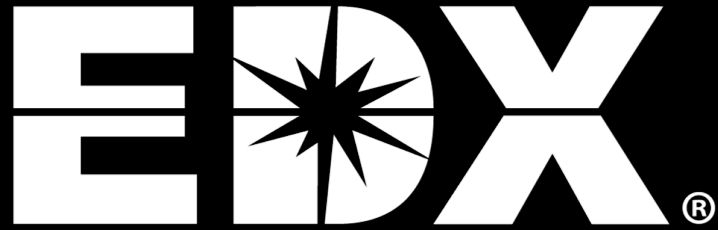# Using Cloud to Host Large Open-Source Datasets

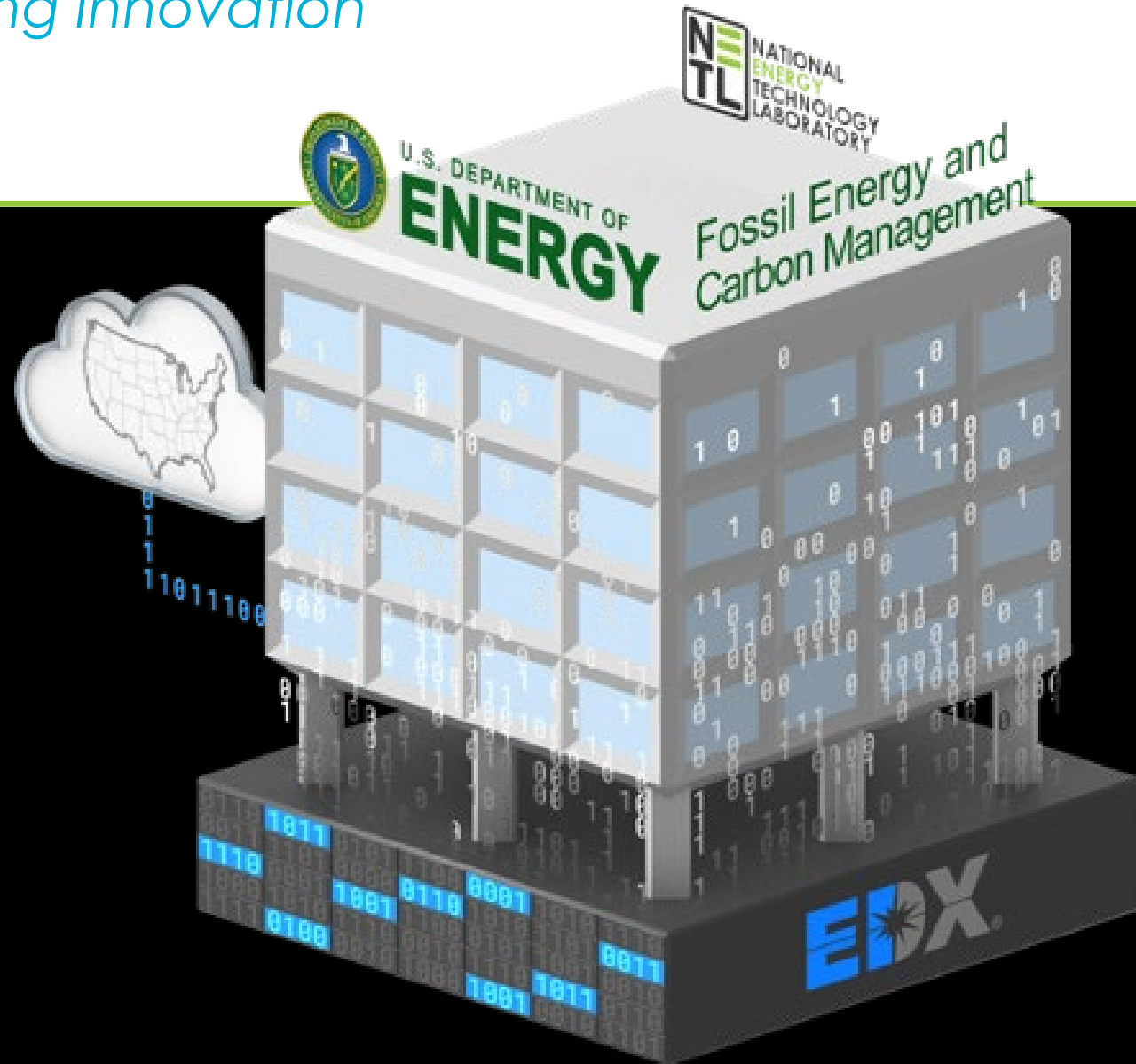*Catalyzing Collaboration & Accelerating Innovation*

