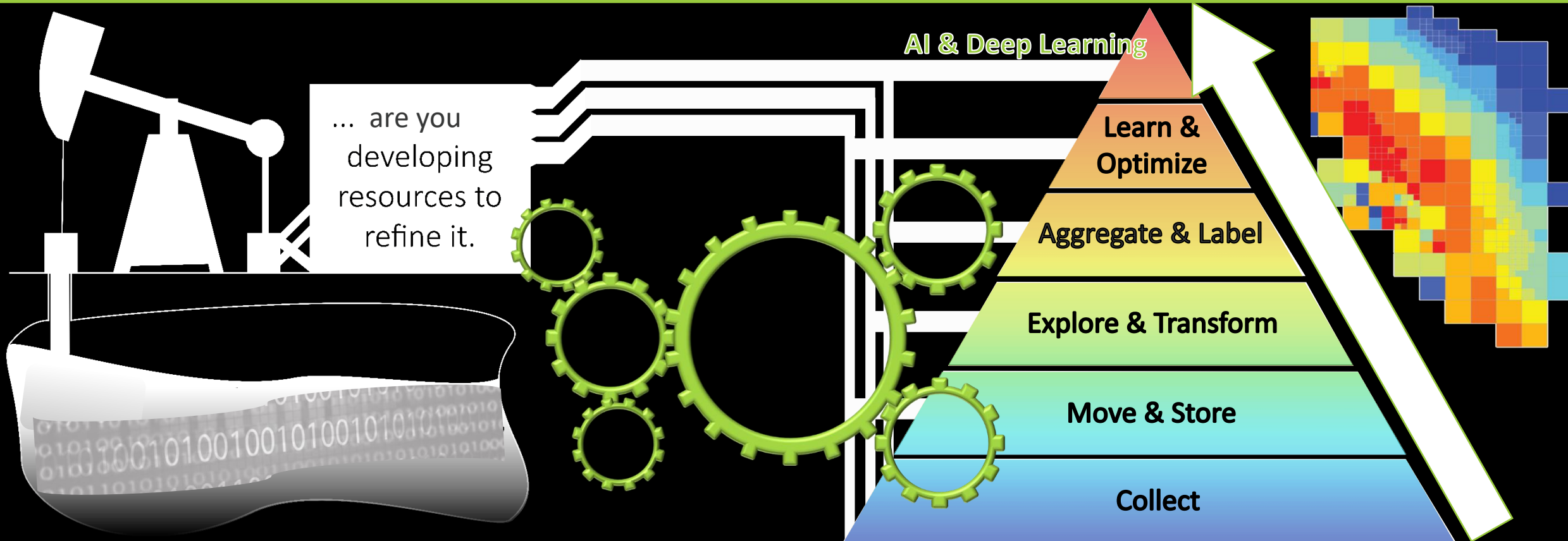


Data is the new oil...



Kelly Rose, HEI Workshop
September 13, 2018

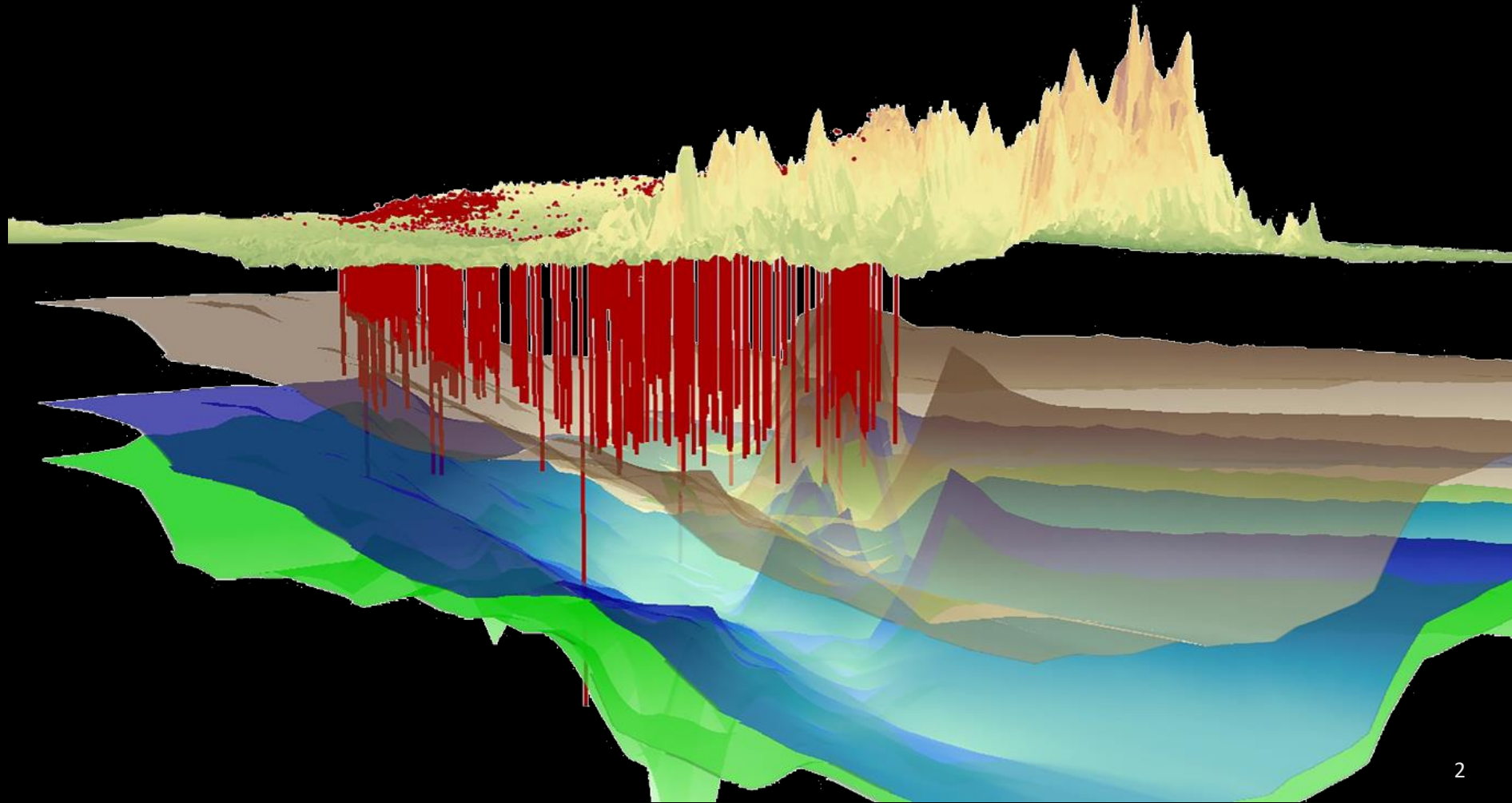


Solutions for Today | Options for Tomorrow

Building a Data Framework for R&D

ML/BD solutions, **tools and capabilities** can be devised or implemented to streamline and automate **data collection, movement, and transformation**

- Collect **disparate datasets and contextual information**
- Incorporate data across **various scales and formats**
- Span surface-subsurface
- **Not all data is equal**, not all data are easy...but there is more out there...put it to work!



NETL's Geo-Data Science:

Inventing Intelligent Solutions to DOE FE Data & R&D Needs

**Data
Discovery**

**Data
Curation**

**Data
Interoperability**

**Data
Analytics &
Visualization**

*Developing & innovating
data, metadata, tools &
approaches to support a
range of user needs*

Data Discovery
& Collaboration
EDXTM
Energy Data eXchange

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

Subsurface
Characterization

SUBSURFACETM
TREND ANALYSIS

Team Digital
Notebook

Subsurface
Databook

Big Data
Geoprocessing
& Analysis

hadoop

Multi-variate
Assessments

**Cumulative
Spatial
Impact
Layers**

Energy Web
Mapping Tool

**GEO
CUBE**

SWIMTM
Weighted Impact
Model

**NATCARB
viewer**

Long-term
Community Data
Management

**VARIABLE
GRID
METHOD[©]**

Spatial
Uncertainty

BLOSUMTM

4D Oil Spill
Modeling

SIMPATM
Improving
Resource
Assessments

CIAMTM

Metocean
Modeling

GLOBAL

REGIONAL

FIELD

MICRO

Scaling the Data Pyramid -

Building Solution for Common Data Challenges

Challenges scientists face in order to effectively use data resources:

Data Access:

~80% loss of published data after 20 yrs

Data Discovery:

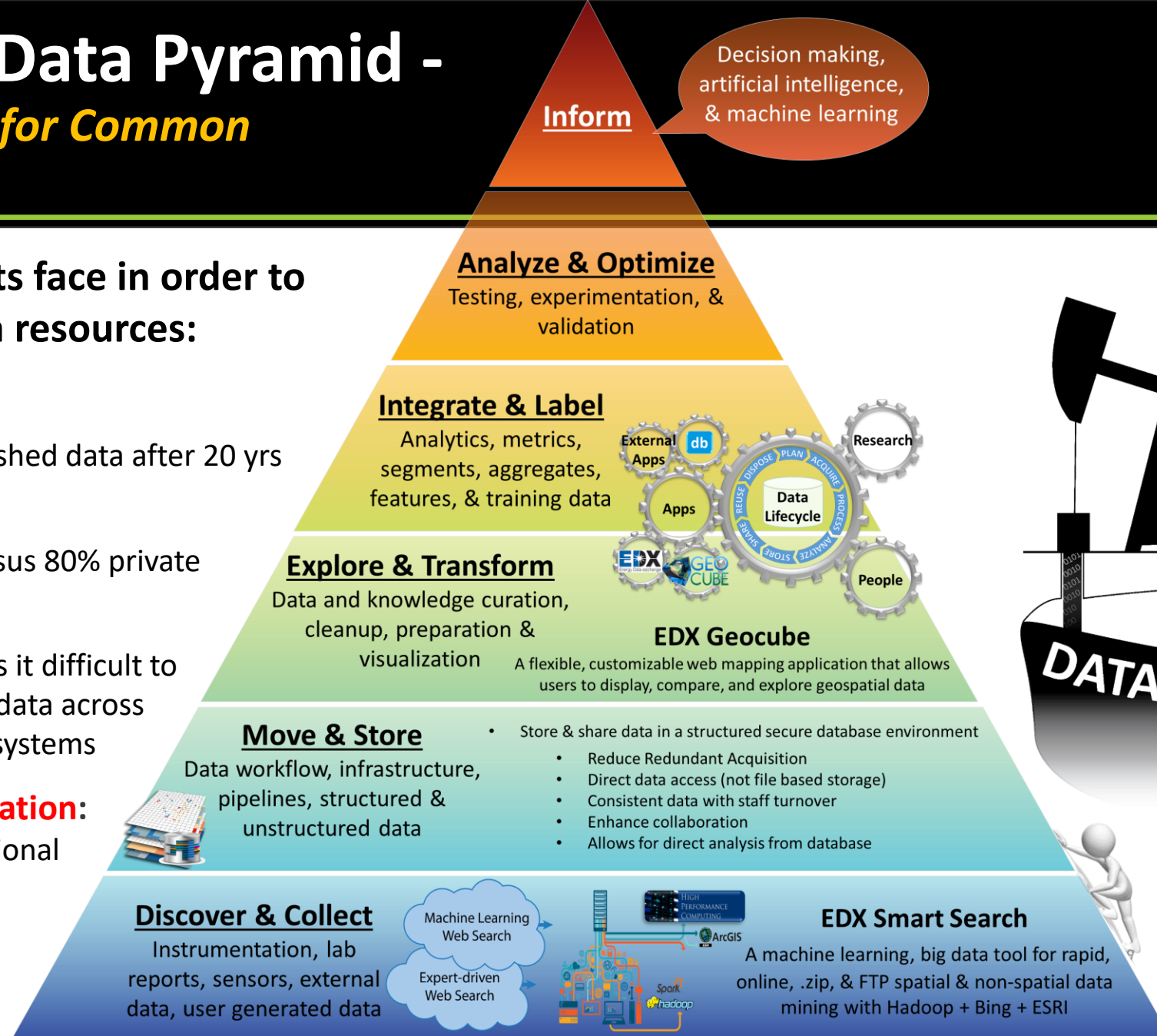
20% public data versus 80% private

Data Interoperability:

