# Improving Electronic Logbook Searches Using Natural Language Processing

Jennefer Maldonado

Collider Accelerator Department, Brookhaven National Laboratory

jmaldonad@bnl.gov

@BrookhavenLab

# Introduction

# What is the elog system?

- The electronic logbook (elog) system is used to record information ranging from meeting notes, to do lists, and critical operations

# Motivation

- The search feature only provides exactly what a user enters, what if there are other entries that do not include those exact words BUT also are related to these things

- Allow for more catering to search filters like date, logbook, etc.

- Eventually provide custom sets of entries based on users' interactions with the system

- Possibly eliminate the need for manual searching

Brookhaven
National Laboratory

# Data

| | ID | Content | Timestamp | Author | ElogID | Tag | Flag |
|---|---|---|---|---|---|---|---|
| **0** | 1 | \<p>bta-th158-ps and bta-qd5-ps both have a sta... | 2013-11-18 20:25:48 | pdyer | 1 | bta | 0 |
| **2** | 3 | NRO wants the same 114 MeV (160 in Booster) se... | 2013-11-18 20:06:38 | NAK | 1 | | 0 |
| **3** | 4 | New 114 MeV Au_Ebis file created. | 2013-11-18 20:00:04 | tape | 1 | | 0 |
| **4** | 5 | It starts out fine thhen fades away | 2013-11-18 18:00:25 | keith | 1 | | 0 |
| **5** | 6 | Entry deleted | 2013-11-18 17:56:49 | anonymous | 1 | | 0 |

- The database includes whether entries are a comment, what book they are in, and time entered
- All elog data is stored in a MySQL database

**Brookhaven** National Laboratory

# Data Processing

- Remove links, numbers, tokenize, lemmatize, lowercase, remove punctuation, also remove any entries with no content in them

```
Service, building, and equipment tour complete.
```

```
[service, building, equipment, tour, complete]
```
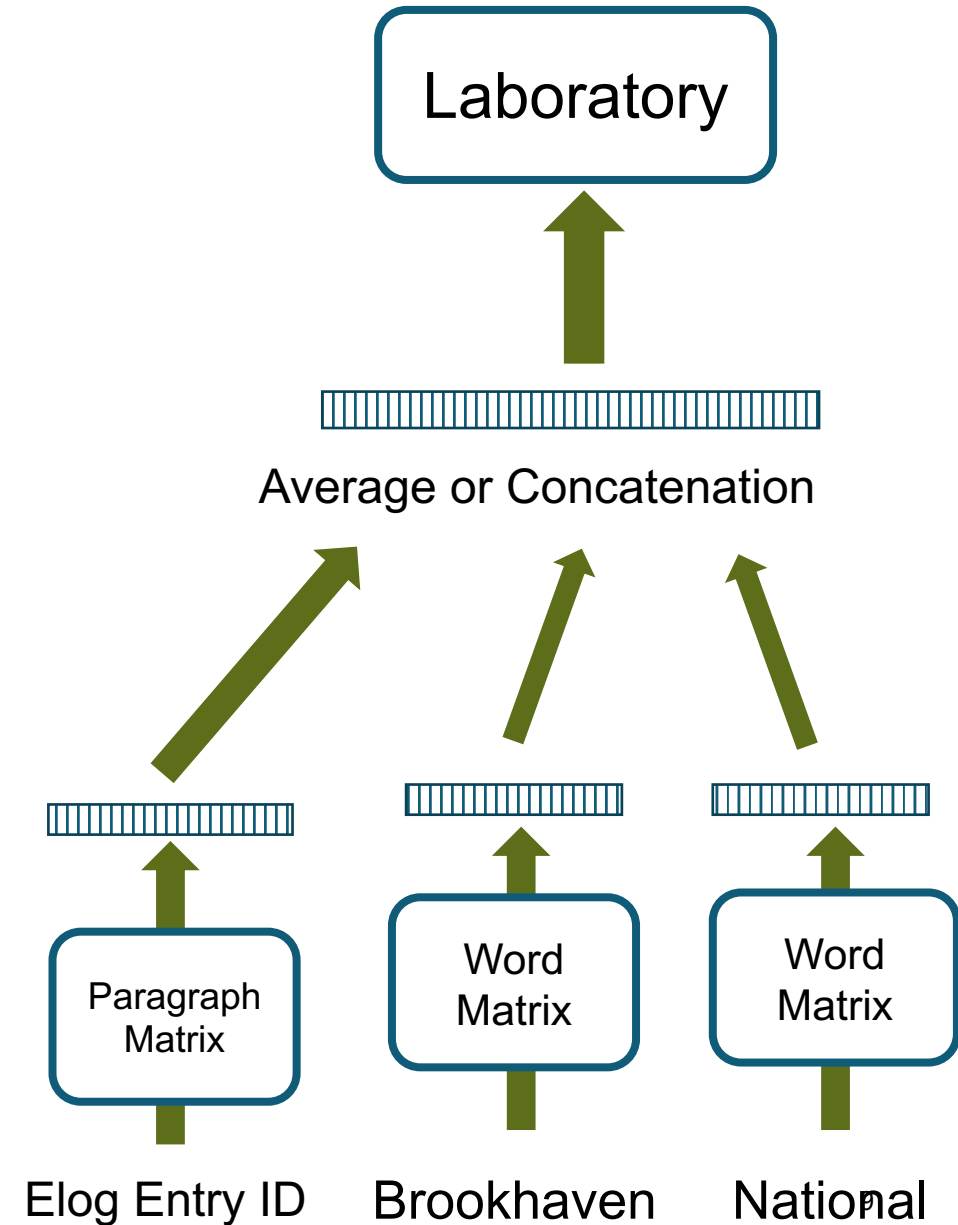
# Models

