

Author(s): Jeremy York, 2010

License: Unless otherwise noted, this material is made available under the terms of the **Creative Commons Attribution–Noncommercial–Share Alike 3.0 License:**
<http://creativecommons.org/licenses/by-nc-sa/3.0/>

We have reviewed this material in accordance with U.S. Copyright Law **and have tried to maximize your ability to use, share, and adapt it.** The citation key on the following slide provides information about how you may share and adapt this material.

Copyright holders of content included in this material should contact open.michigan@umich.edu with any questions, corrections, or clarification regarding the use of content.

For more information about **how to cite** these materials visit <http://open.umich.edu/education/about/terms-of-use>.

Any **medical information** in this material is intended to inform and educate and is **not a tool for self-diagnosis** or a replacement for medical evaluation, advice, diagnosis or treatment by a healthcare professional. Please speak to your physician if you have questions about your medical condition.

Viewer discretion is advised: Some medical content is graphic and may not be suitable for all viewers.

Citation Key

for more information see: <http://open.umich.edu/wiki/CitationPolicy>

Use + Share + Adapt

{ Content the copyright holder, author, or law permits you to use, share and adapt. }



Public Domain – Government: Works that are produced by the U.S. Government. (17 USC § 105)



Public Domain – Expired: Works that are no longer protected due to an expired copyright term.



Public Domain – Self Dedicated: Works that a copyright holder has dedicated to the public domain.



Creative Commons – Zero Waiver



Creative Commons – Attribution License



Creative Commons – Attribution Share Alike License



Creative Commons – Attribution Noncommercial License



Creative Commons – Attribution Noncommercial Share Alike License



GNU – Free Documentation License

Make Your Own Assessment

{ Content Open.Michigan believes can be used, shared, and adapted because it is ineligible for copyright. }



Public Domain – Ineligible: Works that are ineligible for copyright protection in the U.S. (17 USC § 102(b)) *laws in your jurisdiction may differ

{ Content Open.Michigan has used under a Fair Use determination. }



Fair Use: Use of works that is determined to be Fair consistent with the U.S. Copyright Act. (17 USC § 107) *laws in your jurisdiction may differ

Our determination **DOES NOT** mean that all uses of this 3rd-party content are Fair Uses and we **DO NOT** guarantee that your use of the content is Fair.

To use this content you should **do your own independent analysis** to determine whether or not your use will be Fair.



HATHI TRUST

A Shared Digital Repository

Building A Future By Preserving Our Past

The Preservation Infrastructure of
HathiTrust Digital Library

Jeremy York

IFLA 2010

August 15, 2010

Current Partners

- Columbia University
- New York Public Library
- University of California system
- CIC (Committee on Institutional Cooperation)
 - University of Chicago
 - University of Illinois
 - Indiana University
 - University of Iowa
 - University of Michigan
 - Michigan State University
 - University of Minnesota
 - Northwestern University
 - Ohio State University
 - Pennsylvania State University
 - Purdue University
 - University of Wisconsin-Madison
- University of Virginia
- Yale University

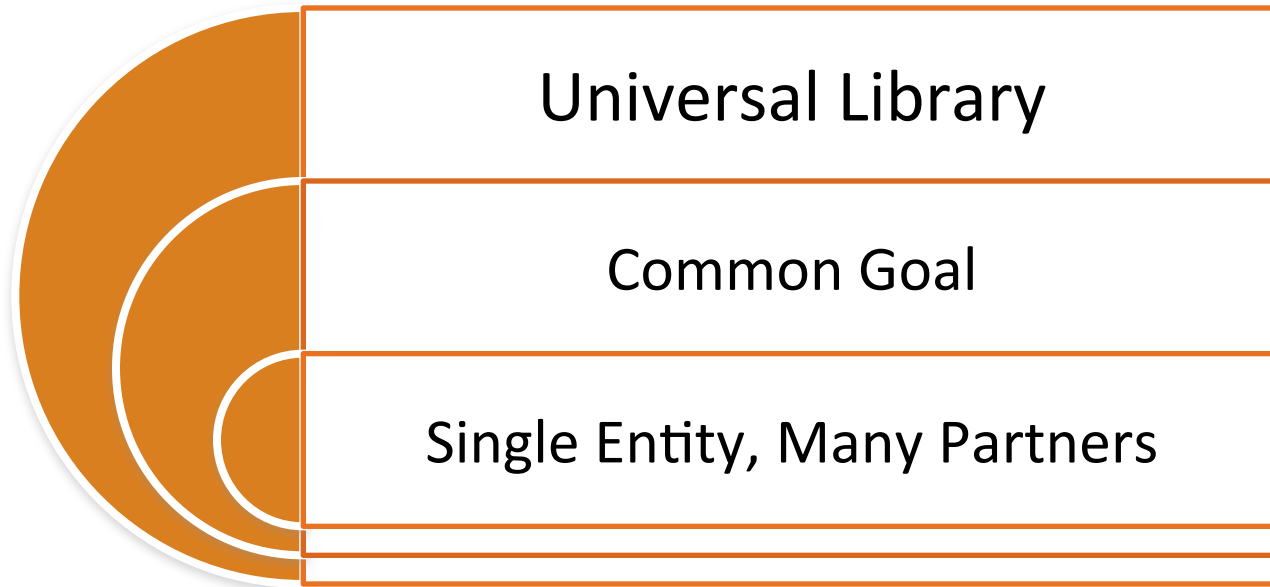


Mission

- To contribute to the common good by collecting, organizing, preserving, communicating, and sharing the record of human knowledge



HathiTrust

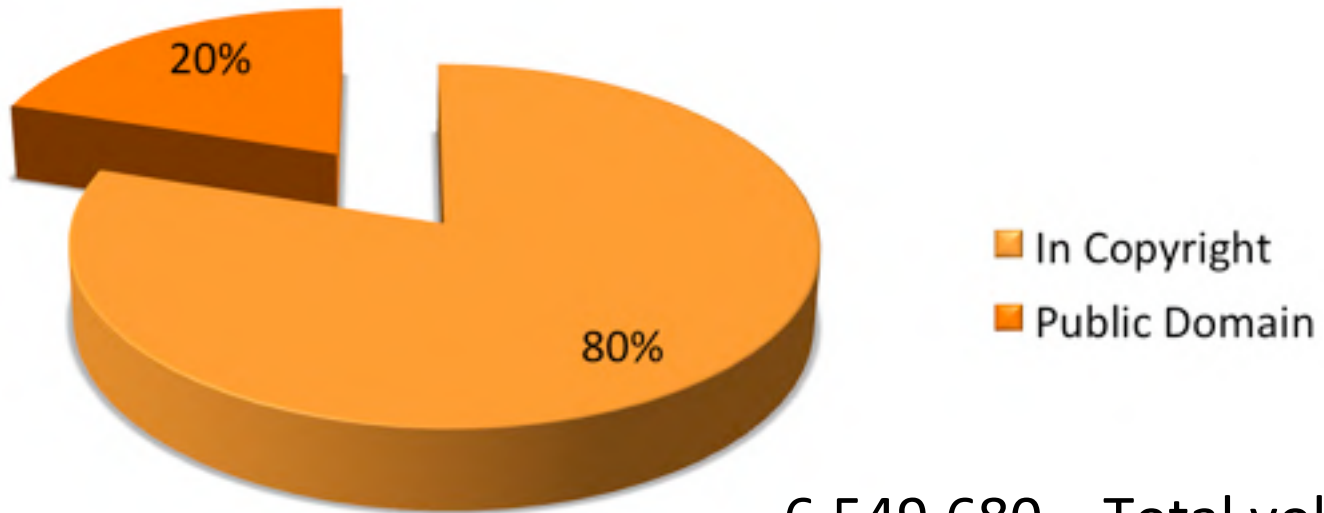


Goals

- Comprehensive collection
- Preservation...with Access
- Shared strategies
 - Collection management, development
 - Preservation
 - Copyright
 - Efficient user services
- Openness