Large variety of data makes it difficult to create, exchange, and use data across different applications and systems

Data Analytics & Visualization:

Require advanced computational capabilities, algorithms, and large data stores to analyses these data





A tool that scans “seed” resources and identifies relevant keywords, then crawls the web and parses the data for integration



Acquisition of disparate data:

- 

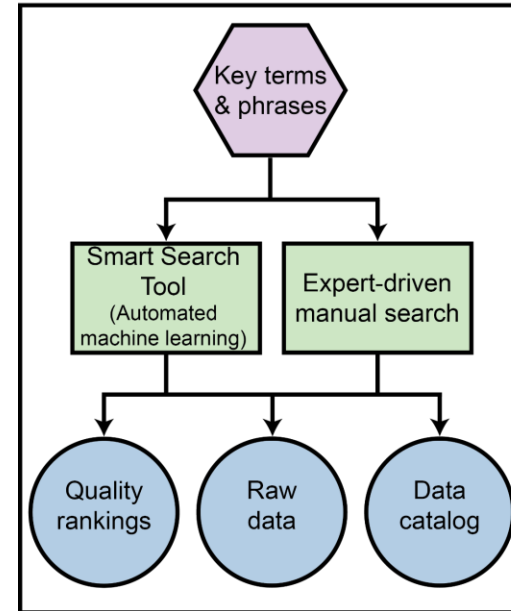
80% Dark Data

Defining a strategy up front

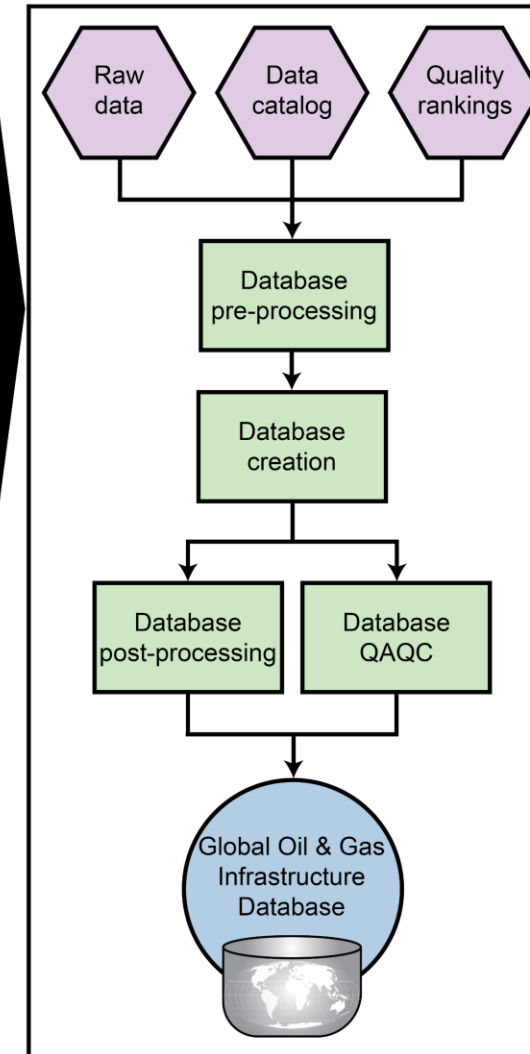
Steps for data:

1. Acquisition,
2. Integration & transformation, and
3. Analytics

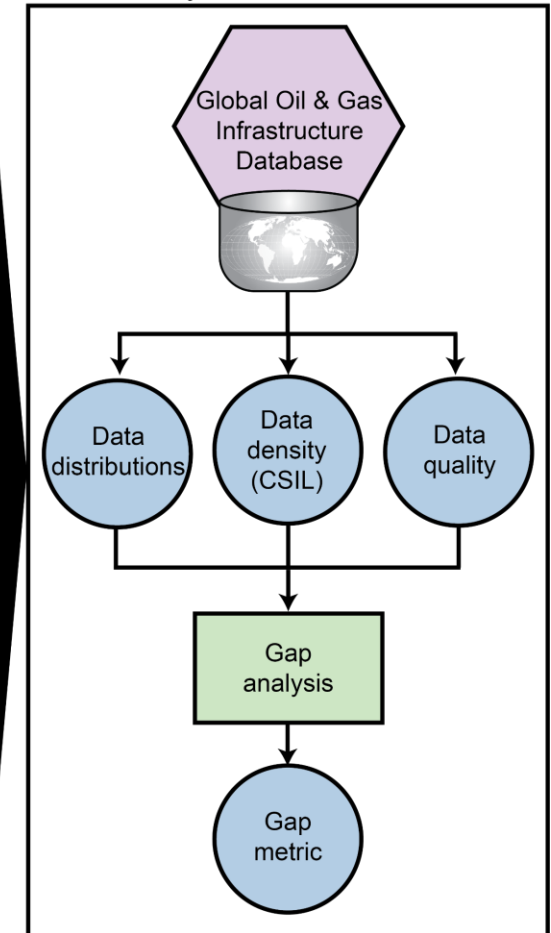
1. Data acquisition



2. Data integration and transformation

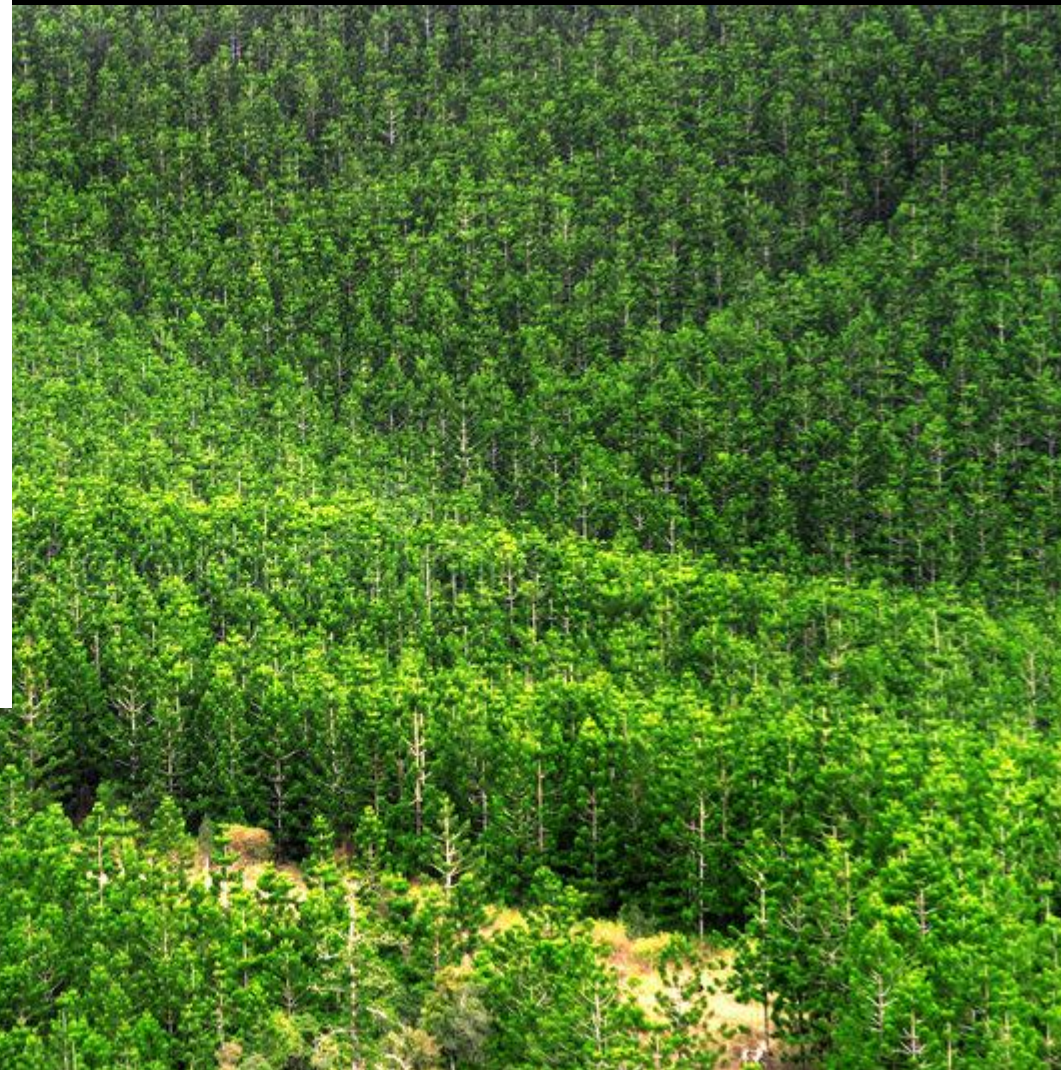
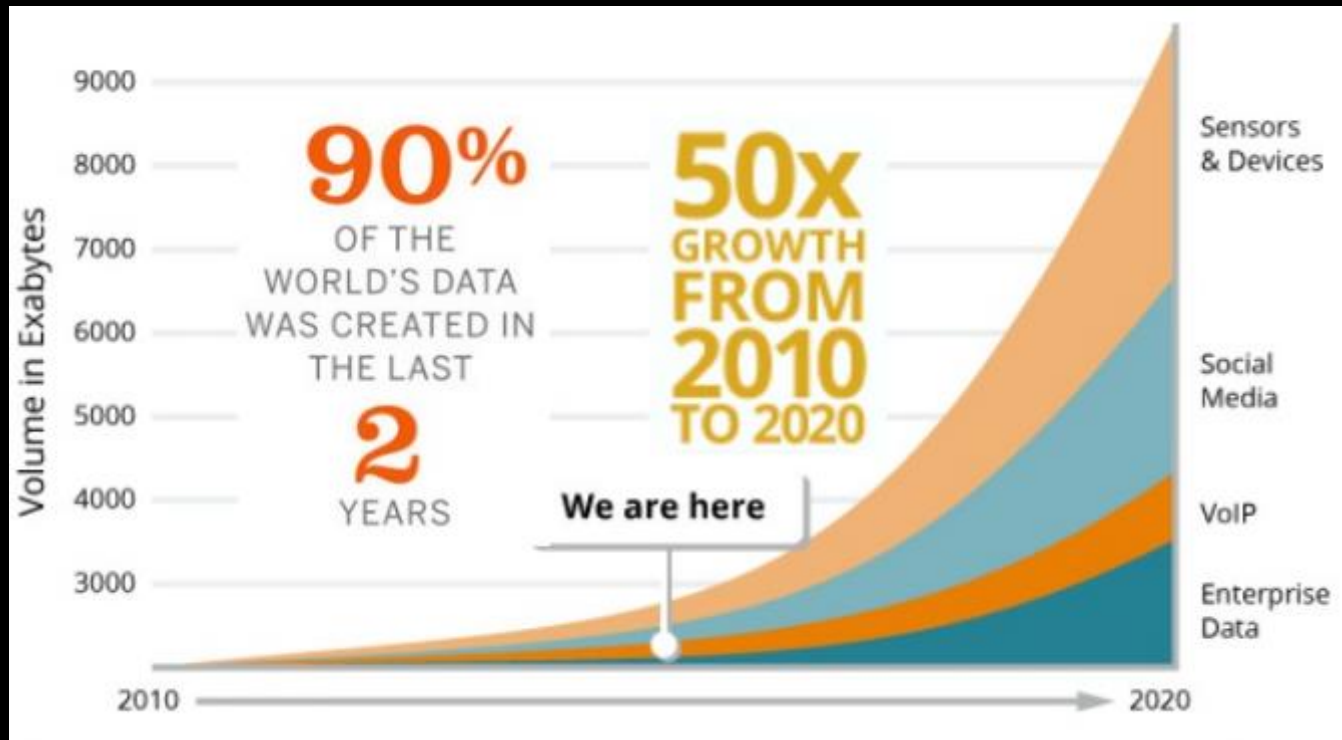


