

OR 2010, The 5th International Conference on Open Repositories
Madrid, July 6-9, 2010

Curation Micro-Services

A Pipeline Metaphor for Repositories

Stephen Abrams *

Patricia Cruse *

John Kunze *

David Minor †

** UC Curation Center, California Digital Library*

† San Diego Supercomputer Center

University of California

University of California Curation Center (UC3)

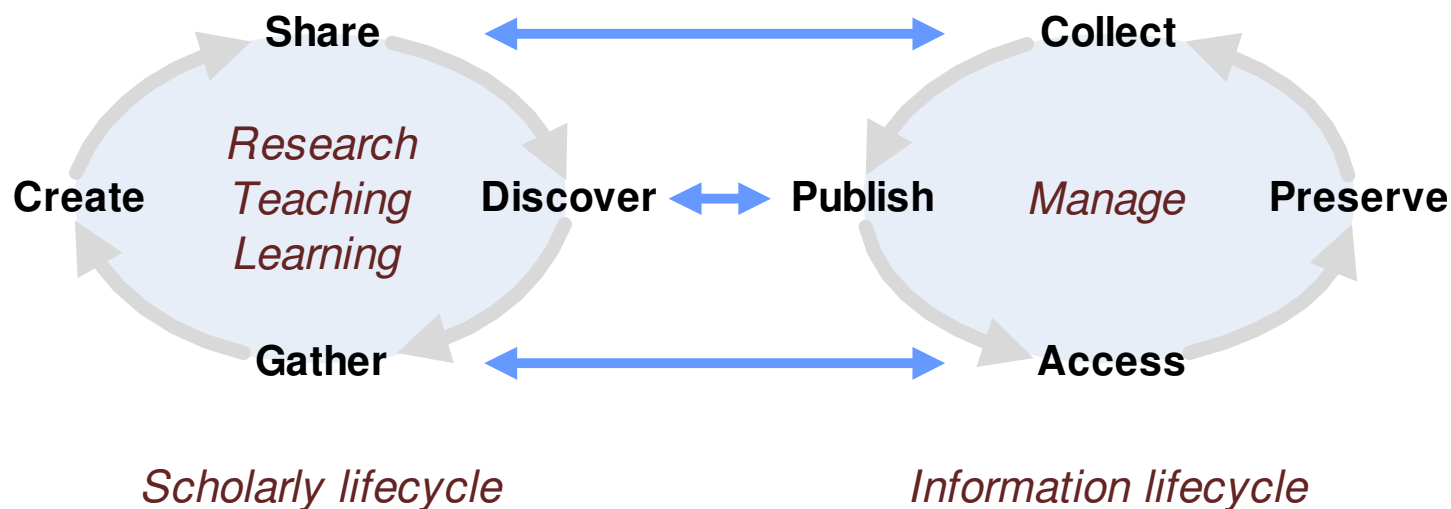
Creative partnership between the CDL, the 10 UC campuses, and peer institutions in the digital curation community

- An evolving community of shared concern and practice
- A means to pool and distribute diverse experience, expertise, and resources
- Robust solutions to counteract inevitable disruptive change



Digital curation

The set of policies and practices focused on *managing and adding value to a body of trusted digital content*, and facilitating the alignment of the scholarly and information lifecycles



Assumptions

Curated content gains

- Safety through redundancy *“Lots of copies keeps stuff safe”*
- Meaning through context *“Lots of description keeps stuff meaningful”*
- Utility through service *“Lots of services keeps stuff useful”*
- Value through use *“Lots of uses keeps stuff valuable”*

