# The Environmental Data Initiative (EDI)

**Supporting Curation and Archiving of Environmental Data**

# Overview

**EDI's mission and approaches - Corinna Gries**

Data repository and publishing - Duane Costa

Data publication workflow support - Colin Smith

Outreach and training - Kristin Vanderbilt
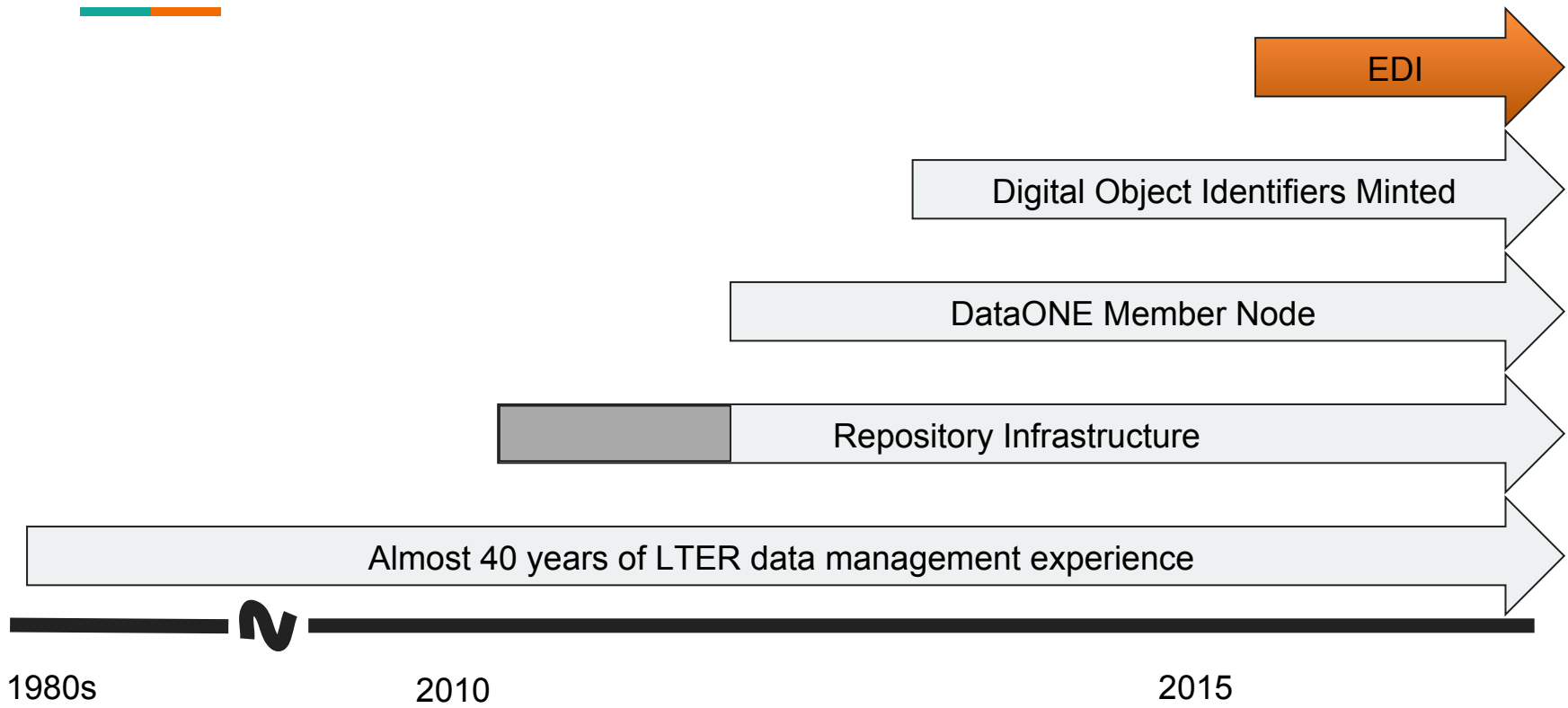
Data integration support - Margaret O'Brien

https://environmentaldatainitiative.org/          https://portal.edirepository.org          @EDIgotdata          edi-got-data

# History - we are standing on the shoulders of giants

EDI

Digital Object Identifiers Minted

DataONE Member Node

Repository Infrastructure

Almost 40 years of LTER data management experience

1980s          2010                              2015

# Mission and Goals

Accelerate the curation, archiving, and dissemination of environmental data

Ensure that environmental data are:

- Deposited into a data repository for long-term preservation and data integrity
- Easily discoverable and seamlessly accessible
- Documented with rich science metadata to enable reuse and integration

# Curation

Data curation support for data providers

- Experienced data managers on staff
- Consultation
- Training
- Data Curation Tool Development

