# chATLAS

## An AI Assistant for the ATLAS Collaboration

DANIEL MURNANE, GABRIEL FACINI,
RUNZE LI & CARY RANDAZZO

AI4EIC WORKSHOP

NOVEMBER 30, 2023

ATLAS EXPERIMENT

Ask me something!

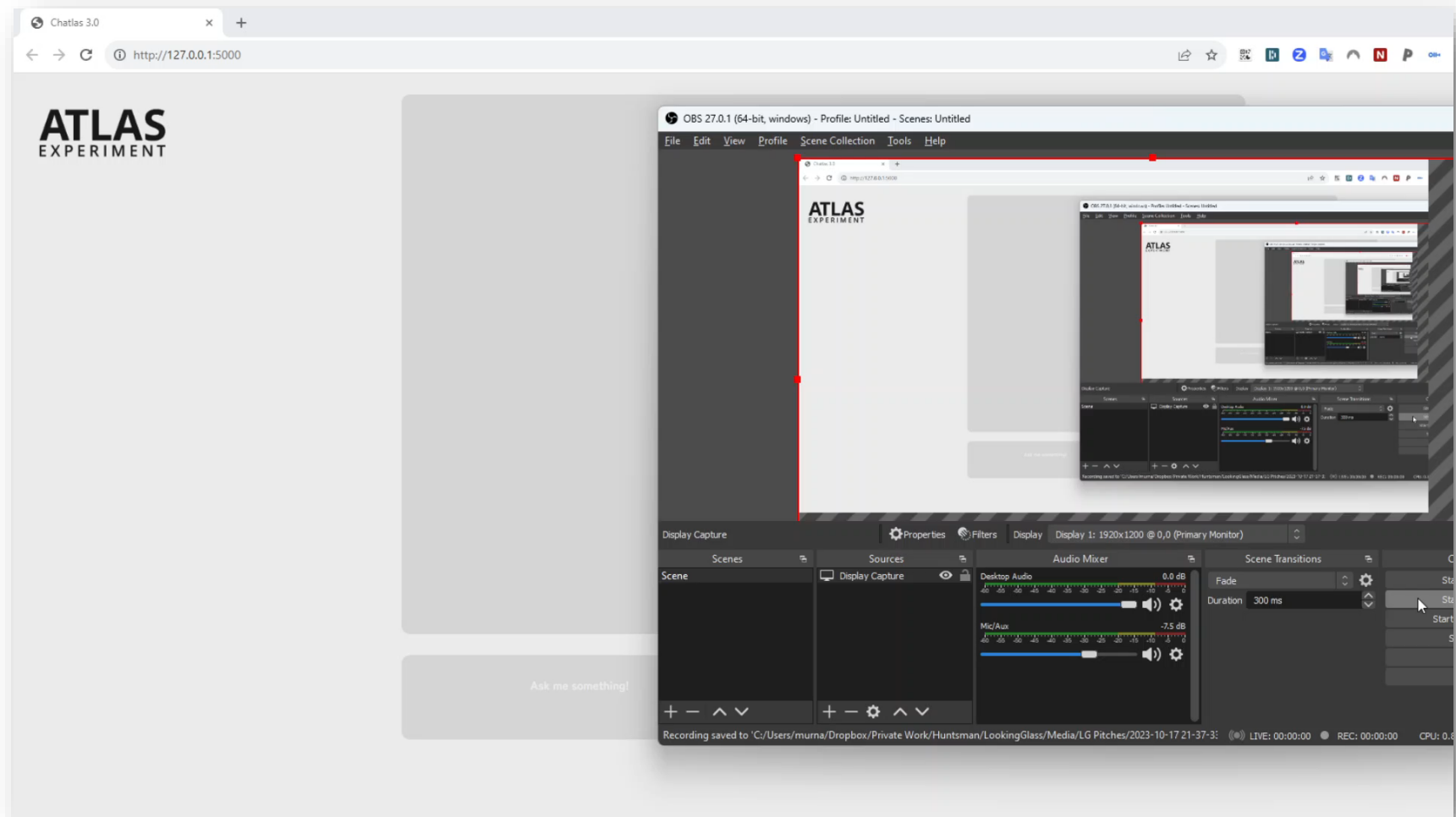BERKELEY LAB  UCL  Yale  LOUISIANA TECH UNIVERSITY

1

# MOTIVATION

- Motivated to have a tool that could be a front-door for...

