

This is a preprint of a paper to be published in the proceedings of RCDL'2013 "Digital Libraries: Advanced Methods and Technologies, Digital Collections" [«Электронные библиотеки: перспективные методы и технологии, электронные коллекции»], held on October 14–17, 2013 in Yaroslavl, Yaroslavl Oblast, Russia.

# A Model for Integrating the Publication and Preservation of Journal Articles

Kevin S. Hawkins  
University of Michigan, Ann Arbor

*There are policy, technical, and workflow gaps in library efforts to preserve online journal literature. Since libraries are increasingly involved in journal publishing, HathiTrust, a shared preservation-quality digital repository, is a natural place to archive and provide access to journal literature to ensure its long-term preservation and discoverability. The U-M Library is funding the creation of mPach, an open-source, end-to-end publishing system in which archiving in HathiTrust happens as a byproduct of publication rather than being carried out after the fact. The architecture of mPach, its envisioned workflow, and plans for creating a shared infrastructure for publishing open-access journals are all summarized.*

## The deficit in journal preservation

Until quite recently, publishers produced documents on physical media, and libraries acquired and preserved copies of these documents. But in the era of the Internet, when publishers host content online, the library's role in acquiring and preserving the content is in jeopardy: without special licensing arrangements such as those often provided by open-access journals, a library has no legal right to make a copy of the content for preservation.

Various business models have evolved to address this situation, especially for journals, which are increasingly available only online. For non-open-access journals, research libraries often negotiate the right to create a digital copy of any content acquired during the period of subscription<sup>1</sup> and make this content available only to their patrons,<sup>2</sup> though few are equipped to provide this kind of restricted access and archiving with integrated browse and search functions. To address the more pressing concern of publishers going out of business without *any* libraries holding a copy of the content, libraries and publishers have collaborated in initiatives like LOCKSS,<sup>3</sup> CLOCKSS,<sup>4</sup> and Portico<sup>5</sup> in order to guarantee that one or more copy of the content will become available if it is no longer available from the publisher. Similarly, the Koninklijke Bibliotheek and Elsevier reached an agreement in 2002 whereby the KB will preserve Elsevier journals under terms similar to those governing journals that use LOCKSS, CLOCKSS, and Portico.<sup>6</sup> Still, there are problems with these models. LOCKSS and

---

<sup>1</sup> Sadie L. Honey. Preservation of electronic scholarly publishing: an analysis of three approaches. *Portal: Libraries and the Academy*, 5(1):59-75, Jan. 2005.

<sup>2</sup> NISO SERU Standing Committee, SERU: A Shared Electronic Resource Understanding: A Recommended Practice of the National Information Standards Organization. National Information Standards Organization (NISO), May 2012.  
[http://www.niso.org/publications/rp/RP-7-2012\\_SERU.pdf](http://www.niso.org/publications/rp/RP-7-2012_SERU.pdf).

<sup>3</sup> Lots of Copies Keeps Stuff Safe. <http://www.lockss.org/>.

<sup>4</sup> CLOCKSS. <http://www.clockss.org/>.

<sup>5</sup> Portico. <http://www.portico.org/>.

<sup>6</sup> National Library of the Netherlands and Elsevier Science make digital preservation history: permanent digital archive assures perpetual accessibility of scientific heritage. August 20, 2002.  
<http://www.kb.nl/en/news/news-archive-2002/national-library-of-the-netherlands-and-elsevier-science-make-digital-preservation-history>.

CLOCKSS use web crawling, which captures only the appearance of webpages but not their underlying structure or search functionality. Portico and the KB, on the other hand, rely on publishers to deliver journal articles in valid file formats, and not just the version first published but also any corrected versions of these articles.

One way to ensure that a library always has access to the latest content is for the library to operate the very system used to publish the journal. A survey in 2010 of a cross-section of North American academic libraries found that, of 144 responding institutions, 43 offered “operational publishing services” to their scholars at the institution.<sup>7</sup> Of these 43 institutions, most host publications using open-source software such as Open Journal Systems (OJS)<sup>8</sup> or DSpace,<sup>9</sup> while about a quarter use Digital Commons,<sup>10</sup> a hosted platform provided by bepress. Unfortunately, all of these platforms deliver to users only those files (primarily PDF files) created and uploaded by a journal editor. Since the library is not in a position to control the software and workflows used to create these files, the library can only provide bitwise preservation of the files, severely hampering future migration of the content.