# Doc2Vec

- Gensim D2V model based on *Distributed Representations of Sentences and Documents* by Le and Mikolov

- Paragraph vectors predict the next word given a sample of words from the paragraph

- Every paragraph is mapped to unique vector which is a column in the paragraph matrix

- Every word in each unique vector also gets mapped to a unique column vector in the word matrix

- Take processed entries and create a list of d2v tagged documents



Laboratory

Average or Concatenation

Paragraph Matrix

Word Matrix

Word Matrix

Elog Entry ID    Brookhaven    National

**Brookhaven**
National Laboratory

# Similar Documents

`polarization for yellow 2h target1 store energy before physics declared yellow beam intensity`

1. Yellow 1 V6 Polarization: -51.53 6.05%
   Yellow 2 H6 Polarization: -51.28 10.83%    **78%**

2. Polarization For Yellow 1 V Target2: 51.44 &plus mn 1.94 Store Energy (254.21) Before Physics Declared, Yellow Beam Intensity: 208.3x10^11    **77%**

3. Yellow 1 V5 Polarization: -56.67  4.87% Yellow 2 H5 Polarization: -59.46  6.09%    **76%**

# Topic Modeling



t-SNE Clustering of 8 LSA Topics

- time monitor loss
- nbsp position blue
- beam energy ebis
- rhic yellow start
- fault vacuum analysis
- ramp blue yellow
- current temperature user
- jet ipm efficiency

**Latent Semantic Analysis**

LSA uses dimension reduction techniques to find meanings and similarities of documents by how frequently words appears in those documents.

