# Linking data repositories - an illustration of agile data curation principles through robust documentation and multiple application programming interfaces

*Karl Benedict[1], Mark Stephen Servilla[2], Kristin Vanderbilt[3], Jonathan Wheeler[4]*

## Problem

The replication of data and associated documentation between systems for data security, enhanced discovery and access, or use on the secondary system is increasingly difficult as the number and variety of data objects increases. The development of data management and archiving systems and workflows that lower the barriers to automated movement of data and documentation between systems can enable efficient and scalable replication/migration of data assets between systems and between systems and users. Conceptually, this model of streamlined data management and curation based on sound system design principles may be described as agile data curation.

When confronted with the planned discontinuation of funding for the Sevilleta LTER research site, the LTER Network Infomation Information System designer (Servilla) and the Sevilleta Information Manager (Vanderbilt) approached the research data services team in the UNM Libraries (Benedict & Wheeler) to discuss the options for replication of the Sevilleta data and metadata assets currently managed within the LTER data management system (known as PASTA) into UNM's institutional repository - LoboVault. The workflow that was developed to complete this replication is an excellent example of how efficient data transfer, management and reuse can be facilitated through the development of independent systems that are based upon design principles that enable well-structured machine-to-machine communication (via well documented Application Programming Interfaces - APIs), support of standards based documentation models, and effective communication between project participants throughout the development, testing, and implementation process. While not specifically intended to do so, the process used (and described in this poster) illustrates a number of core principles of agile data curation.
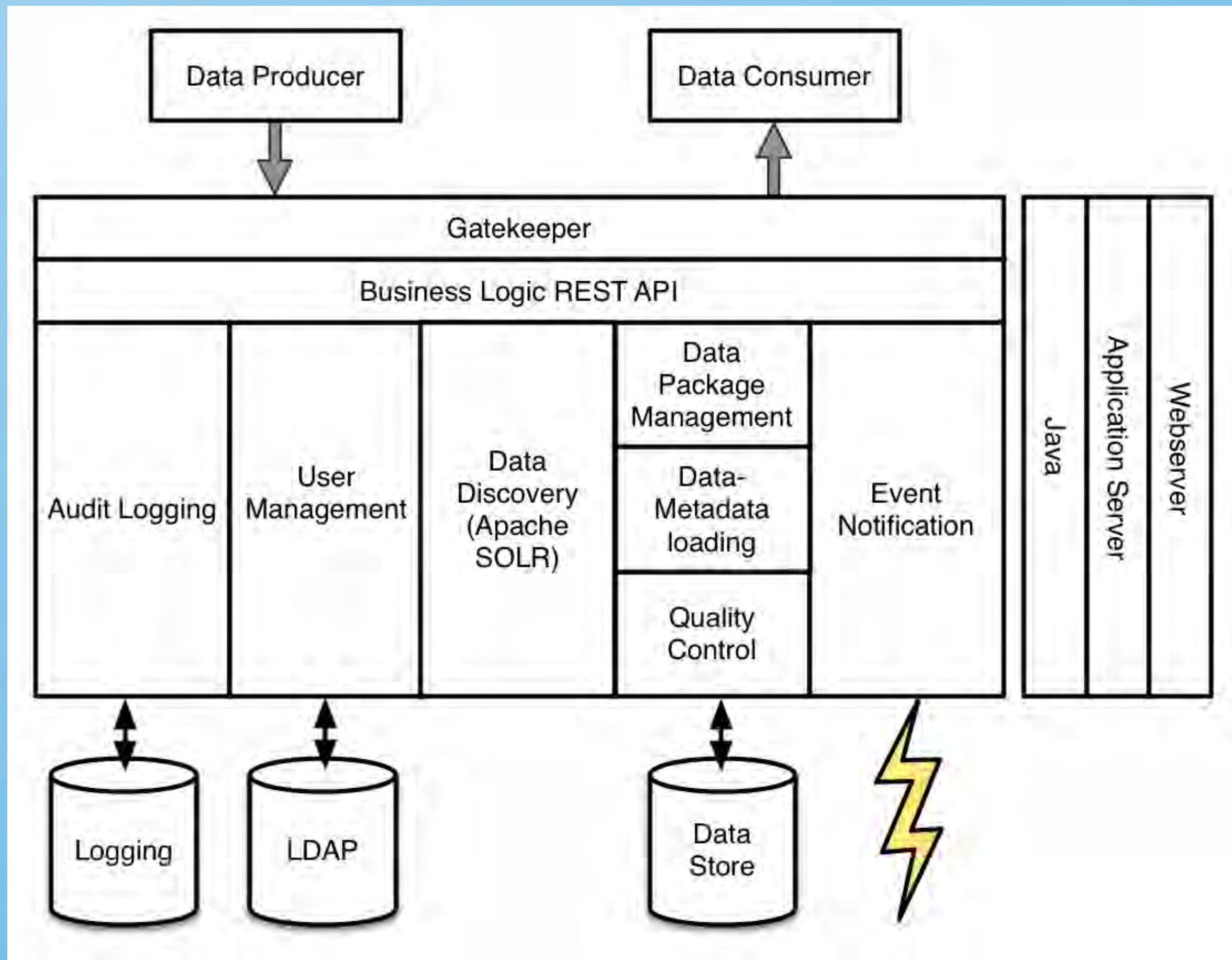
## Integration Workflow Development Process

