



SmartSearch, scalable data search and aggregation in the Cloud for CCS and beyond

Vic Baker, National Energy Technology Laboratory
Kelly Rose, National Energy Technology Laboratory
Paige Morkner, National Energy Technology Laboratory

```
49 set_val
50 plt.
51 4]
4]
size = Run Cell | Run Above
52 #%%# Split-out va
53 array = dataset.values
54 X = array[:,0:4]
55 Y = array[:,4]
56 validation_size = 0.20
57 seed = 7
58 X_train, X_validation, Y_train, Y_validation
```



U.S. DEPARTMENT OF
ENERGY



Multi-Cloud (...really means Multi-Environment)

- “Cloud” is where you’re **not** at
 - Cloud Service Providers
 - On-Prem Compute
 - Edge computing (sensors, cell phone, etc.)
- Leverage APIs for service-to-service communications
- Networking:
 - peered networks
 - commodity networks
 - dedicated/shared interconnects
 - consider egress fees (leverage compute where data lives)
- Utilize cloud-native concepts



Multi-Environment

“Cloud” is compute resources where you’re *not*...



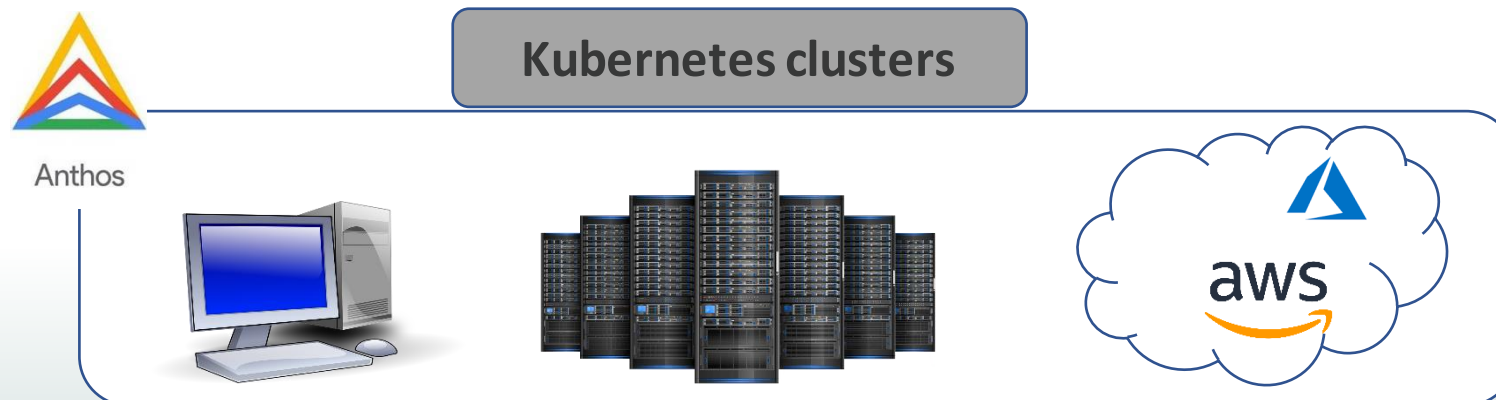
Navigating on-prem to cloud (and local dev)

- One of the biggest challenges: migrating from an “on-prem” development / deployment environment to cloud – HOW???
- Utilize “cloud native” concepts
 - Containerization
 - Microservices
 - Scalability
- Leverage enabling tools:
 - Deployment:
 - Kubernetes
 - Helm
 - Terraform
 - Development:
 - VS Code
 - Containerized environments



Kubernetes

- Supported by all major cloud providers
- Configurable autoscaling based on load thresholds (CPU, RAM usage)
 - Services
 - Nodes
- Managed service deployments via Helm
- Seamless local, on-prem, cloud deployments via kubectl
- Manage multi-environment clusters via GCP Anthos
- Kubeflow



Everyone wants to jump to the **top** of the pyramid...



