APRIL 2024

# HIGH RISK PROJECT

Predicting Patient Encounters Using Local LLMs to Analyze Previous Encounters

**MATT DYL**
MSAI Graduate Student, The University of Texas at Austin

# Project Goals

# What were we trying to accomplish?

Phase 1: Predict future patient encounters using previous medical notes.

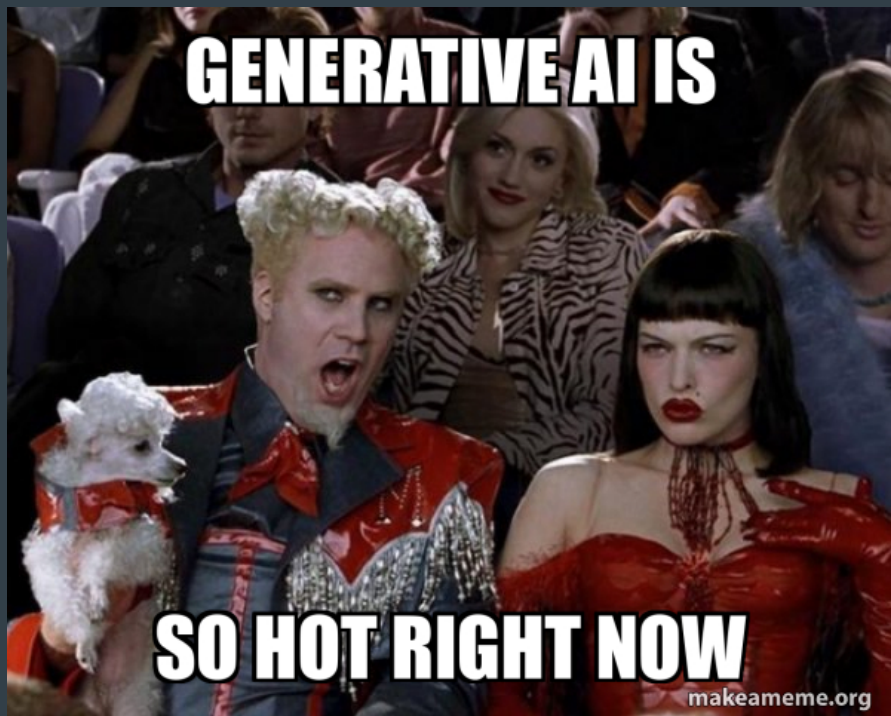Phase 2: Predict what the future encounter would be.

**Why?**

Should I go to the doctor?

Being able to accurately forecast future patient encounters or needs would lead to better utilization of strained medical resources as well as allow for better patient outcomes as care would become more preventative rather than acute.
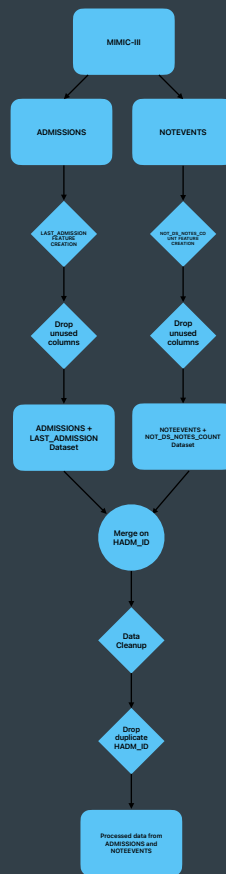
# How were we going to accomplish this?

# The Data

Dataset:
MIMIC-III

Tables:
Admissions
NoteEvents

**The Models**

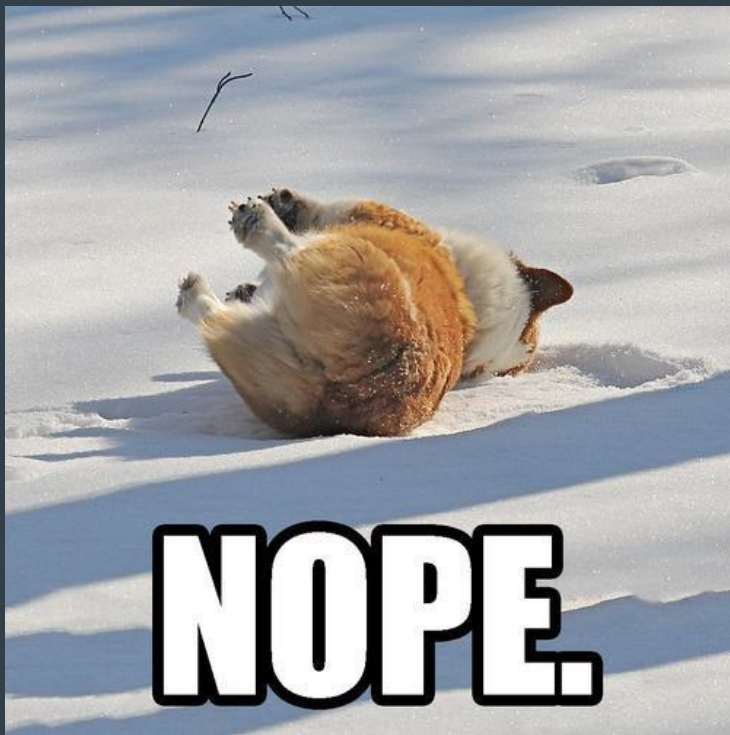nmitchko/medguanaco-65b-GPTQ

microsoft/BioGPT-Large-PubMedQA

- Due to the restrictions and privacy concerns that come with using MIMIC-III we were forced to use local LLM models to try to accomplish our goals
- Both of these models were fine-tuned for medical domain tasks and should give much better results and context vs general purpose AI models

## Combine the two

- By feeding the input text we preprocessed into the model we would format prompts to input into the LLM asking a simple question of if this patient was likely to have a future admission
- We would then be able to train the model by comparing its output to our generated feature of "LAST_ADMISSION" and see if it correctly predicted whether or not the entry was the last admission for the patient
- Upon satisfactory performance of this task our scope would expand to try to predict and classify what the nature of the future admission would be by creating a new feature called "SECOND_TO_LAST_ADMISSION" and retraining the model to see if it could accurately predict the nature of the "LAST_ADMISSION" record

# Were we successful?

# Postmortem

Issues:

LLM models did not want to run on local machine

Feedback loop was excessively long

What we learned:

Local environment matters a lot

For data processing we should offload the processed data into an external file to allow simply reloading the data vs reprocessing it each time

# What would we do different?

**Improvements**

- Be willing to pivot more quickly when an approach is not working

- Use synthetic data instead of MIMIC-III in order to be able to utilize outside computing power