# References (delete or hide for presentation)

- Rubric with tasks - link
- Team project plan - link
- Roadmap - link

# MIDSearch

A RAG based multiplatform search tool for the MIDS program

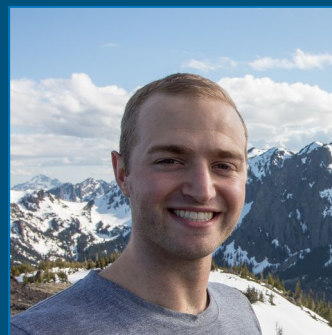Devin Suy, Nadia Tantsyura, Randy Louie, Robert Greer, Thomas Lai

# Team

Devin Suy

Nadia Tantsyura

Thomas Lai

Randy Louie

Robert Greer

# Problem & Motivation

# Market research

## Knowledge bases / Chat apps
**(non-cohesive)**



slack

OneDrive

ATLASSIAN

GitHub

G Suite

Discord

Local storage

## Cohesive-ish apps

Microsoft Copilot

**Recall function**
(releasing 6/18/24!)

glue.

# Market size

~5,000
total current students +
alumni + professors

+200
per
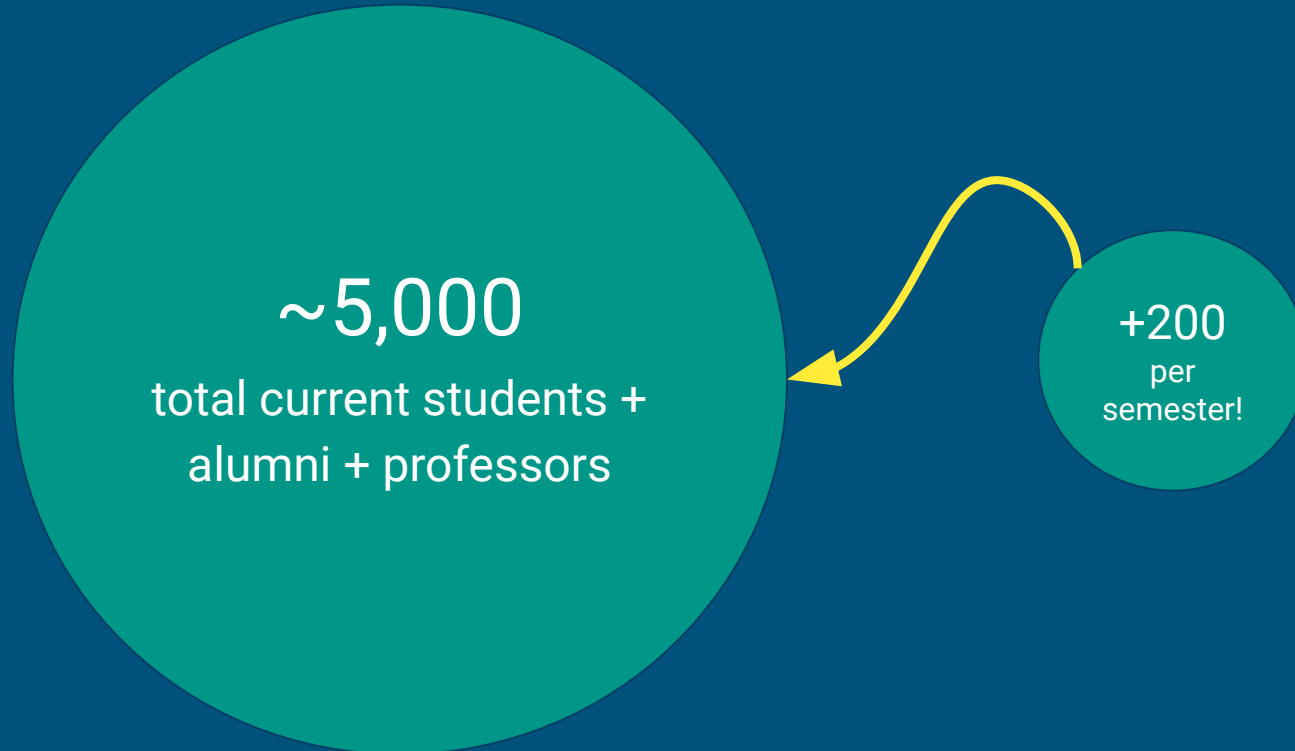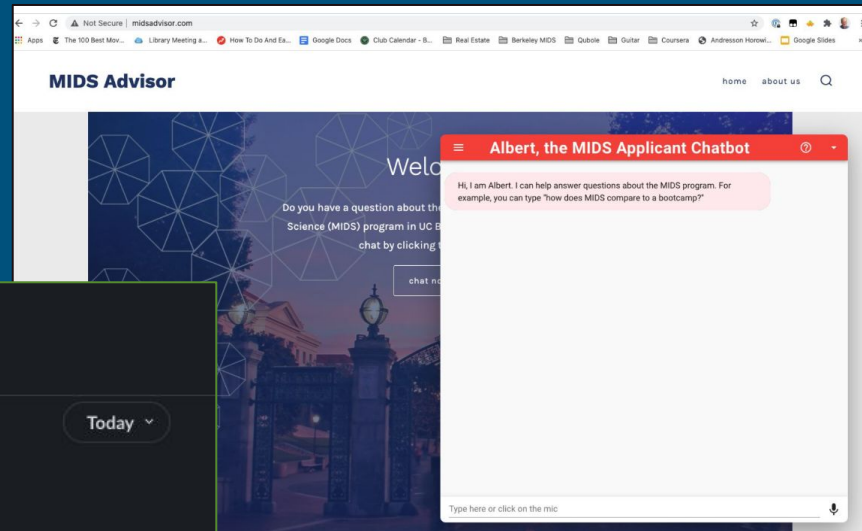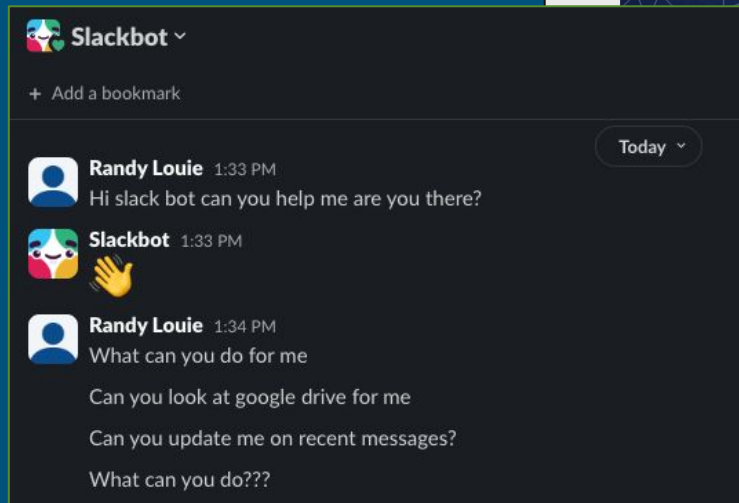semester!