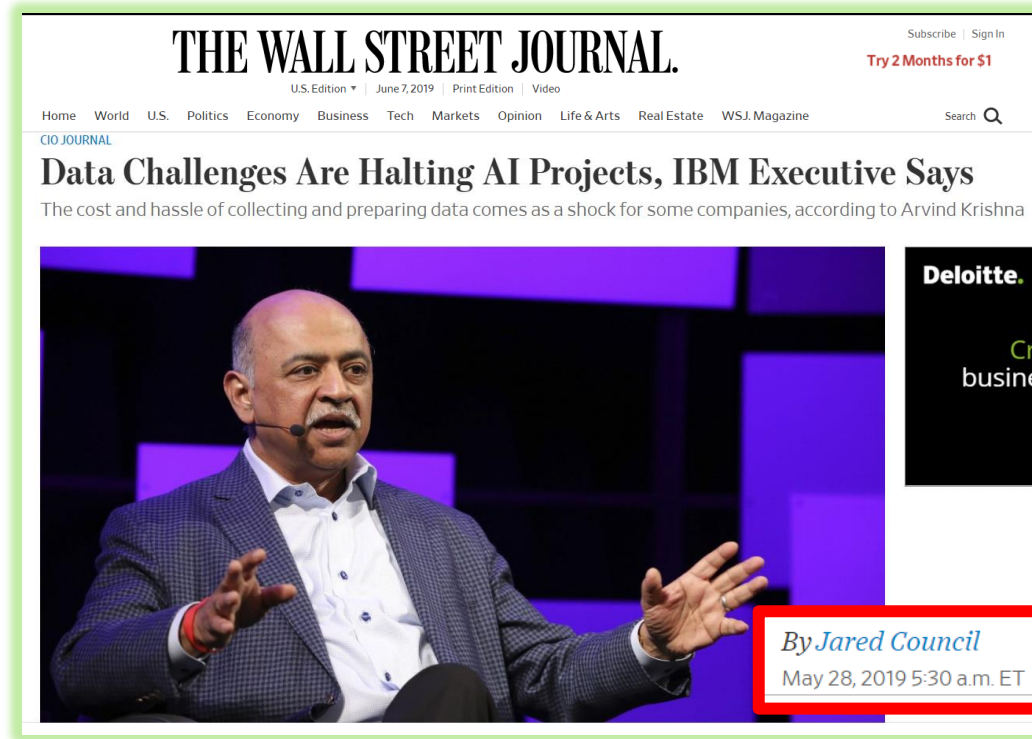
But... scientists still spend nearly 80% of their time **acquiring, cleaning, and organizing** data

The steps at the **bottom** must come first

They **require leadership & support** to ensure a solid data foundation for future R&D

“**Invest 5% of research funds in ensuring data are reusable.** Funders hold the stick: they should disburse no further funding without a data stewardship plan.”

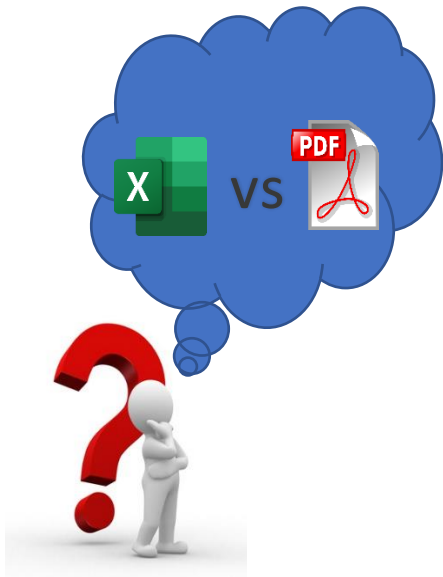
- *Nature*, February 2020



Data are the **energy** for AI and analysis

SmartSearch[®] Conquering the Data Avalanche

- **How do you currently search?**
 - Type in a few keywords
 - Skim the top few results
 - Type in more keywords and try again



- **How do you find and connect to something relevant?**
 - Open a file / web page
 - Read it (skim it)
 - Decide if it's relevant

APPROVED FOR PUBLIC RELEASE

What is SmartSearch?

Problem:

You like these files.

You want to find more data relevant to the content of these files



Solution:

SmartSearch automates data discovery by ...