Content Distribution

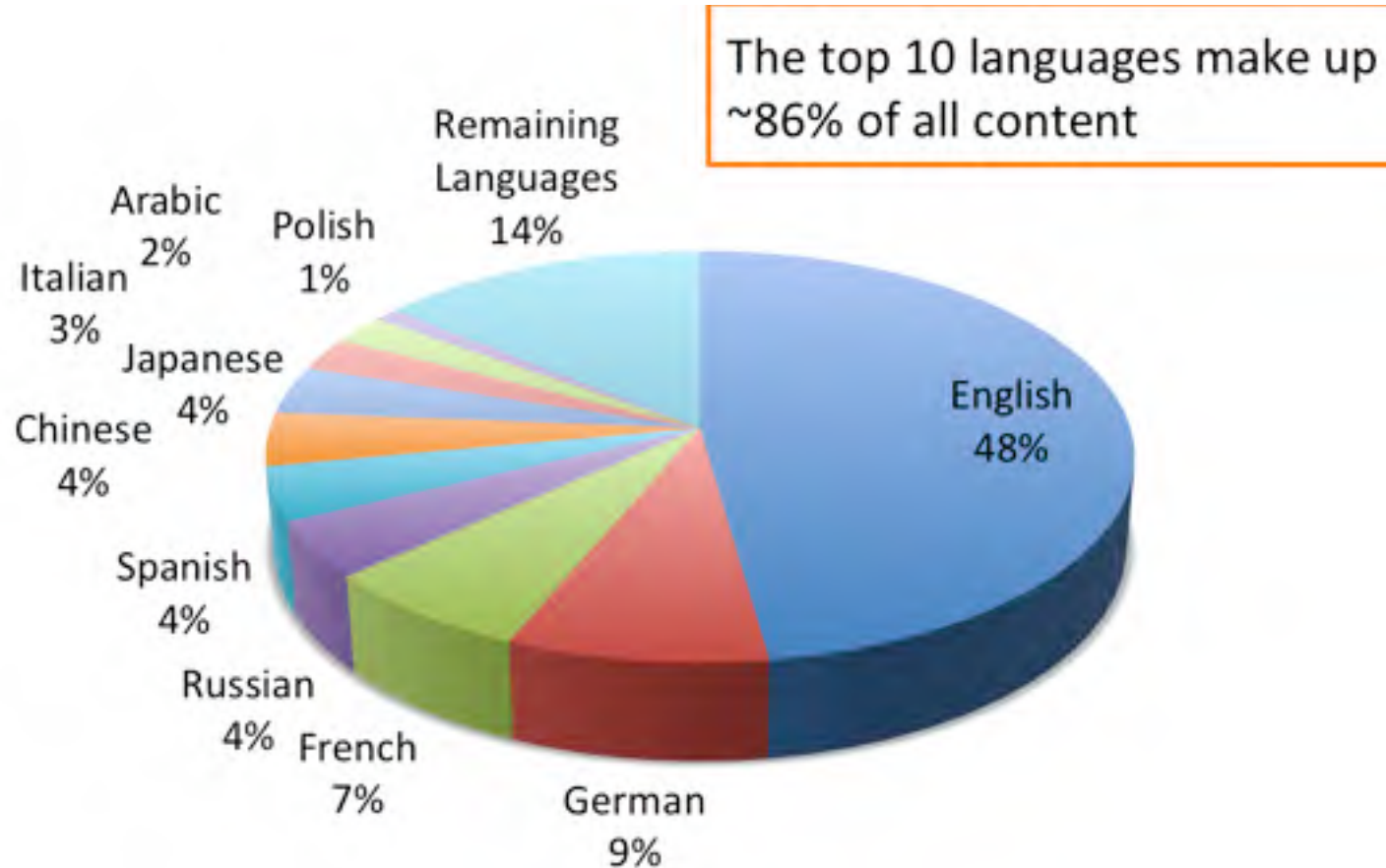


6,549,680 – Total volumes
1,300,896 – Public Domain
3,798,116 Book titles
153,311 Serial titles

* As of August 13, 2010



Language Distribution (1)

