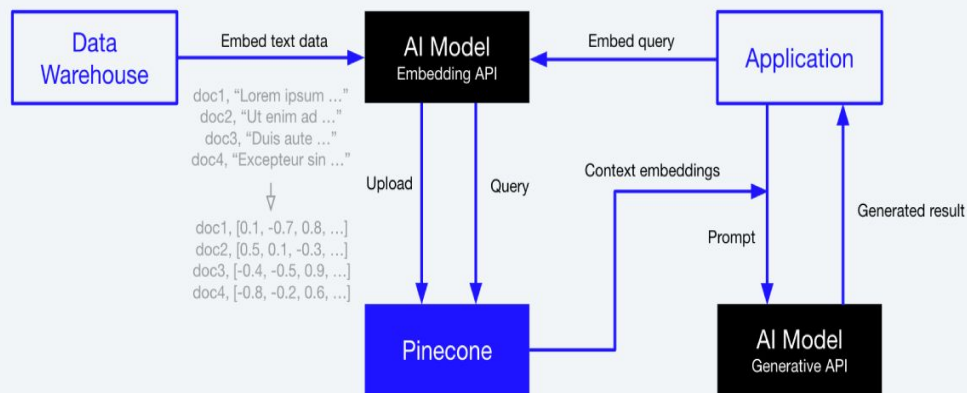  - LLM invocation, retrieval, chat memory

# Langchain 🦜🔗

# Pinecone

- Serves as additional data source for RAG applications
- Data is stored as vectors in indexes
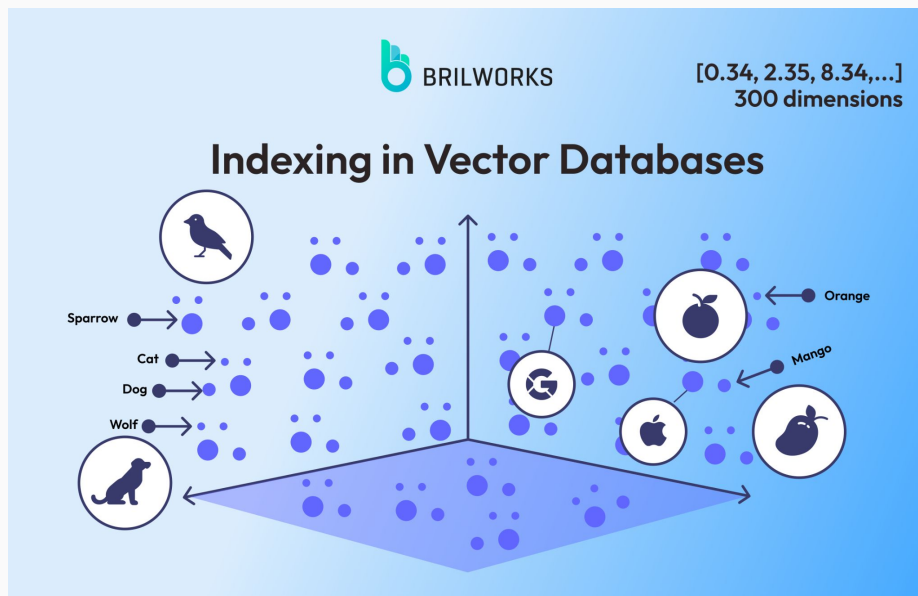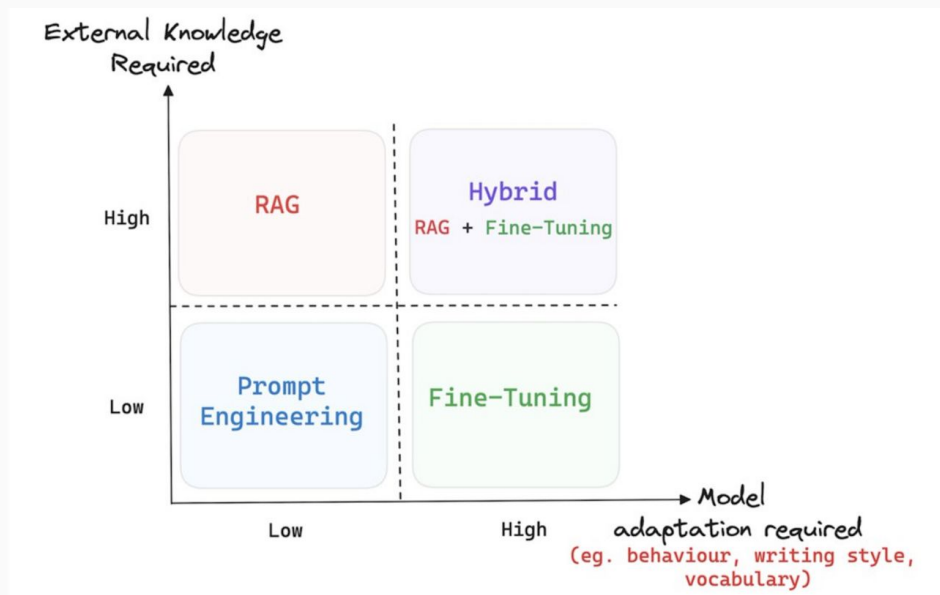- Data needs to be turned into embeddings
- Embedding size is important!