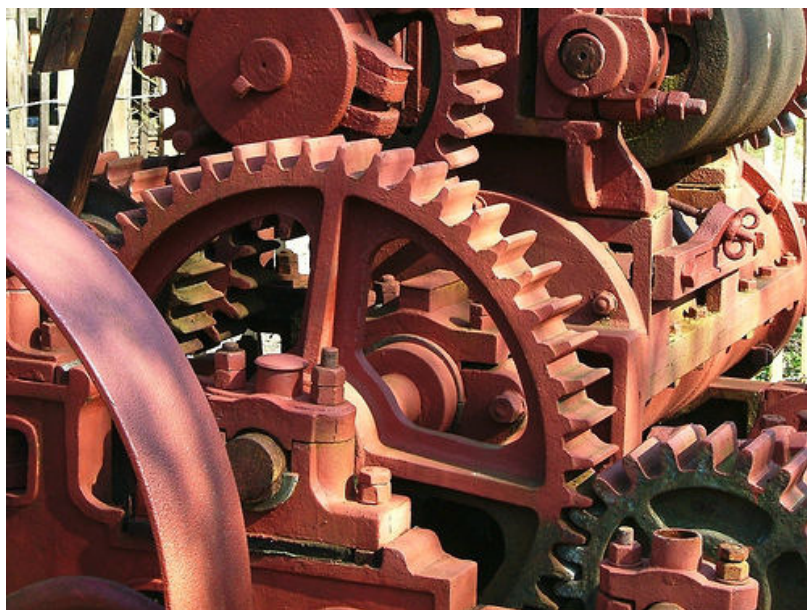
Curation is an outcome, not a place

- Focus on content, not the systems in which that content is managed

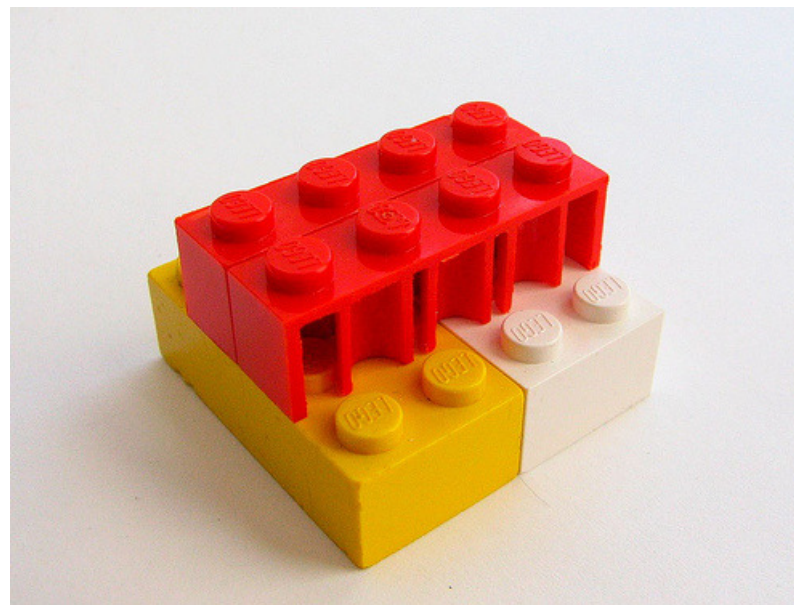
Curation stewardship is a relay

*“A complex system that works is invariably found to have
“Entirely more parts than are necessary for the job”
evolved from a simple system that worked”*