**EDX**®
Energy Data eXchange

Chad Rowan

National Energy Technology Laboratory (NETL)

Computational Science & Engineering

Advanced Computing & Artificial Intelligence

# Disclaimer

# FECM R&D is Impeded by Common Challenges

- Finding and accessing relevant datasets

- Publishing & preserving R&D data products

- Accessing previously developed R&D data

- Sharing secure R&D scale data resources among team

- Collaborating across multi-organizational teams

- Need to access prior R&D data products to accelerate next-generation innovations
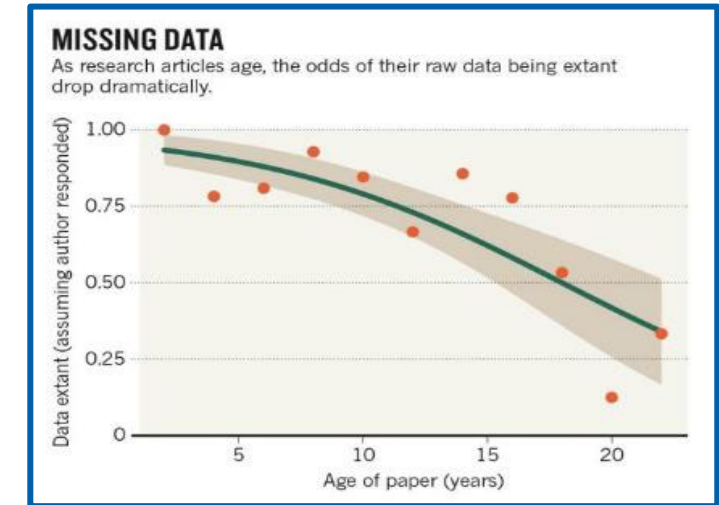


Image from : http://www.nature.com/news/scientists-losing-data-at-a-rapid-rate-1.14416

# The Bigger the Data the Bigger the Challenges

- Datasets are increasingly becoming larger

- Large datasets are more difficult to curate and make publicly available

- Large datasets are not exempt from federal data publishing guidelines

# How has FECM addressed these issues?

**a web-hosted, virtual library and laboratory that supports the NETL/FECM community**

# Advantages of publishing data products



OSTI DOI Number

Data License

Data Citation

Data Access

**Many journals require models, tools and data be publicly available prior to journal publication.**

*can help!*

# Our Big Data Journey

- In 2011, we thought big data was a few gigabytes

- Transferring gigabyte files across the Internet to EDX was slow, but doable

- In the last few years, we have datasets that have grown to over 100 terabytes

- Transferring terabyte files to our on-prem EDX server over the Internet was not feasible

This Photo by Unknown Author is licensed under CC BY-SA

# The Recent Past

Large dataset is mailed to EDX Support on an external drive

Large dataset is uploaded to the Watt Machine Learning Cluster

A researcher requests a copy of the large dataset

The large dataset is transferred to an external drive

The data on the external drive is shipped to the researcher

# We needed a more nimble approach to Big Data

- Rather than upgrade our on-prem hardware we started exploring cloud options

- We were quickly introduced to Cloud Open Data Programs

- All major cloud service providers have Open Data Programs

# What do cloud Open Data Programs provide?

- Free hosting and egress of large, publicly accessible data

- Increased access to ODP hosted datasets

- Faster upload/download speeds

- Access to cloud tools

- Access to cloud compute

- Access to cloud data analytics

# What is the benefit of hosting Carbon Storage data in an ODP?

- ODP is **FREE** for large, public datasets

- Provides **ACCESS** to large, public datasets that were historically difficult to share

- Increases **VISIBILITY** and **DISCOVERABILITY** of large, public datasets

- Facilitates **CLOUD COMPUTE** at the source of the data

# Google's Open Data Program

- After reviewing Open Data Programs from all of the cloud service providers we started our ODP journey with Google

- Google shipped a 300TB transfer appliance to transfer our first collection of large datasets

- The data was transferred to the appliance and sent back to Google where it was uploaded to a GCP bucket

- The EDX Team is currently working with Google engineers to apply metadata, making the dataset discoverable for use/re-use

# What was included in the first ODP package?

✓ 5 datasets

✓ Over 200TBs of data

✓ Over 24M data files

**Total Disk Usage (TBs)**



Bar chart values:
- Illinois Basin Decatur Project (IBDP): 69.1
- Joule Wave Model: 16.7
- MRCSP: 8.1
- MSEEL: 105.8
- Stress in Deep Subsurface: 2

