

I T H A K A

JSTOR | PORTICO | ITHAKA S+R

Digital Preservation Case Studies: Preservation Activities at Portico

Sheila Morrissey

Senior Research Developer, Portico, ITHAKA

UN FAO

Digital Preservation and JHOVE2

Rome

May 24, 2011

“Digital Preservation is Everyone’s Problem ...

*BUT IT ISN'T THE SAME PROBLEM
FOR EVERYONE!!”*



ITHAKA is a not-for-profit organization that helps the academic community use digital technologies to preserve the scholarly record and to advance research and teaching in sustainable ways.

We pursue this mission by providing innovative services that aid in the adoption of these technologies and that create lasting impact.



Ithaka S+R is a research and consulting service that focuses on the transformation of scholarship and teaching in an online environment, with the goal of identifying the critical issues facing our community and acting as a catalyst for change.



JSTOR is a research platform that enables discovery, access, and preservation of scholarly content.



PORTICO

Portico is a digital preservation service for e-journals, e-books, and other scholarly e-content.



PORTICO

Portico is among the largest community-supported digital archives in the world.

Working with libraries, publishers, and funders, we preserve e-journals, e-books, and other electronic scholarly content to ensure researchers and students will have access to it in the future.



PORTICO

An “Insurance Policy” for e-Content

Provide libraries with access to archived content when it becomes lost, orphaned or abandoned (regardless of libraries’ past or current subscription):

Publisher ceases operation

Publisher discontinues title

Publisher drops back file

- Provide libraries with post-cancellation access – if publisher specifically names Portico
- About 90% of titles in Archive are covered by Portico post-cancellation access rights.



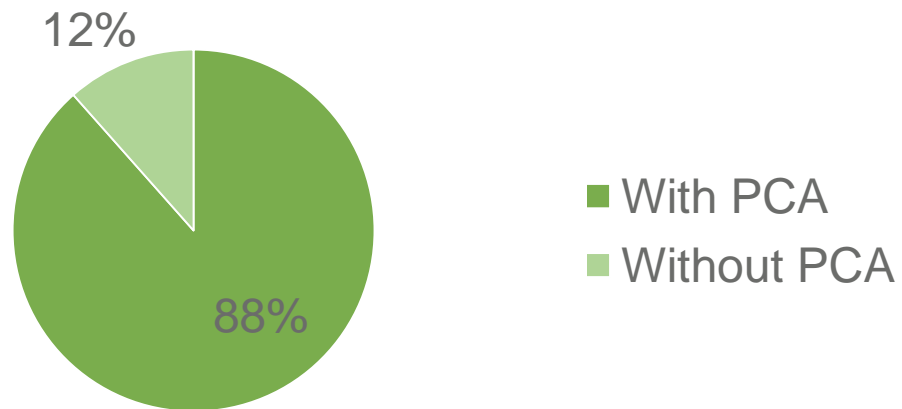
Triggered Content

Title	Trigger Date	Publisher	Holdings Available	Years
<i>Auto/Biography</i>	2008/07	SAGE Publications	v. 12-14	2004-2006
<i>Brief Treatment and Crisis Intervention</i>	2009/04	Press	v. 1-8	2001-2008
<i>Graft</i>	2007/12	SAGE Publications	v. 4-6	2001-2003
<i>Pain Reviews</i>	2009/07	Hodder	v. 5-9	1998-2002

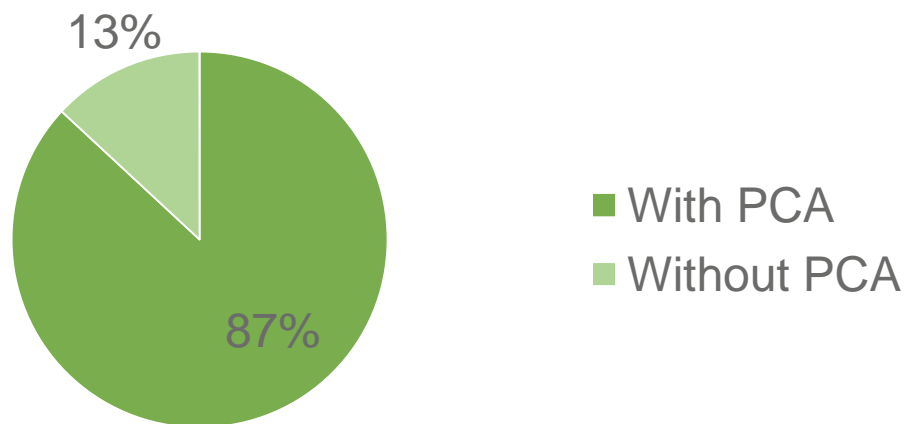
Post-Cancellation Access Requests



E-Journals



E-Books





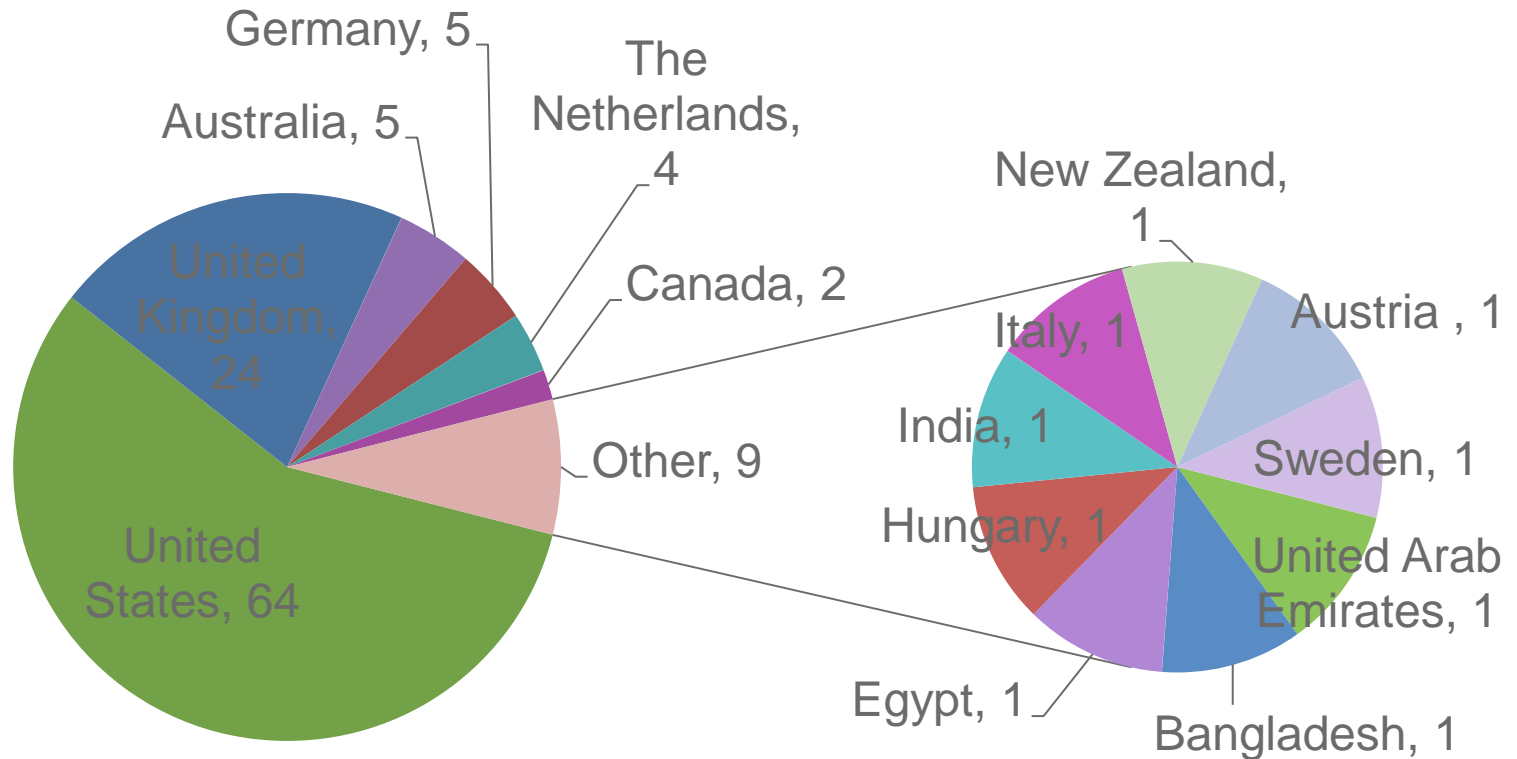
PORTICO

Over 2,000 societies, and associations have committed content to Portico through 122 publishers agreements.