3. Data analytics



Rose et al., 2018

As access to open, authoritative data increases science driven analyses face challenges to **efficiently find**, **integrate** and **use** these resources



- Volume, variety, and velocity of data online is growing... exponentially
- How will you parse the tree from the forest?

Use Case: FTP Data Mining: Hadoop + EDX

- **Problem:**
 - Need to search data in FTP silos (millions of files, spatial and contextual)
- **Solution:**
 - Index FTP silos using Hadoop



NETL's Big Data Discovery Ecosystem (To Date)



Data Collection:

- FTP Recursion
- WWW Crawl

Data Analysis:

- Phrase Generation
- Relevance Analysis
- Geoprocessing

Metastore
(Hive, HBase)

Data Mining Clients



Custom Scripts

Python source scripts used to create, translate, and integrate the GOGI geodatabase.

- Used to check and remove for duplicates
- Conduct language translation to English
- Project spatial data
- Generate updated versions of the geodatabase

```

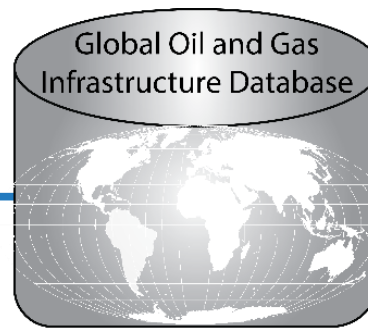
1 #This script will scan through a geodatabase and remove points with identical locations but keep the point with the higher Quality Rating
2 print("Loading arcpy please be patient...")
3 import arcpy
4 import numpy
5 from arcpy import env
6 arcpy.env.overwriteOutput = True
7 arcpy.env.workspace = "C:\\temp\\New Folder\\Final\\V6_Global_Oil_and_Gas_DB_fixpipe.gdb"
8 datasets = arcpy.ListDatasets('','Feature')
9
10 fclass = []
11 allfc = []
12 selAtt = "TempFC_FEAT_SEQ"
13
14
15 for dataset in datasets:
16     fclass.append(arcpy.ListFeatureClasses(feature_dataset=dataset))
17
18 for featureclass in fclass:
19     for feature in featureclass:
20         allfc.append(feature)
21
22 for fc in allfc:
23     print "Processing... ",fc
24     arcpy.MakeFeatureLayer_management(fc,"temp_layer")
25     rows = arcpy.UpdateCursor("temp_layer", "", "", "", selAtt+" D") #use this code to sort descending
26     arcpy.AddField_management("temp_layer", "Remover", "SHORT", "", "")
27     aveHigh = 0.0
28     rowPrev = long(0)
29     for row in rows:
30         if row.QualityRating > rowPrev:
31             row.Remover = 1
32             rowPrev = row.QualityRating
33             row.Update()
34         else:
35             row.Remover = 2
36             row.Update()
37     del rowPrev
38     del row
39     arcpy.DeleteFeatureLayer_management("temp_layer")
40     print "Done with ",fc
41
42 #####Universal attribute Translator
43 #This script lists the attributes of a selected data table and then translates the selected attributes from non-English to english using the Google Translate API
44 #to use this script you must create an account with google and generate an API Key
45 import arcpy
46 import os
47 import unicodedata
48
49 #generate a google api key and place it here
50 googleapikey = 'Place google api key here'
51
52 #os and file setup
53 gdbFiles = []
54 for path, dirs, files in os.walk(current_file_dir):
55     for f in files:
56         if f.endswith(".gdb"):
57             gdbFiles.append(f)
58
59 for i in range(len(gdbFiles)):
60     print i, " ",gdbFiles[i]
61
62 #ask user which GDB to process
63 numselect = int(raw_input("Type in the number of the GDB you want to process: "))
64 catalog = current_file_dir+"\\GDB"+gdbFiles[numselect]
65 print catalog+" selected"
66
67 arcpy.env.workspace = catalog
68 datasets = arcpy.ListDatasets()
69 #List available datasets
70 for i in range(len(datasets)):
71     print i, " ",datasets[i]
72
73 #ask user for input to dataset
74 numselect2 = int(raw_input("Type in the number of the dataset you want to process: "))
75 catalog2 = catalog+"\\Dataset"+datasets[numselect2]
76 print catalog2+" selected."
77 arcpy.env.workspace = catalog2
78
79 #####
80 #This script Reads the data attributes and data location from a CSV file and Populates a Geodatabase with the spatial files found.
81 #It also includes the metadata from the CSV and populates the collected data with these attributes.
82 #Also a log file will be created that lists any errors that the script encounters while running.
83 #####REQUIREMENTS
84 #User must have ARCPY installed
85 #Make sure the following files are in the same
86 #folder as this script:
87 #feature_classes.csv
88 #WGS1984.prj
89 #a CSV file based off the edf catalog
90 #####
91
92 import csv
93 import os
94
95 #os and file setup
96 current_file_dir = os.path.dirname(__file__)
97 print "Opening log.txt file..."
98 features = current_file_dir+"\\feature_classes.csv"
99 sr = current_file_dir + "\\WGS1984.prj" #This is a projection file that will help move all disparate data sets to a single spatial reference
100
101 #Loop to list csv files in current directory
102
103 csvFiles = []
104 for path, dirs, files in os.walk(current_file_dir):
105     for f in files:
106         if f.endswith(".csv"):
107             csvFiles.append(f)
108
109

```

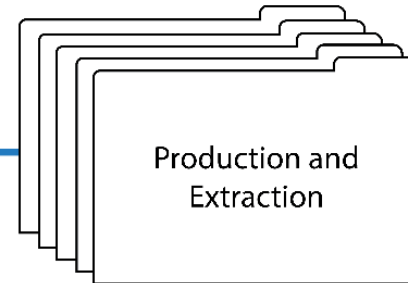

Not all data is equal

- Quality & quantity varies
- But understanding uncertainty and gaps in data is important for data driven analytics, stakeholder decision making, and other needs

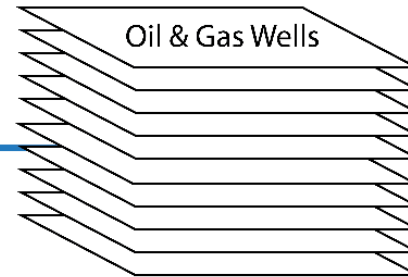
Geodatabase



Feature Datasets



Feature Classes



Features

i.e. records

ID	X	Y
1	80	-111
2	34	65
3	-60	92
4	-49	35
...

Feature Attributes

i.e. fields

ID	Type	Age	Status	...
1	Oil	1990	A	...
2	Oil	1996	PA	...
3	Gas	1982	PA	...
4	Oil	2015	A	...
...

