

Tom Gulbransen - Battelle



NEON Data Delivery & Cyberinfrastructure



neon
Proudly operated by Battelle

(mission, data products, data delivery, processing, architecture, interoperability, improvement plans)



Why is NEON important?

NEON provides a coordinated national system for monitoring a number of critical ecological and environmental properties at multiple spatial and temporal scales.

...transformative science

...workforce development



NEON's 176 Data Products Overlap Multiple Themes

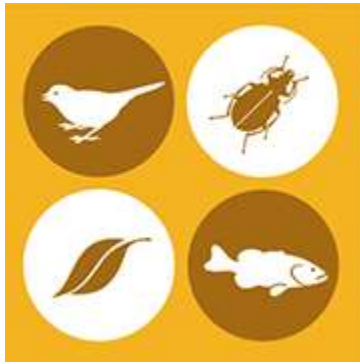
Atmospheric
58



**+11 from
AmeriFlux**

H₂O, CO₂
Heat
Isotopes
Turbulence
Storage
Fluxes

Organismal
51



**+1 from PhenoCam
+3 from MG-RAST**

Abundance
Composition
Pathogens
Phenology
DNA Barcodes
Marker Genes
Metagenomics

Ecohydrology
47



+2 from AeroNet

Water quality
Precipitation
Levels
Discharge
Radiation
Geomorphology

Biogeochemistry
85



Soil conditions
Chemistry
Particulates
Foliar characteristics

Land Related
47



Spectrometry
Hi-Res imagery
LiDAR

NEON CI & Data Science Offerings

National
Ecological
Observatory
Network

Cyberinfrastructure = resources, tools & services, observation to delivery

- of **field observational results**
- of **sampling protocols**
- of **algorithms**

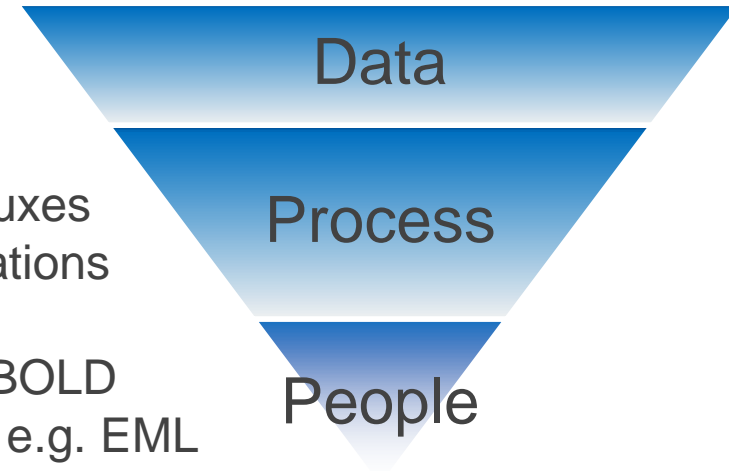
surface atmospheric exchange/fluxes
meteorological temporal interpolations

- of **data science methods**

APIs @ PhenoCam, MG-RAST, BOLD
Ecological semantic conventions e.g. EML
Dictionary-driven generic observational pipeline
Data integration via workflow models

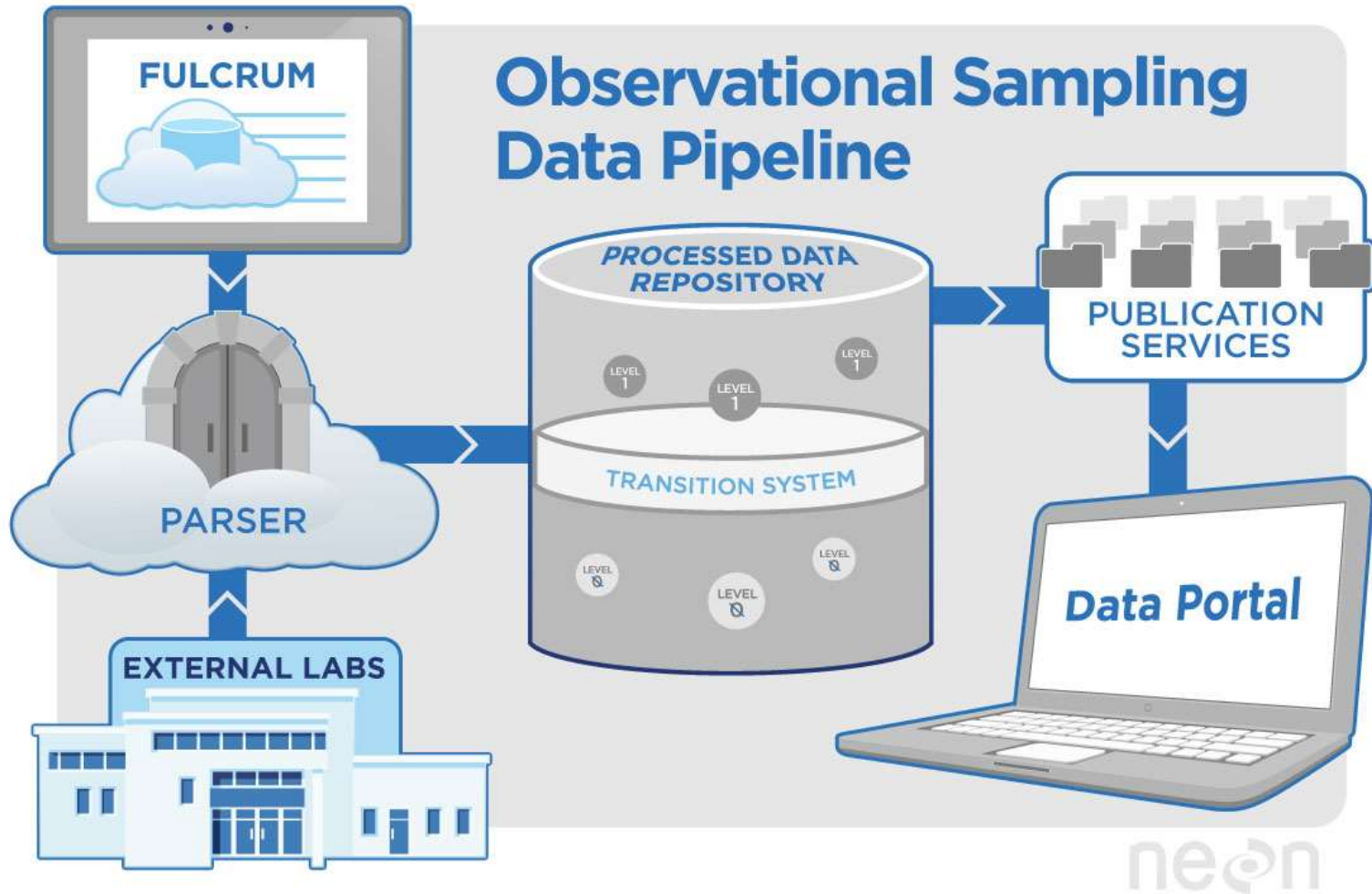
- of **practitioners**

21 Working Groups; ~1,500 unique users & 75k API hits/month
Carpentry workshops & interns
Digital identifiers for data use/provenance
Large data facilities coalitions, e.g. RDA, CDF, ESIP