The integration of the PASTA data packages for the Sevilleta LTER site leveraged the API's of both PASTA and DSpace to automate the process in a way that would meet the immediate needs for the Sevelleta collection while being extensible to all of the collections within PASTA. The specific capabilities of the developed workflow include:

- Leveraging the extensibility of the DSpace metadata model to integrate location data within the PASTA EML metadata into the structured metadata within LoboVault using the Darwin Core extension to the Dublin Core standard
- Maintaining dataset version information provided by PASTA within the LoboVault records
- Conversion of the EML metadata content into corresponding DSpace Dublin Core or Darwin Core elements for search and presentation
- Providing access to the source EML metadata and data objects housed within LoboVault
- Providing an interactive data location map based on information extracted from the source EML metadata

## Lessons Learned

This data and metadata integration project provides both an illustration of many of the agile data curation principles outlined below, but it also highlighted a number of ongoing challenges in addition to the benefits:

Benefits:

- The availability of robust APIs in both PASTA and LoboVault enabled a high degree of automation that provides the promise of significantly streamlined ingest of additional PASTA collections if/when needed
- The detailed and consistent EML metadata standard used in PASTA provided a rich resource of machine-readable documentation that could easily be mapped into LoboVault's metadata model
- The project provided a concrete project around which the LTER and Library data management teams could collaborate - reinforcing the development of a growing research data management community of practice at UNM
- Core PASTA data products are now safely replicated in a complementary system for long-term preservation and access

Challenges:

- The different versioning models implemented within PASTA and LoboVault introduce a step in the processing of materials that hasn't been fully automated yet. This poses a challenge in maintaining synchronization between the systems.
- The communication model based primarily on email exchanges introduced some delays in developing shared understanding that could have been achieved more quickly through in-person interactions.
- Both PASTA and LoboVault support permanent identifiers but the identifiers used in the systems differ in that PASTA uses DOIs and LoboVault (currently) uses Handles. The specification of identifying the canonical identifier remains to be done.
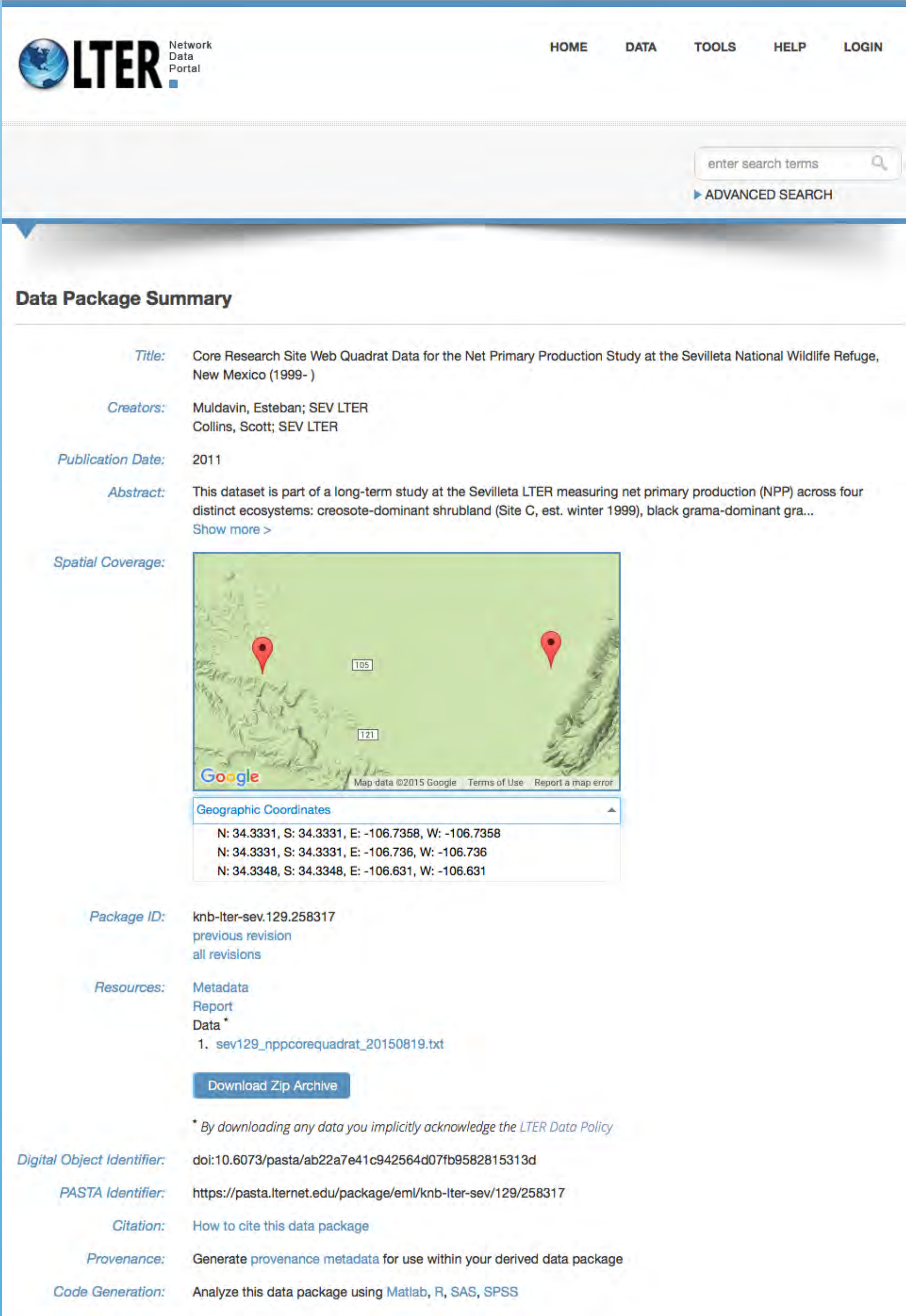