» E-journal titles	12,142
» E-book titles	73,298
» D-collections	39



Portico Participating Publishers



Numbers as of 8/31/2010



Participating Libraries

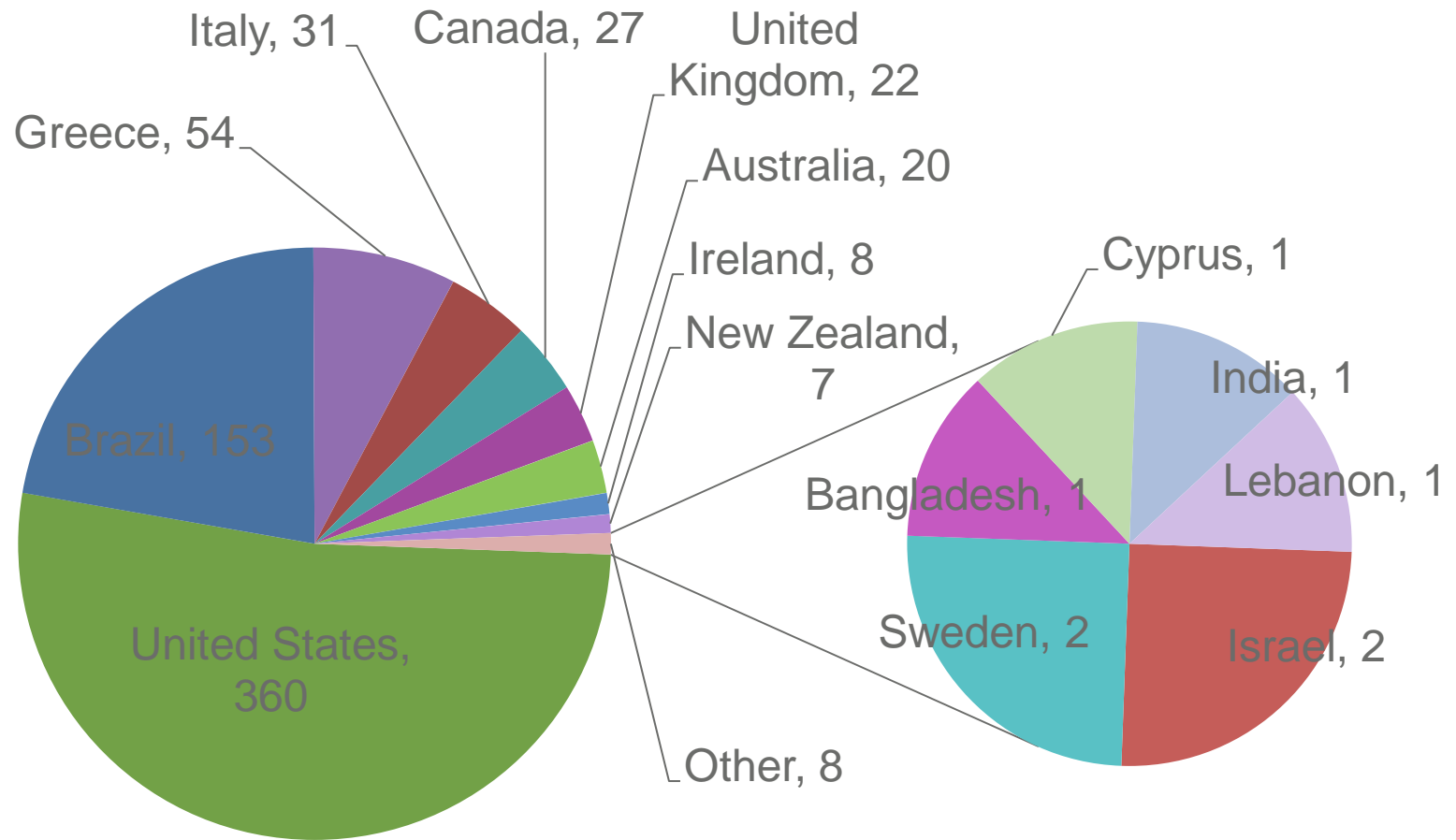
Participating Libraries	
Participating Libraries	690
US Libraries	360
Non-US Libraries	330

Numbers as of 8/31/2010



PORTICO

Portico Participating Libraries



Numbers as of 8/31/2010



TAKE THE LONG VIEW...



Portico Timeline

2002
Launch of
Electronic
Archiving
Initiative
by JSTOR

2005
Portico
signs
initial e-
journal
publishers

2008
Portico
signs
initial e-
book
publishers

2009
Portico
signs
initial d-
collections

2005
Portico
Launched

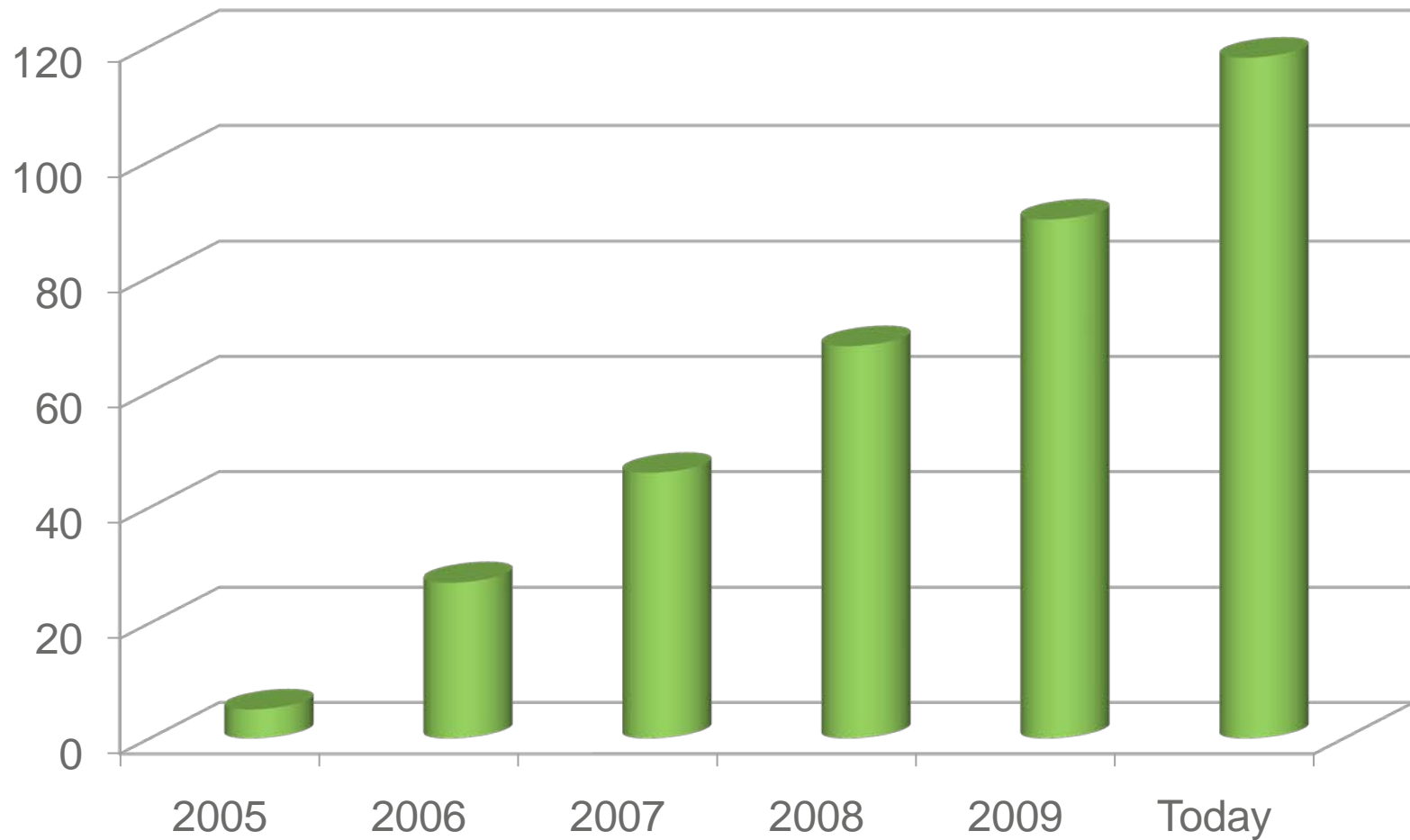
2006
Portico
ingest
initial e-
journal
content
into the
archive