# What has been tried before

- Previous capstone:
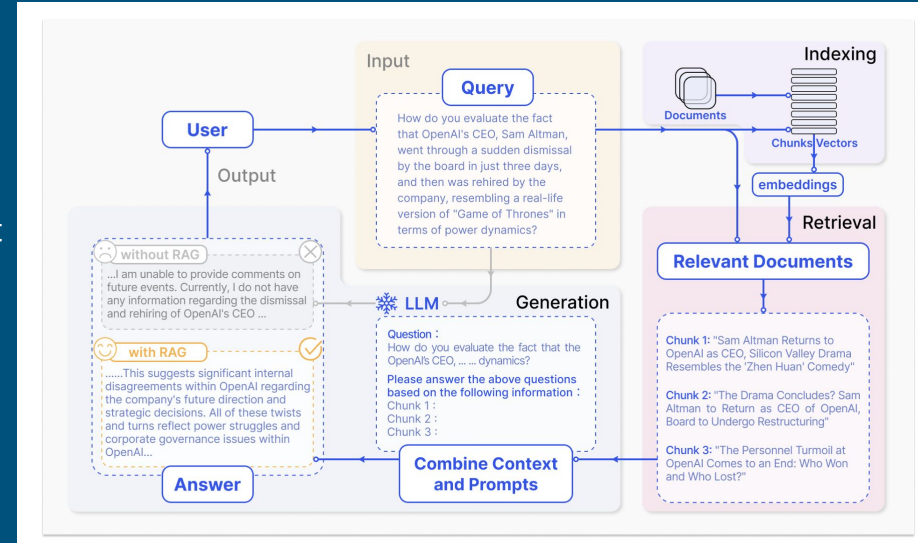  MIDS Applicant Chatbot (Spring 2021)

- Slack bot 👎

# Our Solution and Improvements

RAG based system : Retrieval-Augmented generation is the process of optimizing the output of a large/small language model, so it references an authoritative knowledge base outside of its training data sources before generating a response.