## A higher standard for preservation

Since libraries are increasingly involved in journal publishing, HathiTrust,<sup>11</sup> a shared preservation-quality digital repository, is a natural place to archive and provide access to journal literature to ensure its long-term preservation and discoverability. HathiTrust already archives and provides access to reformatted library holdings, but the University of Michigan Library, a founding member of HathiTrust, sees an opportunity to use HathiTrust for publishing born-digital journals as well. To develop an infrastructure in support of low-cost university-based publishing that addresses the needs and values of both content creators and librarians, the U-M Library is funding the creation of mPach,<sup>12</sup> an open-source, end-to-end publishing system in which the act of publishing and the act of archiving are unified. In other words, archiving in HathiTrust happens as a byproduct of publication rather than being carried out after the fact. mPach leverages existing components of HathiTrust and available open-source software where appropriate.

Archiving is not as simple as saving a copy of a file produced by a journal editor, as OJS and institutional repositories generally do. Instead, the content needs to be stored in a format that allows digital preservation. PDF/A, a non-proprietary variant of the PDF family standardized as ISO 19005, is often suggested for such needs, but even a PDF/A file is poorly suited for use with screen readers for the visually impaired and for any non-paginated display, and is suboptimal even for searching and data mining.

---

<sup>7</sup> James L. Mullins, Catherine Murray-Rust, Joyce L. Ogburn, Raym Crow, October Ivens, Allyson Mower, Daureen Neddill, Mark Newton, Julie Speer, and Charles Watkinson. *Library Publishing Services: Strategies for Success: Final Research Report*. March 2012.  
<http://wp.sparc.arl.org/lps/>.

<sup>8</sup> Open Journal Systems. <http://pkp.sfu.ca/ojs/>.

<sup>9</sup> DSpace. <http://www.dspace.org/>.

<sup>10</sup> Digital Commons. <http://digitalcommons.bepress.com/>.

<sup>11</sup> HathiTrust Digital Library. <http://www.hathitrust.org/>.

<sup>12</sup> mPach. <http://www.lib.umich.edu/mpach>.

Rather than preserving the paginated appearance of a document, the text of the article needs to be stored in a format that reflects its structure and semantics, with associated media in formats that can be preserved and rendered. mPach has developed a specification for journal articles that uses the Journal Article Tag Suite (JATS), an application of NISO Z39.96-2012,<sup>13</sup> for the text and stores this with high-quality versions of media objects and with a METS record containing structural and preservation metadata.

## An overview of mPach

There are three major parts of mPach (see also figure 1), each of which includes components in various stages of development at the time of writing:

- **the peer review and editorial system:** what authors and reviewers interact with
- **Prepper:** what prepares the article for ingest into HathiTrust for archiving and publication
- **modified HathiTrust components:** various modifications to existing components of the HathiTrust environment to support born-digital journal articles

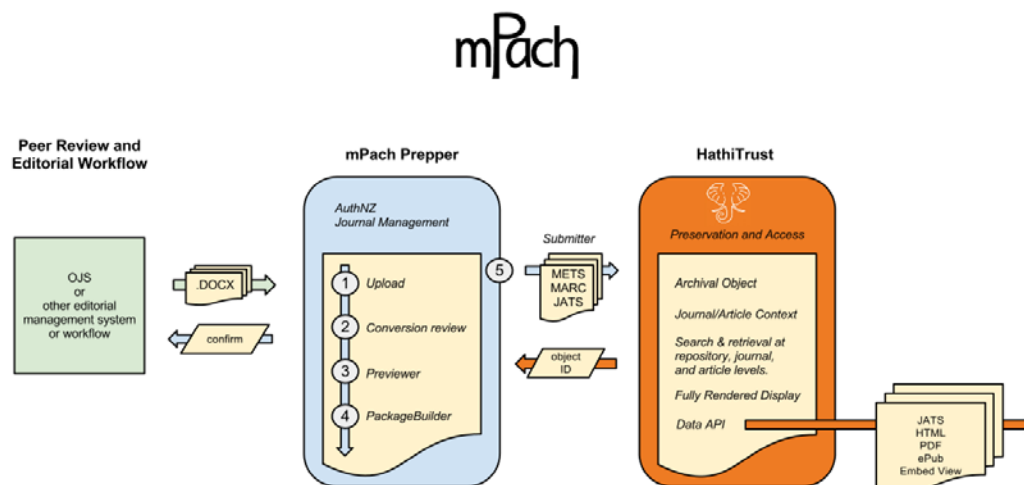


Figure 1: Major parts of mPach

As a modular system, mPach could be used with any peer review and editorial system that is capable of interacting with Prepper; however, the developers have chosen to provide OJS as the default option. Despite having no support for digital preservation, OJS is already widely used for library-based journal publishing, and mPach's integration with this software will allow for a smooth transition of journals already published using OJS into the HathiTrust repository. Integration with mPach requires that manuscripts that reach the "layout" stage in OJS be sent to Prepper, which prepares the HathiTrust Submission Information Package (SIP).