# Archiving

EDI's repository builds on cyber infrastructure developed for LTER

Certified and registered as trustworthy repository (re3data, Nature, ESA, others)

Data are documented in Ecological Metadata Language standard

High quality standards enforced through automated congruence and completeness checking

# Dissemination

DataONE member node

Local search interface

Digital Object Identifier (DOI) through DataCite

Collaborations with FAIR project, journal publishers

Linking of publications and data sets

Documentation of data provenance

# Overview

EDI's mission and approaches - Corinna Gries

**Data repository and publishing - Duane Costa**

Data publication workflow support - Colin Smith

Outreach and training - Kristin Vanderbilt

Data integration support - Margaret O'Brien
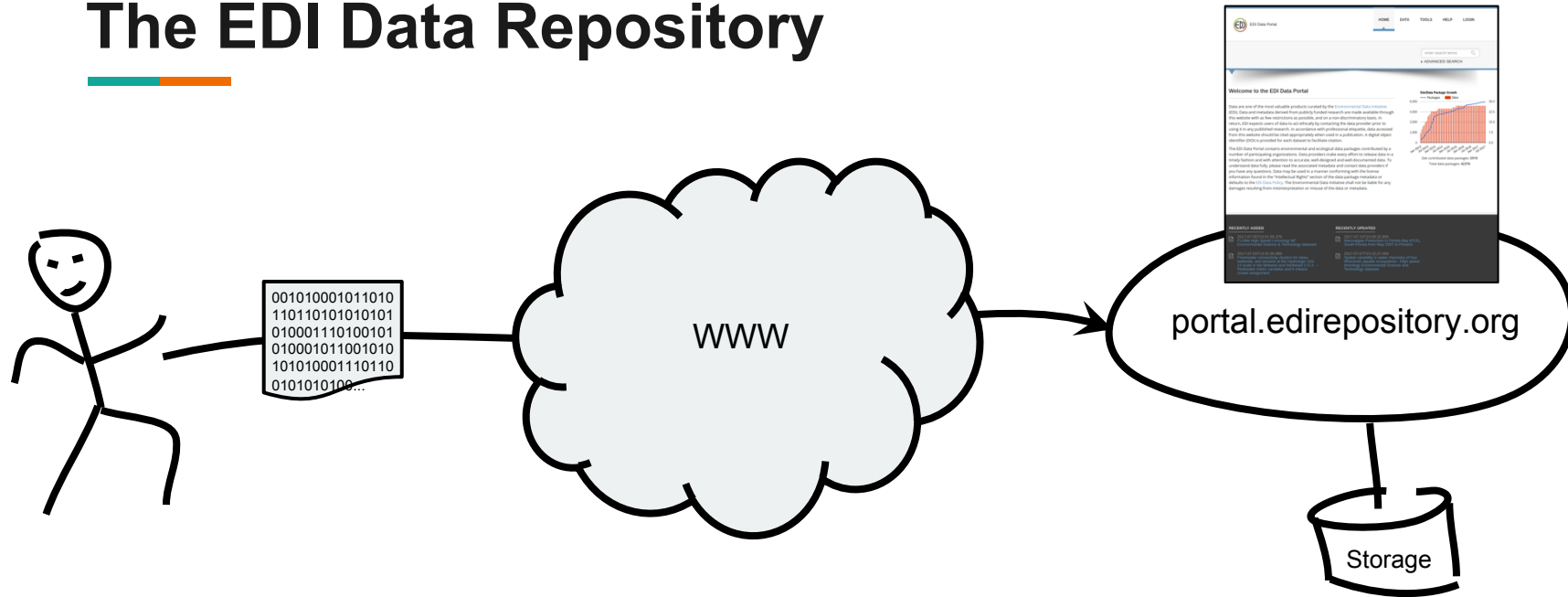
https://environmentaldatainitiative.org/          https://portal.edirepository.org          @EDIgotdata          edi-got-data

