# Advanced Data Extraction to Support a Living Database
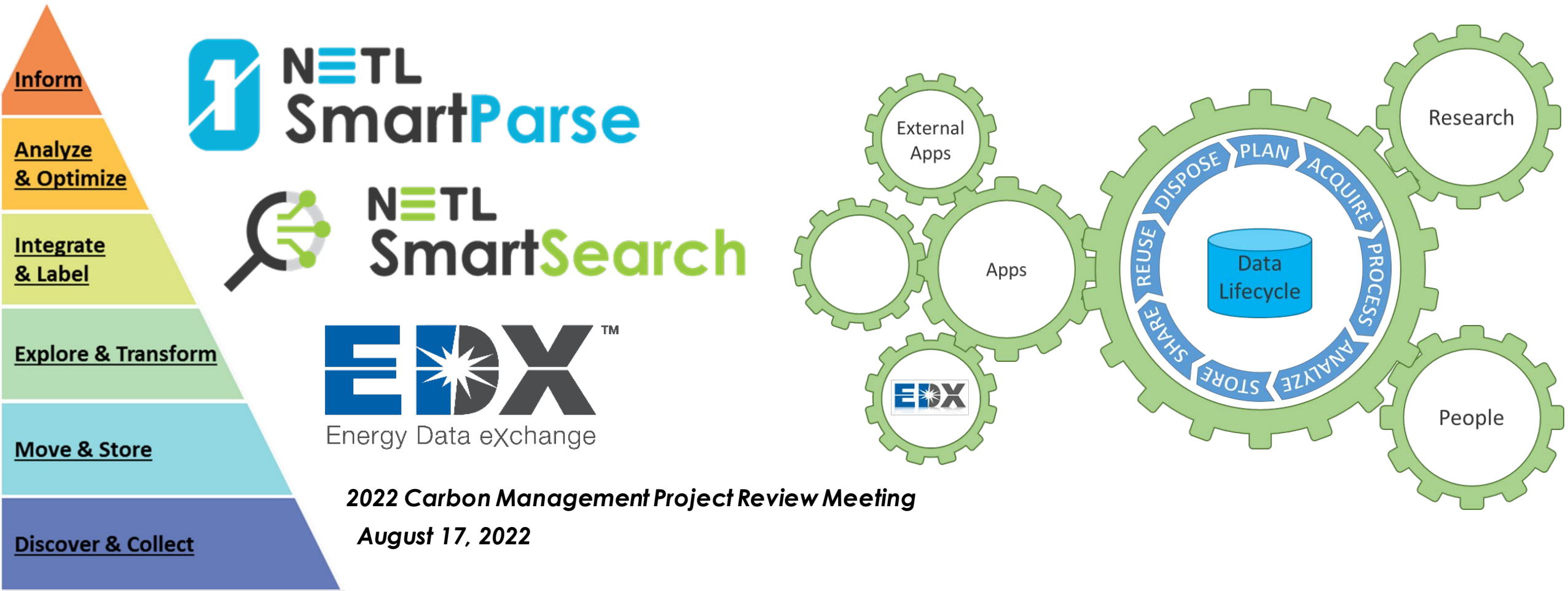
**Michael Sabbatino**

*NETL Support Contractor*
*Research Innovation Center*

*2022 Carbon Management Project Review Meeting*

*August 17, 2022*

# 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 the authors expressed herein do not necessarily state or reflect those of the United States Government or any agency thereof.*

# Author Information

**Michael Sabbatino[1,2], Paige Morkner[1,2], Jennifer Bauer[1], Kelly Rose[1]**

*[1] National Energy Technology Laboratory, 1450 Queen Avenue SW, Albany, OR, 97321 USA*

*[2] NETL Support Contractor, 1450 Queen Avenue SW, Albany, OR, 97321, USA*

*[3] NETL Support Contractor, 3610 Collins Ferry Rd., Morgantown, WV, 26505, USA*

# Research is Data-driven

- **Millions of dollars of research and data are available from carbon storage efforts**
- **How can we preserve and efficiently access those resources to drive the next generation of R&D?**

