



Data Harmonization Facilitates Interoperability and Reuse

Colin Smith¹, Margaret O'Brien², Corinna Gries¹

¹University of Wisconsin, Madison, ²University of California, Santa Barbara



Abstract

Data repositories and research networks publish a diverse array of primary data for meaningful synthesis, integration and future reuse. However, in synthesis research the largest time investment is in cleaning and combining primary datasets until all are completely understood and converted to a usable format. To accelerate this process, EDI defines flexible domain-specific data models, and converts primary data to these models using a lightweight and distributed workflow framework.

Advantages of a Harmonization Workflow

- Original data description and curation practices are maintained
- Workflow framework is repeatable (essential for ongoing datasets)
- Intermediate format is not determined by a single synthesis research question
- Most harmonization steps can be performed by non-specialists

3-Phase Process

Design

- Capture essential attributes from the community
- Consider existing standardization efforts
- Evaluate external vocabularies to disambiguate meaning.



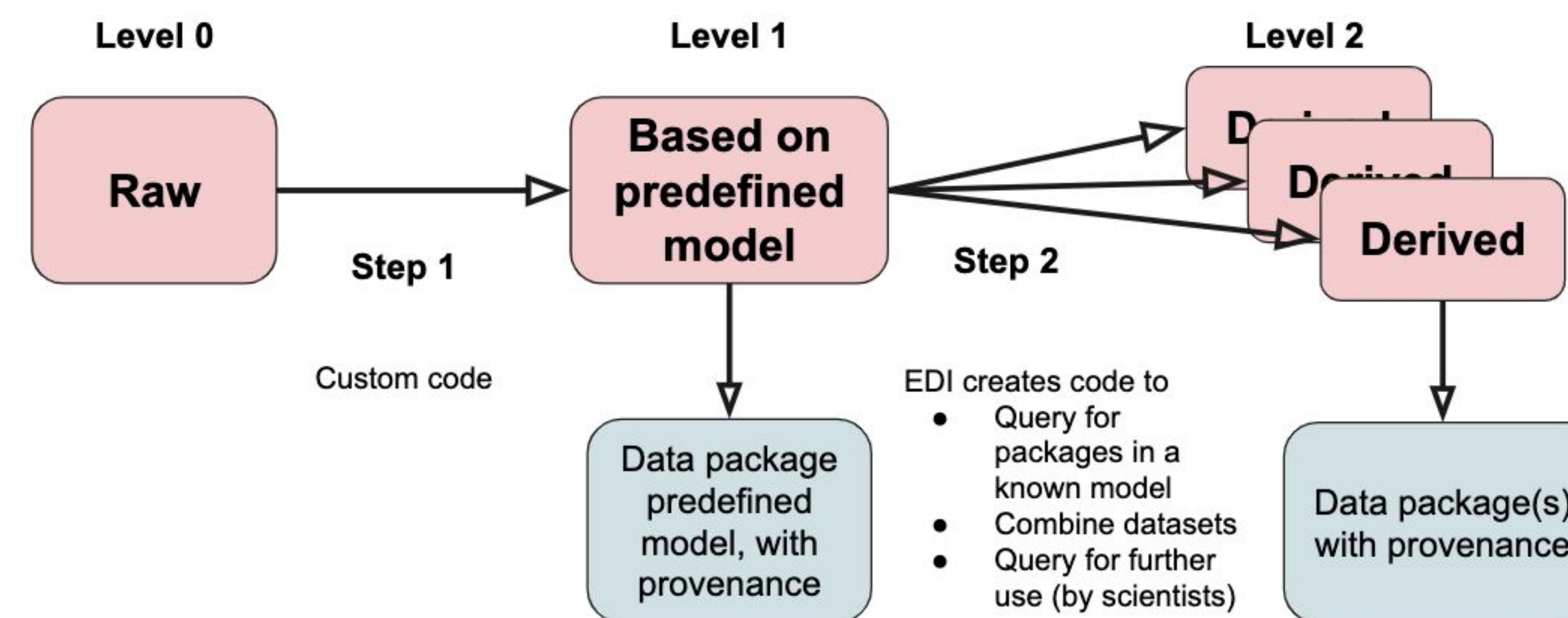
Implementation

- Distribute the data model; convert relevant datasets
- Create
 - templates for building data packages
 - best practice guides for reference
 - software for discovery, exploration, denormalization