# The EDI Data Repository

WWW

portal.edirepository.org

Storage
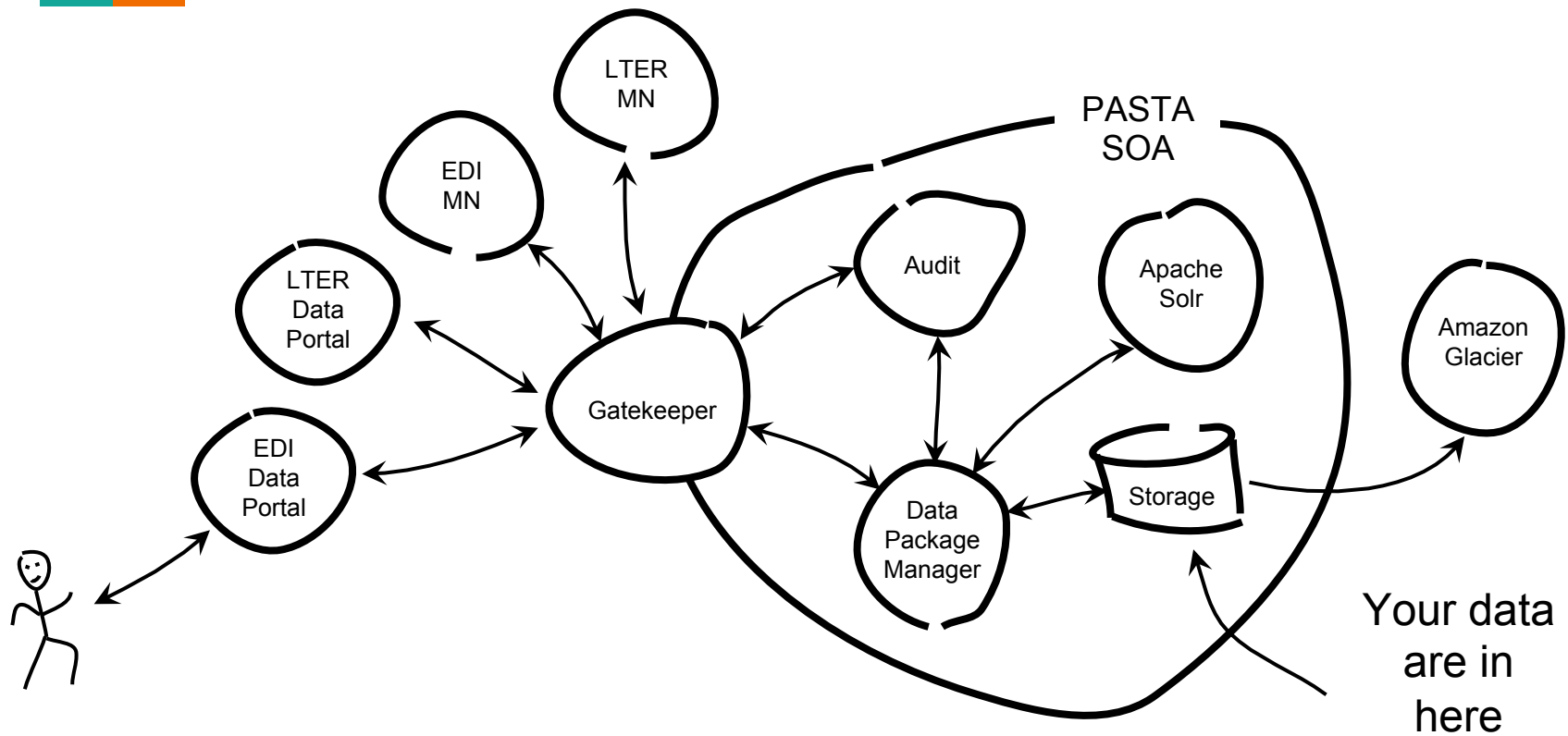
# Components of the EDI Data Repository

# About the EDI Data Repository (1 of 2)

- Community (LTER Network) designed with transparency in mind
- In continuous production since January 2013
- An Internet accessible, open data repository, open source project (GitHub)
- The EDI Data Repository is an *implementation instance* of a PASTA data repository. To put it another way, the EDI Data Repository is *powered by* PASTA software.
  - Metadata-driven (Ecological Metadata Language)
  - Java SE/EE implementation
  - Service Oriented Architecture (SOA) with a RESTful web service API
  - PASTA itself is *not* a website designed for direct human interaction
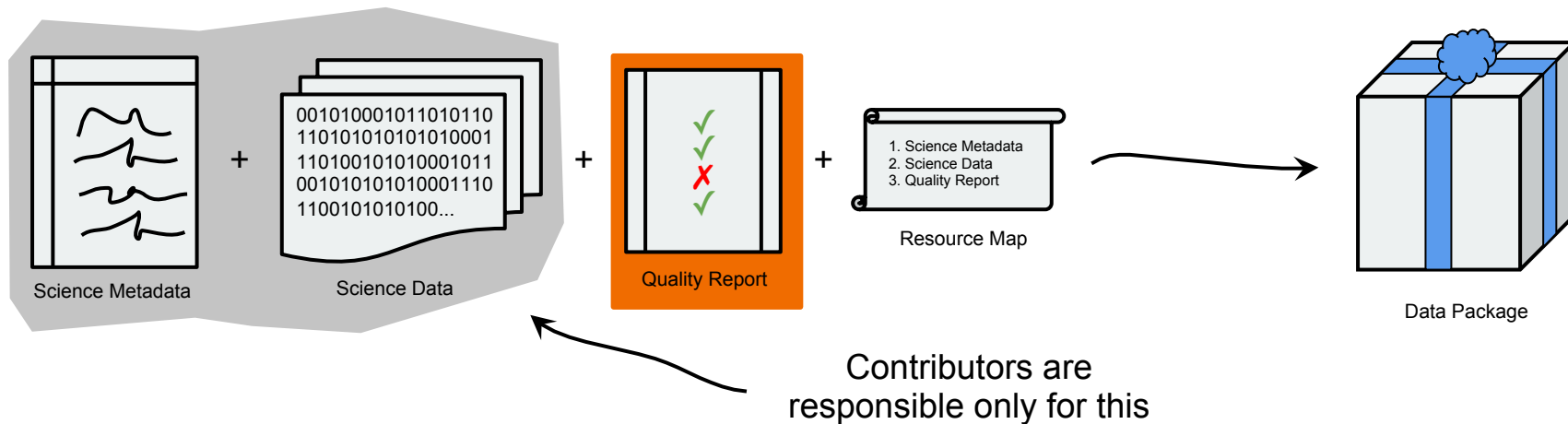  - EDI Data Portal provides a GUI interface to PASTA web services