2009
Portico
ingests
initial e-
book
content
into the
archive

2010
Portico
ingests
initial d-
collection
content



Portico Participating Publishers

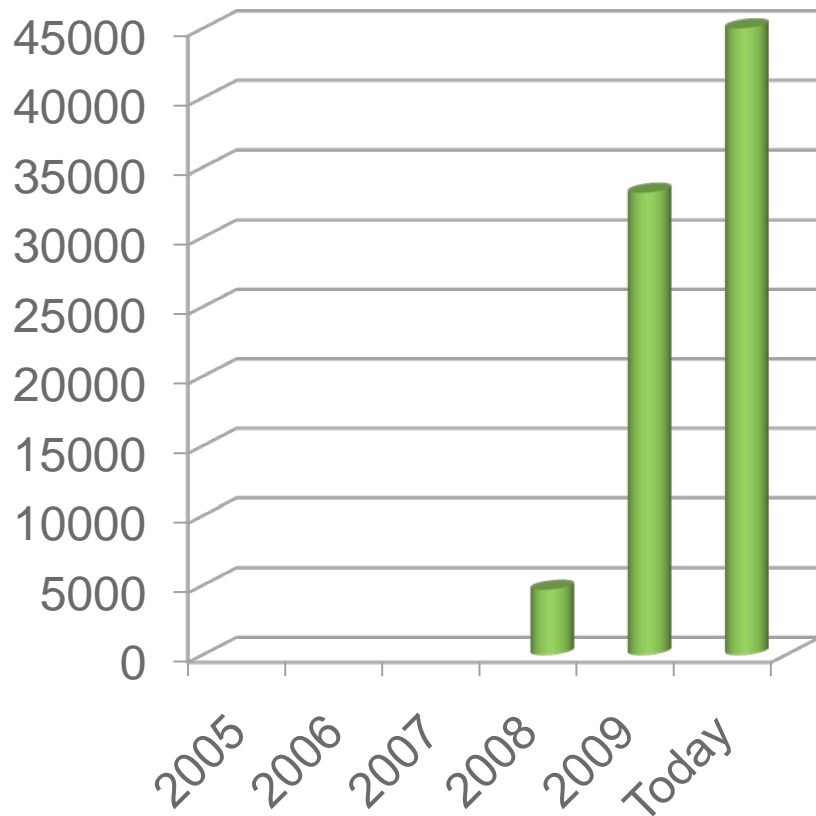


Numbers as of 8/31/2010

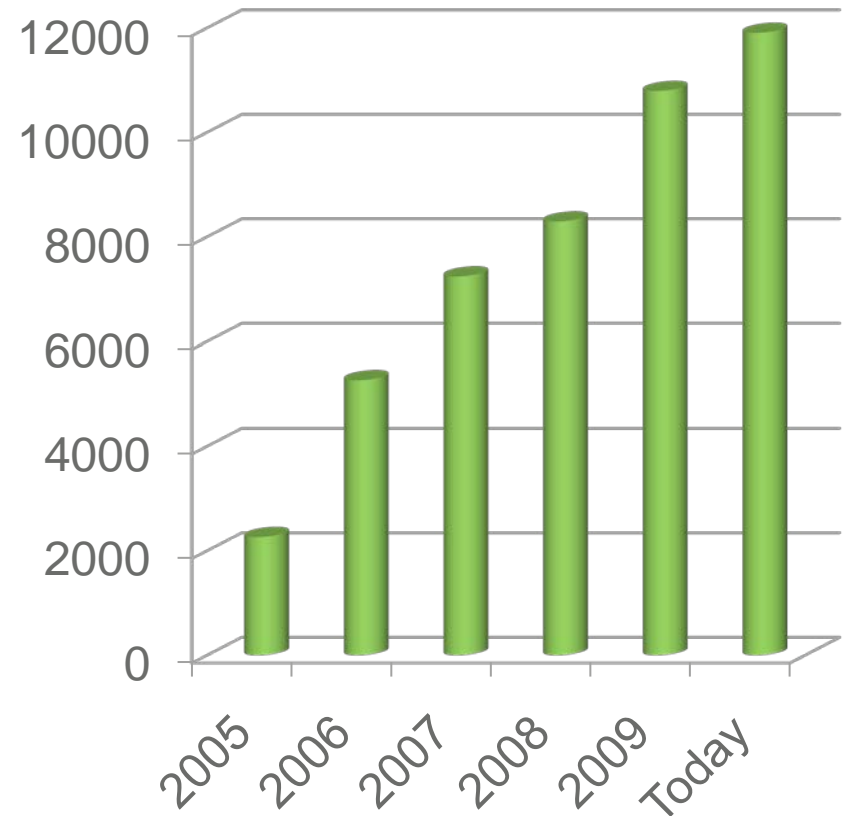


Portico Growth in Participating Titles

Participating E-Books



Participating E-Journals

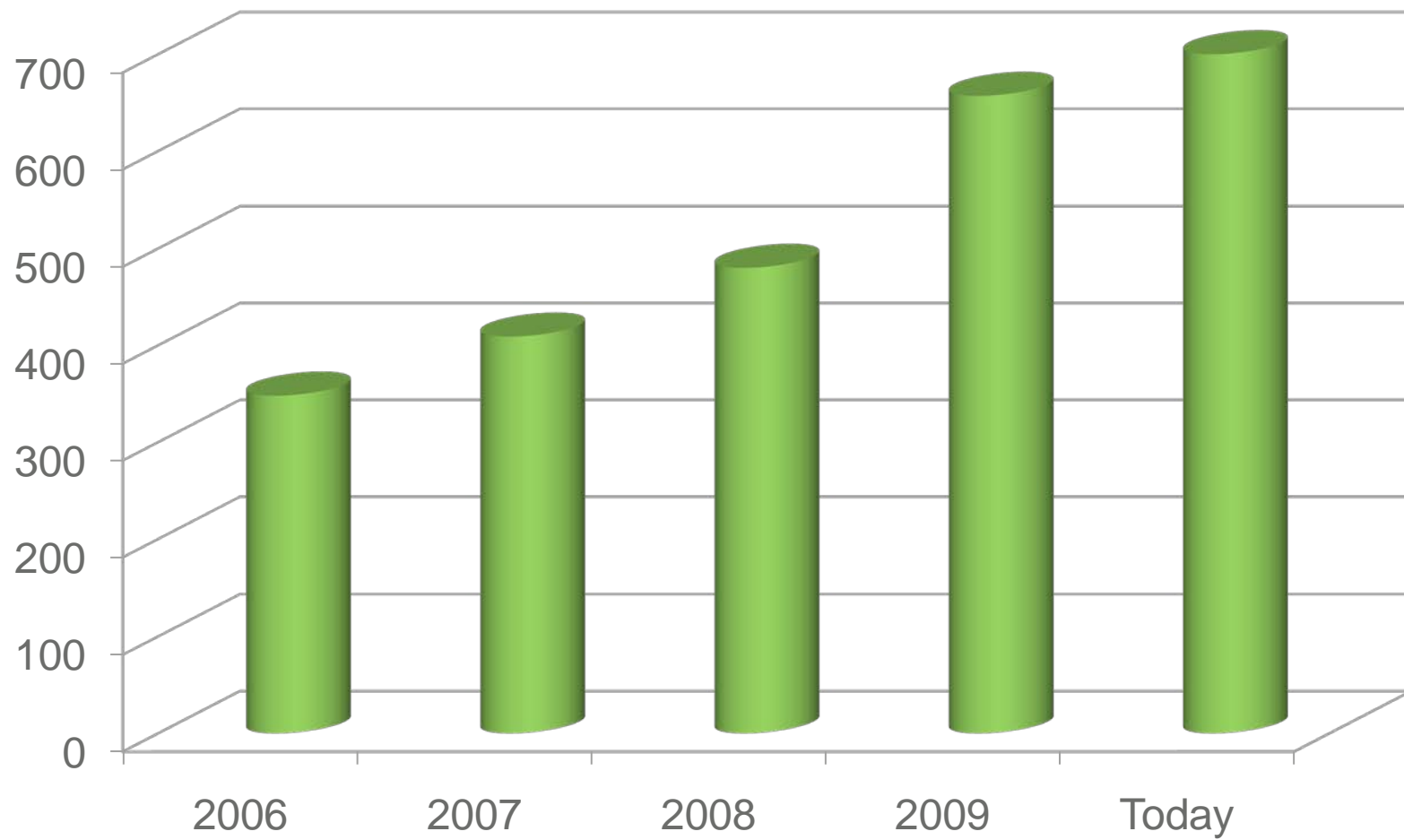


Numbers as of 8/31/2010



PORTICO

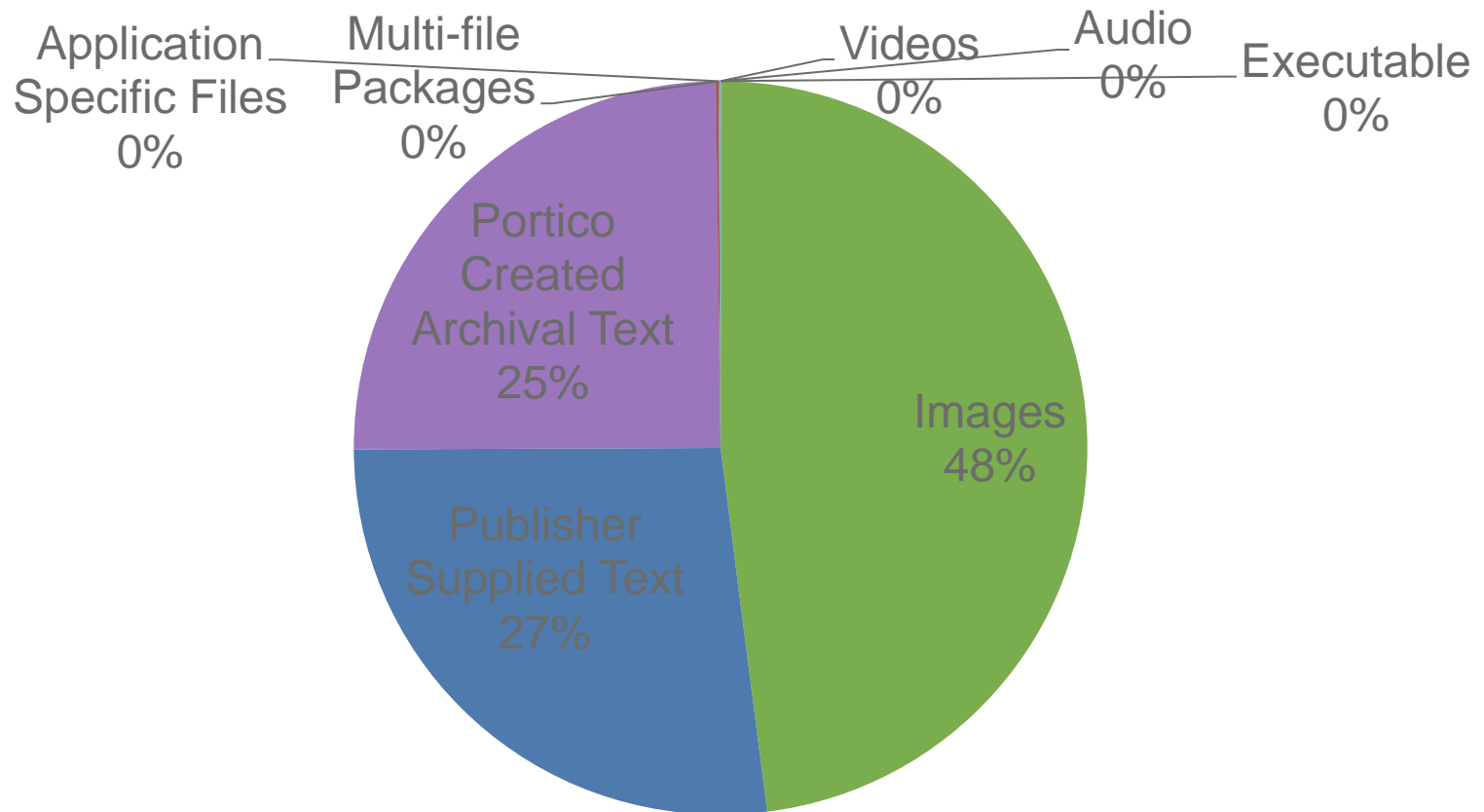
Portico Participating Libraries



Numbers as of 10/29/2010



Types of Files Preserved



Mime Types Preserved

1. application/mathematica
2. application/msword
3. application/octet-stream
4. application/pdf
5. application/postscript
6. application/rtf
7. application/sgml
8. application/vnd.corel-presentations
9. application/vnd.ms-excel
10. application/vnd.ms-htmlhelp
11. application/vnd.ms-powerpoint
12. application/vnd.openxmlformats-officedocument.wordprocessingml.document
13. application/vnd.rn-realmedia
14. application/vnd.wordperfect
15. application/x-asp
16. application/x-gzip
17. application/x-mathcad
18. application/xml
19. application/xml-dtd
20. application/xml-external-parsed-entity
21. application/x-ptc-els-Application Specific Filesset-toc-snippet