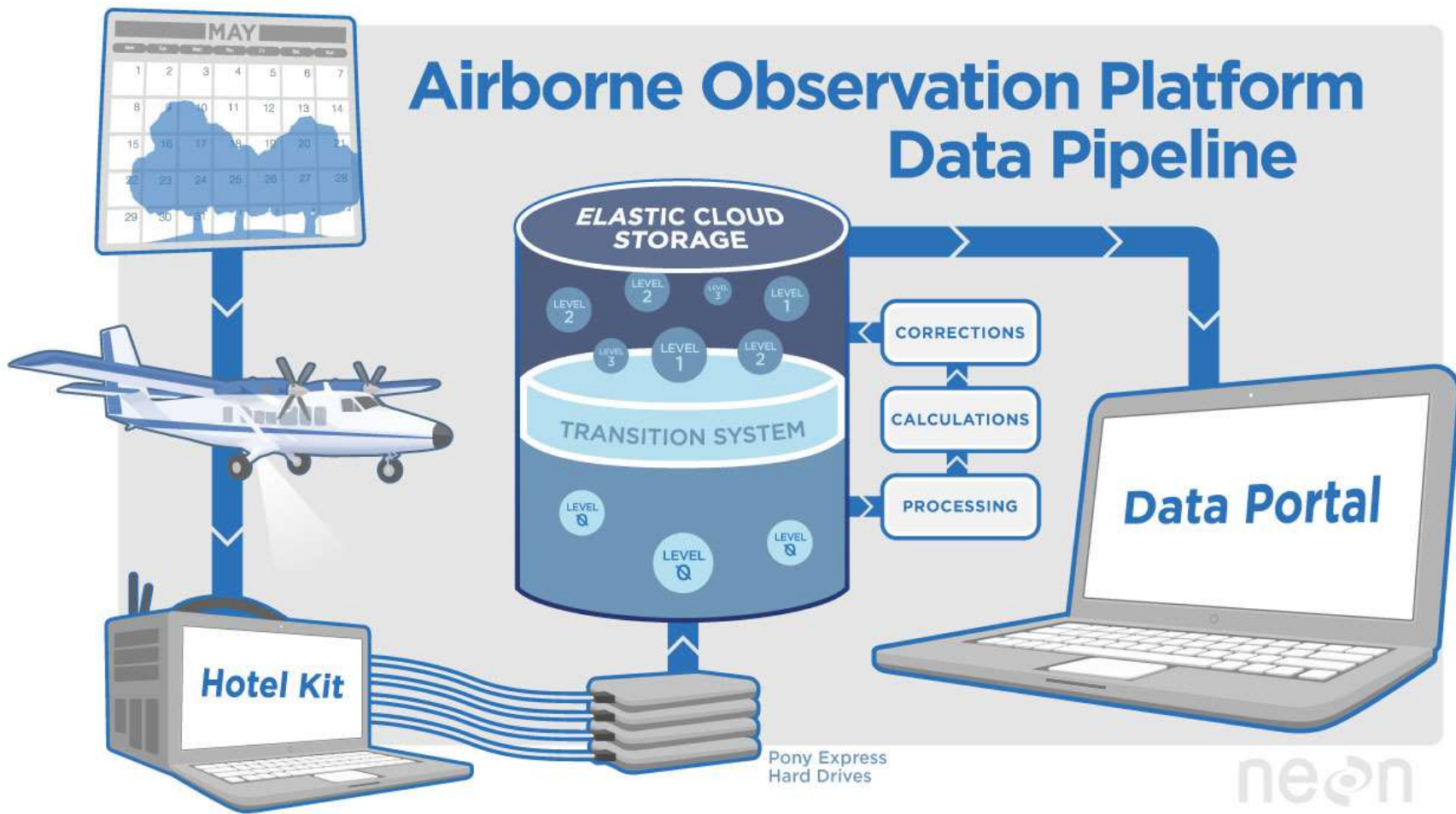


NEON CI Data Delivery

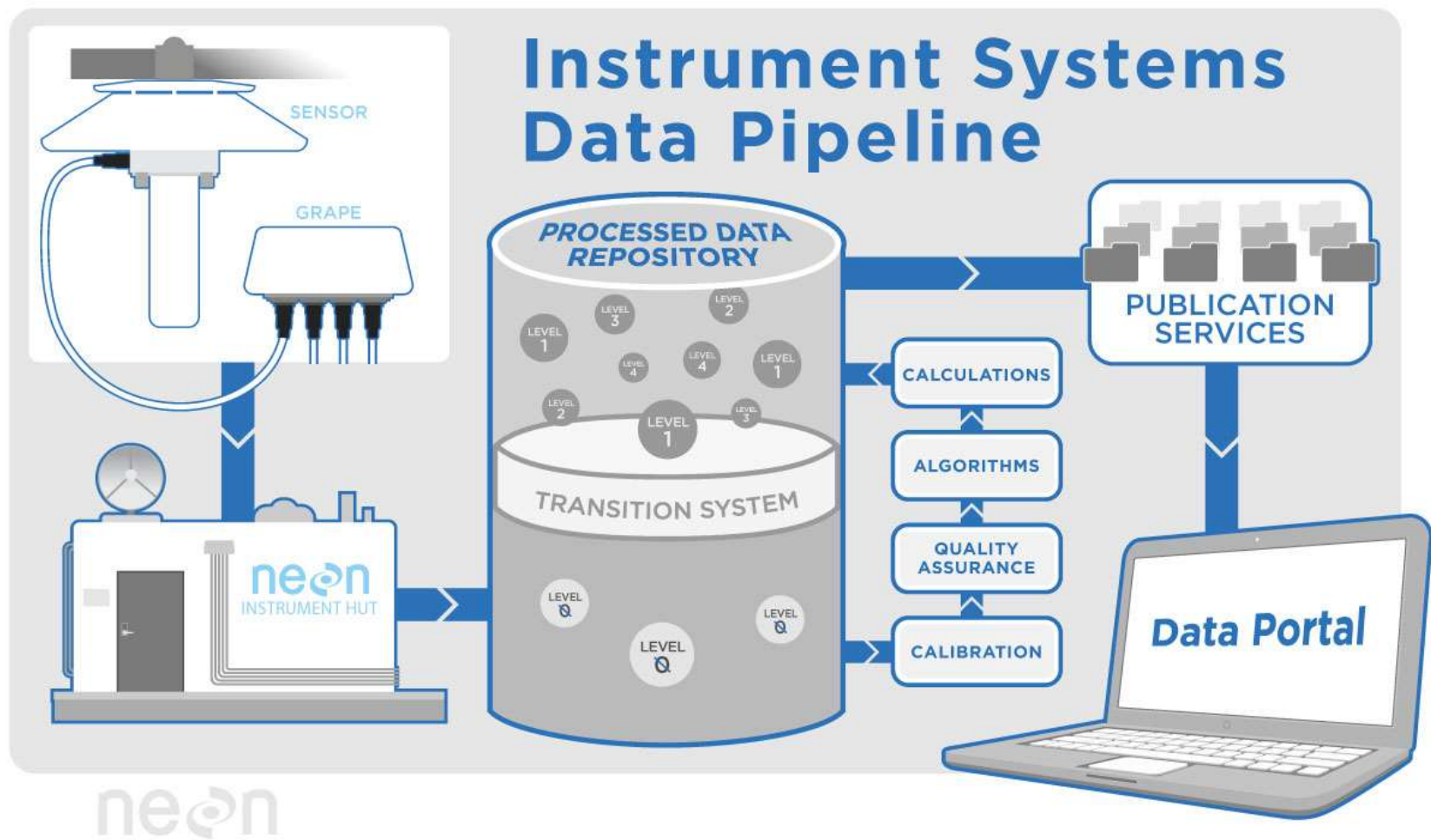
Data Product Type	General Latency from Sampling to Portal	Primary Influence on General Latency	Existing Backlog	Estimated Complete Publication of Backlog
AOP	2 months	Staff during flights	2018 ongoing	Nov2018
AOP legacy		Reprocessing with new algorithms	2013-2016	Dec2018
Land, Water, Soil, Meteorology	15 th of next month	5 day communications completion buffer	2018 ongoing	Monthly
Land, Water, Soil, Meteorology legacy		Volume of prior site-months	2013-2017	Nov2018
Eddy Covariance	15 th of next month	5 day communications completion buffer	2018 ongoing	Monthly
Eddy Covariance legacy		Volume of prior site-months & configurations	2013-2017	Nov2018
OS – field observations	1-4 months (~9-13 for fish & veg structure)	Reviews as season progresses. HQ staff.		Oct2018
OS – domain lab results	1-3 months (6 & 13 for morphosp. & bathy)	Reviews as season progresses. HQ staff.		Sep2018
OS – external lab rolling batches	1-9 months	Reviews as season progresses. Lab deliveries.		Dec2018
OS – one time samples (megapits, soil)		external lab processing - MBL		Dec2018



- 47 terrestrial & 34 aquatic sites, 100s protocols, dozens of labs
- Latencies from 5 to 365+ days

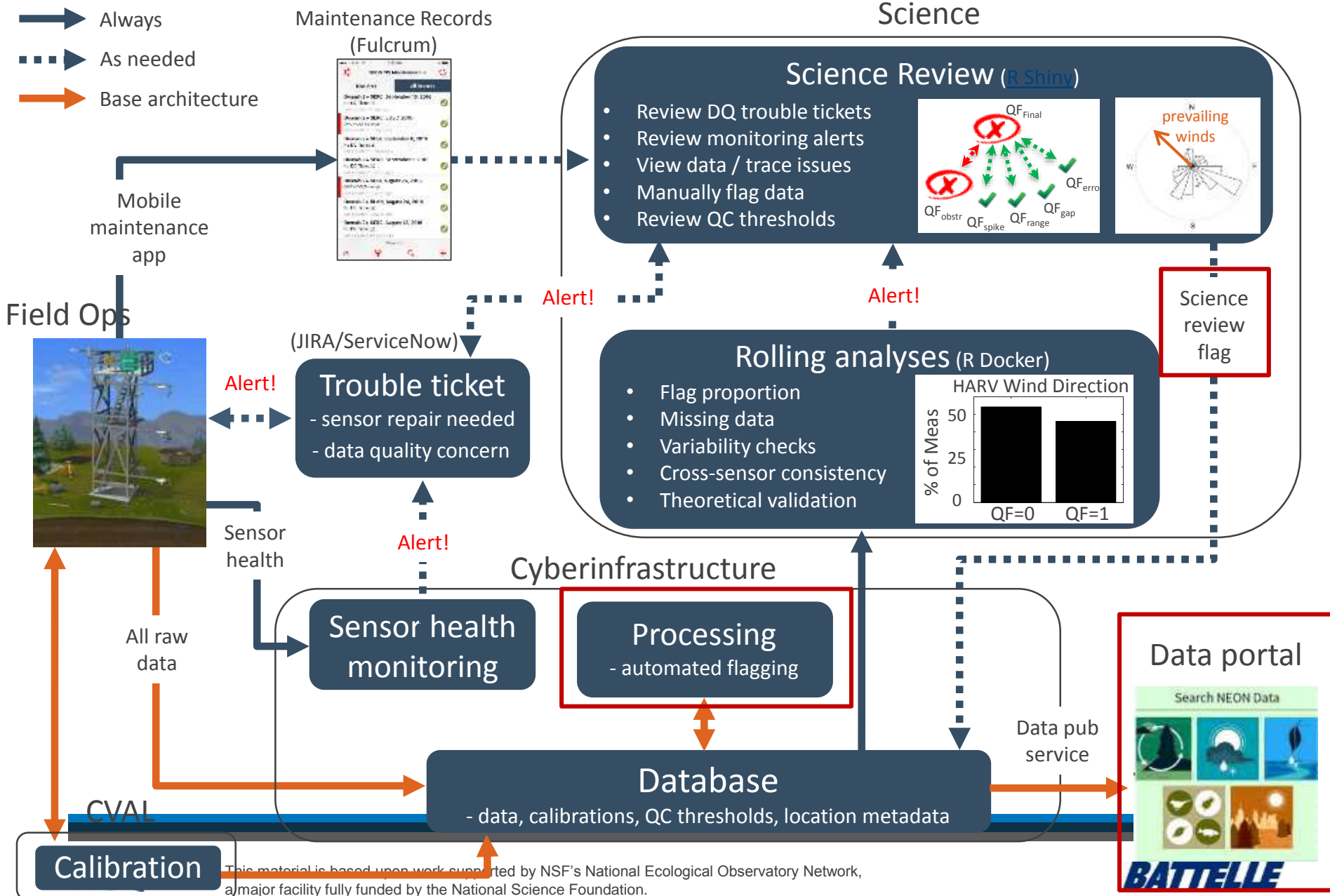


- 19 domains, 100-300km² flown annually at peak greenness
- Images, Spectroscopy (380 to 2500 nm), LiDAR
- 1-500GB/data product, ~2PB/yr