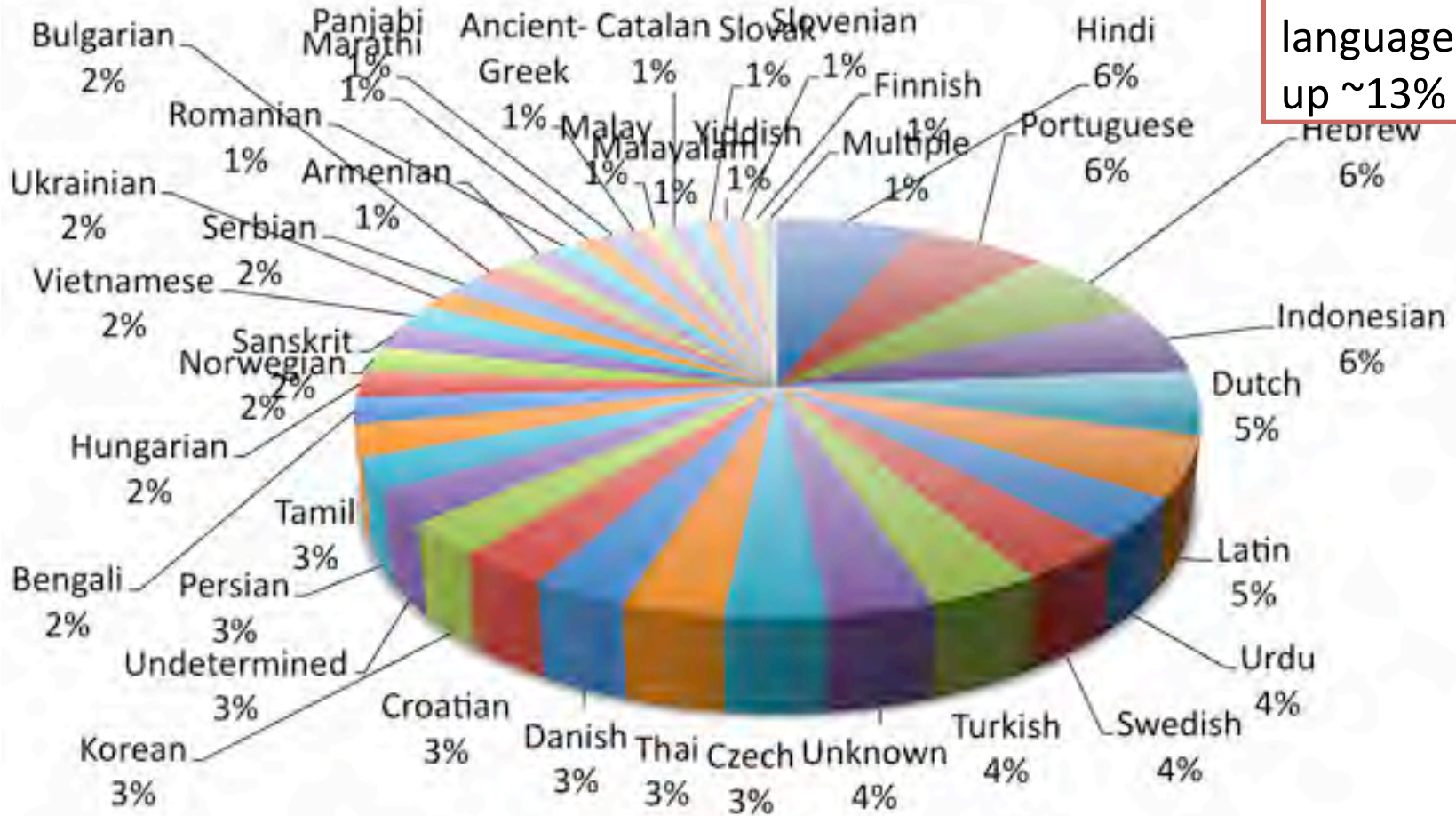


* As of August 13, 2010



Language Distribution (2)

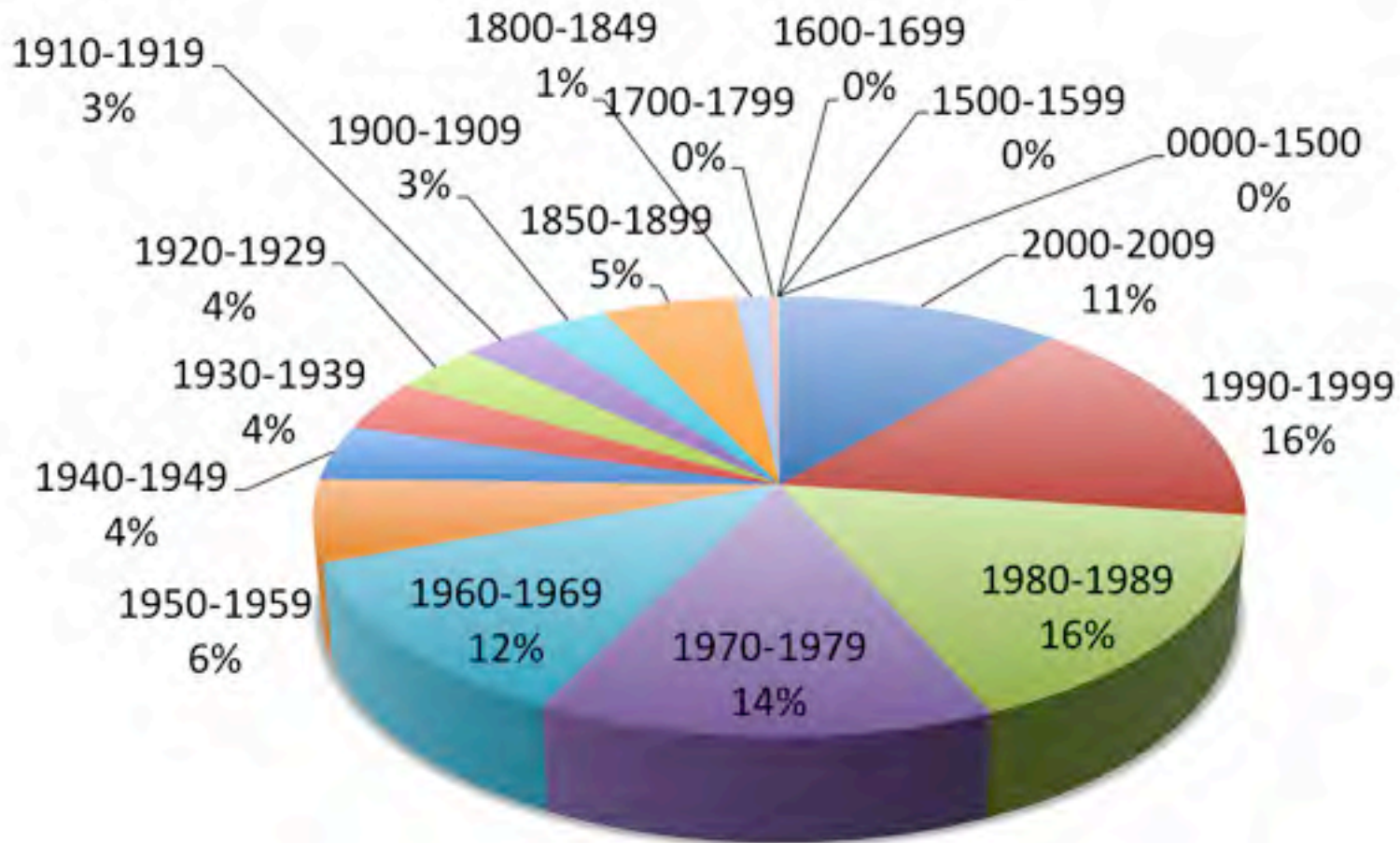
The next 40 languages make up ~13% of total