Prepper provides a user interface for the editor of a journal: a dashboard for administering the journal and putting manuscripts through a production process—akin to composition and typesetting—that prepares all content

<sup>13</sup> Journal Article Tag Suite. <http://jats.nlm.nih.gov/>.

according to the preservation standard developed for mPach content in HathiTrust. Prepper invokes Norm, a Python application developed to convert manuscripts from Office Open XML (“DOCX”) format<sup>14</sup> into XML that conforms to JATS. DOCX is the default option because, like OJS, it is widely used in the editorial process of journals published by libraries. The Prepper interface also guides the staff member through a review of validation errors detected by Norm’s conversion, uploading high-resolution figures, supplying “alt text” for figures, previewing the article as rendered using the default stylesheet (based on the Preview XSLT stylesheets<sup>15</sup>), uploading supplementary material,<sup>16</sup> and submitting for ingest into HathiTrust.

mPach requires a number of significant modifications to HathiTrust components and workflows originally designed to support reformatted print materials. The reading interface in HathiTrust, which previously supported only rendering of digitized page images, renders JATS XML in HTML and allows a user to download a dynamically generated PDF and EPUB, display metadata specific to articles (figure 2), and link to a special “collection” for the journal in HathiTrust’s Collections application<sup>17</sup> that allows for browsing volumes and issues of the journal (figure 3).

---

<sup>14</sup> Office Open XML. Wikipedia. [http://en.wikipedia.org/wiki/Office\\_Open\\_XML](http://en.wikipedia.org/wiki/Office_Open_XML).

<sup>15</sup> NISO Journal Article Tag Set (JATS) version 1.0: Preview XSLT stylesheets. <https://github.com/NCBITools/JATSPreviewStylesheets>.

<sup>16</sup> Recommended Practices for Online Supplemental Journal Article Materials: a recommended practice of the National Information Standards Organization and the National Federation of Advanced Information Services. January 2013. <http://www.niso.org/publications/rp/rp-15-2013>.

<sup>17</sup> Collections. HathiTrust Digital Library. <http://babel.hathitrust.org/cgi/mb>.

The screenshot displays the HathiTrust digital library interface. At the top, there is a navigation bar with links for Home, About, Collections, Help, and Feedback. Below this is the HathiTrust logo and a search bar with options for 'FULL-TEXT' and 'CATALOG'. A search box contains the text 'Search words about or within the items', and there is a 'LOG IN' button. Below the search bar, there are links for 'Advanced full-text search' and 'Search tips', and a 'Full view only' checkbox.

The main content area is divided into a left sidebar and a main article area. The sidebar contains the following sections:

- Journal of Electronic Publishing collection**: A link to go back to the collection.
- JEP logo**: The logo for the Journal of Electronic Publishing, with the text 'the journal of electronic publishing'.
- Journal of Electronic Publishing**: The journal title.
- Vol. 15, No. 1 (Summer 2012)**: The volume and issue information.
- About this journal**: A link to learn more about the journal.
- About this Article**: A section for article details.
  - Title: Refurbishing the Camelot of Scholarship: How to Improve the Digital Contribution of the PDF Research Article
  - Authors: John Willinsky, Alex Garnett, Angela Pan Wong
  - View full catalog record
  - Copyright: (CC) BY
- Get this Article**: Options to download the article in PDF, XML, or EPUB formats.
- Supplemental Materials**: A link to a Data Set (XLS, 35K).
- Add to Collection**: A section for adding the article to a personal collection, including a 'Login' button and a 'Select Collection' dropdown menu.
- Share**: A section for sharing the article, including a 'Permanent link to this article' and a URL: <http://dx.doi.org/0000.0000.0000>.
- Version**: 2012-07-19 16:37 UTC

The main article area contains the following information:

- Article Title**: Refurbishing the Camelot of Scholarship: How to Improve the Digital Contribution of the PDF Research Article
- Authors**: John Willinsky, Alex Garnett, and Angela Pan Wong
- Volume**: 15, Issue 1, Summer 2012
- DOI**: <http://dx.doi.org/10.3998/3336451.0015.102>
- Permissions**: A link to view permissions.
- Abstract**: A paragraph summarizing the article's content, mentioning the Portable Document Format (PDF) and its evolution.
- Introduction**: The beginning of the article's text, discussing the history of the Portable Document Format (PDF) and the Camelot Project.

