

chATLAS

An AI Assistant for the ATLAS Collaboration

DANIEL MURNANE, GABRIEL FACINI,
RUNZE LI, DANIELE DEL SANTO & CARY RANDAZZO

1st Large Language Models in Physics Symposium
February 21, 2024



chATLAS Overview

- Motivation
- History of chATLAS
- Data Gathering
- Infrastructure
- Deployment Result: Chat Interface and Demo
- Roadmap
- Updates, Challenges, and Summary

MOTIVATION

- 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/>

Data Gathering: How Scraping Began

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

Data Gathering: Volume Measurements


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

Data Gathering: Volume Measurements and Tool Selection

Twiki

Over 2000 ATLAS Twiki Topics

Using:
BeautifulSoup and
auth-get-sso-cookie

Twiki > AtlasProtected Web > AtlasPhysics (2024-02-08, MariaCecilia)	
 ATLAS Physics Activity	
Physics coordinators: Monica Dunford, Fabio Cerutti	
Analysis - Getting started and key links	
Getting set-up Data checklist for physics analysis - Starting page for GEANA, data and MC ASO info - Useful FAQ, info on PHYS/PHYS/SLITE, software tutorials, who to ask for help	
Finding your samples Data and Monte Carlo Datasets - Starting point for all data/MC information Derivations MC Sample Request Procedure	
Trigger and Object calibration CP analysis Kick-off work - CP work for your analysis Run 2 CP recommendations - Status and links Run 2 CP recommendations - Status and links Lowest unscaled triggers Luminosity	
High priority CP work Tracking Etanimes Muon Tau JetEtnimes Flavour tagging Simulation PMO Physics validation	
Hepdata and public data ATLAS Hepdata recommendations - helpful tips and scripts to get started	
General links How the ATLAS collaboration works - All helpful ATLAS links in one place! Guidelines to produce event displays - create event displays Run 2 Discrepancies and excesses	
Paper writing links and tips	
Approval of Results - links approval procedure and policies PubComm main page - Contains all instructions, latex templates, etc Conference deadlines for papers and CONF notes - timelines for major conferences ATLAS figure style recommendations (Jun 2014)¶ ATLAS meeting rules (Mar 2014)¶ Physics Office - policy documents, glance-related instructions, egroups naming rules Glance Phase 0 instructions - including who is on what egroup	
Policy Documents ATLAS Policy Documents - includes policy on theses, job talks, etc	
Events	
ATLAS physics workshops - Workshops and guidelines on away meetings ATLAS posters at LHCC meeting, November 27th, 2023	
Physics Organisation	
Physics Coordination Membership Operation Task Planner (OTPP)¶ SCAB - Instructions on SCAB List of qualification tasks Compact list of talks at main meetings (Indicomb)¶	
Task forces	
Combined Performance (CP) Groups	
Eligemma	L. Aperio Bella, K. Lohwasser
Flavour Tagging	F.A. Di Bello, D. Guest
Inner Tracking	C. Giffels, N. Calace
JetEtnimes	R. Camacho Toro, F. Ball
Muon	R. Nikolaidou, S. Angelidakis
Tau	S.M. Farrington, A. De Maria
Physics Analysis (PA) Groups	
B Physics & Light States	A. Cerni, S. Turchetta
Exotics	T. Vazquez Schroeder, D. Hayden
Heavy Ions	A. Angerami, Q. He
Higgs	N. Berger, P. Francavilla
Higgs & Diboron Searches	E. Bross, A. Cortes Gonzalez
Physics Modelling	D. Hirschbuhl, A. Cueto Gomez
Standard Model	P. Sommer, S. Camarda
Supersymmetry	J. Monopoli Berlingieri, S. Alderweireld
Top	A. Knuhl, N.A. Abdah
Upgrade Physics	A. Schwartzman, H. De La Torre Per
List of PA & CP subgroups and subgroup conveners Physics and CP Group Liaisons ATLAS Appointments Database¶ List of CMS conveners¶	
Other Groups related to PC	
Analysis Model Group	E. Kourilis, G. Watts
Boosted Higgs Inco Tagger	A. Coccia, F. Filthaut
Derivation Coordination and Production Team	E. Toro Pastor
Forward proton performance	R. Staszewski, S.E. Clawson
Global particle flow task force	M. Hodgkinson, M. Swiatkowski
Isolation and Fake Forum	N. Bouchie, F. Alonso
Join EFT and Interpretation group	S. Kotler, T. Dado, H. Midher
Luminosity WG (data preparation) (info)	V.S. Lang, E. Torrence
Machine Learning Forum	W.H. Hopkins, D.T. Murnane
MC Production	E.M. Lobodzinska, Y.L. Liu
Physics Office	G. Navarro

Data Gathering: Volume Measurements and Tool Selection

**ATLAS
Software
Docs**

Using:
BeautifulSoup and
auth-get-sso-cookie

Hundreds of ATLAS Software Docs

ATLAS Software Documentation Guides ▾ Analysis SW Tutorial ▾ Other Tutorials ▾ Links ▾ Sea

Athens Developers
Athena Configuration
Trigger Developers
CMake Configuration
Release Coordinators
Merge Request Review Shifters
Building a release
Analysis Tools
W IDE Integrations & VS Code
He Centralized Ntuples Production

Documentation

**ATLAS
EXPERIMENT**

umentation pages.
written and reviewed by experts.

These pages contained structured Software documentation, and exist in addition to the [twiki](#) [ATLAS members only], and are intended to be more authoritative.

For comments and suggestions, please feel free to use the [issue tracker](#).