* As of August 13, 2010



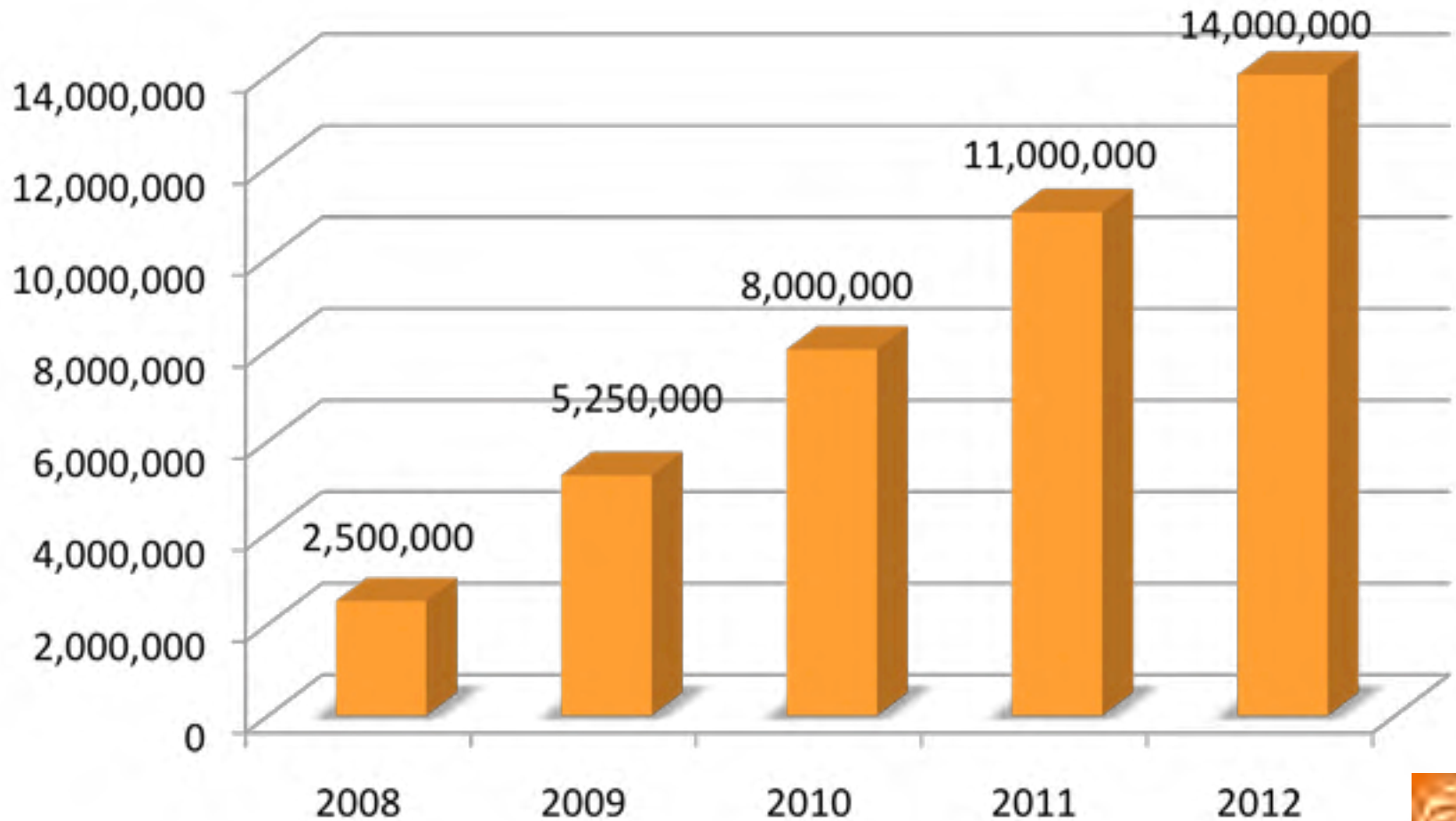
Dates



* As of August 13, 2010



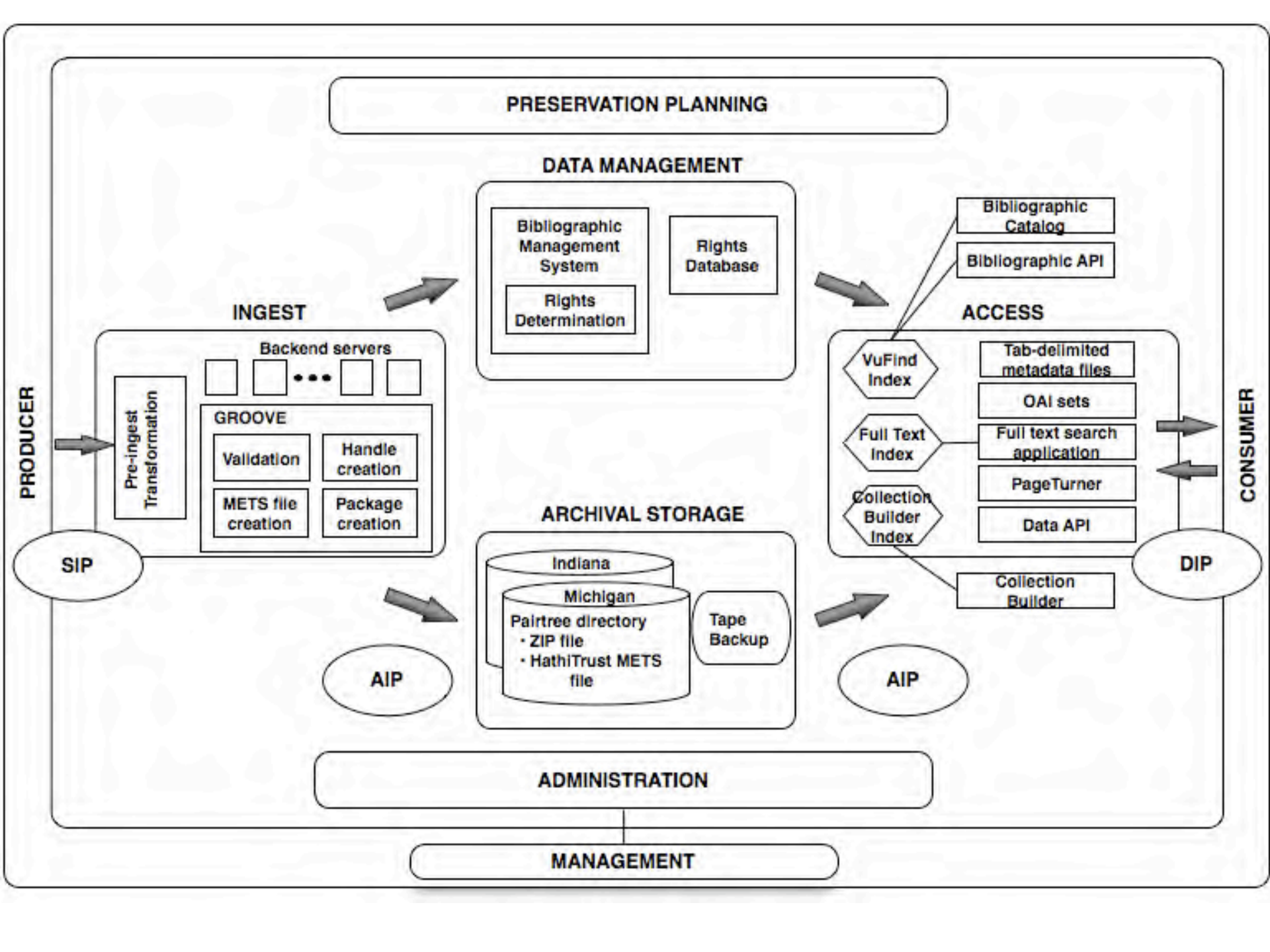
Content Growth



Repository Philosophy/Design

- OAIS/TRAC
- Consistency
- Standardization
- Simplicity (in design, not function)
- Practicality
- Sustainability



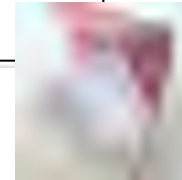
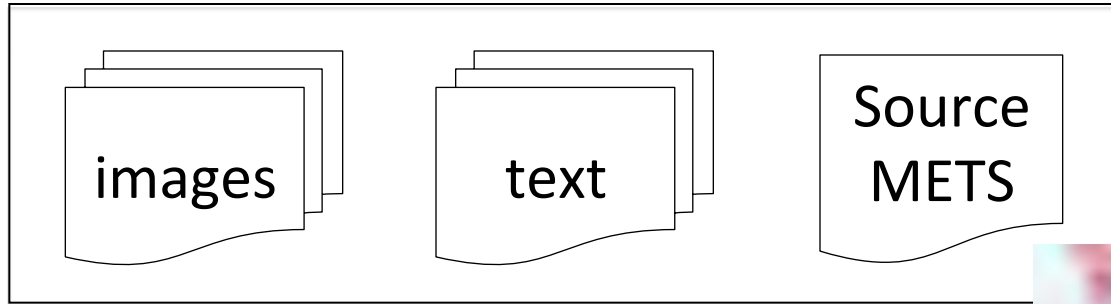
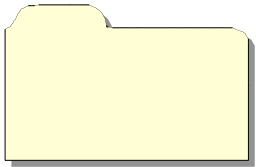


Content

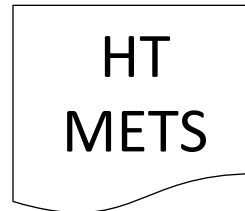
- Largely uniform in technical characteristics
- 4 formats
 - ITU G4 TIFF
 - JPEG2000
 - JPEG
 - Unicode (with and without coordinates)