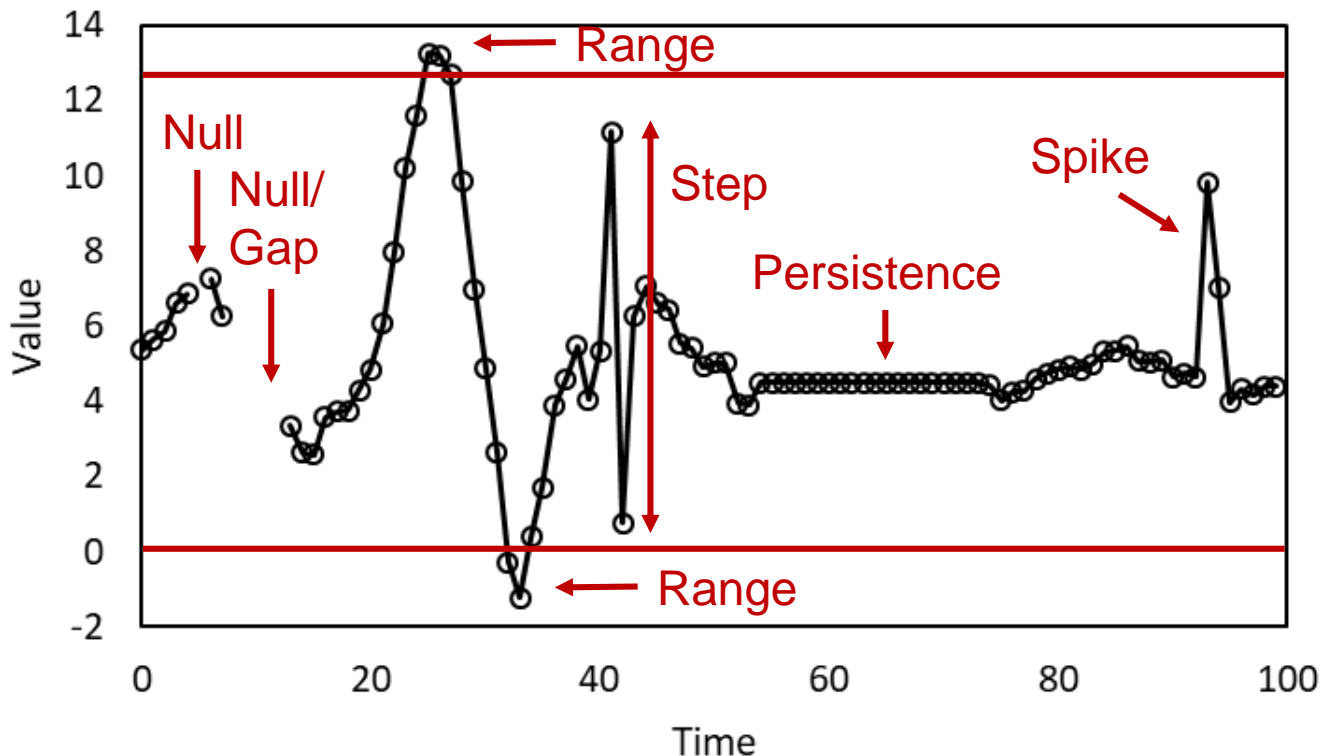
- 81 sites, 100s sensors per site, up to 40Hz
- Sensor→GRAPE→LC→HQ 5day lag, ~5TB/month

Instrumented Systems Data QA/QC Framework



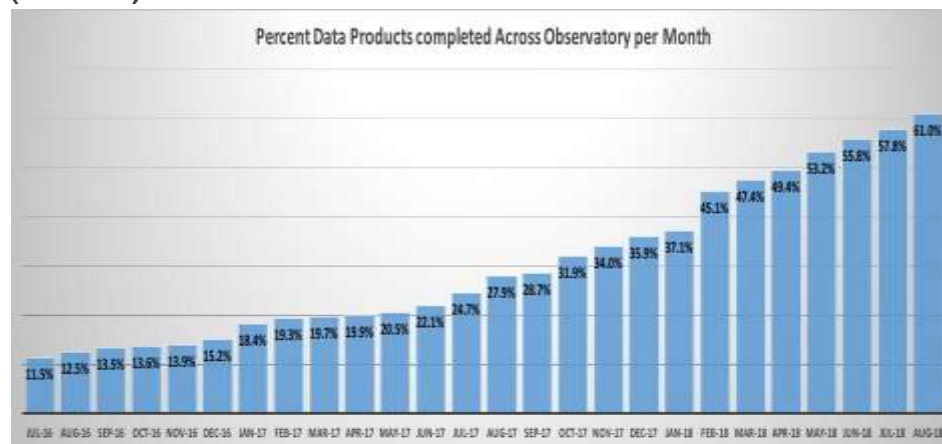
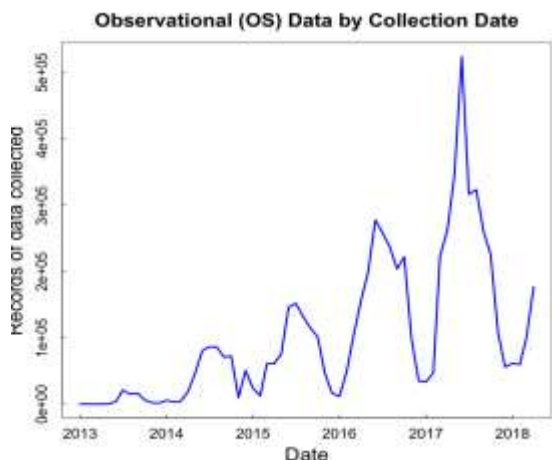
Automated flagging

- Applied on each calibrated data value (native resolution)
- Basic tests: Null, Gap, Range, Step, Spike, Persistence
- Sensor-specific tests: e.g. sensor diagnostic
- If chosen, can remove any data point that fails a test

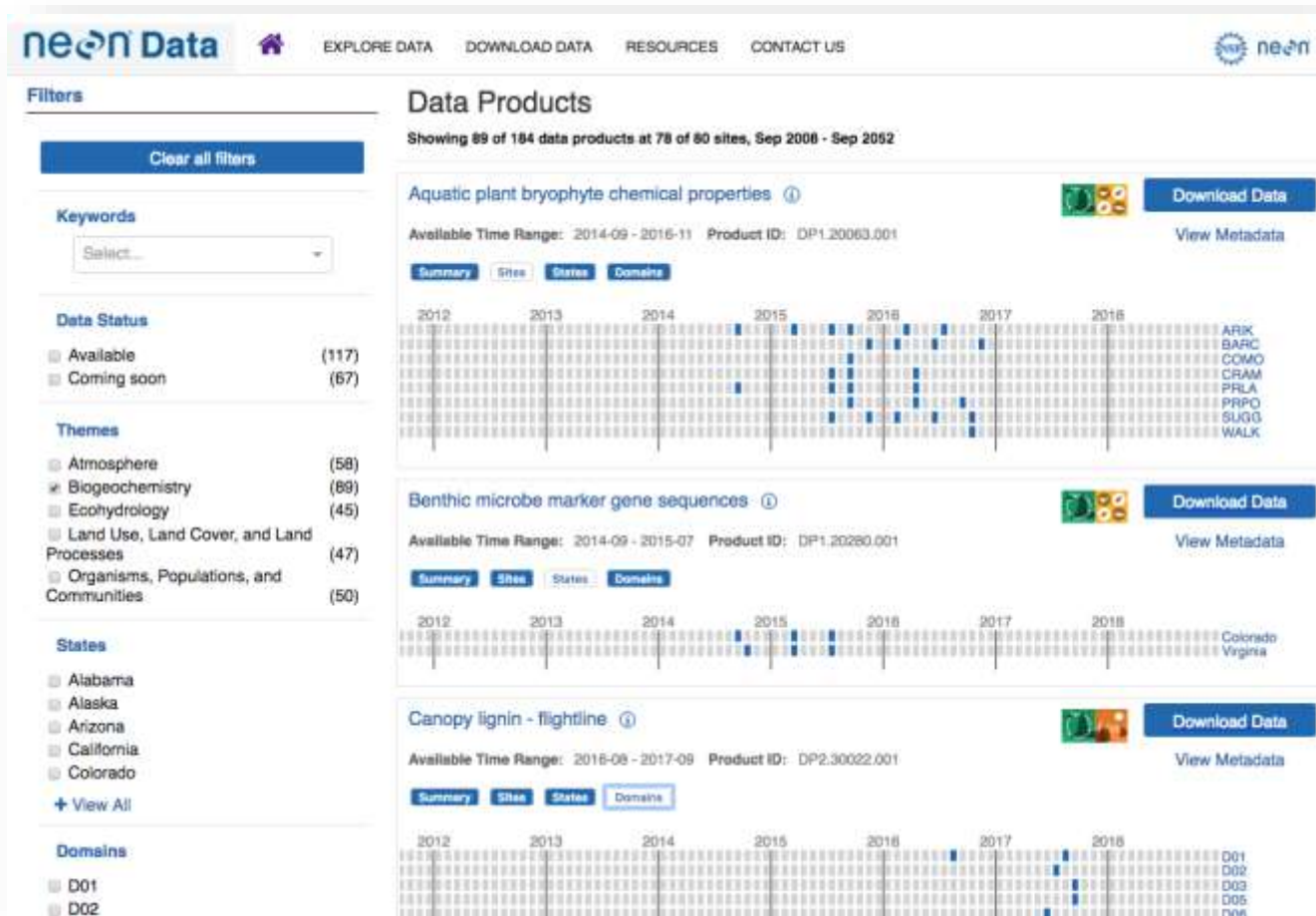


What's in a data package?