Feature Datasets

Transport

Example Feature Classes

Ports, Railways, Pipelines

Facilities/
Installations

LNG, Power Plants, Processing Plants,
Refineries, Stations, Storage, Terminals

Production/
Extraction

Oil and Gas Fields, Platforms, Well Pads,
Underground Storage, Mines, and Wells

Geology

Sedimentary Basins

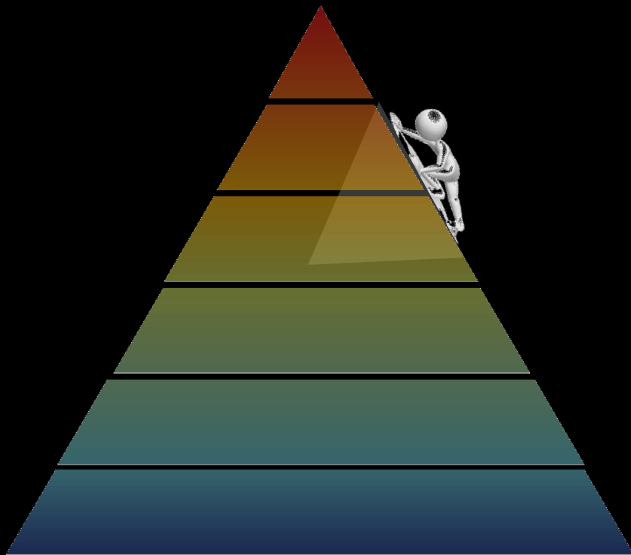
Overall quality score



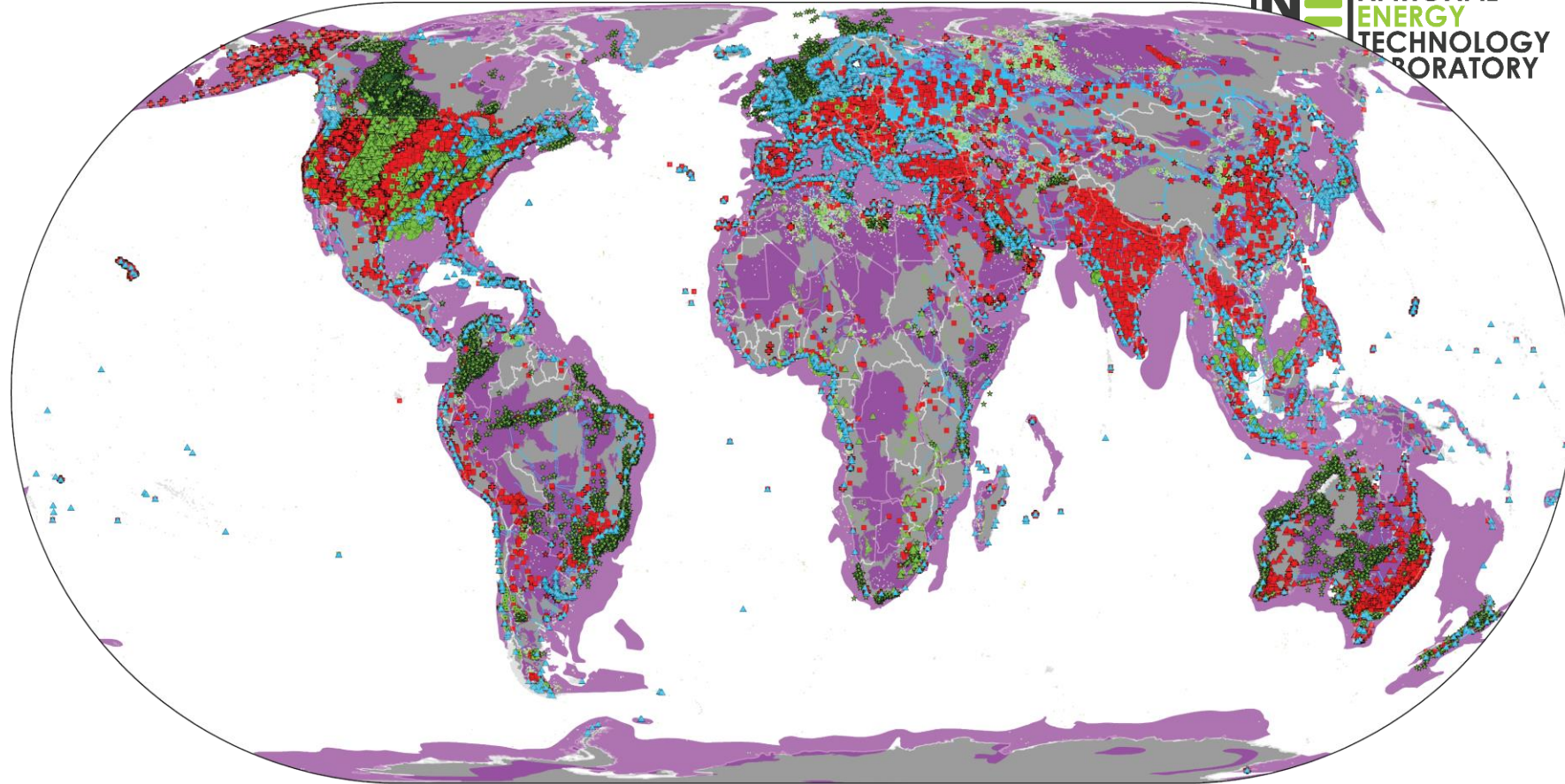
Ranking data quality by
source, spatial and
temporal features helps
with analytics
















Visualization

Tip of the pyramid



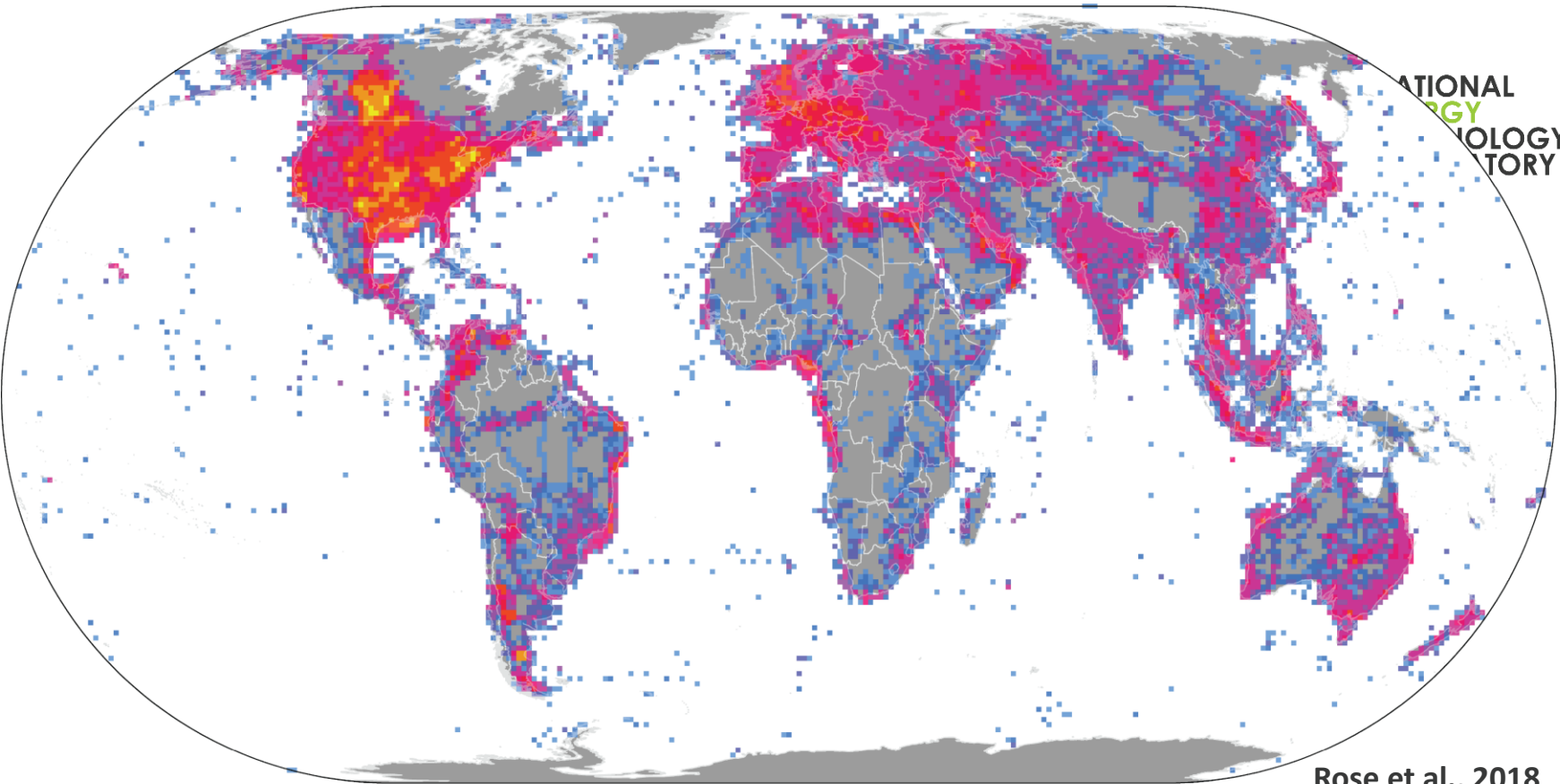
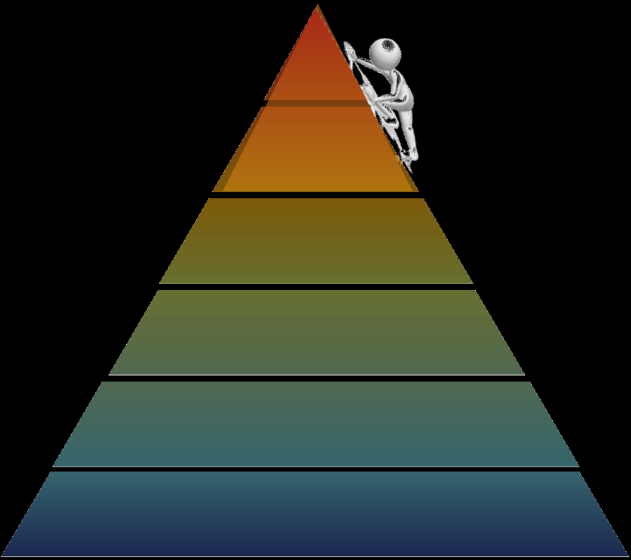
Rose et al., 2018



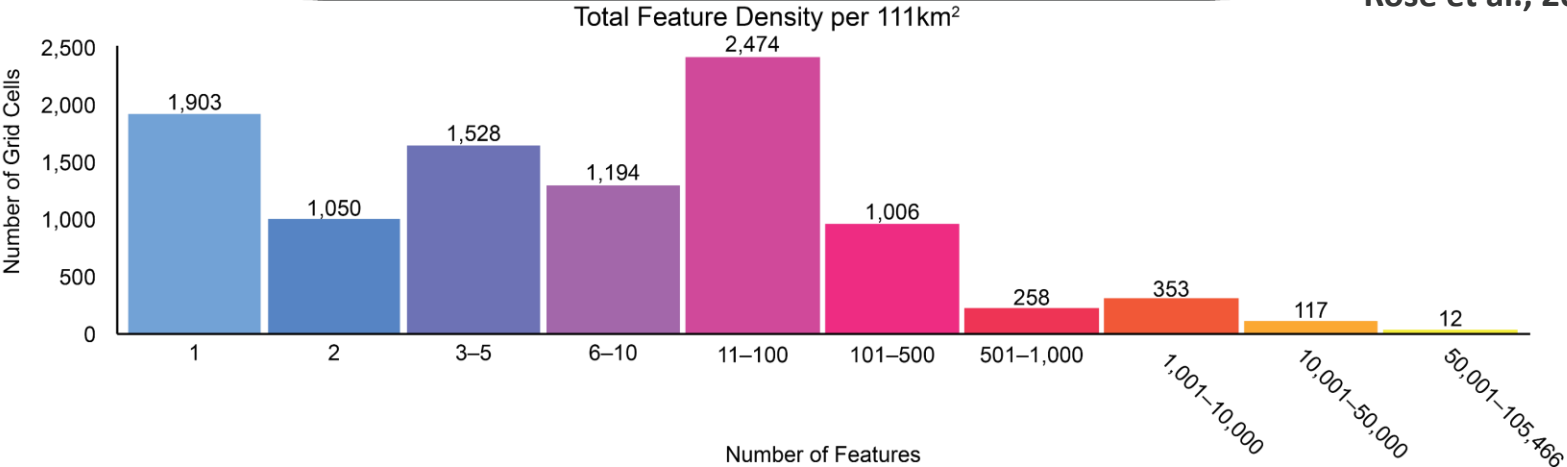
Feature Datasets	Feature Class Symbology and Total Number of Features		
Transport	Ports  n = 3,702	Railways  n = 280,734	Pipelines  n = 94,448
Facilities and Installations	Liquid Natural Gas (LNG)  n = 329	Processing Plants  n = 1,922	Stations  n = 13,876
	Power Plants  n = 14,097	Refineries  n = 2,272	Storage  n = 26,103
Production and Extraction	Fields  n = 25,236	Underground Storage  n = 3,731	Wells  n = 736,476
	Platforms and Well Pads  n = 9,845	Mines  n = 51,602	Well Density 25km ² Grid (represents a total of n = 3,544,809 wells)
Geology	Sedimentary Basins  n = 1,046		