Object Package



malachus,
Flickr.com



Metadata

- Details and specifications at repository level
 - Object specifications / Validation criteria
 - Page-tagging
- Variations at object level
 - Files missing
 - Non-valid files
 - Incorrect file checksums

http://www.hathitrust.org/digital_object_specifications

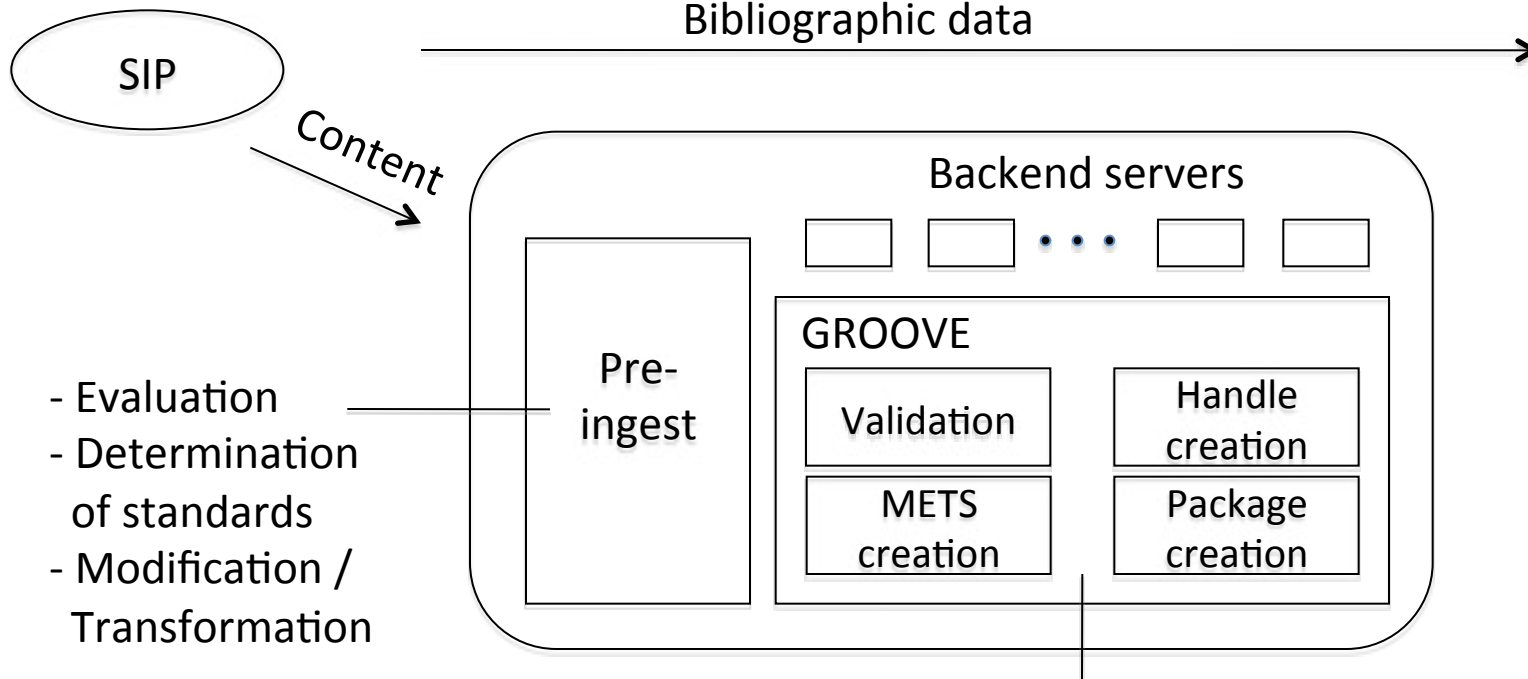


Ingest

- Bibliographic Data
 - Must be present prior to content ingest
 - MARCXML, as complete as possible
- Content
 - Pre-ingest
 - Ingest



Ingest (2)



- Evaluation
- Determination of standards
- Modification / Transformation

- Ensure conformance
 - Barcode
 - Fixity
 - Consistency
 - Well-formedness
 - Prepare archival package



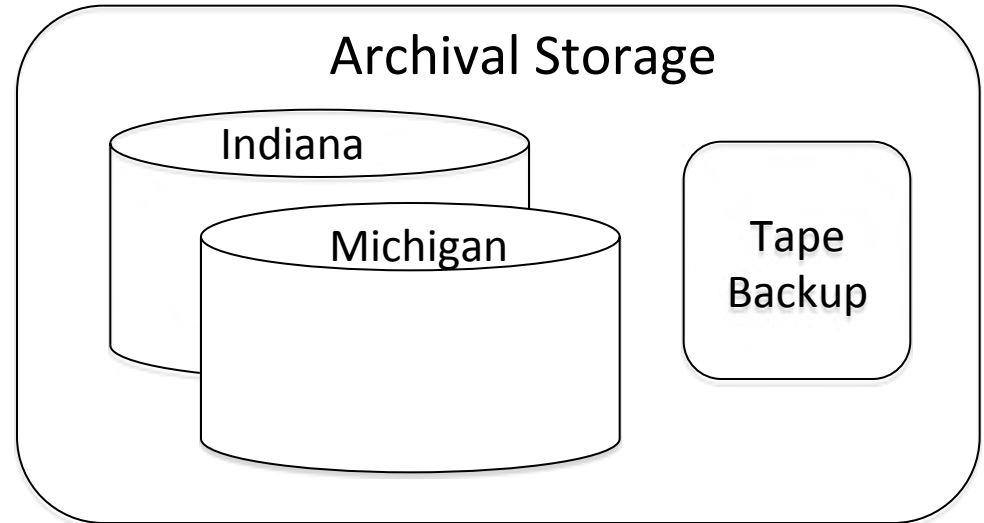
Archival Storage

- Reliability – ensure integrity
- Redundancy – in single and multiple sites
- Scalability – including ease of management
- Accessibility – for repository processes and services
- Platform-independence – for data/object management

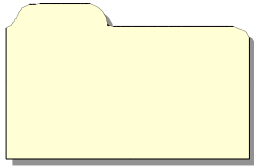


Media & Architecture

- Isilon Systems
- Load balancing and failover
- Ingest at Michigan, replicated to Indiana
- Replacement on 3-4 year cycle