If you want to update or modify the content yourself, please feel free to follow the [contribution guide](#)

Data Gathering: Volume Measurements and Tool Selection













E-group Archive

Atlas forums

<input type="checkbox"/> Archive	E-group name	Description
Category : Computing Documentation and Announcements (5)		
CERN Computing Announcements	hn-atlas-cernCompAnnounce	This forum will contain announcements of changes, outages and other
Documentation and Communication	hn-atlas-docAndCommunication	This forum covers general announcements concerning Documentation
Grid Announcements	hn-atlas-gridAnnounce	This forum contains announcements related to the Grid infrastructure
Releases and Distribution Kit Announcements	hn-atlas-releaseKitAnnounce	This Forum announces the Release Plans and New Releases or patches. Do not reply to the postings, send comments or reports on problems to
Software Developers Announcements	hn-atlas-SWDevelopersAnnounce	Announcements, which are intended for developers of ATLAS offline software. Do not reply to postings in this Forum, send comments or reports on p
Category : Computing Offline Software (28)		
Architecture Team: Core Software Architecture and Design	hn-atlas-SWArchitecture	A forum to follow and participate to the work of the ATLAS Architecture
Athena-ROOT access	hn-atlas-athena-ROOT-access	Developer's discussion of progress in accessing Athena objects from I
Atlantis Event Display	hn-atlas-AtlantisDisplay	Forum for help requests, suggestions, discussions and comments reg
Attfast Support	hn-atlas-attfast-support	Support for Atlas fast simulation users
Attfast Working Group	hn-atlas-attfast2-newsdev-val	This forum is aimed for the development, validation and testing of the
Bugs	hn-atlas-Prelimbugs	Atlas Preliminary/unconfirmed Bugs, Problems, Frustrations, Fixes
Bytestream Initiative	hn-atlas-bytestream-initiative	A finite-lifetime forum to coordinate bytestream infrastructure develop
Digitization Developers	hn-atlas-digitization-developers	This forum is intended for coordinating and discussing digitization soft
Fatras News, Development and Validation	hn-atlas-fatras	This forum is aimed at the developers of Fatras such as the users. In
General Offline Help	hn-atlas-offlineSWHelp	Requests for guidance in writing or using software. Comments on rec. Please report release specific errors and problems to the Forum "Rele
Generator Validation	hn-atlas-generator-validation	Common forum for MC generator responsibilities in ATLAS
Forum for VP1 support	hn-atlas-vp1-help	Forum for VP1 (Virtual Point 1) support
New Job Configuration	hn-atlas-NewJobConfiguration	In this forum we will discuss the migration to Configurables, JobPrope
Offline Commissioning	hn-atlas-offline-commissioning	Offline Commissioning: For discussion about the offline software usec
Offline SW Development Discussions	hn-atlas-offlineSWDevelopment	This group is for the Athena and Package developers to discuss desig
Persistence Help	hn-atlas-persistenceHelp	This forum is a general help list for reading and writing data. It is mea
Physics and Software Validation	hn-atlas-physics-software-validation	Common forum for discussions and announcements concerning Soft
PileUp	hn-atlas-pileup	To discuss and improve pile-up simulation, reconstruction and perform
Reconstruction Bug Monitor	hn-atlas-reco-bug-monitor	Forum for developers who wish to receive mail whenever a Reconstru
Reconstruction Integration	hn-atlas-recoIntegration	This forum is the main discussion channel for the Reconstruction Inte. Although the RIG has a number of defined members, this forum is op. For software development question not particularly related to reconstr
Releases and Distribution Kit Problems	hn-atlas-releaseKitProblem	This Forum is for requests for assistance with problems running the sc

8912 topics x (1-50+ messages) for largest egroup
~ tens of thousands of messages to date

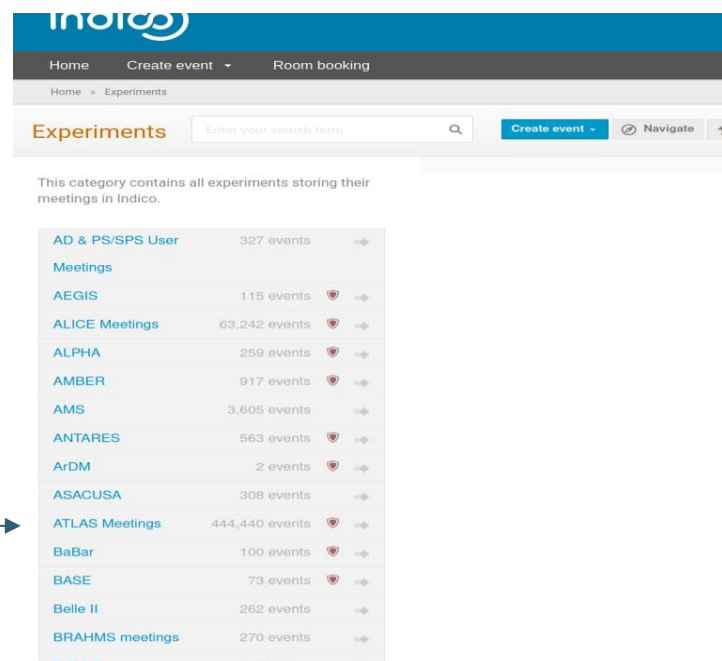
Using:
BeautifulSoup and
Selenium

Name	Size
 hn-atlas-lar-electronic-calibration.txt	304.5 MB
 hn-atlas-triggerReleaseValidation.txt	84.6 MB
 hn-atlas-TDAQCommissioning.txt	81.1 MB
 hn-atlas-tile-performance.txt	68.8 MB
 hn-atlas-dist-analysis-help.txt	39.2 MB
 hn-atlas-muonSWBugMonitor.txt	26.6 MB
 hn-atlas-data-quality-operations.txt	18.1 MB
 hn-atlas-offlineSWHelp.txt	17.4 MB
 hn-atlas-PATHelp.txt	17.2 MB
 hn-atlas-SITInternal.txt	8.2 MB
 hn-atlas-exotics-wg.txt	6.9 MB
 hn-atlas-jetmiss-wg.txt	6.5 MB

Data Gathering: Volume Measurements and Tool Selection

Indico Meetings

Using:
Nougat and Marker



The screenshot shows the Indico website's 'Experiments' page. It features a navigation bar with 'Home', 'Create event', and 'Room booking'. Below the navigation bar is a search bar and a 'Create event' button. The main content area displays a list of experiments with their respective meeting counts. An arrow points from the text '440,440 ATLAS Indico Meeting Events' to the 'ATLAS Meetings' entry in the list.