**Latent Dirichlet Analysis**

LDA utilizes vector representations of the ratio of the counts of words in document data.

t-SNE Clustering of 8 LDA Topics

- pattern beam turn
- mev energy beam
- monitor time loss
- nbsp beam position
- physics dump prepare
- rhic beam blue
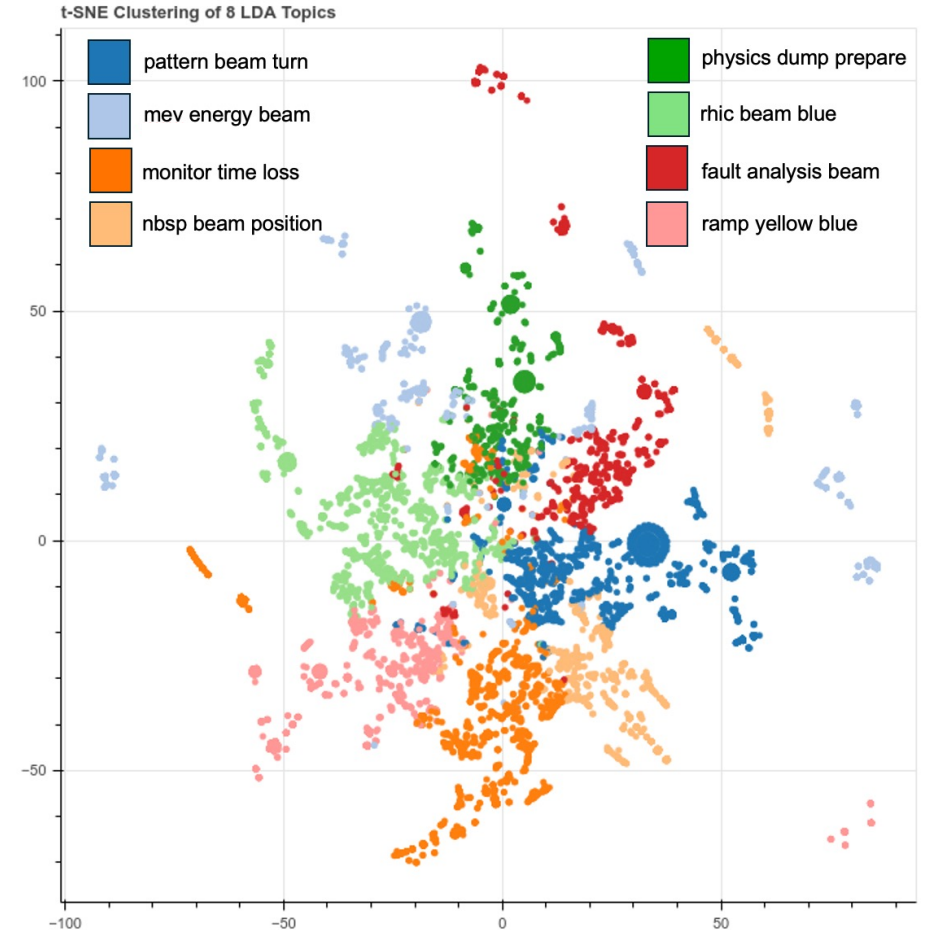- fault analysis beam