# About the EDI Data Repository (2 of 2)

- Content summary
  - 42,000+ data packages
  - ~9TB data volume
- DataONE member node operator for the EDI MN and LTER MN (GMN)
- DOIs minted for all public data packages through DataCite
  - **10.6073/pasta/c868deb66af9c0f72c9672f65b484995**
  - Assigned at the data package level, not assigned to its individual components
- Ensures data integrity through checksums and multiple backup strategies (AWS Glacier)
- Supports upload notification to drive workflows (e.g. Twitter bot)
- Funded by US National Science Foundation - DEB/LTER/ABI and ARRA

# What's A Data Package?

**Data Package** (noun): an assemblage of science metadata (e.g., EML) and one or more science data objects; PASTA data packages include a "quality report" object and are described by package metadata called a "resource map" (i.e., manifest)

# Identifiers in the EDI Data Repository

- PASTA data package identifier
  - edi.12.1 (internal)
  - https://pasta.lternet.edu/package/eml/edi/12/1 (external, public facing)
- DOI
  - doi:10.6073/pasta/43e9a619cbebb98da0011ada25ad5c12
  - https://dx.doi.org/10.6073/pasta/43e9a619cbebb98da0011ada25ad5c12
- DataONE identifier
  - doi:10.6073/pasta/43e9a619cbebb98da0011ada25ad5c12 (data package)
- Identifiers support
  - immutability
  - strong versioning
  - reliable access

# PASTA Congruence/Quality Checking

- Are all required/recommended metadata fields present?

- Do metadata values comply with best practices?

- Is the data available for download?

- Does metadata accurately describe the data, i.e. are the two congruent?

<p style="text-align:center">info, valid, warn, error</p>

**PackageId: knb-lter-nwk.1424.50**

**Report Date/Time: 2017-11-22T20:38:36**

**Dataset Report**

| # | Identifier | Status | Quality Check | Name | Description | Expected | Found |
|---|---|---|---|---|---|---|---|
| 1 | packageIdPattern | valid | Type: metadata<br>System: lter<br>On Failure: error | packageId pattern matches "scope.identifier.revision" | Check against LTER requirements for scope.identifier.revision | 'scope.n.m', where 'n' and 'm' are integers and 'scope' is one of an allowed set of values | knb-lter-nwk.1424.50 |
| | emlVersion | valid | Type: metadata | EML version 2.1.0 or | Check the EML | | eml://ecoinformatics.org/eml 2.1.1 namespace |
| 6 | keywordPresent | warn | Type: metadata<br>System: lter<br>On Failure: warn | keyword element is present | Checks to see if at least one keyword is present | Presence of one or more keyword elements | 0 'keyword' element(s) found |
| 7 | methodsElementPresent | valid | Type: metadata<br>System: lter<br>On Failure: warn | A 'methods' element is present | All datasets should contain a 'methods' element, at a minimum a link to a separate methods doc. | presence of 'methods' at one or more xpaths. | 2 'methods' element(s) found |
| 8 | coveragePresent | warn | Type: metadata<br>System: lter<br>On Failure: warn | coverage element is present | At least one coverage element should be present in a dataset. | At least one of geographicCoverage, taxonomicCoverage, or temporalCoverage is present in the EML. | 0 'coverage' element(s) found |
| 9 | geographicCoveragePresent | info | Type: metadata<br>System: lter | geographicCoverage is present | Check that geographicCoverage exists in EML at the | geographicCoverage at least at the dataset level. | 0 'geographicCoverage' element(s) found |
| 12 | onlineURLs | valid | Type: congruency<br>System: knb<br>On Failure: error | Online URLs | URLs return something | | |
| 13 | integrityChecksum | error | Type: congruency<br>System: lter<br>On Failure: error | Compare the metadata checksum for an entity to the checksum of the downloaded entity | Two possible responses: valid if checksums match; error if checksums do not match. | 915a52bd06ef5730ca5ef33dd359380a59c86ef5 | 815a52bd06ef5730ca5ef33dd359380a59c86ef4 |