Experiment	Events
AD & PS/SPS User	327 events
Meetings	
AEGIS	115 events
ALICE Meetings	63,242 events
ALPHA	259 events
AMBER	917 events
AMS	3,605 events
ANTARES	563 events
ArDM	2 events
ASACUSA	308 events
ATLAS Meetings	444,440 events
BaBar	100 events
BASE	73 events
Belle II	262 events
BRAHMS meetings	270 events
CALET	76 events

440,440 ATLAS Indico
Meeting Events

Data Gathering: Volume Measurements and Tool Selection

**PDF
Plots** **Mattermost** **Jira** **ATLAS
Codebases** **Group
level
Docs**

Measurements Pending - very large

Using:
Pending further
experiments

Data Gathering: Volume Measurements and Tool Selection

CDS Papers & Notes

The screenshot shows the CERN Document Server interface. At the top, there's a navigation bar with 'Search', 'Submit', 'Help', and 'Personalize' buttons. Below this, the breadcrumb trail reads 'Home > CERN Experiments > LHC Experiments > ATLAS'. The main heading is 'ATLAS'. A search bar indicates 'Search 66,465 records for:' with a 'Search' button and links to 'Search Tips' and 'Advanced Search'. Below the search bar is an 'Add to Search' button. The 'Narrow by collection:' section lists various ATLAS collections with checkboxes and counts: ATLAS Papers (1,253), ATLAS Reports (37), ATLAS Conference Notes (1,250), ATLAS Notes (10,087), ATLAS Scientific Notes (69), ATLAS Theses (2,799), ATLAS Conference Slides (11,933), ATLAS Videos (596), ATLAS Footage (0), ATLAS Photos (2,516), ATLAS Event Displays (13), ATLAS eNews (250), ATLAS Preprints (2,578), and ATLAS Internal (36,405). At the bottom, there are links to 'ATLAS Communications (32,265)', 'ATLAS Internal Notes (1,195)', 'ATLAS Publication Drafts (72)', 'ATLAS Plots (136)', 'ATLAS Live News (185)', 'Restricted ATLAS Talks (2,546)', and 'Restricted ATLAS Event Displays (6)'.

Using:
Nougat and Marker

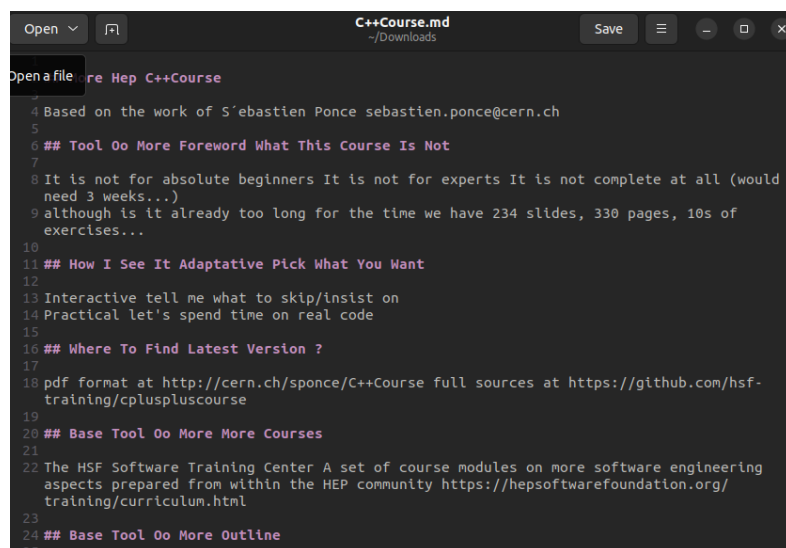
Total: 66,465 records
found

36,405 are Internal
including
communications,
notes, etc.