- ramp yellow blue
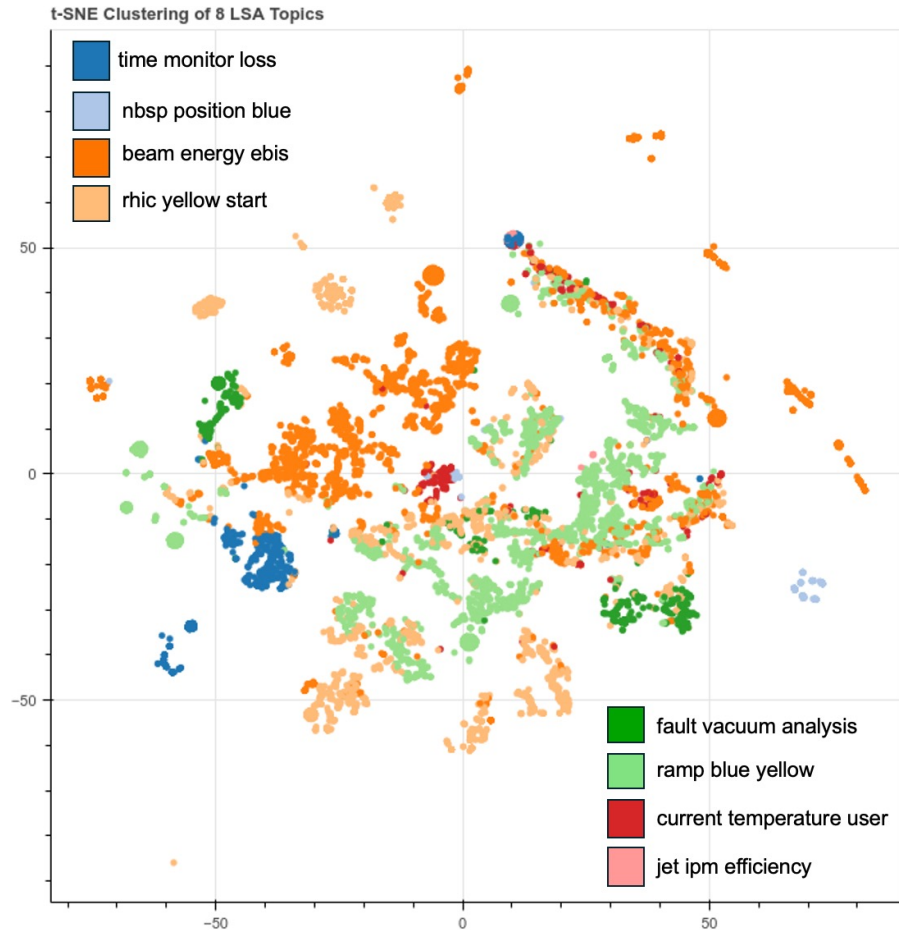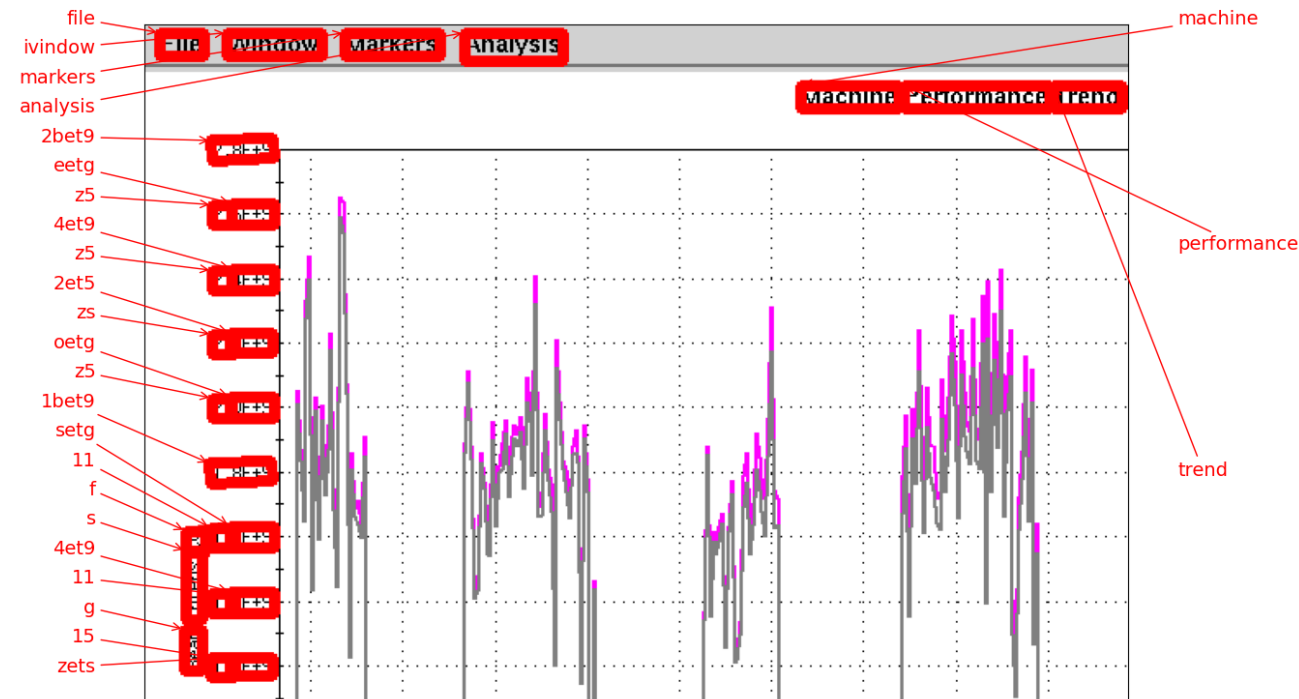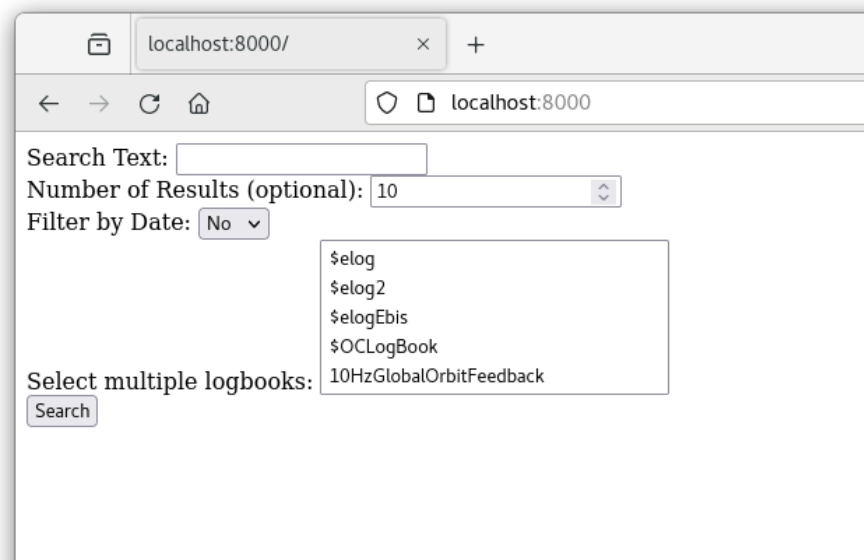
**Brookhaven**
National Laboratory