**Address the needs of the community through advanced AI/ML tools via DOE's virtual data library and laboratory, EDX**

# Carbon Storage Data Lifecycle

**Collection**
- SmartSearch
- Expert-driven research
- EDX submissions

**Metadata development and capture**
- Cataloging
- ReadMe file development
- Natural language processing for keywords, topic modeling, geographic association

**Quality Assessment**
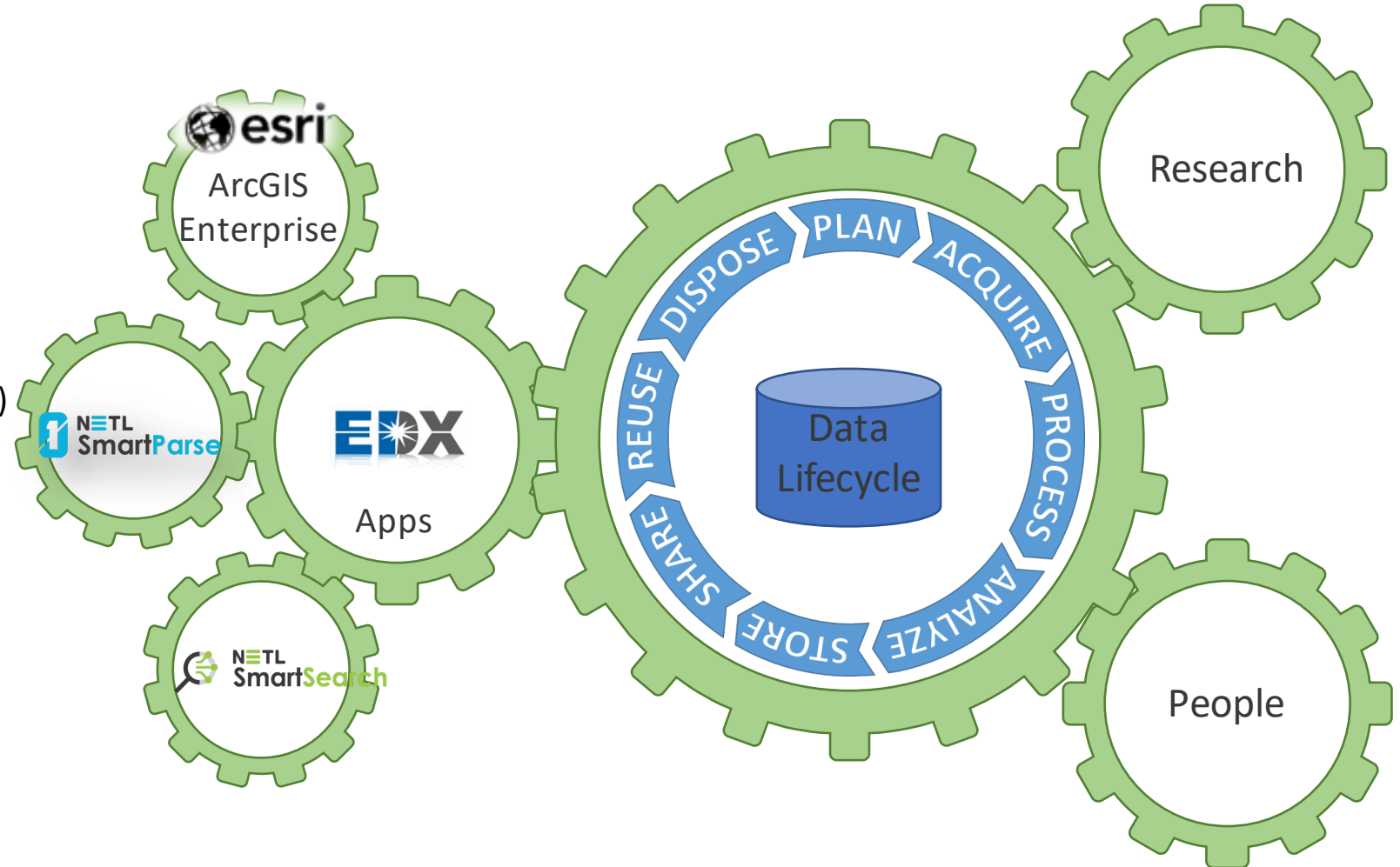- Data ranking
- Data assessment method scoring

**Data Organization and publishing**
- Private workspaces
- Submission packaging
- ArcGIS Enterprise with EDX
  - GeoCube data integration
  - Collaborative Map Applications