*Figure 1*. PASTA Conceptual Diagram



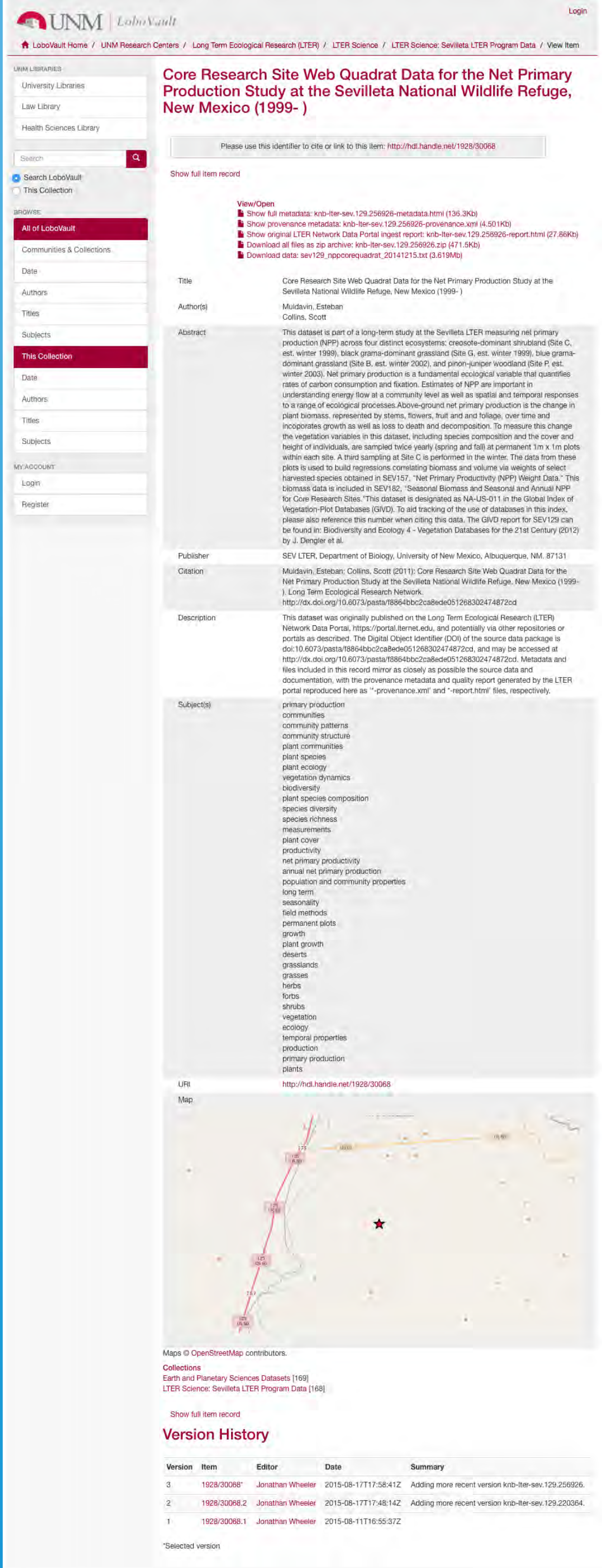