# DEVELOPING A SCIENTIFIC DATA SEARCH ENGINE
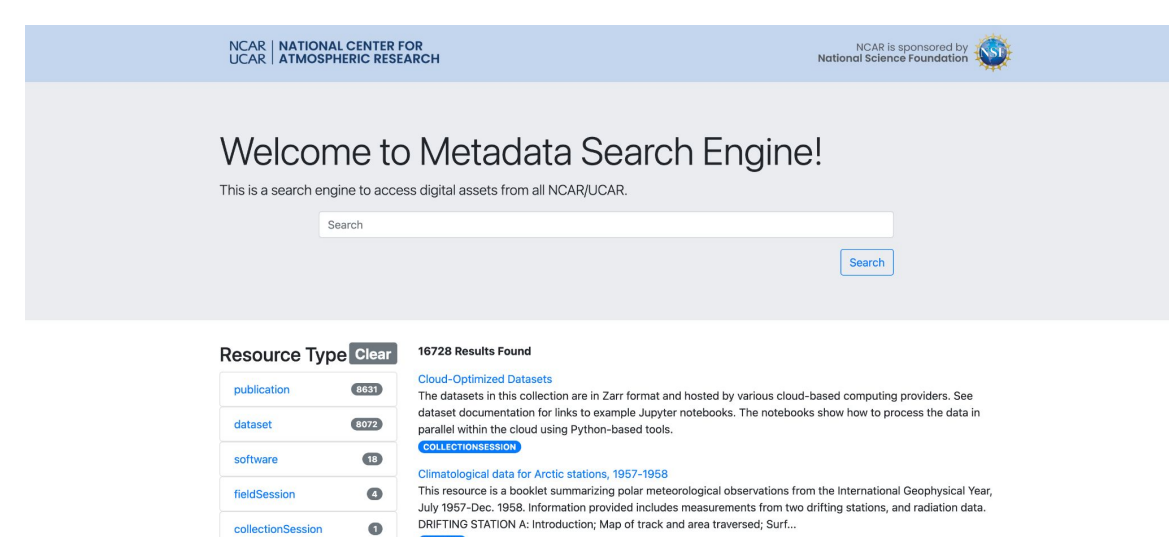
**Teagan Johnson, Sabira Duishebaeva**

**Mentors:** Nathan Hook, Saquib Aziz-Khan, Eric Nienhouse

## INTRO

### NCAR's Search Engine
NCAR has a variety of research labs that produce resources such as datasets, publications, and software. It's important that these resources are accessible for end users which is why NCAR has a search engine. This summer, **we worked on a new search engine** that has been in development for the past two years.

**Metadata**

In the context of the search engine, we refer to these resources as **metadata**.

### The "Problem"
Eventually, this new search engine may be deployed and be used as NCAR's primary search engine. Before it's deployed, there are many features and bugs that need to be addressed. Our job this summer was to **push the search engine closer to deployment.**

## OUR WORK

### Validated Metadata
Implemented a validation feature for the scientific metadata to **ensure search results are accurate.**

### Improved Deletion Efficiency
Decreased the amount of time it takes to delete a file by up to 100%, **further improving the search engine's accuracy.**