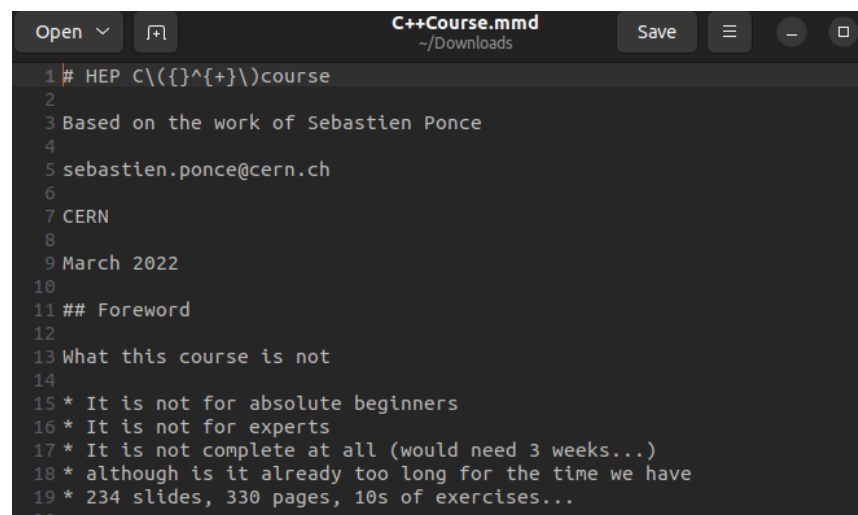
Data Gathering: A Tool Comparison

CDS Papers & Notes

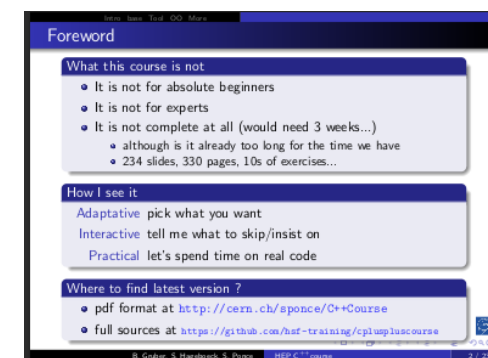
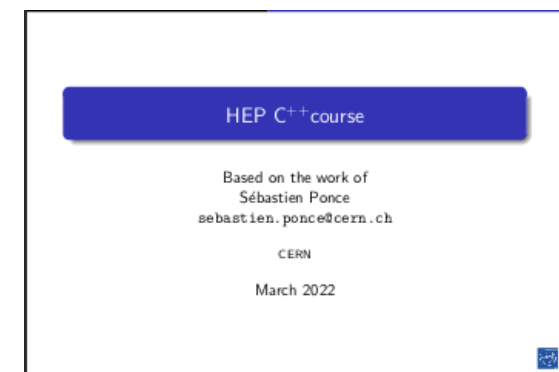
Slides more difficult to process
and largest data in volume



```
Open file: C++Course.mmd
~/Downloads
Save
Open file: re Hep C++Course
4 Based on the work of S'ebastien Ponce sebastien.ponce@cern.ch
5
6 ## Tool Oo More Foreword What This Course Is Not
7
8 It is not for absolute beginners It is not for experts It is not complete at all (would
  need 3 weeks...)
9 although is it already too long for the time we have 234 slides, 330 pages, 10s of
  exercises...
10
11 ## How I See It Adaptative Pick What You Want
12
13 Interactive tell me what to skip/insist on
14 Practical let's spend time on real code
15
16 ## Where To Find Latest Version ?
17
18 pdf format at http://cern.ch/sponce/C++Course full sources at https://github.com/hsf-training/cpluspluscourse
19
20 ## Base Tool Oo More More Courses
21
22 The HSF Software Training Center A set of course modules on more software engineering
  aspects prepared from within the HEP community https://hepsoftwarefoundation.org/training/curriculum.html
23
24 ## Base Tool Oo More Outline
25
```



```
Open file: C++Course.mmd
~/Downloads
Save
1 # HEP C\({}^{+}\)\course
2
3 Based on the work of Sebastien Ponce
4
5 sebastien.ponce@cern.ch
6
7 CERN
8
9 March 2022
10
11 ## Foreword
12
13 What this course is not
14
15 * It is not for absolute beginners
16 * It is not for experts
17 * It is not complete at all (would need 3 weeks...)
18 * although is it already too long for the time we have
19 * 234 slides, 330 pages, 10s of exercises...
20
```



Faster

Marker vs Nougat
(PDF scraping)