Analytics

Tip of the pyramid

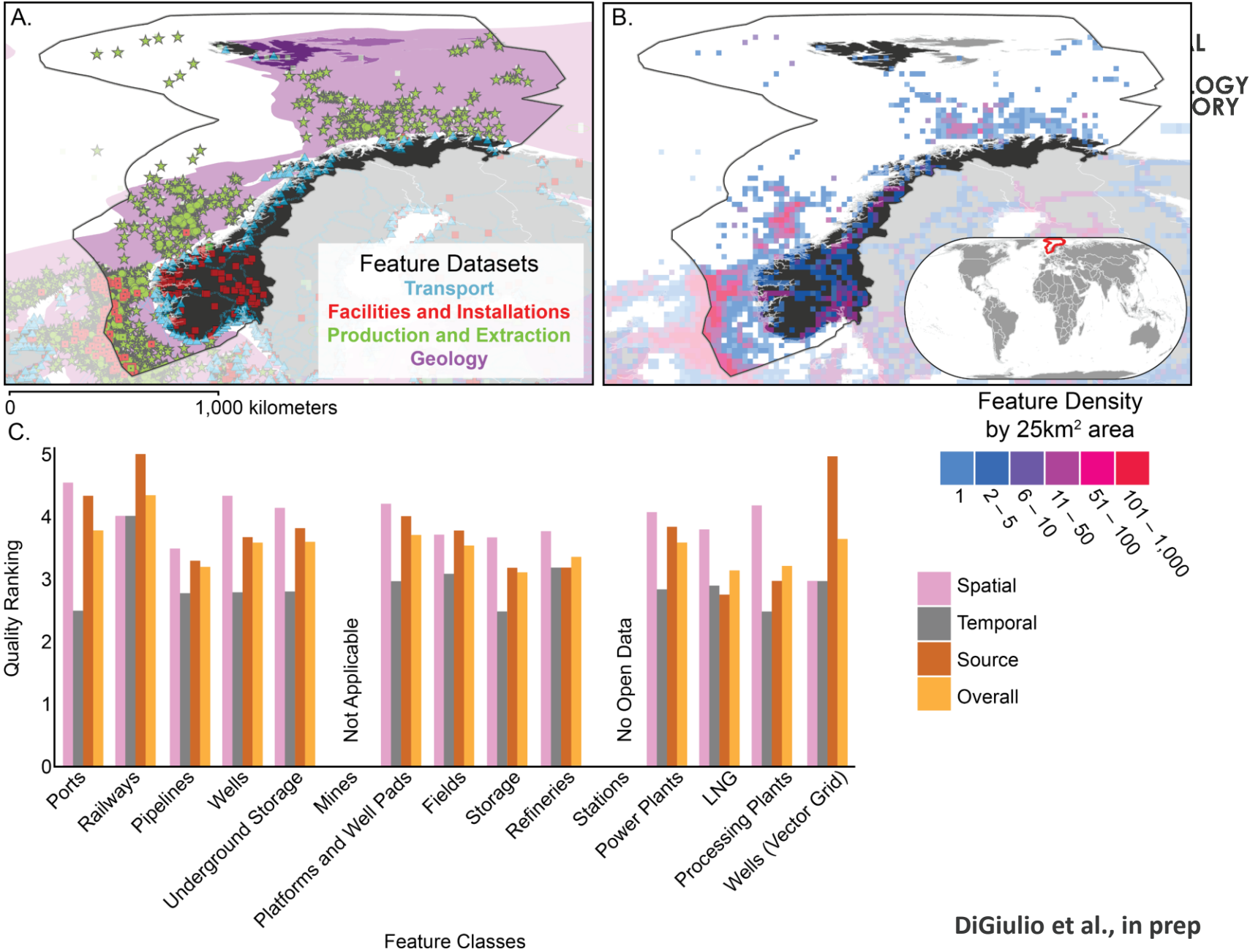
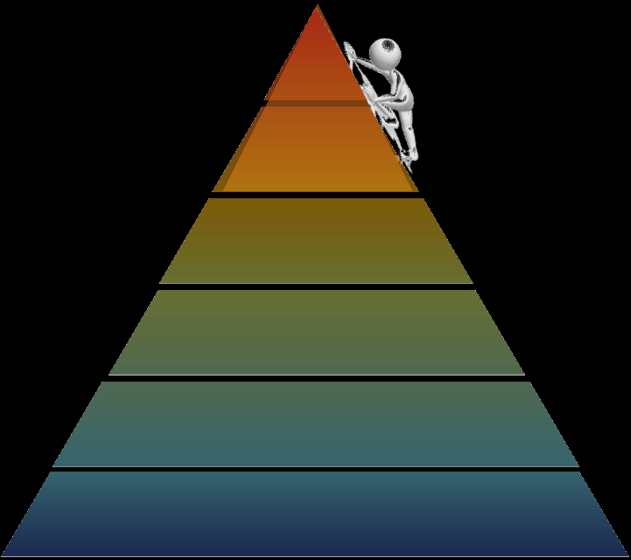


Rose et al., 2018



Analytics!

Tip of the pyramid









Embrace the uncertainty and error in data

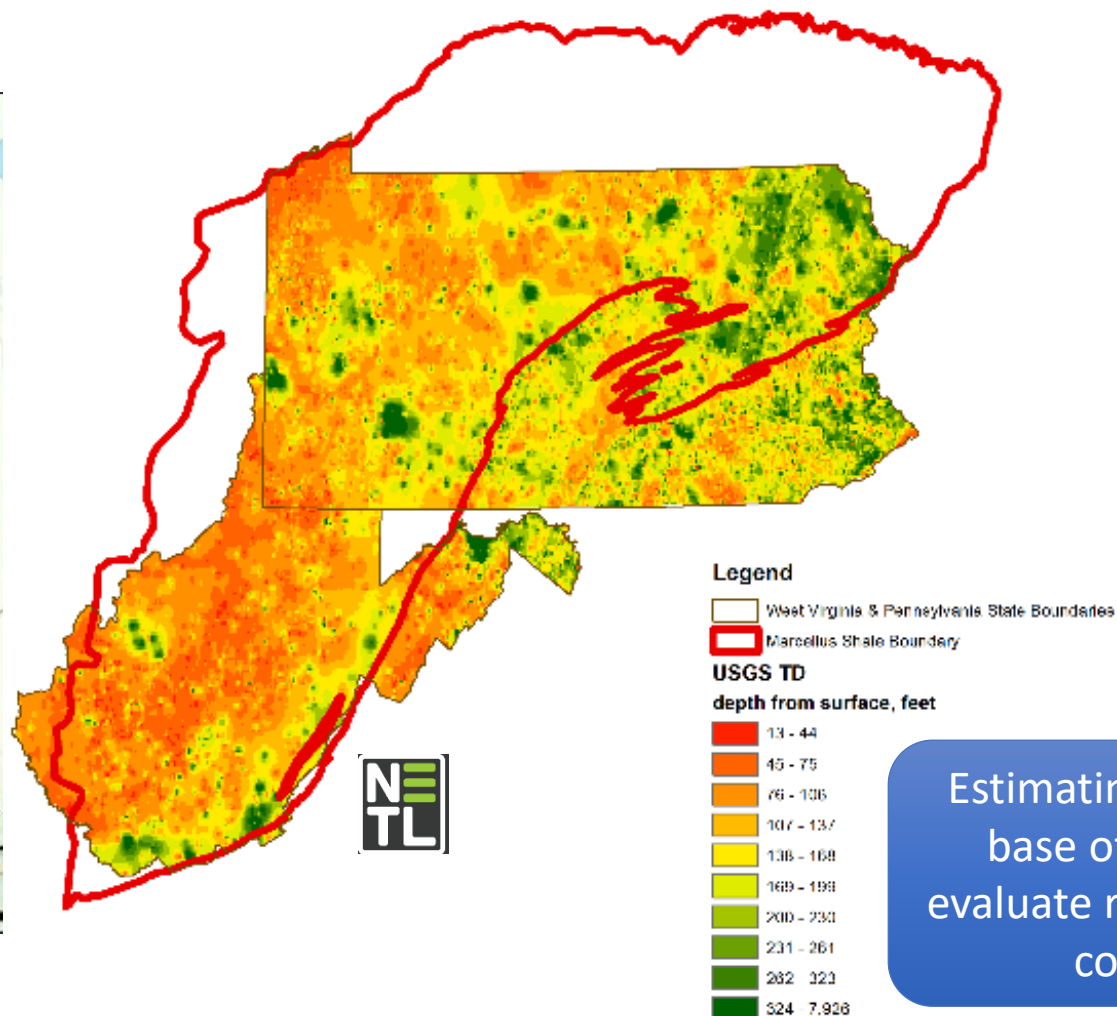
Spatio-temporal data uncertainty information is often **lacking due to difficulties** encountered:

- from the variety of potential sources and definitions,
- visualizing uncertainty, and
- communicating results

Failing to effectively communicate underlying uncertainty can lead to *false conclusions* and *poor decisions* as well as *affect the quality* of current and future research and products

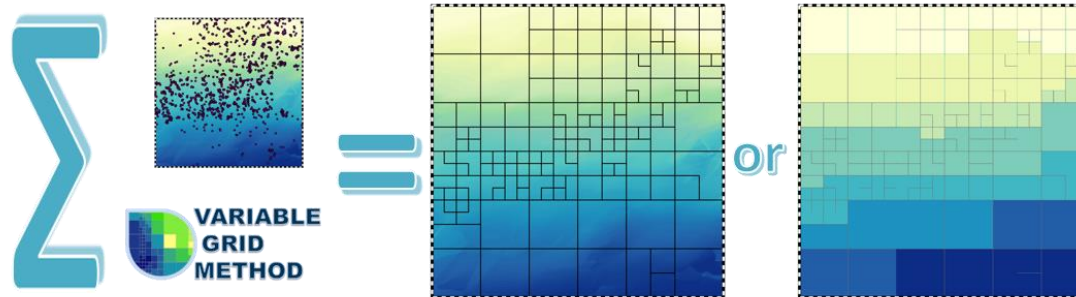
8 PM Sun, Jul 19		75°	75°	Cloudy	0%
9 PM Sun, Jul 19		74°	74°	Cloudy	0%
10 PM Sun, Jul 19		74°	74°	Cloudy	20%
11 PM Sun, Jul 19		74°	76°	Mostly Cloudy	10%
12 AM Mon, Jul 20		74°	75°	Mostly Cloudy	5%
1 AM Mon, Jul 20		73°	75°	Partly Cloudy	5%
2 AM Mon, Jul 20		73°	75°	Cloudy	5%
3 AM Mon, Jul 20		73°	74°	Cloudy	5%
4 AM Mon, Jul 20		73°	74°	Cloudy	5%
5 AM Mon, Jul 20		73°	74°	Cloudy	5%
6 AM Mon, Jul 20		72°	74°	Cloudy	5%

Example, why uncertainty matters



Estimating the depth to the base of groundwater to evaluate risks of groundwater contamination

Allow for simultaneous visualization & quantification of spatial data and uncertainty