### Enabled Search Faceting
Designed a user-facing facet feature that enables filtering search results by various criteria, **improving the search engine's searchability.**
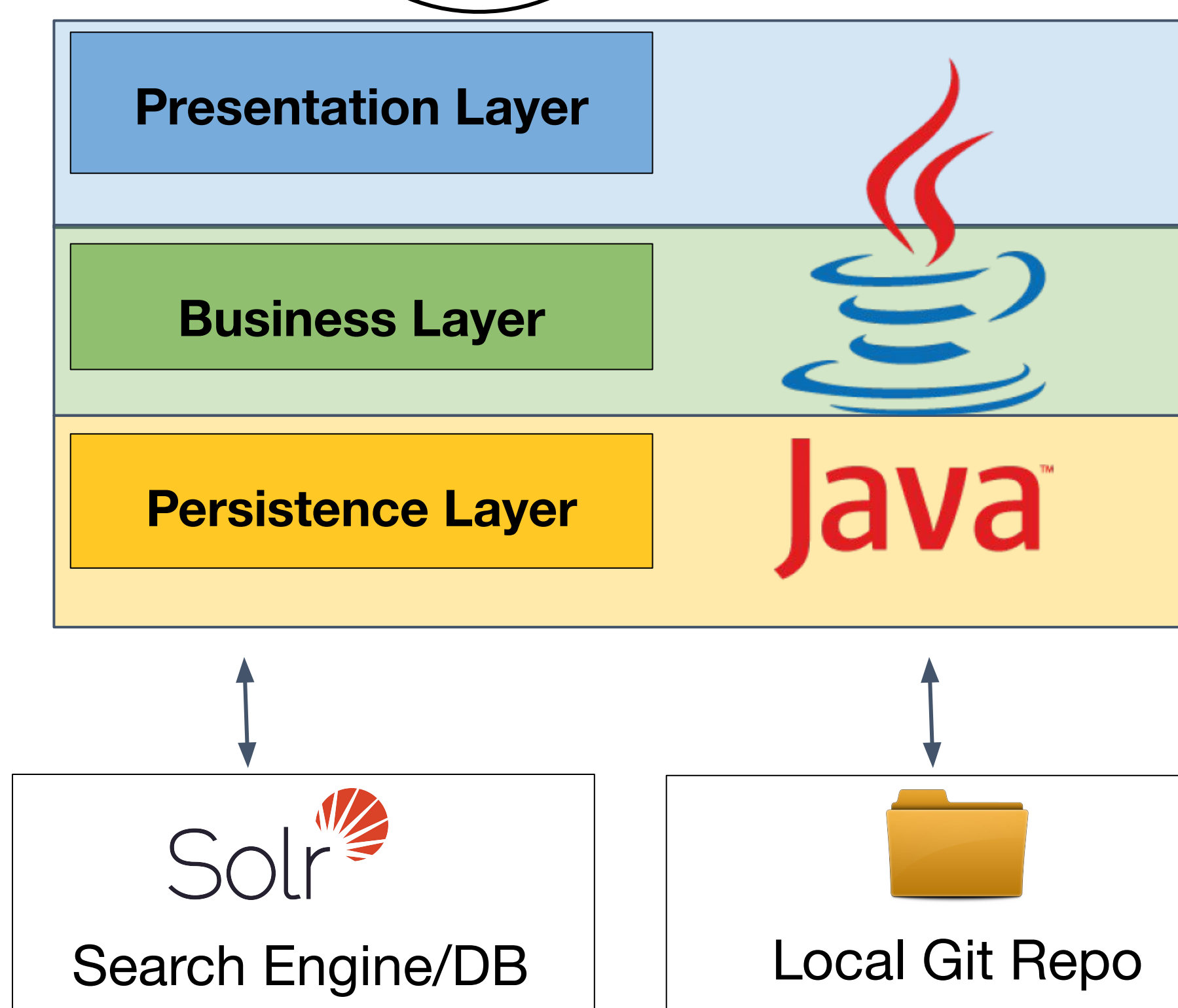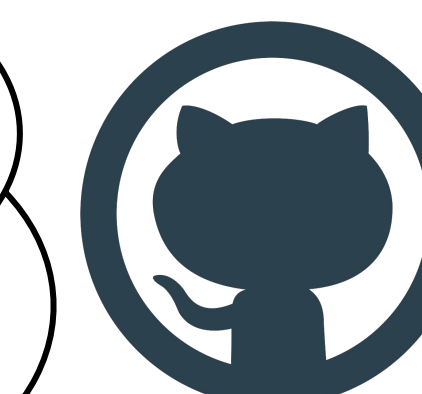
### Implemented Google Indexing
Provided a way for Google to efficiently crawl and index results in the search engine with a sitemap and JSON-Id, **revamping the search engine's findability.**