- Quickly parsing documentation and twiki

- Semantic search and availability of heterogeneous sources of ATLAS information

- Summarizing research

- Connecting the dots between different groups

- Debugging software

- Searching and summarizing JIRA and Glance information
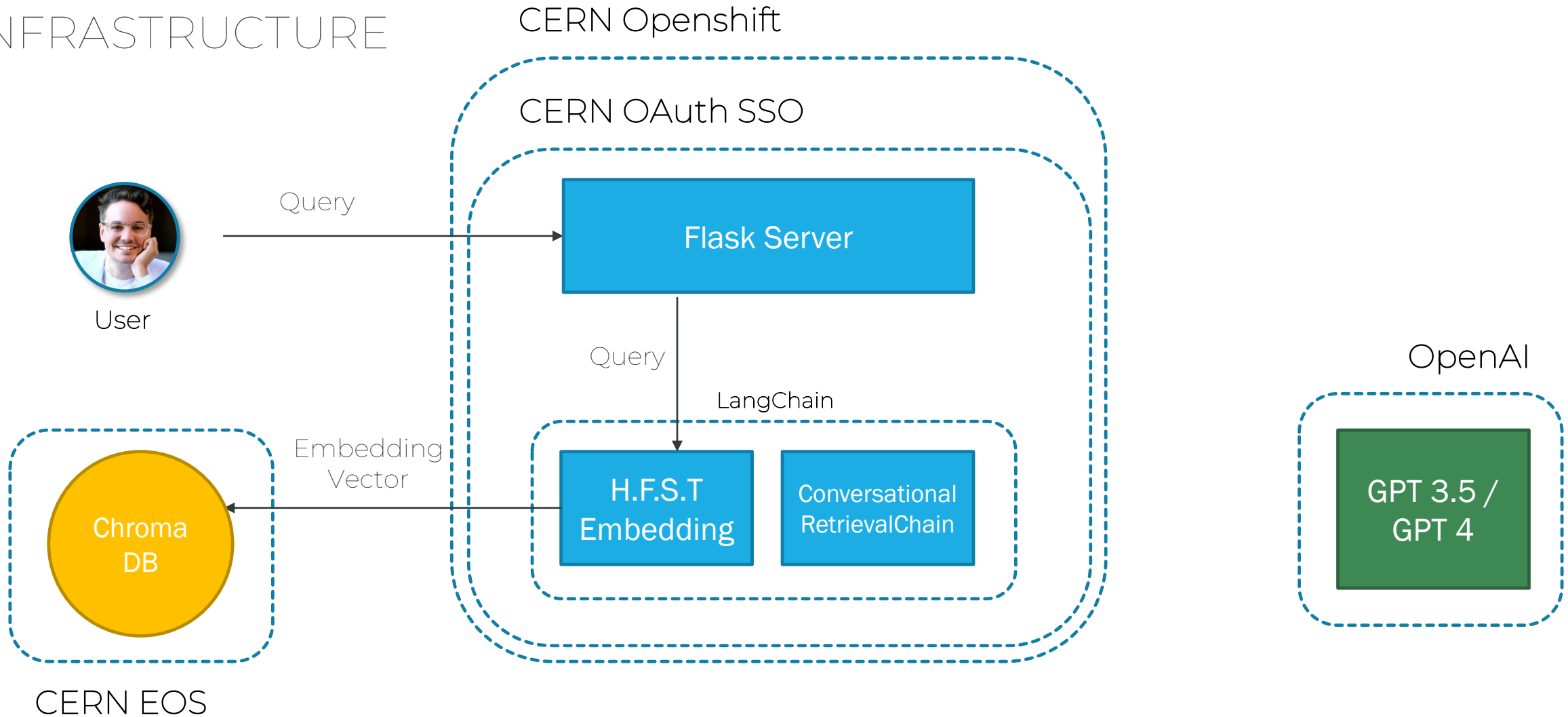
# HISTORY

- In April 2023, initial ATLAS ML Forum meeting to discuss usage of ChatGPT and Github Copilot within ATLAS