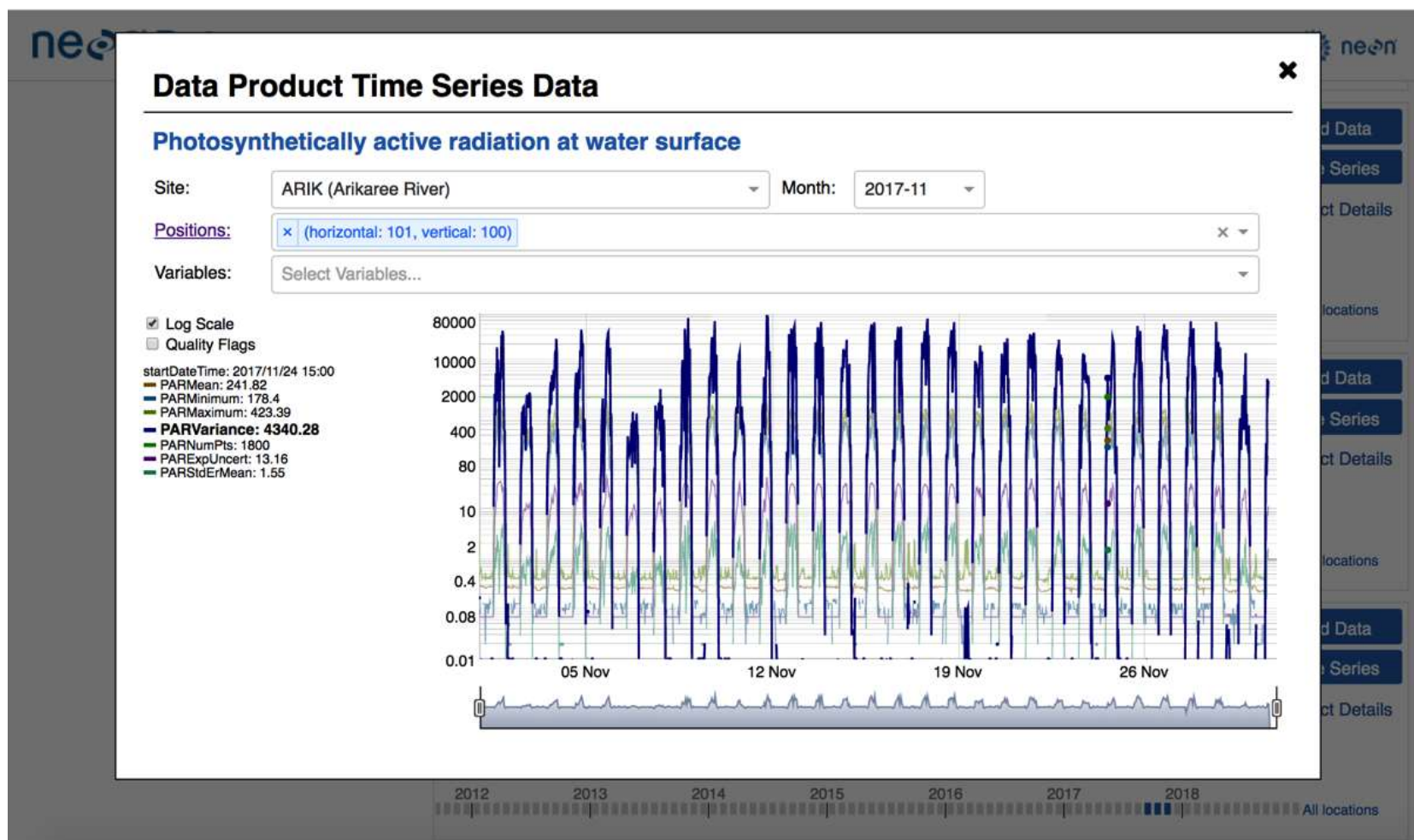
- Monthly or Annual data files with all basic data and quality flags
- CSV (IS, OS), HDF5 (EC, AOP), TIF (AOP), LAS (AOP), FASTA (metagenomics)
- Readme text file – general info about the data product
- Additional quality metrics, external lab data
- User guides & protocols; algorithm & sensor configuration documents
- Variable definitions, validation rules, and sensor positions
- Machine-readable metadata (EML)



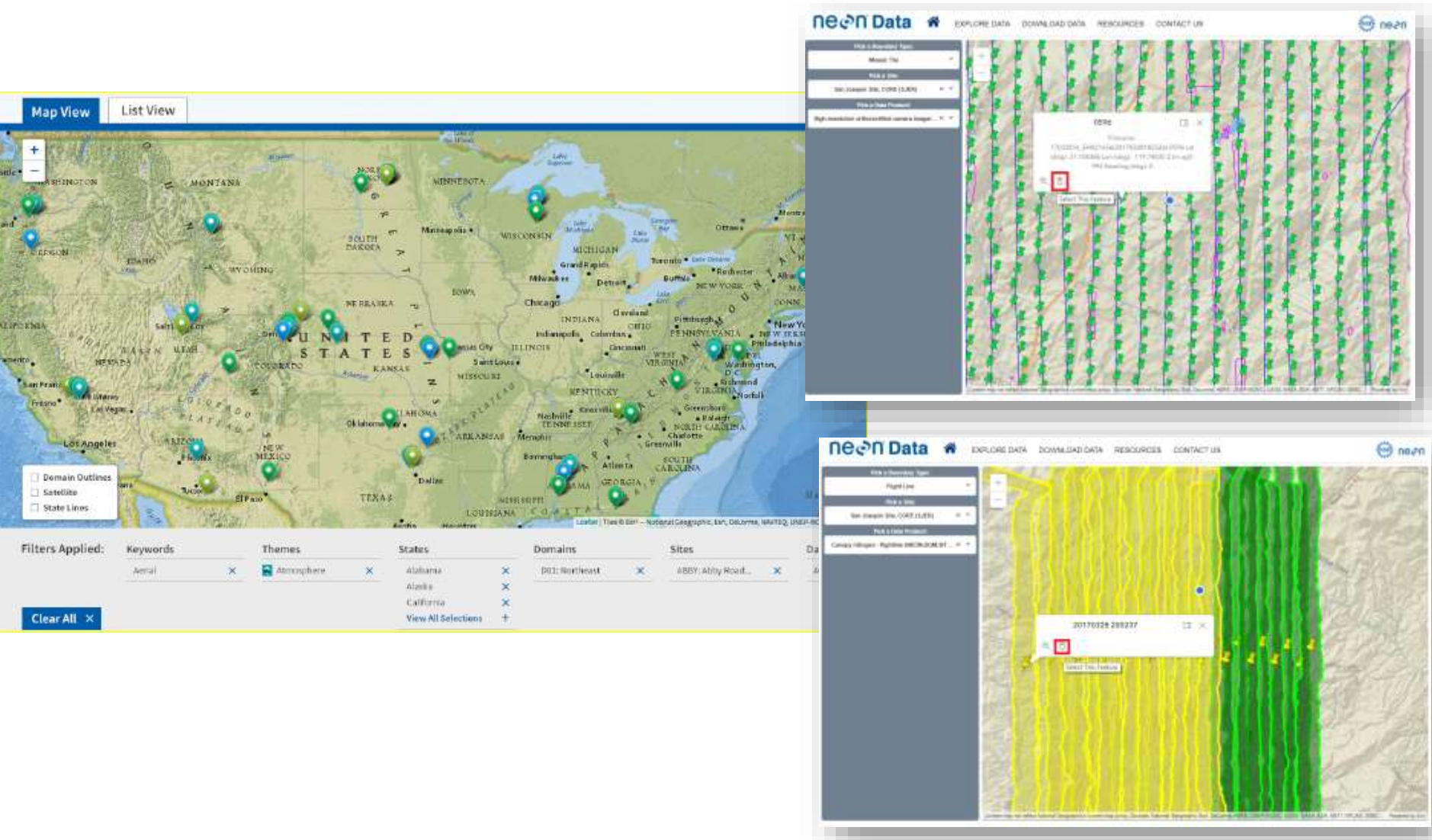
Interactive Browse



Interactive Time Series



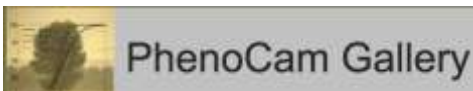
Spatial Viewer (prototypes)



Partner Repositories

- Work with specialized repositories for domain-specific data
- Use their APIs to keep NEON caches in sync
- Field data on portal; specialized data hosted