# Pinecone

# Pinecone

# RAG vs Fine Tuning

- RAG
  - Real time data
  - Ability to store data from different sources (research paper)
  - Factual context
- Fine tuning:
  - Availability of a dataset in a specific format
  - Changes in the "tone"

# "LLMs are frozen in time"

# Dataset

- MedQuad-MedicalQnADataset
    - Q&A pairs
    - 16k + rows

# Demo

# Results

- Verified that the RAG pipeline works ✔️
- LLM has memory ✔️
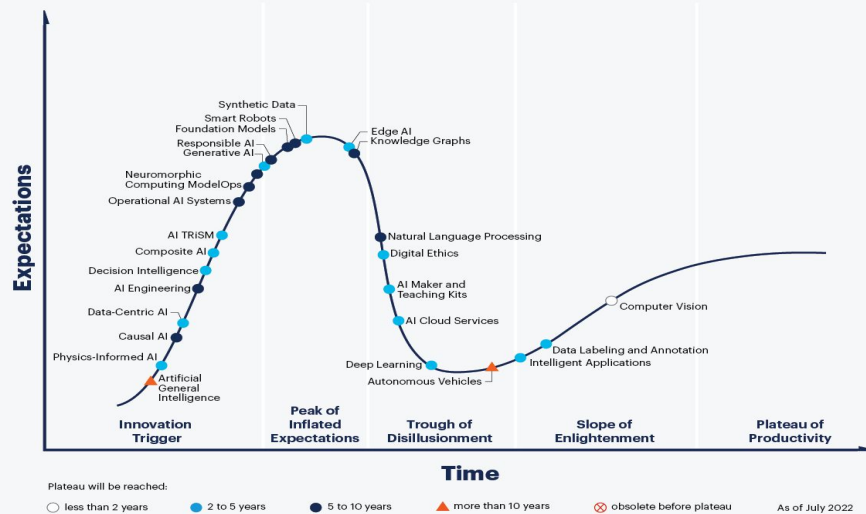- Comprehensive evaluation ❌

# Challenges

- Finding the right dataset (Scale and type)
- Designing user interactions (multi-turn conversations)
- Conditional context retrieval (not every user query requires context)
- Evaluation methods (What is a successful interaction? What do we consider accurate?)
- Additional features like appointment scheduling, speech recognition etc.
- Familiarity with the tech stack
- Unit testing

# Conclusion

Exciting time to be a developer :)



**Hype Cycle for Artificial Intelligence, 2022**

# Thank you!