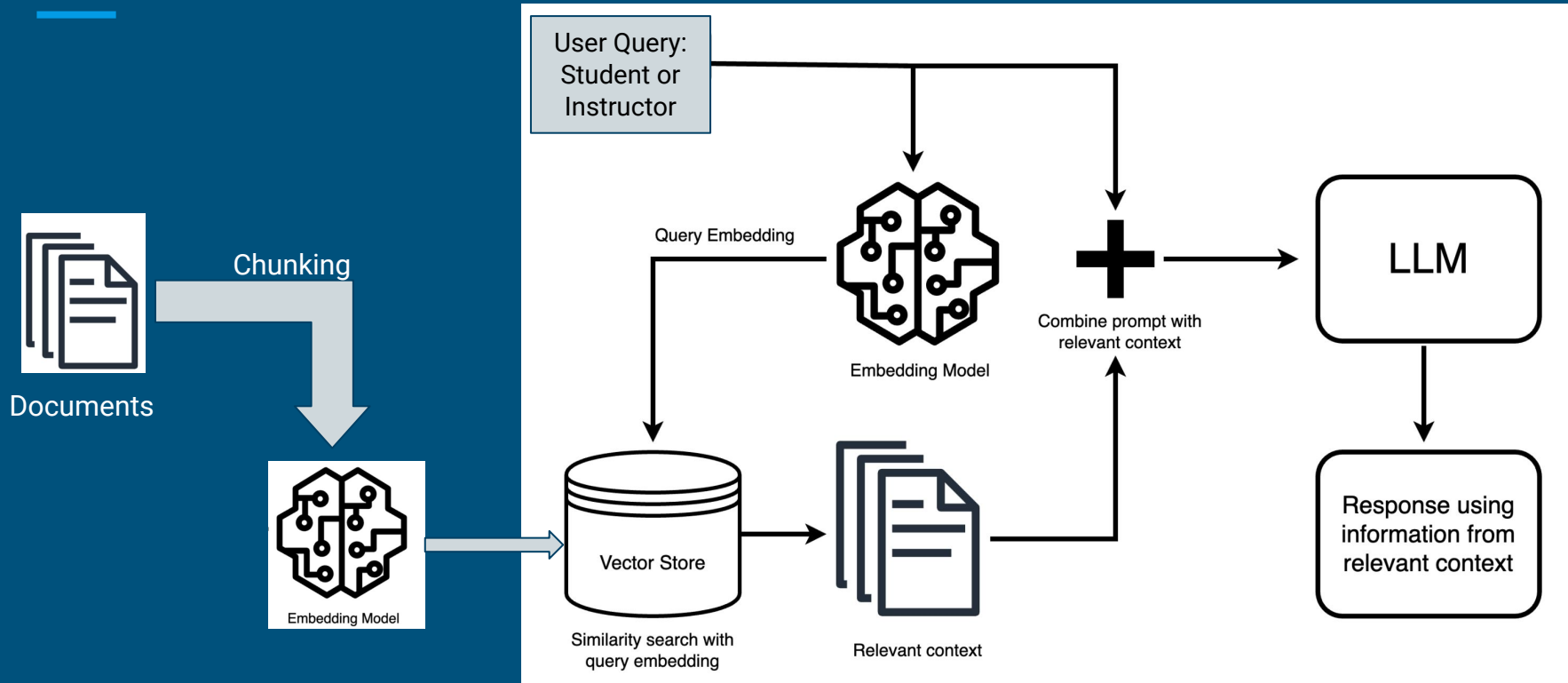Small Language Model : similar to Large language Model but smaller

User type and customers : Students and Instructors.

Goal is to Reduce implementation and training costs, improve reliability, improve security