– William of Occam
– John Gall



© <http://www.flickr.com/photos/elsie/8229790/>



© <http://www.flickr.com/photos/oskay/265899811/>

Approach

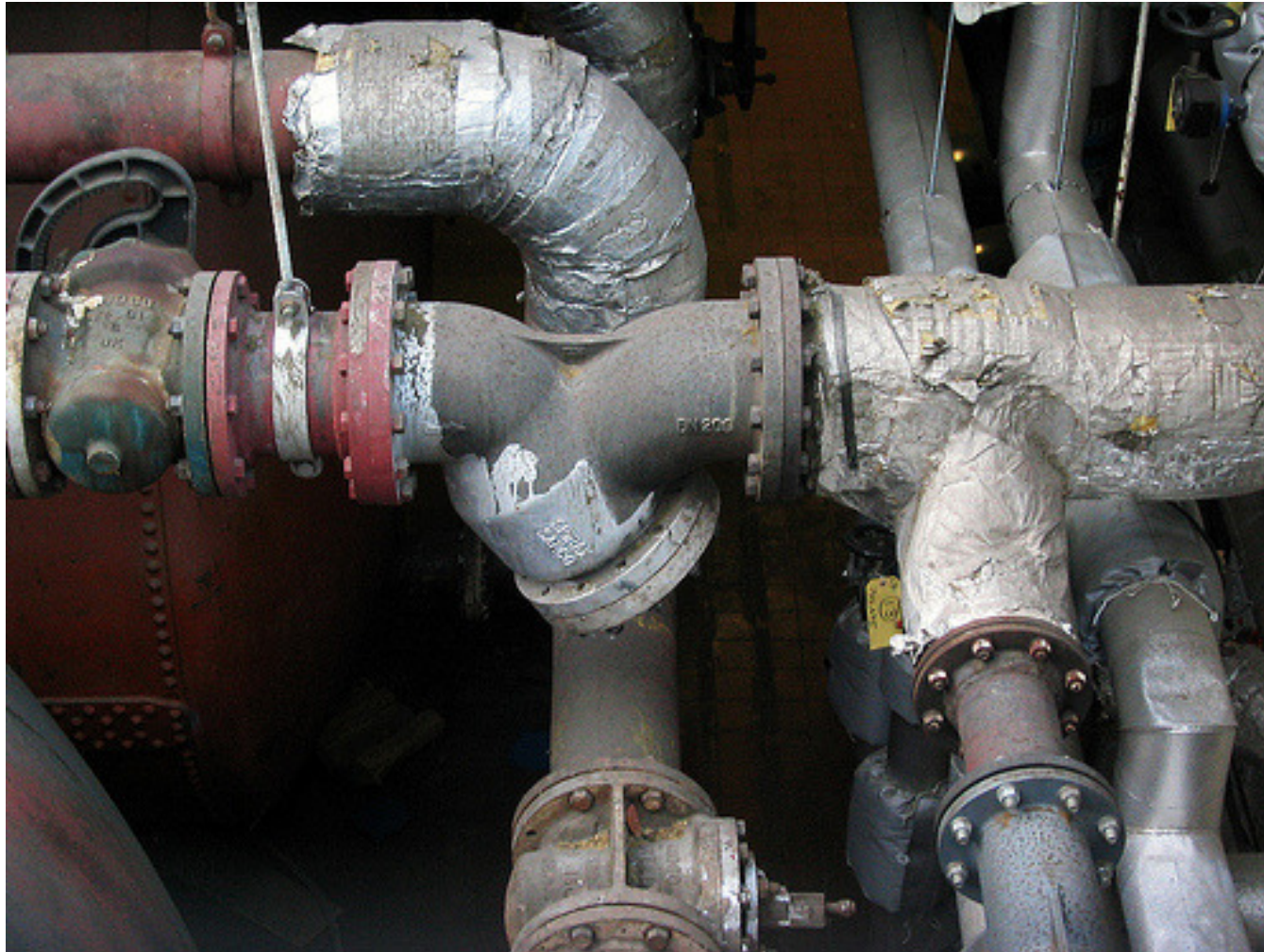
Sound engineering principles (*c.f.* Occam, Gall, and Murphy) suggest

- Favor the small and simple over the large and complex
- Favor the minimally sufficient over the feature laden
- Favor the configurable over the prescribed
- Favor the proven over the (merely) novel

Build complexity through composition, not addition

Approach sufficiency through a series of incrementally necessary steps

The pipeline metaphor



The pipeline metaphor

The pipeline concept grew out of Doug McIlroy's articulation of the "Unix philosophy"

- "Make each program do one thing well"
- "Expect the output of every program to become the input of another, as yet unknown, program"
- "Design and build software ... to be tried early"

McIlroy *et al.*, "Unix time-sharing system forward," *Bell System Technical Journal* 57:6.2 (1978): 1902

Curation micro-services

Devolve curation function into a granular set of independent, but interoperable micro-services

- Since each is small and self-contained, they are collectively easier to develop, maintain, and deploy
- Since the level of investment in any given service is small, they are easier to replace when they have outlived their usefulness
- The scope of each service is limited, but complex behavior can *emerge* from the strategic composition of individual atomistic services
- All user/service (and service/service) interaction through public interfaces