# Expanding Integrations for a Living Database

- **Expanded to Utilize ArcEnterprise GIS data management System hosted on AWS servers**
  - **Store/Share Geospatial Data**
  - **Host New Geocube Web Application**
  - **Custom Interactive Map Applications**
- **Store and Share Data in a Structured Secure Database Environment**
  - Reduce Redundant Acquisition
  - Direct Data Access (not file-based storage)
  - Consistent Data with Staff Turnover
  - Enhanced Collaboration
- **Curation of data and knowledge**
- **Allows Direct Analysis from Database**
- **Available On Research network and Watt ML Cluster**

# Types of Carbon Storage Data

**Text-based Data**
- **Documents**
- **Publications**
- **Power points**
- **Memos**
- **Posters**

**Spatial data:**
- **Shapefiles (field, basin, regional scale)**
- **Datasets**
- **Models**

**Image Extraction:**
- **Documents**
- **Presentations**
- **Maps**
- **Posters**

**Other types of data:**
- **Tools**
- **Applications**
- **APIs**
- **LivingDatabase**



I'M ALWAYS HUNGRY HUNGRY

**EDX** Energy Data eXchange

Feeding the data hippo!

# Advanced AI/ML Tools for CS Data Lifecycle

**Challenge:** Continue increasing available data while enhancing metadata, searchability, and curation.

**Solutions:**

- **Natural Language Processing:**
  - **text-based resource classification, organization, keyword identification**
  - **metadata extraction and preservation**
  - **geographic association (for searchability)**
- **Image Classification and Text Extraction**
  - **Identify images from papers, posters, and documents**
  - **Classify images and extract text**
  - **Extract image metadata**
- **ArcGIS Enterprise on AWS:**
  - **Geographic database development (Geocube)**
  - **Interactive map creation and collaboration**
  - **Integration with EDX**

NETL SmartSearch

NETL SmartParse

Inform — Use of data for site selection and modeling

Analyze & Optimize

Integrate & Label (Analytics, metics, features and training data)

Explore & Transform (Curation, cleanup and visualization)

Move & Store (Collaborative data management)

Discover & Collect (Subsurface and contextual data from various sources)

80% of time is spent acquireing, curating, labeling and organizing data

# Data Cleaning for ML, AI and Spatial Analysis

**Identify data to be collected**

**Includes:**

- **Data, Papers, Catalogs of Data, Online Sources, and Metadata**

**Data collected and processed using Python tools to move, quantify and label data**

# Natural Language Processing (NLP) Unsupervised ML for Document Classification

- Latent Dirichlet allocation **(LDA) model based on corpus of text-based documents**
- Topic names assigned by subject-matter experts
- **Each document is classified** by % of each topic it's associated with
- **Each document has 50+ keywords identified** and can be **associated with metadata on EDX**
- **Parse geographic location to associate with each document** – when possible



Completed on Local Desktop PC

Completed on NETL Watt ML Cluster

Loop Through All Documents in Database
Completed on Local Desktop PC

# Machine Learning Image Data Extraction

- Object Detection Model Development Process
  - Use transfer learning to train object detection model for specific image and data types
  - Detect Graphs, Diagrams, Photos, Maps, and Tables
  - Image Labeling and process Developed with help from Mickey Leland Energy Fellowship

Images and Tables Targeted for Data for Extraction

# Machine Learning Data Extraction
## Utilize Object Identification ML Models to Extract Additional Data

# Machine Learning Data Extraction
## Utilize Object Identification ML Models to Extract Additional Data



Figure 2: Forecast California GHG Emissions by Sector

Source: Modified from (Schiller 2010)

Table 4: Electricity Demand and CO2 Emissions in 2010 and Forecasts for 2050

| | Demand (TWh/year) | Emissions (Mt CO$_2$) |
|---|---|---|
| 2010 | 300 | 100 |
| 2050 Goals | - | 77 |
| 2050 BAS | 1200 | 140 |
| 2050 Scenario* | 500-600 | 60 |