# Data Package
##   Landing Page

Citation suggestion

Provenance information

Code generation



**Data Package Summary**   View Full Metadata

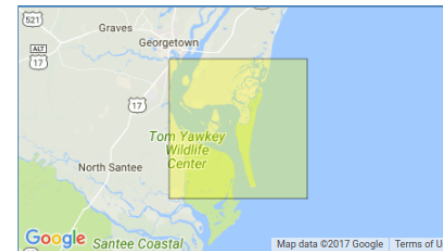| | |
|---|---|
| Title: | Daily Water Sample Nutrient Data for North Inlet Estuary, South Carolina, from 1978 to 1992, North Inlet LTER |
| Creators: | Gaucho, Chase; Conquerors of the Useless |
| Publication Date: | 2006 |
| Citation: | Gaucho C. 2006. Daily Water Sample Nutrient Data for North Inlet Estuary, South Carolina, from 1978 to 1992, North Inlet LTER. Environmental Data Initiative. http://dx.doi.org/10.5072/FK2/f77109ef542e04740af0b4521ebfae08. Dataset accessed 11/15/2017. |
| Abstract: | This data package consists of Daily Water Sample Nutrient Data for North Inlet Estuary, South Carolina, from 1978 to 1992, North Inlet LTER. Its purpose is to establish a long term data base on the nutrient dynamics of a... Show more > |
| Spatial Coverage: | |

N: 33.357     S: 33.1925     E: -79.1002     W: -79.2936

| | |
|---|---|
| Package ID: | edi.98.1 |
| Resources: | Metadata |
| | Report |
| | Data |
| | 1. DailyWaterSample-NIN-LTER-1978-1992  *(902K)* |

[Download Zip Archive]

| | |
|---|---|
| Intellectual Rights: | LTER Network Data Access Requirements The access to all LTER data is subject to requirements set forth by this policy document to enable data providers to track usage, evaluate its impact in the community, and confirm us... Show more > |
| Digital Object Identifier: | doi:10.5072/FK2/f77109ef542e04740af0b4521ebfae08 |
| PASTA Identifier: | https://pasta-s.lternet.edu/package/eml/edi/98/1 |
| Provenance: | Generate provenance metadata for use within your derived data package |
| Code Generation: | Analyze this data package using Matlab, R, SAS, SPSS |

# Overview

EDI's mission and approaches - Corinna Gries

Data repository and publishing - Duane Costa

**Data publication workflow support - Colin Smith**

Outreach and training - Kristin Vanderbilt

Data integration support - Margaret O'Brien

https://environmentaldatainitiative.org/          https://portal.edirepository.org          @EDIgotdata          edi-got-data

# Data submitted to EDI

Long term data sets - regular updating (e.g. LTER, LTREB, OBFS)

Data associated with one project (e.g., MSB, other short term funding)

Data associated with a publication (may be subsets, usually are well curated)

# Data policies in EDI

Data have to be submitted at the same time as metadata

- Except human subject data
- Special access to sensitive data (e.g., endangered species location) can be arranged, but data need to be on EDI server

Embargo time during paper review can be arranged

Data are licensed to be in the public domain or require attribution

# Data publication workflow support

Best practices for organizing, cleaning, and documenting data

Software to help organize, clean, and document data

One-on-one support for any step of this process

# Software to format, clean, and document data

Data formatting tools to help reformat data into a standardized structure

Data cleaning tools

Data documentation tools to help create high quality EML metadata

# Metadata generation

Metadata template, a Microsoft Word document completed by the data provider and converted to EML by EDI's data curation team



EML Assembly Line, a user friendly R library for data providers to generate EML metadata on their own

# Data and metadata upload

Via the data curation team

Via the data provider with a user account

# Overview

EDI's mission and approaches - Corinna Gries

Data repository and publishing - Duane Costa

Data publication workflow support - Colin Smith

**Outreach and training - Kristin Vanderbilt**

Data integration support - Margaret O'Brien

https://environmentaldatainitiative.org/        https://portal.edirepository.org        @EDIgotdata        edi-got-data

# Outreach:

Website:

https://environmentaldatainitiative.org/

Twitter:  @EDIgotdata

# Outreach:

Website:

https://environmentaldatainitiative.org/

Twitter: @EDIgotdata

# Newsletter

## ENVIRONMENTAL DATA INITIATIVE



## Newsletter . October 2017

### HIGHLIGHTS

EDI is happy to announce funding in support of 3 data management internships for the summer of 2018. For more information see here.