## HARVESTER ARCHITECTURE

**The Harvester**

The harvester's general purpose is to **retrieve resources from Github and put them into Solr.** It's three-layered design follows the layered architecture principle.
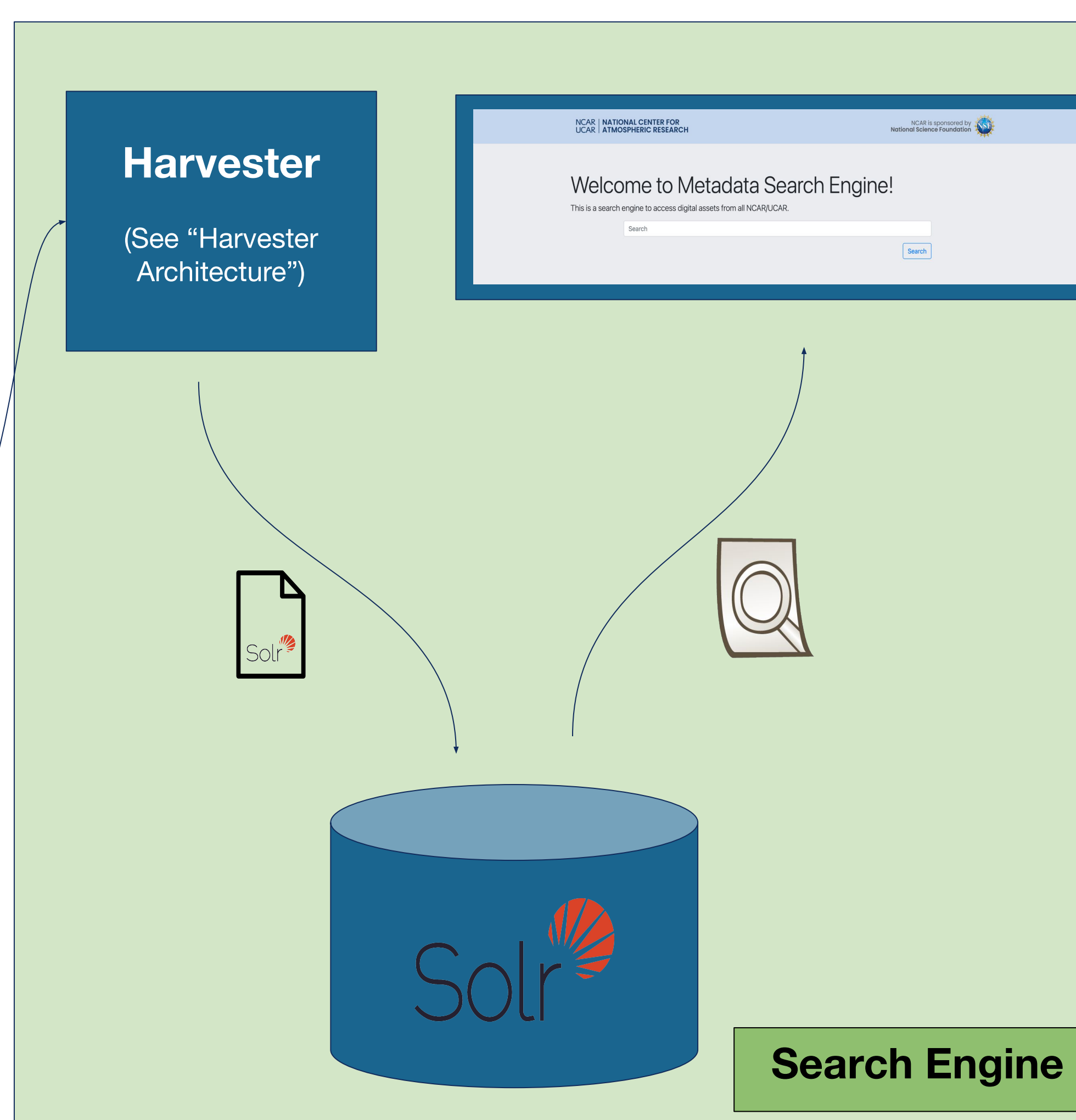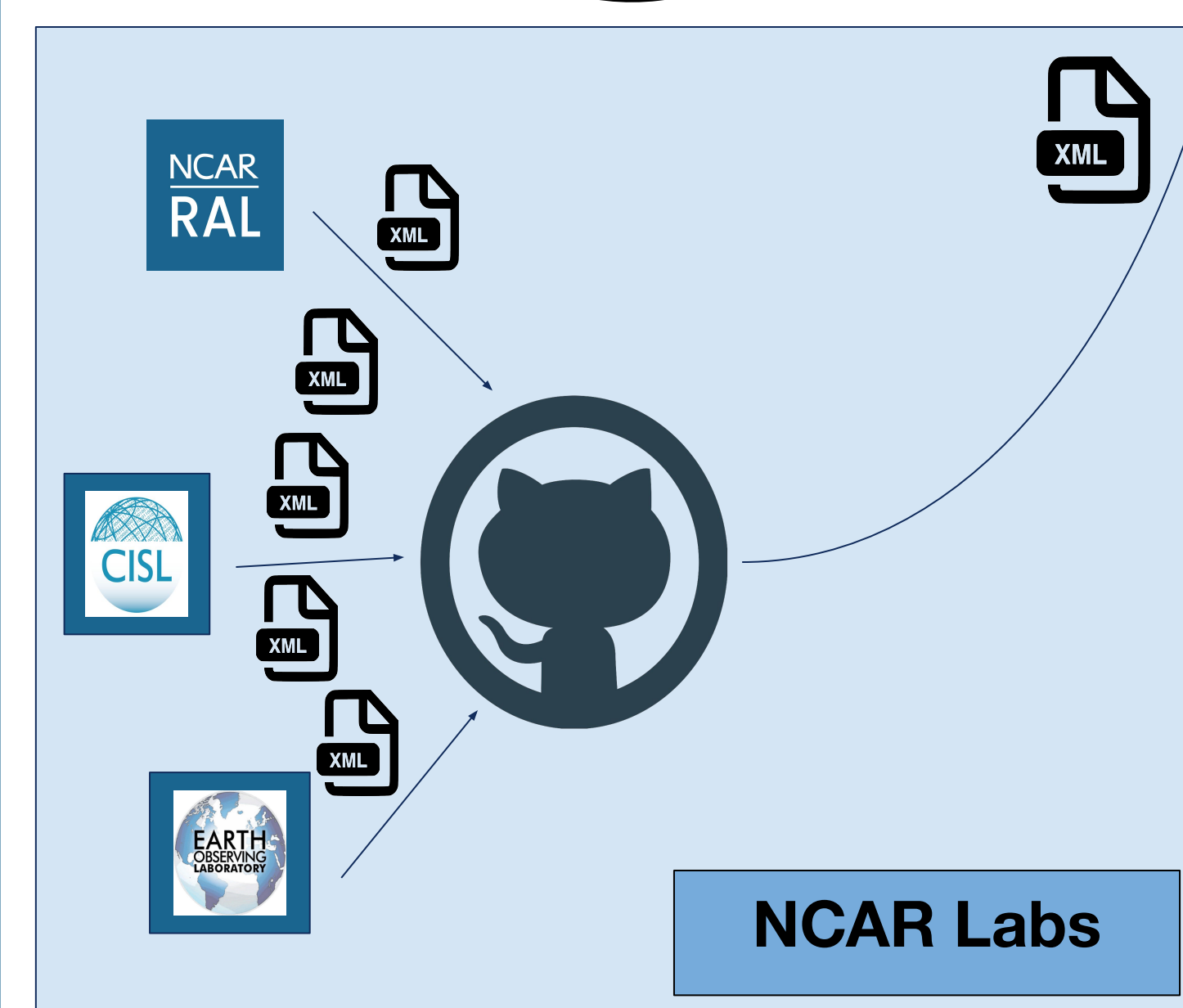
- Presentation Layer
- Business Layer
- Persistence Layer

Search Engine/DB | Local Git Repo

## SOFTWARE

spring boot | Solr | IJ | docker

## DATA FLOW

**XML Files**

Resources produced by the labs are **converted into XML files.** These flow from the labs, to GitHub, to the search engine.

**Harvester**

(See "Harvester Architecture")

Welcome to Metadata Search Engine!

NCAR Labs | Search Engine

## METHODOLOGIES

Scrum | S.O.L.I.D. | Layered Architecture

## CONCLUSION

### Conclusion
We ultimately achieved our goal of progressing the search engine towards deployment by **improving its searchability, findability, and accuracy.** More specifically, we implemented metadata validation, search faceting, efficient deletion, and Google indexing.

### Future Work
Future work includes adding extensions to validation and faceting, implementing autofill and spellcheck algorithms, designing a login feature for the harvester, and more.