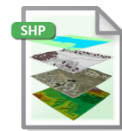
1) **Analyzing** content you like



2) **Finding** new content via www, local, enterprise data stores



3) Telling you **how relevant** the discovered data is to what you like



Input



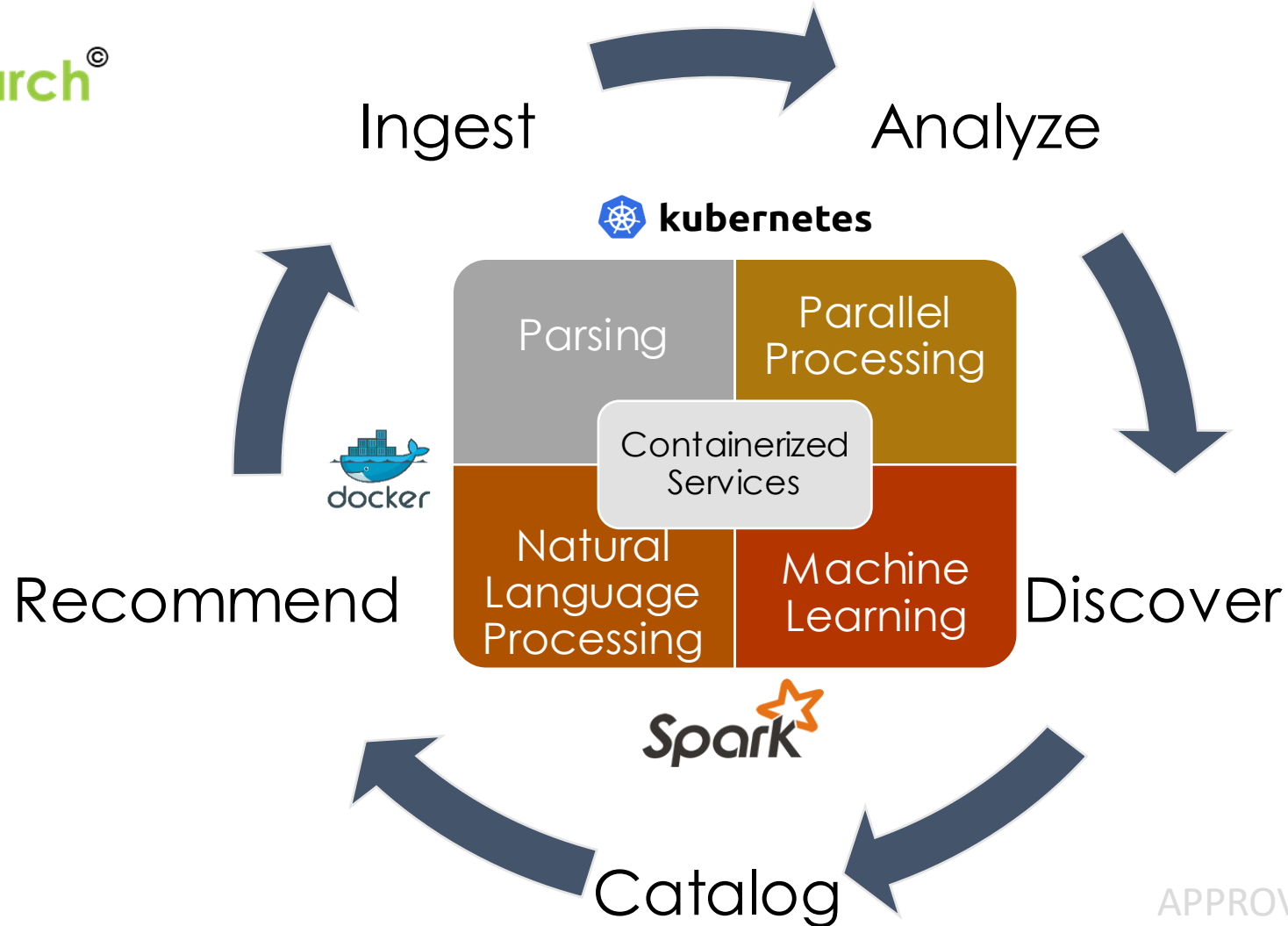
72%



Discovered

APPROVED FOR PUBLIC RELEASE

How Does SmartSearch[®] Work?



APPROVED FOR PUBLIC RELEASE

Benefits of SmartSearch[©]

- **Infinitely Scalable (Automated) Data Discovery**

- Analyze millions+ of files and generate comparison metrics
- Generate topic models, categorization, recommendations
- Desktop, cluster, cloud

- **Treat geospatial data like a document**

- Automatically extract text from geospatial data (shapefiles, geodatabases)
- Compare textual vs geospatial data to identify relevancy

- **Search for meta tags within HTML body of discovered web sites**

- i.e., find map tags

- **Analyze archive files – even archives within archives (zips within zips, etc.)**

- Process every file – docs, spatial, etc.



