# Open Data on AWS

https://opendata.aws

Jed Sundwall, Global Open Data Lead

# Agenda

Overview of Open Data on AWS

How shared data on the cloud can accelerate research

Finding data shared on AWS

Sharing data on AWS

aws

# Why does AWS care about open data?

**Met Office**

Many AWS customers supply data to the public to accelerate research and product development.

**DigitalGlobe®**

Many AWS customers use data shared on AWS to create new products and services.

aws

Sharing data in the cloud lets data users spend more time on data analysis rather than data acquisition.
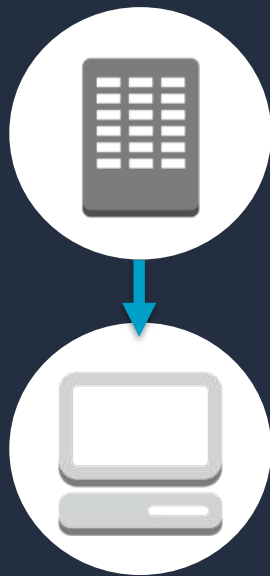
**https://opendata.aws**

aws

"…data must be organized, well-documented, consistently formatted, and error free. Cleaning the data is often the most taxing part of data science, and is frequently **80% of the work**."

— Data Driven by DJ Patil and Hilary Mason

# Undifferentiated heavy lifting

aws

# Advantages of sharing data in the cloud

**Global community of users**

**New services and tools**

**Faster pace of research**
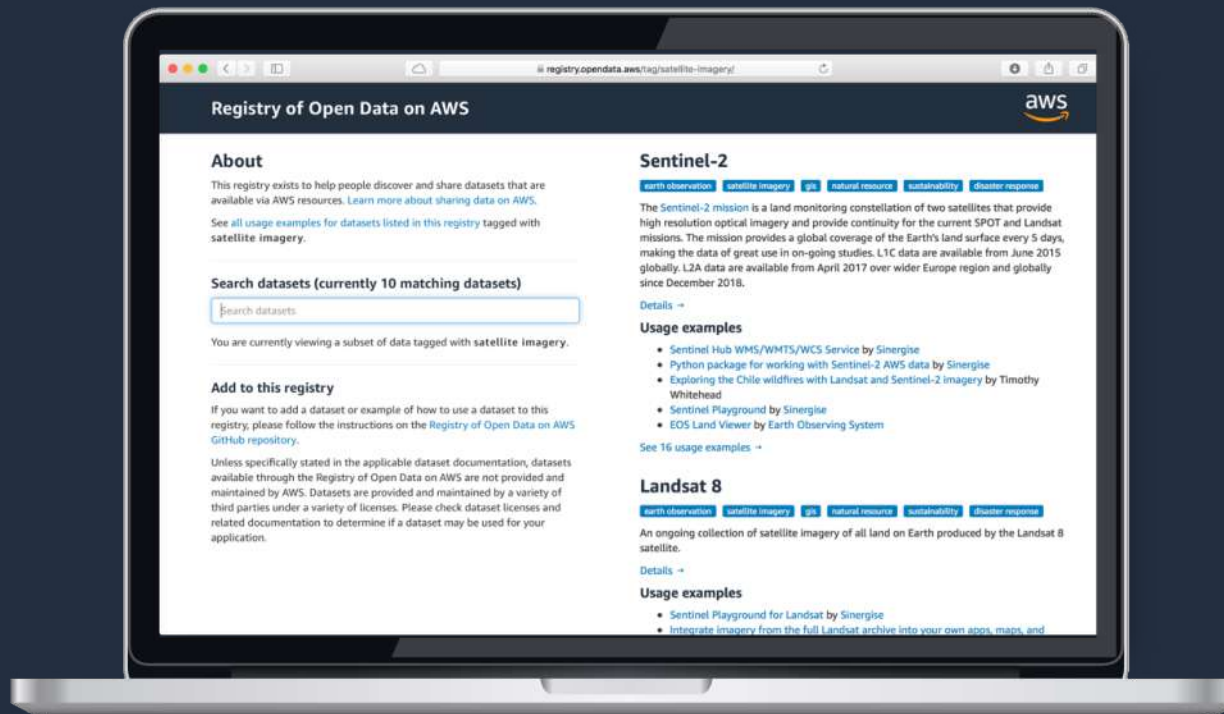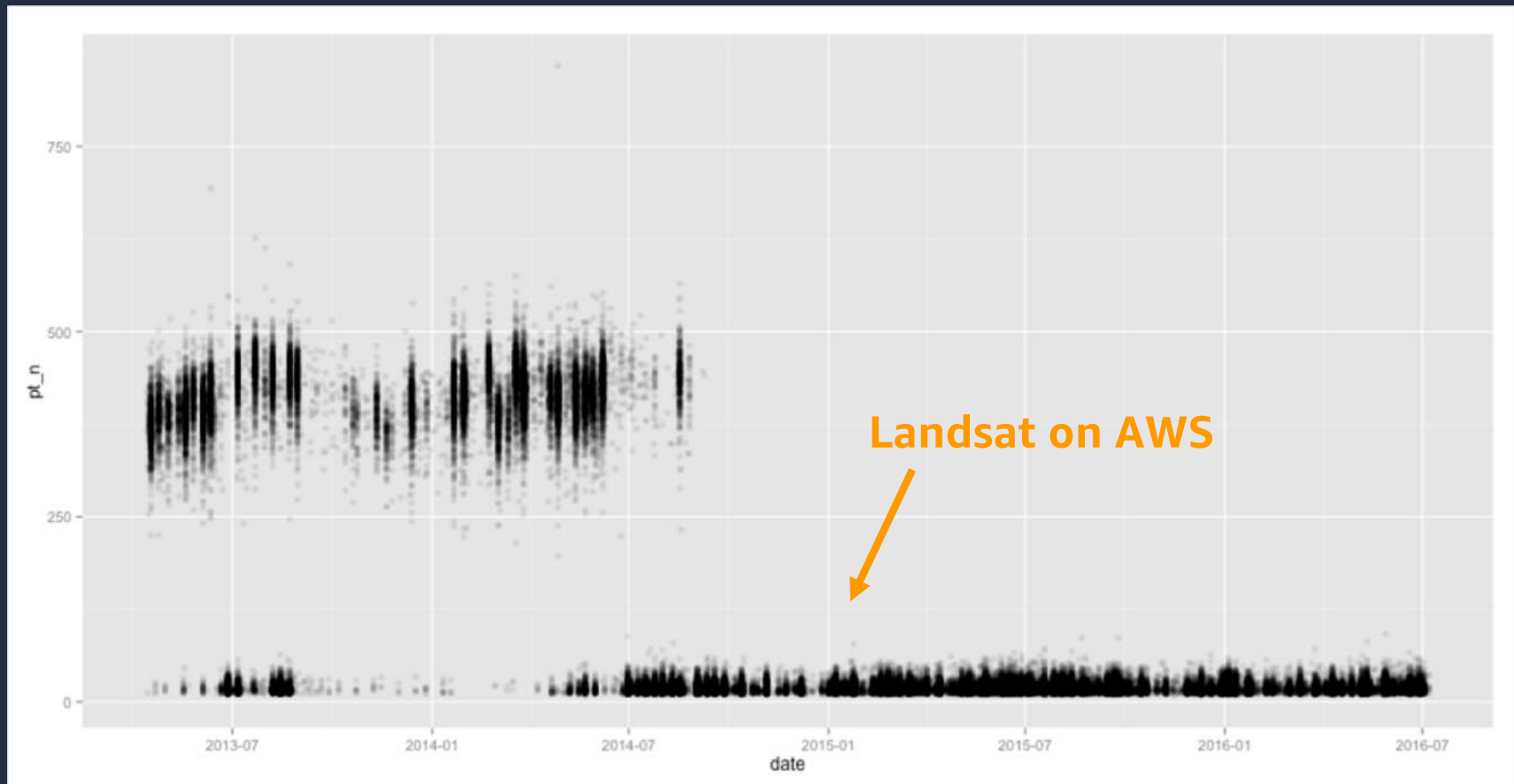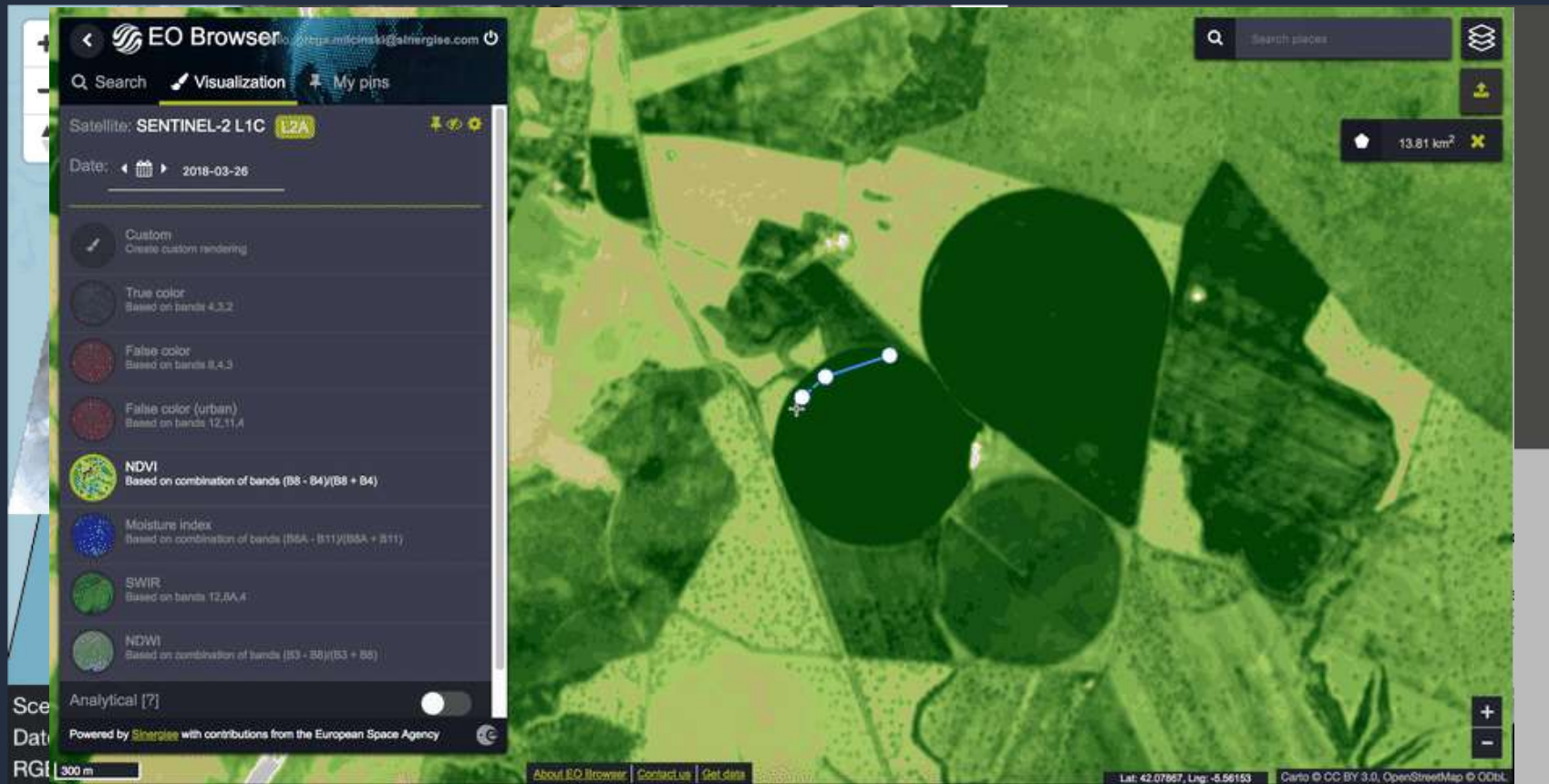
**Lower cost of research**