Source: (Greenblatt and Long 2012)

# Use Case: Abstract White Paper Data Extraction

## NLP and Machine learning to Classify and Text and Images

- 10 topic, 5 topic and variable
- PCA Analysis
- Keywords and Custom Stop Word List

PCA Papers 5 Topics 241 Papers
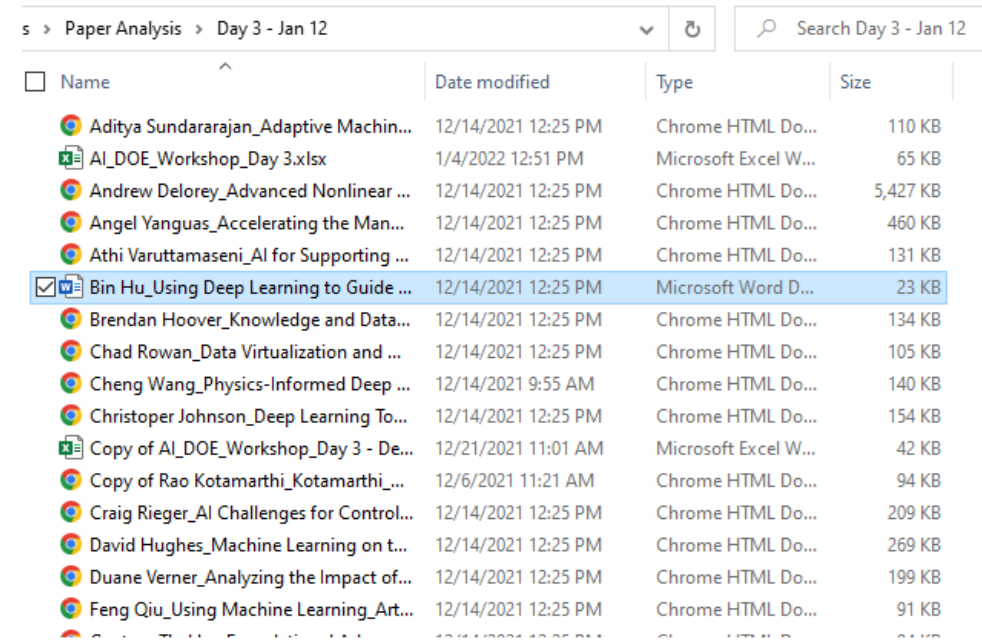
PCA AMO 5 Topics 3037 Documents

# Use Case: Abstract White Paper Data Extraction

## NLP and Machine learning to Classify and Text and Images

- Data Input Zipped Papers, Spreadsheets, and Images
- Process with NLP (Gensim) to Create topic model
- Convert PDF to JPG for preprocessing
- Used trained Yolo model using transfer learning
- Extracted images from papers and classify

# Use Case: Abstract White Paper Data Extraction

## NLP and Machine learning to Classify and Text and Images

- Data Input Zipped Papers, Spreadsheets, and Images
- Process with NLP (Gensim) to Create topic model
- Convert PDF to JPG for preprocessing
- Used trained Yolo model using transfer learning
- Extracted images from papers and classify

**Data Virtualization and Management for Energy R&D**

ROWAN, Chad[1]

ROSE, Kelly[2], BAUER, Jennifer,

BAKER, Vic, JONES, TJ, MCFARLAND, Daniel

1. Maximus, LLC, National Energy Technology Laboratory, 3610 Collins Ferry Road, Morgantown, WV 26507-0880
2. Department of Energy, National Energy Technology Laboratory, 3610 Collins Ferry Road, Morgantown, WV 26507-0880
3. Matric Innovates, 3610 Collins Ferry Road, Morgantown, WV 26507-0880

Curation and access to federally funded research products is key to support the current data revolution, FAIR data practices, and ever-changing landscape of artificial intelligence and machine learning (AI/ML) techniques across the U.S. Department of Energy (DOE). In 2011, the DOE National Energy Technology Laboratory (NETL) began development and maintenance of the Energy Data eXchange (EDX) to address the needs of data management while building the functionality needed to support a virtual laboratory. The motivation of this platform was to address the need for rapid response of data intensive challenges including human/natural disasters and fundamental research.