- In June 2023, presentations on several ongoing works to use LLMs within ATLAS
  - ATLAS-GPT: Daniel Murnane
  - ChATLAS: Gabriel Facini
  - Google Bard + ATLAS: Kaushik De
  - Analysis Description Language + GPT: Gokhan Unel

- Decision made to converge ATLAS-GPT and ChATLAS and create an official prototype

- Fortnightly developer meetings kicked off in August 2023

- Currently approx. seven part-time contributors
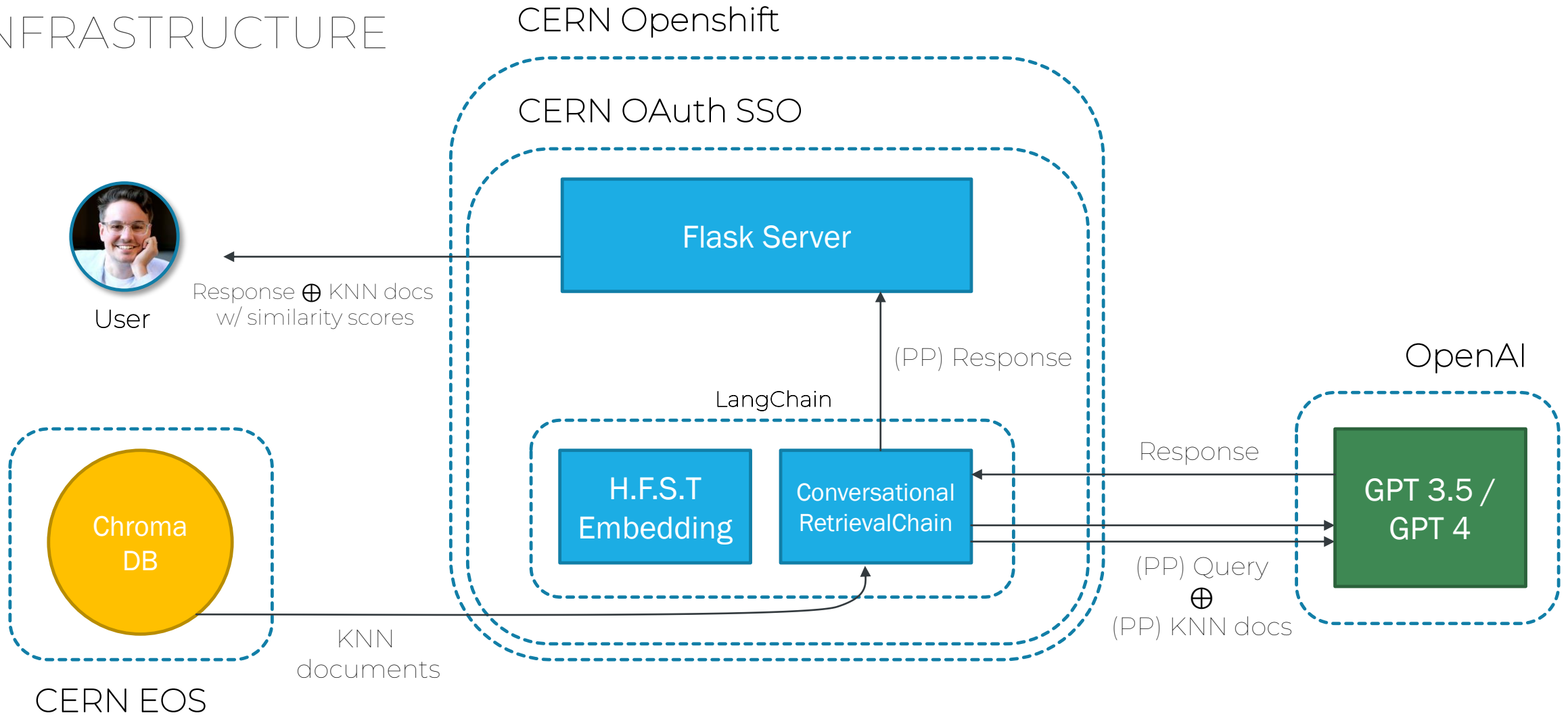
- Launched ATLAS-public demo November 16 https://chatlas-flask-chatlas.app.cern.ch/