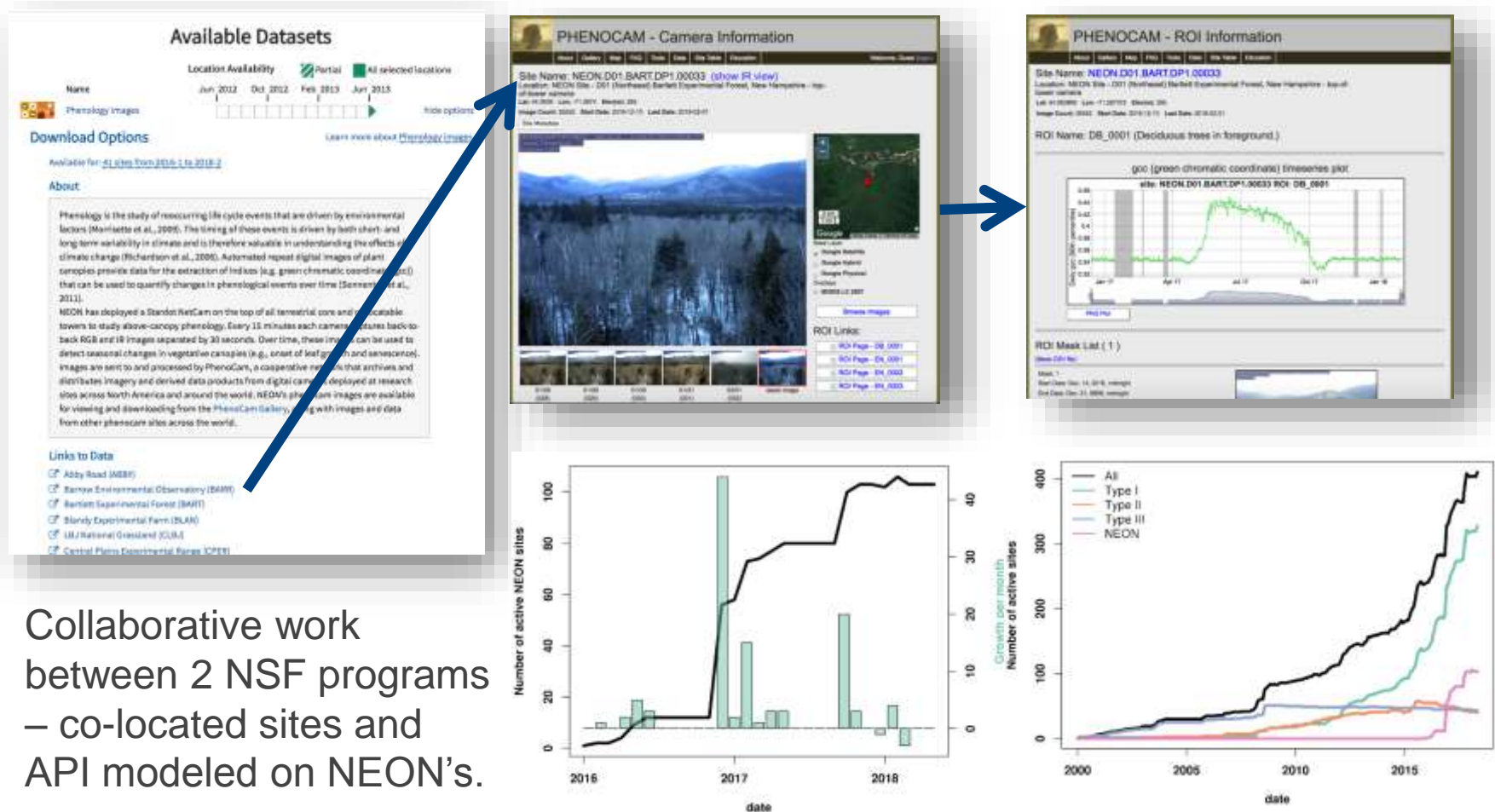
Current



Future



Data Portal & Phenocam Gallery



Collaborative work
between 2 NSF programs
– co-located sites and
API modeled on NEON's.

Data Portal & BOLD

Name

 Ground beetle sequences DNA barcode

Location Availability

Partial	All selected locations		
Jun. 2012	Oct. 2012	Feb. 2013	Jan. 2013

hide options

Download Options

Available for: [3 sites from 2013-7 to 2015-9](#)

Learn more about [Ground beetle sequences DNA barcode](#)

About

This data product contains the quality-controlled laboratory metadata and QA results for NEON's cytochrome oxidase I (COI) barcoding of ground beetles sequences. The DNA barcoding procedure involves

"BOLD Project: Ground beetle sequences DNA barcode" redirects to a page on the BOLD public data portal for the queried data. This is a dynamic link and will automatically update based on the user query.

Links to Data

- BOLD Project: Ground beetle sequences DNA barcode
- BOLD Project: NEON Prototype Ground beetle sequences DNA barcode

Documentation

- include relevant documents for this Data Product

[XML] files for this Data Product are included in all downloads ([more about XML at NCBI](#) and [ESL](#))

Data

- Basic: the core set of attributes.

Format

CSV (Comma Separated Values)
Estimated size 0.429 MB

[DOWNLOAD DATASET](#)

By choosing to download NEON data, you agree to NEON's Terms and Conditions.

BOLD SYSTEMS DATABASE GENTRIFICATION TAXONOMY WORKBENCH RESOURCES LOGIN Q

"BETP" Public Status: [v] **SEARCH**

Specimens: 000 000 000

Sequences: 000 000 000




Combined: 000 000 000

Generate Map: 2

Show Help

Showing Records 1 to 100 Results Per Page: **100** (v)

Page: 5 of 3 records


#1	[IAP361-14 - <i>Pisonia natsumi</i> [COI-SP-058] <i>Taxonomy:</i> Anthrhopoda, Insecta, Coleoptera, Carabidae, Pisonina IdentFacts : NCBI:NCIT Z110.00054 (Lamprell), D10.00081 (Belkin) <i>Catalogue:</i> National Ecological Observatory Network, United States Catalogue : United States, Colorado, Domain 10	
#2	[IAP362-14 - <i>Dicladona punctatula punctatula</i> [COI-SP-058] <i>Taxonomy:</i> Anthrhopoda, Insecta, Coleoptera, Carabidae, Chalcidina IdentFacts : NCBI:NCIT Z110.00058 (Jungbluth), D10.00094 (Belkin) <i>Catalogue:</i> National Ecological Observatory Network, United States Catalogue : United States, Colorado, Domain 10	
#3	[IAP363-14 - <i>Harpalus demetrius</i> [COI-SP-058] <i>Taxonomy:</i> Anthrhopoda, Insecta, Coleoptera, Carabidae, Harpalini IdentFacts : NCBI:NCIT Z110.00058 (Jungbluth), D10.00094 (Belkin) <i>Catalogue:</i> National Ecological Observatory Network, United States	

Results Summary

Found **282** published records, forming **26** taxa clusters, with specimens from 1 country, deposited in 1 institute.

Of these records, **282** have species names, and represent **31** species.

Specimen Distribution



[illegible]

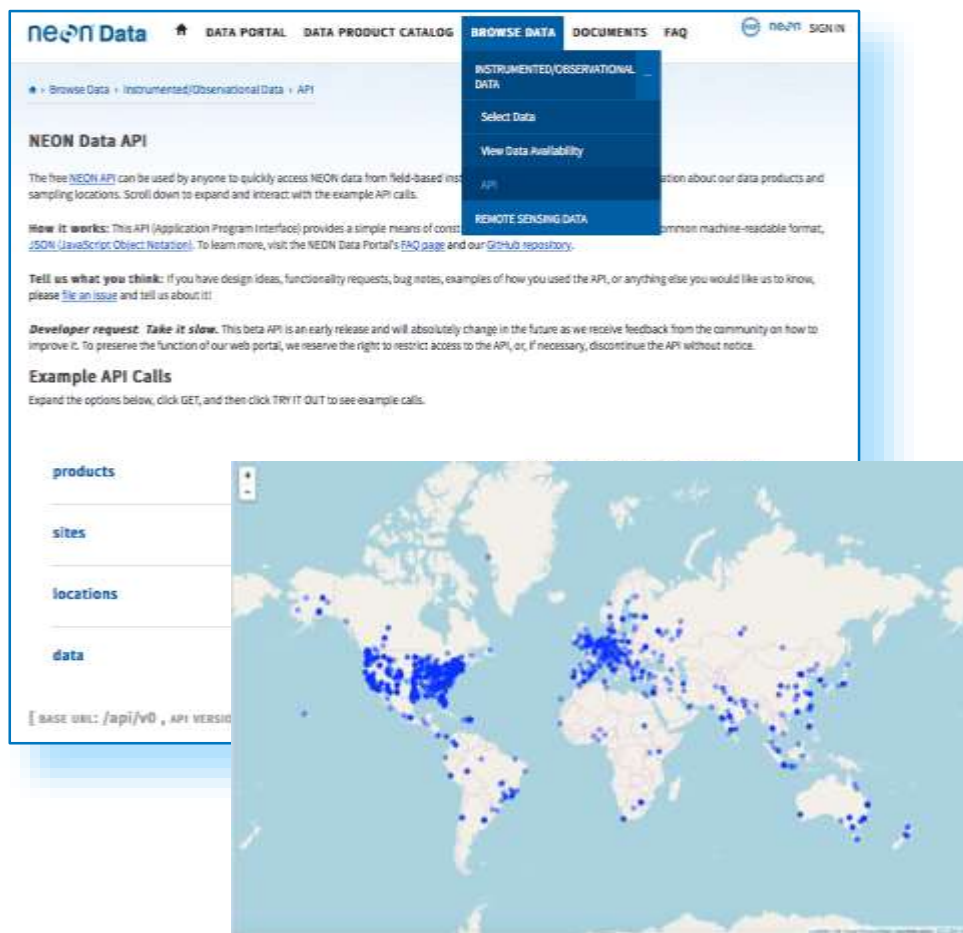
Access to thousands of beetle, mosquito, small mammal, and fish records, from current protocols as well as prototype protocols and sites.