Curation micro-services

	Annotation	<i>of content by consumers</i>
<i>Value</i>	Notification	<i>of new content availability</i>

	Transformation	<i>to create derivatives</i>
	Search	<i>of content and metadata</i>
<i>Service</i>	Index	<i>to enable fast search</i>
	Ingest	<i>of content for curation</i>
<i>Curation</i>	-----	
<i>Preservation</i>	Characterization	<i>to extract content properties</i>
	Inventory	<i>of curated content</i>

	Replication	<i>for safety</i>
	Fixity	<i>to verify bit-level integrity</i>
<i>State</i>	Storage	<i>for long-term retention</i>
	Identity	<i>for long-term reference</i>

Design goals

Policy-neutral, protocol
and platform independent

Linked data

The file system is the
database

- All content and metadata fully expressed in the file system
- Some subset of metadata replicated in databases as an optimization for fast query

Multiple interface modalities

- RESTful HTTP
- Command line
- Procedural (with various language bindings)

Code to interfaces

Principle of least surprise

Appropriate benchmark for
user experience is Flickr

Implementation strategies

Merritt project



- Consolidate management of 140 TB of existing content
- Support centrally-hosted as well as locally-deployed solutions
- Partner with campus data centers for virtualized and cloud provisioning
 - UC Berkeley, SDSC, TACC
- Fine-grained application of modularity and orthogonality
- Open source work products (BSD)
- Exploit agile methods
- Stable URL references



<http://example-store/state/default/1234/3/xyz>

ANVL

ARK

Checkm

CAN

Dflat

EZID

Namaste

Pairtree

ReDD

<http://www.cdlib.org/uc3/curation>

4store

<http://4store.org/>

ERC/Dublin Kernel

<http://dublincore.org/groups/kerne>

ORE resource map

<http://www.openarchives.org/ore>

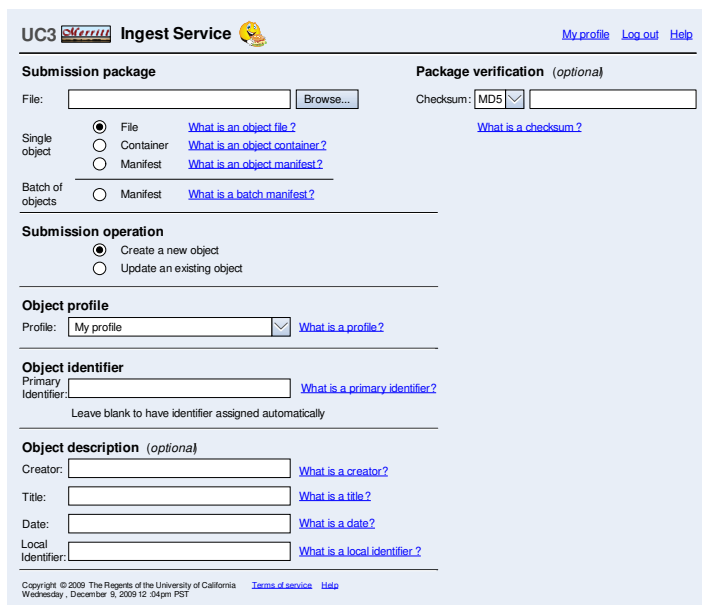
Zookeeper



<http://hadoop.apache.org/zookeepe>
r

Demonstration

A few caveats ... still a work in progress!

- The final interface style sheets are not yet applied
- Inventory service still under development
- More details available in a recent webinar
<https://confluence.ucop.edu/display/Curation/Home>



UC3  Ingest Service  [My profile](#) [Log out](#) [Help](#)

Submission package

File: [Browse...](#)

Package verification (optional)

Checksum: [What is a checksum?](#)

Single object

File [What is an object file?](#)

Container [What is an object container?](#)

Manifest [What is an object manifest?](#)

Batch of objects

Manifest [What is a batch manifest?](#)

Submission operation

Create a new object

Update an existing object

Object profile

Profile: [What is a profile?](#)

