



Ensuring the Quality of Data Packages in the LTER Network Data Management System



Mark Servilla¹, Margaret O'Brien², Duane Costa¹

¹LTER Network Office, University of New Mexico; ²Marine Science Institute, University of California Santa Barbara

Summary

Considerable data analyses use automated workflows to ingest data from public repositories, and rely on data packages of high structural quality. The Long Term Ecological Research (LTER) Network now screens all packages entering its long-term archive to ensure completeness and quality, and requires that metadata and data are structurally congruent.

1. Component of the LTER Provenance Aware Synthesis Tracking Architecture (PASTA)

Operates on data packages described with Ecological Metadata Language (EML), using the EML Data Manager Library (DML), written in Java

Checking is extensible for other data-types and customizable via a template

2. Current Implementation

- The EML metadata specification is widely used to describe environmental data
- The DML code reads EML-described data tables into a relational database for analysis. It also checks features of certain important metadata elements such as title, abstract and data URL.
- Checks are designed to evaluate metadata/data for general use (i.e., other than in a database) and the system is extensible for data types other than tables.

B. Configuration template:

In the template XML, an operator configures the System (scope) and Response status. See `<includeSystem>` `qualityCheck/@system` and `qualityCheck/@statusType`

```
<?xml version="1.0" encoding="UTF-8"?>
<qr:qualityReport
  xmlns:qr="http://ecoinformatics.org/qualityReport"
  xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xsi:schemaLocation="http://ecoinformatics.org/qualityReport http://svn.lternet.edu/svn/NIS/documents/schemas/qualityReport.xsd"
  <createDate>2011-12-01T12:00:00Z</createDate>
  <packageId></packageId>
  <includeSystem>knb</includeSystem>
  <includeSystem>lter</includeSystem>
  <datasetReport>
    <qualityCheck qualityType="metadata" system="lter" statusType="error" >
      <identifier>emlVersion</identifier>
      <name>EML version 2.1.0 or beyond</name>
      <description>Check the EML document declaration for version 2.1.0 or higher</description>
      <expected>eml://ecoinformatics.org/eml-2.1.0 or eml://ecoinformatics.org/eml-2.1.1/expected</expected>
      <found></found>
      <status>metChecked</status>
      <explanation>Validity of this quality report is dependent on this check being valid.</explanation>
      <suggestion>Use an approved namespace.</suggestion>
      <reference></reference>
    </qualityCheck>
    <qualityCheck qualityType="metadata" system="knb" statusType="error" >
      <identifier>schemaValid</identifier>
      <name>Document is schema-valid EML</name>
      <description>Check document schema validity</description>
      <expected>eml://ecoinformatics.org/eml-2.1.0/expected</expected>
      <found></found>
      <status>metChecked</status>
      <explanation>Validity of this quality report is dependent on this check being valid.</explanation>
      <suggestion>Use an approved namespace.</suggestion>
      <reference></reference>
    </qualityCheck>
  </datasetReport>
</qr:qualityReport>
```

A. Check classification

Scope: the domain of the check

- knb, lter, (or other communities)

Type: package component(s) being checked

- metadata, data, metadata:data congruency

Response status:

- info** for information only, does not affect acceptance by system
- Or one which controls system behavior:
 - valid** all check-criteria were met
 - warn** some problem may be present, but data package is acceptable to the system (eg, PASTA)
 - error** data package cannot be accepted