aws

# AWS Public Datasets

https://registry.opendata.aws

aws

# AWS Public Datasets

https://registry.opendata.aws/tag/satellite-imagery

"…data must be organized, well-documented, consistently formatted, and error free. Cleaning the data is often the most taxing part of data science, and is frequently **80% of the work**."

— Data Driven by DJ Patil and Hilary Mason

# Undifferentiated heavy lifting

aws

Graph by Drew Bollinger (@drewbo19) at Development Seed

# Using Serverless to Visualize and Analyze Imagery

aws

# Blue Dot Observatory

https://blue-dot-observatory.com/

"By having the complete Sentinel-2 archive available on AWS and access to on-demand scalable compute resources, the workflow is simple to implement. The cost to process one month of data for about 7,000 bodies of water currently in the system is 6 EUR."

# City of Louisville WARP

"The Waze Analytics Relational-database Platform (WARP) now gives 900+ government partners access to a free platform that uses real-time and historic traffic data to improve mobility, pedestrian and bike safety, road conditions, and emergency response."

aws

# Finding data on AWS

Using the Registry of Open Data on AWS (RODA)

aws

# Registry of Open Data on AWS

https://registry.opendata.aws/

# Registry of Open Data on AWS – Tags

https://registry.opendata.aws/tag/machine-learning

# Registry of Open Data on AWS – Usage examples

https://registry.opendata.aws/usage-examples

# Registry of Open Data on AWS – How to contribute

https://github.com/awslabs/open-data-registry

aws

# Sharing data (on AWS)

What we've learned

aws

# What makes a dataset successful?
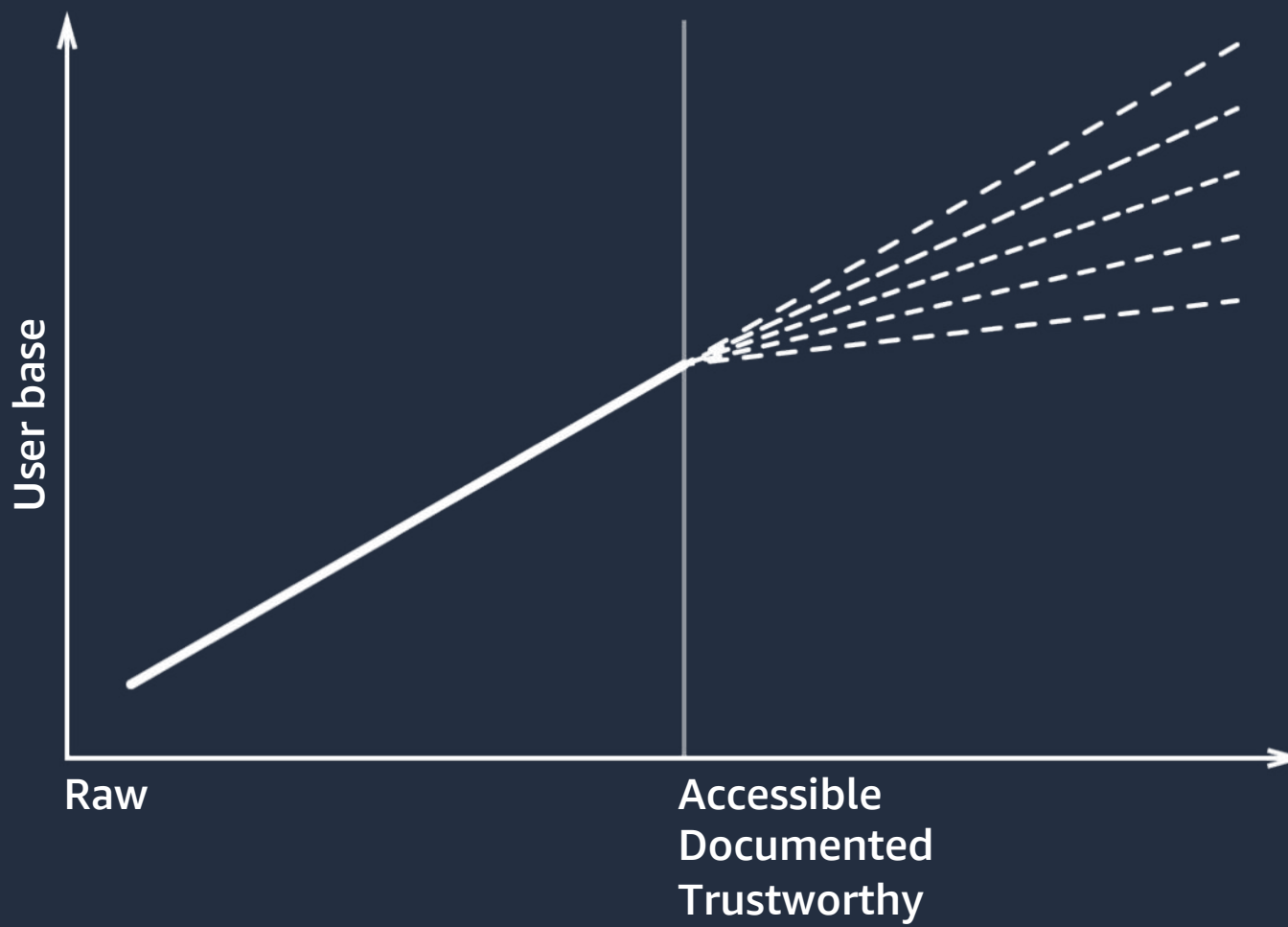## It is treated like a product.

aws

# What makes a dataset successful?
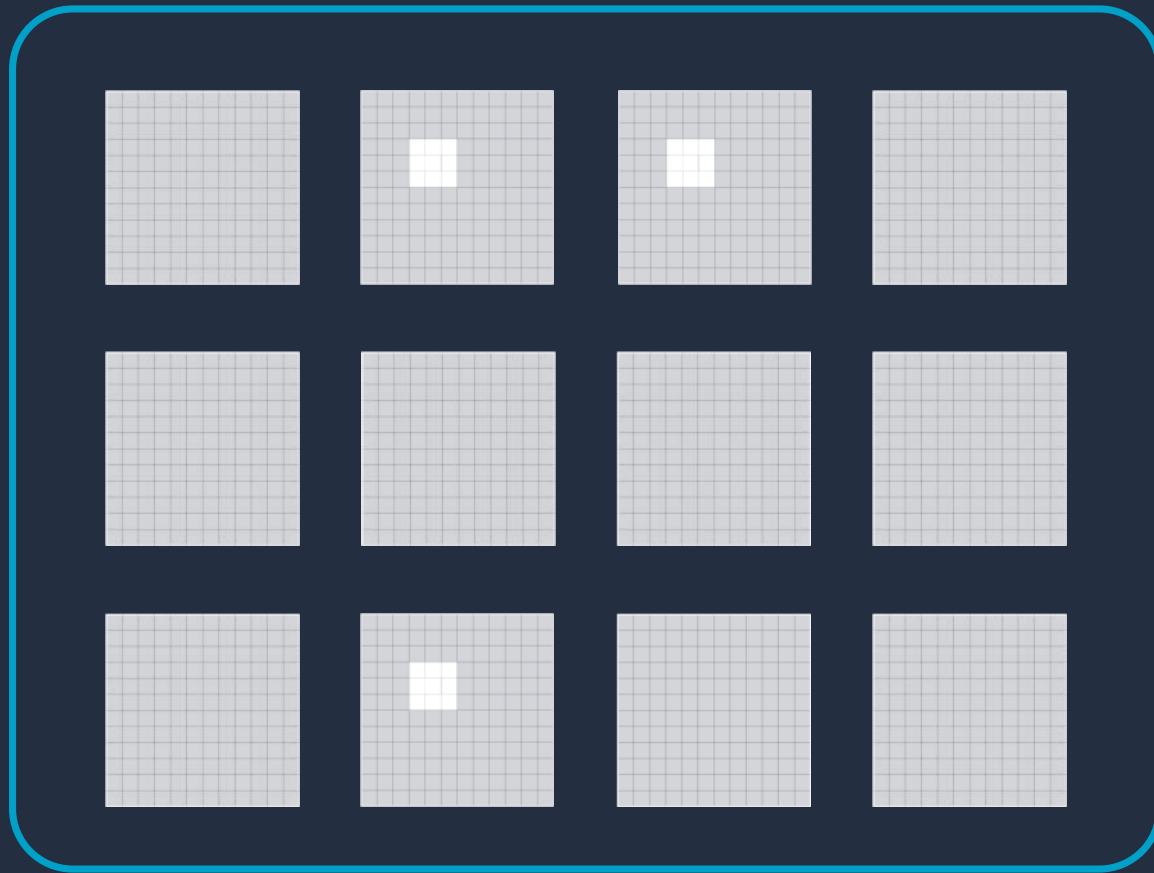## It is treated like a product.
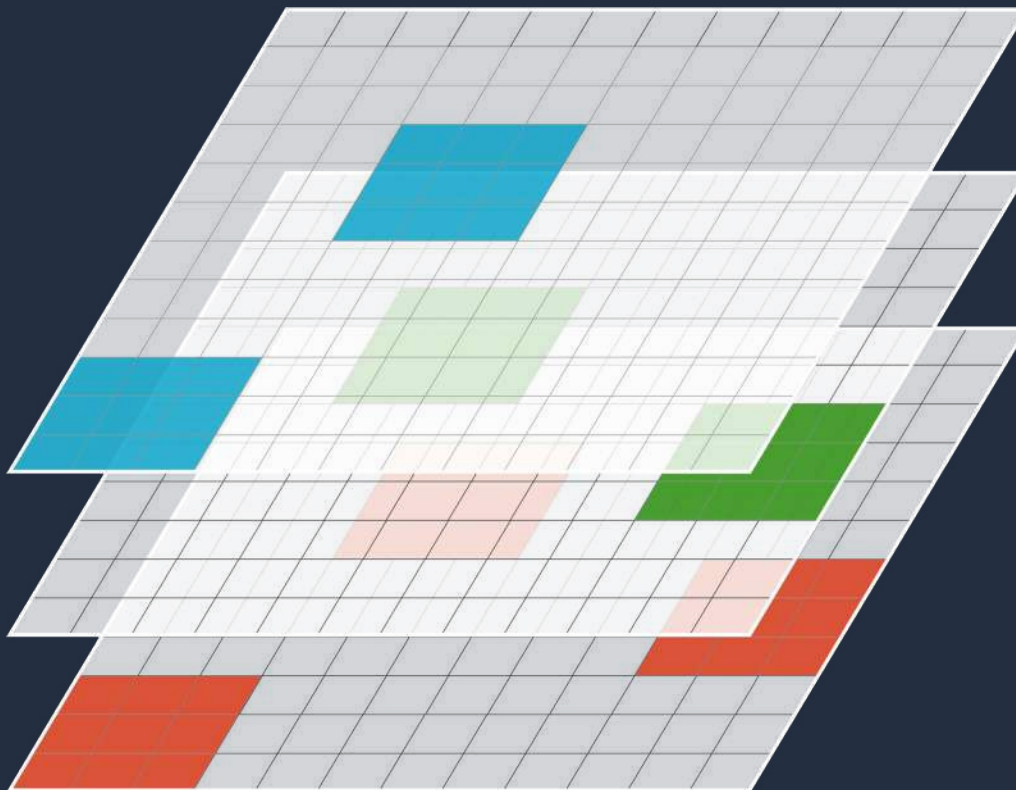## It is optimized for analysis.