MG-RAST is similar in functionality

REST API: Programmatic Access

<http://data.neonscience.org/data-api>

- Open access
- Provides product information, site and within-site information, documents, and data
- Returns information in JSON format
- Returns both zipped data packages and individual data files
- Recently added endpoints for access to taxonomic lists and sample custody histories



NEON CI Software Architecture - Layers

data services, transport, queuing

Publish



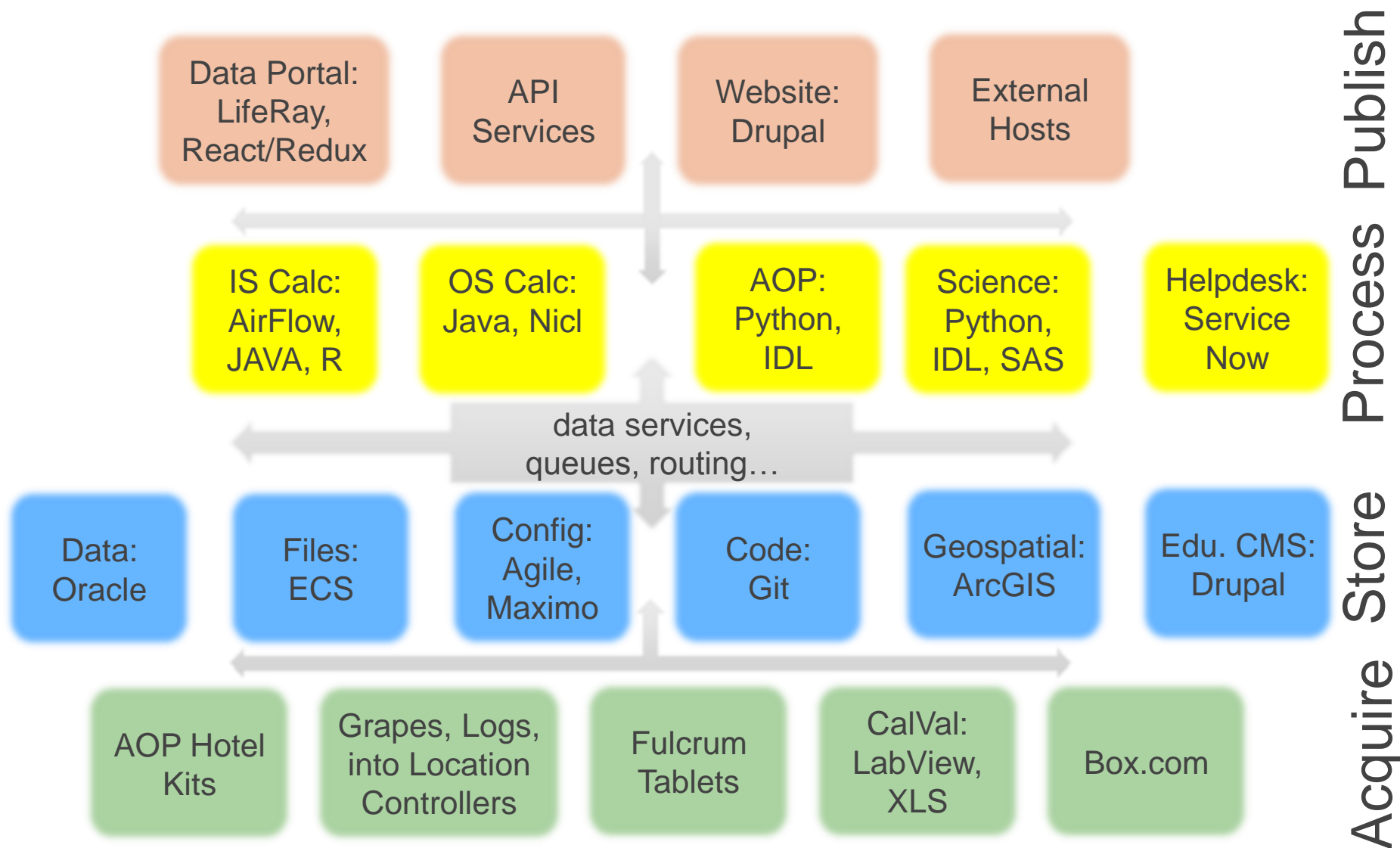
Process

Store

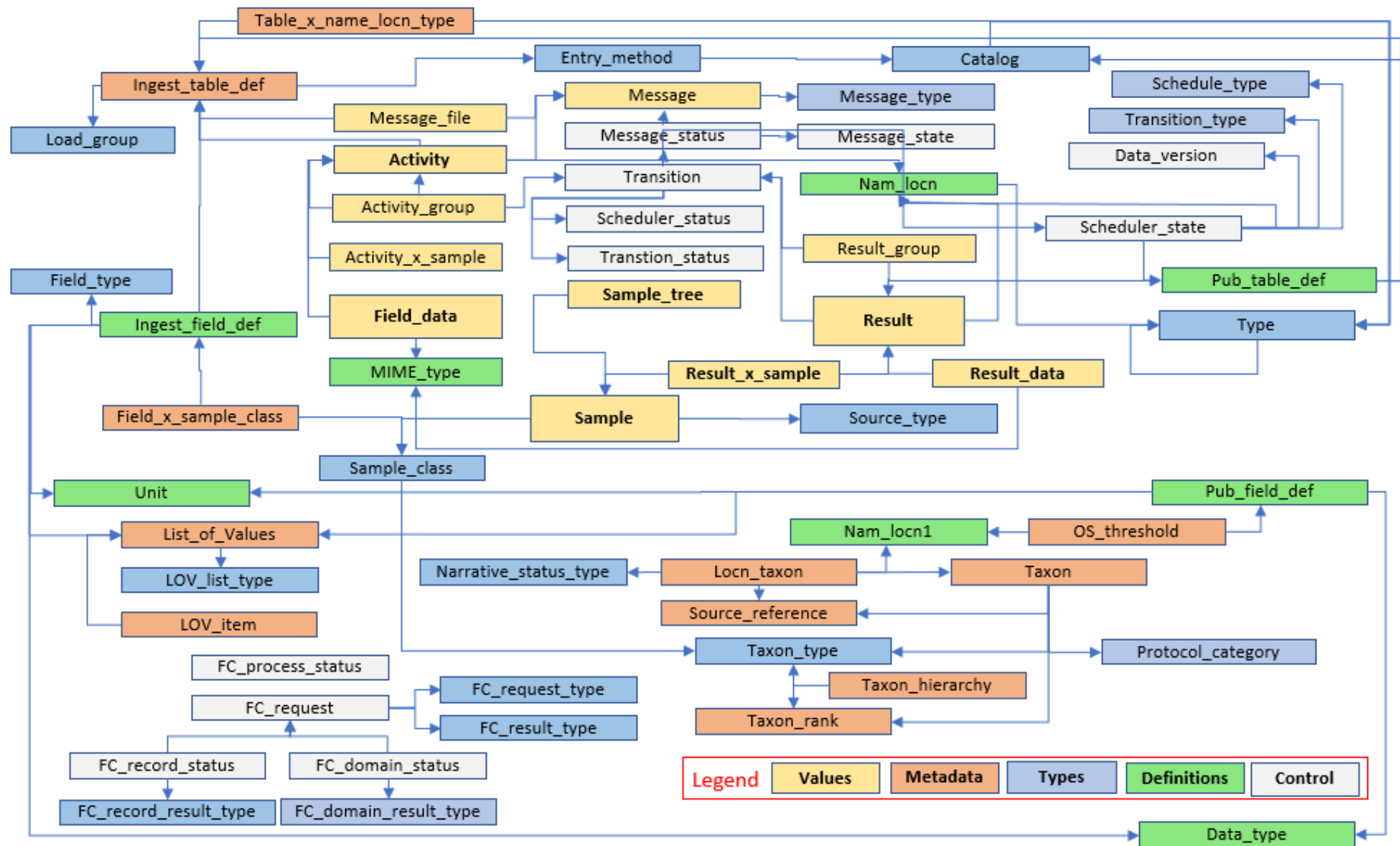
Acquire



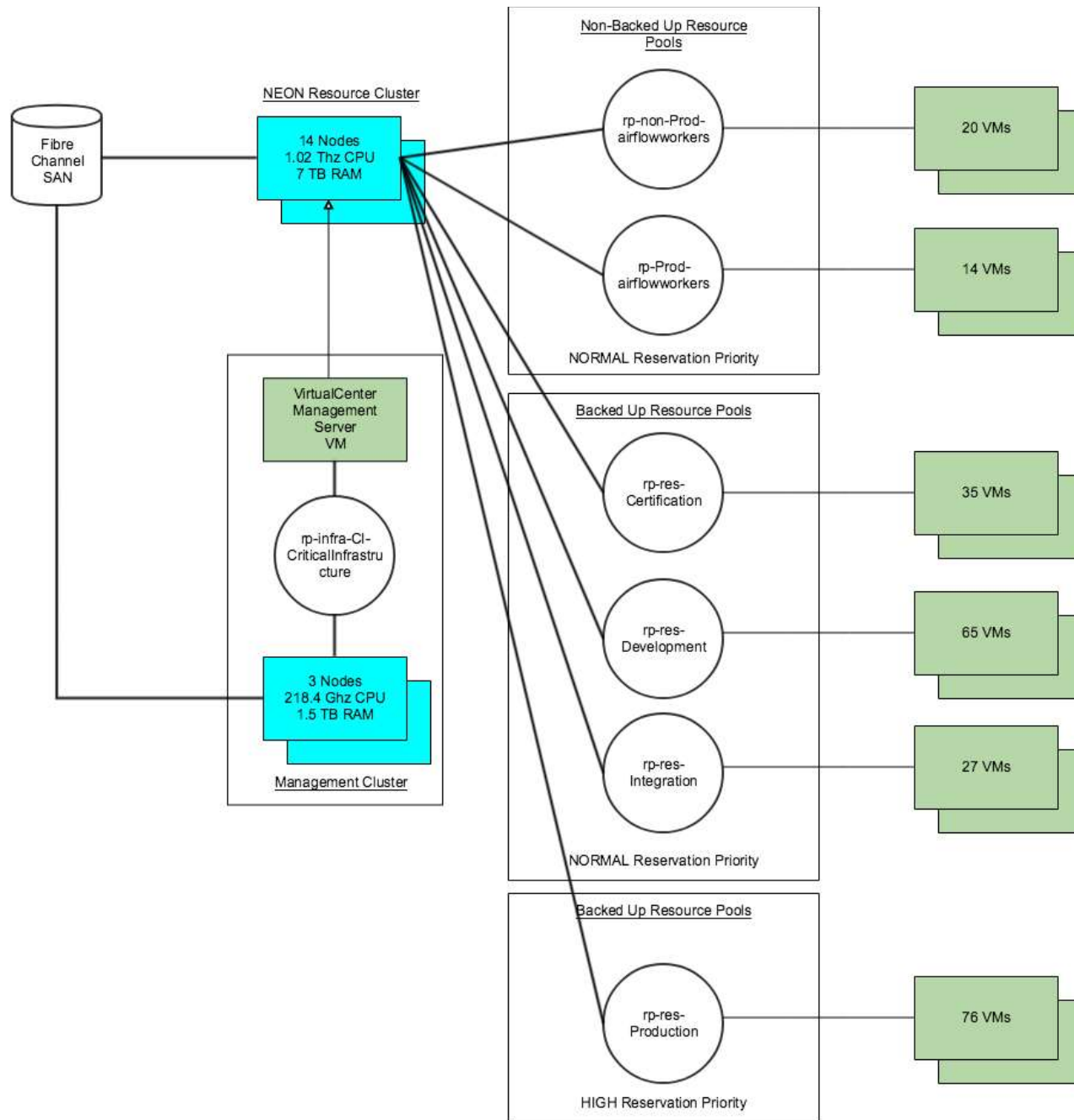
NEON CI Software Architecture - Elements