# Image Processing

Optical character recognition (OCR) is utilized to parse text out of images to be used for natural language processing tasks or many other applications. Keras OCR was used to test text recognition on images attached to elog entries.

# User Interface

# User Interface

- Wanted a web demo for users to test and give feedback

- Created a simple website using FastApi and Jinja templates

# User Interface

- User's can search by elog and date
- Filter number of responses
- This page will be linked to the elog for initial testing
- The current elog system will be evaluated for EIC to see if this functionality will be added to the elog or another tool
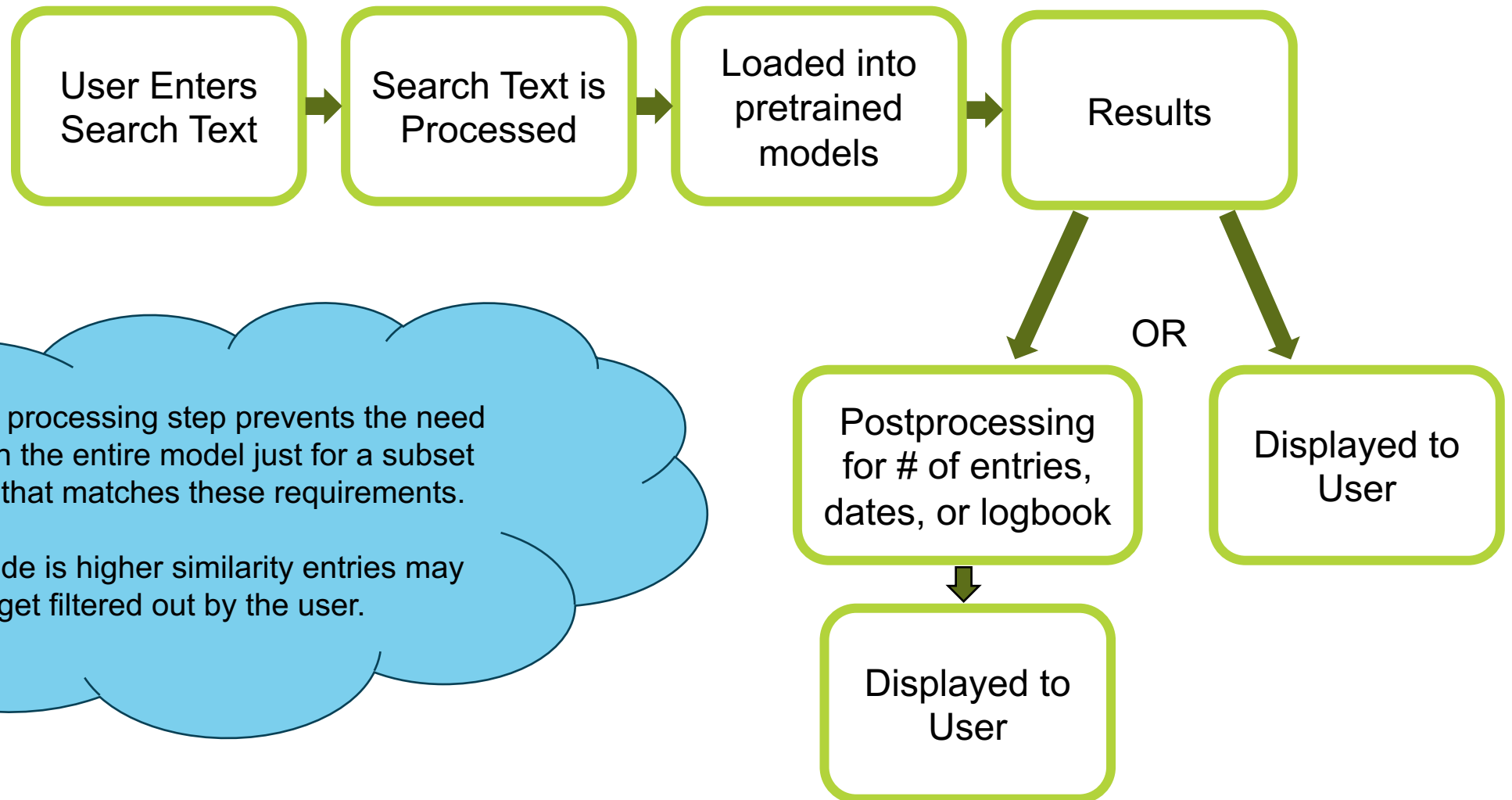


Search Text: blue and yellow
Number of Results (optional): 7
Filter by Date: Yes
Start Date: 01/01/2016
End Date: 02/06/2024

$elog
$elog2
$elogEbis
$OCLogBook
10HzGlobalOrbitFeedback
Select multiple logbooks:
Search

# Web Interface Workflow

User Enters Search Text → Search Text is Processed → Loaded into pretrained models → Results

Results → OR → Postprocessing for # of entries, dates, or logbook OR Displayed to User

Postprocessing for # of entries, dates, or logbook → Displayed to User

The post processing step prevents the need to retrain the entire model just for a subset of data that matches these requirements.

Downside is higher similarity entries may get filtered out by the user.

**Brookhaven** National Laboratory

18

# What's Next?

- Link this demo webpage into the elog system to allow users to test thoroughly
  - Feedback drives new features
- Then decide how to implement directly into the elog (or another tool)
- Student to development reinforcement learning model
  - Improve searches with user interaction!

Brookhaven
National Laboratory

# Thank you!

Thanks to Sam Clark, Wenge Fu, and Seth Nemesure for their knowledge and support.

Brookhaven
National Laboratory