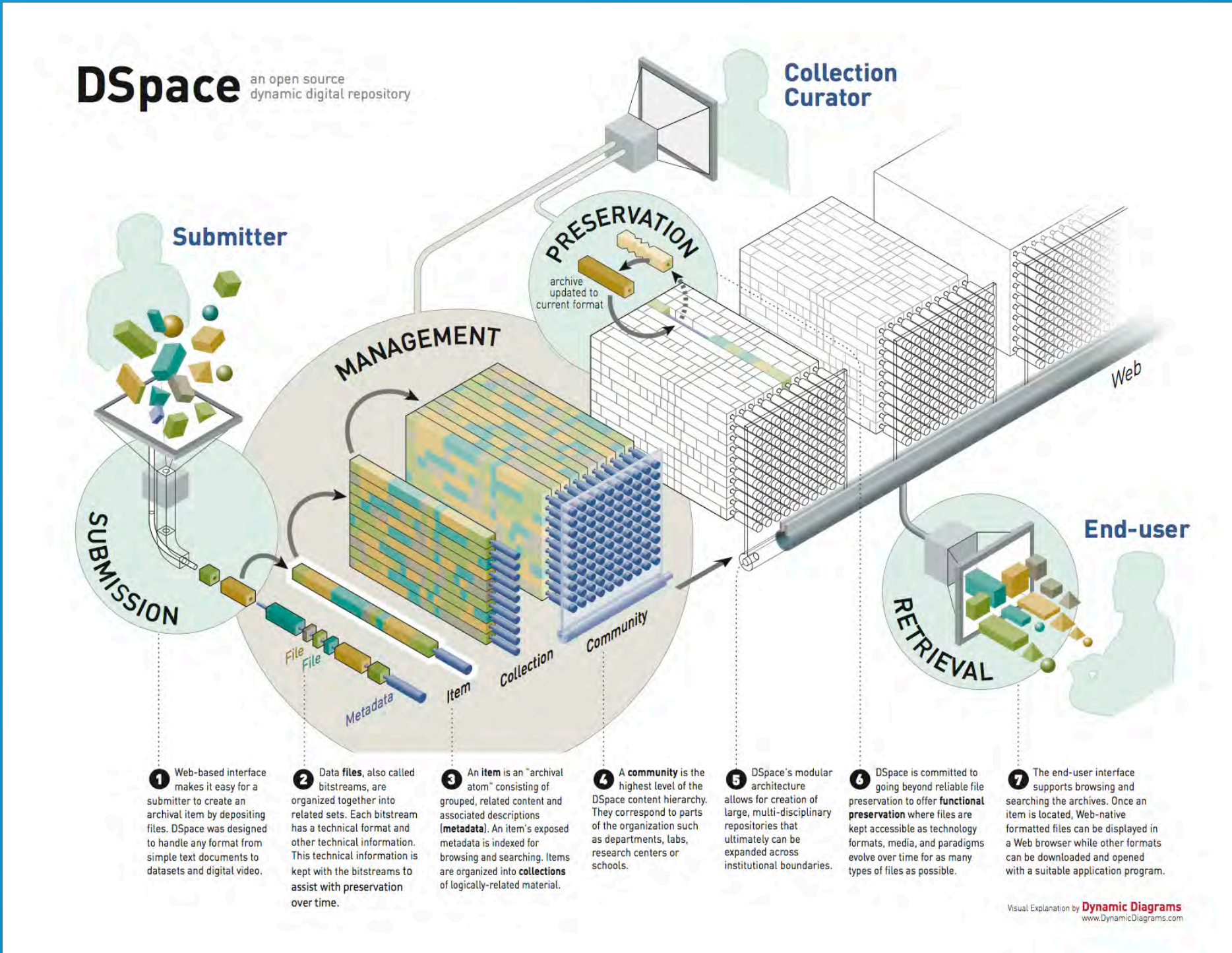*Figure 3*. Sample data package view from the LTER Network Data Portal[14]





