# Developing a Scientific Data Search Engine

## Part 1: Architecture and Technologies

*TEAGAN JOHNSON*[1],
*Sabira Duishebaeva*[1], *Nathan Hook*[2],
*Saquib Aziz-Khan*[2], *Eric Nienhouse*[2]

*SIParCS Intern [1] - NCAR Mentor [2]*

July 27, 2022

# Table of Contents

- Background

- Goals

- Methodologies + Technologies

- Architecture

- Result 1: Efficient Deletion
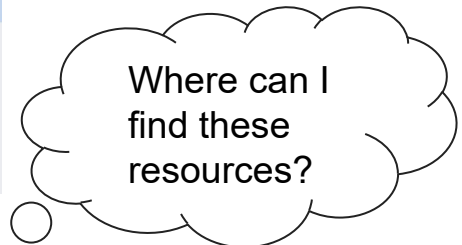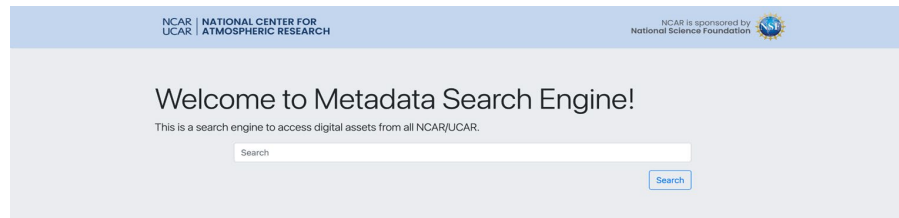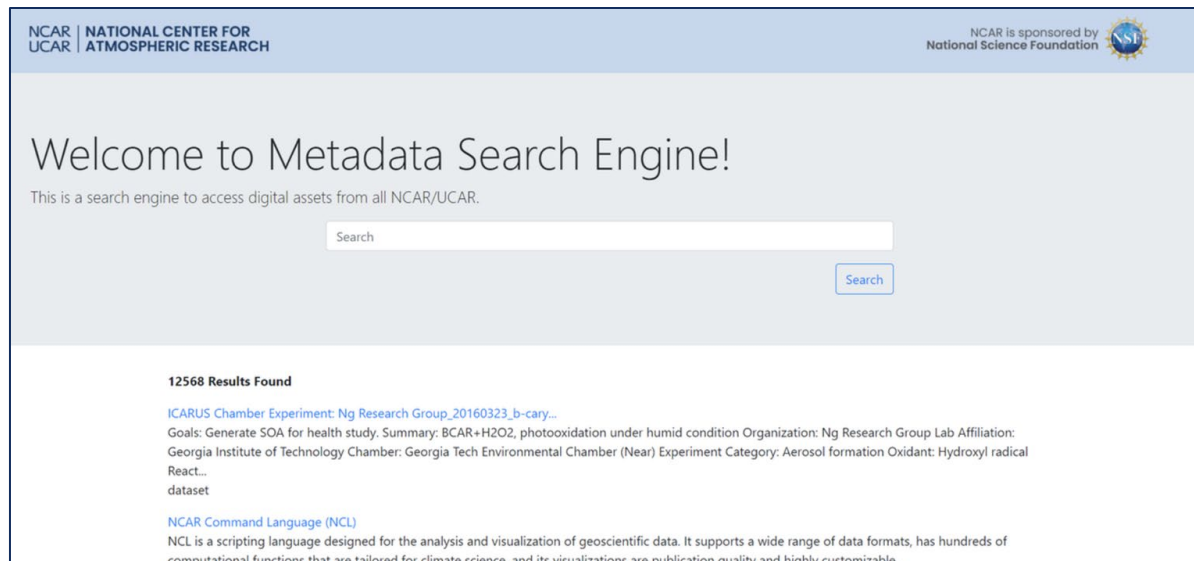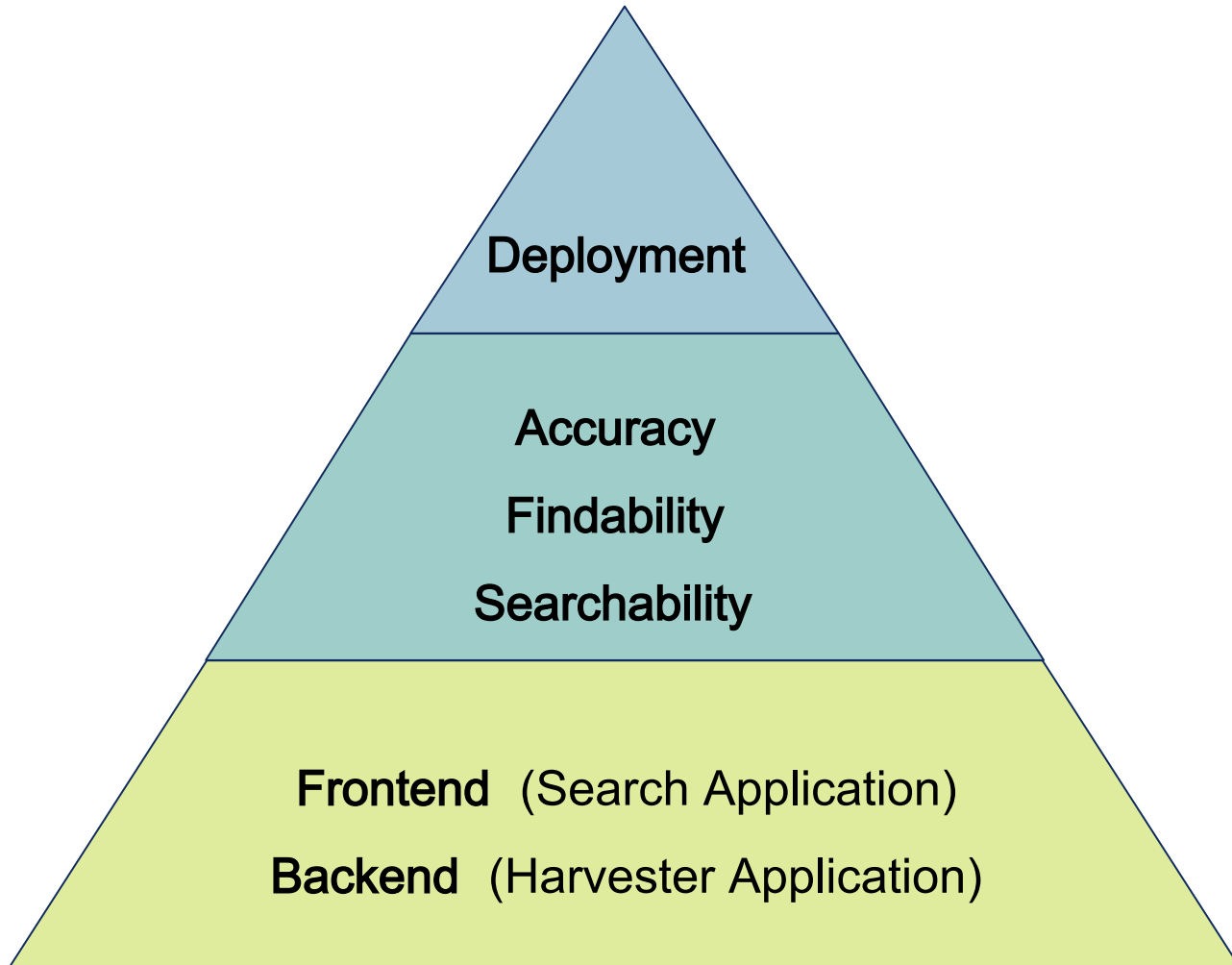
- Acknowledgements