# CURRENT STATUS: DEMO

# CURRENT STATUS: INFRASTRUCTURE

(PP) = (Possibly Processed)



CERN Openshift

CERN OAuth SSO

Query

Flask Server

Query

LangChain

H.F.S.T Embedding

Conversational RetrievalChain

User

Embedding Vector

Chroma DB

CERN EOS

OpenAI

GPT 3.5 / GPT 4

# CURRENT STATUS: INFRASTRUCTURE

(PP) = (Possibly Processed)

CERN Openshift

CERN OAuth SSO

OpenAI

**Flask Server**

Response ⊕ KNN docs
w/ similarity scores

User

(PP) Response

LangChain

**H.F.S.T Embedding**

**Conversational RetrievalChain**

Response

**GPT 3.5 / GPT 4**

(PP) Query
⊕
(PP) KNN docs

**Chroma DB**

KNN documents

CERN EOS

# CURRENT STATUS: INFRASTRUCTURE

(PP) = (Possibly Processed)

CERN Openshift

CERN OAuth SSO

- Currently takes ~3 seconds to respond (mostly waiting on OpenAI)

Query

**Flask Server**

Response ⊕ KNN docs w/ similarity scores

**User**

Query

(PP) Response

LangChain

OpenAI

Embedding Vector

**H.F.S.T Embedding**

**Conversational RetrievalChain**

Response

**GPT 3.5 / GPT 4**

**Chroma DB**

(PP) Query ⊕ (PP) KNN docs

KNN documents

**CERN EOS**

# CURRENT STATUS: DATASETS

- Diagram of all the possible ATLAS datasets and how many we have in the DB

- Add a box around those that are *almost* in the DB

- Stage the progress of each (solid/dashed in-progress/complete, scraped, converted, chunked+embedded)

| Twiki | ATLAS Software Docs | E-group Archive | Indico Meetings PDF Text | CDS Papers & Notes PDF Text |
|---|---|---|---|---|
| JIRA | ATLAS Codebases | Group-level Documentation | Mattermost | PDF Plots |

- - - In-progress
----- Complete
Not yet started
Scraping
Converting
Chunking/Embedding

# DETAILS: SCRAPING

## ATLAS Twiki

- Start with set of "Starting URLs"

- Recursively visit included links

- Find all headers, and visit content below

- Append metadata of twiki (parent structure, date revised, etc.)

## CDS

- Discover whether the CDS paper has a Gitlab latex repo

- If latex exists, pull from repo and (planned) convert to markdown

- (Planned) Use **unstructured** library to parse markdown

- If latex *does not* exist, use **nougat** library to read PDF (including equations) into markdown

## Indico

- Load event list

- Scrape timetable contents (date, title, speaker, etc.)

- (Planned) Pull PDF slide decks and minutes

- (Planned) Parse in the same way as in CDS

# DETAILS: CHUNKING

- (Current) Loop through HTML and Markdown heading sections

- If section exceeds 510 tokens, split with SentenceTransformersTokenTextSplitter

- Pass chunk through HuggingFace's sentence-transformers/all-MiniLM-L6-v2 model

- (Planned) Use built-in **unstructured** library to identify chunks

- Insert chunk into Chroma database, with metadata of file URL, twiki name

# DETAILS: RETRIEVAL

- All handled internally by
  qa = ConversationalRetrievalChain.from_llm(    llm=model,
  retriever=db.as_retriever(),    memory=memory,    verbose=False,)

- LLM Model is GPT-3.5 from OpenAI API, retriever is default Chroma which contains the embedding model, memory is a buffer that retains all previous chat information

- Implicit is that the model aggregates all K-documents with a prompt to produce a **new question** based on the original question and the K-documents