AI/ML in SmartSearch[®]



- **SmartSearch built via cloud native design principles**
- **Combines 'cluster of clusters', Spark, Kubernetes, and data lakes for massively scalable compute infrastructure**
- **Natural Language Processing via SparkNLP**
 - Distributed NLP processing via the Spark framework
 - Implemented via SparkML Pipelines
 - Provides thousands of pretrained models and pipelines (Glove, Bert, Onto, etc.)
 - Custom models can be implemented and trained within same distributed framework
- **Machine Learning via SparkML Pipelines**
 - Distributed ML processing via the Spark framework
 - SmartSearch Recommendation Engine
 - LDA Topic Modeling
 - Named Entity Recognition (NER)
 - Question Answering
 - Summarization

APPROVED FOR PUBLIC RELEASE



U.S. DEPARTMENT OF

ENERGY

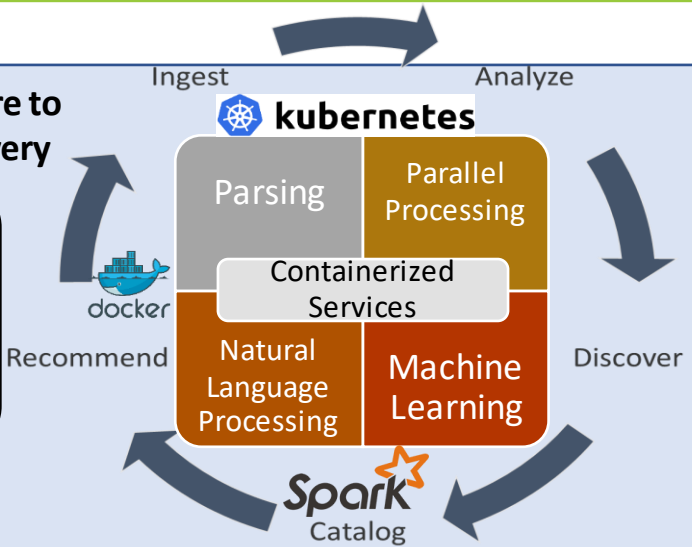
MATRIC

<https://edx.netl.doe.gov/sami>

AI informed approach

Challenge: data infrastructure to AI/ML enhanced data discovery

Employing AI/ML tools to find open resources



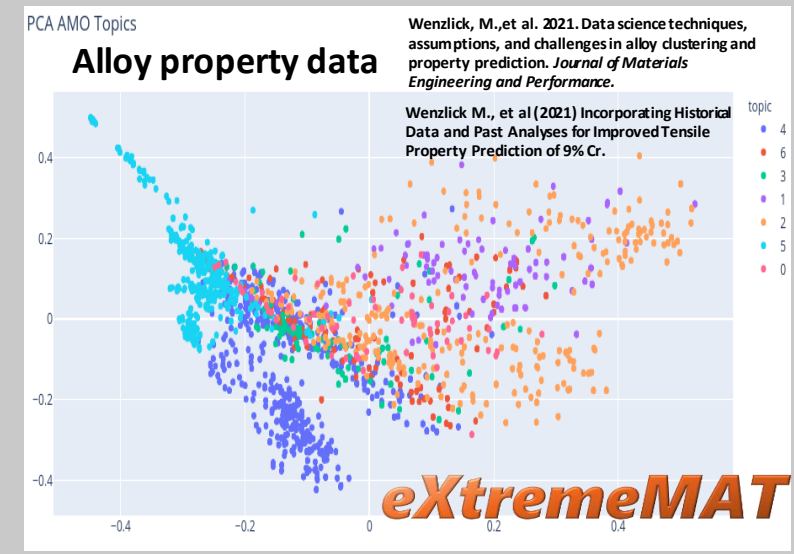
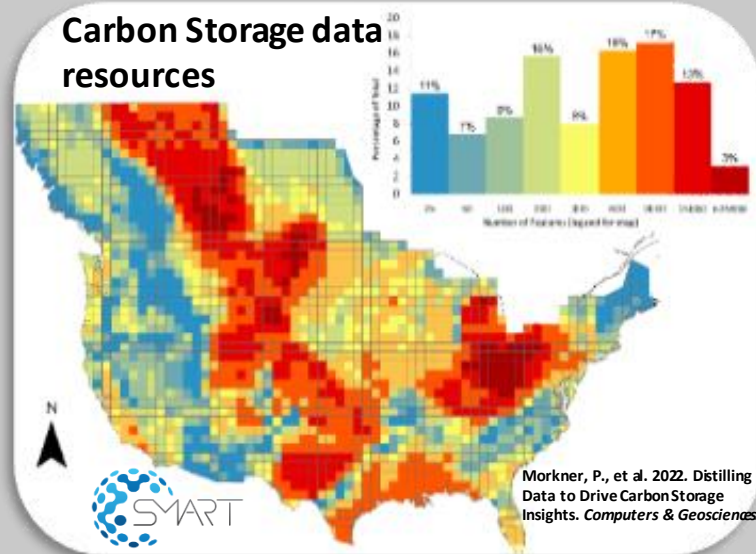
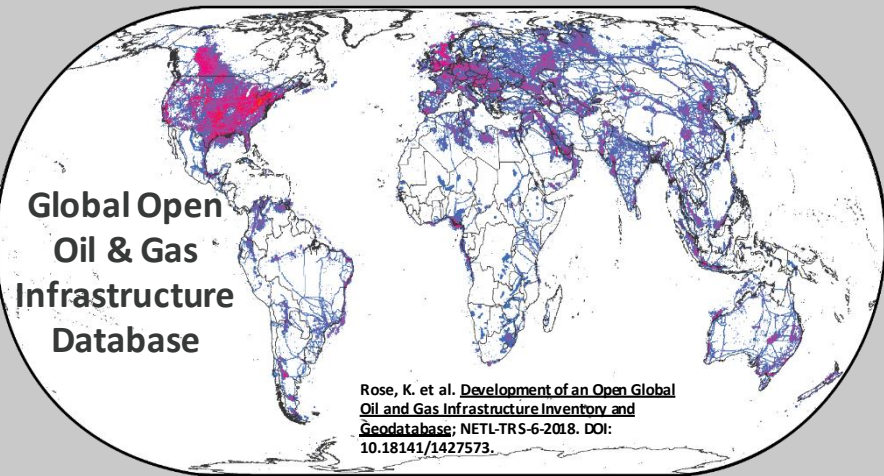
SmartSearch leverages ML+NLP to:

- 1) Analyzing content you like
- 2) Finding new content via www, local, enterprise data stores
- 3) Telling you how relevant the new data is to what you like

Opportunity:

Infinitely scalable to return text, graphical, tabular, image, html, spatial, etc. result

Example applications to date



SmartSearch[®] for building the Carbon Storage Open Database

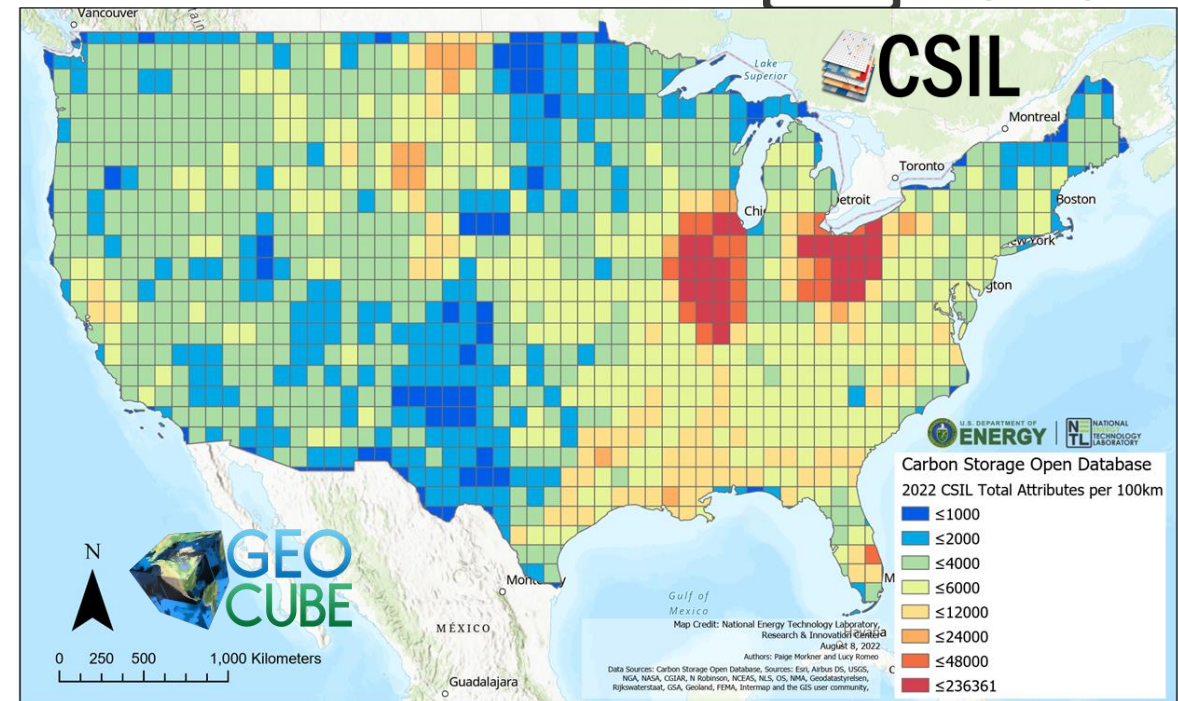
Collection of geospatial data from disparate websites