22. application/x-ptc-els-Application Specific Filesset-toc-xml-snippet
23. application/x-ptc-eps
24. application/x-ptc-exe
25. application/x-ptc-gams
26. application/x-ptc-msoffice
27. application/x-ptc-netlogo
28. application/x-ptc-nexus
29. application/x-ptc-paintshoppro
30. application/x-ptc-r
31. application/x-ptc-stata-Application Specific Files
32. application/x-ptc-stata-program
33. application/x-ptc-tsp
34. application/x-ptc-utf16
35. application/x-ptc-utf8
36. application/x-rar-compressed
37. application/x-sgml-external-parsed-entity
38. application/x-sh
39. application/x-shockwave-flash
40. application/x-tar
41. application/zip
42. audio/mpeg
43. audio/x-ms-wma
44. audio/x-wav
45. image/gif
46. image/jpeg
47. image/png
48. image/tiff
49. image/vnd.adobe.photoshop
50. image/x-ms-bmp
51. image/x-wmf
52. model/vrml
53. text/csv
54. text/html
55. text/plain
56. text/x-c++src
57. text/x-csrc
58. text/x-ptc-iso-8859
59. video/avi
60. video/mp4
61. video/mpeg
62. video/quicktime
63. video/x-flv
64. video/x-ms-wmv

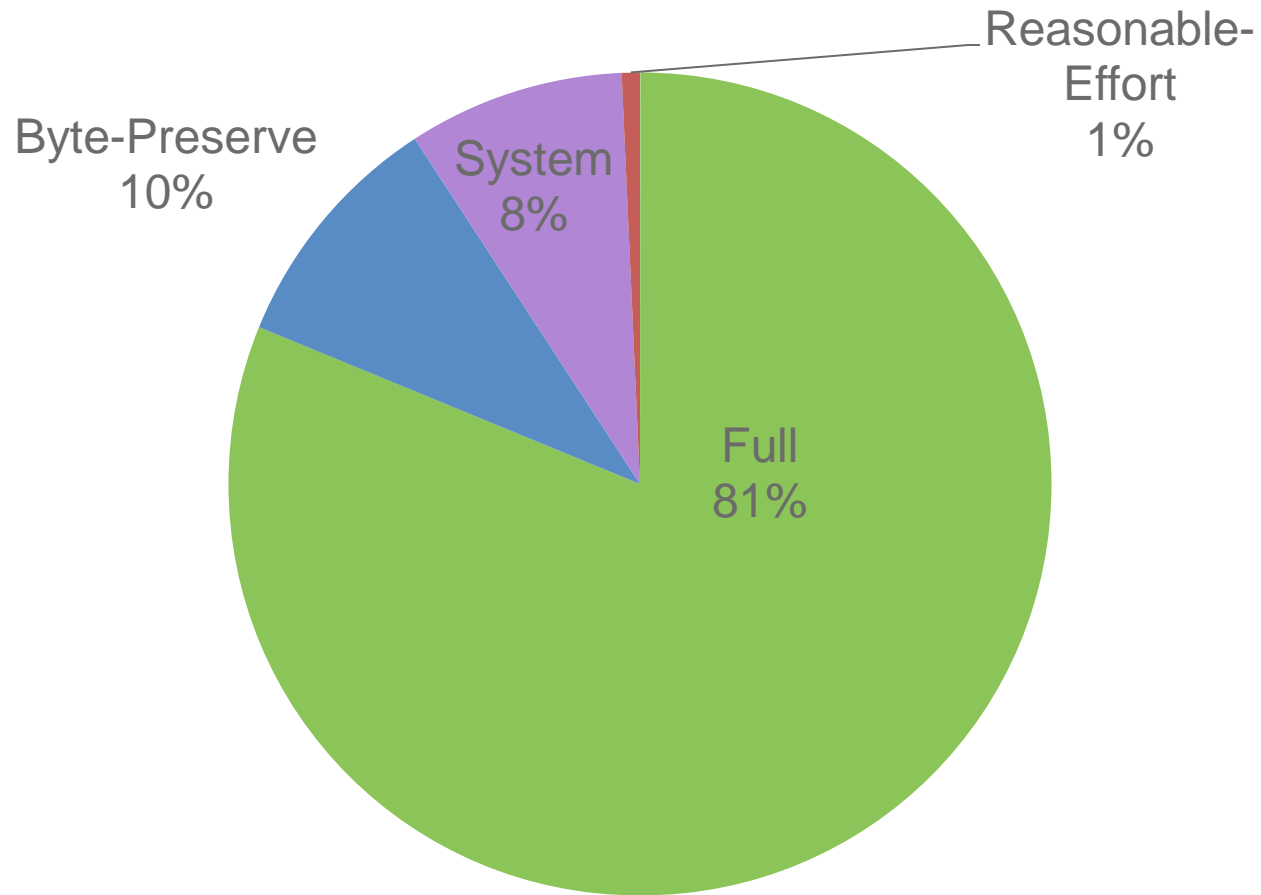


Preservation Levels on Files Preserved

Preservation Level	Files	Percent
Full	142,079,610	81.22%
Byte-Preserve	16,738,528	9.57%
System	14,869,679	8.50%
Reasonable-Effort	1,244,811	0.71%
Total	174,932,628	100.00%