Object identifier

Primary Identifier: [What is a primary identifier?](#)

Leave blank to have identifier assigned automatically

Object description (optional)

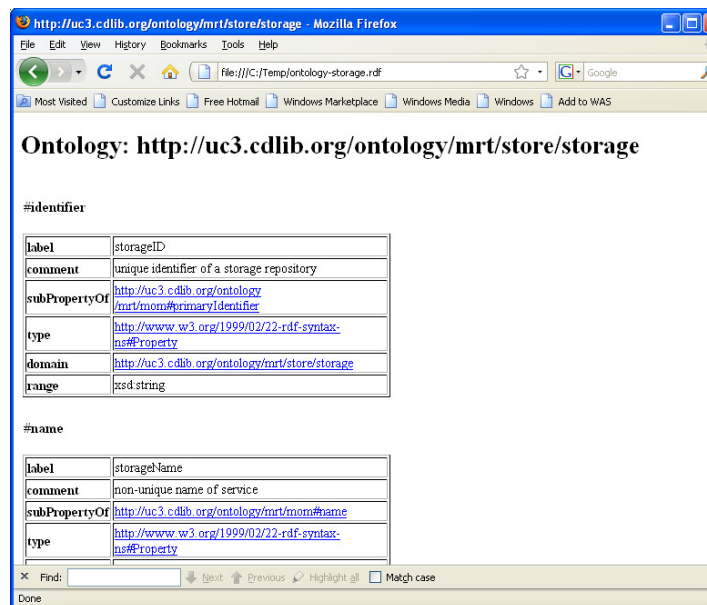
Creator: [What is a creator?](#)

Title: [What is a title?](#)

Date: [What is a date?](#)

Local Identifier: [What is a local identifier?](#)

Copyright © 2009 The Regents of the University of California
 Wednesday, December 9, 2009 12:54pm PST [Terms of service](#) [Help](#)



http://uc3.cdlib.org/ontology/mrt/store/storage - Mozilla Firefox

file:///C:/Temp/ontology-storage.rdf

Most Visited | Customize Links | Free Hotmail | Windows Marketplace | Windows Media | Windows | Add to WAS

Ontology: http://uc3.cdlib.org/ontology/mrt/store/storage

#identifier

label	storageID
comment	unique identifier of a storage repository
subPropertyOf	http://uc3.cdlib.org/ontology/mrt/mom#primaryIdentifier
type	http://www.w3.org/1999/02/22-rdf-syntax-ns#Property
domain	http://uc3.cdlib.org/ontology/mrt/store/storage
range	xsd:string