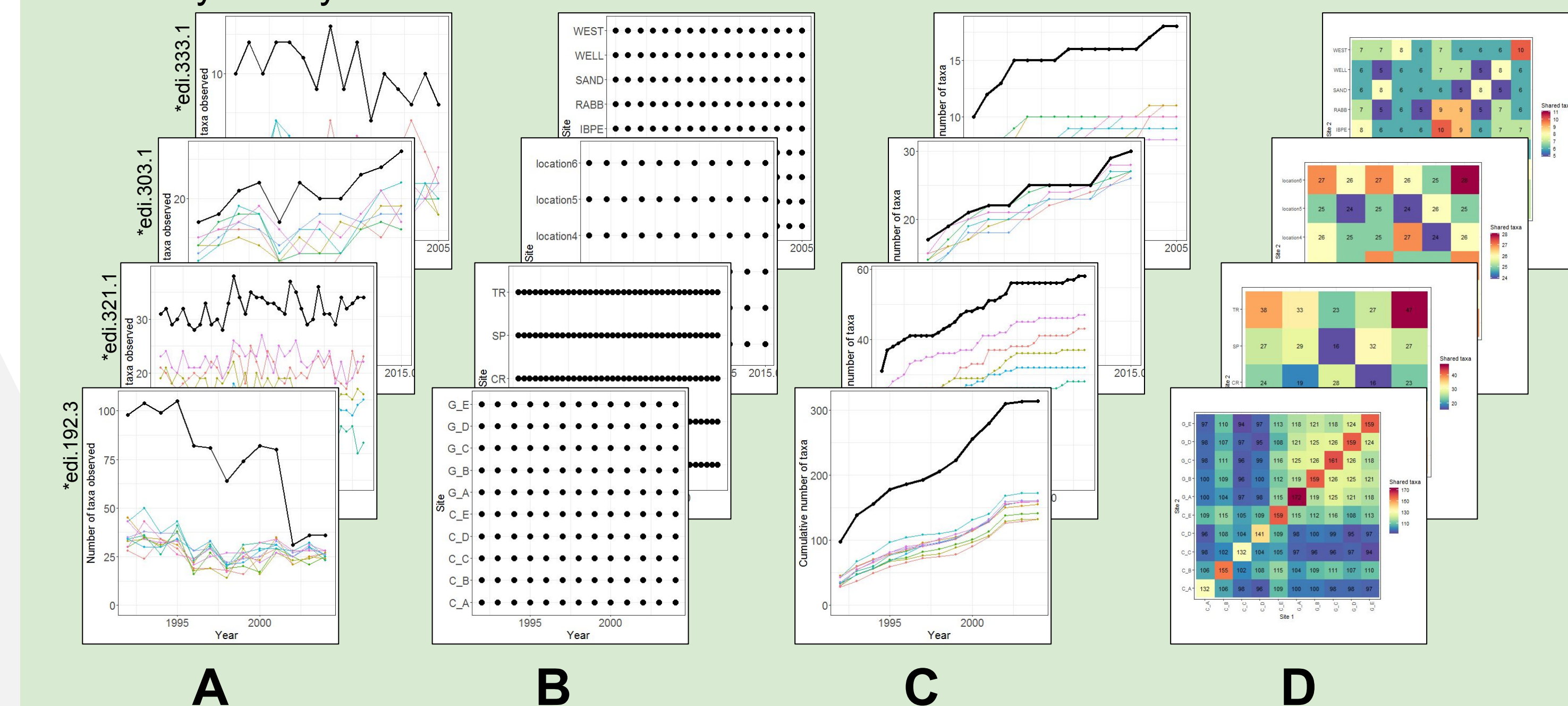
Maintenance

- Build workflows to run when source data are updated
- Automate with event notification services

Harmonization Framework



Harmonized data (Level 1) allows common visualization of many datasets to evaluate features such as number of taxa over time (A), spatio-temporal sampling effort (B), species accumulation curves (C) and species shared among sites (D). These plots are from community survey data harmonized into EDI's ecocomDP model.



*All datasets can be accessed with prefix "https://portal.edirepository.org/nis/mapbrowse?packageid="

Results

Ecological community surveys

Design and Implementation in 2017 & 2018. Today, thousands of records are available as **Level 1**, in EDI design pattern for ecological community data, "ecocomDP". 2019: Maintenance phase.

Metrics of datasets converted to date

non-NEON (N= 70 datasets)					NEON (N = 1)
	N	Min	Max	Median	
Temporal coverage (years)	70	4	70	14	4
Temporal evenness (interval SD)	69	0	10.8	0.05	.93
Geographic coverage (km ² , > 0)	70	1368	1.3 x 10 ¹⁴	1.9 x 10 ⁸	NA
Taxonomic coverage (without OTUs)	69	1	1752	48	1066

Meteorology & hydrology data

Design phase: 2018/2019

- Working group recommends **CUAHSI ODM 1.1** model

Implementation phase: 2019

- Sample datasets converted
- Reference guides curated

Maintenance phase: 2020 (planned)

- **LTER & USFS** early adopters

Use cases

Research synthesis

Up to 80% of synthesis work can be related to cleaning, interpreting and reformatting input data. A harmonization workflow reduces that effort considerably.

Interoperability

A Level 1 data model facilitates integration.

- NEON collaborates with EDI to include NEON biodiversity data in the discovery functionality of ecocomDP's R tools
- EDI is exploring integration with the Popler model (<https://github.com/ropensci/popler>)

Targeted Applications

Harmonization is a necessary step in facilitating contributions to external applications, including

- GBIF, for biodiversity data: <https://www.gbif.org>
- CUAHSI for hydrology data: <https://www.cuahsi.org>

<http://environmentaldatainitiative.org>