Processing: Chunking and Retrieval

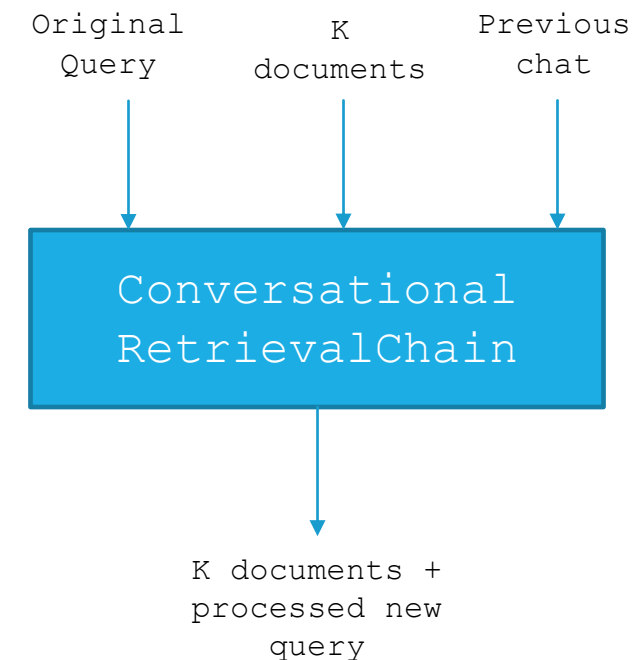
Chunking

- (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

Retrieval

- All handled internally by

```
ga = 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



Summary of Gathering and Processing

Summary includes...

- Diagram of all the possible ATLAS datasets and how many we have in the Database(DB)
- **Chunked & Embedded Datasets** are ready for or have been added to DB
- Stage the progress of each
- Over 37% of textualizable ATLAS datasets (minimum)
- Hardware upgrades were made to scrape in larger volume

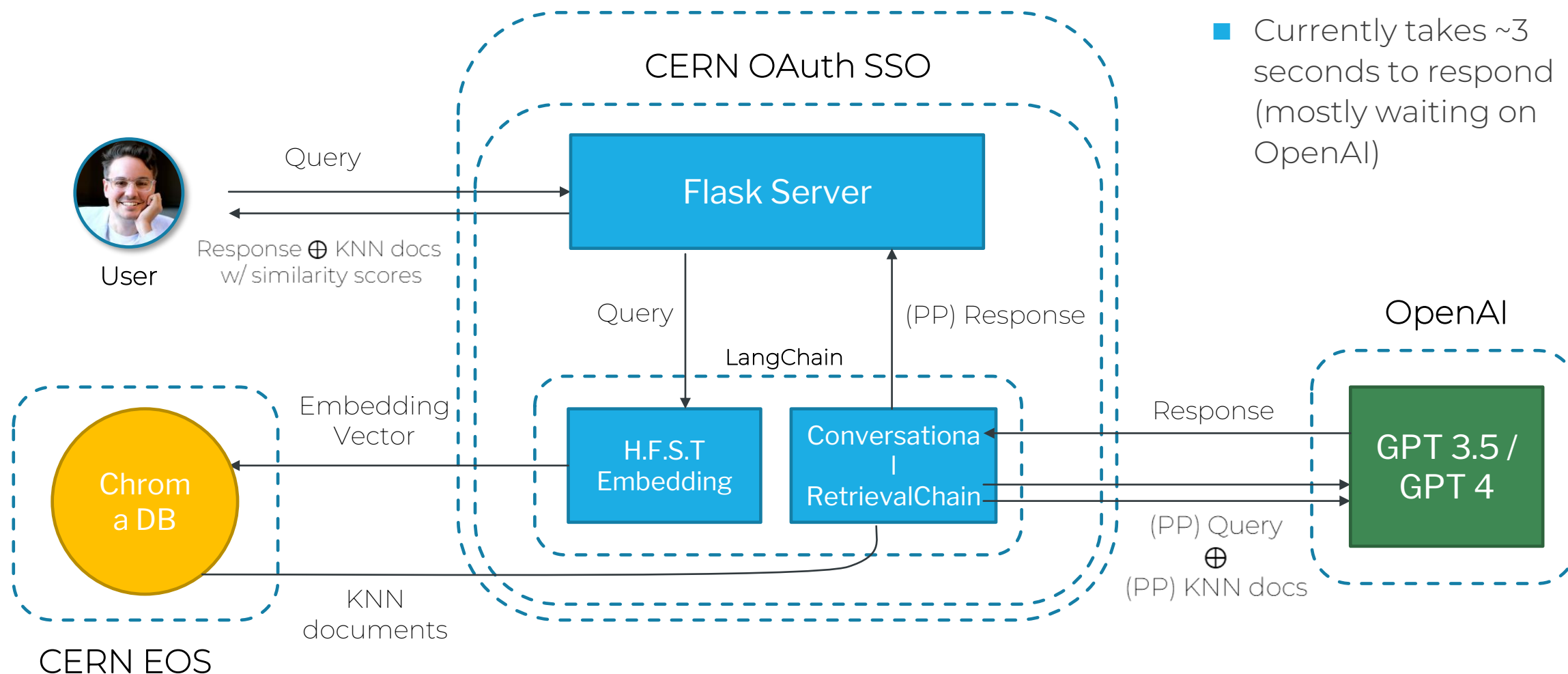
	DB	Twiki	ATLAS Software Docs	E-group Archive	Indico Meetings	PDF Plots	Mattermos t	Jira	ATLAS Codebases	Group level Docs	CDS Papers & Notes
Task											
Scrape		2k+	500+	10k+	~1k+	~1k+					~5k+