EDX has been leveraged significantly by the DOE Office of Fossil Energy and Carbon Management's geospatial and geoscience programs for carbon storage, rare earth elements, unconventional oil and natural gas, and others. It provides users with an online collection of data, capabilities, and resources that advance ongoing research while maintaining the IT and

**U.S. DEPARTMENT OF ENERGY**

## NLP and Machine learning to Classify and Text and Images

- Data Input Zipped Papers, Spreadsheets, and Images
- Process with NLP (Gensim) to Create topic model
- Convert PDF to JPG for preprocessing
- Used trained Yolo model using transfer learning
- Extracted images from papers and classify

# Use Case: Abstract White Paper Data Extraction

**Image Classification Demo/Results**

# Use Case: Abstract White Paper Data Extraction

**Image Classification Demo/Results**

# Use Case: Abstract White Paper Data Extraction

**Image Classification Demo/Results**

| | A | B | C | D | E | F |
|---|---|---|---|---|---|---|
| 1 | Image Type | Number Found | Example Document | Example Page | Image Snapshot | |
| 2 | Diagram | 72 | Akash Dhruv_WhitePaper_P1Q2A1.pdf | 2 |  | |
| 3 | Graph | 28 | Andrew Delorey_Advanced Nonlinear Measurements in the Earth.pdf | 1 |  | |
| 4 | Map | 46 | ed_ScheinkerA_AdaptiveMachineLearningforTime-VaryingParticleAcceler | 2 |  | |
| 5 | Photograph | 21 | Pete Beckman_Intelligent Edge to Predict Extreme Events.pdf | 3 |  | |
| 6 | Table | 122 | J. Darby Smith_smith_inverse_problem_P1Q2A1.pdf | 2 |  | |

# Synergy Opportunities

## Collaborative cross-project technology

- Use material same NLP tech
- Using other NLP Models Louvian Community Detection

# Supporting Data Collection, Curation & Analysis in Other Areas

## Data mining, including…

| Alloy (wt%) | N | C | Mn | Cr | Mo | Ni | Si |
|---|---|---|---|---|---|---|---|
| 316LNSS-7N | 0.07 | 0.027 | 1.7 | 17.53 | 2.49 | 12.2 | 0.22 |
| 316LNSS-11N | | | | 17.62 | 2.51 | 12.27 | 0.21 |
| 316LNSS-14N | 0.14 | 0.025 | | | 2.53 | 12.15 | 0.2 |
| 316LNSS-22N | 0.22 | 0.028 | 1.7 | 17.57 | 2.54 | 12.36 | 0.2 |

*Structured Data*

*Images and Graphs*

| CT C | CS (MPA) | RT, hrs |
|---|---|---|
| 593 | 310.3 | 1.45 |
| 593 | 275.8 | 5.5 |
| 593 | 275.8 | 6.33 |
| 593 | 206.8 | 55 |
| 593 | 171.7 | 357 |
| 593 | 144.8 | 1446 |
| 704 | | 0.37 |
| 704 | 172.4 | 1.5 |
| 704 | 137.9 | 9.5 |
| 704 | 103.4 | 50.5 |
| | 75.8 | 337 |
| 704 | 62.1 | 1227 |
| 816 | 103.4 | 0.75 |
| 816 | 89.6 | 1.87 |
| 816 | 68.9 | 12.75 |
| 816 | 48.06 | 84.3 |
| 816 | 36.5 | 331.8 |
| 816 | 29.0 | 1153 |

*Measurements*

Fig. 1. Orientation imaging micrographs of solution annealed 316LN SS containing nitrogen (wt.%) of (a) 0.07, (b) 0.11, (c) 0.14 and (d) 0.22 N. Nearly equiaxed grains and annealing twins have been observed.