On 11 October 2017, EDI had the first meeting with its Scientific Advisory Board (Peter Arzberger, Nathan Booth, Aaron Ellison, Ian Foster, Rebecca Koskela and Mary Martin).

### FEATURED DATA CONTRIBUTION

Climate history of the northeastern United States during the past 3000 years

Marlon, J., et al. (2017) Environmental Data Initiative, http://dx.doi.org/10.6073/pasta/6ced33c5e07f9fa7f11efb259001bacb.

### UPCOMING EVENTS

For details on our upcoming events see here:

# Newsletter

## ENVIRONMENTAL DATA INITIATIVE

### Newsletter . October 2017

**HIGHLIGHTS**

EDI is happy to announce funding in support of 3 data management internships for the summer of 2018. For more information see here.

On 11 October 2017, EDI had the first meeting with its Scientific Advisory Board (Peter Arzberger, Nathan Booth, Aaron Ellison, Ian Foster, Rebecca Koskela and Mary Martin).

**FEATURED DATA CONTRIBUTION**

Climate history of the northeastern United States during the past 3000 years

Marlon, J., et al. (2017) Environmental Data Initiative, http://dx.doi.org/10.6073/pasta/6ced33c5e07f9fa7f11efb259001bacb.

**UPCOMING EVENTS**

For details on our upcoming events see here:

# Student Internships

- Project leaders apply to host a student to work with particular data sets
- Students will
  - Learn to create metadata, quality control data and archive data packages in the EDI repository
  - Be trained and mentored by EDI
  - Conduct a small data analysis project
- $5000  scholarship for the summer

# Information Management Training: In-Person

- Technical
  - IMs from OBFS, LTREB, LTER (2017)
  - EDI interns/OBFS/LTREB/LTER (2018)
  - Early Career faculty and their lab members (2018)
- Overview
  - GLEON scientists (2017)
  - OBFS site managers (2018)

# Info Management Training: Online

- VTC (Git, EML Creation Using R, PASTA+ REST API)
  - Slides and YouTube videos
- Upcoming: Five Phases of Data Publishing
  - January 30: "What are Clean Data?"
- EDI data managers available to help!

Previous EDI Events

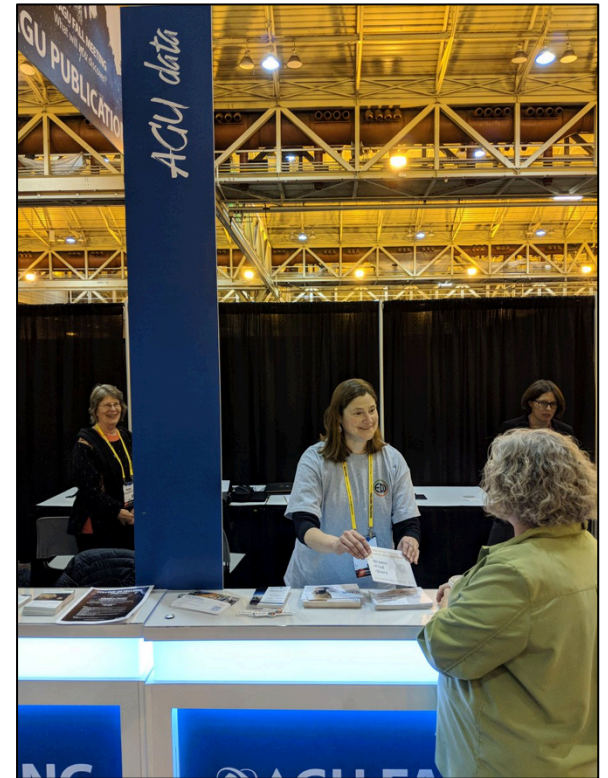The Environmental Data Initiative is on YouTube! Stop by for past tutorials and discussions on information management in the environmental sciences here.

**VTC Discussions &Tutorials**

- VTC tutorial: "DEIMS, the Drupal Ecological Information Management System", 5 December 2017.
- VTC tutorial: "Using checksums to speed up data package uploads", 28 November 2017.
- VTC tutorial (2): "Transform and visualize data in R using the packages tidyr, dplyr and ggplot2", 24 October 2017.
- VTC tutorial (1): "Transform and visualize data in R using the packages tidyr, dplyr and ggplot2", 17 October 2017.
- VTC Discussion: "Data package design, featuring the candidate model for community survey data", September 26, 2017
- VTC Tutorial: "The PASTA+ Rest API" (part 2 of 2), July 18, 2017
- VTC Tutorial: "The PASTA+ Rest API" (part 1 of 2), July 11, 2017
- VTC Tutorial: "Creating EML with R and sharing on GitHub", June 27, 2017
- VTC Tutorial: "Git and GitHub", June 20, 2017
- VTC Tutorial: "R basics", June 13, 2017