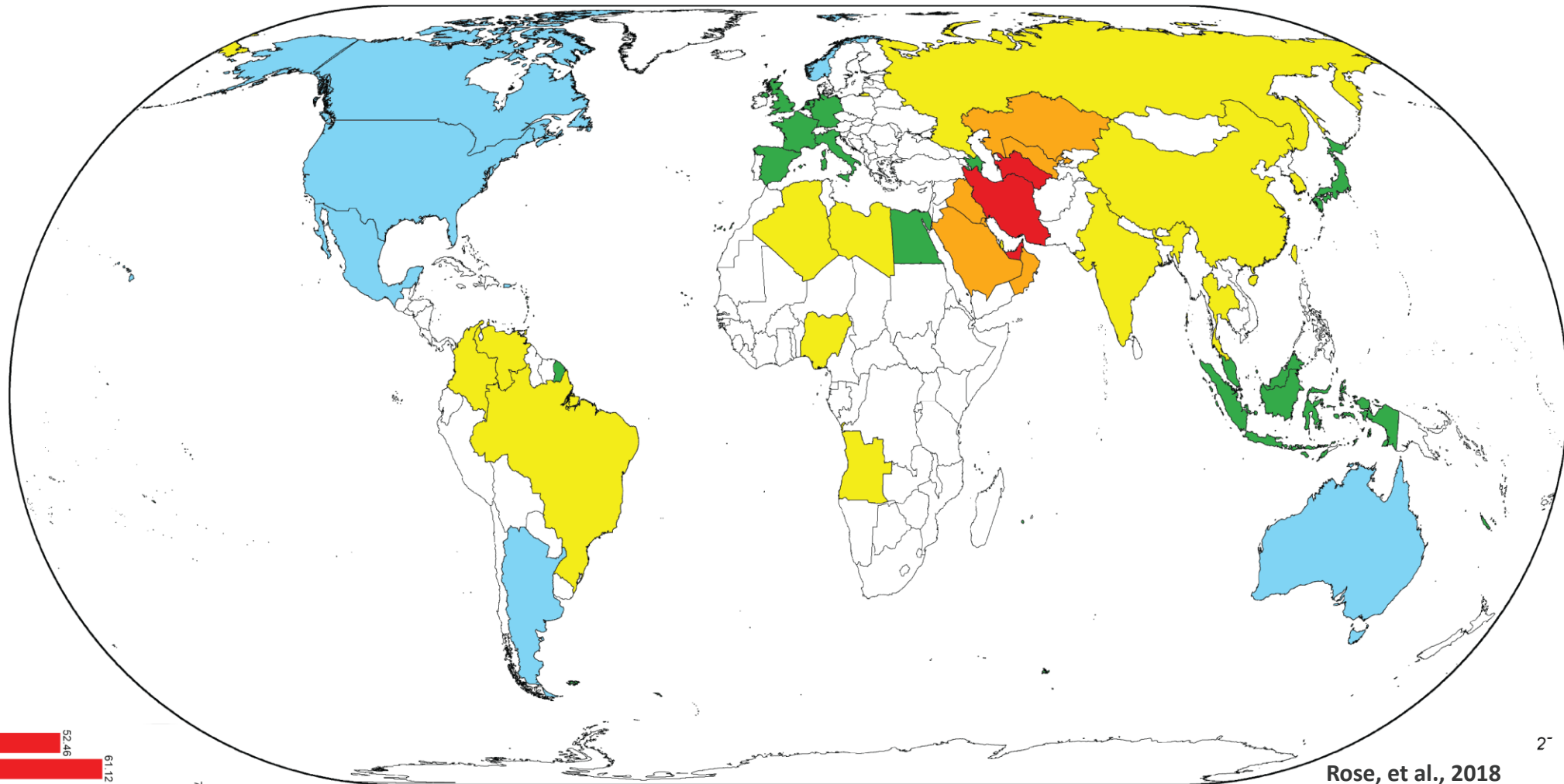
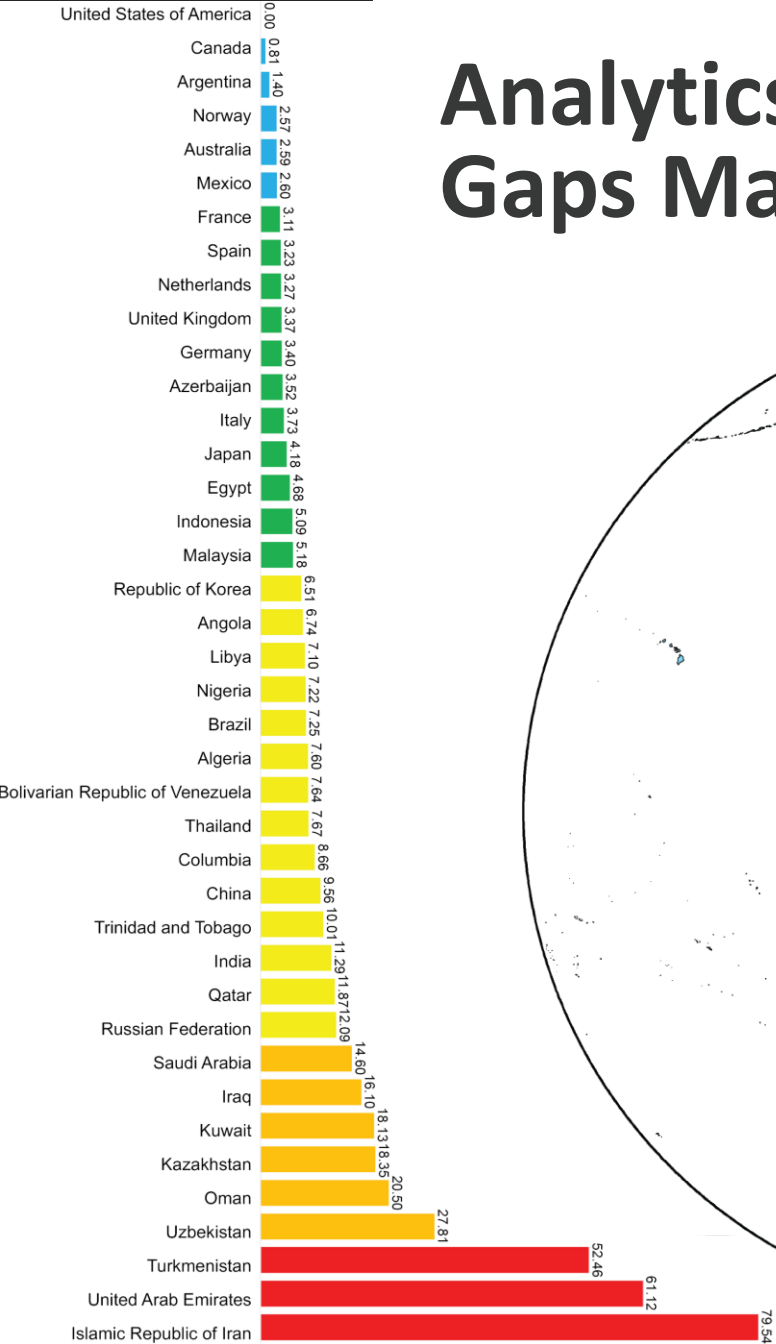
Bauer & Rose, 2015

Communicate data (via colors) and uncertainty (via grid cell size)

Uncertainty Viz/Quant for Spatio-Temporal Analyses Can Improve:

- Resources evaluations
- Impact assessments
- Understanding trends in the data
- Calculating Project Feasibility
- Identifying Knowledge Gaps

Analytics – Gaps Matter



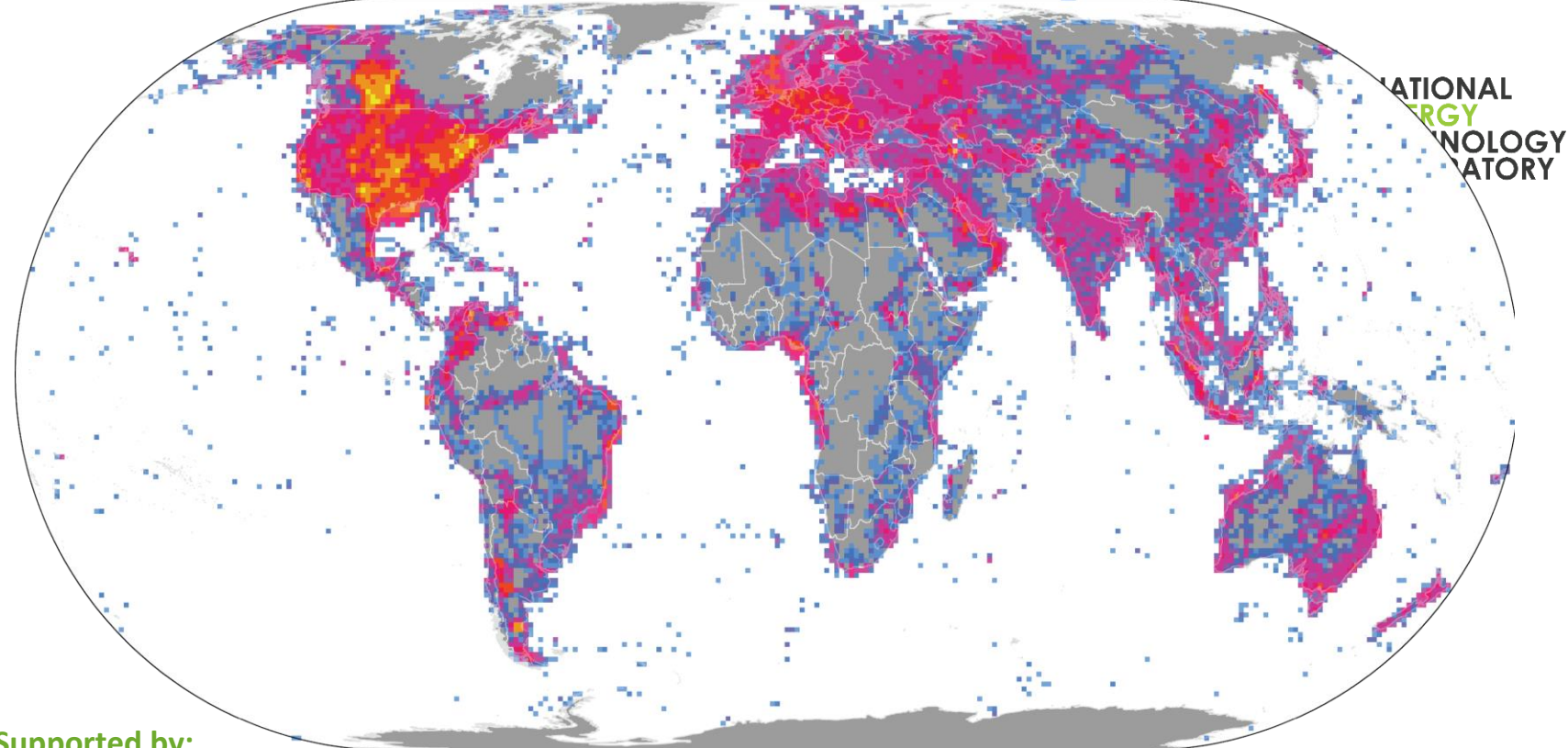
Rose, et al., 2018

27

Global Oil and Gas Infrastructure (GOGI)

- 4 month performance period
- Acquisition of disparate data by country, region, & continent totaling:
 - ~800 datasets
 - 4 million+ of features
 - Attributes some regions/features