# Why are Open Data Programs Free?

- Large datasets are desirable

- CSPs know if they host some of your data for free they have a better chance of hosting your other data at a cost

- CSPs can market cloud tools for compute and visualization that incur a cost

# Open Data Programs facilitate the concept of EDX++
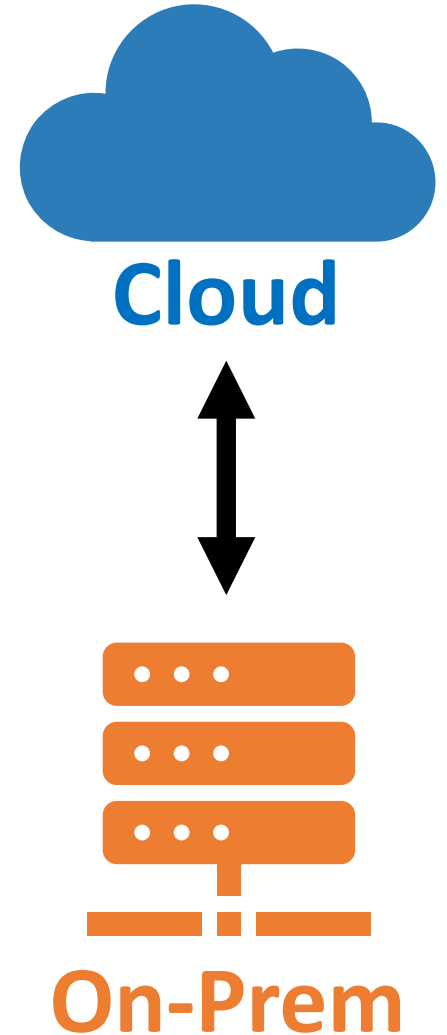
## Freedom for users to use any cloud service providers
- Compute at the source of the data
- Utilize APIs to move data
- Web hosted applications

## Improves flexibility and performance
- Does not limit users to one storage & compute platform
- Compute occurs at the data source

## Resilience
- Redundancy across multiple regions
- Strategic alignment for data transfer and compute across multiple cloud service providers

**Cloud**

**On-Prem**

# ODPs Facilitate FECM R&D Data Use and Reuse

**Carbon Storage Program**

- **Free Hosting and Egress:**
  - Over 200 terabytes of data
  - Over 24 Million data files
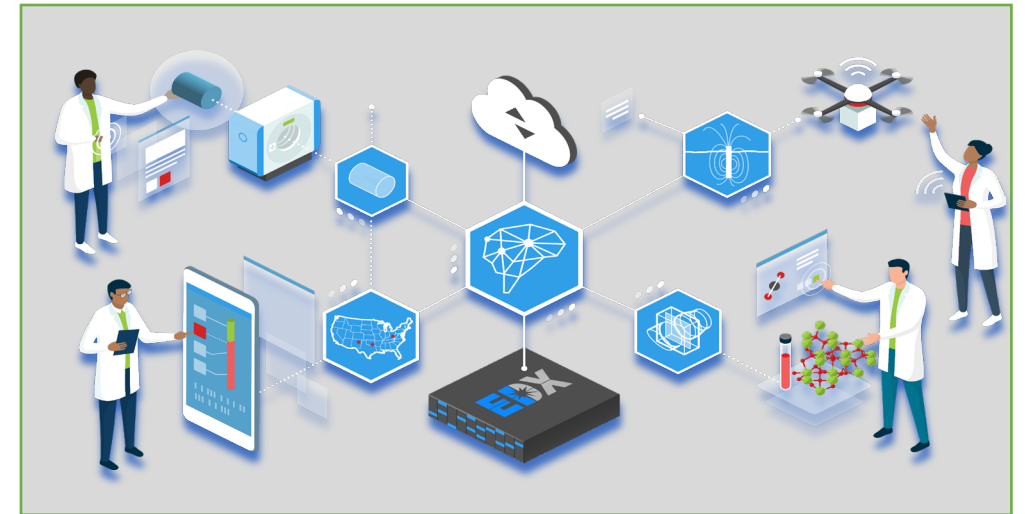  - Supporting use/re-use of current and future Carbon Storage research efforts



**just a few examples**

# The future is now

- Evolving into a multi-cloud solution

- Accelerating AI/ML

- Tackling data compute in the cloud and on-prem

- Improving transfer speed, security, and pipe

# Any questions?

**Key Resources**

- EDX Reference Shelf

- Focused training for research teams (Request Training)

- EDX Training Videos (pre-recorded)

- Robust API Documentation

- **Contact** EDXSupport@netl.doe.gov or SAMI@netl.doe.gov