Original Query          K documents          Previous chat

**Conversational RetrievalChain**

K documents + processed new query

# DETAILS: CHAT INTERFACE

# DETAILS: CHAT INTERFACE

- Version 1.0 contains everything needed to answer a query

- Can cite the top sources used in the response

- Has a quick search for similar sources

Input query

LLM response

Source title

Similarity score
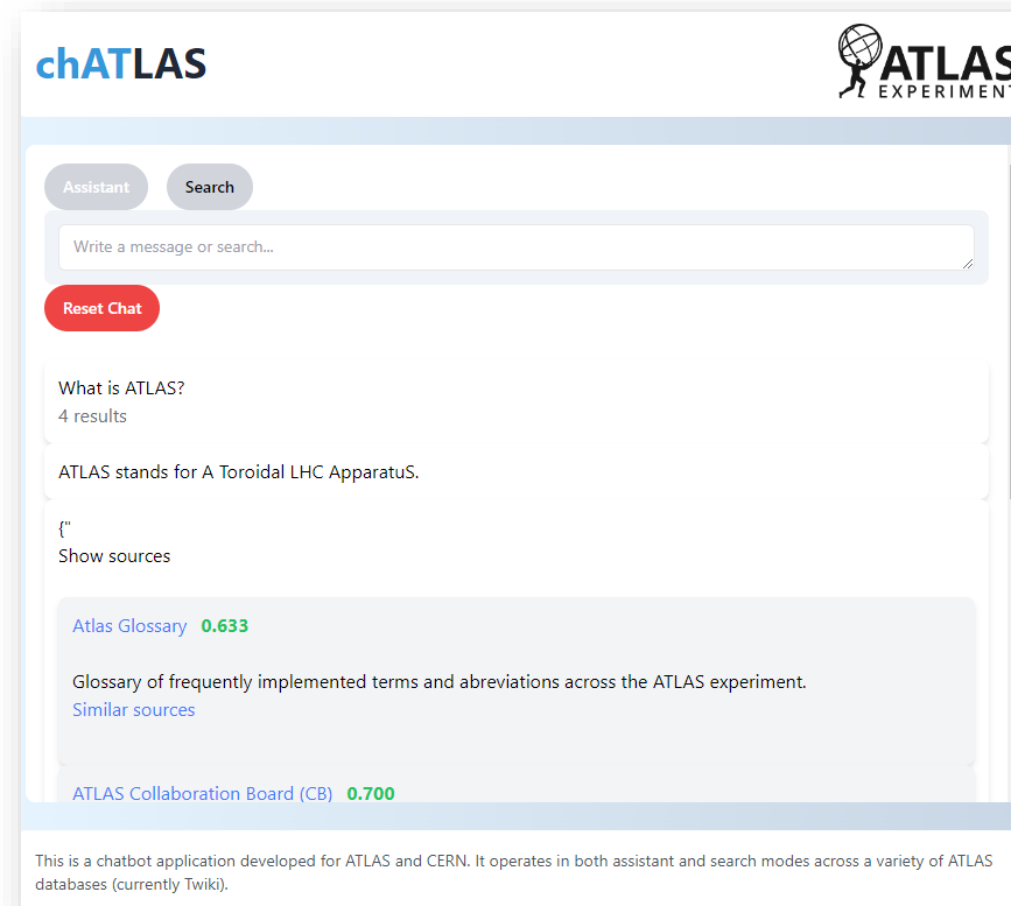
Instant similarity search

Next query entry

# DETAILS: CHAT INTERFACE

- Version 1.0 contains everything needed to answer a query

- Can cite the top sources used in the response

- Has a quick search for similar sources

- Experimenting with a V2.0 appearance that is lighter, and has a dedicated Assistant mode and Search mode

# SHORTCOMINGS AND HURDLES

- Getting the data! Highly heterogeneous file types, many behind authorisation walls, many stale or inaccurate, many requiring high levels of post-processing

- Community solutions could go a long way: Ensure that any experiment/collaboration databases are easily accessible and exportable. All websites should live in a git repo. All publications should be submitted and saved as latex, and compiled separately. All discussion forums should have anonymisation options. This would have saved O(1 year) of data wrangling

- Hallucination is still a very real problem [https://www.arxiv-vanity.com/papers/2311.04348/]

- A high quality AI assistant probably requires fine tuning, which is an expensive task (less in gpu-hours, more in expert-hours)

- Open-source solutions for UI are not particularly flexible. A tool built by+for the scientific community would be **very** useful! Open-source solutions for backend (retrieval, document aggregation) are perfectly fine.