# BACKGROUND

# Initial State

- Java-based search engine that contained a **harvester** and a **search** application
- Numerous features and bugs needed to be addressed **before deployment**

# GOALS

# Goals



Deployment

Accuracy

Findability

Searchability

**Frontend** (Search Application)

**Backend** (Harvester Application)

# METHODOLOGIES

# Agile Scrum

- **Short, iterative cycles** (1 week)

- **Reprioritization** (retro meetings)

- **Communication**

# Design Principles

- SOLID Principles

- Layered Architecture

# TECHNOLOGIES

# ARCHITECTURE

# Resource Representation

Every labs' resources are converted into XML files because XML files are uniform

Properties:
- **Title**
- **ID**
- **Author(s)**
- **Resource Type** (dataset, software, publication, etc.)
- **Lab**
- **Etc.**

Harvester

Welcome to Metadata Search Engine!
This is a search engine to access digital assets from all NCAR/UCAR.

# RESULT 1: DELETION

# The Problem

- If a lab deletes a file, it should be deleted from the search engine **(increase accuracy)**



- Before, to delete a file we'd clear the harvester and re -pull *every* file from Github which was **very inefficient**

Harvester

Welcome to Metadata Search Engine!

This is a search engine to access digital assets from all NCAR/UCAR.

Solr

# Efficient Deletion

- To delete one file it took **11 mins, 36 secs**

- Now it takes LESS THAN A SECOND

# 99%

# ACKNOWLEDGEMENTS

NCAR |
UCAR |

# Thank You

NCAR and CISL

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

**Project Partner:**
- Sabira Duishebaeva

**SIParCS Program Leads:**
- Virginia Do
- Jerry Cyccone
- AJ Lauer
- Francesgladys Pulido

And everyone else involved with the SIParCS program this summer.

# Resources

- https://www.hiclipart.com/ - Most images in presentation
- https://blog.knoldus.com/why_-we-need-solid-principles-and-its-types/ - SOLID Principles
- https://www.youtube.com/c/ncarcgd - CGD logo
- https://mobile.twitter.com/ncar_acom - ACOM, EOL logos
- https://mobile.twitter.com/ncar_cisl CISL logo
- https://www.youtube.com/channel/UCgDMGLn6JKEU87aJJlt8_F-g - HAO logo
- https://mobile.twitter.com/ncar_mmm - MMM logo
- https://www.youtube.com/channel/UCDIFPOu2f7TortTgNYQc6wA - RAL logo
- https://www.scrum.org/resources/scrum_-framework-poster - Scrum Diagram