Aggregation of resources into central database

Cataloging and collection of metadata

Integration of metadata and data into EDX++ framework and online mapping platform

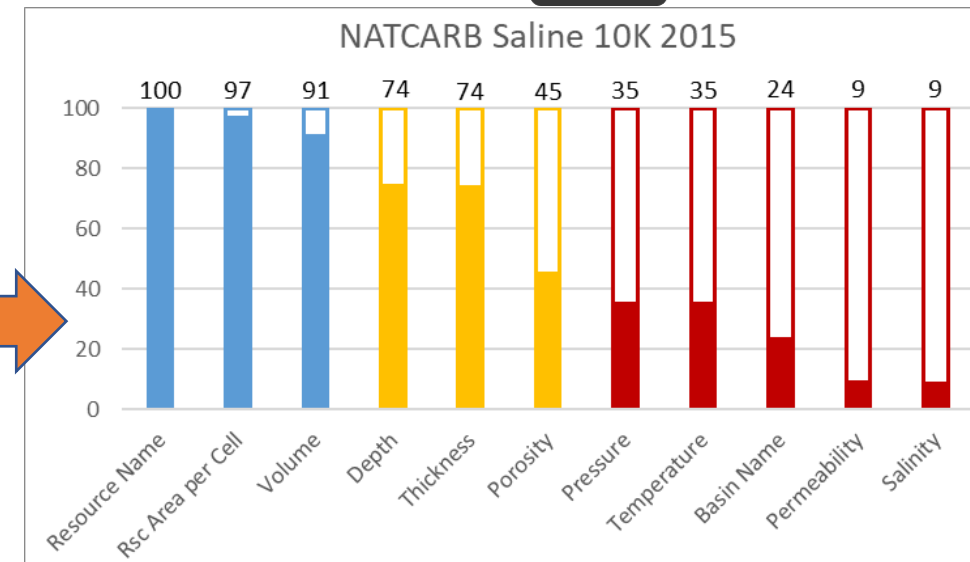


- Use of SmartSearch in finding geospatial data resources to build out collection of 892 geospatial data resources to build the Carbon Storage Open Database on GeoCube
- Morkner, P., et al. 2022. Distilling Data to Drive Carbon Storage Insights. *Computers & Geosciences*.

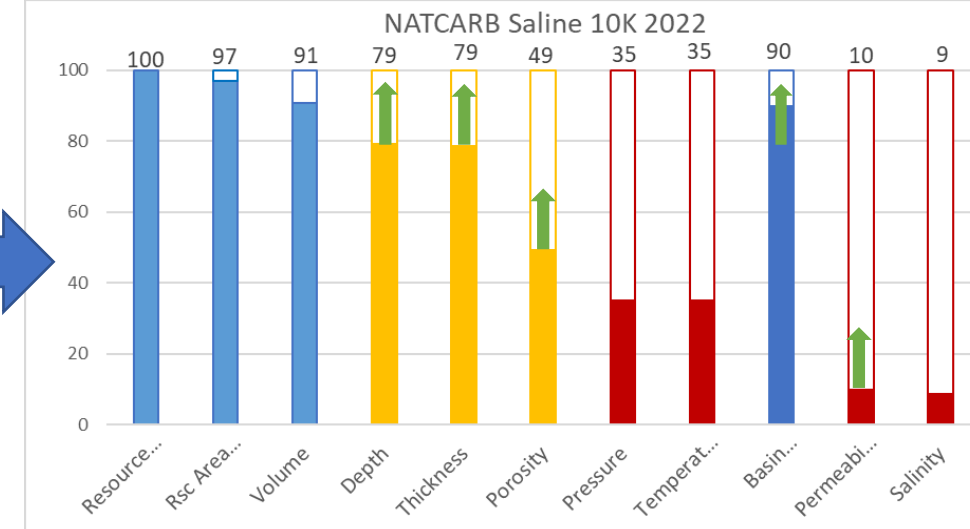
SmartSearch[®] for addressing data gaps in CS data

- SmartSearch was leveraged to find data to meet data gaps in the NATCARB database during EY21
- NATCARB's 2015 Saline 10km layer was used as seed data
- 1000s of results returned above 10% similarity
- Result text bodies were integrated into a postgresQL database where keywords were used to query relevant results
- Researchers integrated thousands of new data results improving data available for:
 - Depth, thickness, porosity, basin name, and permeability