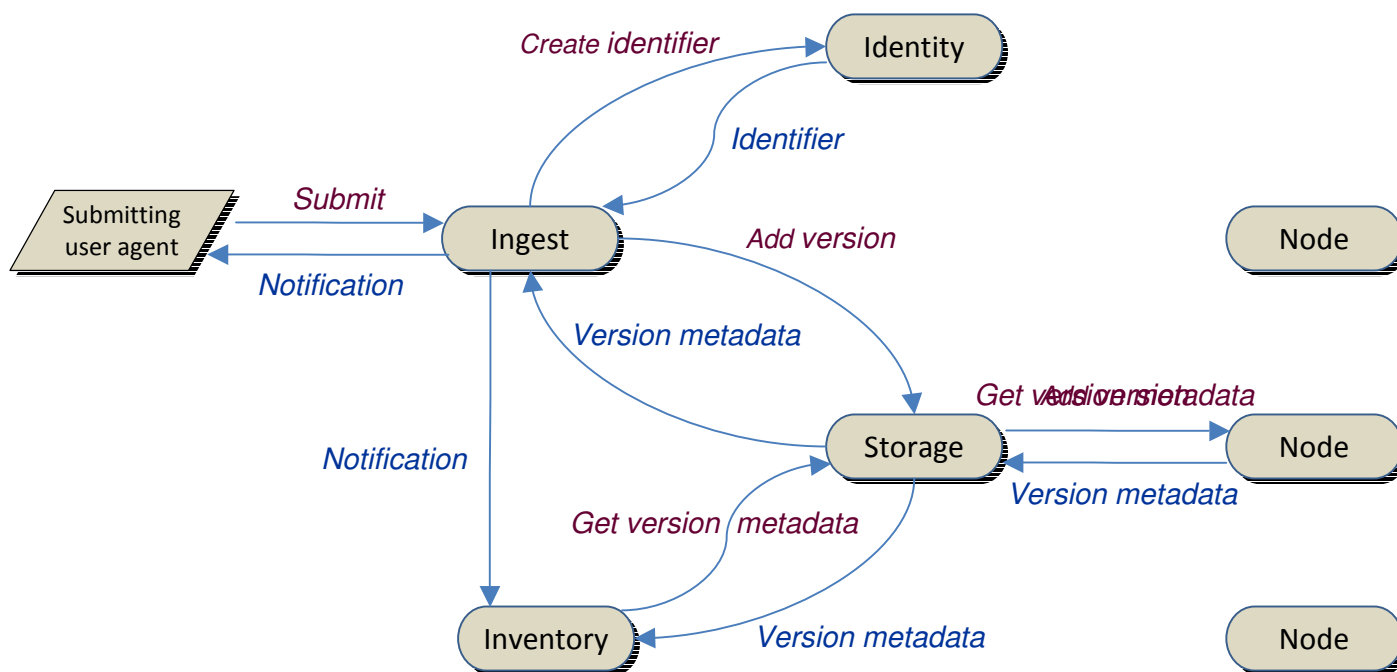
#name

label	storageName
comment	non-unique name of service
subPropertyOf	http://uc3.cdlib.org/ontology/mrt/mom#name
type	http://www.w3.org/1999/02/22-rdf-syntax-ns#Property

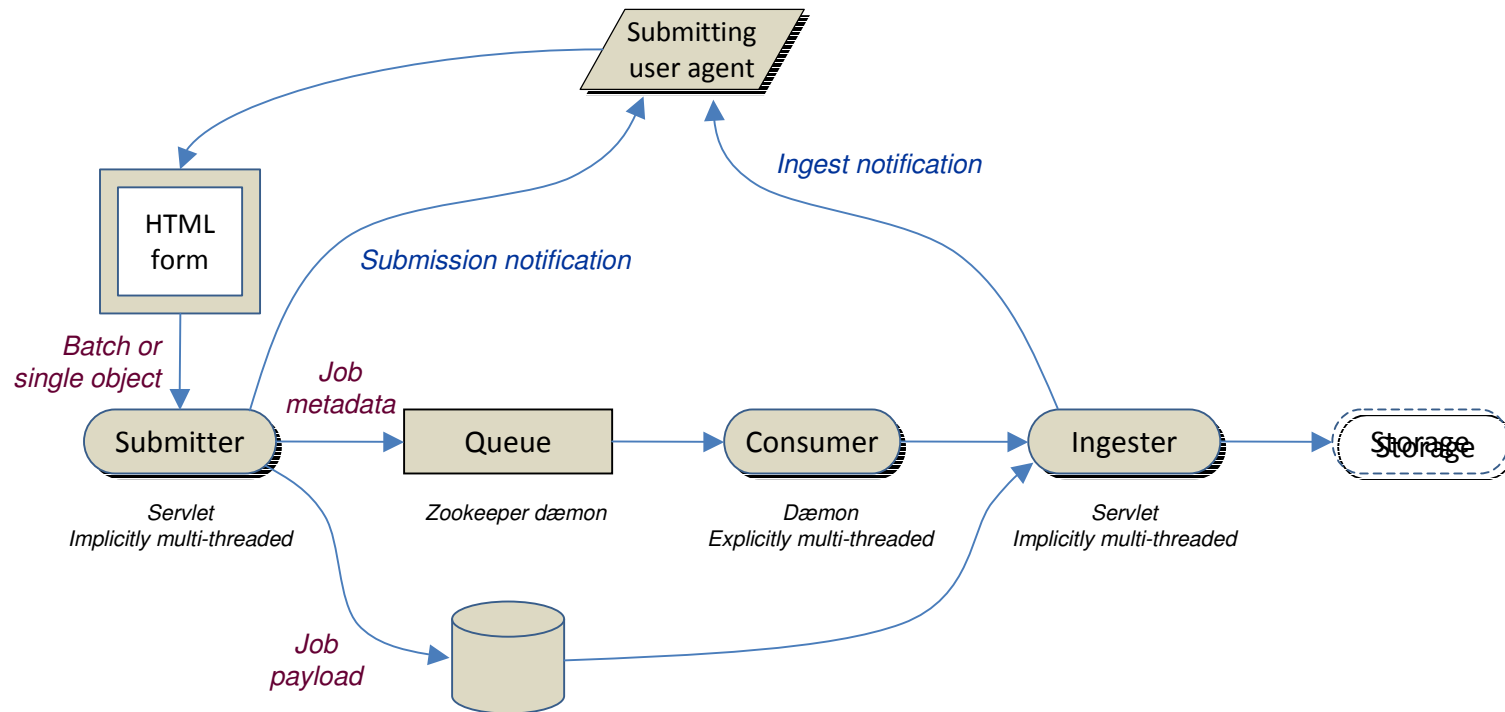
Find: [Next](#) [Previous](#) [Highlight all](#) Match case