EDF used compiled database to inform decision-making about methane emissions



NATIONAL
ENERGY
TECHNOLOGY
LABORATORY

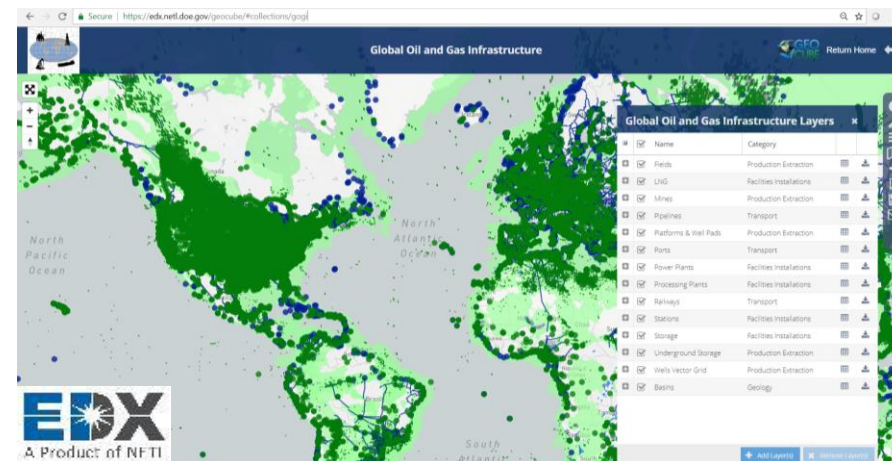
Supported by:



Total Feature Density per 111km²

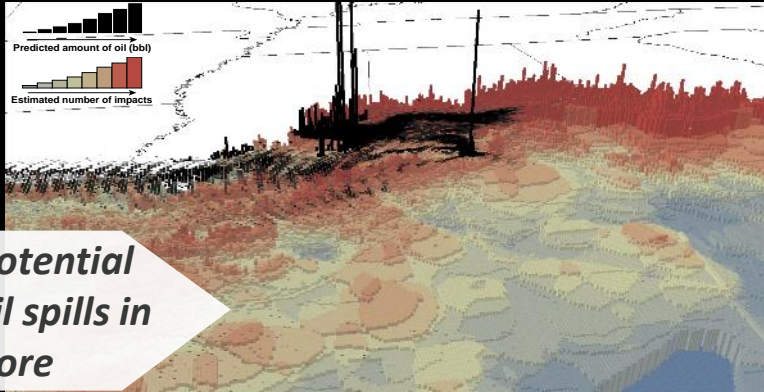
Public Products:

- [Technical report](#)
- [GOGI Database](#)
- [Web mapping application](#)
- Journal pub in prep: Digiulio et al., in prep, *Elementa*

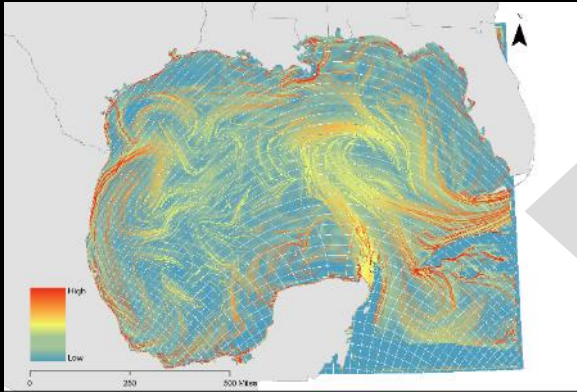


Array of Applications

- Adapted to work with other approaches, tools, and models
- Many data formats
- Multi-scale

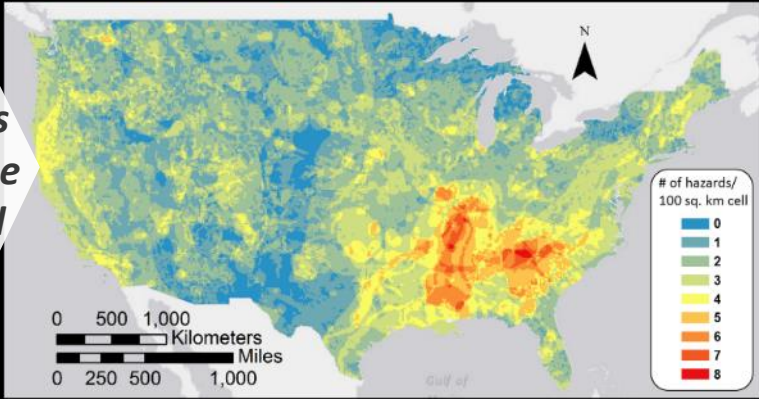


Evaluating potential impacts of oil spills in the US Offshore

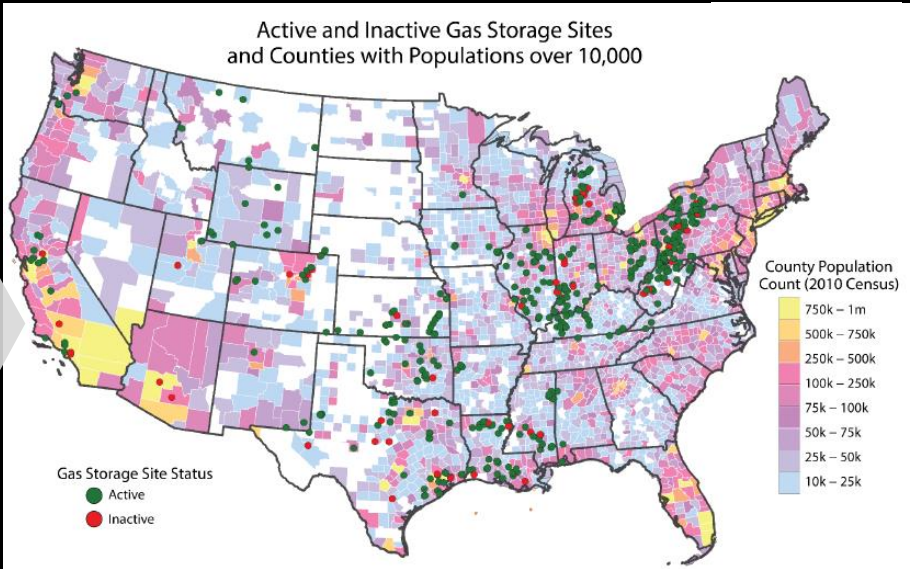
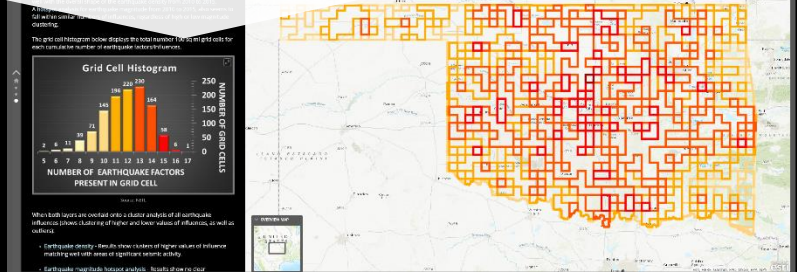


Predicting where oil is likely to go based on oceanographic data

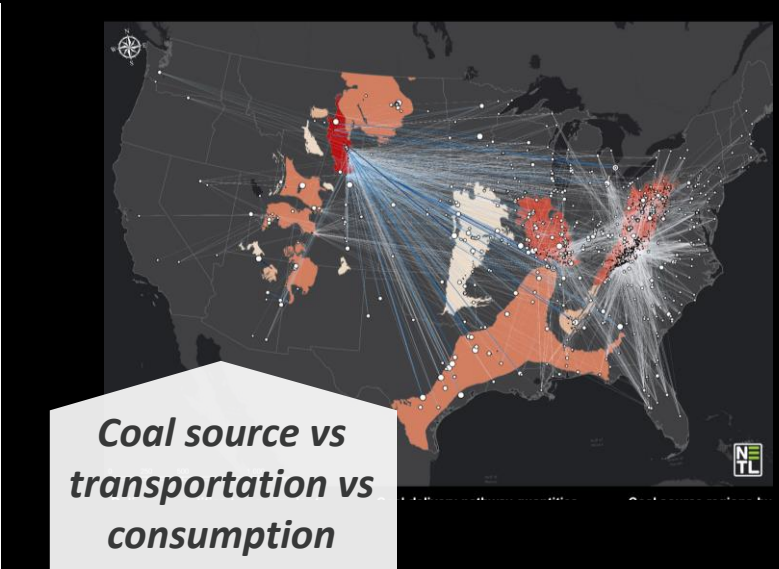
Quantifying pipeline risks from extreme weather and geohazards