- Codebase integration: experiment codebases are huge, not so well-commented, and non-obvious how to chunk. Perhaps an automated commenting algorithm as a pre-process step?
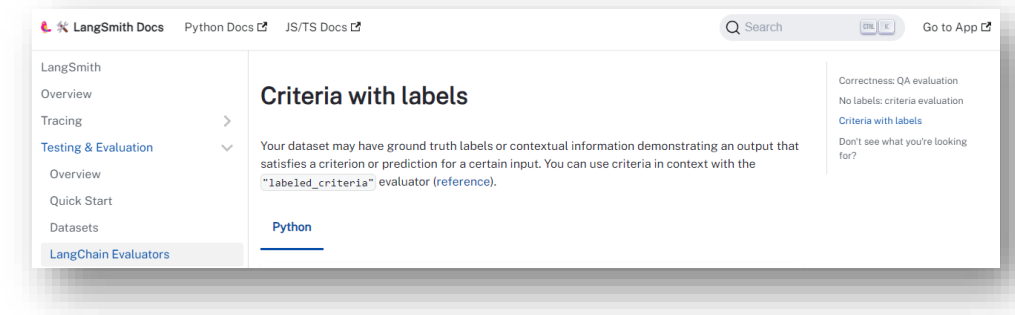
**ROADMAP**

**Datasets**

| | | | | |
|---|---|---|---|---|
| Twiki | ATLAS Software Docs | CDS Papers & Notes PDF Text | ATLAS Codebases | JIRA |
| E-group Archive | Indico Meetings PDF Text | Mattermost | Group-level Documentation | PDF Plots |

**Infrastructure & Models**

| | | | | |
|---|---|---|---|---|
| HuggingFace & OpenAI API | HuggingFace both Embedding & Completion | More sophisticated document aggregation | Vision model for plots and diagrams | Fine-tuned Completion Model |

**Interface & Features**

| | | | | |
|---|---|---|---|---|
| Single-query Assistant & Simple search | Filtered search & Continued-conversation Assistant | Integration with Gitlab and IDEs | Integration with other CERN AI Assistants | Agentic – categorize, review papers, monitor Glance and Jira |

Today     Jan '24     March '24     July '24

chATLAS    BERKELEY LAB

# OPEN QUESTIONS & CONCLUSION

- How to avoid hallucinations? Probably "GPT-5 / Q* / AGI" will make this hurdle irrelevant

- How to best "censor" politically incorrect responses (e.g. which analysis team is the best?)

- How to **measure** the quality of responses – LangSmith AI-assisted evaluators?

- What is the best dataset to gather for fine-tuning?

- How to anonymize email threads and discussion forums?

We are having a lot of fun building this thing from scratch, but if there was an open-source scientific community framework for AI Assistants, it would be even more fun!

# BACKUP

# RETRIEVAL-AUGMENTED GENERATION

# RETRIEVAL-AUGMENTED GENERATION (RAG)



1. Convert document to txt



arXiv:1207.7214v2 [hep-ex] 31 Aug 2012EUROPEANORGANISATION FOR NUCLEARRESEARCH(CERN) CERN-PH-EP-2012-218 Acceptedby: PhysicsLettersB Observation ofa NewParticleinthe Searchfor the Standard Model HiggsBoson withthe ATLAS Detector atthe LHC TheATLAS Collaboration This paper is dedicated tothe memory of our ATLAS colleagues whodid not livetosee the full impact and signicance of their contributions totheex periment. Abstract A search for the Standard Model Higgs boson in proton-proton collisions with the ATLAS detector at the LHC is presented. The datasets used correspond to inte grated luminosities of approximately 4.8fb1collectedats=7TeVin2011and5.8fb1ats=8TeVin2012. Individualsearchesinthe ...

# RETRIEVAL-AUGMENTED GENERATION (RAG)

arXiv:1207.7214v2 [hep-ex] 31 Aug 2012EUROPEANORGANISATION FOR NUCLEARRESEARCH(CERN) CERN-PH-EP-2012-218 Acceptedby: PhysicsLettersB Observation ofa NewParticleinthe Searchfor the Standard Model HiggsBoson withthe ATLAS Detector atthe LHC TheATLAS Collaboration This paper is dedicated tothe memory of our ATLAS colleagues whodid not livetosee the full impact and signicance of their contributions totheex periment. Abstract A search for the Standard Model Higgs boson in proton-proton collisions with the ATLAS detector at the LHC is presented. The datasets used correspond to inte grated luminosities of approximately 4.8fb1collectedats=7TeVin2011and5.8fb1ats=8TeVin2012. Individualsearchesinthe ...

arXiv:1207.7214v2 [hep-ex] 31 Aug 2012EUROPEANORGANISATION FOR NUCLEARRESEARCH(CERN) CERN-PH-EP-2012-218 Acceptedby: PhysicsLettersB Observation ofa NewParticleinthe Searchfor the Standard Model HiggsBoson withthe ATLAS Detector atthe LHC TheATLAS Collaboration This paper is dedicated ...