aws

# The cloud-optimized GeoTIFF
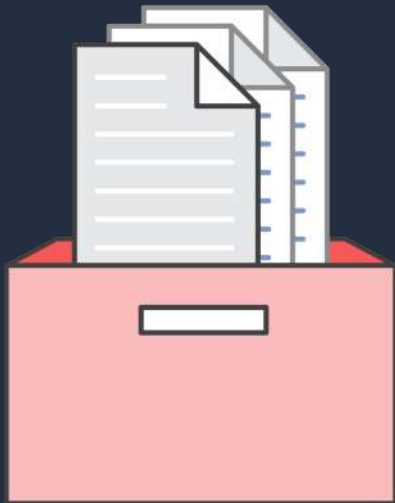


.tar

aws

# The cloud-optimized GeoTIFF

# Patterns

S3 Key Index

External
Index

Internal
Index

aws

# Example: Allen Brain Observatory Key Naming

https://registry.opendata.aws/allen-brain-observatory/

```
visual-coding-2p
├── manifest.json              # used by AllenSDK to look up file paths
├── experiment_containers.csv  # metadata for each container (area, imaging depth,
├── ophys_experiments.csv      # metadata for each experiment session
├── ophys_experiment_data      # traces, running speed, etc per experiment session
│   ├── <experiment_id>.nwb
│   └── ...
├── ophys_experment_analysis   # analysis files per experiment session
│   ├── <experiment_id>_<session_name>.h5
│   └── ...
└── ophys_movies               # motion-corrected video per experiment session
    ├── ophys_experiment_<experiment_id>.h5
    └── ...
```

aws

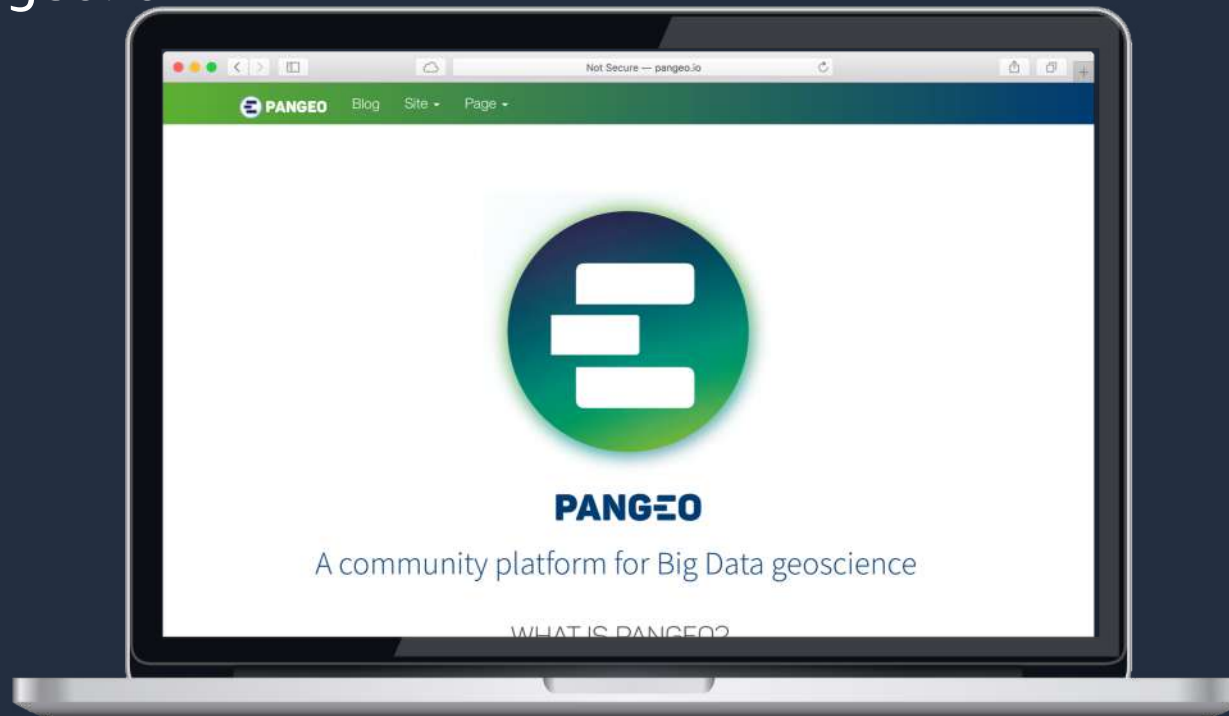# Example: IRS 990 CSV as External Index

**What makes a dataset successful?**
It is treated like a product.
It is optimized for analysis.
There is a community around it.

aws

# Pangeo community platform

http://pangeo.io

aws

# Thank you!

jed@amazon.com

aws