Figure 2: Mockup of an article viewed in HathiTrust's user interface

The screenshot shows the HathiTrust Digital Library interface for the Journal of Electronic Publishing (JEP). At the top, there is a navigation bar with links for Home, About, Collections, Help, and Feedback. Below this is the HathiTrust logo and a search bar with a 'FULL-TEXT' button and a 'LOG IN' button. The main content area is divided into two columns. The left column contains the JEP logo, the journal title 'the journal of electronic publishing', and the full title 'Journal of Electronic Publishing'. It also lists the owner as Michigan Publishing, a description of the journal, and the ISSN 1080-2711. The right column features a search bar for the journal, a 'Find' button, and a list of articles. The articles are sorted by 'Date Descending' and include titles such as 'The Short-Term Influence of Free Digital Versions of Books on Print Sales' by John Hilton, III; David Wiley, 'UP 2.0: Some Theses on the Future of Academic Publishing' by Phil Pochoda, 'Our Book' by Sandra Ordenez, 'Launching (and Sustaining) a Scholarly Journal of the Internet: The International Journal of Baudrillard Studies' by Gerry Coulter, 'Justify Just of Just Justify' by Mohamed Elyaaakoubi; Azzeddine Lazrek, 'XML Production Workflows? Start with the Web' by John W. Maxwell; Meghan MacDonald; Travis Nicolson, et al., and 'Editor's Note' by Judith Axler Turner. The list also shows volumes 11 through 16, with volume 13 expanded to show its issues.

Figure 3: Mockup of a journal viewed in HathiTrust's user interface

Discovery of known items in HathiTrust using metadata like title and author is currently provided for by a catalog of MARC records, with one per item in the repository. For mPach, each article has its own analytic catalog record, tied to a monographic record for the journal as a whole. Finally, the HathiTrust Data API<sup>18</sup> allows for the content of each article to be retrieved for use outside of the native HathiTrust interface.

Note that by policy HathiTrust only closes access to content for legal reasons, not because a rightsholder wants to restrict access. Therefore, mPach only supports the publishing of open-access journals.

<sup>18</sup> HathiTrust Data API. [http://www.hathitrust.org/data\\_api](http://www.hathitrust.org/data_api).

## Workflow

In the typical workflow for publishing a journal using mPach, a journal editor uses OJS to manage submissions, peer review, and the editing process. Once an article reaches the “layout” stage (where a combination of composition and typesetting allows the article to be formatted in a consistent way), the journal editor formats it according to a predefined list of styles in Microsoft Word and submits the article in DOCX to mPach’s Prepper, which guides the editor through conversion to JATS XML, preparation of the SIP, and submission for ingest. Prepper keeps track of articles so that a revised version can be submitted for ingest. Currently the ingest process overwrites any previous version of an item with the same identifier, but eventually HathiTrust will archive past versions and allow users to navigate among them.

## mPach as a shared infrastructure

The U-M Library plans to host the Prepper system, including the submission module, to facilitate authorized deposit of content, and will make this system available for use by organizations wishing to publish journal literature in HathiTrust. The developers envision extending the Norm component to handle OpenDocument (“ODT”)<sup>19</sup> and LaTeX as input formats, each of which is more commonly used in certain communities. Furthermore, if the Book Interchange Tag Suite<sup>20</sup> is adopted as a standard, the mPach architecture might be extended to support monograph publishing. While mPach is currently being developed to meet the needs of the U-M Library, the contribution of the sourcecode to the planned HathiTrust Development Environment should foster contributions from developers not at U-M and therefore lead to the creation of a truly shared infrastructure for publishing open-access scholarly journals.

---

<sup>19</sup> OpenDocument. Wikipedia. <http://en.wikipedia.org/wiki/OpenDocument>.

<sup>20</sup> Book Interchange Tag Suite (BITS) 0.2 DRAFT. <http://jats.nlm.nih.gov/extensions/bits/>.