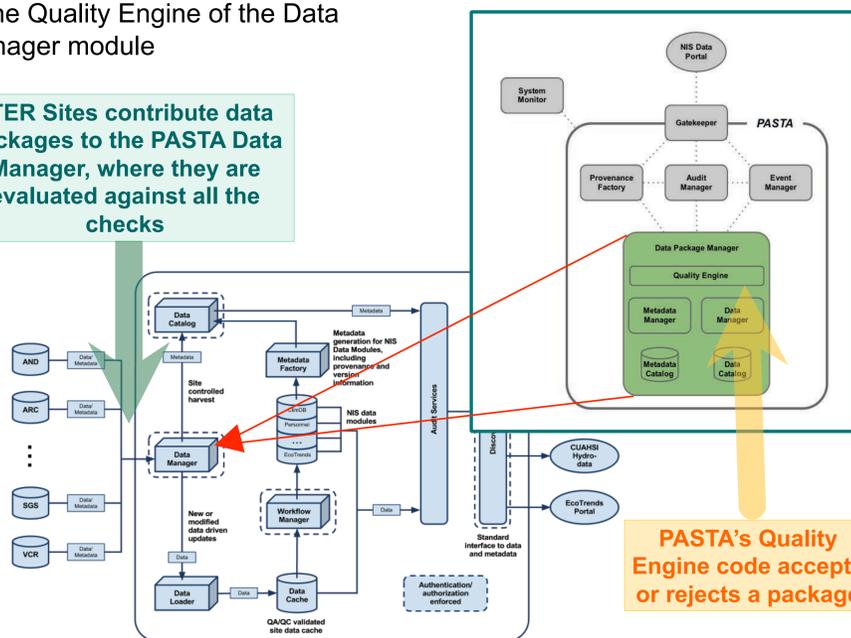
C. Code populates appropriate fields (e.g. `<packageID>`, `<found>` and `<status>`) with results and creates an XML file called a "Quality Report", using the same schema. Reports are included with metadata and data in the package's resource map, and are available to data package users.



1. LTER PASTA Architecture

The quality checks are executed as data packages are inserted into PASTA, in the Quality Engine of the Data Manager module

LTER Sites contribute data packages to the PASTA Data Manager, where they are evaluated against all the checks



PASTA's Quality Engine code accepts or rejects a package

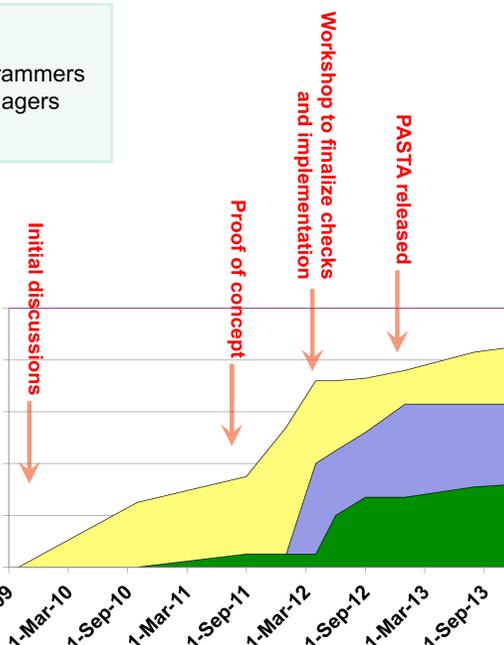
3. Community Involvement

Contributors

LTER Network Office programmers
LTER Site information managers
NCEAS programmers

Progress describing and implementing checks

To date, the community has entered 83 potential checks, and developers have implemented 32. 21 additional checks have been prioritized and fully.



User community requests:

- HTML interface with transformation of Report results
- Checks which return "error" should be implemented first
 - Submitters will know which are most important issues to address
- Code should operate in two modes
 - Evaluate
 - Checking continues after a failure so that a submitter sees as many problems as possible all at once
 - Harvest
 - Checking stops on the first error and the data package is rejected

#	Identifier	Status	Quality Check	Name	Description	Expected	Found	Dataset Report
1	packagePattern	valid	Type: metadata System: lter On Failure: warn	packageId pattern matches "scope identifier" revision	Check against LTER requirements for scope identifier revision	scope identifier revision is one of an allowed set of values	knb-ter-abc-52.1	
2	emlVersion	valid	Type: metadata System: lter On Failure: error	EML version 2.1.0 or beyond	Check the EML document declaration for version 2.1.0 or higher	eml://ecoinformatics.org/eml-2.1.0 or eml://ecoinformatics.org/eml-2.1.1/expected		Validity of this quality report is dependent on this check being valid.
3	schemaValid	error	Type: metadata System: lter On Failure: error	Document is schema-valid EML	Check document schema validity	schema-valid		Validity of this quality report is dependent on this check being valid.
4	parserValid	valid	Type: metadata System: lter On Failure: error	Document is EML, parser-valid	Check document using the EML 5.0 parser	Document is EML, parser-valid		Validity of this quality report is dependent on this check being valid.