# Deep Dive into RAG components

# Roadmap - Yes, it's week 5 already

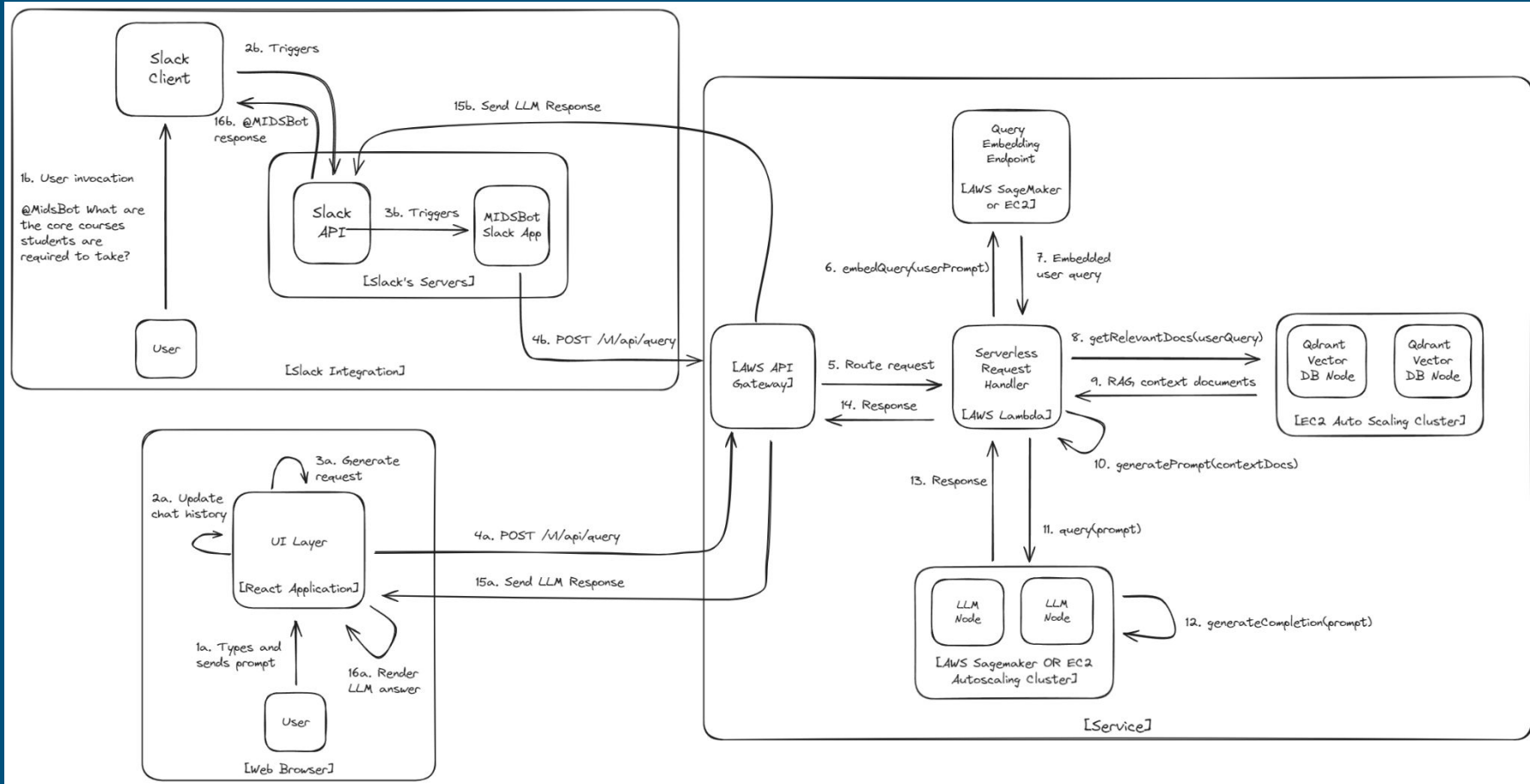| | Week 7 | Week 10 | Week 12 | Week 14 |
|---|---|---|---|---|
| **Language Model** | Baseline LLM selected<br>SLM Settings<br>Default settings implemented | LoRA<br>User Based Prompting<br>Model Compression Techniques | Prompt Tuning<br>Hyperparameter Tuning | API Chaining |
| **Retrieval** | Vector Database<br>Establish Defaults | Search Strategies<br>Reranker | Dynamic Chunking<br>Hybrid Vector Database | Establish Protected Documents |
| **Data** | Collection<br>EDA<br>Survey Users | Generate User Specific Responses<br>PDF Processing | Applying additional datasets | |
| **DevOps** | AWS Initialization | UI Prototype<br>Security Audit | Gather UI Feedback | Host or Containerize |
| **Evaluation** | Research KM KPIs<br>Research NLP Metrics | Apply NLP Experiments | Automated Testing | Human Preference (DPO) |

# MVP Transition to Real Deal

# Datasets

## Data Collection

- **MIDS Intranet:** BeautifulSoup text scraping
- **Bcourses:** PDF processing MIDS Syllabi
- **Other publicly available sites** as needed
- **Slack**: admin export of #mids-class-rec channel
- **LLM Generated Question Answer Pairs**
  - Google Forms surveying students
  - Using openAi to generate questions AND answers based on available documents