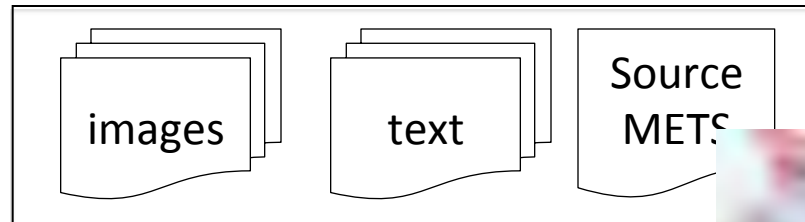


Architecture & Management



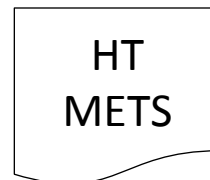
../**uc1**/pairtree_root/b3/54/34/86/b34543486

b34543486.zip



malachus,
Flickr.com

b34543486.mets.xml



Example ids:

wu.89094366434

mdp.39015037375253

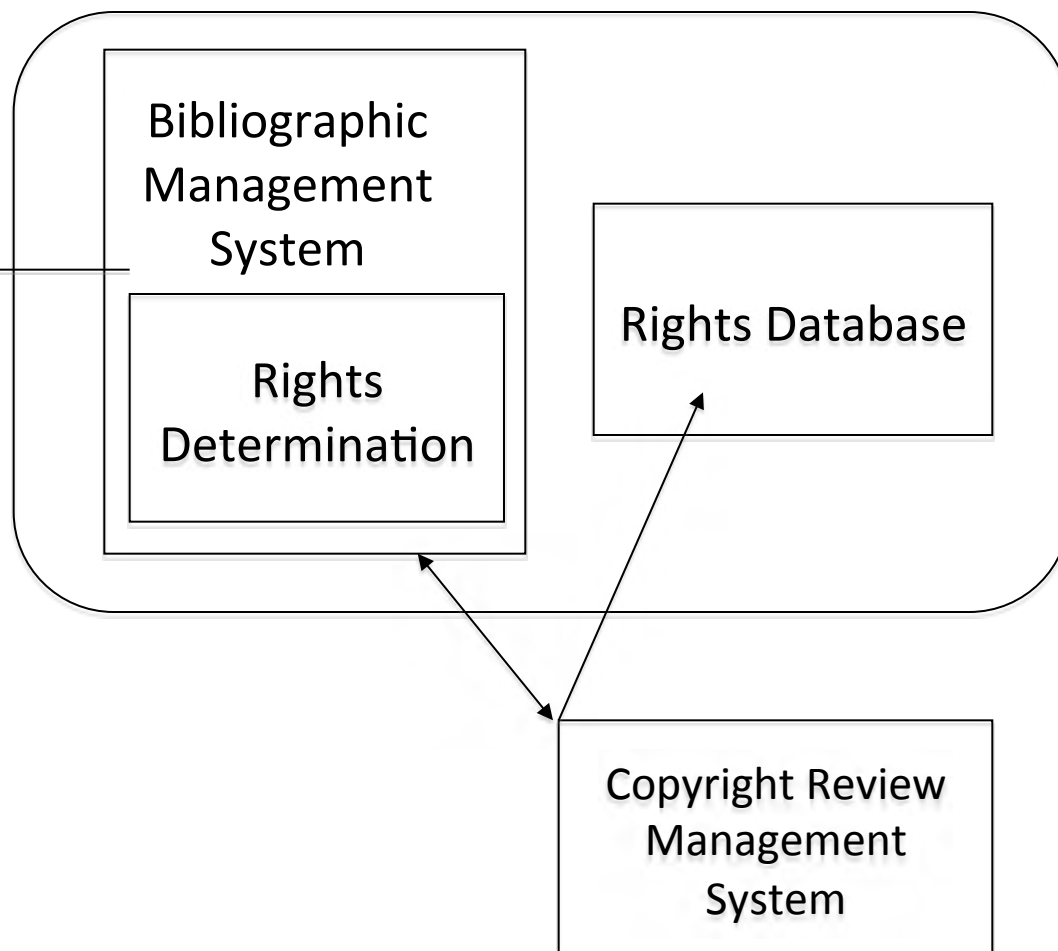
uc2.ark:/1390/t26973133

miua.aaj0523.1950.001



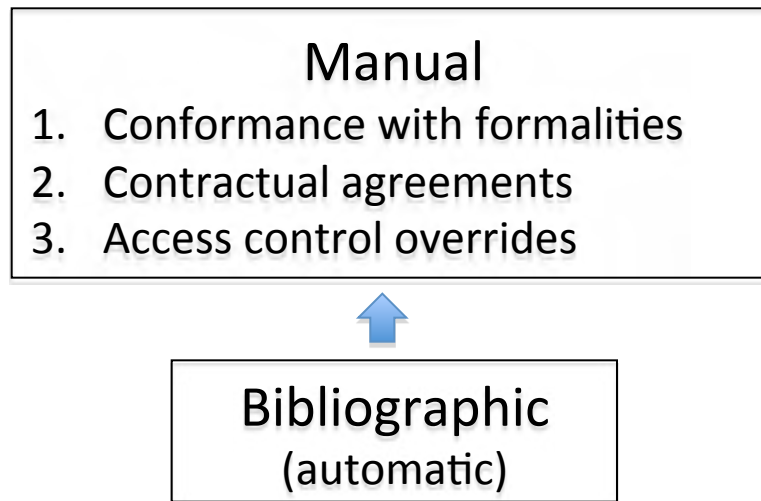
Data Management

- Inventory
- Loading and updating records
- Duplicate detection and collation
- Solr indexes behind VuFind catalog
- Source of information for Access services
- Rights determination (automated and support for manual review)