# Other Outreach:

- National and International Meetings (OBFS, ESA, AGU, ILTER, ESIP)
- Targeted seminars
- Code Repository

# Overview

EDI's mission and approaches - Corinna Gries

Data repository and publishing - Duane Costa

Data publication workflow support - Colin Smith

Outreach and training - Kristin Vanderbilt

**Data integration support - Margaret O'Brien**

https://environmentaldatainitiative.org/          https://portal.edirepository.org          @EDIgotdata          edi-got-data

# Data Integration Support

## Rationale

As archived, primary data are often difficult to use in synthesis  - *even with complete metadata*
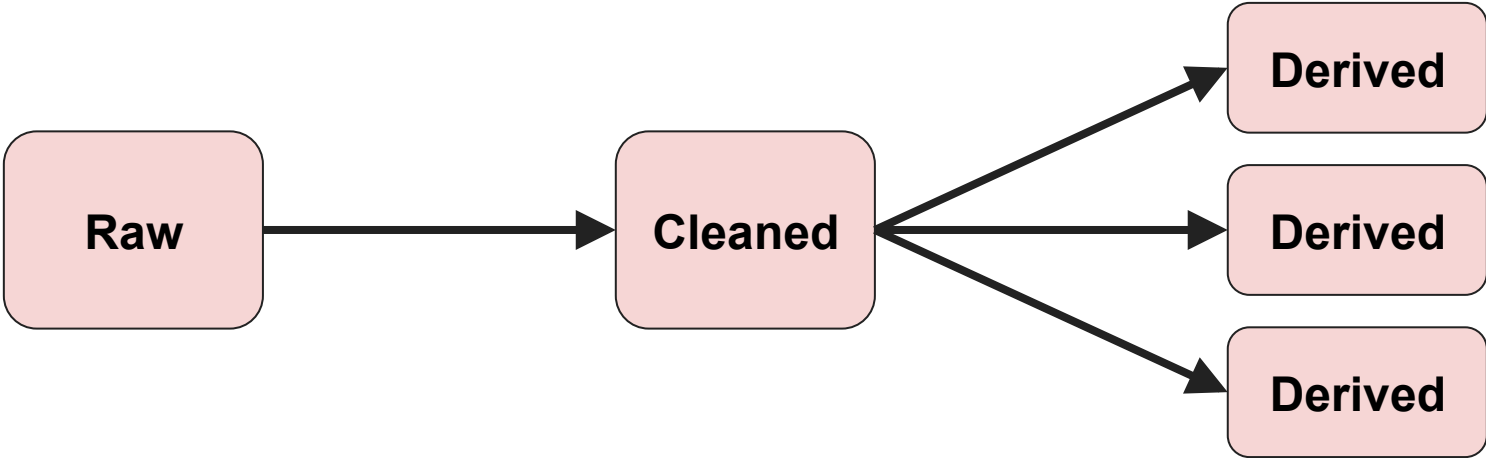
## EDI Provides

Expertise from data managers to define best practices for data package design

## Approach

1) Recommend format and content so analyses can be streamlined
    a) Mechanism for preparers to know
        i) Data elements that are most important
        ii) Layouts that are easiest to use
2) Work with scientists currently engaged in synthesis of primary data

*Template for a process that can be reused in many scientific domains*

# Typical Synthesis Workflow

# Ideal Synthesis Workflow

**Level 0**

**Raw**

Step 1

Custom code

**Level 1**

**Based on predefined model**

Data package predefined model, with provenance

Step 2

EDI creates code to
- Query for packages in a known model
- Combine datasets
- Query for further use (by scientists)

**Level 2**

**Derived**

**Derived**

**Derived**

Data package(s) with provenance

# Objectives - Design Pattern for Level 1 Dataset

Flexible format for a particular domain of data, but usable for multiple types of measurements and synthesis projects
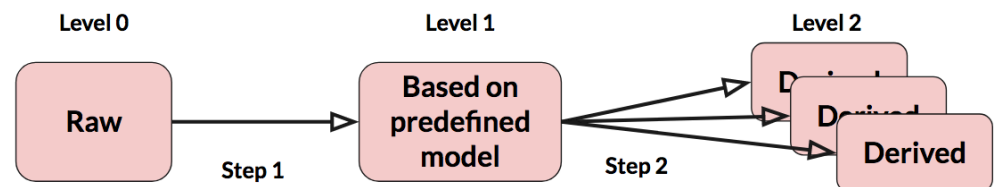
Metadata in EML

Reformat only, no calculations, aggregations or other steps requiring ecological judgement