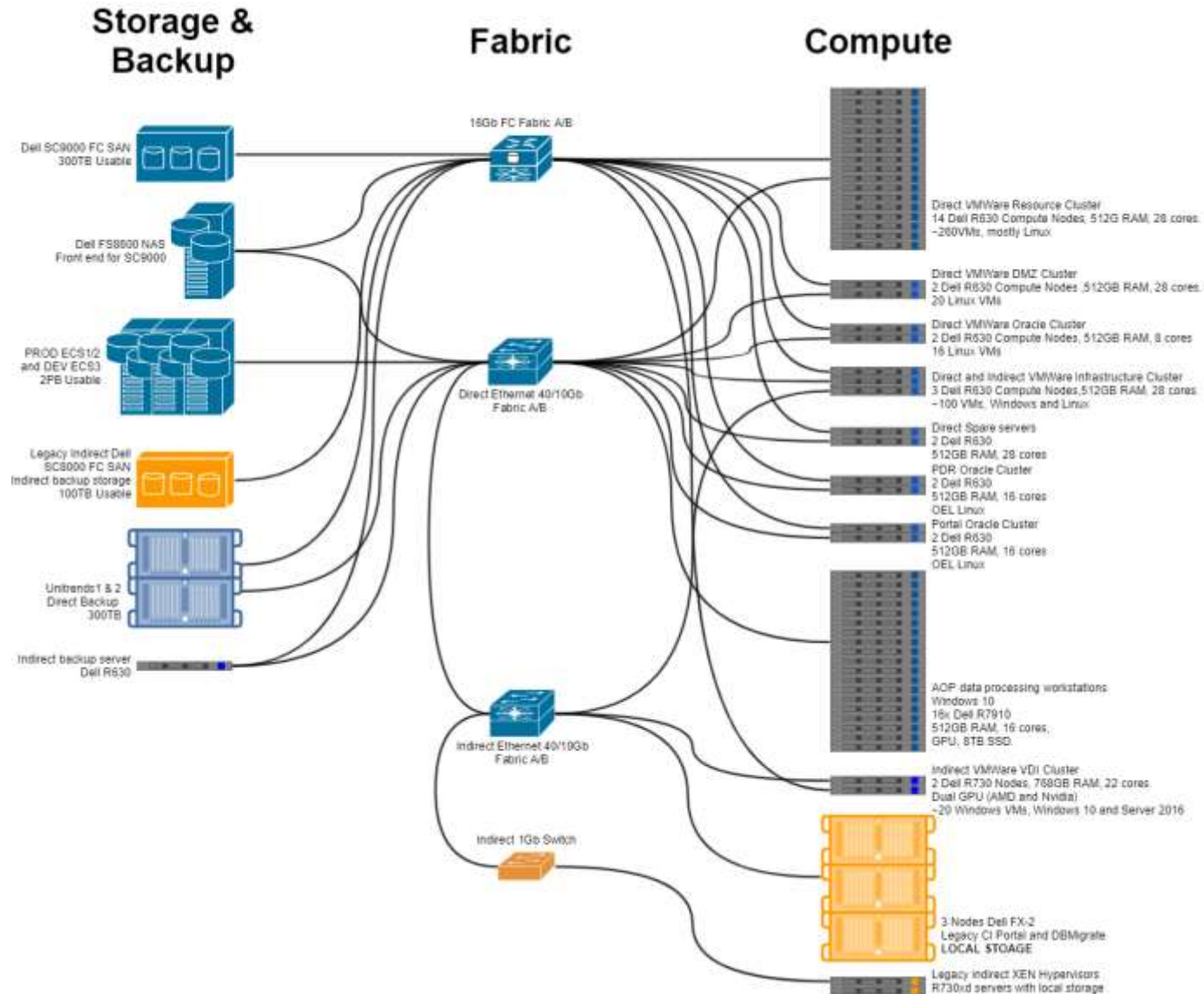
PDR Database – Observations Logical Model



NEON CI Virtual Machine Resource Pools



Denver Datacenter – Configuration Schema



Cyber Security Overview

1. Security Strategy – NIST Cyber Security Framework: Identify, Protect, Detect, Respond, Recover
2. NIST CSF Timeframe – Program in progress – Goal is to address all categories by mid 2019
3. Perimeter Security – All Ingress / Egress Points are Firewall protected
4. Cloud Security – Encrypted Connections (TLS 1.0+) to Cloud Apps
5. Endpoint Security – Anti-malware / Endpoint Controls
6. Vulnerability Remediation – Constant patching schedule
7. Email Security – Microsoft Advanced Threat Protection – Phishing, Malware, Safelink protection
8. User Awareness – All employees are required to complete annual cyber security training

NEON CI & Data Interoperability

1. Across the community of researchers

- Data collection, QA, standardized encoding and uncertainty handling protocols
- Data exchange I/O formats & metadata
- Data processing algorithms, code libraries

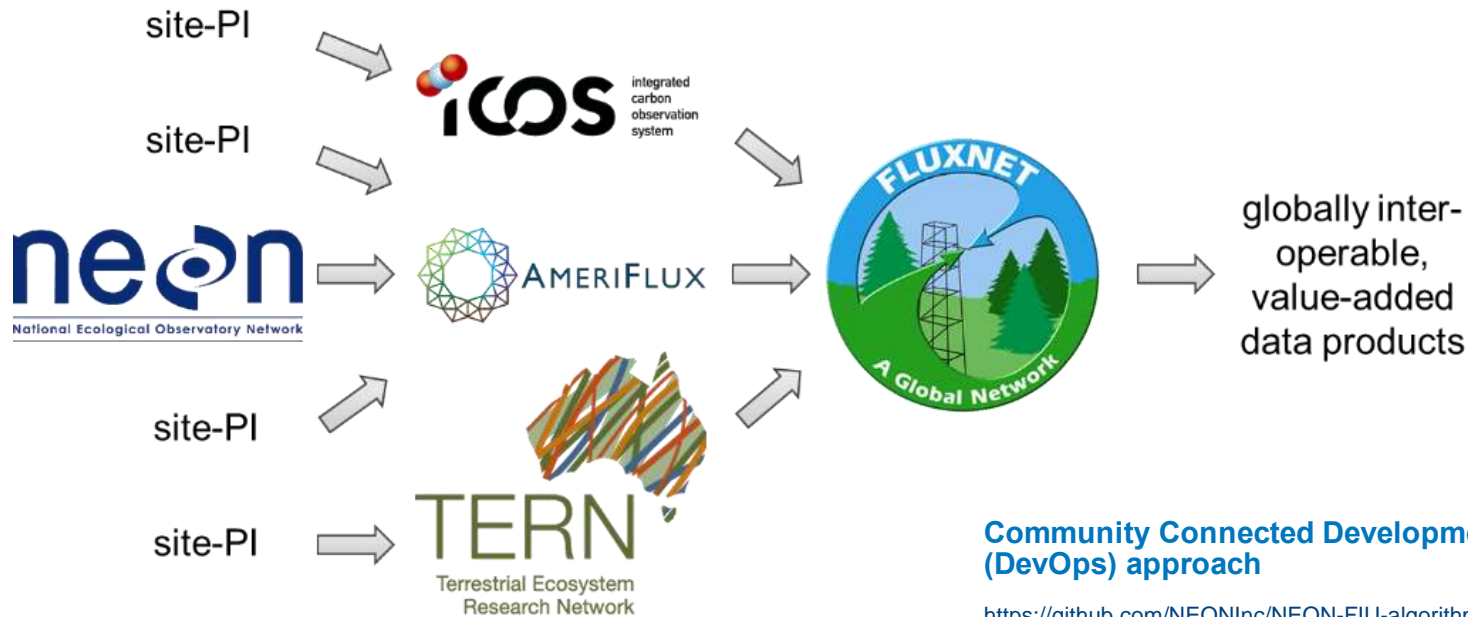
2. Between research facilities

- NEON, AmeriFlux and Fluxnet common methods
- Fully compliant metadata provisioning of NEON results to 6 partner host systems
- Participation in meta-analyses of international CI methods, e.g. workflow objects

3. In active collaboration with aggregators

- DataOne partner and metadata practitioner
- EarthCube, CDF, ESIP, RDA coalition participants

Data Product Collaboration (example)