Convert											
Chunk & Embed											

Not yet started
In Progress
Complete

Current INFRASTRUCTURE

(PP) = (Possibly Processed)

CERN Openshift



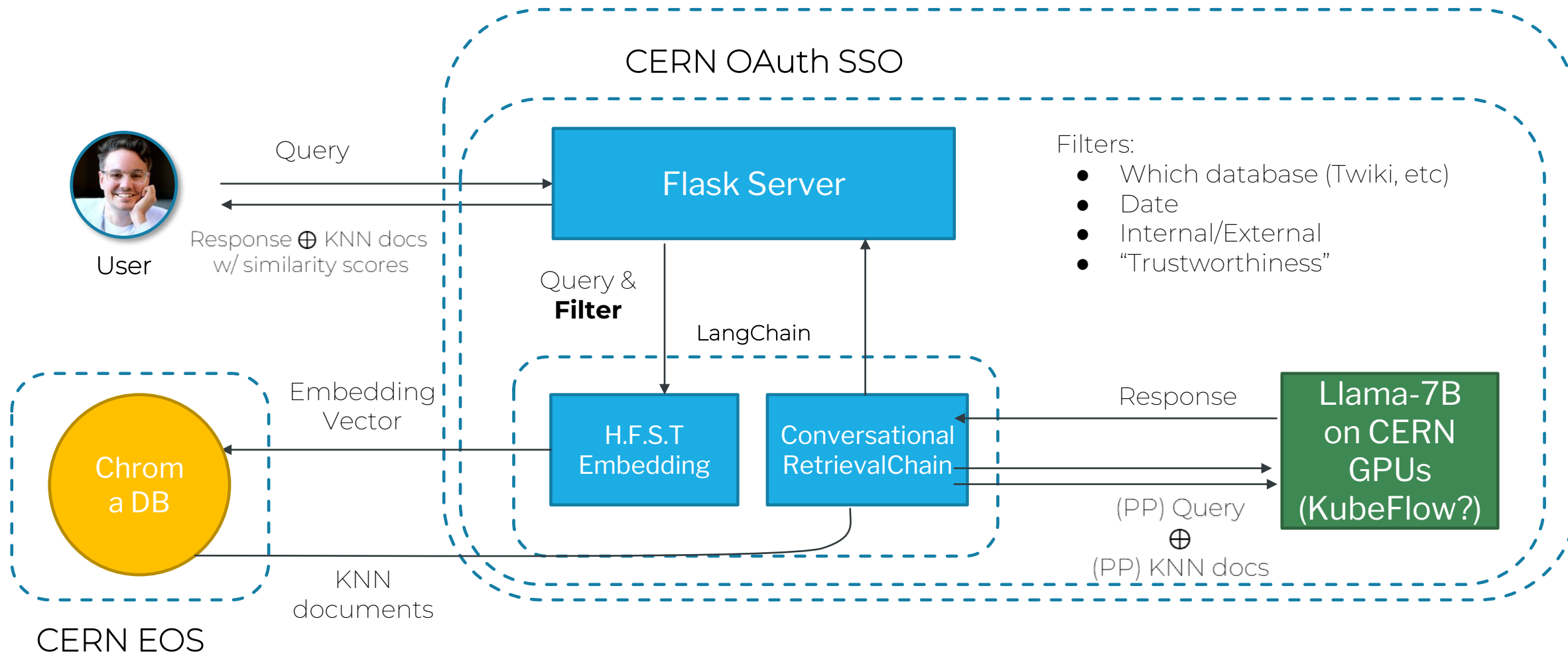
- Currently takes ~3 seconds to respond (mostly waiting on OpenAI)

Planned Infrastructure

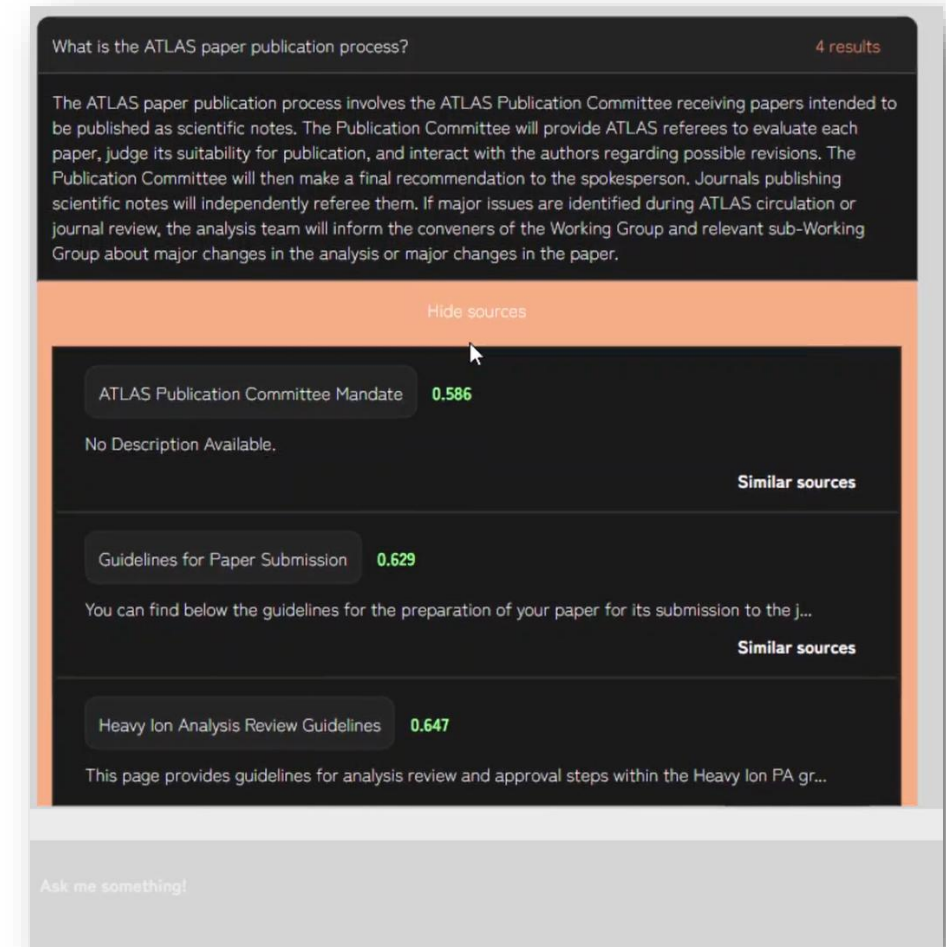
(PP) = (Possibly Processed)

CERN Openshift

CERN OAuth SSO



Chat Interface



Chat Interface

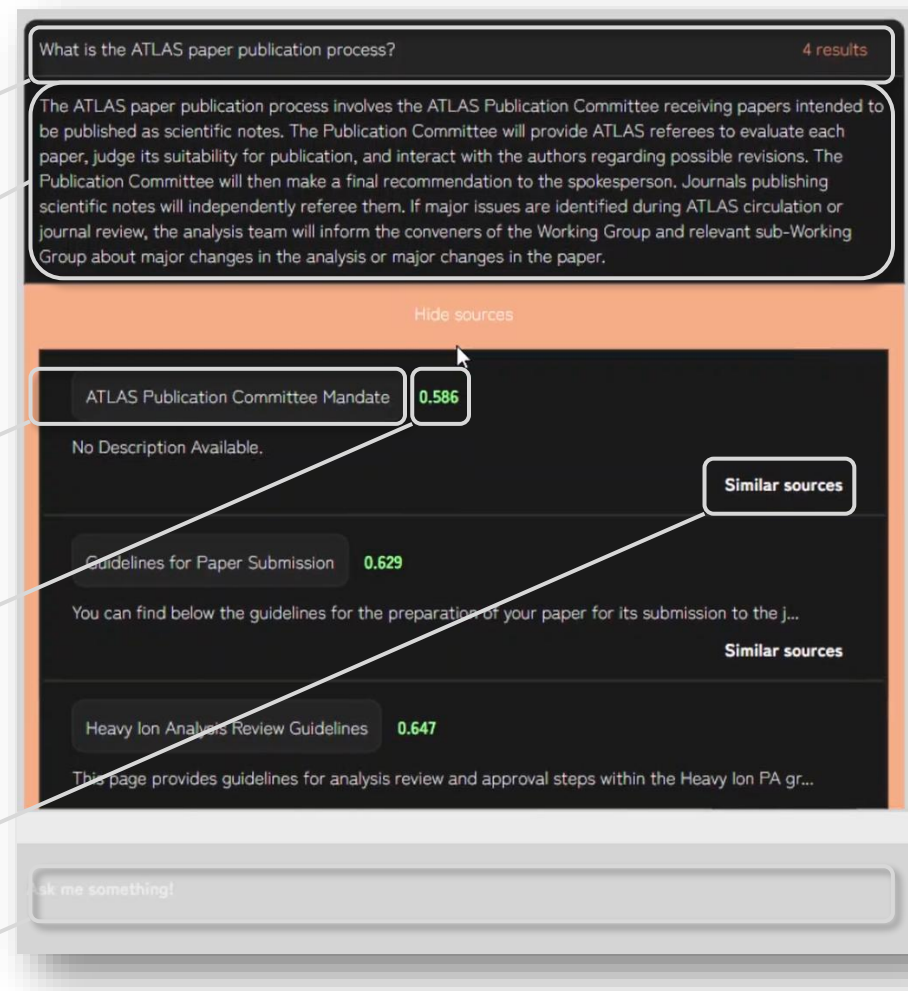
- Version 1.0 contains everything needed to answer a query
- Can cite the top sources used in the response
- Has a quick search for similar sources

Input
query
LLM
response

Source
title

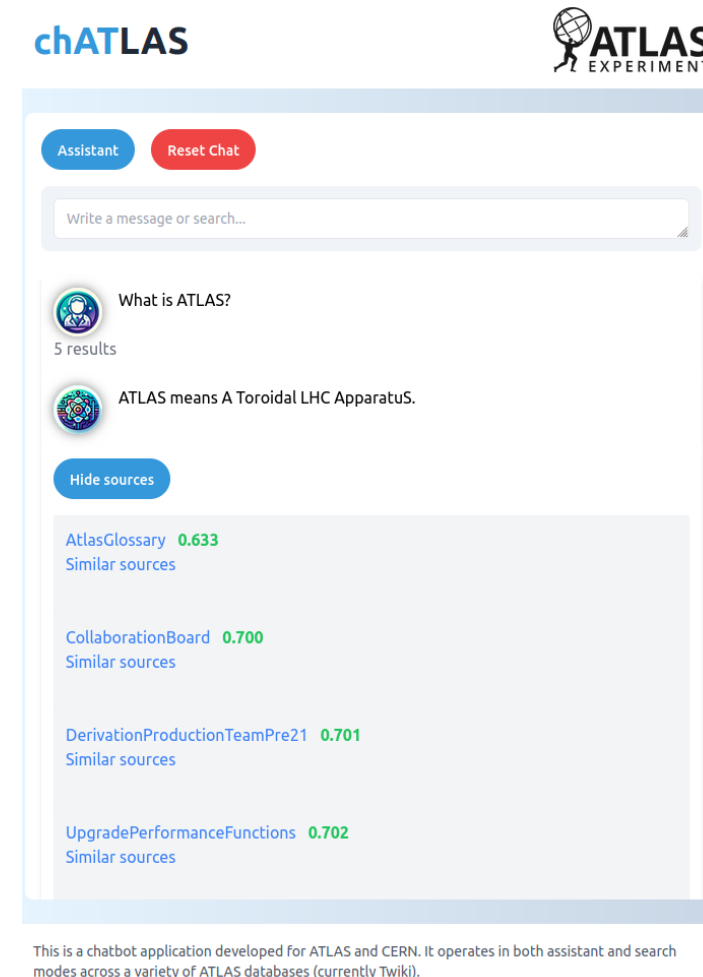
Similarity
score

Instant
similarity
search
Next query
entry



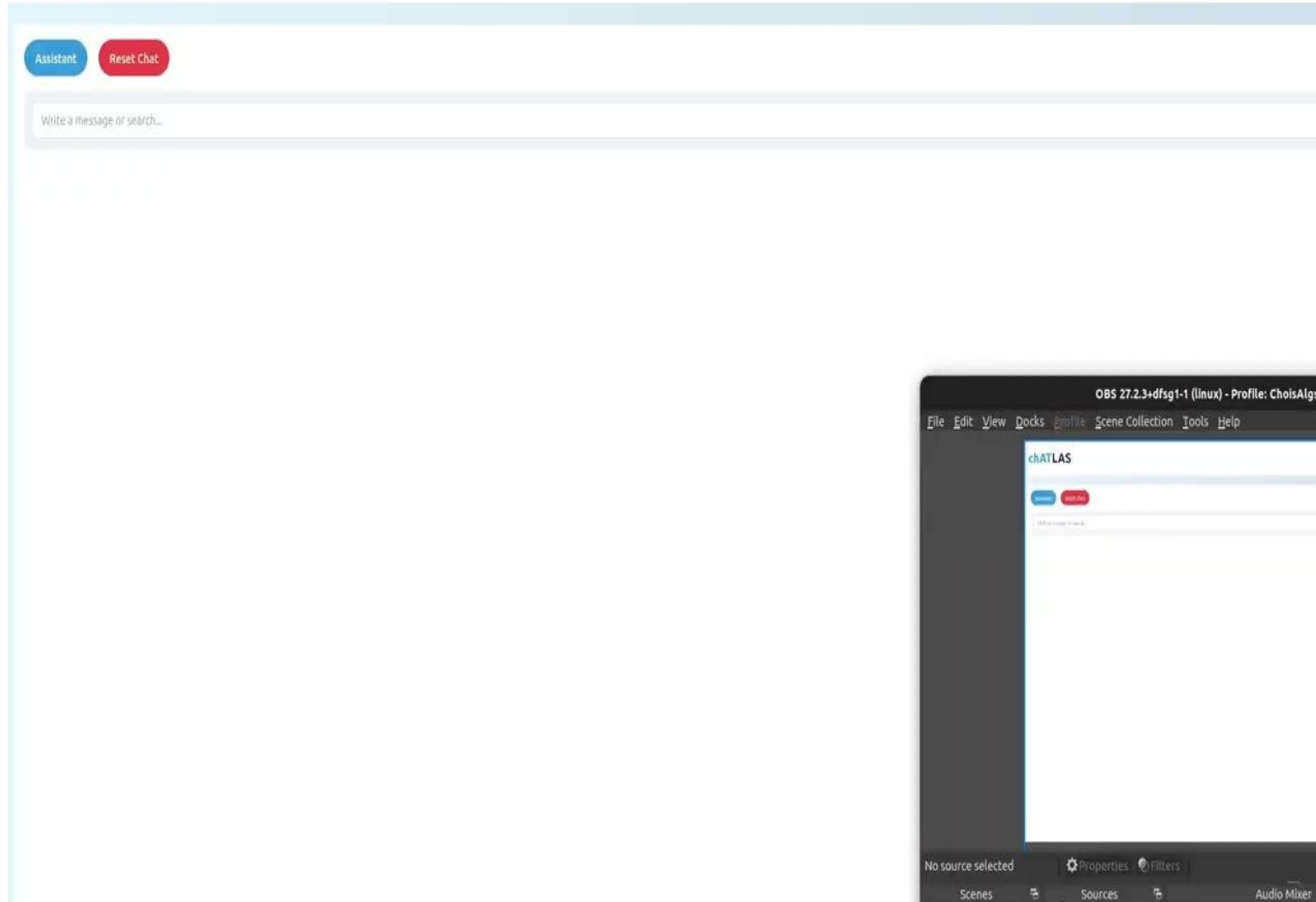
Chat Interface

- 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