Original
NATCARB
Saline 10km

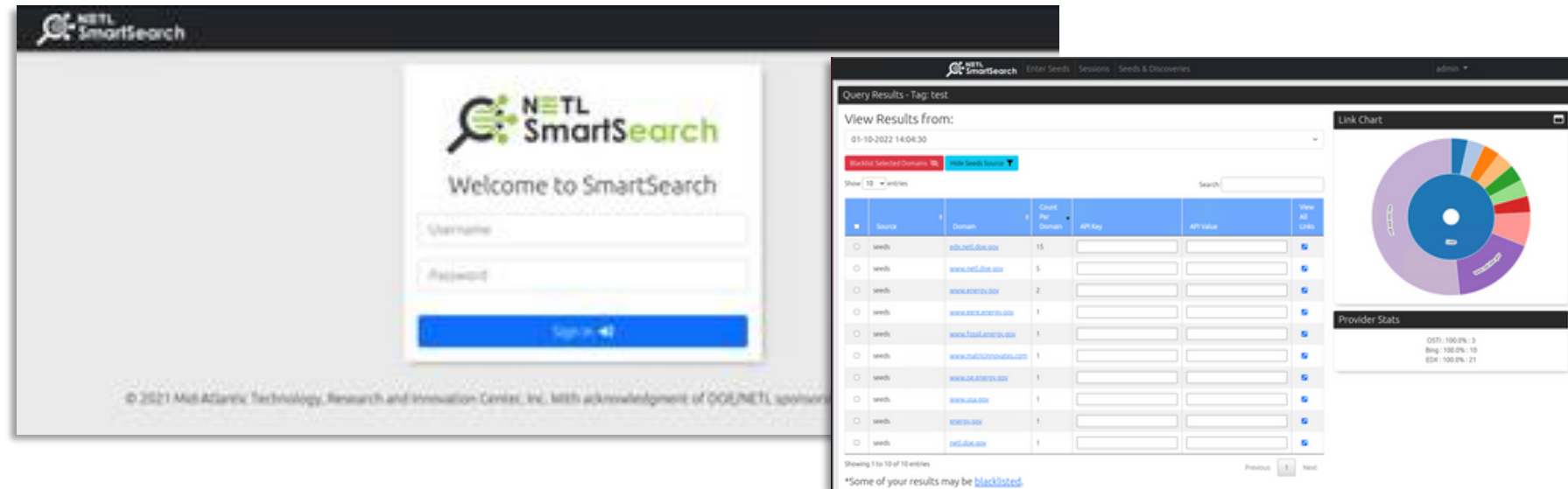


With
additional
data added



Summary

- SmartSearch[®] automates the data discovery process by using a scalable compute environment coupled with NLP and ML
- Will be integrated within EDX++
- SmartSearch supports ongoing research projects
- SmartSearch used to evaluate cloud providers (GCP, AWS, Azure)



Query Results - Tag: test

View Results from:
01-10-2022 14:04:30

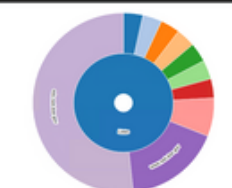
Selected Selected Domains: 10 | Hide Selected Domains

| Source | Domain | Count Per Domain | API Key | API Value | View All Links |
|--------------------------|--------|------------------|---------|-----------|----------------|
| <input type="checkbox"/> | seeds | edu.netl.doe.gov | 15 | | |
| <input type="checkbox"/> | seeds | www.netl.doe.gov | 5 | | |
| <input type="checkbox"/> | seeds | www.edx.gov | 2 | | |
| <input type="checkbox"/> | seeds | www.edx-lab.com | 1 | | |
| <input type="checkbox"/> | seeds | www.edx-lab.com | 1 | | |
| <input type="checkbox"/> | seeds | www.netl.doe.gov | 1 | | |
| <input type="checkbox"/> | seeds | www.edx.gov | 1 | | |
| <input type="checkbox"/> | seeds | www.usa.gov | 1 | | |
| <input type="checkbox"/> | seeds | edx.gov | 1 | | |
| <input type="checkbox"/> | seeds | netl.doe.gov | 1 | | |

Showing 1 to 10 of 10 entries

*Some of your results may be blacklisted

Link Chart



Provider Stats

OSTI: 100.0% / 3
Bib: 100.0% / 10
EDX: 100.0% / 21

Thank you!