Preservation Levels on Files Preserved

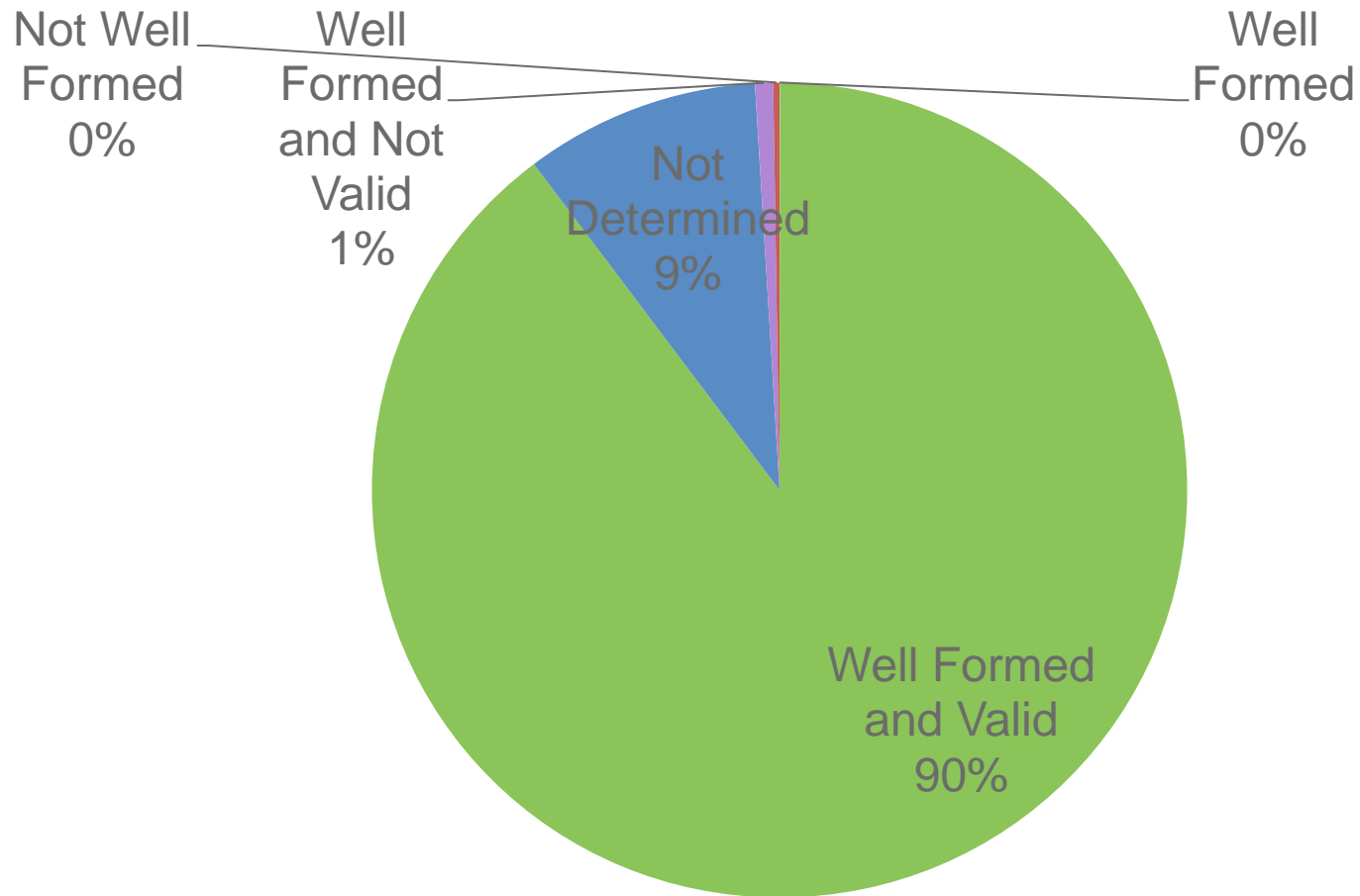


Format Status of Files Preserved

Format Status	Files	Percent
Well Formed and Valid	156,948,510	89.72%
Not Determined	16,304,477	9.32%
Well Formed and Not Valid	1,245,314	0.71%
Not Well Formed	434,074	0.25%
Well Formed	253	0.00%
Total	174,932,628	100.00%



Format Status of Files Preserved



Content Types of Files Preserved

Content Type	Files	Percent
E-journal Files	174,517,812	99.76%
Supplied E-journal Files	304,794	0.17%
E-book Files	108,829	0.06%
Technical Artifact Files	938	0.00%
Business Artifact Files	255	0.00%
Total Files	174,932,628	100.00%



OAIS-compliant repository designed for managed preservation

Key influences:

- » OAIS, GDFR, PRONOM, PREMIS, METS, DC, NLM (JATS), MPEG-21 DIDL, ARK

Key technologies:

- » XML, XML schema, Schematron, JHOVE, NOID
- » Documentum, Oracle, Java, JMS, LDAP
- » Format Registry



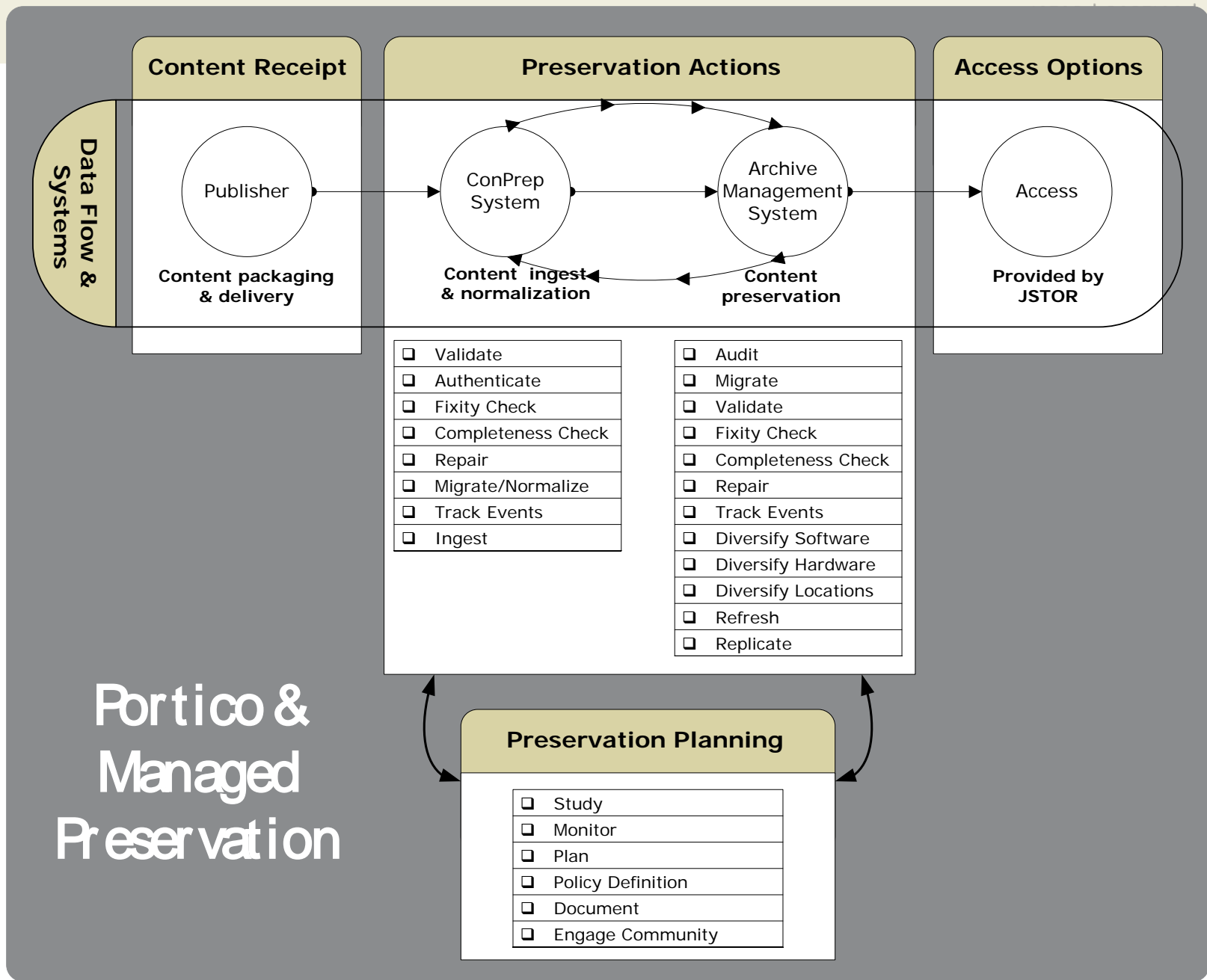
Archive design goals:

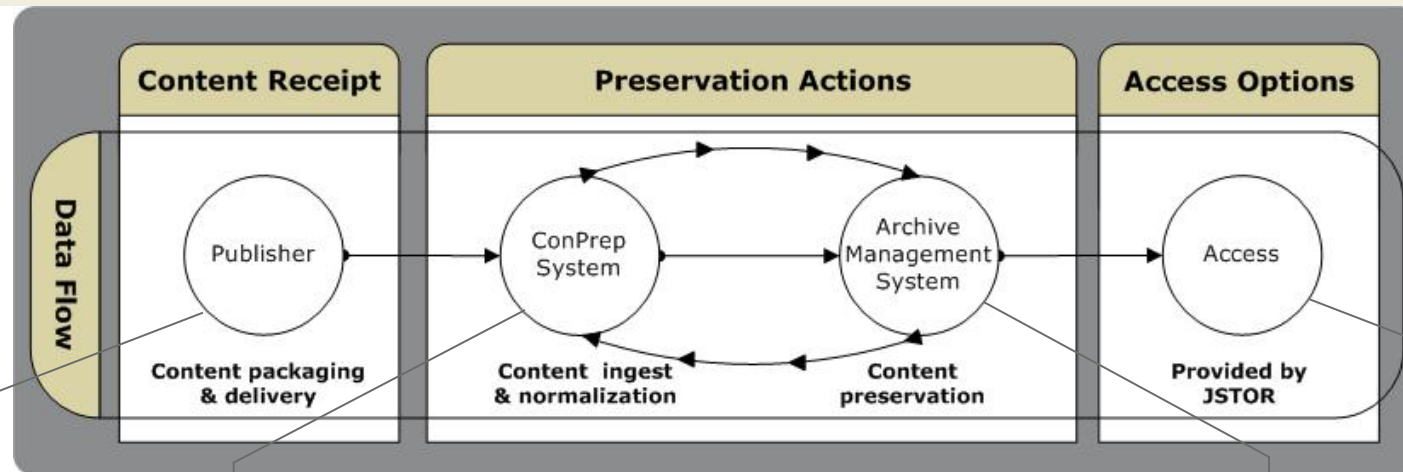
- » Content preserved in application-neutral form using open standards
 - METS, PREMIS, JHOVE
- » A “Bootstrapable Archive”: XML plus Digital Objects
 - Cached in Documentum and Oracle; replicated on file systems

Ingest system design goals:

- » Pluggable tools to facilitate new providers and replacement tools
- » Configurable workflows for different content types
- » Scalable to very high content volumes
- » Built on Documentum workflows







- Publisher supplies XML Source file (including the text, images) and PDF page rendition.
- Best approach for preserving the intellectual content of the article or book.