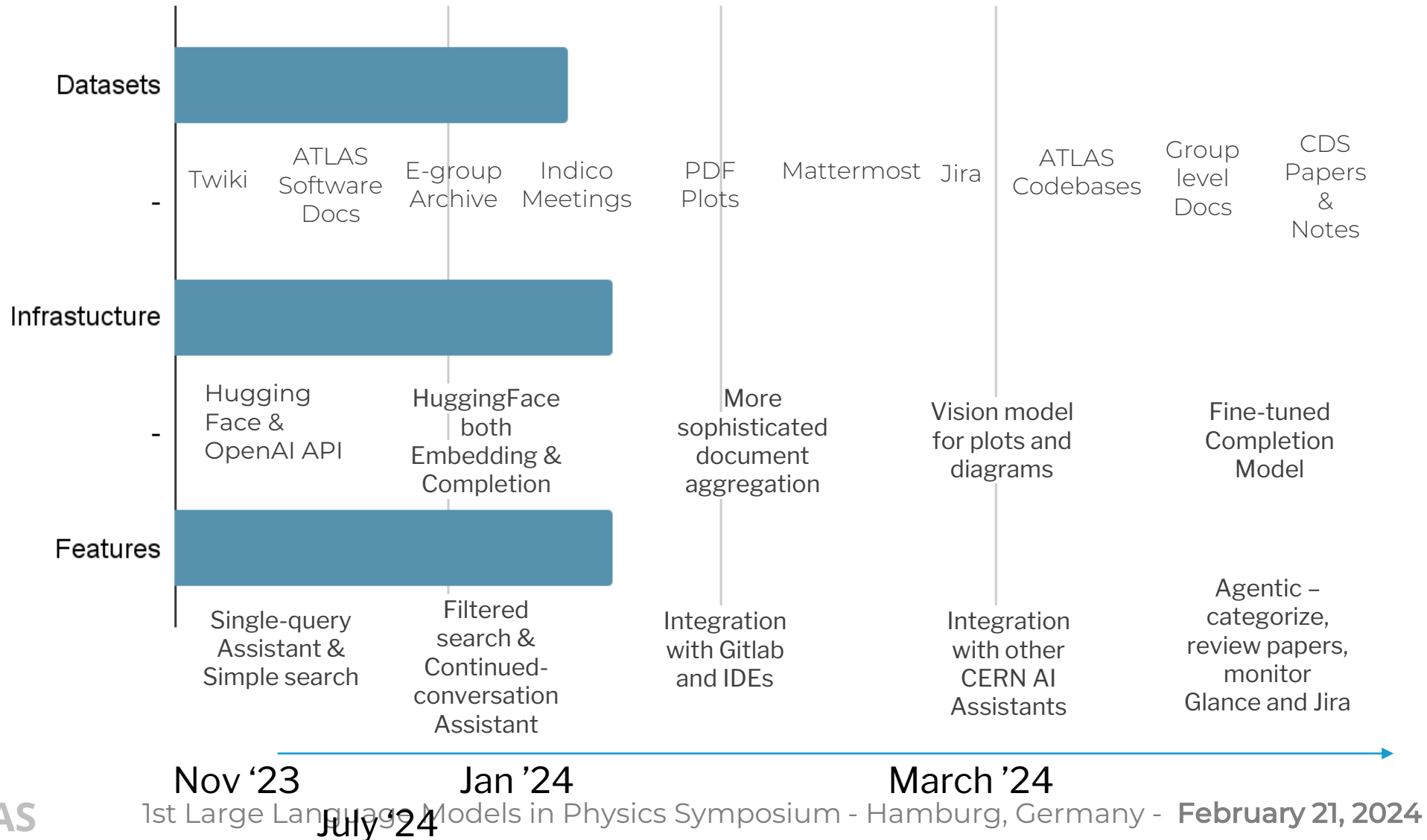


Demo

chATLAS



Roadmap (Milestone Tracking)



Recent Updates and Ongoing Work

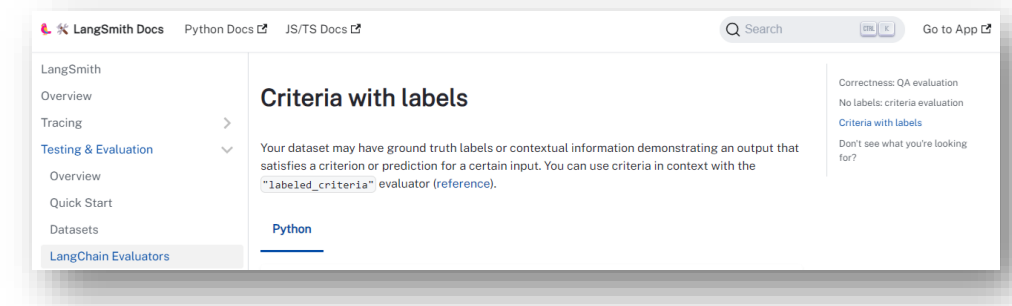
- Studies of the use and performance of different models as well as *response optimization* via *prompt engineering* and *similarity search techniques*
- Indico scraping shown previously with Nougat and Marker
- Video caption search and embedding
- Fine-tuning model research and experiments on either CDS or Twiki
- App development, deployment, documentation, etc
- Studies on how to include the ATLAS Glossary, and other “Dictionary” or “Reference” type documents - first-layer database search, fine-tuning, or even heuristic hand-engineering (“When is the NSW being upgraded?” - search in Glossary for any words present - place them in parentheses = “When is the NSW (New Small Wheel) being upgraded?”)

Challenges and Suggestions

- 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 ~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?

OPEN QUESTIONS

- How to avoid hallucinations? (Integrate latest research?)
- 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? (Metric Development)
- 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!

Presentation Summary

- Our goal is to create a reliable AI assistant across all ATLAS content
- Solved(In Progress) - ATLAS has significant data and presents a logistical challenge that we have largely overcome
- Complete - Implemented good semantic (vector embedded) search
- Complete - Prototype of AI Assistant
- Roadmap - Achieving our goal by Iterative Development
- Feedback - Recent Updates, Challenges, Suggestions, and Open Questions



END