Done

Ingest process flow



Ingest implementation



Development roadmap



<i>First wave</i>	<i>Second wave</i> ✓	<i>Third wave</i>	<i>Fourth wave</i> ✓	<i>Fifth wave</i>	<i>Sixth wave</i> ✓
Identity	Inventory	Index	Search	Notification	Annotation
Storage	Ingest	Fixity	Replication	Characterization	Transformation
Object / collection modeling			Metadata standards		
Authentication / authorization			Semantic interoperability		
Policy / business model development					



Partnerships and collections

UC3 Digital Preservation Repository (DPR)

<http://www.cdlib.org/uc3/dpr.html>



– California Digital Newspaper Collection



– CDL eScholarship publishing



– Media Vault Program



– Minnesota Historical Society



– Open Context



– UCTV



– Water Resource Center Archive



and many others

UC3 Web Archiving Service (WAS)

<http://was.cdlib.org/>



DataCite



<http://datacite.org/>

DataONE



<https://dataone.org/>

Early community reaction

Collaborative development and integration projects with UC3 partners

Independent implementation of key Merritt specifications

- Georg-August-Universität Göttingen – New York University
- HathiTrust / University of Michigan – University of North Texas
- Oxford University

Digital curation group

<http://groups.google.com/group/digital-curation>

Curation micro-services Barcamp, Berkeley, August 16-17

<http://groups.google.com/group/digital-curation/web/curation-technology-sig>

–Possible follow-up events at iPRES (Vienna, September 19-24) and DLF (Palo Alto, November 1-3)

Curation BOF session, Wednesday, 15:45-16:45, Reino Unido A

<http://or10.crowdvine.com/pages/bof>

Summary

The pipeline concept provides a useful metaphor for curation micro-services embodying the following principles

- Modularity
- Orthogonality
- Parsimony
- Evolution
- Emergence



© <http://www.flickr.com/photos/pbogs/2351286100/>

Key activity pattern: *define, decompose, recurse*

Build complexity by composition, not addition

Significantly lower levels of development effort, and higher quality solutions

For more information

Curation micro-services BOF, Wednesday, 15:45-16:45, Reino Unido A
<http://or10.crowdvine.com/pages/bof>

UC Curation Center (UC3)
<http://www.cdlib.org/uc3>

Micro-service specifications
<https://confluence.ucop.edu/display/Curation>

Digital curation group and Barcamp
<http://groups.google.com/group/digital-curation>
<http://groups.google.com/group/digital-curation/web/curation-technology-sig>

UC3 / CDL

Stephen Abrams
Willet

Erik Hetzner

Margaret Low

Mark Reyes

Perry

Patricia Cruse

Greg Janée

David Loy Tracy Seneca

Scott Fisher

John Kunze

Isaac Rabinovitch

Marisa Strong

UC3 / SDSC