- Authenticate: verify that preserved content is what it purports to be.
- Verify format: ensure the file meets syntactic and semantic rules of format specification.
- Repair
- Normalize (XML)
- Create preservation metadata

- Assess archival robustness of file format.
- Migrate files to ensure future usability of content.
- Replicate objects and metadata to protect against bit rot and media deterioration

- Render articles to meet viewing requirements of delivery platform.

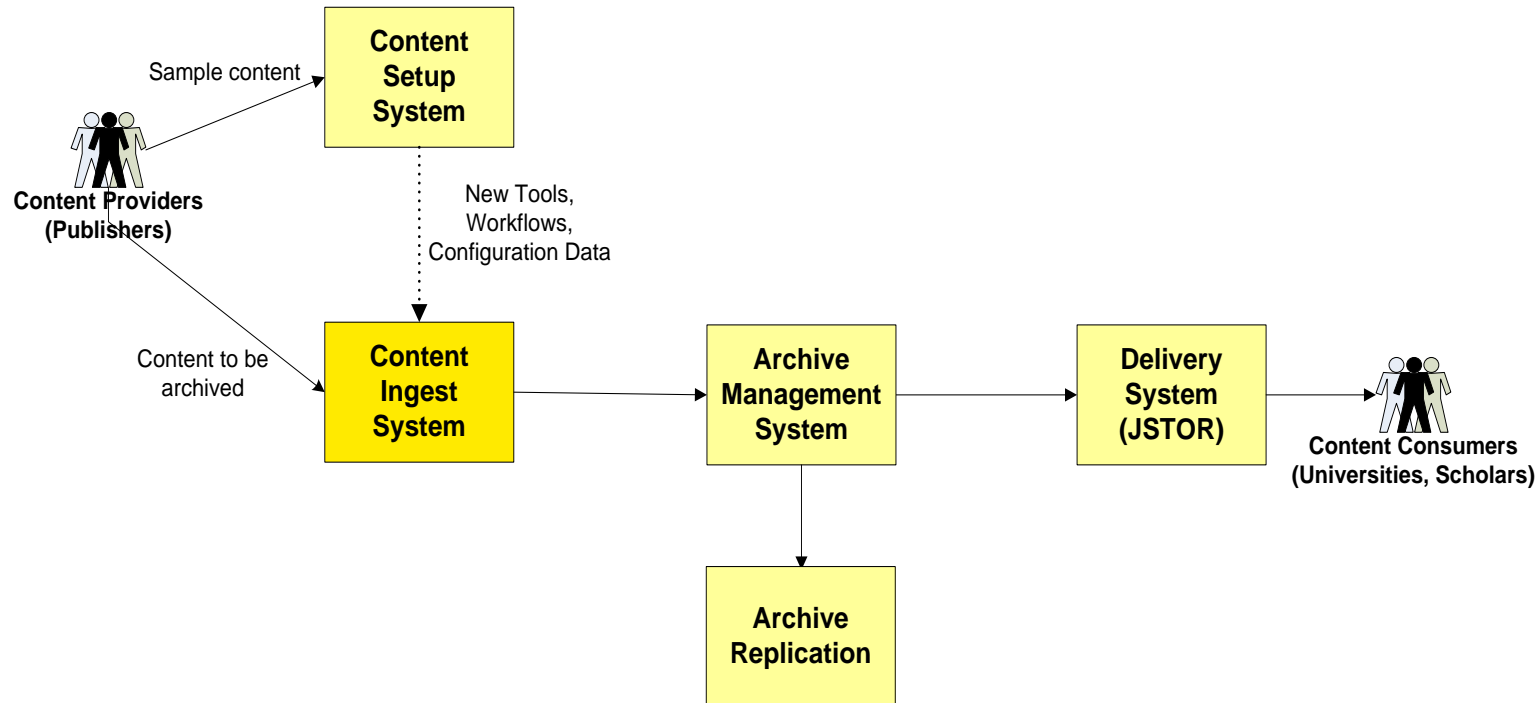


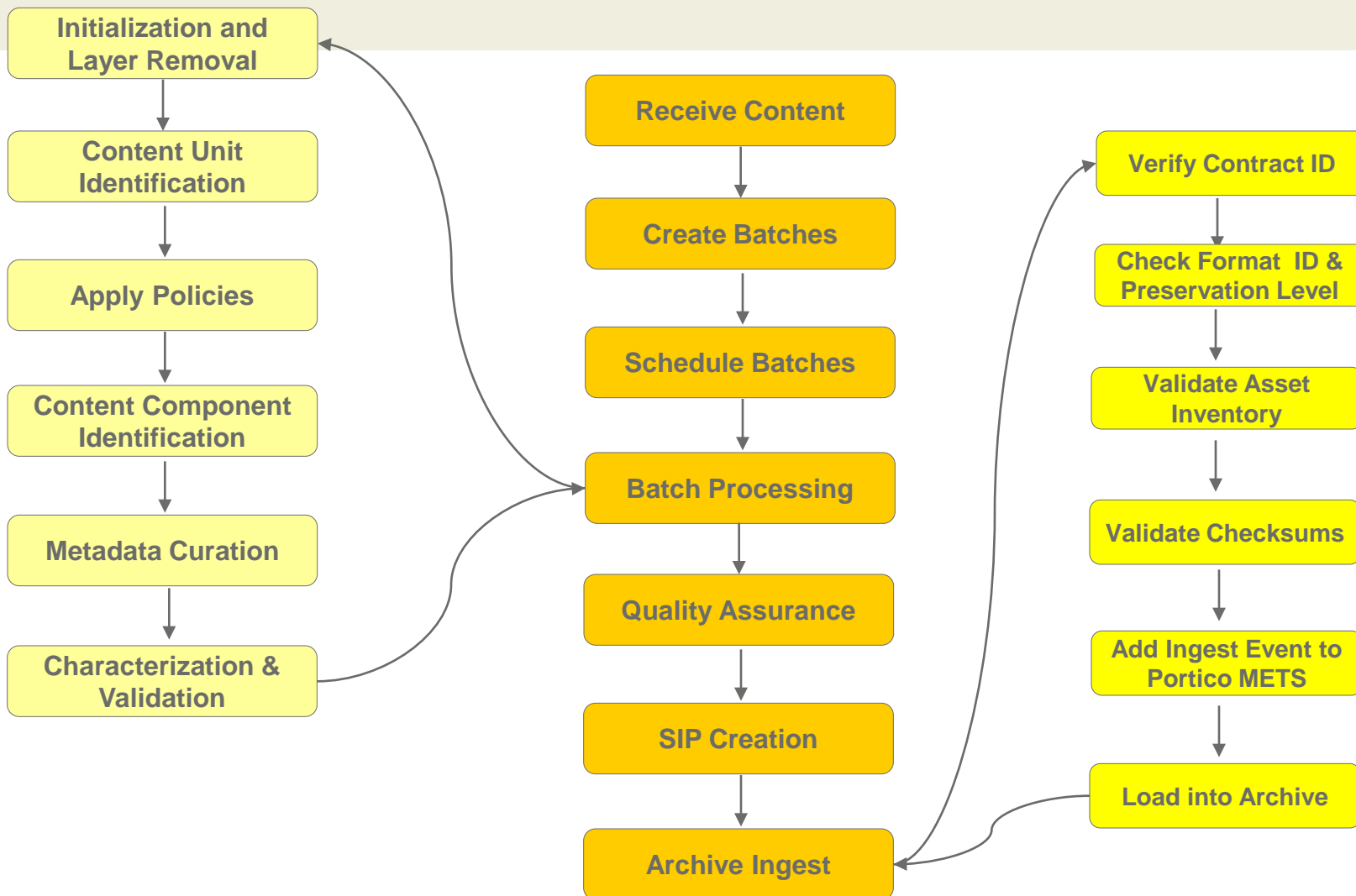
Portico E-Journal/E-Book Preservation Process

- » Interviews with publisher production and technology staff
 - Formats used, production process, content delivered
 - Number of different types of content
 - Updates
 - Supplemental files
- » Large sample data evaluation
- » Formal (written) preservation action plan for each publisher
- » Tool development (as needed per preservation plan)
- » Extensive automated QC during ingest



Portico Systems Overview







- The Problem: Publisher content with arbitrary naming conventions and packaging rules
- The Solution: Profiles containing publisher-specific policies, defaults, and overrides
- The Implementation: XML instances, java regular expression patterns to tokenize complex file names

Design and Coding by Roland Mesde

EXCERPTS FROM SGML TEXT:

The following statistical model was fitted to the data: $$ in which <I>T</I> equals: 1 if Grade 3 or 4 neutropenia was present

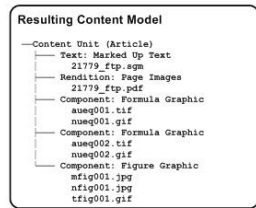
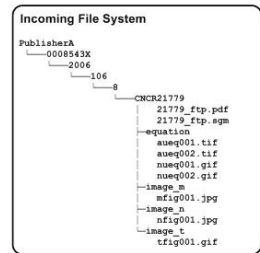
The overall survival for all patients is illustrated in Figure <FIG RREF="fig1"></FIG>.