RESEARCH ARTICLE

**Materials data analytics for 9% Cr family steel**

Vyacheslav N. Romanov ✉, Narayanan Krishnamurthy, Amit K. Verma, Laura S. Bruckman, Roger H. French, Jennifer L.W. Carter, Jeffrey A. Hawk

First published: 15 February 2019 | https://doi.org/10.1002/sam.11406

U.S. Department of Energy. DE-FE0002688S.

Read the full text ›

**Abstract**

A materials data analytics (MDA) method was developed in this study to evaluate publicly available information on mechanical testing and to handle nonlinear relationships and the sparsity in materials data for a given alloy class. The overarching goal is to accelerate the design process for a new Ni-based alloy for fossil energy applications. Data entries in the mechanical testing of iron base alloy compositions, several processing parameters, and several microstructural/tensile mechanical tests selected for this study were arranged in 34 columns with 915 rows. While detailed microstructural information was not available, it is assumed that the compositional space for the 9 to 12% Cr steels is limited such that all data entries have a tempered martensitic microstructure during service. Establishing a hierarchy of first-order trends in the publicly available data requires the MDA to filter out the biases. Complexity of the phase transformations and microstructure evolution in the multicomponent alloys (using 21 chemical elements) with major influence on mechanical

*Contextual knowledge*

MARBN : 9Cr-3W-3Co-VNb, 120 - 150 ppm B & 60 - 90 ppm N
P92 : 9Cr-0.5Mo-1.8W-VNb, 20 ppm B & 500 ppm N

*Test results*

## Move & Convert…

EDX™
Energy Data eXchange

## …& use in predictive analytics for alloy behavior

Actual and predicted creep rupture time using the Gradient Boosted Regression ML Algorithm

**Evaluating machine learning models to:**
- address data gaps
- identify key features in lifetime behavior of the alloy

# Lessons Learned

## Machine Learning, Artificial Intelligence, and Natural Language Processing are Difficult
- Whatever happened to Watson?

## Lack of Labeled Training Data
- Training data is time-consuming to develop and can be costly

## Data availability is limited with Living Database
- Currently deployed on the Research Network
- The database would improve if deployed on a cloud service or other shared environment



THE WALL STREET JOURNAL.

CIO JOURNAL
**Data Challenges Are Halting AI Projects, IBM Executive Says**
The cost and hassle of collecting and preparing data comes as a shock for some companies, according to Arvind Krishna

By *Jared Council*
May 28, 2019 5:30 a.m. ET



*Structured data*

**What Ever Happened to IBM's Watson?**
IBM's artificial intelligence was supposed to transform industries and generate riches for the company. Neither has panned out. Now, IBM has settled on a humbler vision for Watson.

https://www.nytimes.com/2021/07/16/technology/what-happened-ibm-watson.html

# Summary

> **FE and Carbon Storage program investments into data curation and management have led to the development of AI/ML tools and the preservation of millions of dollars of research products which benefits ongoing and future research. This has led to:**

- **A better understanding of CS relevant open- data density and data quality throughout US and Canada**
- **Improved access through the integration of CS data resources on EDX into GeoCube, SmartSearch, and SmartParse (EDX version of NLP tools presented here) for further searchability with spatial searches and keyword searches**
  - Updates to GeoCube for enhanced spatial searchability and integration of modeling tools to come
- **EDX AI/ML data discovery, labeling, integration tool developments trained to support Carbon Storage, SMART-CS, and NRAP**
  - Deployment of AI/ML algorithms to allow on-demand data discovery and integration, ready-made for each end-user needs

# What's next: EDX4CCS

**SmartSearch** NETL

**SmartParse** NETL

**EDX4CCS**

**Data,** Integration, generation, and deployment to feed SMART, NRAP, and regulatory models

**Tools,** Develop and/or integrate the deployment of tools for data interaction and visualization, decision-support such as for pipelines, regulatory permitting, resource characterization, data visualization, and more

**Core CCS EDX DisCO2ver platform,** Broader community virtualized data computing platform, and central EDX CCS data and tool hub