*Figure 4*. Sample data package view from the LoboVault LTER collection[15]



*Figure 2*. DSpace Conceptual Diagram[13]

## Design Goals and Attributes of the Provenance Aware Synthesis Tracking Architecture (PASTA)

The PASTA platform was designed as part of the LTER Network Information System[6] to support the data discovery, access and use needs of the community of Long Term Ecological Research sites and the users of the data that they produce. It has been developed with a well documented Application Programming Interface (API[7]) that enables create, evaluate, read, update, delete, list, and search for data package resources in the PASTA system. In particular it has the following characteristics:

- Use of the Ecological Metadata Language[8] as its metadata standard
- Employs the Metacat XML database to support metadata search
- Uses data warehousing methods to provide a uniform representation of the heterogeneous data managed within the system
- Provides a web interface[9] for browsing, searching, and viewing the contents of the PASTA system

## Mapping of Agile Software Development Principles into Data Curation

| Agile Software Development Principles[5] | Agile Data Curation Principles |
|---|---|
| Our highest priority is to satisfy the customer through early and continuous delivery of valuable software. | Maximize the impact of research data through accelerated capacity for discovery, access and use of valuable data |
| Welcome changing requirements, even late in development. Agile processes harness change for the customer's competitive advantage. | Expect unanticipated needs for and uses of research data (and documentation) and develop flexible systems to support new uses and users without significant modifications |
| Deliver working software frequently, from a couple of weeks to a couple of months, with a preference to the shorter timescale. | Facilitate automated interaction with data and metadata assets through well documented public web services that enable disintermediated use and reuse of research data |
| Business people and developers must work together daily throughout the project. | Data creators and data curators should work closely throughout planning, research and preservation activities to ensure the most efficient and streamlined process |
| Build projects around motivated individuals. Give them the environment and support they need, and trust them to get the job done. | Identify key individuals in a data curation project that have the requisite knowledge and motivation to do the job and get out of their way |
| The most efficient and effective method of conveying information to and within a development team is face-to-face conversation. | Identify the most effective method(s) for maintaining close communication and *use* them |
| Working software is the primary measure of progress. | Delivery, access, use and citation of research data are the primary measures of success |
| Agile processes promote sustainable development. The sponsors, developers, and users should be able to maintain a constant pace indefinitely. | Design principles that enable steady delivery of incremental improvements to research data discovery, access and use should be consistent with a sustainable level of effort and funding from sponsors, data creators and curators, and users |
| Continuous attention to technical excellence and good design enhances agility. | Continuous attention to technical excellence and good design enhances agility |
| Simplicity--the art of maximizing the amount of work not done--is essential. | Start with the basics and only make systems more complex as needed, while maintaining a low bar to entry |
| The best architectures, requirements, and designs emerge from self-organizing teams. At regular intervals, the team reflects on how to become more effective, then tunes and adjusts its behavior accordingly. | Continuously work to develop and evolve a community of data providers, curators and users that all participate in the ongoing evolution of the research data systems that they interact with |

## Design Goals and Attributes of the DSpace Digital Repository Platform (LoboVault)

The University of New Mexico Libraries' DSpace[10] based institutional repository, LoboVault, is an open access portal to the scholarly publications and research data of UNM faculty, graduate student theses and dissertations, and university administrative records. As a general purpose repository, LoboVault has been developed to provide an easy to use access and discovery resource for the full array of scholarly content and data types produced at UNM. The recent implementation of collection-scale batch ingest procedures is further designed to facilitate efficient publication of research output subject to federal publication and data sharing requirements. Relevant features include:

- An extensible metadata model based on qualified Dublin Core[11]
- Metadata federation and integration with external systems via OAI-PMH[12]
- Batch ingest, editing, and cross-registration of content and metadata
- As of DSpace v5, a RESTful API supporting search, create, read, update, and delete functions.
- The development of a *sandbox* version of the repository for development and testing of workflows before execution of bulk ingest processes

1. University of New Mexico, University Libraries - kbene@unm.edu↩
2. University of New Mexico, LTER Network Office↩
3. University of New Mexico, Sevilleta LTER Office↩
4. University of New Mexico, University Libraries↩
5. http://agilemanifesto.org/principles.html↩
6. http://lno.lternet.edu/content/network-information-system↩
7. https://pasta.lternet.edu/package/docs/api↩
8. https://knb.ecoinformatics.org/#external//emlparser/docs/index.html↩
9. https://portal.lternet.edu/nis/home.jsp↩
10. http://www.dspace.org/↩
11. http://dublincore.org/↩
12. http://www.openarchives.org/OAI/openarchivesprotocol.html↩
13. http://www.dspace.org/sites/dspace.org/files/media/DSpace%20Diagram_0.pdf - published by dspace.org and shared under a Creative Commons BY-SA license↩
14. https://portal.lternet.edu/nis/mapbrowse?packageid=knb-lter-sev.129.258317↩
15. https://repository.unm.edu/handle/1928/30068↩

## Acknowledgements