This paper is dedicated tothe memory of our ATLAS colleagues whodid not livetosee the full impact and signicance of their contributions totheex periment. Abstract A search for the Standard Model Higgs boson in proton-proton collisions with the ATLAS detector at the LHC is presented. ...

1. Convert document to txt

2. Split into K-word chunks

# RETRIEVAL-AUGMENTED GENERATION (RAG)

arXiv:1207.7214v2 [hep-ex] 31 Aug 2012EUROPEANORGANISATION FOR NUCLEARRESEARCH(CERN) CERN-PH-EP-2012-218 Acceptedby: PhysicsLettersB Observation ofa NewParticleinthe Searchfor the Standard Model HiggsBoson withthe ATLAS Detector atthe LHC TheATLAS Collaboration This paper is dedicated tothe memory of our ATLAS colleagues whodid not livetosee the full impact and signicance of their contributions totheex periment. Abstract A search for the Standard Model Higgs boson in proton-proton collisions with the ATLAS detector at the LHC is presented. The 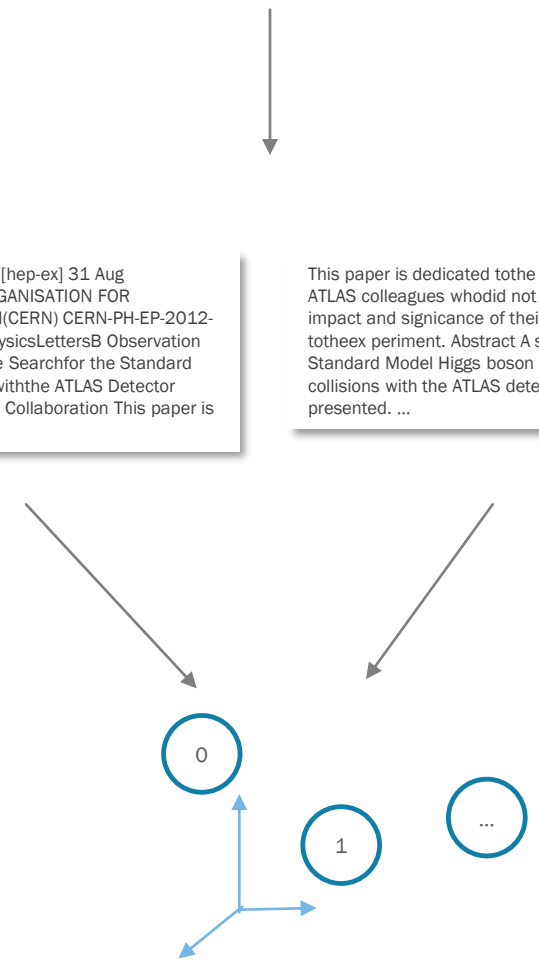datasets used correspond to inte grated luminosities of approximately 4.8fb1collectedats=7TeVin2011and5.8fb1ats=8TeVin2012. Individualsearchesinthe ...