Forecasting induced seismicity risk in Oklahoma

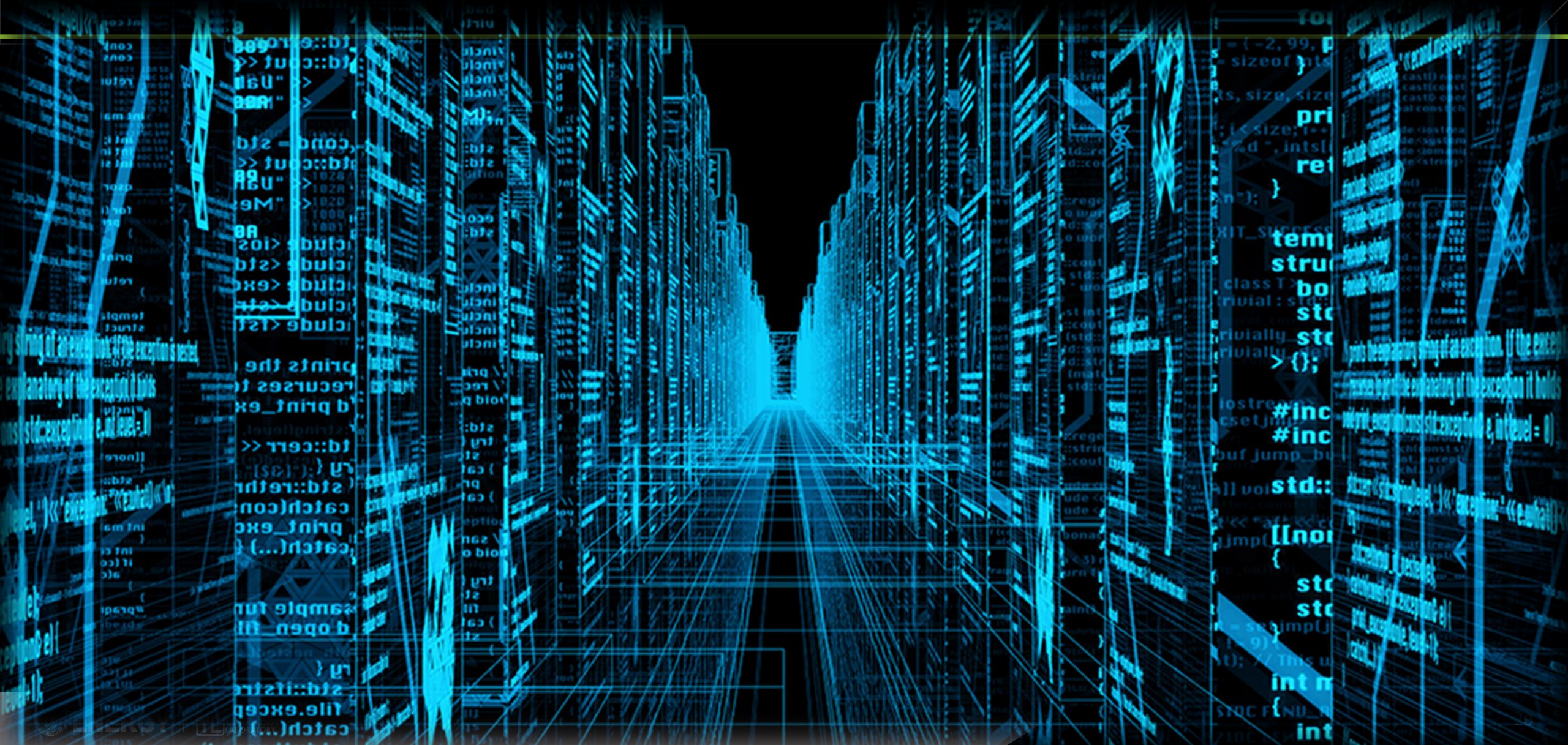


Characterizing gas storage vs population trends



Coal source vs transportation vs consumption

2018 Approach to Publishing R&D



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

A Virtual Library & Laboratory for Energy Science

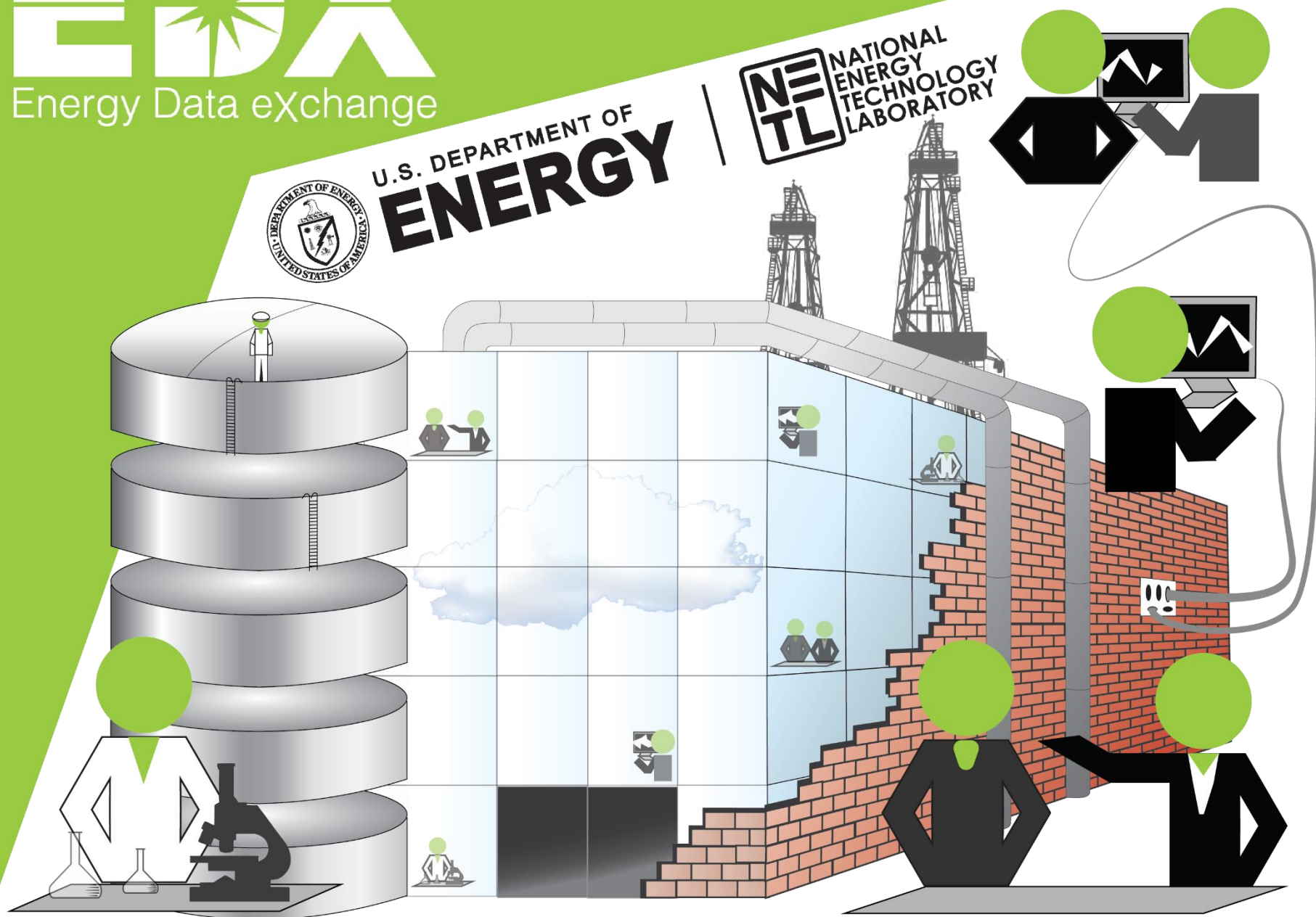
- Virtualizing team analytics
- Continued innovations to connect DOE FE affiliated researchers to online resources (tools, data, etc)
- Publishing data products from FE R&D for public reuse
- A virtual lab/user facility for FE R&D teams collaborate, analyze, and utilize data
- In development since 2011

EDX
Energy Data eXchange



U.S. DEPARTMENT OF
ENERGY

NETL NATIONAL
ENERGY
TECHNOLOGY
LABORATORY



Numerous Data Repositories

Offers opportunities and challenges

<https://www.dataquest.io/blog/free-datasets-for-projects/>

13 SEPTEMBER 2016 / PROJECT

18 places to find data sets for data science projects

This is the fifth post in a series of posts on how to build a Data Science Portfolio. You can find links to the others in this series at the bottom of the post.

If you've ever worked on a personal data science project, you've probably spent a lot of time browsing the internet looking for interesting data sets to analyze. It can be fun to sift through dozens of data sets to find the perfect one, but it can also be frustrating to download and import several csv files, only to realize that the data isn't that interesting after all. Luckily, there are online repositories that curate data sets and (mostly) remove the uninteresting ones.



A screenshot of the NETL's Energy Data eXchange (EDX) website. The header includes the NETL logo and the text 'NETL's Energy Data eXchange'. Below the header is a navigation bar with links for Home, Search, Contribute, Groups, Portfolios, Tools, Workspaces, My EDX, About, and Help. The main content area is titled 'Search Submissions' and features a 'Spatial Search' map of North America. To the right of the map are filters for EDX, EDX (Private), OpenEI, NGDS, and NOAA. Below the map is a search bar with the text 'gogi' and a 'Relevance' dropdown. A message states '2 submissions found for "gogi"'. Below this are two submission cards: 'Global Oil & Gas Features Database' and 'Development of an Open Global Oil and Gas Infra...'. The bottom of the page has a footer with the U.S. Department of Energy and NETL logos.

Show Your Work. Share Your Work. Advance Science.

That's OPEN SCIENCE

ML, NLP, OCR and other tools to resurrect old data

Traditional Approach to Publishing R&D



<http://www.thibaudpoirier.com> Journal manuscripts, books, reports, written results

Data tools can be used to:

- Mine journal/patents other publications
- Convert tables and graphs back to data
- Gather images for analytics
- Scan and characterize documents



Vladimir Fedak [Follow](#)

CEO of IT SvIt since 2005 and don't wanna stop | DevOps & Big Data specialist
Jan 29 · 4 min read

5 Heroic Tools for Natural Language Processing

Big Data analysis is an essential tool for Business Intelligence, and Natural Language Processing (NLP) tools help process a flow of unstructured data from disparate sources.



Never miss a story from Towards Data Science, when you sign up for Medium. [Learn more](#)

[GET UPDATES](#)

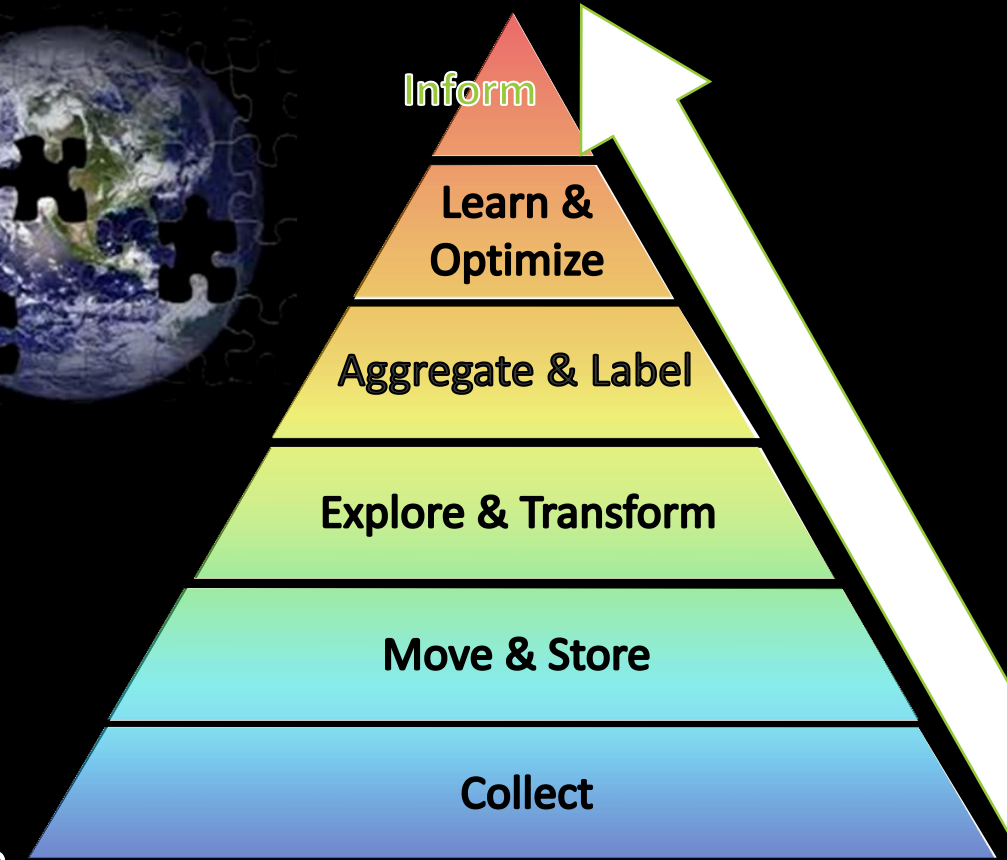
Data driven
science...

...takes a team



Additional Data Challenges & Opportunities

- **Computing science** can help address subsurface systems data challenges
- How do we balance the tug of war between **potential of data to innovate vs. stakeholder concerns?**
- **Is all data equal?** What are the data priorities?
 - Fill in the data puzzle one piece at a time...
- If you don't have the data you want, **are there proxy data** that can fill in the gaps?
- Think about **demonstrating need vs value**
- Error and **uncertainty** are important
- **Incentives** to release data, data citations, journals, and scientific community standards
- **Anonymization** and other big data computing capabilities can help unlock sensitive data to inform



Kelly.rose@netl.doe.gov