<FIG ID="fig1" LOC="FLOAT"><GRAPHIC NAME="fig001"></GRAPHIC><NUMBER><CAPTION></CAPTION></FIG>Overall survival for 160 eligible and evaluable patients with recurrent solid tumors who were enrolled on Children's Cancer Group Study 0962. </CAPTION></FIG>

EXCERPTS FROM NORMALIZED XML TEXT:

The following statistical model was fitted to the data: $$ equals: 1 if Grade 3 or 4 neutropenia was present

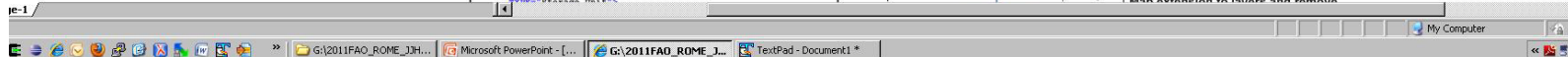
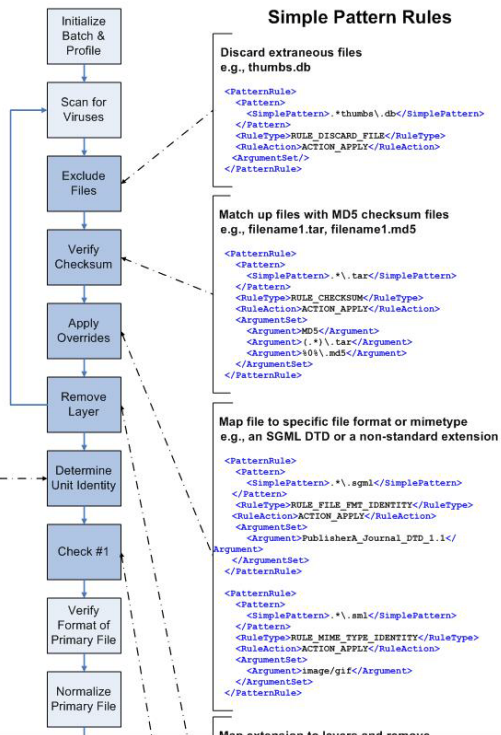
The overall survival for all patients is illustrated in Figure <ref id="FIG1" ref-type="fig"></ref></p>
 <fig fig-type="figure" id="FIG1" position="float"><label>$Figure$</label><caption></caption></fig>Overall survival for 160 eligible and evaluable patients with recurrent solid tumors who were enrolled on Children's Cancer Group Study 0962. </caption><graphic position="anchor" xlink:href="ark:/27927/pc01mabqj"></graphic></fig>



METS Fragment (details simplified)

```

    <div CONTENTID="ark:/27927/pc01mabqj"
    LABEL="Component: Figure Graphic"
    TYPE="Functional Unit">
    <div STATUS="ACTIVE"
    PRESERVATIONLEVEL="FULL"
    CONTENTID="ark:/27927/pc01mabqj"
    LABEL="File"
    PRESERVATIONLEVEL="FULL"
    </div>
    </div>
    
```



Digital Preservation is Everyone's Problem ...

MEANS: "YOU ARE NOT ALONE!!"



Standards and Community Activities

- » NLM DTD Advisory Board
- » NISO Standards Architecture Committee
- » NISO Journal Article Versions Working Group (completed)
- » PREMIS Working Group (completed)
- » Global Digital Format Registry (now UDFR)
- » PEPRS (Piloting an e-journals preservation registry service)
- » DPC (Digital Preservation Coalition)
- » NDSA (National Digital Stewardship Alliance)



Grant-Funded Projects

- » NDIIPP Grant to Portico
- » JISC Digitisation Programme Preservation Study
 - Univ. of London Computer Centre, Portico, and Digital Preservation Coalition
- » IMLS Project on digital book preservation
 - Cornell Univ. and Portico
- » JHOVE2 project (NDIIPP-funded)
 - California Digital Library, Portico, Stanford Digital Repository



Internal Projects

- » E-Book study
- » Library-created content study
- » Portico Preservation Metadata 2.0



The Center for Research Libraries (CRL) conducted a preservation audit of Portico (www.portico.org) between April and October 2009 and hereby certifies Portico as a trustworthy digital repository. CRL has found that Portico's services and operations basically conform to the requirements for a trusted digital repository. The CRL Certification Advisory Panel has concluded that the practices and services described in Portico's public communications and published documentation are generally sound and appropriate to the content being archived and the needs of the CRL community. Moreover the CRL Certification Advisory Panel expects that in the future, Portico will continue to be able to deliver content that is understandable and usable by its designated user community.

This certification is based upon a site visit and sampling of archives content, and upon the review of information gathered by CRL and its Certification Advisory Panel and of documents and documentation provided by Portico. CRL's analysis was guided by the criteria included in the *Trustworthy Repositories Audit and Certification* checklist, and other metrics developed by CRL through its analyses of digital repositories.

CRL conducted its audit with reference to generally accepted best practices in the management of digital systems; and with reference to the interests of its community of research libraries and the practices and needs of scholarly researchers in the humanities, sciences and social sciences in the United States and Canada. The purpose of the audit was to obtain reasonable assurance that Portico provides, and is likely to continue to provide, services adequate to those needs without material flaws or defects and as described in Portico's public disclosures. The CRL audit provides a reasonable basis for these findings.

CRL has assigned Portico the following levels of certification (the numeric rating is based on a scale of 1 through 5, with 5 being the highest level, and 1 being the minimum certifiable level):¹

CATEGORY	PORTICO SCORE
Organizational Infrastructure	3
Digital Object Management	4
Technologies, Technical Infrastructure, Security	4



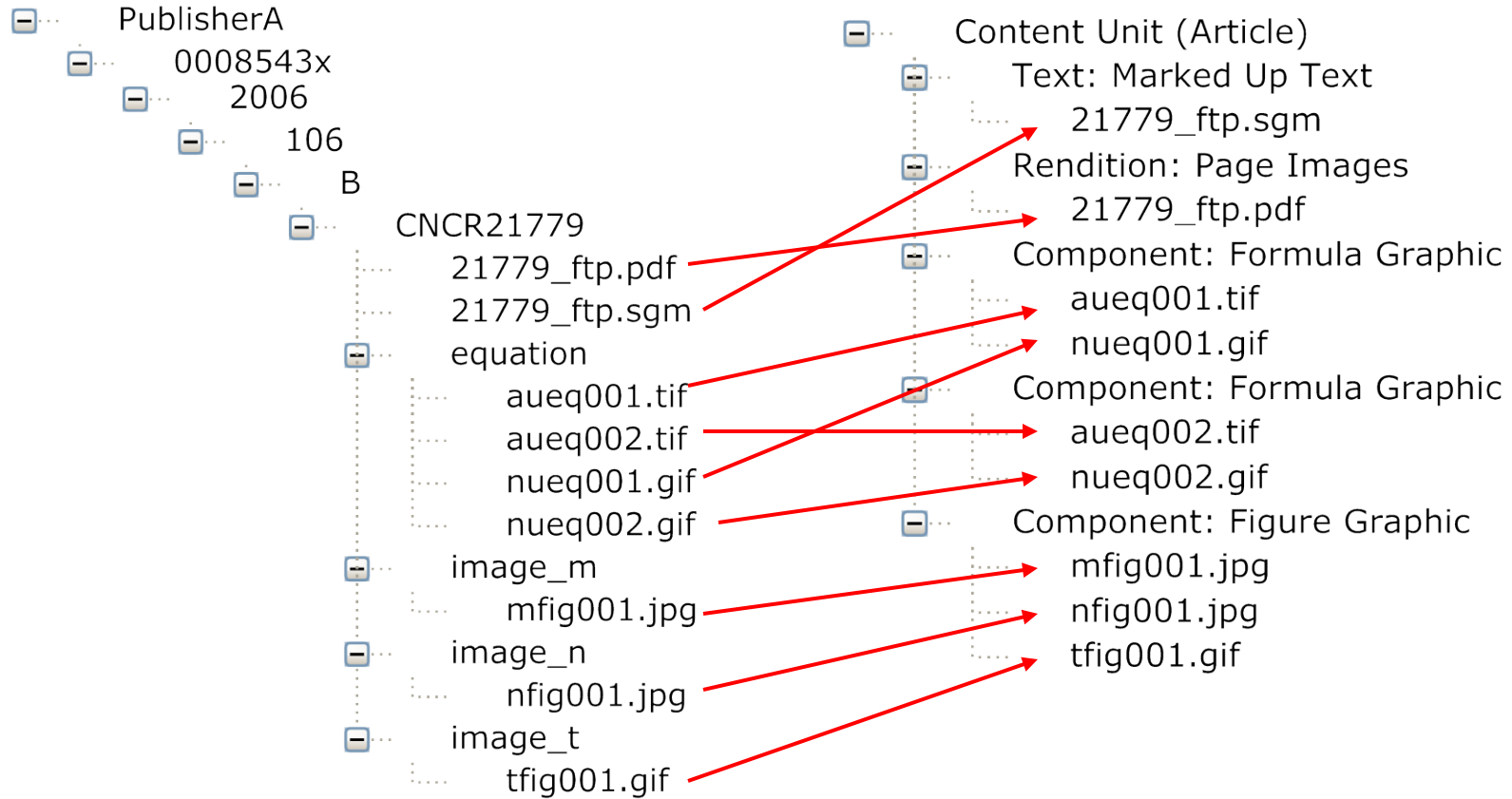
Life is messy...