**Carbon Storage Open Database**

GEO CUBE

EDX®

EDX 4CCS
Energy Data eXchange

U.S. DEPARTMENT OF ENERGY
Fossil Energy and Carbon Management

# Thank you!

# NETL
# RESOURCES

VISIT US AT: www.NETL.DOE.gov

🐦 **@NETL_DOE**

📷 **@NETL_DOE**

f **@NationalEnergyTechnologyLaboratory**

U.S. DEPARTMENT OF **ENERGY**

U.S. DEPARTMENT OF **ENERGY**

# Appendix

- These slides will not be discussed during the presentation, <span style="color:red">but are mandatory.</span>

# Benefit to the Program

- Task 27 supports the development of data, materials, maps, analyses, and figures for the Carbon Storage Atlas, Natcarb Viewer, and Natcarb database. This includes the release of new data insights to the GCS community, through the sixth edition of the Carbon Storage Atlas, and through bi-annual updates to the Natcarb Viewer and Natcarb database.

- Task 28 focuses on addressing CS R&D data curation challenges associated with ingesting, describing, and curating data products from DOE FE to ensure enduring access and more efficient utilization of those resources using AI/ML enhanced approaches to support future CS R&D. Ultimately, this effort will result in tools, data resources, and virtual capabilities for the CSP and community to facilitate efficient CS data discovery, integration, and curation using NETL's EDX

- Use of EDX and development of tools to support the collection, curation, organization, labeling, and publishing of large quantities of data for carbon storage. Whether laboratory, field or computational, CS R&D is both a producer and consumer of data resources (datasets, tools, models, etc.). However, while the volume of open, online data is increasing exponentially, scientists struggle to find, access and make operable data products from previous R&D projects due to insufficient and/or burdensome online data curation tools and outdated techniques.

# Project Overview
## Goals and Objectives

– Funded by DOE as part of Carbon Storage DE FE-1022465, Tasks 27 and 28

– RSS Contract and ITSS contract researchers

– Ongoing performance dates 2018-2022

– Project Participants

- PI: Kelly Rose

- LRST: Paige Morkner, Michael Sabbatino, Andrew Bean, Lucy Romeo, Patrick Wingo

- ITSS: Chad Rowan, TJ Jones, Aaron Barkhurst, Vic Baker

# Organization Chart
# Carbon Storage Data

**Project Partners**
DOE
NETL
RCSPs – Big Sky Carbon Sequestration Partnership, Southwest Partnership, Southeast Regional Carbon Sequestration Partnerhsip, Midwest Regional Carbon Sequestration Partnership, Midwest Geological Sequestration Consortium, Plains CO2 Reduction Partnership.

**Lead Organization**
NETL

**Principal Investigators**
Kelly Rose, Jennifer Bauer

**Task 28**
Curation of Carbon Storage R&D Products Through Advanced Data Computing Solutions

**Lead: Jennifer Bauer**
Contractors: **Chad Rowan, Michael Sabbatino**, Paige Morkner, Andrew Bean, Lucy Romeo, TJ Jones, Aaron Barkhurst, Vic Baker, Other Matric Software Engineers and Developers

**Task 27.0**
Next Generation Development, Deployment, and Modernization of Database, Tools, Online Viewer, and Atlas

**Lead: Jennifer Bauer**
Contractors: **Paige Morkner**, Michael Sabbatino, Patrick Wingo, Andrew Bean, TJ Jones, Aaron Barkhurst, other Matric Software Engineers and Developers
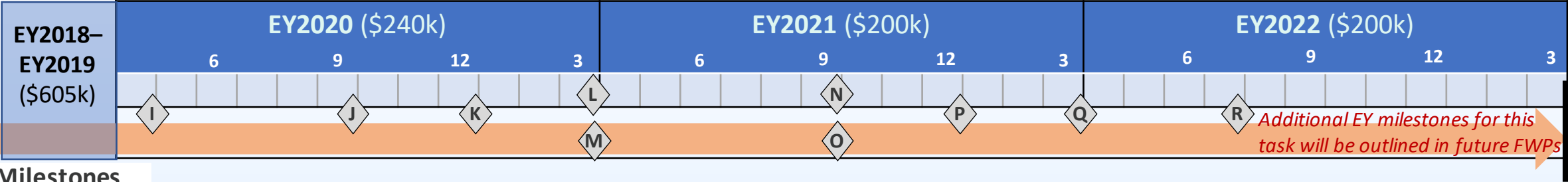
# Task 28.0: Project Timeline Overview