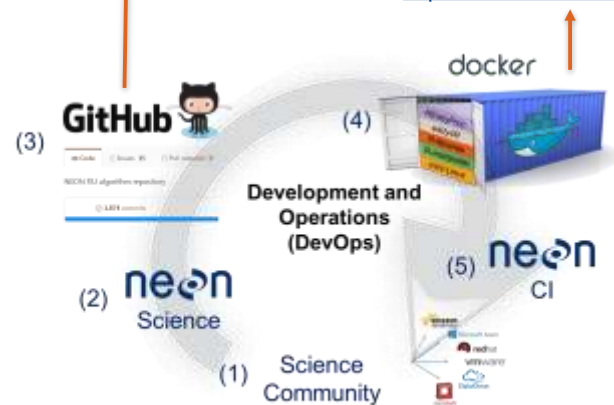


Community Connected Development and Operations (DevOps) approach

<https://github.com/NEONInc/NEON-FIU-algorithm>

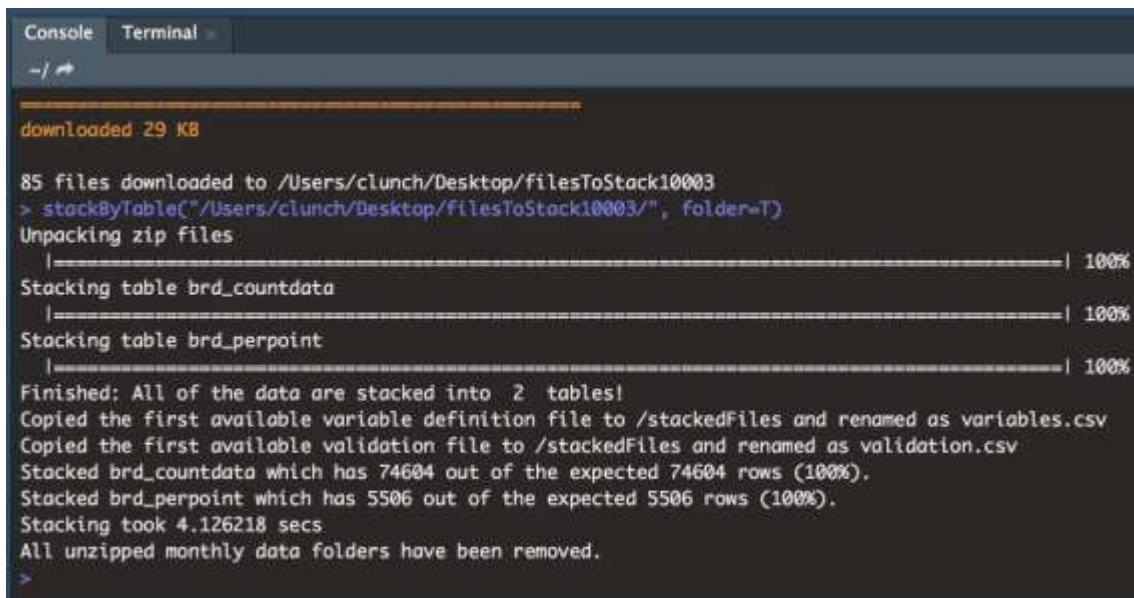
<https://hub.docker.com/u/stefanmet/>

Metzger, et.al. 2017. **eddy4R 0.2.0**: a DevOps model for community-extensible processing and analysis of eddy-covariance data based on R, Git, Docker, and HDF5. *Geosci. Model Dev.*, 10, 3189–3206.



Data Science Tools and Training: NEONScience on GitHub

- Open-source code packages to facilitate NEON data access and use
- Access to code used to generate NEON data products



```
Console Terminal
~/
downloaded 29 KB
85 files downloaded to /Users/clunch/Desktop/filesToStack10003
> stackByTable("~/Users/clunch/Desktop/filesToStack10003/", folder=T)
Unpacking zip files
|=====| 100%
Stacking table brd_countdata
|=====| 100%
Stacking table brd_perpoint
|=====| 100%
Finished: All of the data are stacked into 2 tables!
Copied the first available variable definition file to /stackedFiles and renamed as variables.csv
Copied the first available validation file to /stackedFiles and renamed as validation.csv
Stacked brd_countdata which has 74604 out of the expected 74604 rows (100%).
Stacked brd_perpoint which has 5506 out of the expected 5506 rows (100%).
Stacking took 4.126218 secs
All unzipped monthly data folders have been removed.
>
```

neonUtilities R package:

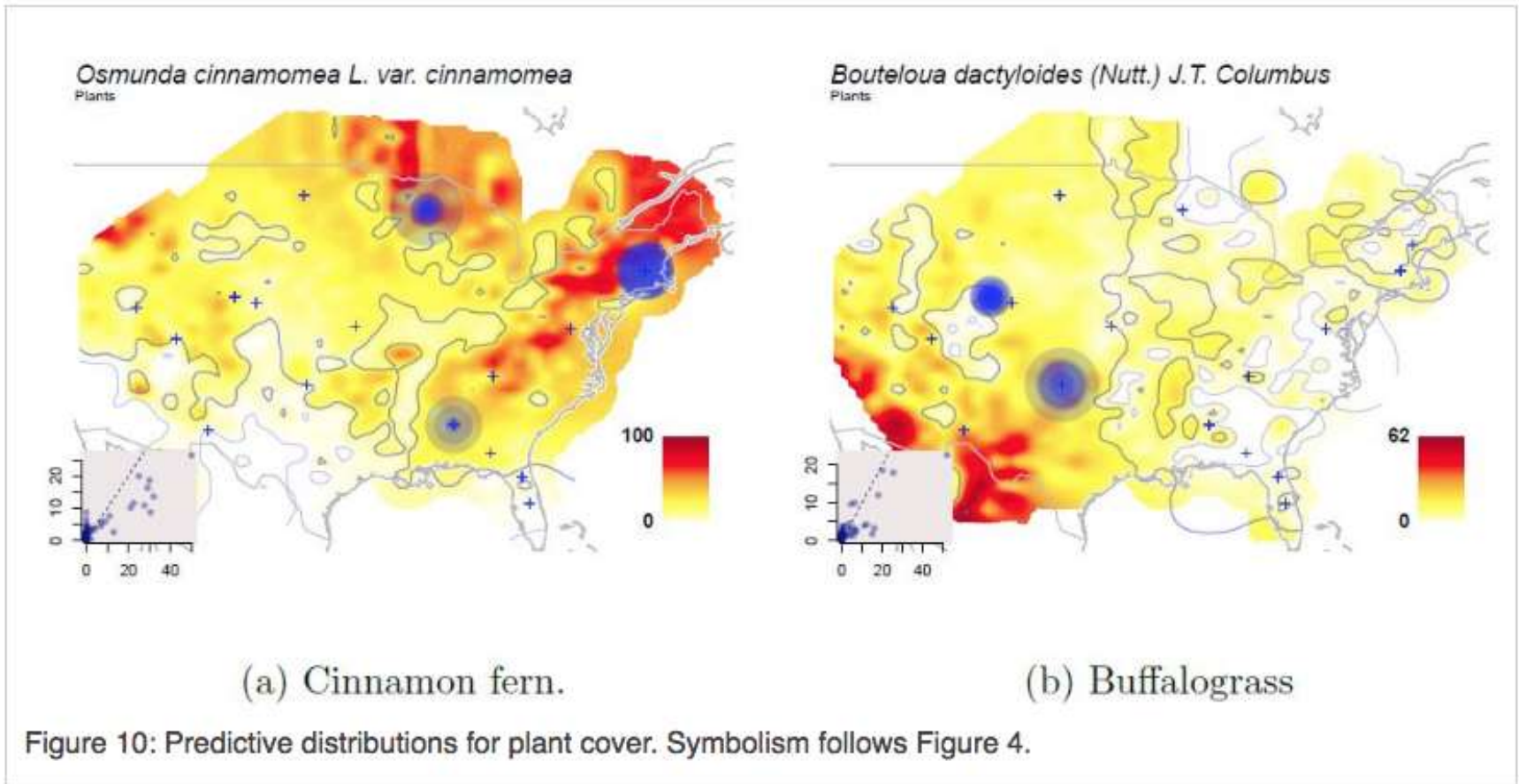
- Download data from NEON API
- Merge data files from Portal or API
- Convert file format for interoperability

Data Science Tools and Training: Tutorials & Workshops

- Advance users' data analysis skill levels
- Online tutorials for self-directed training
- Data Institute on-site at NEON HQ
- Workshops at meetings and conferences
- Includes training in using the [GitHub.com/NEONScience](https://github.com/NEONScience) packages



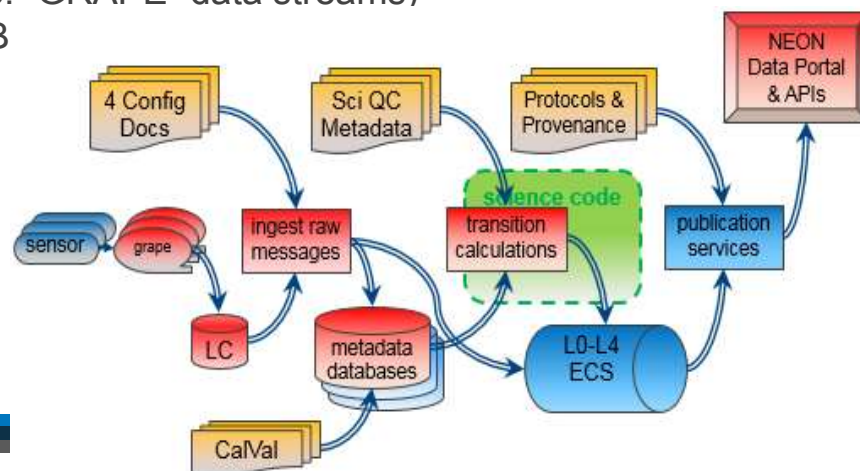
User Community Projects: Predictive Distributions



Clark lab: sites.duke.edu/neon/prediction

NEON CI Improvement Plans

1. Ongoing cyberengineering support
 - e.g. better interface to manage sensor metadata; unexpected data gap checks...
2. Prioritized queue of CI or Data Product enhancements
 - e.g. improve minimum viable products; address tech debt; better algorithms...
3. Service Management continued expansion/integration
 - e.g. more workflow delegation; closer linkage with SOM & monitoring...
4. Asset management tools & methods assessment
 - e.g. mitigate risks of vulnerable sensor installation metadata addressing...
5. Proposed research & development towards VI version2
 - Redesign data acquisition hw&sw (i.e. “GRAPE” data streams)
 - Pilot optimized transition engine & DB
 - Upgrade Data Portal platform & tools
 - Evaluate CI capacity, connectivity & offsite disaster recovery options





neon

Proudly operated by **BATTELLE**