## Data Preparation

**Data Cleansing**
- Remove duplicate content (e.g., repetitive links to Facebook)
- Purge unnecessary content to maintain data relevance and reduce noise

Purge or anonymize all personal information attributes

**Data Partitioning by**
- Secure pages
- Public pages
- Course-specific pages

# User Requirements:

Gather user input to understand their requirements, confirm our assumptions about features, user interface elements,

system usage and system functionality

- Google form
- Open Alpha and Beta to small group of students
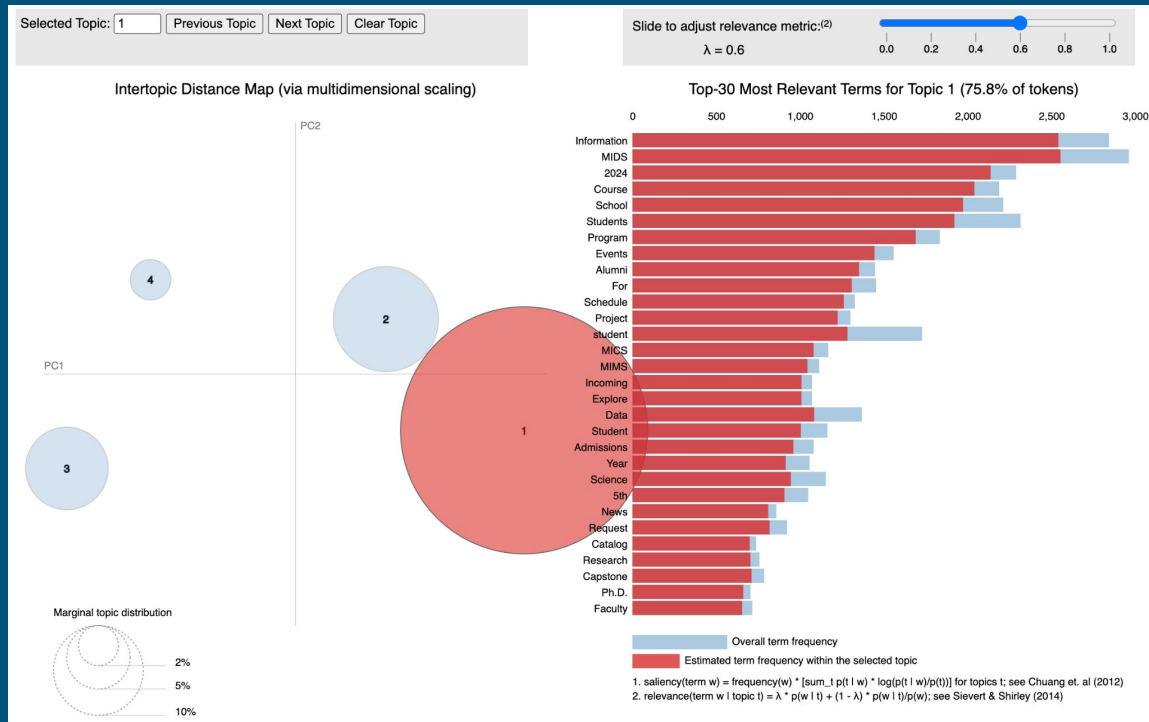- Human feedback via Direct Preference Optimization (DPO)

# EDA

N-gram

Named Entity Recognition

Topic Modeling (LDA)

Topic Complexity

Sentiment Analysis

# Ethics and Privacy

- Use of Slack
  - Terms and Conditions
- LLM considerations
  - Fine tuning on unbalanced dataset
  - Adversarial attacks
  - Guardrails
- Confidential data
  - PII
  - User sensitive information
  - User Privilege Classification on documents for access

# Open Questions and Challenges

- RAG related
  - Generating Reference Answers
  - Model evaluation methods
  - Accessing Slack Data
- Deployment
  - MVP will be hosted locally with a static database
  - Can we build a generative system that run locally for an average users?
- A lot of cool ideas, what are the core features that make the cut?
  - Where will we draw the line ?
  - Not everyone can make varsity

# Conclusion

Knowledge is found in many unstructured sources, we unlock it through a RAG-based approach.

# Notes

Pitching - do we keep SLM? Does that align with BErkeley

- Yes, it minimizes infrastructure and implementation complexity for Berkeley.
-