## Curation of Carbon Storage R&D Products Through Advanced Data Computing Solutions
(PIs: Michael Sabbatino, Jennifer Bauer)



*Additional EY milestones for this task will be outlined in future FWPs*

### Milestones

| Number | Expected Completion Date | Milestone Description |
|---|---|---|
| EY20.28.I | 04/30/2020 | Push to public on EDX appropriate **MGSC** Partnership data products. |
| EY20.28.J | 09/30/2020 | Deploy LivingDatabase beta version capability in EDX, private side, for CS teams (e.g., RCSPs) use and testing. |
| EY20.28.K | 12/31/2020 | Integration of CSP data products that are spatially related through enhanced EDX spatial search and discovery tool on GeoCube. |
| EY20.28.L | 03/31/2021 | Deploy NETL SmartSearch version 2 algorithm in EDX to support automated gathering of open, CS relevant data. |
| EY20.28.M | 03/31/2021 | Deploy LivingDatabase version 1 capability in EDX, private side, for CS teams (e.g., RCSPs) use and testing. |
| EY21.28.N | 09/30/2021 | Develop and test SmartSearch and SmartParse beta integration. |
| EY21.28.O | 09/30/2021 | Complete testing of Living Database dashboard tools. |
| EY21.28.P | 12/31/2021 | Create additional training data for SmartParse image, graph, and table extraction model improvement. |
| EY21.28.Q | 03/31/2022 | Develop beta Living Database user interface and dashboard. |
| EY22.28.R | 07/29/2022 | Ingestion and push to public on EDX appropriate **SW Regional Partnership** data products. |

**Chart Key**

- ◇ Milestone
- ▮ Project Completion
- ▮ Go/No-Go Timeframe

## Key Accomplishments/Deliverables

- 2018–Present, Addition of **Big Sky**, PCOR, Midwest CS Partnership, SECARB, and MGSC data and resources on EDX, for a combined total of 3,037 and 1.64 TB of data
- 2018–2020, Big data computing cluster, Watt, set up and work to directly link EDX with these computing capabilities
- 2019–2021, Test and validate SmartSearch for use with commercial cloud & EDX to evaluate capabilities to assimilate relevant CS data; including work as part of an NDA with Google and collaboration with DOE-HQ OCIO
- 2020–2021, Develop Living Database logic to host and storge large volumes of CS data
- 2021–2022, Deploy beta instance of Living Database front end and dashboard tools
- 2022, Addition of any final RCSP and other CS resources to EDX

## Value Delivered

- **Collecting, curating, and cataloging** data from all regional CS partnerships and open-sources.
- **Developing capabilities** to query curated data.
- **Delivering** EDX's public-private capabilities, including growing access to its **big data computing** cluster and Amazon Web Services (AWS) **cloud services**, seek to facilitate more effective research **for DOE-FE subsurface scientists**.
- **Pairing EDX hosted CS data resources and products with other online capabilities**, data, custom ML algorithms and capabilities to enhance user experience and provide research teams with the resources needed to make subsurface energy research more efficient, reduce redundancy, and drive innovation.

* Task 28.0 is integrating data into an existing tool with no development of a technology. Therefore, no TRL is assigned.

# Bibliography

– List peer reviewed publications generated from the project per the format of the examples below.

- Morkner, P., Bauer, J., Creason, C., Bean, A., and Rose, K., "A Data Quality Assessment Method to Support Carbon Storage," in preparation . Target journal: *Nature Scientific Data.* (Tasks 27.0, 28.0)

- Morkner, P., Creason, C., Sabbatino, M., Wingo, P., DiGiulio, J., Jones, K., Greenburg, R., Bauer, J., and Rose, K., "Distilling Data to Drive Carbon Storage Insights," accepted pending final revisions, *Computers and Geosciences*. (Tasks 27.0, 28.0)

- Barkhurst, A., Morkner, P., Bauer, J., Rose, K. GeoCube, TRS report, in prep, target completion Fall 2021.