Rights Database

- System of precedence

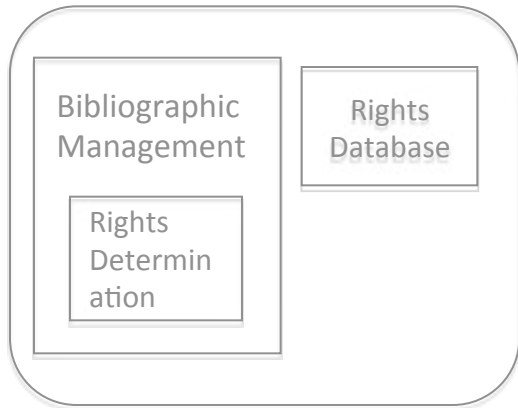


- 9 attributes
- 11 reason codes

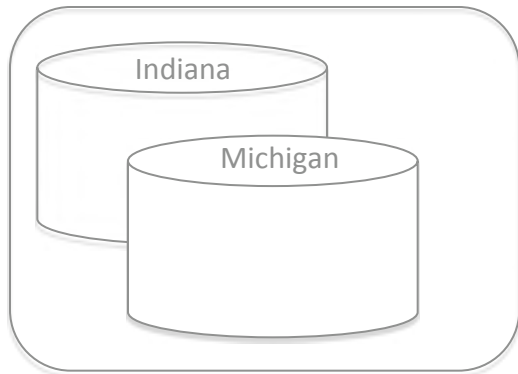


Access

Data Management



Archival Storage



Bibliographic Catalog

Bibliographic API

VuFind Index

Full text Index

Collection Builder Index

Tab-delimited Metadata files

OAI sets

Full text Search application

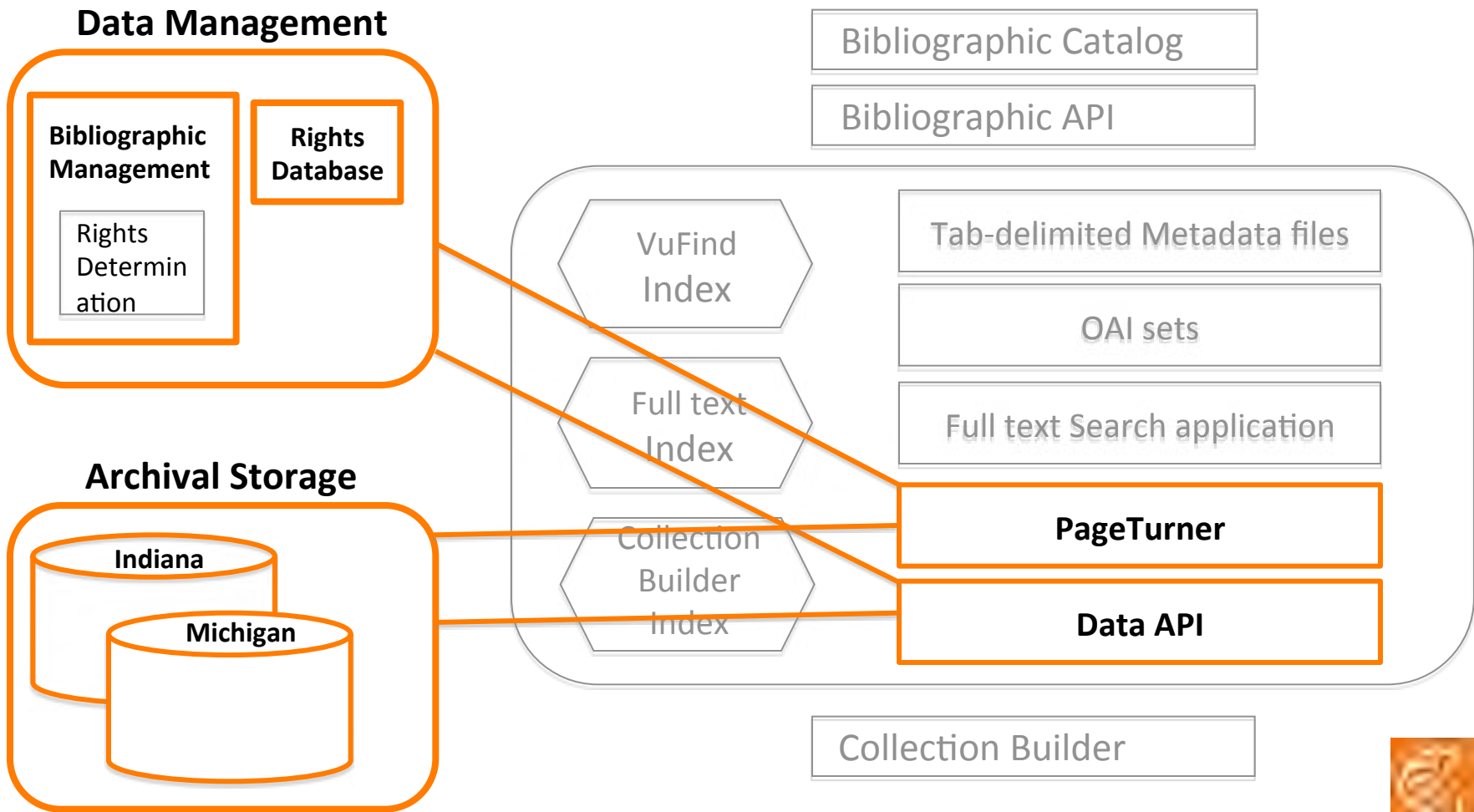
PageTurner

Data API

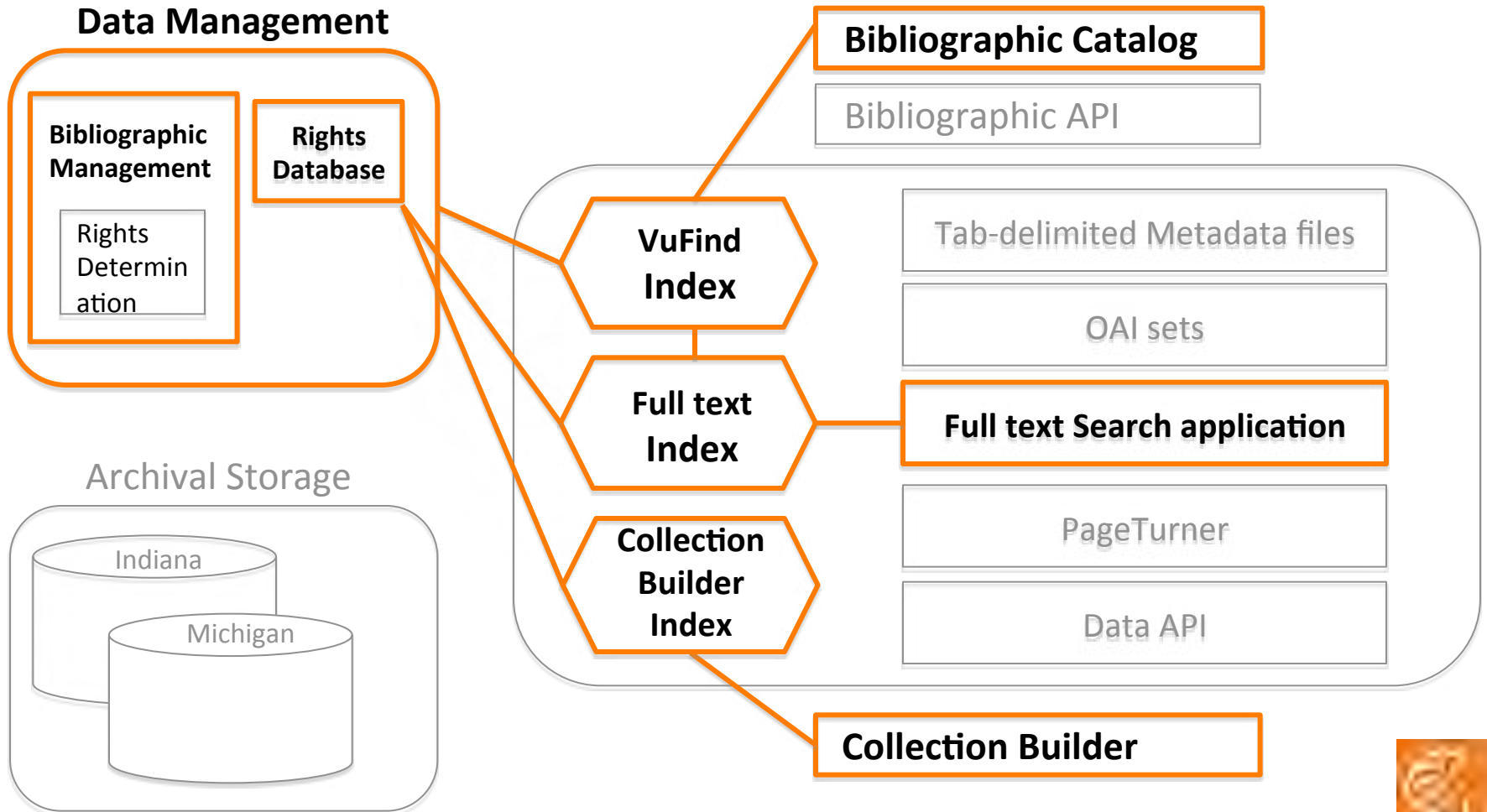