References:

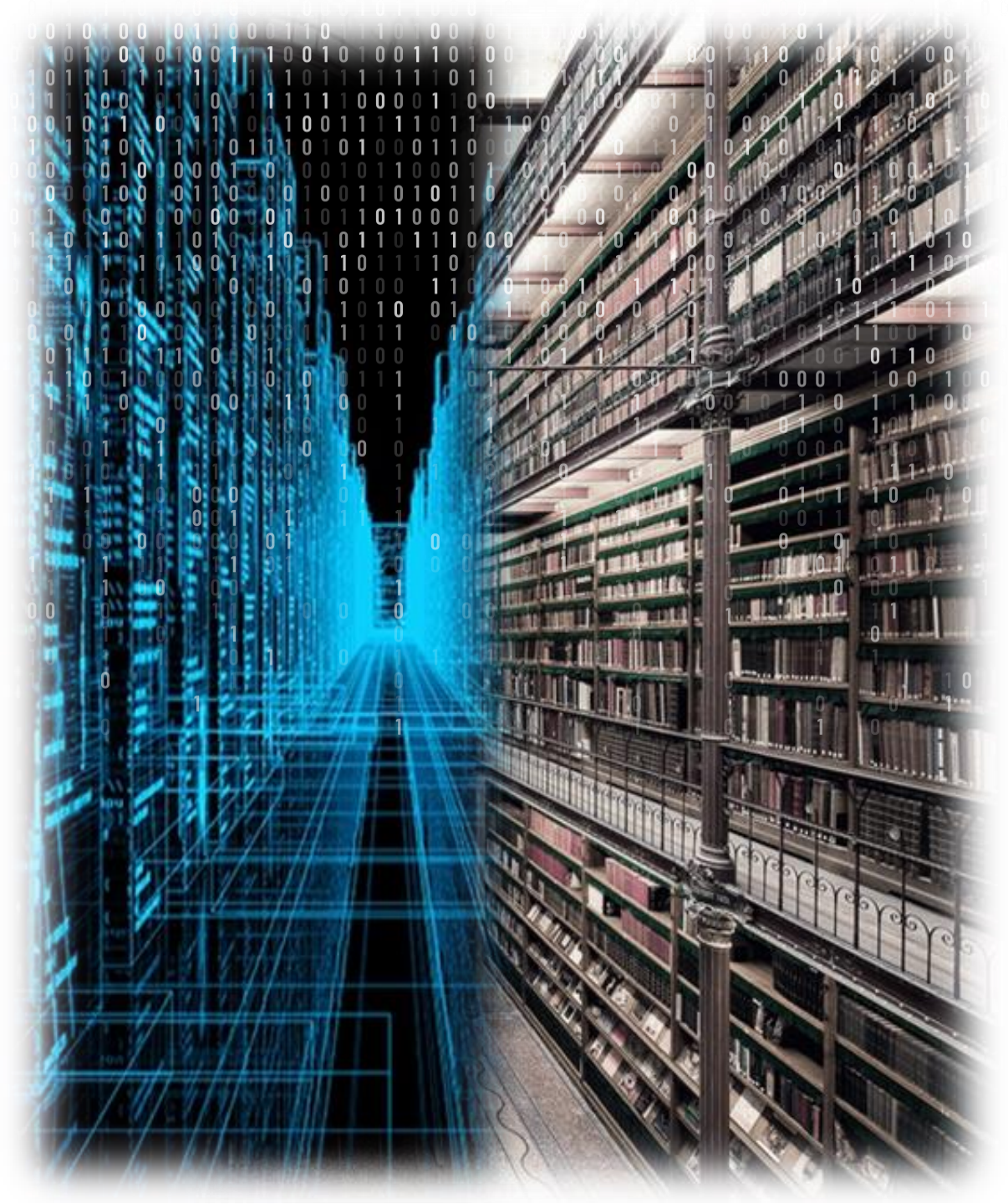
- Morkner, P., et al. 2022. Distilling Data to Drive Carbon Storage Insights. *Computers & Geosciences*.
- Rose, K. et al. [Development of an Open Global Oil and Gas Infrastructure Inventory and Geodatabase](#); NETL-TRS-6-2018. DOI: 10.18141/1427573
- Wenzlick, M., et al. 2021. Data science techniques, assumptions, and challenges in alloy clustering and property prediction. *Journal of Materials Engineering and Performance*.
- Wenzlick M., et al (2021) Incorporating Historical Data and Past Analyses for Improved Tensile Property Prediction of 9% Cr.

CONTACT:

Vic Baker

vic.baker@netl.doe.gov | vic.baker@matricinnovates.com

[SAMI – SAMI \(doe.gov\)](#)





Disclaimer: This project was funded by the United States Department of Energy, National Energy Technology Laboratory, in part, through a site support contract. Neither the United States Government nor any agency thereof, nor any of their employees, nor the support contractor, nor any of their employees, makes any warranty, express or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or any agency thereof. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or any agency thereof.