Original data referenced, one to one relationship with reformatted Level 1 data

Complete; original records can be recreated

Distributed, not complex database application for maintenance
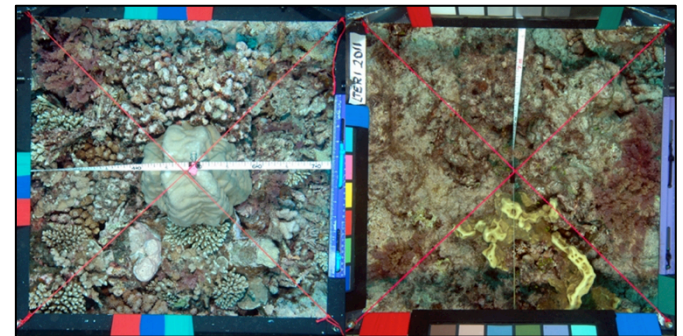
## Basic Process

1. Examine available models currently in use

2. Examine *ad hoc* cleaned (Level 1) data created by a synthesis working group

3. Describe patterns

4. Define design pattern tables and value typing

5. Test model against data of interest

6. Create utility scripts for QC, metadata generation

# Example - Design Pattern for Community Surveys

1. Examine available models currently in use: **Popler, Darwin Core, CUAHSI ODM2**

2. Examine *ad hoc* cleaned (Level 1) data created by a synthesis working group: **LTER Synthesis WG: Synchrony, Metacommunities**

3. Describe patterns: **quite similar to DC-A but with many custom fields, especially for spatial info**

4. Define design pattern tables, typing - "**ecocomDP**"

5. Test model against data of interest

6. Create utility scripts for QC, metadata generation



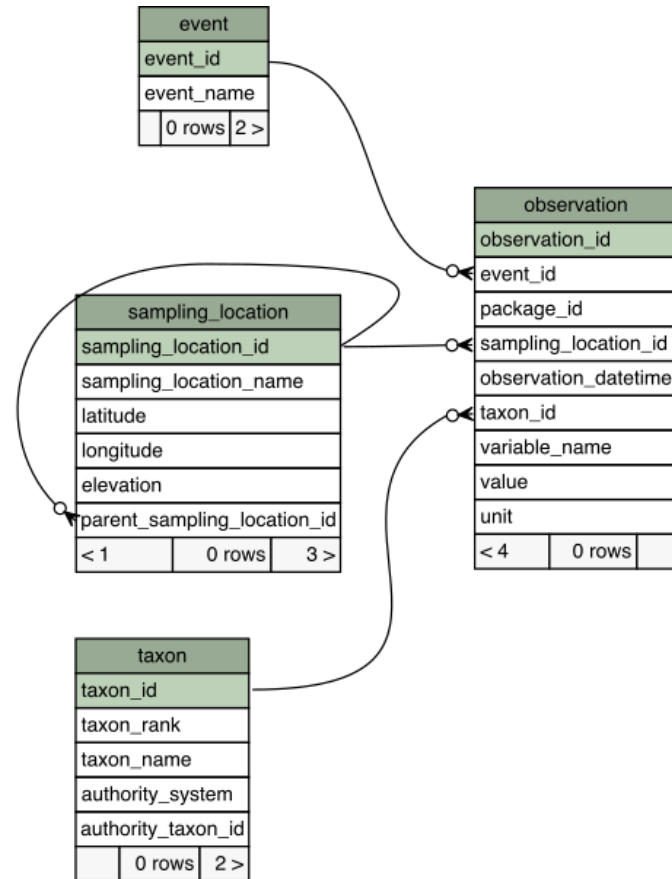*Benthic cover at a Moorea Coral Reef (credit, http://lternet.edu)*

# ecocomDP

**Observation table** for data related to
Count, biomass, abundance, density
Primary organization
Entity, name, value, unit (EAV, U)

**Linked Tables** for
Sampling location
Organism
Event

*Model includes ancillary tables and summary
(https://github.com/EDIorg/ecocomDP)*

# Progress

Creating Level 1
- Usability
- Metadata generation

Using Level 1
- Query scripts



*Woody encroachment at Konza Prarie (credit, http://lternet.edu)*

# Thank you

- Contact - info@environmentaldatainitiative.org

- Website - https://environmentaldatainitiative.org/

- Data portal - https://portal.edirepository.org

- Twitter - @EDIgotdata

- Slack - edi-got-data

- GitHub - https://github.com/EDIorg   https://github.com/PASTAplus/PASTA

- PASTA+ User/Developer Documentation - http://pastaplus-core.readthedocs.io/en/latest/

- Data Package Manager Web Service API - https://pasta.lternet.edu/package/docs/api

- Audit Manager Web Service API - https://pasta.lternet.edu/audit/docs/api