Collection Builder



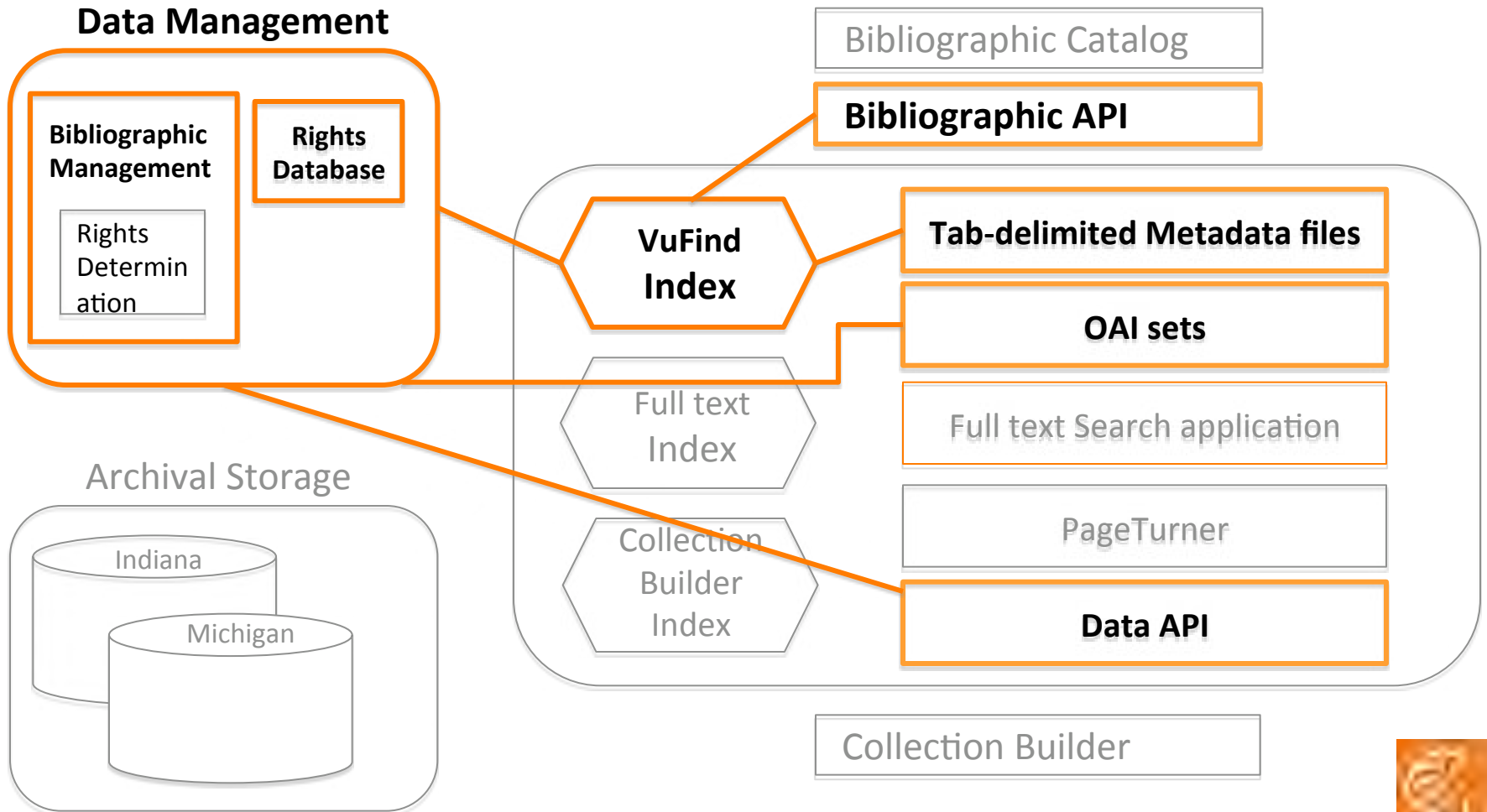
Content Access

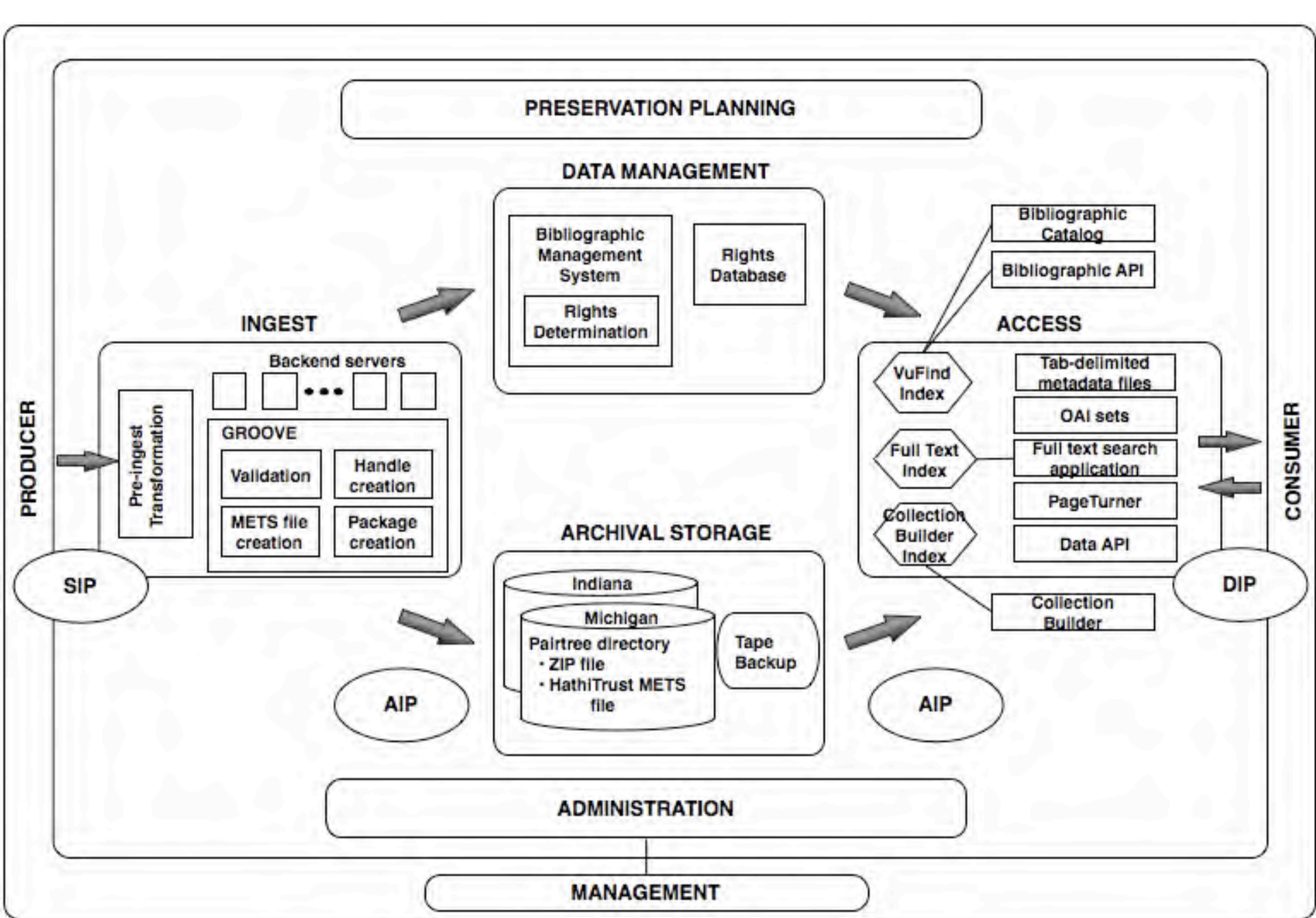


Search and Aggregation Access



Metadata Access





Thank you!

hathitrust-info@umich.edu



Additional Source Information

for more information see: <http://open.umich.edu/wiki/CitationPolicy>

Slide 16, Image 11: malachus, Flickr.com

Slide 22, Image 11: malachus, Flickr.com, <http://www.flickr.com/photos/malachus/5152200478/>

Slide 29, Image 0: Source Undetermined