The Challenges of Curating a Research Data Set From the World's Most Complex Machine

Eric Andersen and Jeff Banning, Pacific Northwest National Laboratory

The US Department of Energy's Office of Electricity has supported research for big data analytics of phasor measurement unit (PMU) data as a step toward the development of real-time early warning tools, operator decision support tools, and potential PMU-based automated controls for the world's most complex machine: the U.S. bulk electric system. The goal of DOE's Funding Opportunity Announcement (FOA) 1861 was to explore the use of artificial intelligence (AI) tools on time-synchronized telemetry data (synchrophasors), to confirm and improve existing knowledge, and to discover new insights and tools for better electric grid operation and management. This FOA offered pre-packaged datasets to the FOA awardees to develop Al tools and capabilities. Pacific Northwest National Laboratory (PNNL) was responsible for acquiring the electric utility data to be provided to awardees of the FOA call and to provide support with the execution of the projects and their outcomes.













Identify **Utilities** Establish NDA's

Acquire Data & **Event Logs**

Load Raw Data onto Cloud

Clean Up and **Anonymize Data**

Data Delivered to Awardees for Use

Path Forward:

set

set

Additional

There have been

for access to the

anonymized data

numerous requests

NDAs with the Data

Providers are being

additional use of the

synchrophasor data

If and when we are

able to provide the

set for others to use,

have to establish an

anonymized data

the data users will

NDA with PNNL

is being added,

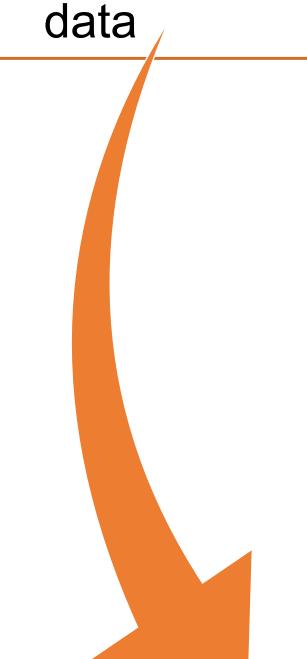
where possible

modified to enable

anonymized data

Data Request:

- Geographically dispersed utilities from each of the 3 U.S.-based electric grid interconnections
- Two years of



Interconnection

Eastern

ERCOT

Western

Challenges:

- NDAs with Data Providers
 - Legal negotiations
 - Providing assurance of data protection
 - Providing details of the anonymization process before we had access to the data

Potential Remedies:

- Get a better understanding of utility data sharing risks and mitigation strategies
- Improve the value proposition for utility participation

Dataset Range

2016-01-01 - 2017-12-31

2018-07-21 - 2019-08-24

2016-01-01 - 2017-12-31

Summary of Utility Synchrophasor Data Contributed

Total

Challenges:

- Pulling archive data was a heavy lift for some Data Providers
- Duration and age of data requested resulted in some providers having to obtain their data from third party archives
- Exported data came from a wide variety of archive tools (commercial, open-source, and custom) and were in different formats
- Event logs were all unique and lacked common taxonomy, and appear to be created manually

Potential Remedies:

- Move toward standardizing the data retention and archive processes
- Automate the creation of event logs and improve their consistency

Total

PMUs

250

221

514

Challenges:

- Inconsistency of
 - data between providers (e.g., different sequences and phases)
 - Extracted data formats
 - UTC timestamp formats
 - significant digits
- Extensive manipulation was required to aggregate data
- Data quality varied by provider
- Duplicate data
- Anonymization required to remove topology metadata

Improve

Raw Data Size

Received (TB)

38.6

10.6

19.0

68.2

- consistency across
- archive processes
- what data is archived
- UTC formats and timestamps
- Data quality

Potential Remedies:

- utilities



www.pnnl.gov



13

Number of Data

Providers (Utilities)