

Requirements

1. **What are the names and NetIDs of all your team members? Who is the captain? The captain will have more administrative duties than team members.**
 - a. Captain: Hannah Benig (NetID: hhlím2)
2. **What topic have you chosen? Why is it a problem? How does it relate to the theme and to the class?**
 - a. I have chosen the Intelligent Browsing Topic.
 - b. Part 1- intelligent assistant answering questions from recent relevant news articles. Currently, if you are to ask a question in the google search engine, you may receive an answer to your question but will more likely receive a list of documents and then you need to read the documents to deem if they are relevant and find the answer to your question. In this project I plan to create a web application that uses an API to gather recent news articles and then index the gathered documents using a vector database so that users can query the application about recent news events. The output will be a direct answer to the question, the source(s) where the answer was found and the cosine similarity between the query and returned answer. This project will relate to the class in that I will be creating a collection language model (week 4) using recent news articles and will calculate the cosine similarity (week 1) between the query and returned answer.
 - c. Part 2 - Implementing pseudo feedback (week 5) and returning summaries of the top 10 documents for a given topic from the day. A user might need to know updates on a general topic but not have time to browse or read every news article before they head to work. This summarization tab will give them a summary of the top 10 (determined by pseudo feedback) articles for the day updating the user on the news and allowing them to choose which documents they need to read further.
3. **Briefly describe any datasets, algorithms or techniques you plan to use**
 - a. I will use the NewsCatcherAPI to gather news articles (<https://www.newscatcherapi.com/>)
 - b. I will use Langchain to create the vector database
 - c. I will subset a publicly available large language model from Huggingface with news articles to retrieve an answer to the user's query
 - d. I will determine the best answer using cosine similarity
4. **How will you demonstrate that your approach will work as expected?**
 - a. Part 1 - I will create a list of 5-10 sample questions with known answers and will determine that the approach has worked if 75% correctness has been achieved.
 - b. Part 2 - if the summary is sufficient in length and content I will deem my approach successful.
5. **Which programming language do you plan to use?**
 - a. Python
6. **Please justify that the workload of your topic is at least $20 \cdot N$ hours, N being the total number of students in your team. You may list the main tasks to be completed, and the estimated time cost for each task.**
 - a. $N=1$, Total Workload ~20 hours
 - b. Tasks to be completed and estimated time of completion
 - i. Understand how the NewsCatcherAPI works and set it up to gather articles and metadata based on a query (2 hours)
 - ii. Intelligent assistant (part 1)
 1. Setup vector database using Langchain (2 hours)
 2. Create a set of evaluation questions then choose best base LLM from huggingface that we will subset with the news articles based on the answers to

- evaluation questions and ease of use (plan to try out Llama-2, Flan UL2, and roBERTa to see which works best) (5 hours)
- 3. Set up LLM to receive context from news articles (1 hour)
- 4. Fine tune LLM parameters (temperature, top-p, top-k, etc) to find which combination gives best answer (3 hours)
- 5. Write cosine similarity function to calculate between query and answer (2 hours)
- iii. Article summary page (part 2)
 - 1. Experiment with prompt engineering to determine which instructions give the best summary (1 hour)
 - 2. Write summarization script (2 hours)
 - 3. Design webpage (1 hour)
- iv. Create streamlit web app for frontend and package everything up for grading/evaluation (4 hours)