1. Convert document to txt

2. Split into K-word chunks

3. Embed each chunk with HuggingFace SentenceTransformer model into 1536-dimensional vector space

arXiv:1207.7214v2 [hep-ex] 31 Aug 2012EUROPEANORGANISATION FOR NUCLEARRESEARCH(CERN) CERN-PH-EP-2012-218 Acceptedby: PhysicsLettersB Observation ofa NewParticleinthe Searchfor the Standard Model HiggsBoson withthe ATLAS Detector atthe LHC TheATLAS Collaboration This paper is dedicated ...
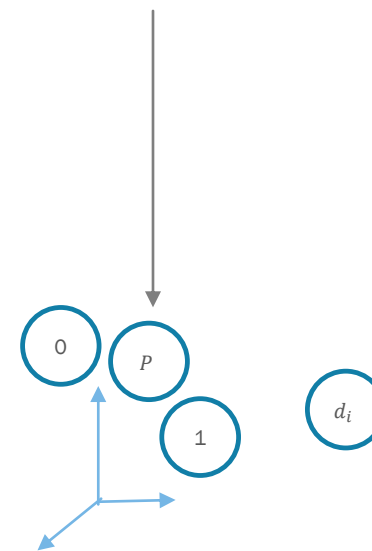
This paper is dedicated tothe memory of our ATLAS colleagues whodid not livetosee the full impact and signicance of their contributions totheex periment. Abstract A search for the Standard Model Higgs boson in proton-proton collisions with the ATLAS detector at the LHC is presented. ...

0

1

...

# HOW IT WORKS

1. Convert document to txt

2. Split into K-word chunks

3. Embed each chunk with HuggingFace SentenceTransformer model into 1536-dimensional vector space

4. Embed prompt $\boldsymbol{P}$ in this space

Q In the Higgs observations tudy, what was the discovery significance?

# HOW IT WORKS

1. Convert document to txt

2. Split into K-word chunks

3. Embed each chunk with HuggingFace SentenceTransformer model into 1536-dimensional vector space

4. Embed prompt $P$ in this space

5. Nearest neighbour search (KNN, $K \leq 10$, $P \cdot d_i \leq 0.5$)

Q   In the Higgs observations tudy, what was the discovery significance?

# HOW IT WORKS

1. Convert document to txt

2. Split into K chunks

3. Embed each chunk with HuggingFace SentenceTransformer model into 1536-dimensional vector space
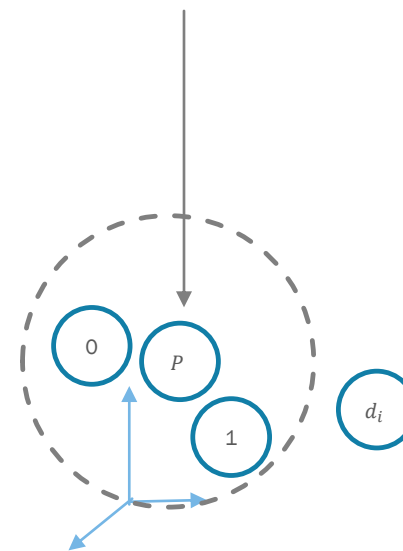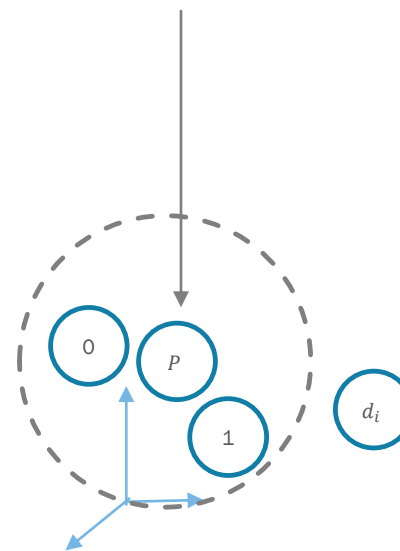
4. Embed prompt $P$ in this space

5. Nearest neighbour search (KNN, $K \leq 10$, $\mathrm{P} \cdot d_i \leq 0.5$)

6. Append all matching sections to prompt and query GPT-4

arXiv:1207.7214v2 [hep-ex] 31 Aug 2012EUROPEANORGANISATION FOR NUCLEARRESEARCH(CERN) CERN-PH-EP-2012-218 Acceptedby: PhysicsLettersB Observation ofa NewParticleinthe Searchfor the Standard Model HiggsBoson withthe ATLAS Detector atthe LHC TheATLAS Collaboration This paper is dedicated ...

$\oplus$

This paper is dedicated tothe memory of our ATLAS colleagues whodid not livetosee the full impact and signicance of their contributions totheex periment. Abstract A search for the Standard Model Higgs boson in proton-proton collisions with the ATLAS detector at the LHC is presented. ...

$\oplus$

Q   In the Higgs observations tudy, what was the discovery significance?