... FOR EVERYONE!!



Incoming File System

Resulting Content Model



Content isn't perfect

- » Must have policies and workflow for invalid data
- » There are degrees of “badness”
- » Strict format validity does not equate to usefulness or usability
 - E.g., Well-formed but not valid PDF
 - E.g., Valid PDF with bad embedded font
 - E.g., Invalid JPEG

Content creation practices change over time

- » Publishers (content providers) aren't consistent
- » Or don't warn you that they are changing something
- » Defensive programming required

Software isn't perfect

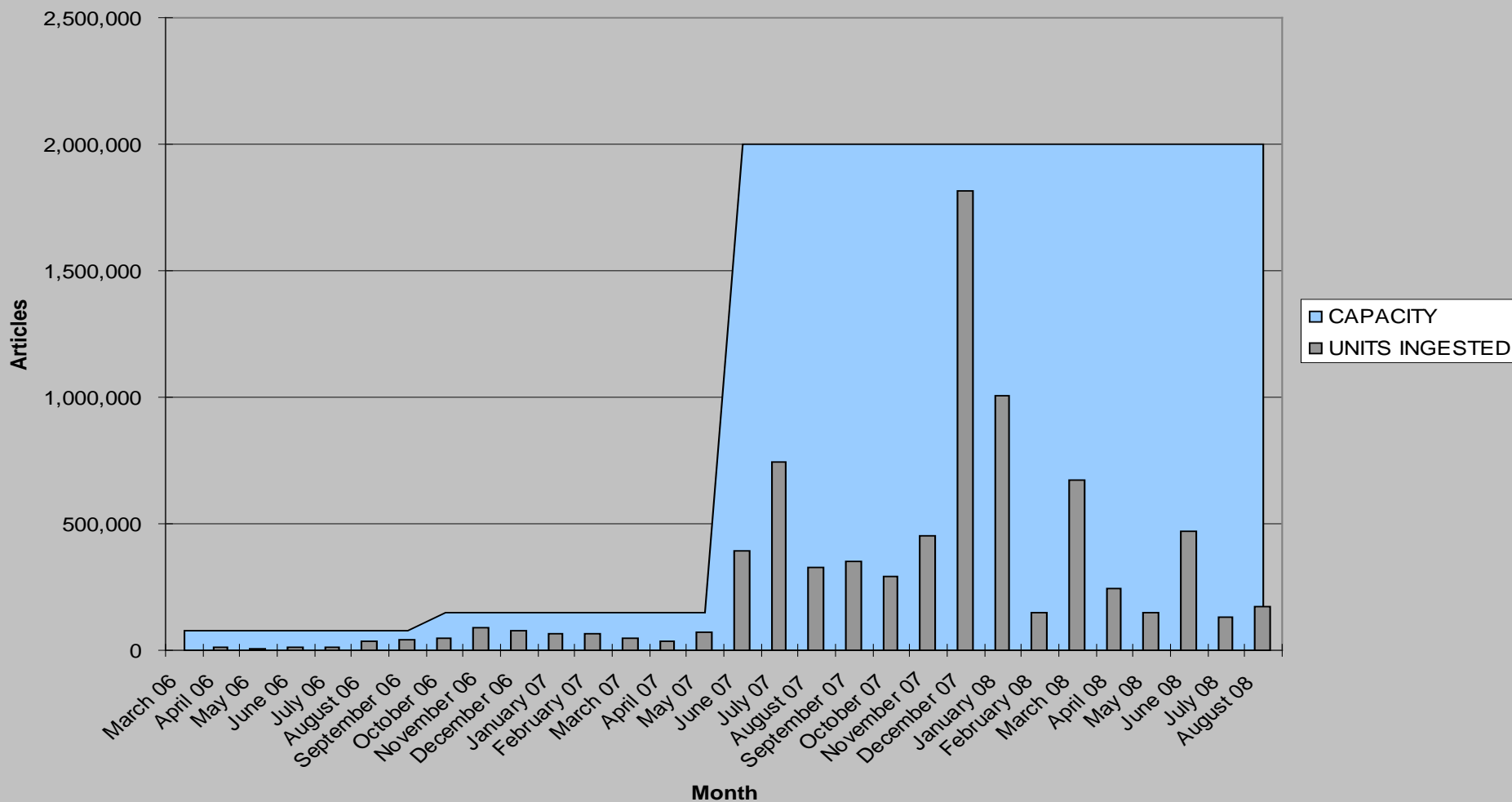
- » Assume that there will be internal failures
- » Reversibility and audit trail are essential



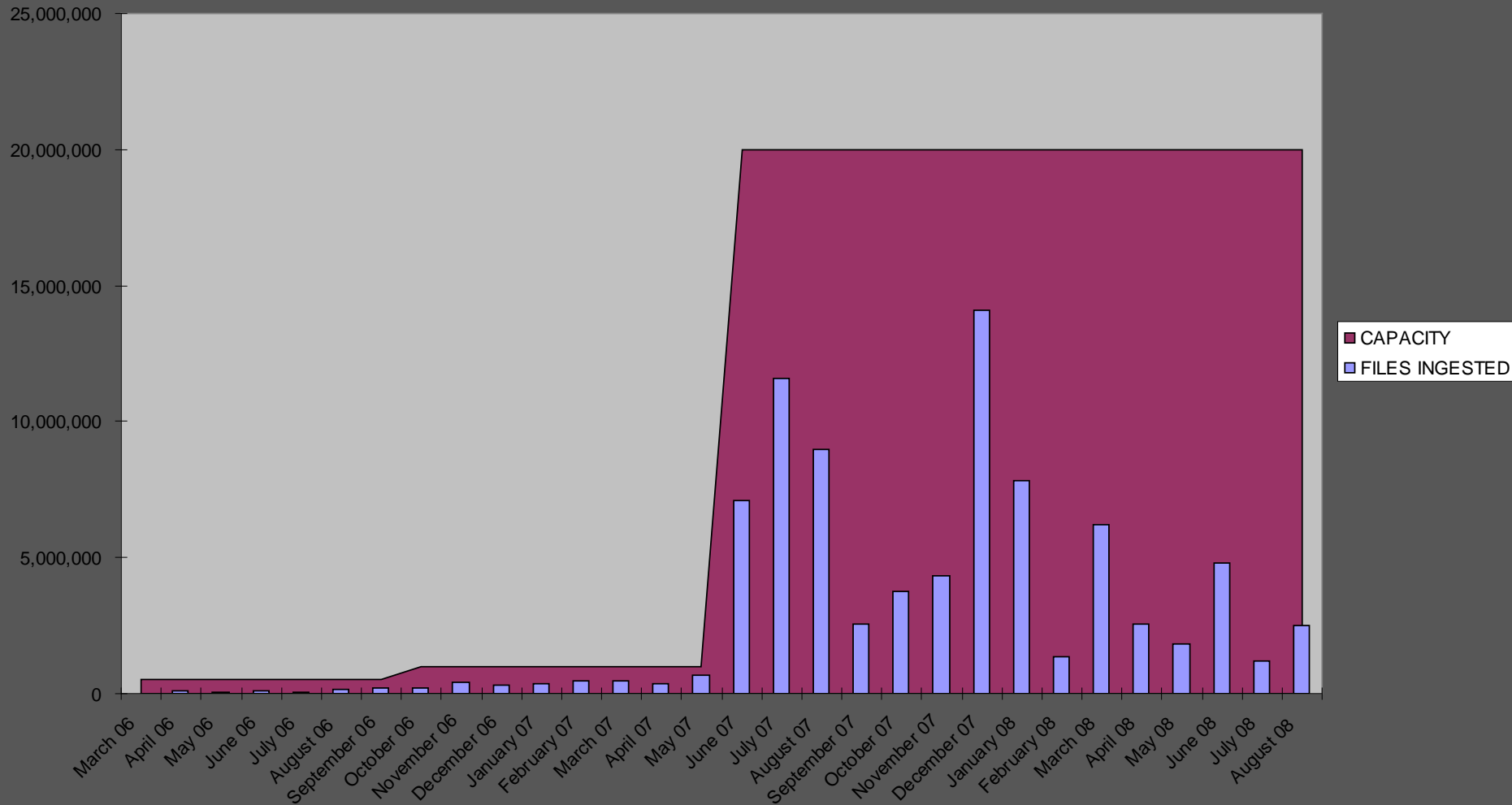
- » Scale up from 900K articles/year to 10 million articles/year
- » Involved changes to
 - Software
 - Hardware
 - Procedures
- » Testing, tuning
 - How many threads?
 - Good data, bad data
 - More batches? Bigger batches?
 - Long-running tests
- » Side effects
 - Loaders
 - Cleanup
 - Logging
 - User interface
 - Storage backup and recovery



Monthly Article Ingest versus Capacity



Month File Ingest Versus Theoretical Capacity



■ CAPACITY
■ FILES INGESTED

- » Portico Web Site
 - Portico TRAC self-audit
 - Portico Policy Documents
- » CRL Audit Report
- » NEH Report on Preservation of Books and Other Digital Content
- » Blue Ribbon Task Force on Sustainable Digital Preservation and Access
- » NLM Journal Archiving and Interchange Tag Suite



PORTICO

THANK YOU.

Sheila Morrissey
sheila.morrissey@ithaka.org

<http://www.portico.org>



PORTICO